



ML-Based NLP and Medical Text Corpora

Borchert, Dr. Schapranow
Data Management for Digital Health
Winter 2023

Agenda

Pillars of the Lecture

Medical Use Cases



Biology Recap



Oncology



Nephrology



Infectious
Diseases

Technology Foundation



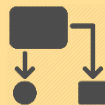
Data
Sources



Data
Formats



Processing and
Analysis



Software
Architectures

Machine Learning

Data



Refine

Evaluate



Prediction +
Probability

ML and Corpora

Data Management for
Digital Health, Winter
2023
2

Agenda

Pillars of the Lecture

Medical Use Cases



Biology Recap



Oncology



Nephrology



Infectious
Diseases

Technology Foundation



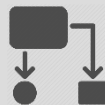
Data
Sources



Data
Formats



Processing and
Analysis



Software
Architectures

Machine Learning

Data



Refine

Evaluate



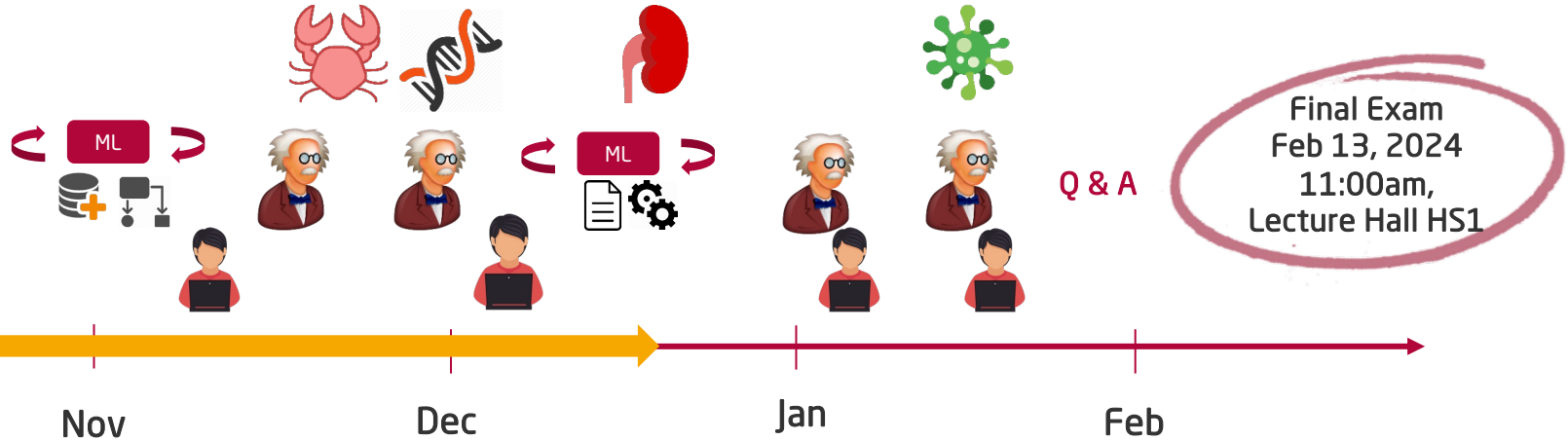
Prediction +
Probability

ML and Corpora

Data Management for
Digital Health, Winter
2023

3

Lecture Schedule



- Lecture Kickoff
- Actors in Healthcare
- Digital Health Data

- Machine Learning (ML) Foundations
- Use Case Oncology
- Biology Recap

- Natural Language Processing
- Use Case Nephrology & Intensive Care
- Supervised ML & Deep Learning

- Use Case Infectious Diseases
- Unsupervised ML

ML and Corpora

Data Management for
Digital Health, Winter
2023

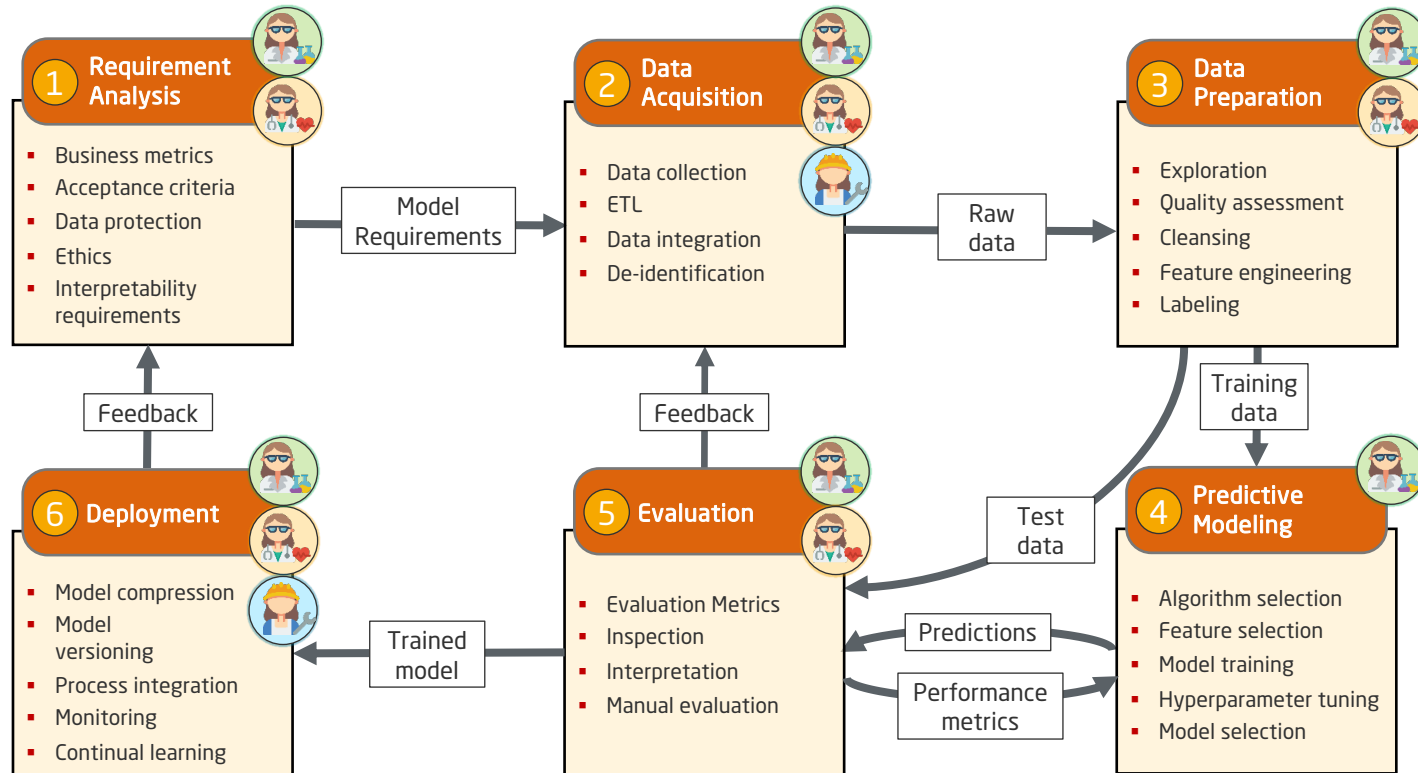
Agenda

- Acquisition and Preparation of Text Corpora
- Feature Extraction from Texts
- Neural Networks
- NER as Sequence Tagging
- Evaluation of NER Models

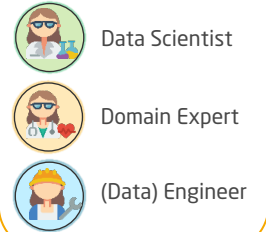
Text Data & NLP

Data Management for
Digital Health, Winter
2023

Revisiting the Process Model for ML in Digital Health



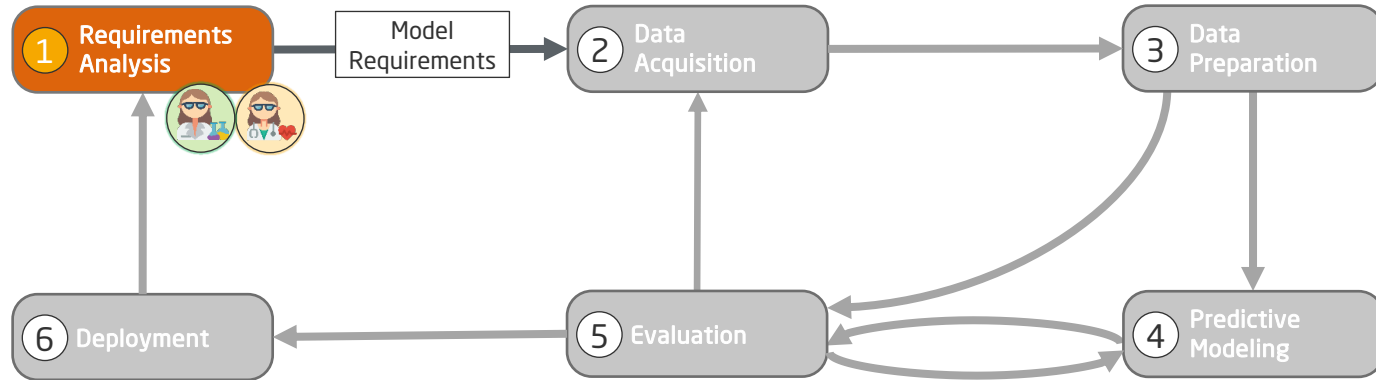
Roles



ML and Corpora

Data Management for Digital Health, Winter 2023
6

Requirements Analysis



Roles



Data Scientist



Domain Expert



(Data) Engineer

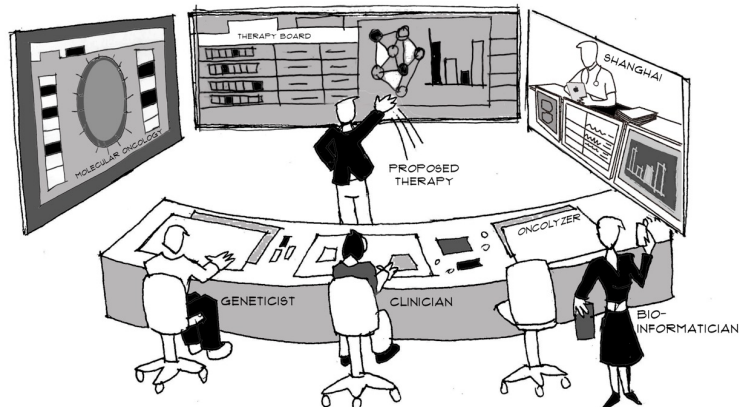
ML and Corpora

Data Management for
Digital Health, Winter
2023

NLP Use Case: Molecular Tumor Board

Patient data

- Genomic data
- Imaging data
- Clinical data
 - Patient history
 - Pathology reports
 - Radiology reports
 - ...



Medical knowledge

- Variant / gene databases
- Drug databases
- Cellular pathways
- Clinical Trials
- Clinical Guidelines
- Research publications

2008-15



Patient: Jimmy Smith
 Date of birth: 6/6/1972 Sex: Male
 Biopsy Date: 1/3/2008
 Doctor: Jennifer Tabernackie

Part A: LEFT MAXILLARY SOFT TISSUE

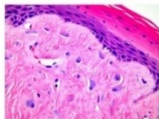
Gross Description:
 Submitted is formalin fixed tissue, measuring 1.6x1.4x1.0cm, stated to be from the left maxilla. The specimen consists of multiple pieces of brown soft tissue. Sections multiple. All submitted. Also submitted is a tooth, no sections taken.

Microscopic Description:

Multiple sections show keratinic, stratified squamous epithelium covering a core of dense and cellular fibrous connective tissue. Numerous enlarged stellate-shaped fibroblasts, some containing multiple nuclei, are seen in the lesional stroma.

Diagnosis: **Fibroma, giant cell type**

ICD: 210.4
 CPT: 88305



6.25.	Evidence-based Recommendation	2017
Grade of Recommendation B	After complete removal of a traditional serrated adenoma or sessile serrated adenoma, the follow-up should be the same as for classic adenomas.	
Level of Evidence 3b	Sources: [576, 617, 618]	
	Consensus	

ML and Corpora

Data Management for Digital Health, Winter 2023

8



Some Problems of Rule- and Dictionary-Based NER

- **Ambiguity:** e.g., gene names can be extremely weird (*a, white, swiss cheese, upregulated during skeletal muscle growth 5*)
- **Emerging terminology:** e.g., new drugs are approved all the time, have different trade names and only numeric codes when still experimental
- **Complex entities:** e.g., population and outcome statements come with a variety of modifiers and attributes that need to be detected

Disorders

Trial population and design

Genes / Proteins / Variants

Methods: In this open-label, phase 3 trial, we enrolled 665 patients with BRAF V600E-mutated metastatic colorectal cancer who had had disease progression after one or two previous regimens. Patients were randomly assigned in a 1:1:1 ratio to receive encorafenib, binimetinib, and cetuximab (triplet-therapy group); encorafenib and cetuximab (doublet-therapy group); or the investigators' choice of either cetuximab and irinotecan or cetuximab and FOLFIRI (folinic acid, fluorouracil, and irinotecan) (control group). The primary end points were overall survival and objective response rate in the triplet-therapy group as compared with the control group. A secondary end point was overall survival in the doublet-therapy group as compared with the control group. We report here the results of a prespecified interim analysis.

Drugs

Outcomes

Results: The median overall survival was 9.0 months in the triplet-therapy group and 5.4 months in the control group (hazard ratio for death, 0.52; 95% confidence interval [CI], 0.39 to 0.70; $P < 0.001$). The confirmed response rate was 26% (95% CI, 18 to 35) in the triplet-therapy group and 2% (95% CI, 0 to 7) in the control group ($P < 0.001$). The median overall survival in the doublet-therapy group was 8.4 months (hazard ratio for death vs. control, 0.60; 95% CI, 0.45 to 0.79; $P < 0.001$). Adverse events of grade 3 or higher occurred in 58% of patients in the triplet-therapy group, in 50% in the doublet-therapy group, and in 61% in the control group.

Rule- or ML-based Information Extraction?

	Pros	Cons
Rule-based	<ul style="list-style-type: none">• Declarative• Easy to comprehend• Easy to maintain• Easy to incorporate domain knowledge• Easy to trace and fix the cause of errors	<ul style="list-style-type: none">• Heuristic• Requires tedious manual labor
ML-based	<ul style="list-style-type: none">• Trainable• Adaptable• Reduces manual effort	<ul style="list-style-type: none">• Requires labeled data• Requires retraining for domain adaptation• Requires ML expertise to use or maintain• Opaque

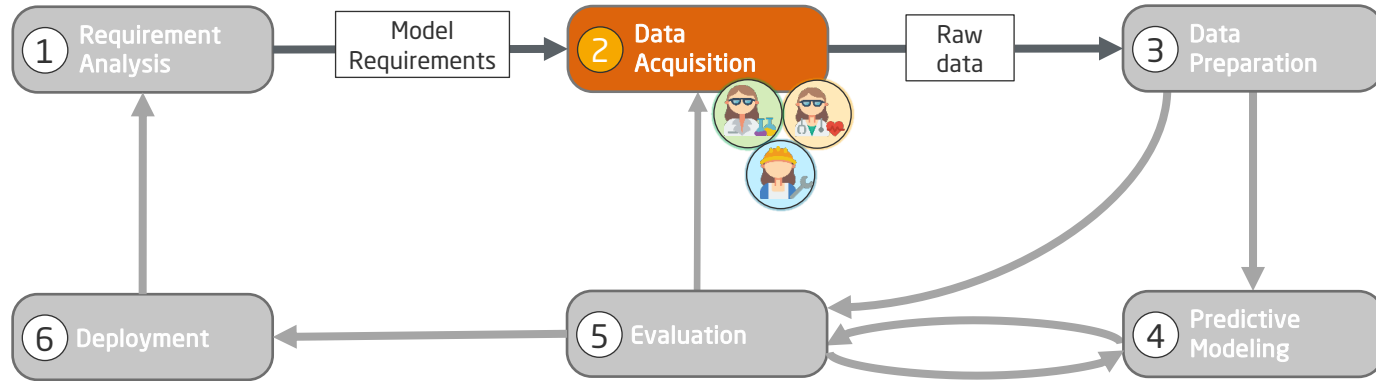
CAVE:

- from 2013 (ML-based back then != today)
- by IBM researchers (IBM sells / sold rule-based IE software)

ML and Corpora

Data Management for
Digital Health, Winter
2023
10

Data Acquisition



Roles



Data Scientist



Domain Expert



(Data) Engineer

ML and Corpora

Data Management for
Digital Health, Winter
2023

11

- Real-world **gold standard** text corpora necessary for development and evaluation of NLP methods
 - Available for general domains (Wikipedia, Twitter)
 - Many (biomedical) datasets based on MEDLINE abstracts
 - Few datasets (MIMIC, i2b2 / n2c2) of clinical text available for English language
- Issues:
 - Corpora for **low-resource languages** hardly available
 - **Data protection** for clinical text (every researcher is creating their own in-house corpus)
 - **Licensing** issues with full-text scientific articles
 - Tedious **annotation** process



HARVARD
MEDICAL SCHOOL

BLAVATNIK INSTITUTE
BIOMEDICAL INFORMATICS

n2c2 NLP Research Data Sets

Unstructured notes from the Research Patient Data Repository at Partners Healthcare.

[Need help? Contact us!](#)

Description

The majority of these Clinical Natural Language Processing (NLP) data sets were originally created at a former NIH-funded National Center for Biomedical Computing (NCBC) known as i2b2: Informatics for Integrating Biology and the Bedside.

- 2006 - Deidentification & Smoking
- 2008 - Obesity
- 2009 - Medication
- 2010 - Relations
- 2011 - Coreference
- 2012 - Temporal Relations
- 2014 - Deidentification & Heart Disease

<https://portal.dbmi.hms.harvard.edu/projects/n2c2-nlp/>

ML and Corpora
Data Management for
Digital Health, Winter
2023
12

Data Availability Problem for Non-English Language Communities



Authentic Clinical Texts

English

- Discharge Summaries
- Radiology Reports
- Pathology Reports



De-identified / Synthetic Clinical Texts



Non-individual Medical Texts



German

- Discharge Summaries
- Radiology Reports
- Pathology Reports

- Very few / small (e.g., GraSCCo, BRONCO)

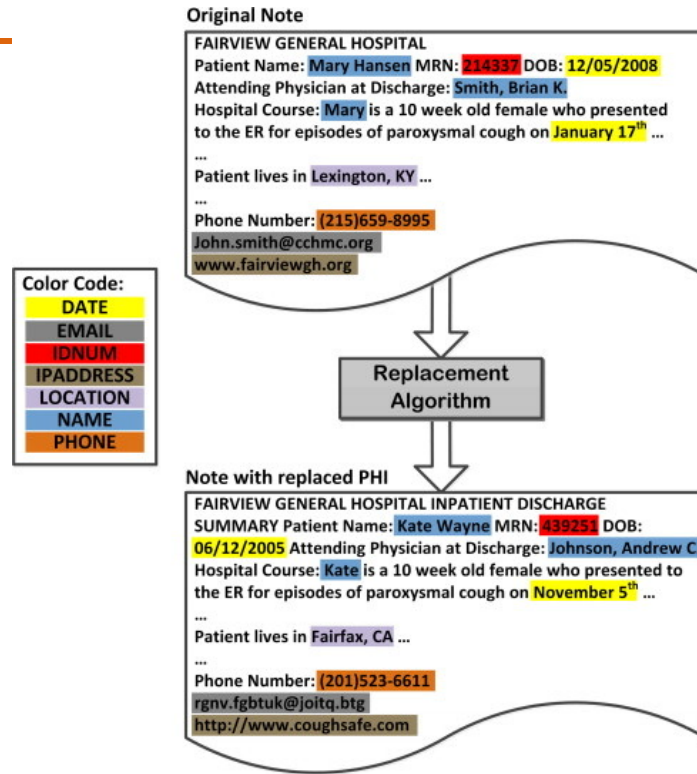
- Guidelines
- Textbooks

ML and Corpora

Data Management for
Digital Health, Winter
2023
13

De-identification of Clinical Text

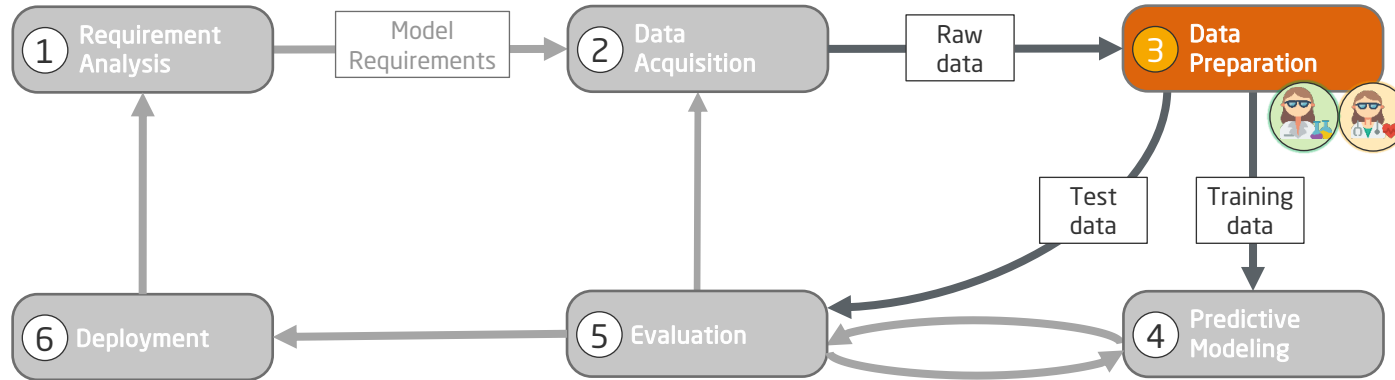
- De-identification of clinical text is tricky
 - US: 18 PHI (protected health information) categories shall be removed
 - EU: not well defined (GDPR)
- Automatic de-identification is challenging NLP task itself
 - e.g. [i2b2 Challenge 2014](#)
 - Similar problem to NER
 - Possible with relatively high accuracy (F1 scores > .90)
→ recall below 100% acceptable?



ML and Corpora

Data Management for
Digital Health, Winter
2023
14

Data Preparation



Roles



Data Scientist



Domain Expert



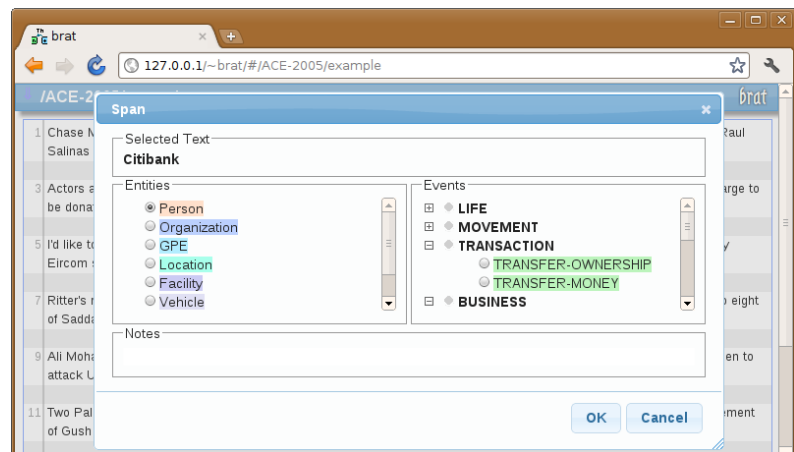
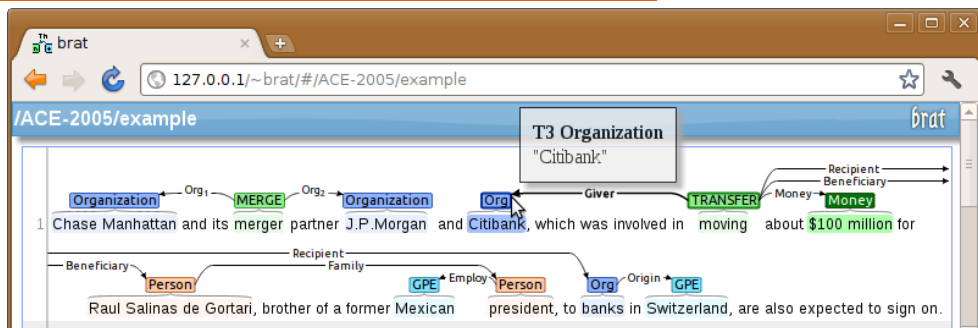
(Data) Engineer

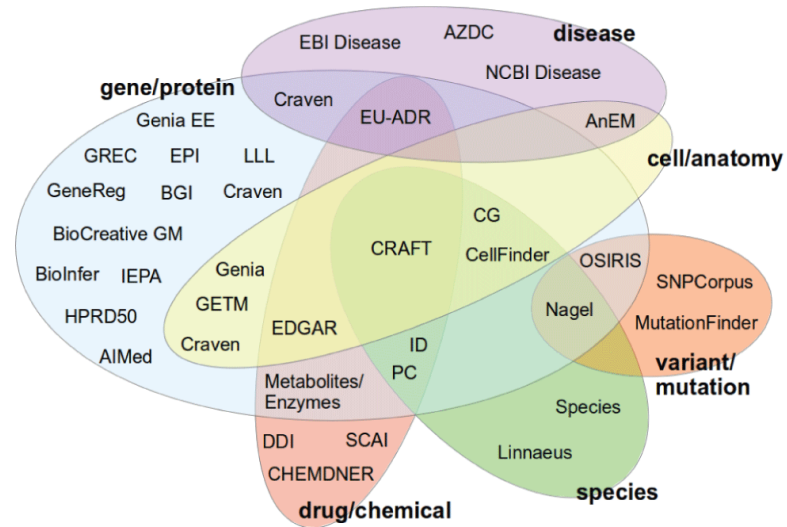
ML and Corpora

Data Management for
Digital Health, Winter
2023
15

Annotation of Text Corpora

- Requires domain expertise
- Usually **multiple annotators**
- Need clearly defined **annotation guidelines**
- Crowdsourcing (Mechanical Turk, Upwork) is an option
 - Many noisy labels can be better for ML than few good labels
 - Gold standard evaluation data still necessary
- Popular visual online annotation tools: BRAT, INCEpTION





ML and Corpora

Data Management for
Digital Health, Winter
2023
17

NCBI Disease Corpus (Dogan et al. 2014)

- 793 MEDLINE abstracts
- Annotations:
 - Disease Mentions (6892)
 - MeSH / OMIM Concepts (790)
 - 14 annotators
 - Agreement F1-Score around 0.9
- Entity-level F1-Score:
 - upon publication: 63.7%
 - State-of-the-art: 89.1%
(BioBERT)

PMID:10633128 **Friedreich ataxia: an overview.**
Publication: Journal of medical genetics; 2000 Jan ; 37(1) 1-8
[CompositeMention](#) [Modifier](#) [SpecificDisease](#) [DiseaseClass](#) [Clear](#) [Reset](#)

TITLE:
Friedreich ataxia: an overview.

ABSTRACT:
Friedreich ataxia, an **autosomal recessive neurodegenerative disease**, is the most common of the **inherited ataxias**. The recent discovery of the gene that is mutated in this condition, FRDA, has led to rapid advances in the understanding of the pathogenesis of **Friedreich ataxia**. About 98% of mutant alleles have an expansion of a GAA trinucleotide repeat in intron 1 of the gene. This leads to reduced levels of the protein, frataxin. There is mounting evidence to suggest that **Friedreich ataxia** is the result of accumulation of iron in mitochondria leading to excess production of free radicals, which then results in cellular damage and death. Currently there is no known treatment that alters the natural course of the disease. The discovery of the **FRDA** gene and its possible function has raised hope that rational therapeutic strategies will be developed..

[Highlight]: [FRDA](#) | [inherited ataxias](#) | [autosomal recessive neurodegenerative disease](#) | [Friedreich ataxia](#)

Type	Mention	Annotator 1	Annotator 2	Nomenclature	Delete
SpecificDisease	Friedreich ataxia	D005621	-TheSame-	CTD Disease	Delete
SpecificDisease	Friedreich ataxia	D005621	-TheSame-	CTD Disease	Delete
DiseaseClass	autosomal recessive neurodegenerative disease	D020271	-TheSame-	CTD Disease	Delete
SpecificDisease	inherited ataxias	D013132	D020754	CTD Disease	Delete
SpecificDisease	Friedreich ataxia	D005621	-TheSame-	CTD Disease	Delete
SpecificDisease	Friedreich ataxia	D005621	-TheSame-	CTD Disease	Delete
Modifier	FRDA	D005621	-TheSame-	CTD Disease	Delete

[Save Annotation Results](#) [Save & Export Annotation Results](#)

German Clinical Text BRONCO (Kittner et al. 2021)

- First publicly available, de-identified German clinical text corpus (200 oncological discharge summaries)
- Interesting de-identification procedure:
 - “First, the corpus was completely **manually deidentified**. This process was **confirmed by Charité and UKT data protection officers**.”
 - “Second, we only annotate and publish certain sections of the discharge summaries, **avoiding all sections containing mostly biographic information**”
 - “Third, we **shuffled all sentences in the 2 subcorpora** to blur their order and relationships.”

1	2010 Erstdiagnose Aderhautmelanom rechts	DIAGNOSIS [R][C69.3]
2	cMRT: keine zerebralen Metastasen	DIAGNOSIS [negative][C79.3]
3	Die naechste Ausbreitungsdiagnostik (CT und MR-Oberbauch) wurde fuer den 1.03.2023 terminiert.	TREATMENT [possibleFuture][3-804]
4	Z.n. radikaler Lymphadenektomie rechte Axilla, Level I-III 08/2005	TREATMENT [R][5-404.03#]1 TREATMENT [R][5-404.03#]
5	Manifestationen: pulmonal, fragl. ossaer	DIAGNOSIS [C78.0] DIAGNOSIS [speculative][C79.3] DIAG [speculative][C79.3]
6	Am 7.04.2134 erfolge die komplikationslose Nivolumab-Infusion.	TREATMENT [6-008.m] MEDICATION [L01XC17]
7	Beginn Chemotherapie nach dem GeT-Schema Zyklus 1.	TREATMENT [8-542] TREATMENT [6-001.1] MEDICATION [L01BC05] MEDICATION [L01AB02]
8	Im CT hatte sich eine hochgradig HCC suspekta Laesion im Lebersegment VI gezeigt.	TREAT [3-207]

ML and Corpora

Data Management for
Digital Health, Winter
2023
19

- 1.87 M tokens from 30 German clinical guidelines in oncology
- Currently the largest, **shareable** corpus of German medical text
(<https://www.leitlinienprogramm-onkologie.de/projekte/ggponc-english/>)
- Available through Data Use Agreement (i.e., non-commercial use only)

8.1.4. UICC-Stadium II

8.5.	Evidenzbasierte Empfehlung	2017
Empfehlungsgrad 0	Bei Patienten mit einem kurativ resezierten Kolonkarzinom im Stadium II kann eine adjuvante Chemotherapie durchgeführt werden.	
Level of Evidence 1b	Quellen: [880-884]	
	Starker Konsens	

Hintergrund

Der Nutzen einer **adjuvanten Therapie** im **UICC Stadium II** ohne Risikofaktoren liegt absolut zwischen 2-5% im 5-Jahresüberleben. Studien und gepoolte Analysen von Studien bei Patienten mit einem **Kolonkarzinom** im **Stadium II** fand sich kein signifikanter **Überlebensvorteil** durch eine **postoperative adjuvante Chemotherapie** [880-883]. Die gepoolte Analyse von 7 randomisierten Studien, die eine adjuvante Chemotherapie mit einer alleinigen Operation verglichen, zeigte in der univariaten Analyse nur eine signifikante Verbesserung für das **5 Jahres krankheitsfreie Überleben (DFS)** (72 versus 76%; $p=0.049$), aber nicht für das **5 Jahres Gesamtüberleben** (80 versus 81%; $p=0.1127$) im **Stadium II**, wobei sich die Einzelstudien deutlich in den Therapiemodi unterschieden und kleine Patientenzahlen einschlossen [875]. In einem Kollektiv von 43.032 Medicare Empfängern (> 66 Jahre; **Stadium III** 18.185 Pat., 57 % erhielten eine **adjuv. Chemotherapie**) besaßen im **Stadium II** 6.234 Pat. kein (19 % mit **adjuv. Chemotherapie** und 18.613 Pat. mindestens ein prognostisch ungünstiges Kriterium (21 % mit adjuv.

GGPONC Annotation & Distribution

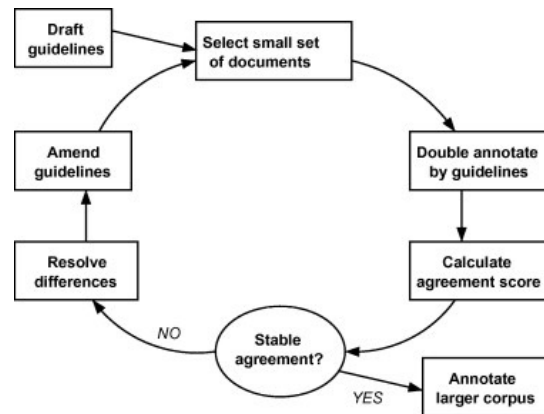
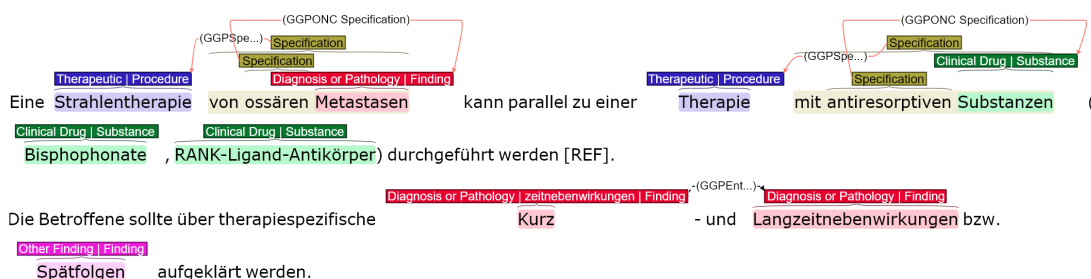
■ Annotation team:

- 7 annotators (medical students)
- 1 curator (medical doctor)

■ 6 months / 1200 hours

■ > 200k entities

■ Iterative guideline refinement until stable agreement was reached ($\gamma = .94$)



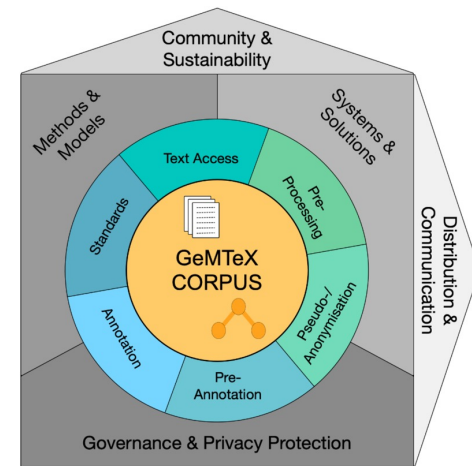
ML and Corpora

Data Management for
Digital Health, Winter
2023
21

GeMTeX

German Medical Text Corpus

- Goal: Create a large annotated text corpus of **German medical texts from routine patient care**
- Building upon German **Medical Informatics Initiative** (MII) infrastructure:
 - Broad consent
 - Data integration centers (DIC) at six university hospitals
- Automated + manual de-identification
- **Deep semantic annotation:**
 - Entities and grounding (SNOMED CT)
 - Relations (e.g., temporal, causal)
 - Domain-specific annotations (oncology, cardiology, neurology, pharmacology)



ML and Corpora

Data Management for
Digital Health, Winter
2023
22

Feature Engineering for ML-based Named-Entity Recognition

- Training data := Sequences of tokens
- Labels := Sequences of I/O/B tags
- Simple approach to ML-based NER

- Turn token into feature vector v

- **Classify** each token as I/O/B

- $f(v(\text{"BRAF"})) = B$

- $f(v(\text{"V600E"})) = I$

- $f(v(\text{"mutation"})) = O$

$\mathbf{x} = (x_1, \dots, x_n)$, e.g. $x_1 = \text{"BRAF"}$, $x_6 = \text{"mutation"}$

$\mathbf{y} = (y_1, \dots, y_n)$, e.g. $y_1 = \text{"B"}$, $y_6 = \text{"O"}$

BRAF V600E is a driver mutation found in multiple tumor types

B I 0 0 0 0 0 0 0 0 0

- What do you think could be useful features for NER on the token level?

ML and Corpora

Data Management for
Digital Health, Winter
2023
23

- Simple feature: **Identity** of current token
- **One-hot-encoding** : create a feature vector with length equal to size of vocabulary and set 1 if word == ith element else 0
- Vocabulary size is typically restricted to most common k tokens
- Special feature for **unknown tokens** (unknown tokens are often named entities!)
- Resulting feature vectors are very **sparse** (almost all elements are 0)

$$v(\text{"apple"}) \rightarrow \begin{pmatrix} 1 \\ 0 \\ \dots \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$v(\text{"cancer"}) \rightarrow \begin{pmatrix} 0 \\ 0 \\ \dots \\ 0 \\ 1 \\ 0 \\ \dots \\ 0 \end{pmatrix}$$

Word Shape Features

- **Word shape** features := abstract letter patterns
- Prefixes
- Suffixes (e.g., word ends with “-itis” → disease)
- character n-grams

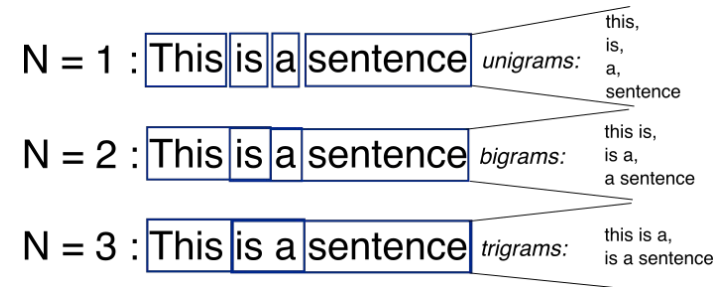
<i>Feature</i>	<i>About feature</i>	<i>Regular Expression</i>
Initcaps	First letter is in uppercase	[A-Z]\w+
Initcapsalpha	First letter is in uppercase. Second letter is in lowercase.	[A-Z][a-z]\w+
All caps	All letters are in uppercase	[A-Z]+
Caps mix	Mixture of uppercase and lowercase letters	[A-Za-z]+
Has digit	Protein name has a number in the middle	\w+[0-9]\w+
Single digit	Ranges from 0-9	[0-9]
Double digit	Two digit numbers	[0-9][0-9]
Natural numbers	Any natural number	[0-9]+
Real numbers	Decimal numbers /numbers with comma	[-0-9]+[.][0-9.]+
Has dash	Protein name with dash in middle	\w+ - \w+
Init dash	First character is a dash	- \w+
End dash	Last character is a dash	\w+ -
Alphanumeric (starts with alphabet)	Combination of alphabets and numbers. First character is an alphabet.	\w+ [A-Za-z] \w+ [0-9] \w+
Alphanumeric (starts with number)	Combination of alphabets and numbers. First character is a number.	\w+ [0-9] \w+ [A-Za-z] \w+
Roman letter	Any roman letter	[IVXDLCM]+
Has Roman	Any roman letter in the middle	\w+ [IVXDLCM]+ \w+
Greek	Any Greek letter	\w+ [αβγÖE]
Has Greek	Any Greek letter in middle	\w+ [αβγÖE] \w+
Punctuation	Protein name with punctuation	

ML and Corpora

Data Management for
Digital Health, Winter
2023
25

More Token Features

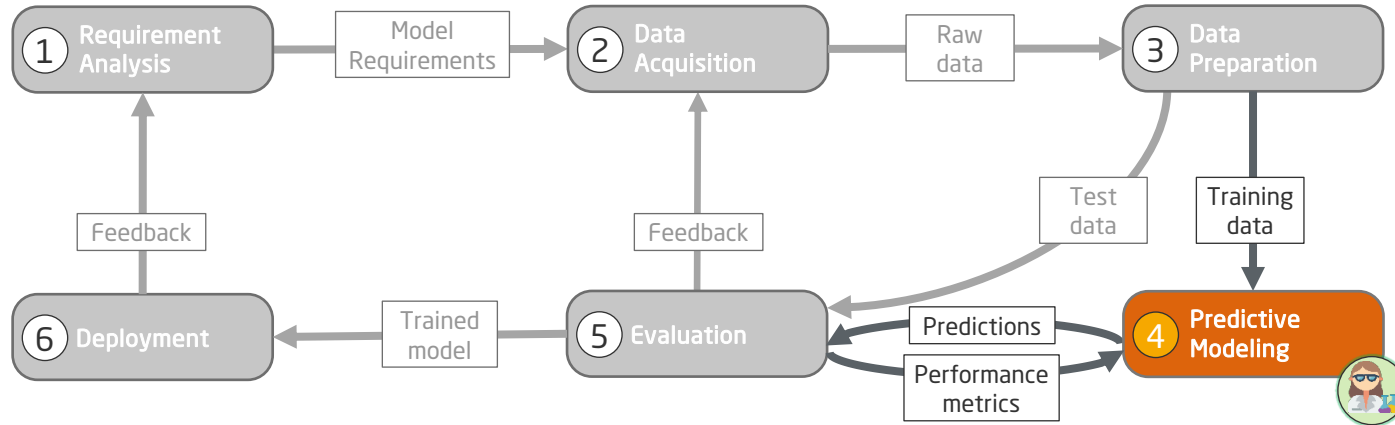
- Dictionary match
- Part-of-speech tags (e.g., is the word a proper noun)
- Context features
 - **Distributional semantics**: “a word is characterized by the company it keeps” (J.R. Firth, 1957)
→ Surrounding tokens can be more important than token itself
 - word **n grams** (common: trigrams with current, left and right word)



ML and Corpora

Data Management for
Digital Health, Winter
2023
26

Predictive Modeling



Roles



Data Scientist



Domain Expert

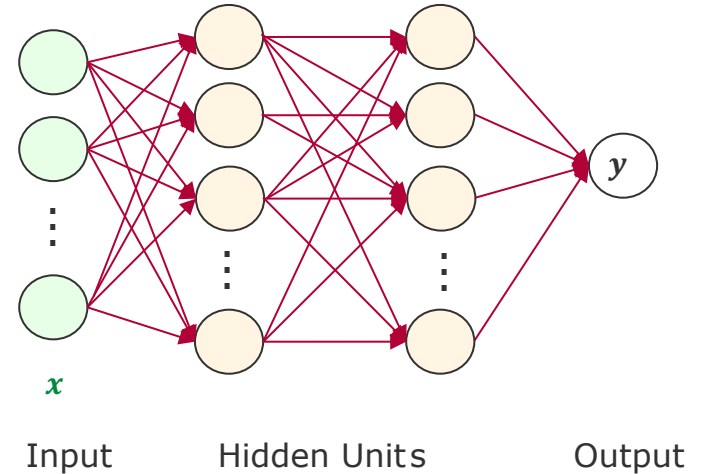


(Data) Engineer

ML and Corpora

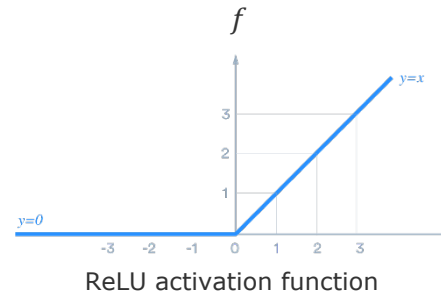
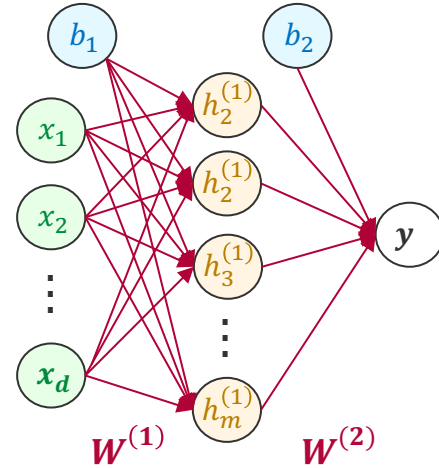
Data Management for
Digital Health, Winter
2023
27

- Class of machine learning models that are loosely inspired by neurons in the brain
- Connections between units have learnable weights
- Hidden units can learn **higher-level representations** of the input



Feed Forward Neural Networks

- Simplest version of a NN (also called Multilayer Perceptron)
- Representations are learned in *hidden layers* (h)
 - $\mathbf{h}^{(1)} := f(\mathbf{W}^{(1)T} \mathbf{x} + \mathbf{b}_1)$
 - $\mathbf{y} := \sigma(\mathbf{W}^{(2)T} \mathbf{h}^{(1)} + \mathbf{b}_2)$ (binary classification)
 - $\mathbf{y} := \text{softmax}(\mathbf{W}^{(2)T} \mathbf{h}^{(1)} + \mathbf{b}_2)$ (multiclass classification)
- Note: Logistic Regression is equivalent to neural network without hidden layers!
- 2-Layer MLP is more expressive than Logistic Regression because of **non-linear** activation function f
- Common choice for f : Sigmoid or ReLU
- **Deep Learning** = Neural Networks with many layers

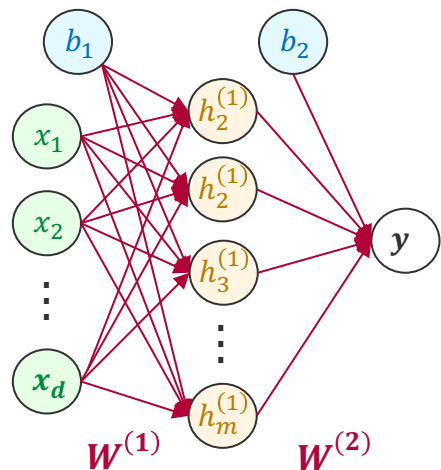


ML and Corpora

Data Management for
Digital Health, Winter
2023
29

Backpropagation

Input →



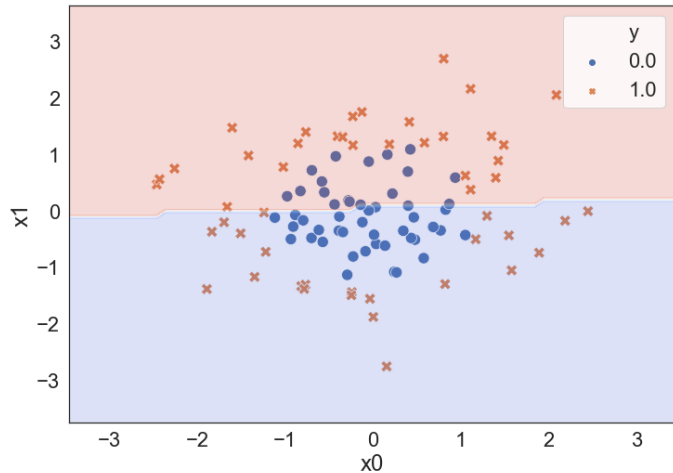
Training := Iteratively adapting weights such that error is decreased

↔ **Error** y_{true}

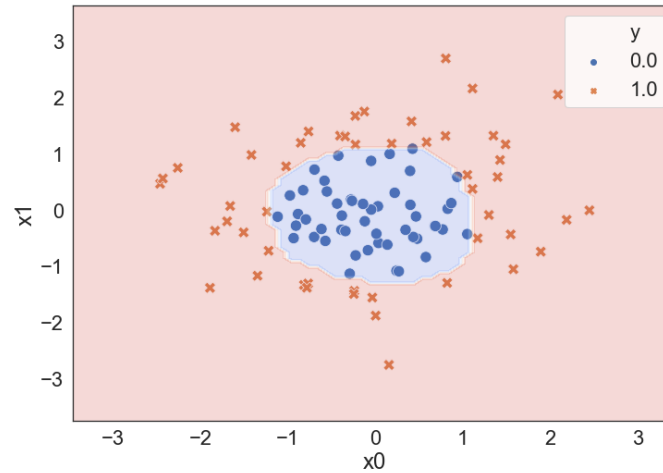
$$\Delta w_{ij} = -\eta \frac{\partial E}{\partial w_{ij}}$$

← **Gradient**

Neural Networks Can Learn Non-Linear Functions

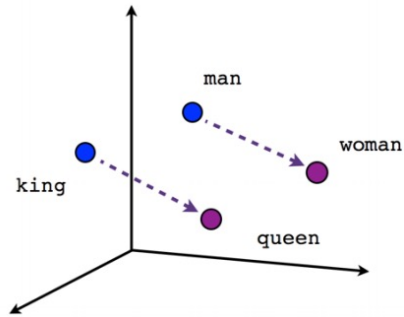


Decision boundary learned by
Logistic Regression

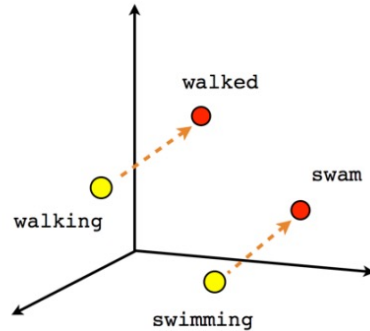


Decision boundary learned by
Neural Network with 1 hidden layer

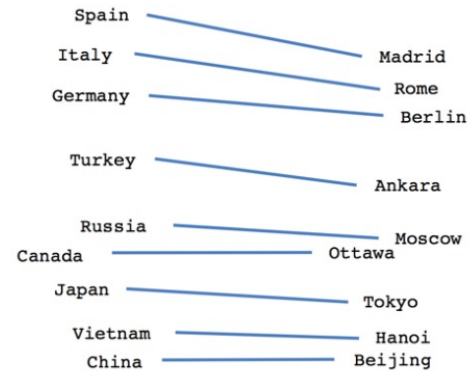
(Shallow) Semantics in Word Embedding Space



Male-Female



Verb tense



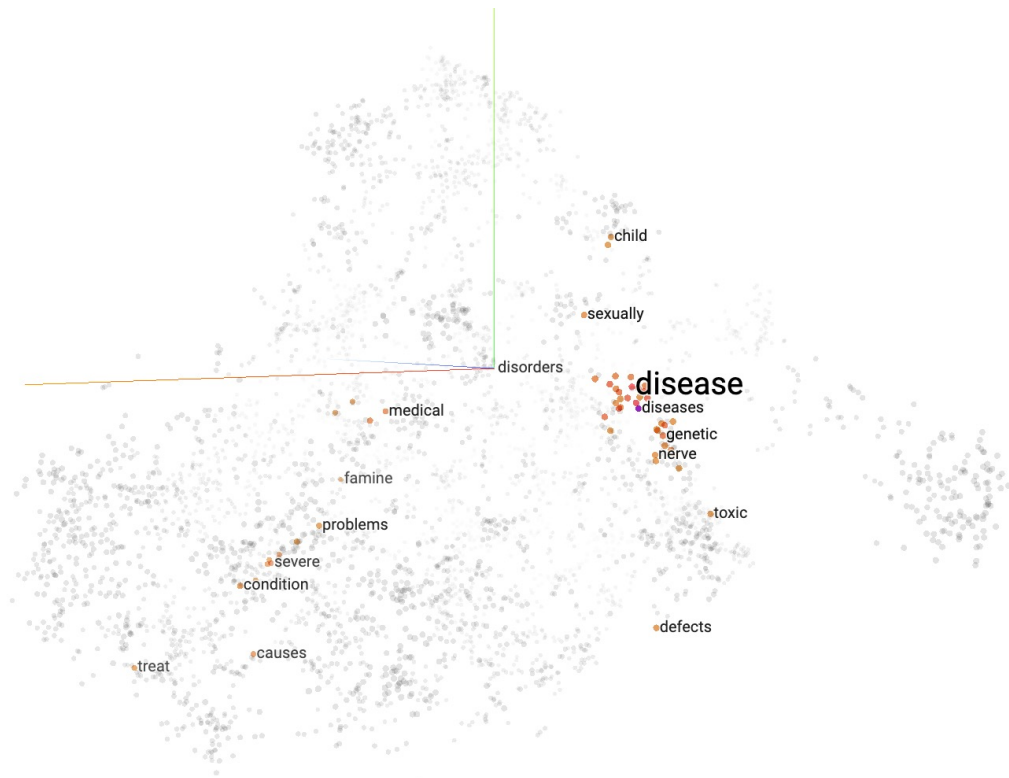
Country-Capital

ML and Corpora

Data Management for
Digital Health, Winter
2023
33

Word Embedding Demo

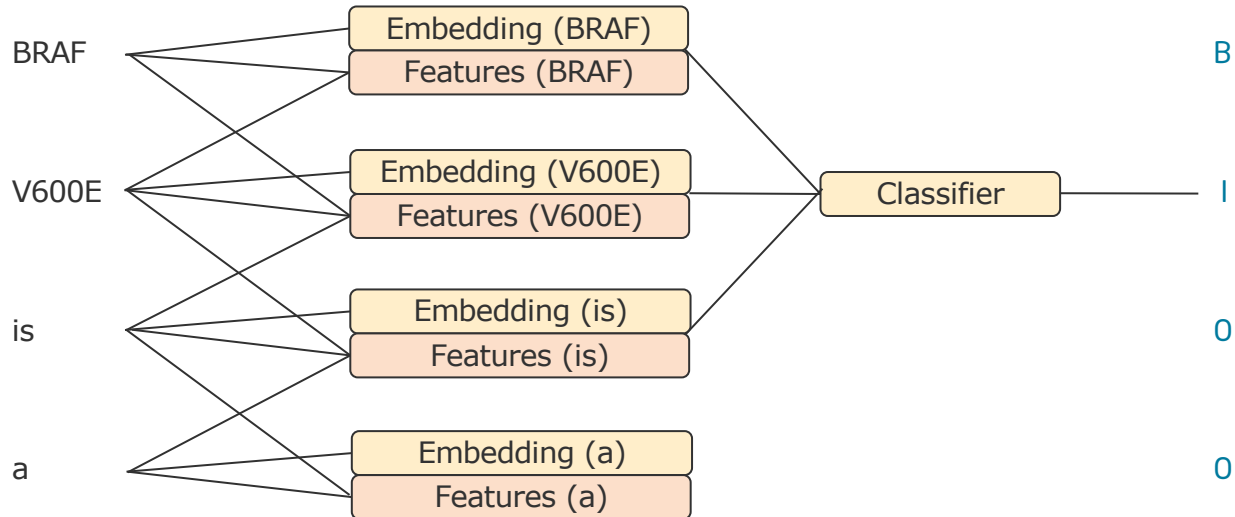
<https://projector.tensorflow.org/>



ML and Corpora

Data Management for
Digital Health, Winter
2023

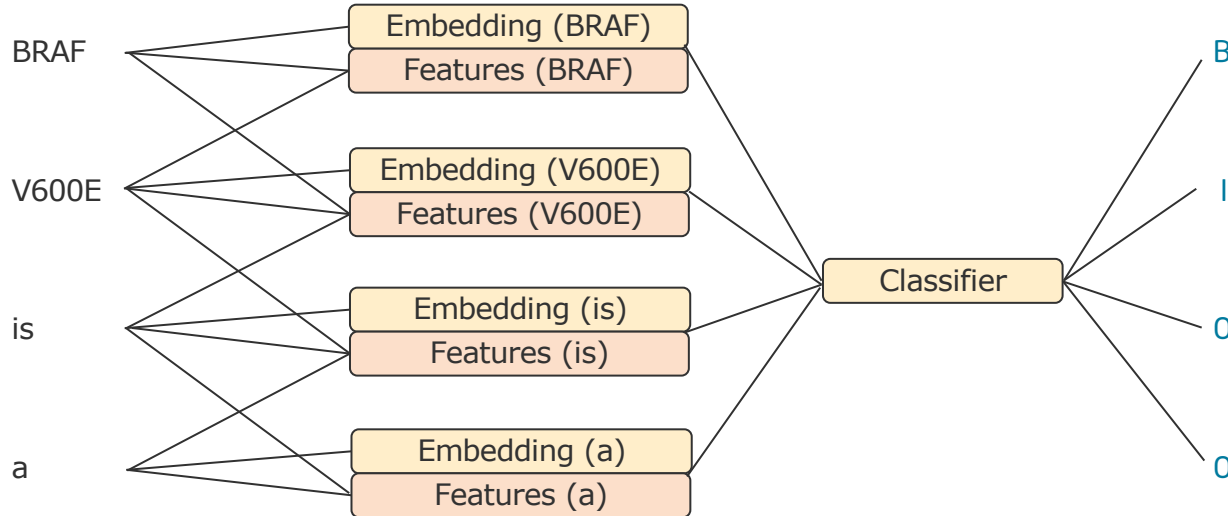
Word Embedding as Features



ML and Corpora

Data Management for
Digital Health, Winter
2020
35

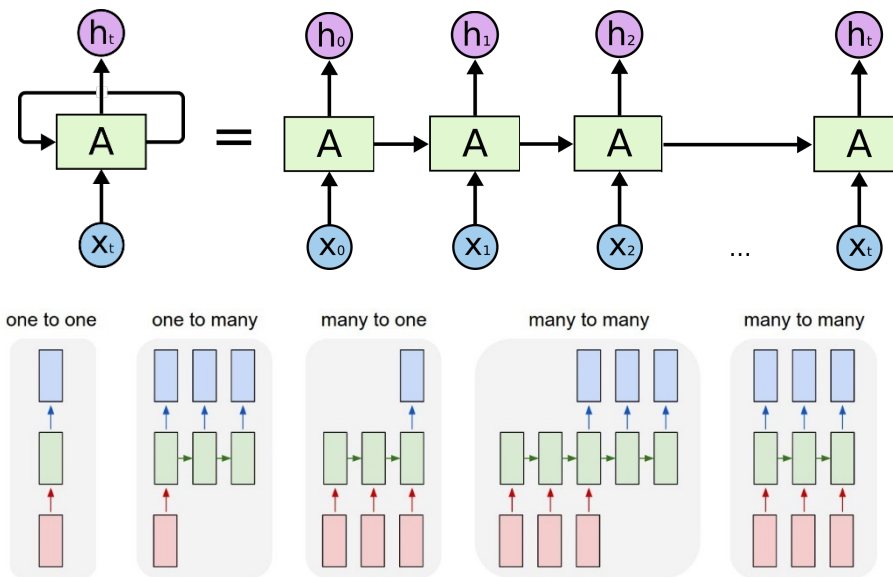
- Classifying single tokens assumes **statistical independence of adjacent labels** → usually wrong
- Idea: predict whole sequence of labels y from sequence of inputs x (**Structured Prediction**)



- ... in practice we need to impose some restrictions to keep the problem tractable

Recurrent Neural Networks

- Special type of neural network with recurrent connections → variable number of computational steps
- Various types of RNN cells exist (denoted by **A** in the picture)
- Enable **end-to-end learning** := directly map inputs (sentences) to outputs without explicitly solving intermediate steps (feature engineering, linguistic analysis)
- Computationally much more demanding than linear models

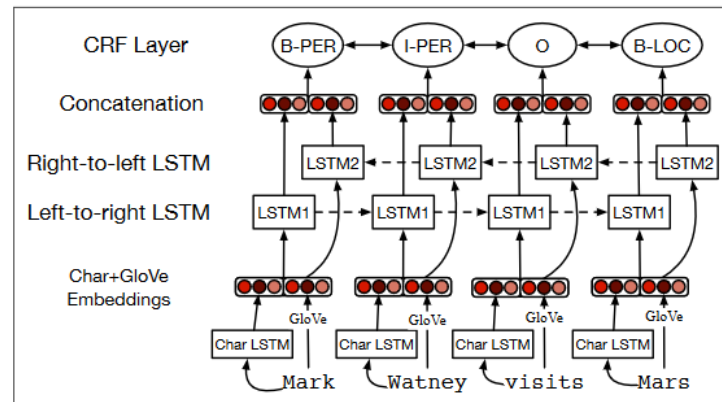


ML and Corpora

Data Management for
Digital Health, Winter
2023
37

Deep Learning for Sequence-to-Sequence Problems

- State-of-the-art approach for NER used to be **bi-LSTM-CRF**:
 - **Word embeddings** pretrained on large, unlabelled corpus (such as PubMed abstracts)
 - **Bidirectional Long-Short-Term-Memory (LSTM)** neural network (special form of recurrent neural network) for feature learning
 - **CRF output layer** to predict sequence of labels from LSTM representations
- Recent superseded by **Transformer**-based architectures, yet, still widely in use
- Sequence-to-sequence architectures also applicable to other problems, such as translation or summarization

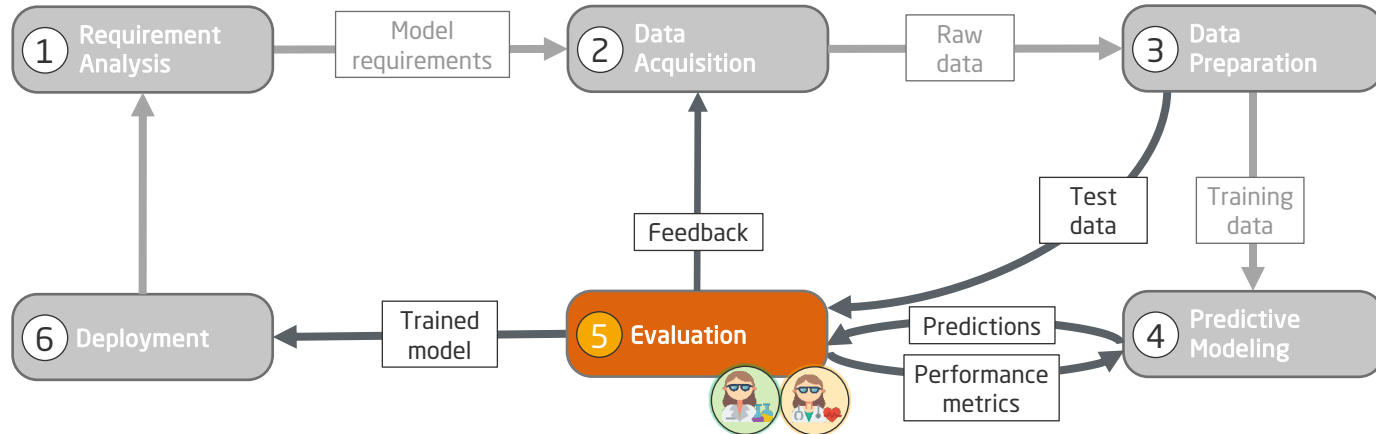


Martin, James H., and Daniel Jurafsky. *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition*. 3rd ed. Draft <https://web.stanford.edu/~jurafsky/slp3/>

ML and Corpora

Data Management for
Digital Health, Winter
2023
38

Evaluation



Roles



Data Scientist



Domain Expert



(Data) Engineer

ML and Corpora

Data Management for
Digital Health, Winter
2023
39

- To evaluate performance of NER and other NLP systems, performance is compared against **gold standard**
- Typical evaluation measures precision, recall, F1 score
→ need to be adapted to account for partial matches

The patient underwent a CT scan in April which did not reveal lesions in his liver.

B I B B ← Ground Truth

B B B ← NER Result

- Strict evaluation (entity level): $TP = 1, FP = 1, FN = 3, TN = 12$
- Loose evaluation (token level): $TP = 2, FP = 1, FN = 2, TN = 12$
- Many possible ways to weight partial matches

Error Types in Strict Evaluation

Token	Paracetamol	for	migraine
True label	Medication	0	Diagnosis
Predicted label	Medication	Diagnosis	
Counted errors	$2 = 1FP + 1FN$		

Boundary Error

Token	Paracetamol	for	migraine
True label	Medication	0	Diagnosis
Predicted label	Medication	0	Medication
Counted errors	$2 = 1FP + 1FN$		

Labeling Error

Token	Paracetamol	for	migraine
True label	Medication	0	Diagnosis
Predicted label	Medication	Medication	Diagnosis
Counted errors	1FP		

False Positive

Token	Paracetamol	for	migraine
True label	Medication	0	Diagnosis
Predicted label	Medication	0	0
Counted errors	1FN		

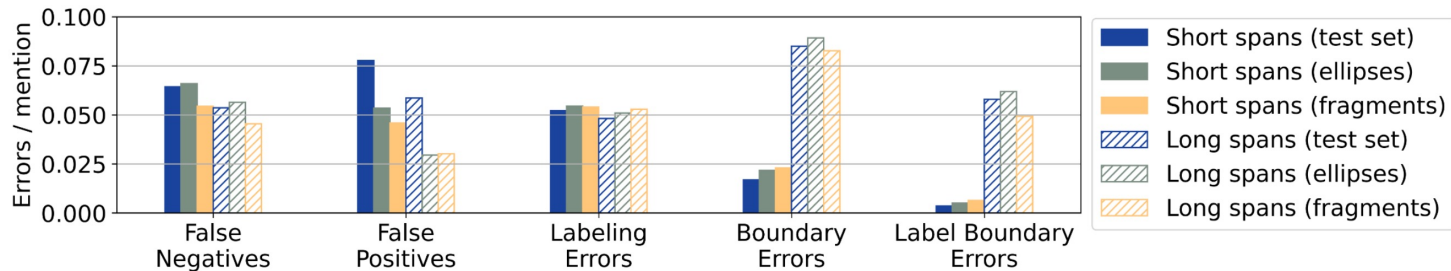
False Negative

- Traditional metrics punish labelling errors, boundary errors, and label boundary errors multiple times
- Alternative scores account for partial overlap, e.g., $1LE = 1BE = 1LBE = 0.5FP + 0.5FN$ (cf. Ortmann, 2022 and others)

GGPONC: Evaluation and Error Analysis

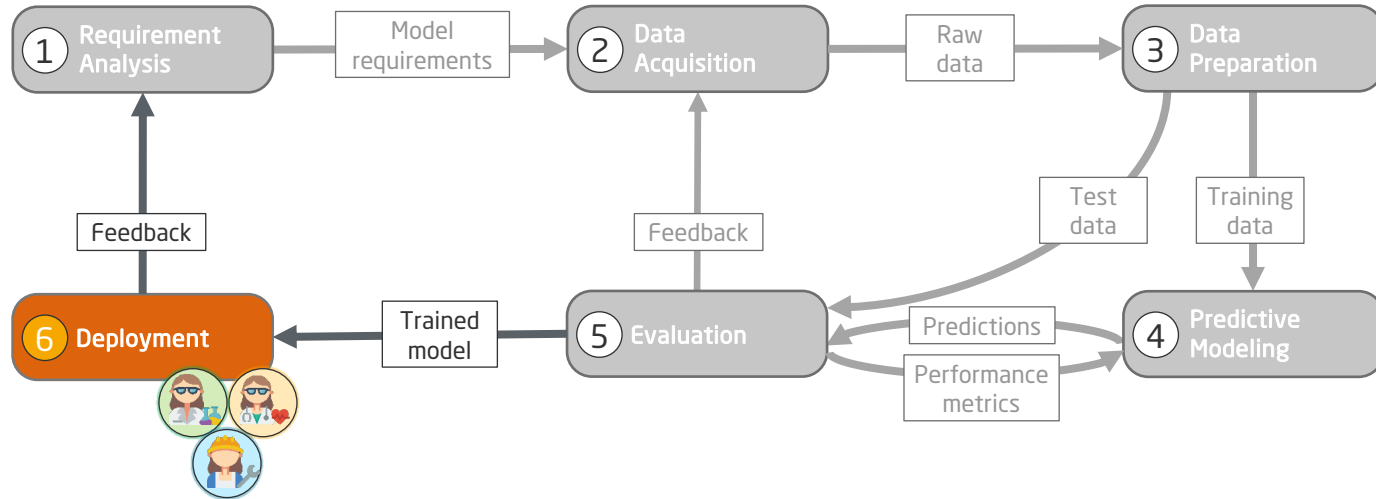
- Good performance on short / coarse annotations
- Low F1 score for long / fine annotations
- Why? → Error Analysis

F1 Score	Coarse	Fine
Short	.89	.86
Long	.75	.72



ML and Corpora

Deployment



Roles



Data Scientist



Domain Expert

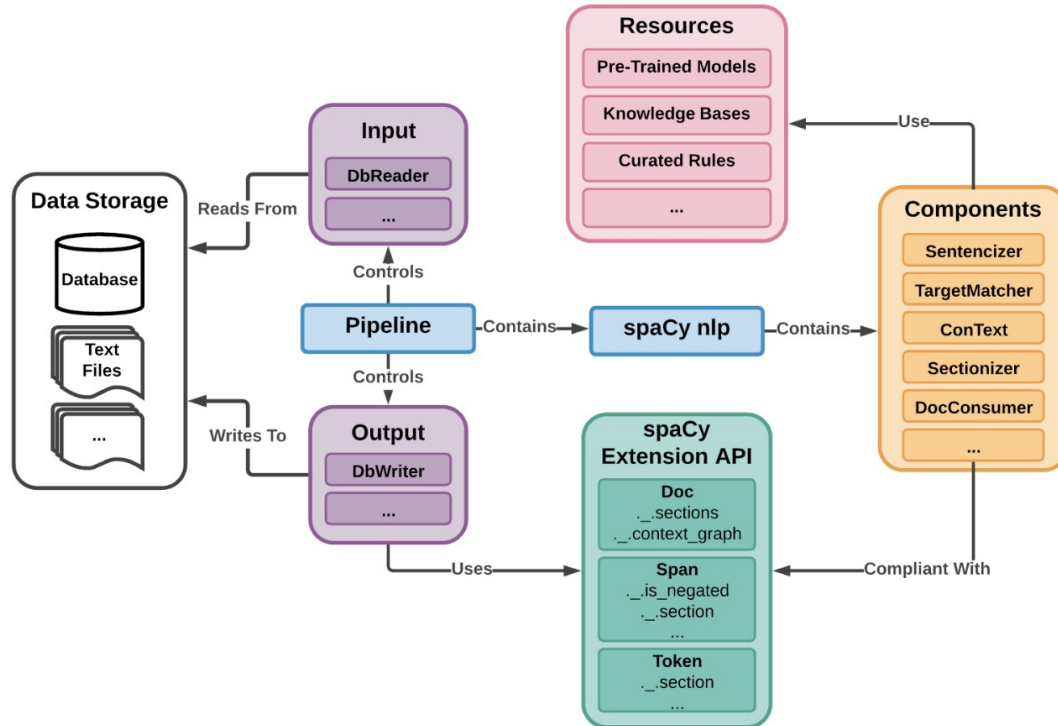


(Data) Engineer

ML and Corpora

Data Management for
Digital Health, Winter
2023

43

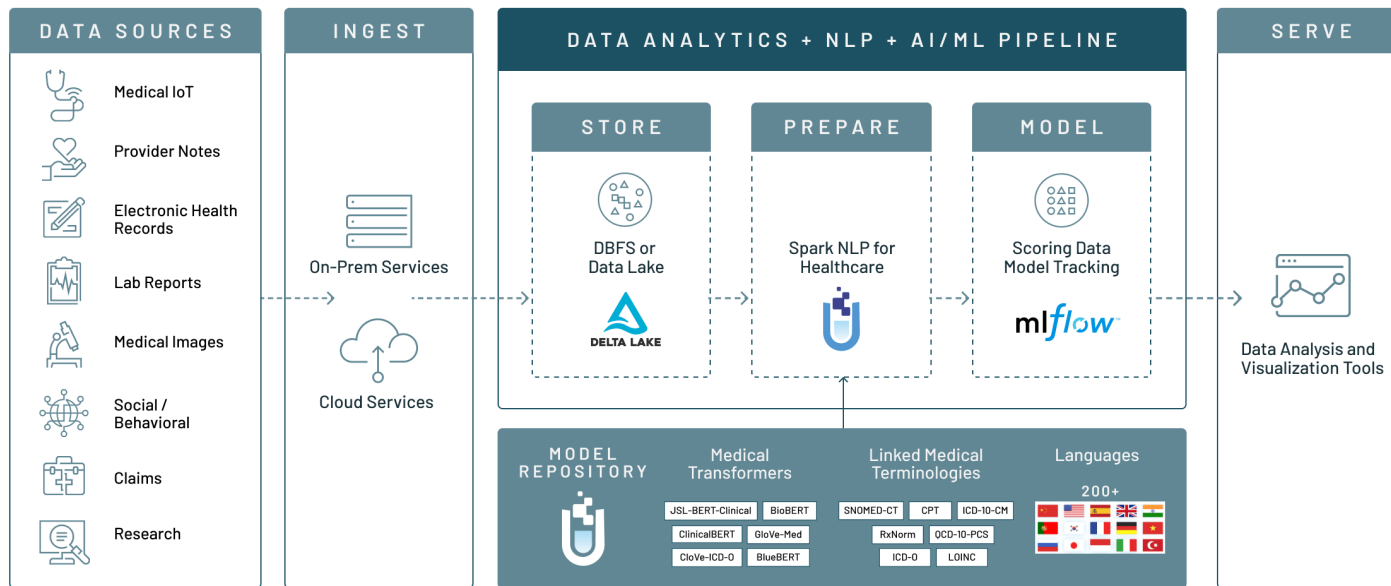


ML and Corpora

Data Management for
Digital Health, Winter
2023

Commercial Solutions for Large-scale Clinical NLP

databricks LAKEHOUSE PLATFORM



ML and Corpora

Data Management for
Digital Health, Winter
2023
45

What to Take Home



- Clinical text corpora (de-identification, annotation, availability)
- Evaluation of NER results
- Features for ML-based NLP algorithms
- Neural networks and neural architectures for Named Entity Recognition

ML and Corpora

Data Management for
Digital Health, Winter
2023
46