

PDF hosted at the Radboud Repository of the Radboud University Nijmegen

The following full text is a publisher's version.

For additional information about this publication click this link.

<https://hdl.handle.net/2066/216864>

Please be advised that this information was generated on 2021-06-25 and may be subject to change.

Comparing Rule-based, Feature-based and Deep Neural Methods for De-identification of Dutch Medical Records

Jan Trienes
Nedap Healthcare
Groenlo, Netherlands
jan.trienes@nedap.com

Dolf Trieschnigg
Nedap Healthcare
Groenlo, Netherlands
dolf.trieschnigg@nedap.com

Christin Seifert
University of Twente
Enschede, Netherlands
c.seifert@utwente.nl

Djoerd Hiemstra
Radboud University
Nijmegen, Netherlands
djoerd.hiemstra@ru.nl

ABSTRACT

Unstructured information in electronic health records provide an invaluable resource for medical research. To protect the confidentiality of patients and to conform to privacy regulations, de-identification methods automatically remove personally identifying information from these medical records. However, due to the unavailability of labeled data, most existing research is constrained to English medical text and little is known about the generalizability of de-identification methods across languages and domains. In this study, we construct a varied dataset consisting of the medical records of 1260 patients by sampling data from 9 institutes and three domains of Dutch healthcare. We test the generalizability of three de-identification methods across languages and domains. Our experiments show that an existing rule-based method specifically developed for the Dutch language fails to generalize to this new data. Furthermore, a state-of-the-art neural architecture performs strongly across languages and domains, even with limited training data. Compared to feature-based and rule-based methods the neural method requires significantly less configuration effort and domain-knowledge. We make all code and pre-trained de-identification models available to the research community, allowing practitioners to apply them to their datasets and to enable future benchmarks.

KEYWORDS

natural language processing, machine learning, privacy protection, medical records

1 INTRODUCTION

With the strong adoption of electronic health records (EHRs), large quantities of unstructured medical patient data become available. This data offers significant opportunities to advance medical research and to improve healthcare related services. However, it has to be ensured that the privacy of a patient is protected when performing secondary analysis of medical data. This is not only an ethical prerequisite, but also a legal requirement imposed by privacy legislations such as the US Health Insurance Portability and Accountability Act (HIPAA) [13] and the European General Data Protection Regulation (GDPR) [9]. To facilitate privacy protection, de-identification has been proposed as a process that removes or masks any kind of protected health information (PHI) of a patient such that it becomes difficult to establish a link between an individual and the data [20]. What type of information constitutes PHI is in part defined by privacy laws of the corresponding country.

For instance, the HIPAA regulation defines 18 categories of PHI including names, geographic locations, and phone numbers [14]. According to the HIPAA safe-harbor rule, data is no longer personally identifying and subject to the privacy regulation if these 18 PHI categories have been removed. As the GDPR does not provide such clear PHI definitions, we employ the HIPAA definitions throughout this paper.

As most EHRs consist of unstructured, free-form text, manual de-identification is a time-consuming and error-prone process which does not scale to the amounts of data needed for many data mining and machine learning scenarios [7, 21]. Therefore, automatic de-identification methods are desirable. Previous research proposed a wide range of methods that make use of natural language processing techniques including rule-based matching and machine learning [20]. However, most evaluations are constrained to medical records written in the English language. The generalizability of de-identification methods across languages and domains is largely unexplored.

To test the generalizability of existing de-identification methods, we annotated a new dataset of 1260 medical records from three sectors of Dutch healthcare: elderly care, mental care and disabled care (Section 3). Figure 1 shows an example record with annotated PHI. We then compare the performance of the following three de-identification methods on this data (Section 4):

- (1) A rule-based system named DEDUCE developed for Dutch psychiatric clinical notes [19]
- (2) A feature-based Conditional Random Field (CRF) as described in Liu et al. [17]
- (3) A deep neural network with a bidirectional long short-term memory architecture and a CRF output layer (BiLSTM-CRF) [3]

We test the transferability of each method across three domains of Dutch healthcare. Finally, the generalizability of the methods is compared across languages using two widely used English benchmark corpora (Section 5).

This paper makes three main contributions. First, our experiments show that the only openly available de-identification method for the Dutch language fails to generalize to other Dutch medical domains. This highlights the importance of a thorough evaluation of the generalizability of de-identification methods. Second, we offer a novel comparison of several state-of-the-art de-identification methods both across languages and domains. Our experiments show that a popular neural architecture generalizes best even when limited amounts of training data are available. The neural method only considers word/character sequences which we find to be sufficient and more robust across languages and domains compared to the structural features employed by traditional machine learning approaches. However, our experiments also reveal that the neural method may

Medische overdracht Datum	26-04-2017	DATE	(patiënt nr. 64088	ID)		
Instelling	Duinendaal	CARE INSTITUTE				
Datum verrichting	24-04-2017	DATE	Tijdstip 23:45			
S regel: VG ALS Heeft sonde deze is eruit, alle medicatie al gekregen. Familie is boos, dhr heeft last van slijmvorming. Is hier iets aan te doen?						
O regel: NV						
E regel: Slijmvorming						
P regel: Nu niet direct op te lossen.						
ICPC code A45.00 (Advies/observatie/voorlichting/dieet)						
Patiënt Dhr.	Jan P. Jansen	NAME	(M),	06-11-1956	DATE	Arts
	J.O. Besteman	NAME	Adres	Wite Mar 782 Kamerik	ADDRESS	
Verrichting	Telefonisch consult ANW (t:	06-7802651	PHONE/FAX)			
===== English Translation =====						
Medical transfer date	26-04-2017	DATE	(patient no. 64088	ID)		
Institution	Duinendaal	CARE INSTITUTE				
Date	24-04-2017	DATE	Time 23:45			
Subjective (S): VG ALS got feeding tube removed, already received all medication. Family is upset, Mr. suffers from increased mucus formation. Can anything be done about that?						
Objective (O): NV						
Evaluation (E): Mucus formation						
Plan (P): Cannot be solved immediately.						
ICPC code A45.00 (Advice/observation/information/diet)						
Patient Mr.	Jan P. Jansen	NAME	(M),	06-11-1956	DATE	Doctor
	J.O. Besteman	NAME	Address	Wite Mar 782 Kamerik	ADDRESS	
Provided phone consult ANW (t:	06-7802651	PHONE/FAX)				

Figure 1: Excerpt of a medical record in our dataset with annotated protected health information (PHI). Sensitive PHI was replaced with automatically generated surrogates.

still experience a substantially lower performance in new domains. A direct consequence for de-identification practitioners is that pre-trained models require additional fine-tuning to be fully applicable to new domains. Third, we share our pre-trained models and code with the research community. The creation of these resources is connected to a significant time effort and requires access to sensitive medical data. We anticipate that this resource is of direct value to text mining researchers.

This work was presented at the first Health Search and Data Mining Workshop (HSDM 2020) [8]. The implementation of the de-identification systems, pre-trained models and code for running the experiments is available at: github.com/nedap/deidentify.

2 RELATED WORK

Previous work on de-identification can be roughly organized into four groups: (1) creation of benchmark corpora, (2) approaches to de-identification, (3) work on languages other than English, and (4) cross-domain de-identification.

Various English benchmark corpora have been created including nursing notes, longitudinal patient records and psychiatric intake

notes [21, 29, 31]. Furthermore, Deléger et al. [4] created a heterogeneous dataset comprised of 22 different document types. Contrary to the existing datasets which only contain records from at most two different medical institutes, the data used in this paper was sampled from a total of 9 institutes that are active in the Dutch healthcare sector. The contents, structure and writing style of the documents strongly depend on the processes and individuals specific to an institute which contributes to a heterogeneous corpus.

Most existing de-identification approaches are either rule-based or machine learning based. Rule-based methods combine various heuristics in form of patterns, lookup lists and fuzzy string matching to identify PHI [11, 21]. The majority of machine learning approaches employ feature-based CRFs [1, 12], ensembles combining CRFs with rules [30] and most recently also neural networks [5, 18]. A thorough overview of the different de-identification methods is given in Meystre [20]. In this study, we compare several state-of-the-art de-identification methods. With respect to rule-based approaches, we apply DEDUCE, a recently developed method for Dutch data [19]. To the best of our knowledge, this is the only openly available de-identification method tailored to Dutch data. For a feature-based machine learning method, we re-implement the token-level CRF by Liu et al. [17]. Previous work on neural de-identification used a BiLSTM-CRF architecture with character-level and ELMo embeddings [5, 15]. Similarly, we use a BiLSTM-CRF but apply recent advances in neural sequence modeling by using contextual string embeddings [3].

To the best of our knowledge, we are the first study to offer a comparison of de-identification methods across languages. With respect to de-identification in languages other than English, only three studies consider Dutch data. Scheurwegs et al. [27] applied a Support Vector Machine and a Random Forest classifier to a dataset of 200 clinical records. Menger et al. [19] developed and released a rule-based method on 400 psychiatric nursing notes and treatment plans of a single Dutch hospital. Tjong Kim Sang et al. [33] evaluated an existing named entity tagger for the de-identification of autobiographic emails on publicly available Wikipedia texts. Furthermore, de-identification in several other languages has been studied including German, French, Korean and Swedish [22, 26].

With respect to cross-domain de-identification, the 2016 CEGS N-GRID shared task evaluated the portability of pre-trained de-identification methods to a new set of English psychiatric records [29]. Overall, the existing systems did not perform well on the new data. Here, we provide a similar comparison by cross-testing on three domains of Dutch healthcare.

3 DATASETS

This section describes the construction of our Dutch benchmark dataset called NUT (Nedap/University of Twente). The data was sampled from 9 healthcare institutes and annotated for PHI according to a tagging scheme derived from Stubbs and Uzuner [31]. Furthermore, following common practice in the preparation of de-identification corpora, we replaced PHI instances with realistic surrogates to comply with privacy regulations. To compare the performance of the de-identification methods across languages, we use the English i2b2/UTHealth and the nursing notes corpus [21, 31]. An overview of the three datasets can be found in Table 1.

Table 1: Overview of the datasets used in this study.

Datset	NUT	i2b2 [31]	Nursing [21]
Language	Dutch	English	English
Domain(s)	elderly, mental and disabled care	clinical	clinical
Institutes	9 (3 per domain)	2	1
Documents	1260	1304	2434
Patients	1260	296	148
Tokens	448,795	1,057,302	444,484
Vocabulary	25,429	36,743	19,482
PHI categories	16	32	10
PHI instances	17,464	28,872	1779
Median PHI/doc.	9	18	0

3.1 Data Sampling

We sample data from a snapshot of the databases of 9 healthcare institutes with a total of 83,000 patients. Three domains of healthcare are equally represented in this snapshot: elderly care, mental care and disabled care. We consider two classes of documents to sample from: surveys and progress reports. Surveys are questionnaire-like forms which are used by the medical staff to take notes during intake interviews, record the outcomes of medical tests or to formalize the treatment plan of a patient. Progress reports are short documents describing the current conditions of a patient receiving care, sometimes on a daily basis. The use of surveys and progress reports differs strongly across healthcare institute and domain. In total, this snapshot consists of 630,000 surveys and 13 million progress reports.

When sampling from the snapshot described above, we aim to maximize both the variety of document types, and the variety of PHI, two essential properties of a de-identification benchmark corpus [4]. First, to ensure a wide variety of document types, we select surveys in a stratified fashion according to their type label provided by the EHR system (e.g., intake interview, care plan, etc.). Second, to maximize the variety in PHI, we sample medical reports on a patient basis: for each patient, a random selection of 10 medical reports is combined into a patient file. We then select patient files uniformly at random to ensure that no patient appears multiple times within the sample. Furthermore, to control the annotation effort, we impose two subjective limits on the document length. A document has to contain at least 50 tokens, but no more than 1000 tokens to be included in the sample. For each of the 9 healthcare institutes, we sample 140 documents (70 surveys and 70 patient files), which yields a total sample size of 1260 documents (see Table 1).

We received approval for the collection and use of our dataset from the ethics review board of our institution. Due to privacy regulations, the dataset constructed in this paper cannot be shared.

3.2 Annotation Scheme

Since the GDPR does not provide any strict rules about which types of PHI should be removed during de-identification, we base our PHI tagging scheme on the guidelines defined by the US HIPAA regulations. In particular, we closely follow the annotation guidelines and the tagging scheme used by Stubbs and Uzuner [31] which

Table 2: PHI tags used to annotate our dataset (NUT). The tagging scheme was derived from the i2b2 tags.

Category	i2b2 [31]	NUT
Name	Patient, Doctor, Username	Name Initials
Profession	Profession	Profession
Location	Room, Department Hospital, Organization	Internal Location Hospital, Organization Care Institute
	Street, City, State, ZIP, Country	Address
Age	Over 90, Under 90	Age
Date	Date	Date
Contact	Phone, FAX, Email URL, IP	Phone/FAX, Email URL/IP
IDs	SSN, 8 fine-grained ID tags	SSN, ID
Other	Other	Other

consists of 32 PHI tags among 8 classes: *Name*, *Profession*, *Location*, *Age*, *Date*, *Contact Information*, *IDs* and *Other*. The *Other* category is used for information that can be used to identify a patient, but which does not fall into any of the remaining categories. For example, the sentence “*the patient was a guest speaker on the subject of diabetes in the Channel 2 talkshow.*” would be tagged as *Other*. It is worth mentioning that this tagging scheme does not only capture direct identifiers relating to a patient (e.g., name and date of birth), but also indirect identifiers that could be used in combination with other information to reveal the identity of a patient. Indirect identifiers include, for example, the doctor’s name, information about the hospital and a patient’s profession.

We made two adjustments to the tagging scheme by Stubbs and Uzuner [31]. First, to reduce the annotation effort, we merged some of the 32 fine-grained PHI tags to a more generic set of 16 tags (see Table 2). For example, the fine-grained location tags *Street*, *City*, *State*, *ZIP*, and *Country* were merged into a generic *Address* tag. While this simplifies the annotation process, it complicates the generation of realistic surrogates. Given an address string, one has to infer its format to replace the individual parts with surrogates of the same semantic type. We address this issue in Section 3.4. Second, due to the high frequency of care institutes in our dataset, we decided to introduce a separate *Care Institute* tag that complements the *Organization* tag. This allows for a straightforward surrogate generation where names of care institute are replaced with another care institute rather than with more generic company names (e.g., Google).

3.3 Annotation Process

Following previous work on the construction of de-identification benchmark corpora [4, 31], we employ a double-annotation strategy: two annotators read and tag the same documents. In total, 12 non-domain experts annotated the sample of 1260 medical records independently and in parallel. The documents were randomly split into 6 sets and we randomly assigned a pair of annotators to each set. To ensure that the annotators had a common understanding of

Table 3: Distribution of PHI tags in our dataset. The inter-annotator agreement (IAA) as measured by the micro-averaged F1 score is shown per category.

PHI Tag	Count	Frac. (%)	IAA
Name	9558	54.73	0.96
Date	3676	21.05	0.86
Care Institute	997	5.71	0.52
Initials	778	4.45	0.46
Address	748	4.28	0.75
Organization	712	4.08	0.38
Internal Location	242	1.39	0.29
Age	175	1.00	0.39
Profession	122	0.70	0.31
ID	114	0.65	0.43
Phone/Fax	97	0.56	0.93
Email	95	0.54	0.94
Hospital	92	0.53	0.42
Other	33	0.19	0.03
URL/IP	23	0.13	0.70
SSN	2	0.01	0.50
Total	17,464	100	0.84

the annotation instructions, an evaluation session was held after each pair of annotators completed the first 20 documents.¹ In total, it took 77 hours to double-annotate the entire dataset of 1260 documents, or approximately 3.7 minutes per document. We measured the inter-annotator agreement (IAA) using entity-level F1 scores.² Table 3 shows the IAA per PHI category. Overall, the agreement level is fairly high (0.84). However, we find that location names (i.e., care institutes, hospitals, organizations and internal locations) are often highly ambiguous which is reflected by the low agreement scores of these categories (between 0.29 and 0.52).

To improve annotation efficiency, we integrated the rule-based de-identification tool DEDUCE [19] with our annotation software to pre-annotate each document. This functionality could be activated on a document basis by each annotator. If an annotator used this functionality, they had to review the pre-annotations, correct potential errors and check for missed PHI instances. During the evaluation sessions, annotators mentioned that the existing tool proved helpful when annotating repetitive names, dates and email addresses. Note that this pre-annotation strategy might give DEDUCE a slight advantage. However, the low performance of DEDUCE in the formal benchmark in Section 5 does not reflect this.

After annotation, the main author of this paper reviewed 19,165 annotations and resolved any disagreements between the two annotators to form the gold-standard of 17,464 PHI annotations. Table 3 shows the distribution of PHI tags after adjudication. Overall the adjudication has been done risk-averse: if only one annotator identified a piece of text as PHI, we assume that the other annotator

missed this potential PHI instance. In addition to the manual adjudication, we performed two automatic checks: (1) we ensured that PHI instances occurring in multiple files received the same PHI tag, and (2) any instances that were tagged in one part of the corpora but not in the other were manually reviewed and added to the gold-standard. We used the BRAT annotation tool for both annotation and adjudication [28].

3.4 Surrogate Generation

As the annotated dataset consists of personally identifying information which is protected by the GDPR, we generate artificial replacements for each of the PHI instances before using the data for the development of de-identification methods. This process is known as surrogate generation, a common practice in the preparation of de-identification corpora [32]. As surrogate generation will inevitably alter the semantics of the corpus to an extent where it affects the de-identification performance, it is important that this step is done as thoroughly as possible [36]. Here, we follow the semi-automatic surrogate generation procedure that has been used to prepare the i2b2/UTHealth shared task corpora. Below, we summarize this procedure and mention the language specific resources we used. We refer the reader to Stubbs et al. [32] for a thorough discussion of the method. After running the automatic replacement scripts, we reviewed each of the surrogates to ensure that continuity within a document is preserved and no PHI is leaked into the new dataset.

We adapt the surrogate generation method of Stubbs et al. [32] to the Dutch language as follows. A list of 10,000 most common family names and given names is used to generate random surrogates for name PHI instances.³ We replace dates by first parsing the format (e.g., “12 nov. 2018” → “%d %b. %Y”),⁴ and then randomly shifting all dates within a document by the same amount of years and days into the future. For addresses, we match names of cities, streets, and countries with a dictionary of Dutch locations,⁵ and then pick random replacements from that dictionary. As Dutch ZIP codes follow a standard format (“1234AB”), their replacement is straightforward. Names of hospitals, care institutes, organizations and internal locations are randomly shuffled within the dataset. PHI instances of type *Age* are capped at 89 years. Finally, alphanumeric strings such as *Phone/FAX*, *Email*, *URL/IP*, *SSN* and *IDs* are replaced by substituting each alphanumeric character with another character of the same class. We manually rewrite *Profession* and *Other* tags, as an automatic replacement is not applicable.

4 METHODS

This section presents the three de-identification methods and the evaluation procedure.

4.1 Rule-based Method: DEDUCE

DEDUCE is an unsupervised de-identification method specifically developed for Dutch medical records [19]. It is based on lookup tables, decision rules and fuzzy string matching and has been validated on a corpus of 400 psychiatric nursing notes and treatment

¹We include the annotation instructions that were provided to the annotators in the online repository of this paper. The instructions are in large parts based on the annotation guidelines in Stubbs and Uzuner [31].

²It has been shown that the F-score is more suitable to quantify IAA in sequence-tagging scenarios compared to other measures such as the Kappa score [4].

³See www.naamkunde.net, accessed 2019-12-09

⁴Rule-based date parser: github.com/nedap/dateinfer, accessed 2019-12-09

⁵See openov.nl, accessed 2019-12-09

Table 4: Features used by the CRF method. The features are identical to the one by Liu et al. [17], but we exclude word-representation features.

Group	Description
Bag-of-words (BOW)	Token unigrams, bigrams and trigrams within a window of $[-2, 2]$ of the current token.
Part-of-speech (POS) BOW + POS	Same as above but with POS n-grams. Combinations of the previous, current and next token and their POS tags.
Sentence	Length in tokens, presence of end-mark such as ‘, ‘?, ‘!’ and whether sentence contains unmatched brackets.
Affixes	Prefix and suffix of length 1 to 5.
Orthographic	Binary indicators about word shape: is all caps, is capitalized, capital letters inside, contains digit, contains punctuation, consists of only ASCII characters.
Word Shapes	The abstract shape of a token. For example, “7534-Df” becomes “####-Aa”.
Named-entity recognition (NER)	NER tag assigned by the spaCy tagger.

plans of a single hospital. Following the authors’ recommendations, we customize the method to include a list of 1200 institutions that are common in our domain. Also, we resolve two incompatibilities between the PHI coding schemes of our dataset and the DEDUCE output. First, as DEDUCE does not distinguish between hospitals, care institutes, organizations and internal locations, we group these four PHI tags under a single *Named Location* tag. Second, our *Name* annotations do not include titles (e.g., “Dr.” or “Ms.”). Therefore, titles are stripped from the DEDUCE output.

4.2 Feature-based Method: Conditional Random Field

CRFs and hybrid rule-based systems provide state-of-the-art performance in recent shared tasks [29, 30]. Therefore, we implement a CRF approach to contrast with the unsupervised rule-based system. In particular, we re-implement the token-based CRF method by Liu et al. [17] and re-use a subset⁶ of their features (see Table 4). The linear-chain CRF is trained using LBFSGS and elastic net regularization [37]. Using a validation set, we optimize the two regularization coefficients of the L_1 and L_2 norms with a random search in the \log_{10} space of $[10^{-4}, 10^1]$ with 250 trials. We use the *CRFSuite* implementation by Okazaki [23].

4.3 Neural Method: BiLSTM-CRF

To reduce the need for hand-crafted features in traditional CRF-based de-identification, recent work applies neural methods [5, 15, 18]. Here, we re-implement a BiLSTM-CRF architecture with contextual string embeddings, which has recently shown to provide

⁶We disregard word-representation features as Liu et al. [17] found that they had a negative performance impact.

state-of-the-art results for sequence labeling tasks [3]. Hyperparameters are set to the best performing configuration in Akbik et al. [3]: we use stochastic gradient descent with no momentum and an initial learning rate of 0.1. If the training loss does not decrease for 3 consecutive epochs, the learning rate is halved. Training is stopped if the learning rate falls below 10^{-4} or 150 epochs are reached. Furthermore, the number of hidden layers in the LSTM is set to 1 with 256 recurrent units. We employ locked dropout with a value of 0.5 and use a mini-batch size of 32. With respect to the embedding layer, we use the pre-trained GloVe (English) and fasttext (Dutch) embedding on a word-level, and concatenate them with the pre-trained contextualized string embeddings included in Flair⁷ [2, 10, 24].

4.4 Preprocessing and Sequence Tagging

We use a common preprocessing routine for all three datasets. For tokenization and sentence segmentation, the spaCy tokenizer is used.⁸ The POS/NER features of the CRF method are generated by the built-in spaCy models. After sentence segmentation, we tag each token according to the Beginning, Inside, Outside (BIO) scheme. In rare occasions, sequence labeling methods may produce invalid transitions (e.g., $O- \rightarrow I-$). In a post-processing step, we replace invalid $I-$ tags with $B-$ tags [25].

4.5 Evaluation

The de-identification methods are assessed according to precision, recall and F1 computed on an entity-level, the standard evaluation approach for NER systems [34]. In an entity-level evaluation, predicted PHI offsets and types have to match exactly. Following the evaluation of de-identification shared tasks, we use the micro-averaged entity-level F1 score as primary metric [30].⁹

We randomly split our dataset and the nursing notes corpus into training, validation and testing sets with a 60/20/20 ratio. As the i2b2 corpus has a pre-defined test set of 40%, a random set of 20% of the training documents serves as validation data. Finally, we test for statistical significance using two-sided approximate randomization with $N = 9999$ [35].

5 RESULTS

In this section, we first discuss the de-identification results obtained on our Dutch dataset (Section 5.1). Afterwards, we present an error analysis of the best performing method (Section 5.2). This section is concluded with the benchmark for the English datasets (Section 5.3) and the cross-domain de-identification (Section 5.4).

5.1 De-identification of Dutch Dataset

Both machine learning methods outperform the rule-based system DEDUCE by a large margin (see Table 5). Furthermore, the BiLSTM-CRF provides a substantial improvement of 10% points in recall over the traditional CRF method, while maintaining precision. Overall, the neural method has an entity-level recall of 87.1% while achieving

⁷github.com/zalando-research/flair, accessed 2019-12-09

⁸spacy.io, accessed 2019-12-09

⁹De-identification systems are often also evaluated on a less strict token-level. As a system that scores high on an entity-level will also score high on a token-level, we only measure according to the stricter level of evaluation.

Table 5: Evaluation summary: micro-averaged scores are shown for each dataset and method. Statistically significant improvements over the score on the previous line are marked with \blacktriangle ($p < 0.01$), and \circ depicts no significance. The rule-based method DEDUCE is not applicable to the English datasets.

Method	NUT (Dutch)			i2b2 (English)			Nursing Notes (English)		
	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1
DEDUCE	0.807	0.564	0.664	-	-	-	-	-	-
CRF	0.919\blacktriangle	0.775\blacktriangle	0.841\blacktriangle	0.952	0.796	0.867	0.914	0.685	0.783
BiLSTM-CRF	0.917 \circ	0.871\blacktriangle	0.893\blacktriangle	0.959\blacktriangle	0.869\blacktriangle	0.912\blacktriangle	0.886 \circ	0.797\blacktriangle	0.839\blacktriangle

Table 6: Entity-level precision and recall per PHI category on the NUT dataset. Scores are compared between the rule-based tagger DEDUCE [19] and the BiLSTM-CRF model. The *Named Loc.* tag is the union of the 4 specific location tags which are not supported by DEDUCE. Tags are ordered by frequency with location tags fixated at the bottom.

PHI Tag	BiLSTM-CRF		DEDUCE	
	Prec.	Rec.	Prec.	Rec.
Name	0.965	0.956	0.849	0.805
Date	0.926	0.920	0.857	0.441
Initials	0.828	0.624	0.000	0.000
Address	0.835	0.846	0.804	0.526
Age	0.789	0.732	0.088	0.122
Profession	0.917	0.262	0.000	0.000
ID	0.800	0.480	0.000	0.000
Phone/Fax	0.889	1.000	0.929	0.812
Email	0.909	1.000	1.000	0.900
Other	0.000	0.000	0.000	0.000
URL/IP	1.000	0.750	0.750	0.750
Named Loc.	0.797	0.659	0.279	0.058
Care Institute	0.686	0.657	n/a	n/a
Organization	0.780	0.522	n/a	n/a
Internal Loc.	0.737	0.509	n/a	n/a
Hospital	0.778	0.700	n/a	n/a

a recall of 95.6% for names, showing that the neural method is operational for many de-identification scenarios. In addition, we make the following observations.

Neural method performs at least as good as rule-based method. By inspecting the model performance on a PHI-tag level, we observe that the neural method outperforms DEDUCE for all classes of PHI (see Table 6). Only for the *Phone* and *Email* category, the rule-based method has a slightly higher precision. Similarly, we studied the impact of the training data set size on the de-identification performance. Both machine learning methods outperform DEDUCE even with as little training data as 10% of the total sentences (see Figure 2). This suggests that in most environments where training data are available (or can be obtained), the machine learning methods are to be preferred.

Rule-based method can provide a “safety net.” It can be observed that DEDUCE performs reasonably well for names, phone numbers, email addresses and URLs (see Table 6). As these PHI

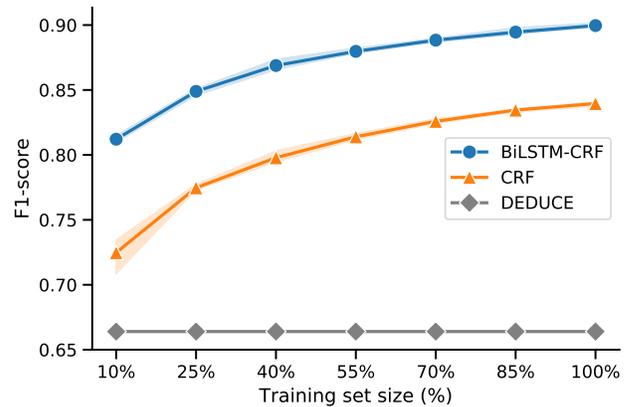


Figure 2: Entity-level F1-score for varying training set sizes. The full training set (100%) consists of all training and validation sentences in NUT (34,714). The F1-score is measured on the test set. For each subset size, we draw 3 random samples and train/test each model 3 times. The lines show the averaged scores along with the 95% confidence interval. The rule-based tagger DEDUCE is shown as a baseline.

instances are likely to directly reveal the identity of an individual, their removal is essential. However, DEDUCE does not generalize beyond the PHI types mentioned above. Especially named locations are non-trivial to capture with a rule-based system as their identification strongly relies on the availability of exhaustive lookup lists. In contrast, the neural method provides a significant improvement for named locations (5.8% vs. 65.9% recall). We assume that word-level and character-level embeddings provide an effective tool to capture these entities.

Initials, IDs and professions are hard to detect. During annotation, we observed a low F1 annotator agreement of 0.46, 0.43, and 0.31 for initials, IDs and professions, respectively. This shows that these PHI types are among the hardest to identify, even for humans (see Table 3). One possible cause for this is that IDs and initials are often hard to discriminate from abbreviations and medical measurements. We observe that the BiLSTM-CRF detects those PHI classes with high precision but low recall. With respect to professions, we find that phrases are often wrongly tagged. For example, colloquial job descriptions (e.g., “works behind the cash desk”) as opposed to the job title (e.g., “cashier”) make it infeasible

to tackle this problem with lookup lists, while a machine learner likely requires more training data to capture this PHI.

5.2 Error Analysis on Dutch Dataset

To gain a better understanding of the best performing model and an intuition for its limitations, we conduct a manual error analysis of the false positives (FPs) and false negatives (FNs) produced by the BiLSTM-CRF on the test set. We discuss the error categorization scheme in Section 5.2.1 and present the results in Section 5.2.2.

5.2.1 Error Categorization. We distinguish between two error groups: (1) modeling errors, and (2) annotation/preprocessing errors. We define modeling errors to be problems that can be addressed with different de-identification techniques and additional training data. In contrast, annotation and preprocessing errors are not directly caused by the sequence labeling model, but are issues in the training data or the preprocessing pipeline which need to be addressed manually. Inspired by the classification scheme of Deroncourt et al. [5], we consider the following sources of modeling errors:

- **Abbreviation.** PHI instances which are abbreviations or acronyms for names, care institutes and companies. These are hard to detect and can be ambiguous as they are easily confused with medical terms and measurements.
- **Ambiguity.** A human reader may be unable to decide whether a given text fragment is PHI.
- **Debatable.** It can be argued that the token should not have been annotated as PHI.
- **Prefix.** Names of internal locations, organizations and companies are often prefixed with articles (i.e., “de” and “het”). Sometimes, it is unclear whether the prefix is part of the official name or part of the sentence construction. This ambiguity is reflected in the training data which causes the model to inconsistently include or exclude those prefixes.
- **Common Language.** PHI instances consisting of common language are hard to discriminate from the surrounding text.
- **Other.** Remaining modeling errors that do not fall into the categories mentioned above. In those cases, it is not immediately apparent why the misclassification occurs.

Preprocessing errors are categorized as follows:

- **Missing Annotation.** The text fragment is PHI, but was missed during the annotation phase.
- **Annotation Error.** The annotator assigned an invalid entity boundary.
- **Tokenization Error.** The annotated text span could not be split into a compatible token span. Those tokens were marked as “Outside (O)” during BIO tagging.

We consider all error categories to be mutually exclusive.

5.2.2 Results of Error Analysis. Table 7 summarizes the error analysis results and shows the absolute and relative frequency of each error category. Overall, we find that the majority of modeling errors cannot be easily explained through human inspection (“Other reason” in Table 7). The remaining errors are mainly caused by ambiguous PHI instances and preprocessing errors. In more detail, we make the following observations:

Abbreviations are the second most common cause for modeling errors (13.9% of FNs, 9.7% of FPs). We hypothesize that more

Table 7: Summary of the manual error analysis of false negatives (FNs) and false positives (FPs) produced by the BiLSTM-CRF. All error categories are mutually exclusive.

Category	FNs ($n = 469$)		FPs ($n = 288$)	
	Count	Part	Count	Part
<i>Model Errors</i>				
Abbreviation	65	13.9%	28	9.7%
Ambiguity	15	3.2%	7	2.4%
Debatable	7	1.5%	4	1.4%
Prefix	10	2.1%	10	3.5%
Common language	35	7.5%	9	3.1%
Other reason	275	58.6%	159	55.2%
<i>Annotation/Preprocessing Errors</i>				
Missing Annotation	-	-	33	11.5%
Annotation Error	21	4.5%	18	6.3%
Tokenization Error	41	8.7%	20	6.9%
Total	469	100%	288	100%

training data will likely not in itself help to correctly identify this type of PHI. It is conceivable to design custom features (e.g., based on shape, positioning in a sentence, presence/absence in a medical dictionary) to increase precision. However, it is an open question how recall can be improved.

PHI instances consisting of common language are likely to be wrongly tagged (7.5% FNs, 3.1% FPs). This is caused by the fact that there are insufficient training examples where common language is used to refer to PHI. For example, the organization name in the sentence “Vandaag heb ik Beter Horen gebeld” (Eng: “I called Beter Horen today”) was incorrectly classified as non-PHI. Each individual word, and also the combination of the two words, can be used in different contexts without referring to PHI. However, in this specific context, it is apparent that “Beter Horen” must refer to an organization.

A substantial amount of errors is due to annotation and preprocessing issues. Annotation errors (4.5% FNs, 6.3% FPs) can be resolved by correcting the respective PHI offsets in the gold standard. Tokenization errors (8.7% FNs, 6.9% FPs) need to be fixed through a different preprocessing routine. For example, the annotation <DATE 2016>/<DATE 2017> should have been split into [2016, /, 2017] with BIO tagging [B, O, B]. However, the spaCy tokenizer segmented this text into a single token [2016/2017]. In this case, entity boundaries do no longer align with token boundaries which results in an invalid BIO tagging of [O] for the entire span. **Several false positives are in fact PHI and should be annotated.** The model identifies several PHI instances which were missed during the annotation phase (11.5% of the FPs). Once more, this demonstrates that proper de-identification is an error-prone task for human annotators.

5.3 De-identification of English Datasets

When training and testing both machine learning methods on the English i2b2 and the nursing notes datasets, we can observe that the BiLSTM-CRF significantly outperforms the CRF in both cases

Table 8: Summary of the transfer learning experiment on our Dutch dataset. Each method is trained on data of one care domain and tested on the other two domains. All scores are micro-averaged entity-level F1.

Method	Training Domain		
	Elderly	Disabled	Mental
DEDUCE	0.683	0.565	0.675
CRF	0.414	0.697	0.719
BiLSTM-CRF	0.775	0.775	0.839

(see Table 5). Similar to our Dutch dataset, the neural method provides an increase of up to 11.2% points in recall (nursing notes) while the precision remains relatively stable. This shows that the neural method has the best generalization capabilities even across languages. More importantly, it does not require the development of domain-specific lookup lists or sophisticated pattern matching rules. To put the results into perspective: the second-highest ranked team in the i2b2 2014 challenge used a sophisticated ensemble combining a CRF with domain-specific rules [30]. Their system obtained an entity-level F1 score of 0.9124 which is on-par with the performance of our neural method that requires no configuration. We can expect that the performance of the neural method further improves after hyperparameter optimization. Finally, note that both machine learning methods can be easily applied to a new PHI tagging scheme, whereas rule-based methods are limited to the PHI definition they were developed for.

5.4 Cross-domain De-identification

In many de-identification scenarios, heterogeneous training data from multiple medical institutes and domains are rarely available. This raises the question, how well a model that has been trained on a homogeneous set of medical records generalizes to records of other medical domains. We trained the three de-identification methods on one domain of Dutch healthcare (e.g., elderly care) and tested each model on the records of the remaining two domains (e.g., disabled care and mental care). We followed the same training and evaluation procedures described in Section 4.5. Table 8 summarizes the performance of each method on the different tasks.

Again, the neural method consistently outperforms the rule-based and feature-based methods in all three domains which suggests that it is a fair default choice for de-identification. This is underlined by the fact that the amount of training data is severely limited in this experiment: each domain only has 420 documents of which 20% of the records are reserved for testing. Interestingly, DEDUCE performs rather stable and even outperforms the CRF within the domain of elderly care.

Given an ideal de-identification method, one would expect that performance on unseen data of a different domain is similar to the test score obtained on the available (homogeneous) data. Table 9 shows a performance breakdown for each of the three testing domains for the neural method. It can be seen that in 4 out of 6 cases, the test score in a new domain is lower than the test score obtained on the in-domain data. The largest delta of the observed in-domain test score (disabled care, 0.919 F1) and the performance in

Table 9: Detailed performance analysis of the BiLSTM-CRF method in the transfer learning experiment. In-domain test scores are shown on the diagonal. All scores are micro-averaged entity-level F1.

Test Domain	Training Domain		
	Elderly	Disabled	Mental
Elderly	0.746	0.698	0.703
Disabled	0.796	0.919	0.879
Mental	0.744	0.806	0.871

the transfer domain (elderly care, 0.698 F1) is 0.221 in F1. This raises an important point when performing de-identification in practice: while the neural method shows the best generalization capabilities compared to the other de-identification methods, the performance can still be significantly lower when applying a pre-trained model in new domains.

5.5 Limitations

While the contextual string embeddings used in this paper have shown to provide state-of-the-art results for NER [3], transformer-based architectures for contextual embeddings have also gained significant attention (e.g., BERT [6]). It would make an interesting experiment to benchmark different types of pre-trained embeddings for the task of de-identification. Furthermore, we observe that the neural method provides strong performance even with limited training data (see Figure 2). It is unclear what contribution large pre-trained embeddings have in those scenarios which warrants an ablation study testing different model configurations. We leave the exploration of those ideas to future research.

6 CONCLUSION

This paper presents the construction of a novel Dutch dataset and a comparison of state-of-the-art de-identification methods across Dutch and English medical records. Our experiments show the following. (1) An existing rule-based method for the Dutch language does not generalize well to new domains. (2) If one is looking for an out-of-the-box de-identification method, neural approaches show the best generalization performance across languages and domains. (3) When testing across different domains, a substantial decrease of performance has to be expected, an important consideration when applying de-identification in practice.

There are several directions for future work. Motivated by the limited generalizability of pre-trained models across different domains, transfer learning techniques can provide a way forward. A preliminary study by Lee et al. [16] shows that they can be beneficial for de-identification. Finally, our experiments show that phrases such as professions are among the most difficult information to de-identify. It is an open challenge how to design methods that can capture this type of information.

REFERENCES

- [1] John S. Aberdeen, Samuel Bayer, Reyhan Yeniterzi, Benjamin Wellner, Cheryl Clark, David A. Hanauer, Bradley Malin, and Lynette Hirschman. 2010. The

- MITRE Identification Scrubber Toolkit: Design, training, and assessment. *I. J. Medical Informatics* 79, 12 (2010), 849–859.
- [2] Alan Akbik, Tanja Bergmann, and Roland Vollgraf. 2019. Pooled Contextualized Embeddings for Named Entity Recognition. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 724–728. <https://doi.org/10.18653/v1/N19-1078>
 - [3] Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual String Embeddings for Sequence Labeling. In *Proceedings of the 27th International Conference on Computational Linguistics*. Association for Computational Linguistics, 1638–1649. <https://www.aclweb.org/anthology/C18-1139>
 - [4] Louise Deléger, Qi Li, Todd Lingren, Megan Kaiser, Katalin Molnár, Laura Stoutenborough, Michal Kouril, Keith Marsolo, and Imre Solti. 2012. Building Gold Standard Corpora for Medical Natural Language Processing Tasks. In *AMIA 2012, American Medical Informatics Association Annual Symposium, Chicago, Illinois, USA, November 3-7, 2012*.
 - [5] Franck Dernoncourt, Ji Young Lee, Özlem Uzuner, and Peter Szolovits. 2017. De-identification of patient notes with recurrent neural networks. *JAMIA* 24, 3 (2017), 596–606.
 - [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *Computing Research Repository* arXiv:1810.04805 (2018). <http://arxiv.org/abs/1810.04805>
 - [7] Margaret Douglass, Gari D. Clifford, Andrew Reisner, George B. Moody, and Roger G. Mark. 2004. Computer-assisted de-identification of free text in the MIMIC II database. In *Computers in Cardiology, 2004*. IEEE, 341–344.
 - [8] Carsten Eickhoff, Yubin Kim, and Ryen White. 2020. Overview of the Health Search and Data Mining (HSDM 2020) Workshop. In *Proceedings of the Thirteenth ACM International Conference on Web Search and Data Mining (WSDM '20)*. ACM, New York, NY, USA. <https://doi.org/10.1145/3336191.3371879>
 - [9] GDPR. 2016. Regulation on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (Data Protection Directive). *Official Journal of the European Union* L119 (2016), 1–88.
 - [10] Edouard Grave, Piotr Bojanowski, Prakhar Gupta, Armand Joulin, and Tomas Mikolov. 2018. Learning Word Vectors for 157 Languages. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)*. <https://www.aclweb.org/anthology/L18-1550>
 - [11] Dilip Gupta, Melissa Saul, and John Gilbertson. 2004. Evaluation of a Deidentification (De-Id) Software Engine to Share Pathology Reports and Clinical Documents for Research. *American Journal of Clinical Pathology* 121, 2 (2004), 176–186.
 - [12] Bin He, Yi Guan, Jianyi Cheng, Keting Cen, and Wenlan Hua. 2015. CRFs Based De-identification of Medical Records. *Journal of Biomedical Informatics* 58, S (2015), S39–S46.
 - [13] HIPAA. 1996. Health Insurance Portability and Accountability Act. *Public Law* 104-191 (1996).
 - [14] HIPAA. 2012. *Guidance regarding methods for de-identification of protected health information in accordance with the Health Insurance Portability and Accountability Act (HIPAA) Privacy Rule*. Retrieved December 09, 2019 from <https://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-identification/index.html>
 - [15] Kaung Khin, Philipp Burckhardt, and Rema Padman. 2018. A Deep Learning Architecture for De-identification of Patient Notes: Implementation and Evaluation. *Computing Research Repository* arXiv:1810.01570 (2018). <http://arxiv.org/abs/1810.01570>
 - [16] Ji Young Lee, Franck Dernoncourt, and Peter Szolovits. 2018. Transfer Learning for Named-Entity Recognition with Neural Networks. In *Proceedings of the 11th Language Resources and Evaluation Conference*. Miyazaki, Japan, 4470–4473. <https://www.aclweb.org/anthology/L18-1708>
 - [17] Zengjian Liu, Yangxin Chen, Buzhou Tang, Xiaolong Wang, Qingcai Chen, Haodi Li, Jingfeng Wang, Qiwen Deng, and Suisong Zhu. 2015. Automatic de-identification of electronic medical records using token-level and character-level conditional random fields. *Journal of Biomedical Informatics* 58 (2015), S47–S52.
 - [18] Zengjian Liu, Buzhou Tang, Xiaolong Wang, and Qingcai Chen. 2017. De-identification of Clinical Notes via Recurrent Neural Network and Conditional Random Field. *Journal of Biomedical Informatics* 75, S (2017), S34–S42.
 - [19] Vincent Menger, Floor Scheepers, Lisette Maria van Wijk, and Marco Spruit. 2018. DEDUCE: A pattern matching method for automatic de-identification of Dutch medical text. *Telematics and Informatics* 35, 4 (2018), 727–736.
 - [20] Stephane M. Meystre. 2015. De-identification of Unstructured Clinical Data for Patient Privacy Protection. In *Medical Data Privacy Handbook*, Aris Gkoulalas-Divanis and Grigorios Loukides (Eds.). Springer International Publishing, 697–716.
 - [21] Ishna Neamatullah, Margaret M. Douglass, Li-Wei H. Lehman, Andrew T. Reisner, Mauricio Villarroel, William J. Long, Peter Szolovits, George B. Moody, Roger G. Mark, and Gari D. Clifford. 2008. Automated de-identification of free-text medical records. *BMC Med. Inf. & Decision Making* 8 (2008), 32.
 - [22] Aurélie Névéol, Hercules Dalianis, Sumithra Velupillai, Guergana Savova, and Pierre Zweigenbaum. 2018. Clinical Natural Language Processing in languages other than English: opportunities and challenges. *Journal of Biomedical Semantics* 9, 1 (2018), 12:1–12:13.
 - [23] Naoaki Okazaki. 2007. *CRFsuite: a fast implementation of Conditional Random Fields (CRFs)*. Retrieved December 09, 2019 from <http://www.chokkan.org/software/crfsuite/>
 - [24] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word Representation. In *Empirical Methods in Natural Language Processing (EMNLP)*. 1532–1543. <https://doi.org/10.3115/v1/D14-1162>
 - [25] Nils Reimers and Iryna Gurevych. 2017. Optimal Hyperparameters for Deep LSTM-Networks for Sequence Labeling Tasks. *Computing Research Repository* arXiv:1707.06799 (2017). <http://arxiv.org/abs/1707.06799>
 - [26] Phillip Richter-Pechanski, Stefan Riezler, and Christoph Dieterich. 2018. De-Identification of German Medical Admission Notes. *Studies in health technology and informatics* 253 (2018), 165–169.
 - [27] Elyne Scheurwegs, Kim Luyckx, Filip Van der Schueren, and Tim Van den Bulcke. 2013. De-Identification of Clinical Free Text in Dutch with Limited Training Data: A Case Study. In *Proceedings of the Workshop on NLP for Medicine and Biology associated with RANLP 2013*. 18–23. <https://www.aclweb.org/anthology/W13-5103>
 - [28] Pontus Stenetorp, Sampo Pyysalo, Goran Topić, Tomoko Ohta, Sophia Ananiadou, and Jun'ichi Tsujii. 2012. BRAT: A Web-based Tool for NLP-assisted Text Annotation. In *Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics (EACL '12)*. Association for Computational Linguistics, Stroudsburg, PA, USA, 102–107. <https://www.aclweb.org/anthology/E12-2021>
 - [29] Amber Stubbs, Michele Filannino, and Özlem Uzuner. 2017. De-identification of psychiatric intake records: Overview of 2016 CEGS N-GRID shared tasks Track 1. *Journal of Biomedical Informatics* 75 (2017), S4–S18.
 - [30] Amber Stubbs, Christopher Kotfila, and Özlem Uzuner. 2015. Automated systems for the de-identification of longitudinal clinical narratives: Overview of 2014 i2b2/UTHealth shared task Track 1. *Journal of Biomedical Informatics* 58 (2015), S11–S19.
 - [31] Amber Stubbs and Özlem Uzuner. 2015. Annotating longitudinal clinical narratives for de-identification: The 2014 i2b2/UTHealth corpus. *Journal of Biomedical Informatics* 58 (2015), S20–S29.
 - [32] Amber Stubbs, Özlem Uzuner, Christopher Kotfila, Ira Goldstein, and Peter Szolovits. 2015. Challenges in Synthesizing Surrogate PHI in Narrative EMRs. In *Medical Data Privacy Handbook*, Aris Gkoulalas-Divanis and Grigorios Loukides (Eds.). Springer International Publishing, 717–735.
 - [33] Erik Tjong Kim Sang, Ben de Vries, Wouter Smink, Bernard Veldkamp, Gerben Westerhof, and Anneke Sools. 2019. De-identification of Dutch Medical Text. In *2nd Healthcare Text Analytics Conference (HealTAC2019)*. Cardiff, Wales, UK.
 - [34] Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 Shared Task: Language-independent Named Entity Recognition. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003 - Volume 4 (CONLL '03)*. Association for Computational Linguistics, Stroudsburg, PA, USA, 142–147. <https://doi.org/10.3115/1119176.1119195>
 - [35] Alexander Yeh. 2000. More Accurate Tests for the Statistical Significance of Result Differences. In *Proceedings of the 18th Conference on Computational Linguistics - Volume 2 (COLING '00)*. Association for Computational Linguistics, Stroudsburg, PA, USA, 947–953. <https://doi.org/10.3115/992730.992783>
 - [36] Reyvan Yeniterzi, John S. Aberdeen, Samuel Bayer, Benjamin Wellner, Lynette Hirschman, and Bradley Malin. 2010. Effects of personal identifier resynthesis on clinical text de-identification. *JAMA* 17, 2 (2010), 159–168.
 - [37] Hui Zou and Trevor Hastie. 2005. Regularization and variable selection via the Elastic Net. *Journal of the Royal Statistical Society, Series B* 67 (2005), 301–320.