A Short Introduction to NILE

Sheng Yu^{1*} and Tianxi Cai²

¹Partners HealthCare Personalized Medicine, Brigham and Women's Hospital & Harvard Medical School, Boston MA

²Department of Biostatistics, Harvard School of Public Health, Boston MA

Abstract

In this paper, we briefly introduce the Narrative Information Linear Extraction (NILE) system, a natural language processing library for clinical narratives. NILE is an experiment of our ideas on efficient and effective medical language processing. We introduce the overall design of NILE and its major components, and show the performance of it in real projects.

1 Introduction

The electronic medical record (EMR) is a rich source of clinical information that can substantially support biomedical research and healthcare improvement [1, 2]. In addition to the structured data such as billing codes, EMR consists of free text clinical narrative recorded by physicians. The unstructured nature of the narrative data necessitates the need for natural language processing (NLP) technology to unlock the information they contain for further analysis. Clinical narratives also possess unique linguistic features both in syntax and semantics, making generic NLP programs and modules not directly applicable to the interpretation of EMR. For instance, a question mark, while normally seen as the end of a sentence, is usually placed before a disease name in clinical narratives, meaning "to diagnose"; a plus sign, when used after a lab test, means "positive" or "above normal". Medical language processing (MLP) is a subfield of NLP that addresses to these distinct features in clinical text.

The work flow in most MLP systems has at least two stages: named entity recognition and semantic analysis. Named entity recognition (NER), also known as entity identification or entity extraction,

^{*}Corresponding author. Email: syu7@partners.org

is the process of locating atom elements in text and classifying them into predefined categories, such as disease, symptom, quantity, time, etc. Multitudes of approaches exist, ranging from a simple and straightforward string matching to a subsystem that matches phrases of interest (such as noun phrases) against a dictionary after doing part-of-speech (POS) tagging and shallow parsing to identify certain types of phrases.

The semantic analysis step uses the result of the NER step to interpret the meaning of the clinical text. Common analyses about a mention of a named entity (disease, symptom, medication, operation, etc.) include whether it is positive, negative, conditional, or whether it is a speculation, a general discussion, past medical history, or family history. Common analyses about the state of a disease/problem include determining its location, severity, and chronicity. There are plenty of other specialized analyses. Examples are determining the smoking state of a patient [3], value extraction from semi-structured EMR [4, 5], extraction of drug and food allergies [6] and medication information [7]. Again, there are many ways to achieve the goal. In addition to using grammar and statistical models, it is also popular to use simple rules (e.g., any mention within a certain range after "no" is negative) and pattern matching, both of which have been proved effective in practice. [8, 9]

The Narrative Information Linear Extraction (NILE) System is our experiment of several ideas on MLP. One of the ideas is that a straightforward string matching for longest left-most phrases (similar to the behavior of DFA regex engines) should work great for NER in English clinical narratives. Another idea is that semantic analysis can be performed effectively and efficiently in most cases by scanning and interpreting the sentence linearly with finite-state machines, which is similar to how humans read. If these ideas work, then one may prefer such programming paradigm for MLP for its coding simplicity and execution speed.

2 How NILE Works

The general work flow of NILE is similar to many other MLP systems. It has three stages: preprocessing, named entity recognition, and semantic analysis. Figure 1 shows some of the major components of NILE with illustration on the processing flow.

2.1 Preprocessing

The preprocessing stage includes sentence boundary detection and tokenization. Section recognition is currently not a component of NILE, and should to be done outside the library.

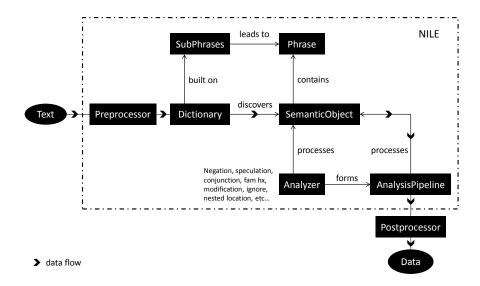


Figure 1: NILE Flow Chart

2.2 Named Entity Recognition

NILE uses a generalized prefix tree as its dictionary data structure, with tokens (mostly, words) at intermediate nodes and full phrase information at leaf nodes. The matching algorithm reads a sentence as a series of tokens, and matches the longest phrase (by the number of tokens) from the left. For instance, when NILE sees "patient had a heart attack in 2006..." it identifies "heart attack" rather than "heart", and goes on to find the next phrase starting from "in". The algorithm is capable of handling prefix and suffix sharing, so it can interpret "no mediastinal, hilar, or axillary lymphadenopathy" as "no mediastinal lymphadenopathy, hilar lymphadenopathy, or axillary lymphadenopathy", and "right upper, middle, and lower lobes" as "right upper lobe, right middle lobe, and right lower lobe".

Conventionally, with general text, like web pages, such string matching would have high recall and low precision, due to the ambiguity of word sense. However, when limited to the context of medicine, the ambiguity has substantially reduced, making such string matching approach attractive, since its speed is among the fastest that one could have. We use hash maps for token look ups in the prefix tree, so the dictionary size hardly affect the look up speed, and the processing time is majorly proportional to the length of the input text.

One the other hand, the commonly adopted NLP approaches using POS tagging + shallow parsing (or deep parsing) could have problems. Clinical narratives contain plenty of grammar errors, and physicians frequently use all kinds of shorthands (e.g. "dx" for "diagnose" or "diagnosis", "tx" for

"treat" or "treated", "hx" for "history", etc.) and acronyms, which are not in regular dictionaries, making POS tagging and the following parsing difficult and resulting in low recall. Since such approaches do not have a clear advantage over string matching and are expensive in terms of computational cost and training cost, we would just go string matching for the speed.

The identified words and phrases are encapsulated in what the program calls "semantic objects". A semantic object contains the text of the identified phrase, the concept code of the phrase (which is shared among inflections and synonyms, e.g., "pulmonary embolism", "pulmonary emboli", and "PE" all have the same code), its semantic role, and its location in the sentence. It also contains fields to be filled by later semantic analyses, such as references to other semantic objects that modify it, and attributes such as certainty (present/absent/unclear) and experiencer (self/family). The semantic role tells the role or function of the phrase in the sentence. Categories of roles include grammatical words, meaning cues, and medical terms. Grammatical words are words that have little lexical meaning, but serve to help express meanings of the other words in the sentence. Some of the semantic roles in this category are pronouns, conjunctions, prepositions, link verbs, and auxiliary verbs. Comma"," is also a role in this category as it has important functions in grammar, Meaning cues are words and phrases that tell us how to interpret the sentence meaning. Roles of this category include confirmation cues, negation cues, speculation cues, ignore cues (to ignore patterns such as "assess for PE", "in case of PE", and "PE study" that do not indicate presence of absence of pulmonary embolism), and etc. The medical terms are concepts that are related to diagnosis or treatment, such as fact, modifier, and anatomical location, where fact can be disorder, finding, procedure, test, substance (like drugs), etc. Although locations are a type of modifiers, NILE distinguishes them from the others in order to use dedicated location analyzers.

NILE has a built-in basic dictionary of grammatical words and meaning cues that are common to all applications. The medical terms need to be populated by the application, and it is generally recommended to load all the terms that may appear in the notes of the target field, which benefits the following semantic analysis. All terms in the dictionary are easily customizable. NILE allows term sense ambiguity, i.e. a term can have multiple concept codes. However, in the current implementation, a term can only have one semantic role. Eventually, the output of NER is a list of semantic objects, where the order of semantic objects is the same as the left-to-right order of the identified terms in the sentence.

2.3 Semantic Analysis

Humans read sentences from left to right linearly to understand the meaning. We want to mimic this process with finite-state machines. Finite-state machines are a type of simple programs. In theory, they are not powerful enough to understand all languages [10], but we imagine they are good

Combination	Semantic Role
found (used alone)	confirmation cue
$\mathrm{not} + \mathrm{found}$	backward negation cue
have + found	confirmation cue
haven' t + found	negation cue
have + been + found	backward confirmation cue
$\frac{1}{1}$ haven't + been + found	backward negation cue
is + found	backward confirmation cue
isn't + found	backward negation cue

Table 1: Combined Effects of Semantic Roles

enough to extract the information interesting to MLP, such as presence and modification relations, from clinical text.

NILE does semantic analysis by sending the NER result of a sentence to a pipeline of analyzers. Each analyzer is a finite-state machine and focuses on a single task. The state typically includes what the machine has just read and an interpretation mode, such as Negation On/Off. The analyzer reads the semantic objects one by one and switches the state accordingly, generally by the semantic role of the term, the distance from the previous term, and occasionally the actual text. Recall that the identified terms include grammatical words, which are hints to how to interpret the text. Also recall that clinical texts contain plenty of grammar errors, thus we avoid using strict grammar. The result is that we use a mixture of grammar and pattern analysis.

In this section we introduce some analyzers of NILE to show that this proposed methodology is very capable at capturing information in MLP tasks.

2.3.1 Presence

Analysis of presence is for determining whether a condition is present, a drug is prescribed/used, and assigns Yes, No, or Maybe to the presence attribute of each fact semantic object. It combines negation and certainty analyses, but differs slightly in that it reports the presence at the time when writing the document. For example, "the tumor was removed" means that the tumor is not present, though it existed. The analyzer uses meaning cues and grammatical words to determine the meaning and its scope. It also looks at combinations of meaning cues and grammatical words, and can merge phrases and modify semantic roles accordingly. For example, originally in the dictionary, "found" is a past participle for confirmation, "not" is a negation cue, "is" and "been" are positive link verbs, "is not" and "isn't" are negative link verbs, "have" is an positive auxiliary verb, and "have not" and "haven't" are negative auxiliary verbs. Table 1 shows the effects when these phrases appear together.

Example Sentence	NILE Return
No filling defects are seen to suggest	filling defects: NO
pulmonary embolism.	pulmonary embolism: NO
No filling defects are seen, suggesting	filling defects: NO
pulmonary embolism.	pulmonary embolism: YES
No filling defects are seen and it suggests	filling defects: NO
$pulmonary\ embolism.$	pulmonary embolism: YES
No filling defects are seen to suggest	filling defects: NO
$pulmonary\ embolism.$	pulmonary embolism: NO
No filling defects are seen to suggest	filling defects: NO
pulmonary embolism.	pulmonary embolism: YES

Table 2: NILE Negation Analysis

The presence analyzer can correctly understand the meaning of the sentences most of the times, including subtle differences as shown in Table 2 (though the example sentences may not make sense clinically). It also understands whether negation applies to a fact or a certain aspect of the fact. For instance, in "no change in the pleural effusion", the negation only applies to "change" but not to "pleural effusion", where "change" is a recognizable term with a semantic role that represents a fact attribute.

Overall, the presence analyzer covers most of the ways that physicians express negation and speculation. There are certain patterns the current implementation does not cover. For instance, it does not do cross-sentence inference, such as in "The study is of good diagnostic quality for pulmonary embolus and the vein thrombosis. There is no evidence of either", it cannot apply the negation from the second sentence to "pulmonary embolus" and "vein thrombosis" in the first sentence. Also, it would be better to refine the result by adding a time property, as illustrated by the above tumor example.

2.3.2 Location

Location analysis is another sophisticated component of NILE. Descriptions of anatomical locations usually involve compounded nested modifications, which can be expressed in very diverse ways. The simplest patterns of nested modifications are "location A location B" (e.g. "left upper lobe arteries") and "location A of/in/within/... location B" (e.g. "arteries of the left upper lobe"). NILE treats the former as "location B (location A)" and the latter as "location A (location B)", where the inside of the parentheses modifies the outside. These patterns can be chained. For instance, "location A location B location C" is interpreted as "location C (location B (location A))". However, it becomes difficult when conjunction comes in. Take a look at the following sentence, which is an excerpt from a CT pulmonary angiography (CTPA) report that we have processed:

"There are segmental and subsegmental filling defects in the right upper lobe, superior segment of the right lower lobe, and subsegmental filling defect in the in the anterolateral segment of the left lower lobe pulmonary arteries. [sic]"

Apparently, this sentence has problems. It is better to use "and" to replace the first comma, because the "and" after the second comma does not connect to a third location modifier of the "filling defects" of the first line; instead, it is another "filling defect". And "in the in the" is obviously a typo. Despite of these, this sentence is a good example that includes multiple entities, each entity having multiple location modifiers, some of which being nested and some being not. The location analyzer reads the semantic objects in sequence, looks at their semantic roles, and uses an empirical rule to determine the modification relations. Also, as mentioned before, NILE is designed to tolerate certain level of grammar problems in clinical text. For the above sentence, NILE returned:

```
filling defects: YES (right upper lobe; superior segment (right lower lobe); segmental; subsegmental),
```

filling defect: YES (segment (pulmonary arteries (left lower lobe)); subsegmental), which is the same as the intended meaning, except that "anterolateral segment" was not in our dictionary for that project.

This location analyzer is written for general purpose. There are also times when domain knowledge is needed. For instance, consider the following sentences:

- 1. There are filling defects in the right upper lobe and superior segment of the right lower lobe.
- 2. There are filling defects in the anterior basal segment and superior segment of the right lower lobe.

These two sentences have identical structures, but the first one means

```
filling defects: YES (right upper lobe; superior segment (right lower lobe)), while the second means
```

filling defects: YES (anterior basal segment (right lower lobe); superior segment (right lower lobe)).

We know that they should be such because we know that: (i) the superior segment is comparable to the anterior basal segment, and both are in the right lower lobe; and (ii) the superior segment is not of the same level as the right upper lobe. If one wishes to write a dedicated analyzer for a domain, such free-form knowledge would be hard to fit into statistical learning models, but would be easily coded if using the proposed analysis methodology.

2.3.3 Other Analyses

NILE uses a generic analyzer to handle all other modifications and the results are not nested. Another important analyzer determines if some part of the sentence should be ignored by looking for ignore cues such as "exam", "evaluate", and "diff diagnosis". The default scope to ignore is the whole sentence, but certain semantic objects such as confirmation and negation cues can modify that. Another analyzer checks if the sentence contains words like "mother", "uncle", etc (there is a semantic role for relatives). If it does, every fact in the sentence will be labeled as family history. New analyzers are actively being added to the library as new need arises.

In the end, after all the analyses are done, NILE returns only the semantic objects that are facts, which may contain other semantic objects like locations and modifiers in their attributes.

3 Preliminary Results

NILE has been used in several research occasions. In a study to detect pulmonary embolism from CTPA reports [11], a statistical algorithm using information extracted by NILE achieved AUC = 0.998 (cross-validated) in classifying PE present/absent, and AUC = 0.986 in detecting subsegmental PE. Other similar studies that used NILE includes classifying whether patients have rheumatoid arthritis and congestive heart failure.

The processing speed of NILE is fast enough for real time analysis. It took 3098 ms (2361 ms inside NILE) to analyze 10330 CTPA reports (an 18 megabyte text file on a local drive). In another project, NILE processed 21 million clinical notes from a remote SQL Server database in 21 hours, or 3.6 ms per note (this is the overall time that includes database querying, data transfer over the network and decryption, additional text cleaning, and writing to the local hard drive, which made up the majority of the processing time). The former job was done on Windows 7 64-Bit machine using a single thread of an Intel Core i7-2600 CPU @ 3.40GHz. The latter job was done on a CentOS 6.5 (64-Bit) virtual machine using a single core of Intel Xeon X5450 @ 3.00GHz. Both used Java HotSpot 64-Bit Server VM.

The accuracy and speed demonstrate the effectiveness and efficiency of the library. In addition, the subsegmental PE detection task could not be done without the nesting modification structure of the locations, which is a unique feature of NILE and an illustration of the merit of the methodology it uses for semantic analysis. In conclusion, we consider this experiment of MLP ideas a success, and will further develop NILE following this paradigm.

References

[1] Isaac S Kohane. Using electronic health records to drive discovery in disease genomics. *Nature Reviews Genetics*, 12(6):417–428, 2011.

- [2] Stéphane M Meystre, Guergana K Savova, Karin C Kipper-Schuler, and John F Hurdle. Extracting information from textual documents in the electronic health record: a review of recent research. *IMIA Year Book of Medical Informatics*, 35:128–44, 2008.
- [3] Özlem Uzuner, Ira Goldstein, Yuan Luo, and Isaac Kohane. Identifying patient smoking status from medical discharge records. *Journal of the American Medical Informatics Association*, 15(1):14–24, 2008.
- [4] Alexander Turchin. Identification of clinical characteristics of large patient cohorts through analysis of free text physician notes. Master's thesis, Massachusetts Institute of Technology, 2005.
- [5] Krister J Kristianson, Henrik Ljunggren, and Lars L Gustafsson. Data extraction from a semi-structured electronic medical record system for outpatients: a model to facilitate the access and use of data for quality control and research. *Health Informatics Journal*, 15(4):305–319, 2009.
- [6] Richard H Epstein, Paul St Jacques, Michael Stockin, Brian Rothman, Jesse M Ehrenfeld, and Joshua C Denny. Automated identification of drug and food allergies entered using non-standard terminology. *Journal of the American Medical Informatics Association*, 20(5):962–968, 2013.
- [7] Hua Xu, Shane P Stenner, Son Doan, Kevin B Johnson, Lemuel R Waitman, and Joshua C Denny. Medex: a medication information extraction system for clinical narratives. *Journal of the American Medical Informatics Association*, 17(1):19–24, 2010.
- [8] Wendy W Chapman, Will Bridewell, Paul Hanbury, Gregory F Cooper, and Bruce G Buchanan. A simple algorithm for identifying negated findings and diseases in discharge summaries. *Journal of Biomedical Informatics*, 34(5):301–310, 2001.
- [9] Carol Friedman, Philip O Alderson, John HM Austin, James J Cimino, and Stephen B Johnson. A general natural-language text processor for clinical radiology. *Journal of the American Medical Informatics Association*, 1(2):161–174, 1994.
- [10] Noam Chomsky and George A Miller. Finite state languages. *Information and Control*, 1(2):91–112, 1958.
- [11] Sheng Yu, Kanako K Kumamaru, Elizabeth George, Ruth M Dunne, Arash Bedayat, Matey Neykov, Andetta R Hunsaker, Karin E Dill, Tianxi Cai, and Frank J Rybicki. Classification of ct pulmonary angiography reports by presence, chronicity, and location of pulmonary embolism with natural language processing. *Journal of Biomedical Informatics*, page http://dx.doi.org/10.1016/j.jbi.2014.08.001, 2014.