

NATURAL LANGUAGE PROCESSING FOR DATA EXTRACTION IN CLINICAL TRIALS

Sri Nitya Yanamala

M.Sc, Biochemistry, student at ClinoSol Research, Hyderabad, India

Uma L

B. Pharmacy, student at ClinoSol Research, Hyderabad, India

Annotation: The growing recognition of the significance of integrating Natural Language Processing (NLP) methods into clinical informatics research has spurred transformative advances in recent years. Typically, clinical NLP systems undergo development and evaluation based on annotations at the word, sentence, or document levels. These annotations model specific attributes and features, encompassing document content (e.g., patient status or report type), document section types (e.g., current medications, past medical history, or discharge summary), named entities and concepts (e.g., diagnoses, symptoms, or treatments), and semantic attributes (e.g., negation, severity, or temporality). The utilization of NLP in clinical informatics facilitates a nuanced understanding of medical texts, enabling the extraction of valuable information from diverse sources. This abstract highlights the evolving landscape of clinical NLP, emphasizing its role in capturing and analyzing intricate details within healthcare documents. The continued integration and refinement of NLP methods hold great promise for advancing clinical informatics research, ultimately contributing to enhanced information extraction and decision-making processes in the medical domain.

Key words: Natural Language Processing, clinical informatics, information extraction, named entities, semantic attributes, document section types, healthcare documents, clinical NLP systems, medical text analysis, document-level annotations, transformative advances, decision-making processes.

Introduction: Appropriate utilization of large data sources such as Electronic Health Records (eHealth records or EHR) databases could have a dramatic impact on health care research and delivery. Owing to the large amount of free text documentation now available in EHRs, there has been a concomitant increase in research to advance Natural Language Processing (NLP) methods and applications for the clinical domain. The field has matured considerably in recent years, addressing many of the challenges identified by Chapman et al, and meeting the recommendations by Friedman et al.

For example, the above include recommendations to address the key challenges of limited collaboration, Lack of shared resources and evaluations-approaches of crucial tasks, such as de-identification, recognition and classification of medical concepts, semantic modifiers, and temporal information. these challenges have been addressed by the organization of several shared tasks. These include the Informatics for Integrating Biology and the Bedside challenges, the Conference and Labs of the Evaluation Forum (CLEF) eHealth challenges, and the Semantic Evaluation (SemEval) challenges. These efforts have enabled a valuable platform for international NLP method development.

Furthermore, the development of open-source NLP software specifically tailored to clinical text has led to increased adoptability. Such NLP software include the *clinical Text Analysis Knowledge Extraction System (cTAKES)* and *Clinical Language Annotation, Modeling, and Processing Toolkit (CLAMP)*, information extraction and retrieval infrastructure solutions such as SemEHR, as well as general purpose tools such as the the general architecture for text engineering (GATE) and Stanford CoreNLP. New initiatives, such as the Health Natural Language Processing (hNLP) Center,⁵ also aim to facilitate the sharing of resources, which would enable further progress through availability, transparency, and reproducibility of NLP methodologies.

In recent years, the field of mental health has shown a burgeoning increase in the use of NLP strategies and methods, mainly because most clinical documentation is in free-text, but also arising from the increasing availability of other types of documents providing behavioral, emotional, and cognitive indicators as well as cues on how patients are coping with different conditions and treatments. Such texts sources include social media and online fora as well as doctor-patient interactions and online therapy, to mention a few examples. However, although there have been a few shared tasks related to mental health the field is still narrower than that of biomedical or general clinical NLP.

The maturity of NLP method development and state-of-the-art results have led to an increase in successful deployments of NLP solutions for complex clinical outcomes research. However, the methods used to evaluate and appraise NLP approaches are somewhat different from methods used in clinical research studies, although the latter often rely on the former for data preparation and extraction. There is a need to clarify these differences and to develop novel approaches and methods to bridge this gap.

This paper stems from the findings of an international one-day workshop in 2017. The objective was to explore these evaluation issues by outlining ongoing research efforts in these fields, and brought together researchers and clinicians working in the areas of NLP, informatics, mental health, and epidemiology.

Unlocking the Potential of Natural Language Processing in Clinical Trial Data Extraction

Clinical trials stand as the cornerstone of medical research, playing a pivotal role in advancing our understanding of the safety and efficacy of new treatments. These trials, meticulously designed and conducted, serve as the bedrock upon which medical advancements are built. However, the landscape of clinical trials is evolving rapidly, marked by an exponential growth in both the volume and complexity of clinical data. In this ever-expanding sea of information, the need to efficiently extract, analyze, and interpret valuable insights becomes paramount. This is where Natural Language Processing (NLP) emerges as a revolutionary tool, poised to transform the landscape of clinical trial data analysis.

At its core, NLP is a subfield of artificial intelligence (AI) that focuses on enabling machines to comprehend, interpret, and generate human-like language. In the context of clinical trials, NLP proves to be a game-changer by offering sophisticated techniques to process and derive meaning from the vast amount of textual data generated throughout the trial lifecycle. From patient records and electronic health records (EHRs) to clinical narratives and medical literature, NLP has the potential to unlock a treasure trove of information that might otherwise remain buried in unstructured text.

The complexity of clinical trial data is multifaceted, ranging from the inclusion of diverse data types (such as structured, semi-structured, and unstructured data) to the intricacies of medical terminology, patient narratives, and healthcare jargon. NLP excels in handling this complexity, providing a set of tools and techniques that go beyond traditional methods of data extraction. By harnessing the power of machine learning, deep learning, and linguistic analysis, NLP algorithms can discern patterns, relationships, and trends within clinical text that may elude manual review or conventional data processing methods.[1]

One of the key areas where NLP demonstrates its transformative potential is in the extraction of structured information from unstructured clinical text. Clinical documents, such as medical records and trial reports, often contain a wealth of unstructured information that is critical for understanding patient outcomes, treatment efficacy, and adverse events. NLP algorithms can systematically extract and organize this information, converting it into a structured format that can be easily integrated into databases, analyzed for trends, and utilized for evidence-based decision-making.

Furthermore, NLP's application extends to the identification and classification of entities and concepts within clinical text. Named Entity Recognition (NER) algorithms, a subset of NLP, can identify and categorize entities such as diagnoses, symptoms, medications, and procedures mentioned in clinical documents. This capability not only facilitates efficient information retrieval but also enhances the accuracy and granularity of data analysis. The ability to precisely identify and categorize medical entities is particularly valuable in the context of clinical trials, where detailed and accurate information is crucial for regulatory compliance, safety monitoring, and overall trial success.

In the realm of clinical narrative analysis, NLP proves to be a valuable ally. Patient narratives, often embedded in clinical notes or adverse event reports, provide qualitative insights into the patient experience and treatment impact. Analyzing these narratives manually is labor-intensive and time-consuming, making it impractical for large-scale clinical trials. NLP algorithms, however, can sift through these narratives, identifying sentiment, capturing nuances, and extracting valuable information regarding treatment efficacy, patient-reported outcomes, and adverse events. This automated analysis not only accelerates the review process but also unveils valuable qualitative data that might otherwise be overlooked.

Moreover, NLP contributes significantly to the automation of adverse event detection and monitoring. Adverse events, ranging from mild side effects to severe complications, are meticulously documented throughout the course of a clinical trial. NLP algorithms can systematically review and analyze these reports, identifying potential adverse events, assessing their severity, and correlating them with specific treatments or patient demographics. By automating this process, NLP not only expedites safety monitoring but also enhances the overall efficiency of clinical trial management.

The impact of NLP extends beyond data extraction and analysis; it has the potential to revolutionize the entire clinical trial lifecycle. In the pre-trial phase, NLP can streamline the identification of eligible patients by mining EHRs and other healthcare databases. By automating the screening process, NLP accelerates patient recruitment, a critical aspect that often determines the success or failure of a trial. The ability to efficiently identify and recruit suitable participants not only reduces trial timelines but also enhances the generalizability and external validity of study findings.

During the trial, NLP aids in real-time data monitoring and quality assurance. By continuously analyzing incoming data, NLP algorithms can identify inconsistencies, outliers, and potential data entry errors. This proactive approach to data quality ensures the integrity of trial results and provides researchers with the confidence that the data collected is accurate and reliable. Additionally, NLP can contribute to adaptive trial design by dynamically analyzing emerging data trends, allowing researchers to make informed adjustments to study protocols in real-time.

In the post-trial phase, NLP facilitates the synthesis and dissemination of trial results. Automating the extraction of key findings, statistical outcomes, and patient demographics from trial reports expedites the publication process. Moreover, NLP supports the creation of structured summaries that can be readily integrated into clinical trial registries, databases, and publications. This not only enhances the accessibility of trial results but also contributes to the transparency and reproducibility of research findings.[2]

While the potential benefits of NLP in clinical trials are vast, challenges and considerations accompany its integration. The variability in language usage, diverse documentation styles, and the dynamic nature of medical terminology pose challenges for NLP algorithms. Additionally, ensuring the privacy and security of patient data is paramount, requiring robust measures for de-identification and compliance with data protection regulations.

Clinical Trial Data Extraction of the Challenge

Clinical trials generate vast amounts of unstructured data in the form of medical records, physician notes, lab reports, and more. This data often resides in disparate formats and systems, making it arduous for researchers and analysts to extract meaningful insights efficiently. Traditional methods of data extraction involve manual review, which is not only time-consuming but also prone to human error.

The valuable insights embedded in the vast and complex landscape of clinical trial data positions NLP as a formidable ally in the pursuit of efficient, accurate, and meaningful research outcomes. From structured information extraction to entity recognition, narrative analysis, and adverse event monitoring, NLP not only addresses the challenges of data extraction but also redefines the entire clinical trial lifecycle.

As we delve deeper into the potential applications of NLP in clinical trials, it becomes evident that the extraction of structured information from unstructured text is a fundamental aspect where NLP excels. Clinical documents, often laden with unstructured information, can be systematically processed by NLP algorithms, converting this wealth of data into a structured format. This structured data, in turn, becomes a valuable resource for researchers, regulators, and healthcare professionals, facilitating evidence-based decision-making, trend analysis, and seamless integration into databases.

Named Entity Recognition (NER) is a specific area within NLP that plays a pivotal role in identifying and classifying entities and concepts within clinical text. Whether it's diagnoses, symptoms, medications, or procedures, NER algorithms can meticulously categorize these entities, contributing to the granularity and accuracy of data analysis. This capability becomes particularly crucial in the context of clinical trials, where precise and detailed information is essential for regulatory compliance, safety monitoring, and overall trial success.[3]

The analysis of clinical narratives, enriched with patient experiences and qualitative insights, is another domain where NLP showcases its prowess. Patient narratives, often buried in clinical notes or adverse event reports, provide a qualitative dimension to the data. NLP algorithms can navigate through these narratives, capturing sentiments, nuances, and essential information related to treatment efficacy, patient-reported outcomes, and adverse events.

How NLP Transforms Data Extraction

Natural Language Processing, a branch of artificial intelligence, equips machines with the ability to understand, interpret, and generate human language. In the context of clinical trials, NLP acts as a powerful tool for processing unstructured text, extracting pertinent information, and converting it into structured, actionable data.

1. Text Mining and Information Extraction

NLP techniques enable the extraction of key information from textual data. By employing methods like named entity recognition (identifying entities such as drugs, diseases, or demographics), relationship extraction, and semantic parsing, NLP algorithms can effectively identify and categorize critical data elements within clinical trial documents.

2. Streamlining Data Abstraction

One significant challenge in clinical trials is abstracting relevant data from a multitude of sources. NLP facilitates the automated extraction of pertinent information from various documents, allowing for streamlined and accurate data abstraction processes.[4]

3. Supporting Decision Making

NLP doesn't just extract data; it transforms it into actionable insights. By converting unstructured text into structured data, NLP empowers researchers and clinicians to make informed decisions based on comprehensive and organized information.

Applications in Clinical Trials

Adverse Event Monitoring: Automated adverse event detection and monitoring can be enhanced through NLP, enabling quicker identification of potential safety concerns.

Protocol Analysis: NLP can assist in analysing trial protocols, comparing them to similar studies, and extracting crucial information for optimizing study design and execution.

Overcoming Challenges

While NLP offers immense potential, it faces challenges in the clinical trial domain, including privacy concerns, data variability, and the need for domain-specific customization. Ensuring the accuracy and reliability of NLP models in understanding medical jargon and context remains a significant hurdle.

The Future Outlook

As technology advances and NLP algorithms become more sophisticated, the future of data extraction in clinical trials appears promising. Integration of machine learning, improved models trained on larger datasets, and collaborative efforts between data scientists, clinicians, and domain experts will drive further innovation in leveraging NLP for efficient and accurate data extraction in the realm of clinical research.

Natural Language Processing stands as a beacon of hope in the realm of clinical trial data extraction, holding the key to unlocking invaluable insights from the vast sea of unstructured information. As it evolves, NLP will continue to reshape the landscape of medical research, enhancing our capacity to derive actionable knowledge from the troves of clinical data available to us.

Practical Steps of the NLP-Featured Approach to Investigator Recruitment

The healthcare sector has always been of particular interest to data scientists. Many consider it a near-perfect domain to showcase NLP's value. By various estimates, 80% of medical data (i.e., from medical records, imaging devices, sensors, wearables, health documents, and articles) remains unlabelled and untapped after it was created. However, all this unstructured data when sorted, labelled, and cleared has an enormous potential to disrupt clinical research.

Modern NLP techniques help to process and analyse clinical documentation, extract the required information, and automate much of the work that researchers previously had to do themselves. Some of the techniques that have proved to be especially effective and time-saving are:[5]

- **Named entity recognition** identifies patterns, doctors' names, phones, locations, drug components and other entities and objects that may be of interest. For example, it can locate the most frequently mentioned doctors' names and the attributes of certain specified parameters.
- **Semantic parsing** produces precise meaning representations from unstructured clinical trial data. Broadly speaking, it converts natural language utterances into logical forms. Applied in practice, it helps to classify investigators and patients and label the relationships between them.

- **Topic modelling** helps to conduct topic segmentation and recognition. It allows researchers to automatically define what topics were used and what text segments concern a specific case.
- **Keyword extraction** aids with the extraction of essential information from unstructured articles and publications. It saves considerable time for the professionals conducting the trials.
- **Text summarisation** is employed to analyse clinical trial data and summarise it according to different abstracts or a particular query.
- **Relationship extraction** is a technique that extracts semantic relationships between two or more entities, for example, between article authors, doctors, clinics, diagnosis, medications. Different relationships can be extracted depending on the researcher's goals.

Conclusion

In this article, we have explored some of the key features of NLP in Healthcare, which will help to understand the complex healthcare text data. We also implemented scispaCy and spaCy and constructed a simple custom NER model through a pre-trained NER model and rule-based matcher. While we have only covered one NER model, numerous others are available, and a vast amount of additional functionality to discover. Within the scispaCy framework, there are numerous additional techniques to explore, including methods for detecting abbreviations, performing dependency parsing, and identifying individual sentences. The latest trends in NLP for healthcare include the development of domain-specific models like BioBERT and ClinicalBert and using large language models like GPT-3. These models offer a high level of accuracy and efficiency, but their use also raises concerns about bias, privacy, and control over data.

References

1. Chaudhary N, Weissman D, Whitehead KA. mRNA vaccines for infectious diseases: principles, delivery and clinical translation. *Nat Rev Drug Discov.* 2021;20(11):817-838. doi:10.1038/s41573-021-00283-5
2. Wolfson B, Franks SE, Hodge JW. Stay on target: reengaging cancer vaccines in combination immunotherapy. *Vaccines (Basel).* 2021;9(5):509. doi:10.3390/vaccines9050509
3. Sioutos N, de Coronado S, Haber MW, Hartel FW, Shaiu WL, Wright LW. NCI thesaurus: a semantic model integrating cancer-related clinical and molecular information. *J Biomed Inform.* 2007;40(1):30-43. doi:10.1016/j.jbi.2006.02.013
4. Pardi N, Hogan MJ, Weissman D. Recent advances in mRNA vaccine technology. *Curr Opin Immunol.* 2020;65:14-20. doi:10.1016/j.coi.2020.01.008
5. Aramaki E, Wakamiya S, Yada S, Nakamura Y. Natural language processing: from bedside to everywhere. *Yearb Med Inform.* 2022;31(1):243-253. doi:10.1055/s-0042-1742510