

Article

Not peer-reviewed version

Advancements in NLP for Clinical Data Extraction from Electronic Health Records

[Ahmod Hasan Siddiky](#) *

Posted Date: 4 March 2025

doi: [10.20944/preprints202503.0133.v1](https://doi.org/10.20944/preprints202503.0133.v1)

Keywords: Natural Language Processing; Electronic Health Records; Clinical Data Extraction; Deep Learning; Machine Learning; Rule-Based Systems



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Advancements in NLP for Clinical Data Extraction from Electronic Health Records

Ahmod Hasan Siddiky

Electrical and Electronics Engineering, Northern university Bangladesh, Dhaka, 1230, Bangladesh; nakibhassan304@gmail.com

Abstract: Electronic Health Records (EHRs) store vast amounts of clinical data in both structured and unstructured formats. Extracting meaningful clinical insights from EHRs is essential for improving patient care, enabling precision medicine, and supporting clinical research. Natural Language Processing (NLP) models have shown significant promise in automating the extraction of clinical information from unstructured text in EHRs. This paper reviews recent advancements in NLP techniques, including rule-based methods, machine learning, and deep learning approaches, to enhance clinical data extraction. We discuss challenges such as data heterogeneity, privacy concerns, and model interpretability and explore potential solutions. The paper also highlights emerging trends, such as self-supervised learning and multimodal integration, that are shaping the future of clinical NLP.

Keywords: Natural Language Processing, Electronic Health Records, Clinical Data Extraction, Deep Learning, Machine Learning, Rule-Based Systems

1. Introduction

Electronic Health Records (EHRs) are a cornerstone of modern healthcare systems, serving as comprehensive digital repositories of patient information, including demographics, diagnoses, medications, lab results, and clinical notes. While structured data can be easily queried, the majority of clinically relevant data resides in unstructured formats such as physician notes, radiology reports, and discharge summaries. Extracting actionable insights from these texts requires sophisticated NLP techniques. Recent advancements in NLP, particularly in deep learning, have revolutionized the field, enabling more accurate and efficient extraction of clinical data. This paper reviews these advancements, discusses challenges, and outlines future directions for clinical NLP.

2. Methods for Clinical Data Extraction

2.1. Rule-Based Approaches

Rule-based systems were among the earliest methods for clinical data extraction, relying on lexicons, ontologies, and regular expressions to identify medical terms and concepts. Examples include MedLEE (Medical Language Extraction and Encoding System) and cTAKES (clinical Text Analysis and Knowledge Extraction System) [1]. MedLEE, for instance, uses a combination of syntactic and semantic rules to extract and encode clinical information from free text. While effective for specific tasks, rule-based systems require extensive domain expertise and are limited by their inability to generalize across diverse datasets and documentation styles. They also struggle with scalability and adaptability to new domains or languages [2].

2.2. Machine Learning-Based Approaches

Supervised machine learning techniques, such as Support Vector Machines (SVM), Decision Trees, and Random Forests, have been widely applied to classify clinical text and extract relevant entities. These methods rely on annotated datasets to train models for tasks such as named entity recognition (NER) and relation extraction. The MIMIC-III (Medical Information Mart for Intensive Care III) dataset,

a publicly available critical care database, has been instrumental in advancing these approaches [3]. However, the performance of machine learning models is often constrained by the need for large, high-quality annotated datasets, which are expensive and time-consuming to produce.

2.3. Deep Learning Approaches

Deep learning models, particularly Transformer-based architectures like BERT (Bidirectional Encoder Representations from Transformers) and its clinical variants (e.g., ClinicalBERT, BioBERT), have significantly improved clinical NLP tasks. These models leverage contextual embeddings to enhance entity recognition, relation extraction, and document classification. For example, ClinicalBERT, a version of BERT fine-tuned on clinical notes from the MIMIC-III dataset, has demonstrated superior performance in predicting hospital readmissions and extracting medical concepts [4–8]. Similarly, BioBERT, pretrained on biomedical text, has been adapted for clinical tasks, achieving state-of-the-art results in NER and relation extraction [9]. Fine-tuning these models on domain-specific corpora has proven effective in capturing medical terminology and contextual nuances.

3. Challenges and Considerations

3.1. Data Heterogeneity

EHRs are collected from multiple healthcare institutions with varying documentation styles, terminologies, and formats. This heterogeneity poses significant challenges for NLP models, which often struggle to generalize across datasets. Pretraining models on diverse datasets, such as MIMIC-III and the n2c2 (National NLP Clinical Challenges) corpora, has been shown to improve robustness [10]. Additionally, efforts to standardize clinical terminologies, such as SNOMED CT and LOINC, can help mitigate this issue.

3.2. Privacy and Ethical Considerations

The sensitive nature of clinical data necessitates strict adherence to privacy regulations such as HIPAA (Health Insurance Portability and Accountability Act) and GDPR (General Data Protection Regulation). Federated learning, a decentralized approach that allows models to be trained across multiple institutions without sharing raw data, has emerged as a promising solution [11]. Differential privacy, which adds noise to data to protect individual identities, has also been explored to ensure compliance with privacy regulations while maintaining model performance [12].

3.3. Model Interpretability

The black-box nature of deep learning models poses challenges in clinical decision support, where interpretability is critical for gaining trust from healthcare providers. Techniques such as attention visualization and explainability frameworks like SHAP (SHapley Additive exPlanations) have been employed to improve transparency [13]. For instance, SHAP values can be used to quantify the contribution of individual features to model predictions, enabling clinicians to understand and validate the model's decisions.

4. Future Directions

Recent advancements in self-supervised learning and few-shot learning offer promising avenues for reducing dependence on large annotated datasets. Self-supervised models, such as GPT-3 and its successors, can generate high-quality embeddings from unlabeled text, reducing the need for manual annotation [14]. Few-shot learning techniques, which enable models to generalize from a small number of examples, are particularly valuable in clinical settings where annotated data is scarce. Additionally, integrating multimodal learning approaches that combine textual data with imaging, genomic, and sensor data can further enhance clinical information extraction and enable more comprehensive patient profiling [15].

5. Conclusions

NLP models play a crucial role in extracting clinical data from unstructured EHRs, facilitating enhanced patient care and clinical research. While significant progress has been made, challenges related to data heterogeneity, privacy, and interpretability must be addressed to enable broader adoption in real-world healthcare settings. Future research should focus on improving model adaptability and robustness while ensuring compliance with ethical and regulatory standards. Emerging trends such as self-supervised learning, few-shot learning, and multimodal integration hold great promise for advancing the field of clinical NLP.

References

1. C. Friedman *et al.*, "MedLEE: A medical language extraction and encoding system," *Proceedings of AMIA Annual Fall Symposium*, pp. 19–23, 1995.
2. G. K. Savova *et al.*, "Mayo clinical text analysis and knowledge extraction system (ctakes): Architecture, component evaluation, and applications," *Journal of the American Medical Informatics Association*, vol. 17, no. 5, pp. 507–513, 2010.
3. A. E. Johnson *et al.*, "MIMIC-III, a freely accessible critical care database," *Scientific Data*, vol. 3, no. 1, p. 160035, 2016.
4. K. Huang *et al.*, "ClinicalBERT: Modeling clinical notes and predicting hospital readmission," *arXiv preprint arXiv:1904.05342*, 2019.
5. M. Arifuzzaman, M. J. U. Chowdhury, I. Ahmed, M. N. A. Siddiky, and D. Rashid, "Heart disease prediction through enhanced machine learning and diverse feature selection approaches," in *2024 IEEE 10th International Conference on Smart Instrumentation, Measurement and Applications (ICSIMA)*, 2024, pp. 119–124.
6. S. Sakib, M. N. A. Siddiky, M. Arifuzzaman, and M. J. U. Chowdhury, "Optimizing facial recognition: An analytical comparison of traditional and deep learning approaches," in *2024 International Conference on Data Science and Its Applications (ICoDSA)*, 2024, pp. 271–276.
7. S. D. Tusu, S. Chowdhury, M. J. Uddin Chowdhury, R. M. Pir, N. A. Alam, M. R. Rahman, M. N. A. Siddiky, M. E. Rahman *et al.*, "Advancing chronic kidney disease prediction through machine learning and deep learning with feature analysis," *Frontiers in Health Informatics*, vol. 13, no. 3, 2024.
8. S. M. K. Pathan, S. B. Imran, M. M. S. Iqbal, M. E. Rahman, M. N. A. Siddiky, M. R. Rahman, M. R. Hasan, N. L. Dey, and M. S. Hossain, "Comparative analysis of machine learning models for predicting healthcare traffic: Insights for optimized emergency response," *Magna Scientia Advanced Research and Reviews*, vol. 12, no. 2, pp. 54–61, 2024. [Online]. Available: <https://doi.org/10.30574/msarr.2024.12.2.0175>
9. J. Lee *et al.*, "Biobert: A pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
10. O. Uzuner *et al.*, "The 2010 i2b2/va challenge on concepts, assertions, and relations in clinical text," *Journal of the American Medical Informatics Association*, vol. 18, no. 5, pp. 552–556, 2011.
11. N. Rieke *et al.*, "The future of digital health with federated learning," *NPJ Digital Medicine*, vol. 3, no. 1, pp. 1–7, 2020.
12. M. Abadi *et al.*, "Deep learning with differential privacy," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, 2016, pp. 308–318.
13. S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems*, vol. 30, 2017.
14. T. Brown *et al.*, "Language models are few-shot learners," *Advances in Neural Information Processing Systems*, vol. 33, pp. 1877–1901, 2020.
15. Y. Wang *et al.*, "Multimodal learning for clinical decision support: Challenges and opportunities," *Journal of Biomedical Informatics*, vol. 126, p. 103982, 2022.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.