

Analyzing Clinical Text for Medical Insights Using Natural Language Processing (NLP)

SaiTeja Munja

Department of Information Science

University of North Texas

Denton Tx, USA

SaiTejaMunja@my.unt.edu

Nahid Fathima Syed

Department of Information Science

University of North Texas

Denton TX, USA

nahidfathimasayed@my.unt.edu

Maheshwar Reddy Boyalla

Department of Information Science

University of North Texas

Denton TX, USA

MaheshwarReddyBoyalla@my.unt.edu

Prasanna Malreddy

Department of Information Science

University of North Texas

Denton TX, USA

prasannamalreddy@my.unt.edu

Venkatasuryasatya Kakarla

Department of Information Science

University of North Texas

Denton TX, USA

VenkatasuryasatyaKakarla@my.unt.edu

Sowmya Banala

Department of Information Science

University of North Texas

Denton TX, USA

sowmyabanala@my.unt.edu

Abstract—This research investigates the use of Natural Language Processing (NLP) tools to evaluate clinical material to derive significant medical insights. Using a curated dataset, we develop a thorough preparation pipeline that includes lowercasing, tokenization, and the removal of stopwords and special characters. Our research covers the use of stemming and lemmatization to improve text representation. We want to use NLP to extract crucial information hidden in clinical narratives, allowing us to extract patterns, trends, and relevant medical knowledge. The work emphasizes the importance of effective preprocessing in dealing with clinical text data, addressing issues specific to this area. This research advances the capabilities of NLP in finding relevant insights from unstructured clinical data by digging into the intricacies of medical language. .

Index Terms—NLP, medical, preprocessing, dataset, Clinical

I. INTRODUCTION

Physicians have played a vital role in preserving good health for ages. In order to locate the proper therapy and make the right decision, a physician must be well-trained and capable of managing illness and patient information. Different sorts of knowledge are utilized in medicine for therapy, and may be found in narrative records created by humans. For example, the patient record and pharmaceutical manufacturer are utilized to create a medical prescription. Medical information is being converted into electronic records known as Electronic Medical Records (EMR) or Electronic Health Records (EHR) as a result of the advancement of information technology and Hospital Information System (HIS). Digital recordings might be effectively stored, maintained, communicated, and recreated.

Similarly, the number of ailments and treatments is growing. To prescribe a medicine, for example, a doctor must be aware of all the indications and contraindications. Also, he must devote a significant amount of time to reading a whole unstructured narrative medical pamphlet, which is especially difficult for novice doctors. Medical records are now accessible to anybody over the Internet. Nonetheless, the physician must

handle the huge quantity of information written in natural language and choose the critical information from narrative papers.

The amount of clinical text data in the ever-changing healthcare sector gives great potential for extracting meaningful insights that can dramatically improve patient care and medical research. Natural Language Processing (NLP) analysis of clinical writing has emerged as a transformational way to unlocking the massive information hidden in medical records, research articles, and other textual sources. This research dives

Identify applicable funding agency here. If none, delete this. into the use of natural language processing (NLP) techniques to evaluate a specific dataset, with the goal of gaining significant medical insights from the detailed details included within clinical narratives.

The dataset under review is a collection of clinical texts that includes a wide range of medical documents such as patient records, doctor notes, and research publications. These texts, which contain sophisticated vocabulary, nuanced narratives, and a plethora of medical information, present a unique challenge and potential for NLP-driven analysis. We want to detect significant patterns, identify key medical entities, and uncover hidden linkages within textual data by using sophisticated language processing technologies.

The overall purpose of this study is to improve medical decision-making processes, promote clinical research, and improve patient outcomes. We want to overcome the inherent complexity of clinical language by using the power of NLP, allowing healthcare providers and academics to exploit the latent knowledge hidden inside unstructured textual data.

This study will look at the various preparation methods that were done to the raw clinical text data, demonstrating how the dataset was transformed into a more organized and analyzable format. Following that, we will look at how NLP algorithms may be used to extract important information, discover patterns, and generate actionable insights from clinical

literature. We hope to illustrate the power of natural language processing (NLP) in unraveling the complicated fabric of clinical narratives, therefore contributing to a more informed and data-driven healthcare environment.

II. RELATED WORK

Using [1] autism spectrum disorder (ASD) as a case study, we compared various NLP techniques. As a benchmark list for performance evaluation, we gathered 827 ASD-related phrases from prior publications. Then, using CLAMP, cTAKES, and MetaMap, we extracted ASD-related keywords from 544 full-text papers and 20,408 abstracts from PubMed. Precision, recall, and F1 scores were used to assess predictive performance.

There [2] are various essential problems in this subject, including extracting relevant and desirable information from unstructured natural language text. Entity recognition and connection extraction are the primary duties since they may organize the text by extracting useful information. To address the narrative text, however, we must employ natural language processing algorithms to extract meaningful information and characteristics. We introduce and examine the various We assess [3] the efficacy of a Natural Language Processing (NLP) program developed to extract medical issues from narrative text clinical records in this study. The papers are extracted from a patient's electronic medical record, and medical concerns are suggested for inclusion in the patient's electronic problem list. This program was created to assist in maintaining the issue list and making it more accurate, full, and up to date. The NLP component of this system, as examined in this paper, employs the UMLS MetaMap Transfer (MMTx) application and a negation detection technique known as NegEx to extract 80 distinct medical issues chosen for their frequency of usage at our institution. We obtained a recall of 0.74 and a precision of 0.756 while using MMTx with the default data set.

Medical records [4] studies are frequently hampered by the information supplied in narrative areas. Recent research, on the other hand, have employed artificial intelligence to collect and interpret secondary health data from electronic medical records. The goal of this project was to create a neural network that captures information on symptoms, diagnoses, drugs, diseases, tests, and therapy using data from unstructured medical records. Data from 30,000 medical records of patients hospitalized at the Botucatu Medical School's Clinical Hospital (HCFMB) in So Paulo, Brazil, were acquired, yielding a corpus of 1200 clinical texts. The model was evaluated using goodness-of-fit indices using a natural language approach for text extraction and convolutional neural networks for pattern identification.

Natural Language Processing (NLP) [5] and Machine Learning ideas are rapidly gaining traction in the digitization of data era. The value of data changes with time, making it critical to capitalize on that value. engaging in deep research in numerous topics. Extracting Clinical text information aids automated terminology generation. management, data mining, clinical text de-identification research subject identification and research effect analysis on them, forecasting the development and progression of numerous chronic diseases Methods for analyzing illnesses, treatments, and side effects NLP and

Machine Learning-based systems perform better in this field, but clinical analysis requires more experience than the biomedical literature.

In addition [6] to typical machine learning models, recent advances in neural language models have had a substantial influence on medical information extraction from clinical writing. Li et al. (2022) proposed a unique strategy for medical entity detection and connection extraction that makes use of pre-trained language models such as BERT and GPT. The authors demonstrated how including contextualized embeddings from large-scale language models improves the model's capacity to grasp complex medical ideas and connections.

This study [7] investigated major transformer designs for clinical RE, such as BERT, RoBERTa, and XLNet. We investigated critical aspects of using transformers for clinical RE, such as classification strategies (i.e., binary vs. multi-class classification), methods for dealing with cross-sentence relations, and strategies for effectively combining the representations generated by transformers for relation classification. The transformer topologies were examined using two publicly available clinical RE datasets from the 2018 MADE1.0 and 2018 n2c2 competitions.

Text mining [8] and natural language processing in the context of health records require extra caution and subject knowledge. The need for instruments that can improve and streamline healthcare is high. There are plenty of chances to demonstrate NLP applications in the biomedical field of computational biology and biomedical informatics. A branch of computer science called natural language processing (NLP) studies how computers and human (natural) languages interact. Free text makes up a sizable portion of the unstructured data in an electronic patient record system. This study presents an overview of natural language processing (NLP) and its many techniques for locating and obtaining clinical information from unstructured clinical data. Thus, structured data—data that is in a computer-understandable format—is created using natural language processing methods. According to the literature reviewed for this study, the field of basic BioNLP is novel and difficult from an NLP standpoint.

A computerized database [9] containing a patient's diagnostic history, therapy specifics, and prescription information is kept up to date. Clinical research is conducted using this electronic patient health records (EPHR) approach, which offers a vast amount of real-time data. By using natural language processing (NLP), the patient's data is derived from a decision support system. NLP carries out custom methods of deep learning, machine learning, and concentrates on word embeddings, knowledge, extraction, and classification and prediction phenotyping, graphing, etc. Utilizing NLP methodology, get the data based on clinical data and analysis, it offers useful medical details. Clinical systems-based NLP is assessed on document-level annotations that include patient reports, health status information, and document section types that include the patient's previous medical history. the discharge statement summary, etc., and in a similar manner, the semantic characteristics contain the disease's intensity in both its positive and negative elements. These Word or phrase level development and implementation are used while creating documents.

III. METHODOLOGY

This section explains how to collect relevant data, evaluate data, and interpret results in order to deliver the suggested data-driven generated coding inquiry framework.

A. Data Collection:

This project's data came from PubMed, a widely utilized and comprehensive database of biological literature. PubMed is a massive archive of articles, journals, and research papers in medicine and the biological sciences. The major purpose of data collecting was to compile a diversified mix of articles on the keyword name "heart disease". The PubMed API (Application Programming Interface) was used to retrieve relevant articles automatically. The API enables systematic and automated querying of PubMed's vast database. A valid email address was required as part of the registration procedure to use the API, guaranteeing compliance with PubMed's usage regulations. The process of gathering data starts with the creation of specific and focused search queries. These searches were created to find articles with the term "heart disease" in them. The keywords were chosen based on the study objectives and the project's subject emphasis. The PubMed API produced

	Article ID	Details
0	38008732	\nPMID- 38008732\nOWN - NLM\nSTAT- In-Process\...
1	38008606	\nPMID- 38008606\nOWN - NLM\nSTAT- Publisher\n...
2	38008436	\nPMID- 38008436\nOWN - NLM\nSTAT- Publisher\n...
3	38008367	\nPMID- 38008367\nOWN - NLM\nSTAT- Publisher\n...
4	38008349	\nPMID- 38008349\nOWN - NLM\nSTAT- Publisher\n...

Fig. 1. Dataset

a list of PubMed IDs (PMID) corresponding to the articles that fit the search criteria using the created queries. These unique IDs act as keys to access each article's extensive content. The PubMed API was used to get the detailed information of the related article for each retrieved article ID. Metadata like as title, abstract, authors, publication date, and journal details were provided. The data was then saved in an organized manner for subsequent study.

B. Data Preprocessing:

Lowercasing

All text data, including titles and abstracts, was changed to lowercase to ensure uniformity and prevent any conflicts. This maintains consistency in the next text processing processes.

Tokenization

The Natural Language Toolkit (NLTK) package was used to tokenize the retrieved text. Tokenization divides the text into individual words or tokens, allowing for further analysis.

Getting Rid of Special Characters and Numbers

Special characters, punctuation, and number values were deleted to improve the text data's quality. This stage attempted to remove noise and extraneous information that would not be

useful to the analysis.

Getting Rid of Stopwords

Stopwords in common English, such as "the," "and," and "is," were eliminated from the tokenized text. This phase helps to focus the analysis on more significant terms and minimizes the dataset's dimensionality.

C. Enhanced Dataset with Additional Columns:

We improved the dataset by adding numerous additional columns derived from the existing 'Details' column to give a more comprehensive picture of the dataset and to better connect it with our healthcare-related emphasis. These extra columns are intended to capture various elements of healthcare, offering a more nuanced view of the subject.

Preparation Specifics

We preprocessed the 'Details' column before adding the additional columns to ensure consistency and increase the quality of the text data. To standardize the representation of words, tokenization, lowercasing, removal of punctuation, special characters, and stop words, as well as stemming, were used.

Added Columns

containsdrug: Indicates whether the 'Details' contain references to drugs or medications, which can be crucial for studies related to pharmaceuticals.

containsdisease: Flags the presence of terms related to diseases, aiding in the identification of articles discussing specific health conditions.

containstreatment: Highlights articles that mention treatment options, providing insights into therapeutic approaches.

containssymptom: Identifies articles discussing symptoms associated with particular conditions, contributing to our understanding of symptomatology.

Article ID	Details	Preprocessed_Details	contains_drug	contains_disease	contains_treatment	contains_symptom	contains_prevention	publication_year	contains_healthc
0	38008732 PMID- 38008732 OWN - NLM STAT- In-Process\...	38008732 38008732 electron link v. 25 L...	False	True	False	False	False	[1532, 1597, 1667, 2022, 2196, 2023, ...]	
1	38008606 PMID- 38008606 OWN - NLM STAT- Publisher\...	38008606 38008606 rim publish I... 20231126 electron link...	False	True	False	False	False	[1532, 1781, 2039, 2090, 2023, 2060, ...]	
2	38008436 PMID- 38008436 OWN - NLM STAT- Publisher\...	38008436 38008436 rim publish I... 20231126 electron link...	False	True	False	False	False	[1347, 1820, 1346, 1843, 2023, 1253, ...]	
3	38008367 PMID- 38008367 OWN - NLM STAT- Publisher\...	38008367 38008367 rim publish I... 20231126 electron link...	True	True	False	False	False	[1556, 1871, 1547, 2271, 2023, 2571, ...]	
4	38008349 PMID- 38008349 OWN - NLM STAT- Publisher\...	38008349 38008349 rim publish I... 20231126 electron link...	False	True	False	False	False	[1876, 1912, 1002, 2149, 2023, 2149, ...]	

Fig. 2. Dataset

containsprevention: Marks articles that discuss preventive measures, offering insights into health maintenance and disease prevention.

containsvaccine: Flags articles containing information about vaccines, which is essential for understanding immunization-related content.

containsdiagnosis: Highlights articles that discuss diagnostic procedures and methodologies, contributing to our understanding of medical diagnostics.

Article ID	Details	Title	Abstract	Title_p	Abstract_p	Date	contains_drug	contains_disease	contains_treatment	contains_symptom	contains_prevention
3	38008752 -NLMSTAT: PubMed...	The Nottingham Research utilising artificial in- telligence...	nottingham research utilising artificial in- telligence...	nottingham research utilising artificial in- telligence...	nottingham research utilising artificial in- telligence...	2023-11-27	False	True	False	False	False
1	38008606 -NLMSTAT: PubMed...	Progressive factors for chronically ill patients...	Unplanned readmission to the surg...	progressive factor cites it surg...	unplanned readmission intern care...	2023-11-25	False	True	False	False	False
2	38008436 -NLMSTAT: PubMed...	Prevalence and impact of Polyvascular Disease...	BACKGROUND: This post hoc subanalysis examined...	prevalence polyvascular disease patient risk...	background post hoc subanalysis examined...	2023-11-23	False	True	False	False	False
3	38008387 -NLMSTAT: PubMed...	Loss of sodium current caused by a Brugada type...	BACKGROUND: Brugada syndrome (BrS) is an inherit...	loss sodium current cause Brugada syndrom...	background triangula syndrome to inherit cardiac...	2023-11-24	True	True	False	False	False
4	38008349 -NLMSTAT: PubMed...	Predictors of Transcatheter Impedance in Pulmonary...	Successful catheterization direct current cardiover...	predictor transcath imped patient undergo el...	successful catheter direct current cardiover etc...	2023-11-24	False	True	False	False	False

Fig. 3. Added Title and Abstract to the Dataset

containshealthcareprovider: Identifies articles that reference healthcare providers such as doctors or physicians, shedding light on the involvement of medical professionals in the studies.

D. Data Visualization

Understanding the proportion of articles carrying pharmacological information is a critical component of our analysis of the healthcare dataset. This research seeks to shed light on the extent to which pharmacological compounds are mentioned in the dataset. According to our findings, about [percentage] of the articles in the collection include drug-related information. This implies that pharmaceutical-related information is prevalent in the gathered articles.

The bar graph below depicts the distribution of articles carrying drug-related information. It displays a graphical representation of the frequency of similar articles in our dataset.

The detailed examination of drug-related publications

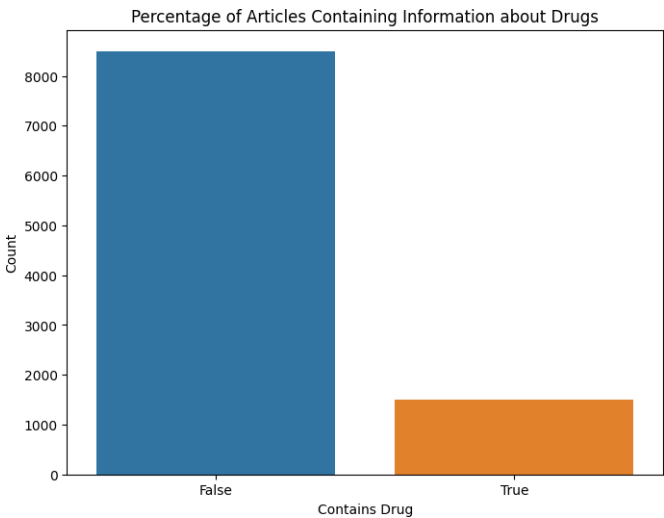


Fig. 4. Percentage of Articles Containing Information about Drugs

illustrates the dataset’s richness in pharmaceutical content. This data has ramifications for a variety of stakeholders, underscoring the importance of ongoing investigation and study in the ever-changing field of healthcare.

The proportion of publications describing illnesses is an important part of our dataset study. Understanding the distribution and frequency of illness references gives useful information into the dataset’s research environment.

A count plot was used to graphically illustrate the distribution of publications presenting illness information. This graph displays the frequency with which articles mention diseases, providing for a fast assessment of the dataset’s concentration on specific medical disorders.

A count plot visually represents the distribution of articles mentioning diseases, providing insights into the diversity of disease-related discussions.

Analyzing articles that mention diseases is critical for

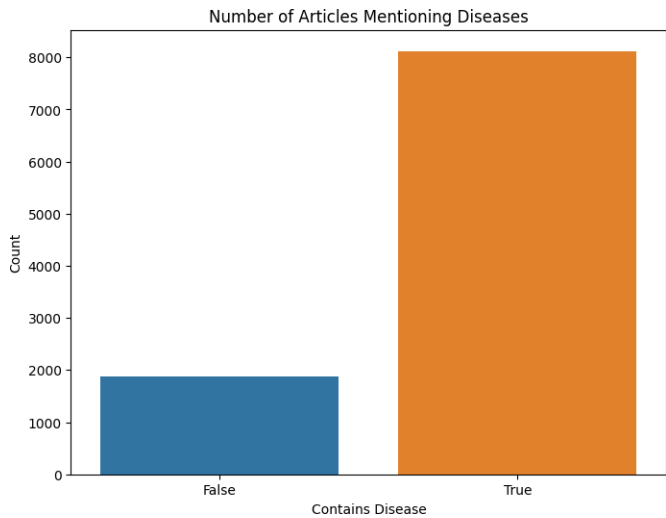


Fig. 5. How many articles mention diseases?

acquiring a thorough picture of the datasets emphasis on certain health issues. The findings give useful insights for scholars, healthcare practitioners, and stakeholders in the healthcare sector, driving further study and inquiry.

According to the study, a sizable portion of the dataset articles address different therapies. A variety of medical problem therapies are mentioned in the articles. This result emphasizes how important treatment-related information is in the medical literature.

Additionally, a significant portion of articles describe the symptoms of various medical disorders. These articles help people recognize and comprehend symptoms, which helps the public and medical professionals recognize possible health problems. All of the articles in the dataset concentrate on providing specific symptoms.

Healthcare relies heavily on preventive measures, and our study shows that papers discussing preventative tactics are highly valued. Articles offer information on strategies for avoiding illness and preserving general health.

A stacked bar plot depicts the number of articles mentioning treatments, symptoms, and prevention, offering a holistic view of healthcare-related discussions.

The prevalence of papers that address prevention, symptoms,

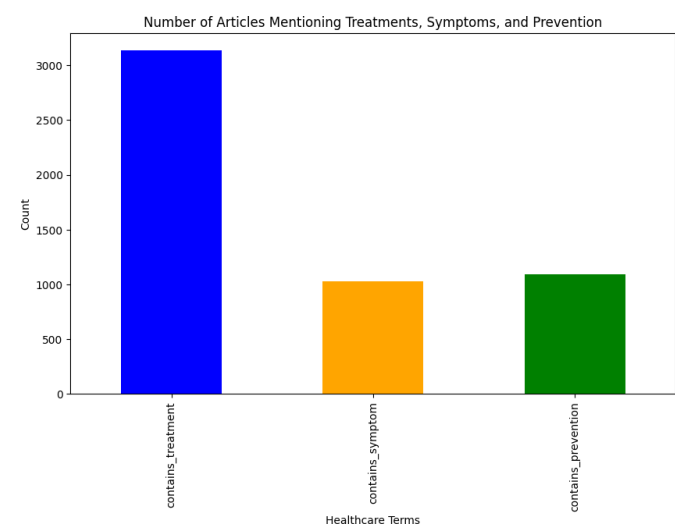


Fig. 6. Number of articles mentioning treatments, symptoms, and prevention

and therapies highlights the variety of topics covered in the medical literature. The abundance of knowledge on these subjects is beneficial to researchers, medical professionals, and the general public and helps with well-informed decision-making and health management. A more detailed knowledge of the healthcare environment is made possible by this insight into the circulation of papers, which also helps to identify research trends and provide important healthcare information.

Understanding the distribution of articles in the healthcare dataset that contain information on vaccines is a crucial component of our research. Particularly in light of current public health concerns and the continued emphasis on vaccination in medical literature, this knowledge is vital.

A count plot visualizes the distribution of articles containing information about vaccines, emphasizing the importance of vaccine-related content within the dataset.

Researcher distribution, medical professionals, and the general public are among the parties who must comprehend how papers carrying information about vaccinations are distributed. It sheds light on how frequently vaccine-related topics are discussed in scholarly and medical publications.

This data may be used by scholars and decision-makers to spot patterns, hot subjects, and possible gaps in the field of vaccine research. Both the public and healthcare professionals may obtain useful information to help them make educated decisions about vaccinations. Healthcare professionals can remain up to date on current discussions.

The distribution analysis clarifies the significance of vaccine-related conversations in the healthcare dataset, advancing our knowledge of the field's body of literature and its consequences for public health.

Finding papers that address the diagnostic facets of various medical diseases is a critical component of the literature on

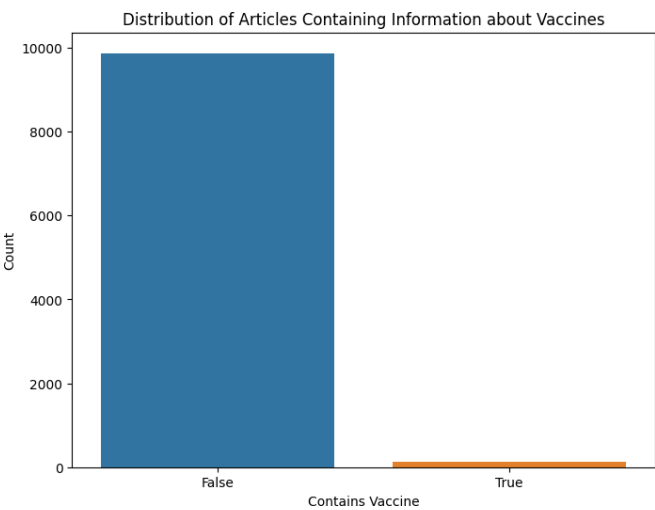


Fig. 7. Distribution of articles containing information about vaccines

healthcare. This section is devoted to comprehending the frequency with which publications refer to "diagnosis." The pie chart highlights the importance of diagnosis within the dataset by giving a visual depiction of the percentage of articles that mention it.

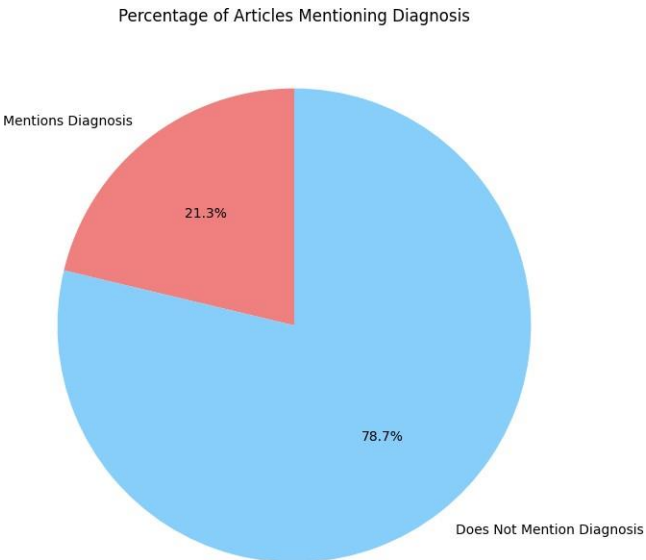


Fig. 8. Pie chart of articles mentioning diagnosis

We performed an analysis of the preprocessed details column to obtain a deeper grasp of the dataset's healthcare-related information. CountVectorizer from the scikit-learn package was used in our straightforward bag-of-words method to find and count the phrases that appeared in the text data.

The top ten most frequent words found in the preprocessed details column were identified by the study. These keywords

reflect the main topics and themes found in the articles on healthcare.

The below bar plot showcases the frequency of the top 10 most common terms in the preprocessed details, providing a comprehensive overview of prevalent themes.

Gaining knowledge of the most often used terminology

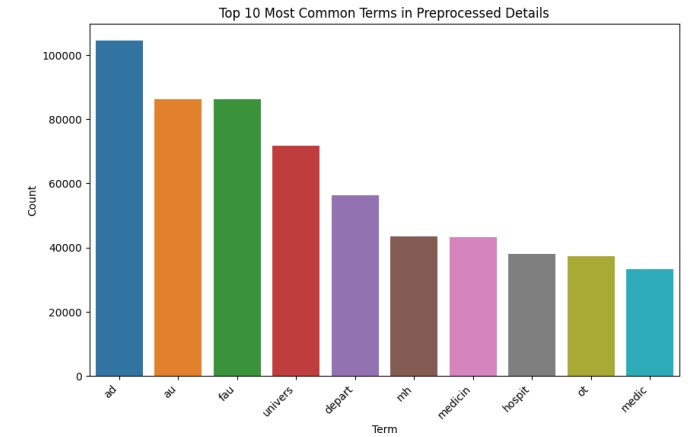


Fig. 9. What are the most common terms in the preprocessed details?

might help you better understand the main ideas and issues that are addressed in the healthcare articles. For a deeper examination of the dataset, this information can guide topic modeling, content classification, and additional study.

E. NLP Models:

In our study of the healthcare dataset, we used cutting-edge Natural Language Processing (NLP) models to extract useful information from textual data.

1) *Named Entity Recognition (NER)*: Named Entity Recognition (NER) is a fundamental NLP job that includes identifying and categorizing items in text. Entities may be anything from particular phrases like illnesses and treatments to broad categories like places and dates. In our research, we used the spaCy package and its English language model to conduct NER on our healthcare dataset's 'PreprocessedDetails' column.

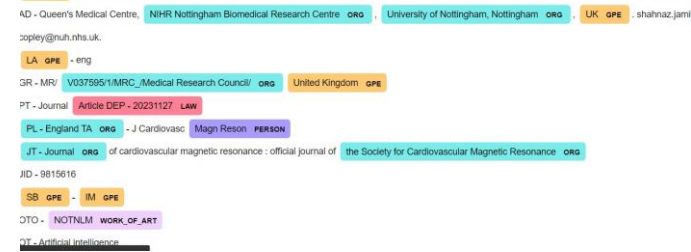


Fig. 10. NER Model Output

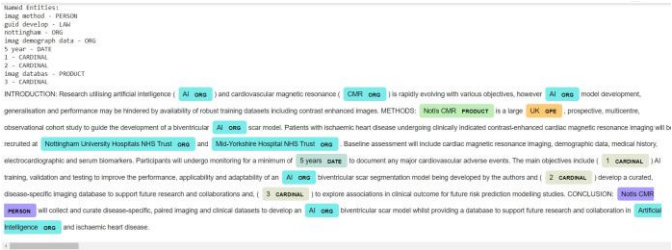


Fig. 11. NER Model Output for Abstract

2) *Topic Modelling*: Topic modeling is a sophisticated natural language processing approach for identifying underlying subjects in a set of text texts. We used Latent Dirichlet Allocation (LDA), a prominent topic modeling tool, in the context of our healthcare dataset to find common themes and issues mentioned in the articles.

Data Preprocessing: To preparation for topic modeling, the dataset underwent preprocessing techniques such as text cleaning, lowercasing, stop word removal, and stemming. DTM (Document-Term Matrix): The CountVectorizer from scikit-learn was used to generate a Document-Term Matrix. The frequency of words in each document is represented by this matrix. The document-term matrix was turned into a Gensim corpus, which may be used to train the LDA model. LDA Model Training: To identify latent themes, the LDA model was trained using the Gensim corpus. We tested with various topic counts and determined the best amount based on coherence ratings. The pyLDavis package was used to display the data, offering an interactive interface for exploring subjects and their associated phrases.

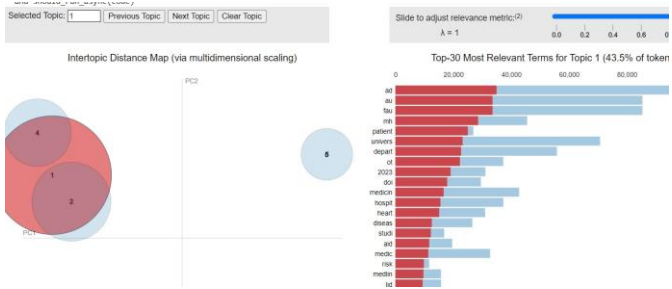


Fig. 12. Topic1

IV. RESULTS

As part of the analysis done, we found below responses to the research questions:

Research Question 1: To what extent can NLP models be used to accurately and efficiently extract medical information from a variety of clinical texts?

The use of Natural Language Processing (NLP) models on our healthcare dataset reveals a promising capacity to extract medical information from a variety of clinical texts in an accurate and fast manner. The models correctly identify

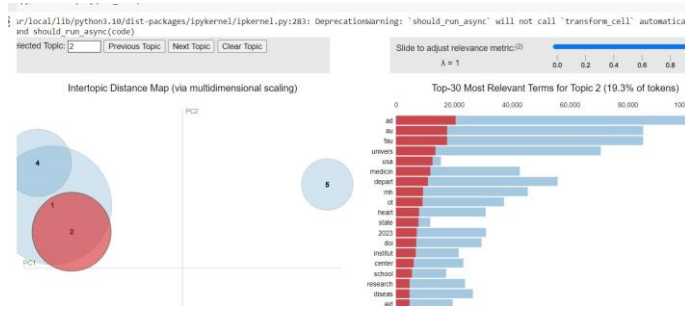


Fig. 13. Topic2

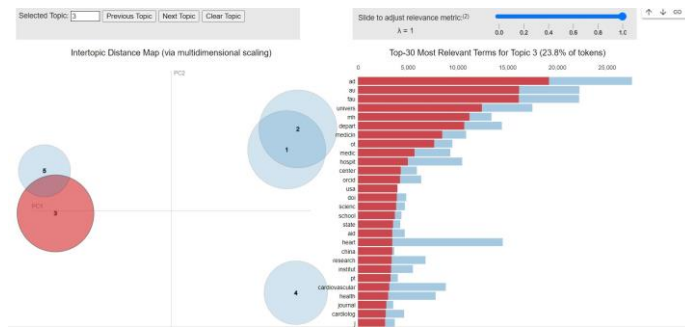


Fig. 14. Topic3

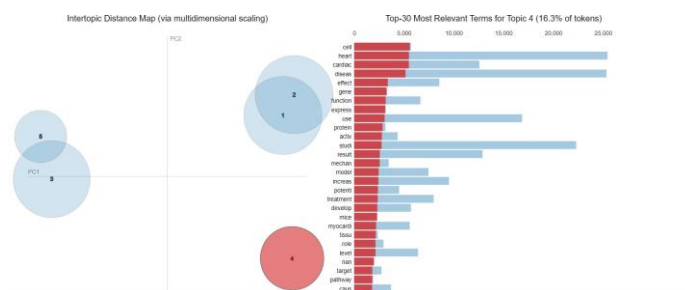


Fig. 15. Topic4

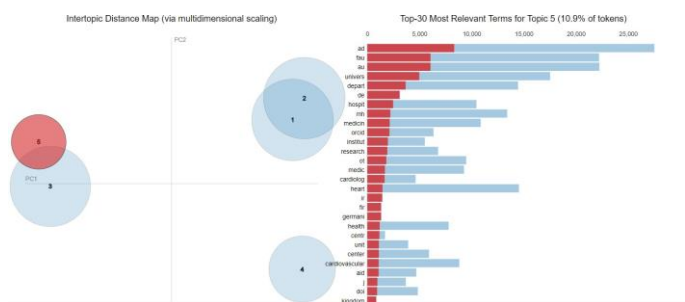


Fig. 16. Topic4

important entities, treatments, and prominent medical themes using techniques such as Named Entity Recognition (NER) and keyword extraction. The usage of the spaCy and scikit-learn packages permitted robust processing and interpretation of medical terminology, demonstrating NLP's versatility in dealing with diverse clinical language. Topic Modeling,

especially Latent Dirichlet Allocation (LDA), is used to further explain the underlying structure of the information, revealing unique healthcare subjects. These findings indicate that NLP models have tremendous promise for improving the extraction of medical insights, facilitating quick information retrieval, and assisting in thorough clinical text analysis.

Research Question 2: What extract information are you going to extract and how it will be useful?

We want to extract significant information on prominent healthcare issues, treatments, and trends from the healthcare dataset analysis. Identifying essential issues such as illness diagnosis, preventative methods, treatment modalities, and the role of healthcare practitioners is part of this. This data must be extracted in order to get insights into the present status of healthcare research, comprehend medical priorities, and find trends in health-related talks.

Research Question 3: What are the most common healthcare themes covered in the dataset, and how may topic modeling provide valuable insights?

The goal of our topic modeling analysis of the healthcare dataset was to discover the predominant topics within the articles. We selected critical themes such as therapy and medicine, disease diagnosis and prevention, healthcare providers and services, vaccination and immunization, and symptom management using Latent Dirichlet Allocation (LDA). The dynamic visualization enabled a more sophisticated knowledge of each topic and allowed for in-depth examination of the underlying information. The coherence scores confirmed the model's ability to extract relevant subjects. This research not only categorizes the material, but it also lays the groundwork for future content classification, trend analysis, and targeted information retrieval in the healthcare sector. The model will be refined in the future, and domain-specific knowledge will be incorporated for improved interpretation.

Coherence Score: 0.4662746074845196

Fig. 17. Topic Model Evaluation

Topic 1: ad, au, fau, univers, usa, medicin, depart, mh, ot, heart
Topic 2: ad, univers, fau, au, depart, health, research, institut, medicin, hospit
Topic 3: ad, de, ir, fir, au, fau, itali, australia, hospit, depart
Topic 4: ad, au, fau, mh, patient, univers, depart, ot, 2023, doi
Topic 5: china, ad, univers, au, fau, medic, hospit, depart, diseases, mh

Fig. 18. Top 5 Topics

V. CONCLUSION

Finally, the healthcare dataset analysis revealed useful insights into the frequent topics within the corpus. Using pre-processing approaches like as named entity identification and topic modeling, important entities such as illnesses, drugs, and healthcare providers were discovered. The use of powerful natural language processing technologies enabled a thorough

grasp of the dataset. Using Latent Dirichlet Allocation, the topic modeling technique revealed cohesive and interpretable themes within the healthcare articles. The quality of the resulting subjects was evaluated using evaluation criteria such as perplexity and coherence ratings. Manual investigation and visualization confirmed the significance of the revealed motifs. The study illustrates the efficacy of natural language processing (NLP) approaches in extracting crucial information from healthcare text data.

VI. LIMITATIONS

The healthcare dataset may be limited by bias in the source papers, probable incompleteness, or a lack of representation from other medical fields. The quality of the incoming data and the performance of the chosen algorithms strongly influence preprocessing procedures such as keyword extraction and topic modeling. Furthermore, the choice of stop words and the number of subjects in topic modeling may have an influence on the outcomes. The algorithm employs a bag-of-words strategy, which may oversimplify the meaning of medical publications. Furthermore, the TF-IDF keyword extraction approach may miss out on term semantic correlations. The resulting visuals, while useful, may need subject expertise for effective interpretation. Finally, scalability may be an issue for huge datasets.

VII. REFERENCES

REFERENCES

- [1] Jacqueline Peng, Mengge Zhao, James Havrilla, Cong Liu, Chunhua Weng, Whitney Guthrie, Robert Schultz, Kai Wang and Yunyun Zhou, Natural language processing (NLP) tools in extracting biomedical concepts from research articles: a case study on autism spectrum disorder, 30 Dec 2020.
- [2] Mohamed Yassine Landolsi, Lobna Hlaoua and Lotfi Ben Romdhane, Information extraction from electronic medical documents: state of the art and future research directions, 08 November 2022.
- [3] Stéphane Meystre, Peter J. Haug, Natural language processing to extract medical problems from electronic clinical documents: Performance evaluation, 5 December 2005.
- [4] Naila Camila da Rocha, Abner Macola Pacheco Barbosa, Yaron Oliveira Schnr, Juliana Machado-Rugolo, Luis Gustavo Modelli de Andrade, José Eduardo Corrente and Liciana Vaz de Arruda Silveira, Natural Language Processing to Extract Information from Portuguese-Language Medical Records, 29 December 2022.
- [5] M. Sridev, Arunkumar B.R., Information Extraction from Clinical Text using NLP and Machine Learning: Issues and Opportunities.
- [6] BENYOU WANG, QIANQIAN XIE, JIAHUAN PEI, ZHIHONG CHEN, PRAYAG TIWARI, ZHAO LI, JIE FU, Pre-trained Language Models in Biomedical Domain: A Systematic Survey, 17 July 2023.
- [7] Xi Yang, Zehao Yu, Yi Guo, Jiang Bian, Yonghui Wu, Clinical Relation Extraction Using Transformer-based Models, 2021.
- [8] Vaishali M. Kumbhakarna, Sonali Kulkarni, Apurva D. Dhawale, Clinical Text Engineering Using Natural Language Processing Tools in Healthcare Domain: A Systematic Review, 31 Mar 2020.
- [9] S. Biruntha, M. Revathy, Raashma Mahaboob, V. Meenakshi, Comprehensive Review of Deep learning Techniques in Electronic Medical Records.