**1) In the article "Challenges and opportunities beyond structured data in the analysis of electronic health records," the two most popular sources of information in an unstructured data format are:**

**Clinical Text:** This refers to narrative information that is kept in electronic health records (EHRs), such as clinical notes, surgical records, discharge summaries, radiology reports, and pathology reports. Important details about patient care, medical procedures, and diagnosis can be found in clinical texts. Because clinical literature is unstructured, the context might vary, and there are grammar and spelling errors, ambiguities, and abbreviations present, it can be difficult to analyze it.

**Medical Images:** Medical images are a key source of unstructured data in EHRs. X-rays, MRIs, CT scans, and other sorts of medical pictures are included in this category. Analysis of medical photographs is crucial for diagnosis and therapy planning since they offer vital visual information about a patient's condition. Though difficult and requiring specialized algorithms and methods, deciphering and extracting information from medical images is a major area of study for healthcare informatics.

**2) The article's main ideas and justifications for why clinical data is frequently restricted are as follows:**

**Data Sensitivity and Ethical Issues:** The sensitivity of the data in EHRs is one of the main justifications for restricting clinical data. These records contain comprehensive and frequently private information regarding the diseases, therapies, and medical histories of specific patients.

**Data security:** Clinical data is a valuable target for cyberattacks and data breaches, according to data security experts. Healthcare organizations need to take strict security precautions to shield patient data from theft, unauthorized access, and data tampering.

**Informed Consent:** Patients must give their informed consent before their data may be utilized for research, according to ethical norms. Patients have the right to be informed about how their information will be used and to decide whether it should be shared with researchers.

**Legal and Regulatory Constraints**; The use and sharing of healthcare data is regulated by a number of laws and regulations, including the Health Insurance Portability and Accountability Act (HIPAA) in the US.

**Institutional Policies:** Healthcare facilities and research organizations may have their own rules and regulations regarding data access and sharing. These regulations frequently serve to safeguard patient interests, uphold confidence, and guarantee data security.

I feel overall, the justifications for limiting healthcare data are usually seen as legitimate and essential for safeguarding patient privacy, ensuring data security, and abiding by ethical and legal requirements. Finding a balance between data privacy and promoting research and innovation, however, is a subject of constant discussion in the medical and scientific fields. In an effort to solve some of these issues while still protecting patient interests, efforts are being made to provide safe and privacy-preserving techniques for data sharing and analysis.

**3) These methods are intended to take useful information from unstructured data and turn it into something that can be analyzed. Listed below are a few of the methods from the article:**

**Natural Language Processing (NLP):** A branch of artificial intelligence called "natural language processing" (NLP) is concerned with how computers and human language interact. Clinical text data in EHRs is processed and analyzed using NLP algorithms.

**Deep learning and machine learning:** To analyze unstructured data, such as clinical writing and photos, machine learning methods, including deep learning, are used. Convolutional neural networks (CNNs) for images and recurrent neural networks (RNNs) for text are two examples of deep learning models that have been used to extract features and patterns from unstructured data.

**Radiomics:** The study of radiographic images, such as X-rays, CT scans, and MRIs, with a focus on their quantitative aspects is known as radiomics. These characteristics can offer useful information for planning a diagnosis and a course of treatment. The extraction of numerous texture, shape, and intensity-based elements from images is a key component of radionic analysis.

**Text Mining:** Finding patterns and insights in clinical text data is done using text mining tools. This could entail looking for patterns in clinical notes, discovering connections between symptoms and diagnoses, or organizing unstructured data into organized categories for simpler analysis.

**Synthetic Data Generation:** The article recommends the development of machine learning techniques that can produce clinically pertinent synthetic data in order to address issues with data quality and accessibility.

**Techniques for Protecting Privacy**: Because clinical data is sensitive, privacy-protecting strategies including deidentification and pseudonymization of clinical language are suggested as potential remedies. These techniques strive to preserve the research and analytical value of unstructured material while removing or encrypting personally identifiable information from it.

**Data Quality Improvement:** The article recognizes the difficulty of improving data quality. To increase the dependability of unstructured data, techniques for enhancing the completeness, conformance, and plausibility of the data are required. Procedures for data standardization, validation, and cleansing may be required.

**4) Here are some common scenarios where data types may be mistaken for each other:**

**Structured vs. Semi-Structured Data**: Typically, structured data is arranged into rows and columns according to a preset schema, which facilitates searching and analysis. However, unlike relational databases, semi-structured data lacks the strict structure that characterizes structured data. Because semi-structured data may occasionally contain metadata or tags that offer some level of organization, semi-structured data may occasionally resemble structured data.

**Semi-Structured vs. Unstructured Data:** Due to the fluidity of its structure, semi-structured data can occasionally be confused with unstructured data. Semi-structured data lacks the rigid structure of structured data, yet it may include some organization or tags..

**Unstructured Data with Metadata:** Unstructured data with associated information or tags that offer context or a certain amount of organization, such as text documents, is referred to as unstructured data. An example of metadata might be the author, creation date, or keywords of a text file

**Data that is completely unclassified or categorized can be encountered in a variety of contexts:**

**Raw sensor data**: Data produced by sensors, such as Internet of Things (IoT) devices or scientific instruments, may not be naturally categorized or structured. To be useful, this raw sensor data frequently needs to be processed and contextualized. Web scraping: When gathering information from the internet via web scraping, the information could be presented in an unstructured way without any established categories.

**User-Generated Content**: User-generated content frequently lacks a standardized structure or classification on social media platforms, forums, or comment sections. It displays the variety and impromptu nature of user contributions.

**Audio and video streams**: Continuous data streams without intrinsic classification may be produced by live audio and video streams, such as those from security cameras or webcams. These streams need to be interpreted and examined.

**Free-Form Text:** Textual information might be completely unstructured and uncategorized, such as handwritten notes or personal journals. Manual annotation and interpretation are frequently needed.

## 5) Improving Data Quality and Accessibility of Unstructured Data:

**Synthetic Data Produced by Machine Learning:** The article makes the case that machine learning techniques can be used to produce synthetic data that closely resembles genuine clinical data while yet safeguarding patient privacy. Without disclosing private patient information, research and analysis can be done using this artificial data. This method aids in addressing privacy and data quality issue.

**Privacy-Preserving Techniques:** Deidentification and pseudonymization are privacy-preserving procedures that are used to replace or obliterate personally identifiable information (PII) in clinical text and other unstructured data. Without compromising patient privacy, deidentified data can be shared more broadly for research purposes.

**Enhancing Data Quality:** The essay emphasizes the value of raising the completeness, conformance, and believability of data. Unstructured data can be improved in quality and made more analytically accessible by using techniques to standardize and normalize it.

When it comes to unstructured data sources that could be used in market analysis, patient-generated data sentiment analysis may be particularly intriguing. To ascertain public mood and opinions regarding healthcare providers, treatments, drugs, and medical devices, may entail mining social media posts, patient forums, and online reviews. In order to make wise judgments and spot market trends, healthcare organizations, pharmaceutical firms, and other stakeholders in the healthcare sector may benefit greatly from the analysis of this unstructured data. In order to better design healthcare services and goods, it can be helpful to understand the satisfaction, issues, and preferences of patients from unstructured sources. When using patient data for analysis, it's imperative to treat it carefully and respect privacy laws.

**6).** Because structured data is arranged in rows and columns, searching and analyzing it is simple. It conforms to a data schema and can be mapped into predefined fields. It is often controlled using Structured Query Language (SQL), and examples include Excel spreadsheets and **structured data**.

**Traditional data** are unable to store unstructured data because they lack a set structure. Examples include text in emails, images, movies, and posts on social media, among others. Unstructured data was difficult to manage and analyze before the emergence of artificial intelligence and machine learning.

Between structured and unstructured data is **semi-structured data**. Although it doesn't have a strict structure like relational databases, it does have some constant features. Examples are emails, which comprise both structured and unstructured data such as timestamps, sender information, and metadata.

Consequently, these are the distinctions between structured, unstructured, and semi-structured data.

**References & Citation:**

-Tayefi, M., Ngo, P., Chomutare, T., Dalianis, H., Salvi, E., Budrionis, A., & Godtliebsen, F. (2021). Challenges and opportunities beyond structured data in analysis of electronic health.. *Wiley Interdisciplinary Reviews: Computational Statistics*, *13*(6), e1549.

-Marr, B. (October 2019). What's the Difference between structured, semi-structured, and unstructured data.. Retrieved June 3, 2022.

-https://blog.hubspot.com/marketing/semi-structured-data