

Data Analysis Report

Topic: Cybersecurity Events on Healthcare Sector

Prepared for:
HEAL Security: Junior Data Analyst

Prepared by:
Aida Martirosyan

December 3, 2023

Contents

1	Introduction	2
2	Data Summary	3
2.1	Dataset Overview	3
2.2	Basic Statistics	3
2.3	Data Preprocessing	3
2.4	Data Columns Descriptions	3
3	Insights and observations	5
3.1	Visuals and Interpretations	5
3.2	Interesting Findings	7
4	Challenges for integration	8
5	Recommendations	10
6	Conclusion	11

The consequences of cybersecurity events in the healthcare sector extend beyond mere data breaches. They can disrupt essential healthcare services, compromise patient care, and lead to financial losses for both healthcare providers and individuals. The interconnectedness of healthcare systems also raises concerns about the potential for widespread impacts, making it imperative to understand, analyze, and address these cybersecurity challenges.

In the following sections, the report includes detailed data analysis, summarizing key attributes, exploring noteworthy patterns, and addressing challenges related to data integration. The findings from this analysis will inform recommendations aimed at enhancing cybersecurity measures and fortifying the healthcare sector against potential threats.

2 Data Summary

2.1 Dataset Overview

The dataset consists of cybersecurity events, capturing a diverse range of incidents affecting covered entities and business associates. Each entry in the dataset represents a unique case, detailing the nature of the breach, the entities involved, and the impact on individuals' protected health information (PHI).

Let us embark on this data-driven journey to unravel insights that will contribute to the ongoing efforts to bolster cybersecurity defenses in the healthcare domain.

2.2 Basic Statistics

Dataset Size:

- The data has 1000 entries (rows).
- There are 10 columns in total.

Data Types:

- The data types of the columns are as follows:
 - **object**: 9 columns are string or categorical data. However, note that **breach_submission_date** should be changed into date.
 - **int64**: 1 column is numeric and includes data in a 64-bit integer type.

Missing Values:

- Columns with missing values (NaN or null) are visible. For instance, **"web_description"** has a relatively low count of non-null entries (4), suggesting that there are many missing values (996) in this column. There are 4 missing values in **"state"** column and just 1 in **"covered_entity_type"**. Overall problematic one is **"web_description"** column as nearly all values are null.

2.3 Data Preprocessing

Data preprocessing is a crucial step before analysis to ensure the accuracy and reliability of insights. In this process, I converted the **"breach_submission_date"** from a string to a datetime format to facilitate temporal analyses. Addressing missing values, the **"state"** and **"covered_entity_type"** columns, with very few missing entries, were imputed using the mode for enhanced completeness. However, in the case of the **"web_description"** column, after careful examination and considering that only 4 descriptions are present out of 1000 entries, I decided to drop this column. The decision is based on the observation that the essential information from **"web_description"** is already captured in other columns, such as **"individuals_affected"** and **"name_of_covered_entity"**, however I am using its information for defining the data description.

2.4 Data Columns Descriptions

1. status

Indicates whether the record is "current" or "archive." In the context of a data breach dataset, "Current" status may indicate active or recent security incidents that are still being addressed or investigated. "Archive" status may represent historical security incidents that have been resolved, closed, or are no longer considered active. In this column there 2 unique values.

2. name_of_covered_entity

The name of the covered entity involved in the security incident. Typically, it refers to the name of the organization or entity that experienced the data breach. In this column there are 917 unique values.

3. individuals_affected

The number of individuals affected by the security incident. In this column there are 869 unique values.

4. breach_submission_date

The date when the breach was submitted or reported. In the context of data security and privacy, the term "breach" refers to an incident where unauthorized access, disclosure, or acquisition of sensitive information occurs. A data breach can involve various types of information, including personal data, financial records, or other confidential data. In this column there are 415 unique values.

5. web_description

Textual descriptions related to the security incidents and their impacts. In this column there are 4 unique values.

6. location_of_breached_information

The location where the breached information was stored (e.g., Network Server). It provides information about the specific location or system where the breached information was stored at the time of the security incident. This information is crucial for understanding the scope and potential impact of the breach. In this column there are 22 unique values.

7. state

The state associated with the covered entity or business associate. In this column there are 50 unique values.

8. covered_entity_type

The type of entity affected, such as "Healthcare Provider" or "Business Associate" or 'Health Plan' or 'Healthcare Clearing House.' In this column there are 4 unique values.

9. type_of_breach

Describes the type of breach, such as "Hacking/IT Incident" or "Unauthorized Access/Disclosure." In this column there are 5 unique values.

10. business_associate_present

Indicates whether a business associate was present during the security incident (Yes/No). In the context of healthcare and data security, a business associate refers to an individual or entity that performs certain functions or activities on behalf of a covered entity. In this column there are 2 unique values.

Data Overview: The dataset comprises information on 1000 cybersecurity events, with details such as the status of incidents, entities affected, the number of individuals impacted, breach submission dates, and various breach-related attributes. The entities involved include diverse types such as healthcare providers, business associates, health plans, and healthcare clearing houses. The breaches encompass a range of incident types, including hacking incidents, unauthorized access, and disclosure events. Notably, the dataset exhibits diversity in terms of geographical locations, with incidents reported across 50 states. The insights gained from the "web_description" column assisted in understanding the context, leading to the decision to drop this column due to a significant number of missing values. The comprehensive analysis of these columns sets the stage for deeper exploration and identification of patterns, trends, and potential areas of concern in the cybersecurity landscape.

3 Insights and observations

3.1 Visuals and Interpretations

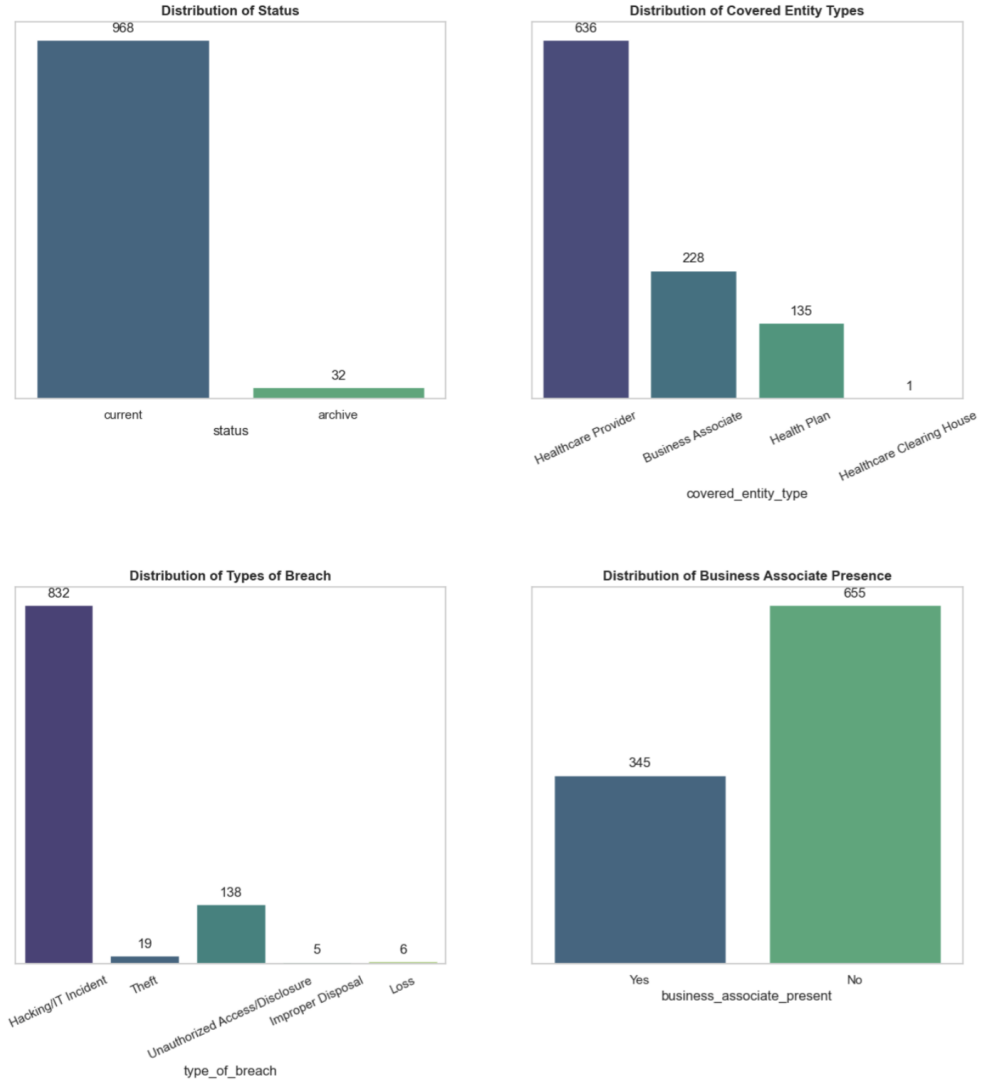


Figure 1: Exploratory Data Analysis.

Figure 1 presents visualizations of key categorical columns. The distribution of **status** reveals that the majority of records are labeled as "current," indicating ongoing or recent security incidents, while only 32 records are classified as "archive," signifying resolved or historical incidents.

Moving on to **covered_entity_type**, the graph illustrates that a significant portion of entities affected by incidents are Healthcare Providers. Business Associates follow with approximately one-third of the frequency, and organizations associated with health, Health Plan, represent another substantial segment.

The barplot for **type_of_breach** highlights that the predominant incidents involve hacking, characterized by unauthorized access to data systems or computers. This underscores the prevalence of cybersecurity threats involving unauthorized digital intrusions.

The distribution of the **business_associate_present** column is depicted in the last barplot. While the exact definition of a business associate is not explicitly provided, analysis of the available descriptions suggests a role related to data security. In most cases, the entities affected by breaches appear to lack business associates, implying potential gaps in data security measures.

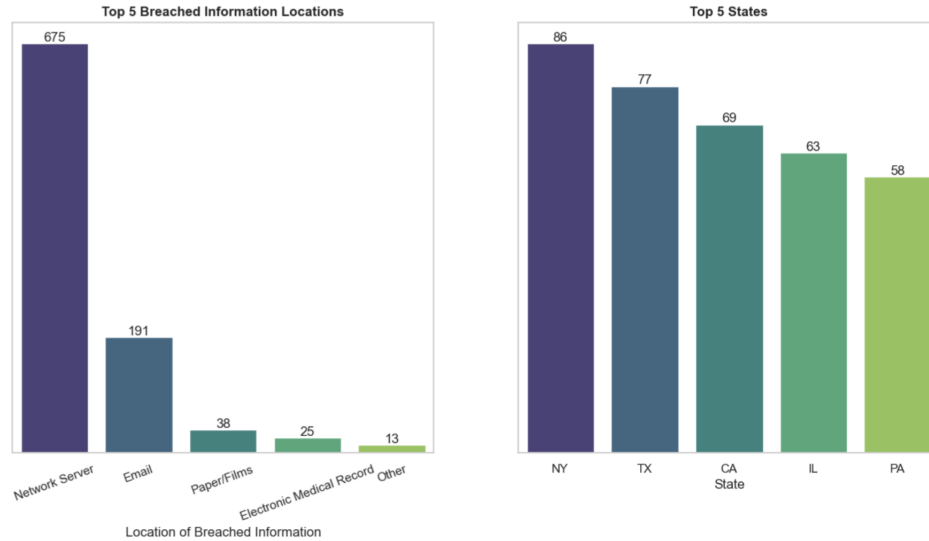


Figure 2: Exploratory Data Analysis.

Figure 2 shows visualizations of the top 5 categories of other categorical columns. The distribution of `location_of_breached_information` reveals that the majority of incidents are connected to "Network Server," "Email," "Paper/Films," "Electronic Medical Record."

The barplot for `state` highlights that the predominant incidents took place in New York, Texas, California, Illinois, and Pennsylvania. This underscores that quite a number of cybersecurity happen in the US.

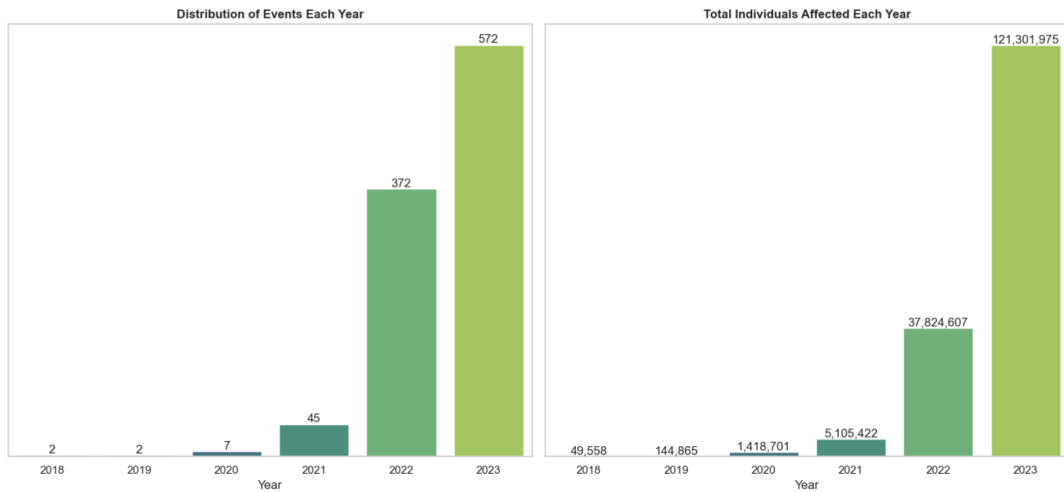


Figure 3: Exploratory Data Analysis.

Figure 3 presents visualizations of key numerical columns through barplots. Examining the distribution of the `breach_submission_date` column's years reveals that the majority of records occurred in the recent years, specifically in 2021, 2022, and 2023.

Moving on to the subsequent graph, it illustrates the yearly aggregate of affected individuals. The observed trend aligns logically with the increase in the number of cases during recent years, resulting in a higher count of affected individuals. Regrettably, the magnitude of individuals impacted is substantial, particularly noteworthy is the figure of nearly 121 million affected individuals reported in the year 2023.

3.2 Interesting Findings

1. Escalating Cybersecurity Incidents:

Over the past few years, there has been a noticeable surge in cybersecurity incidents. A substantial portion of these incidents is current, indicating an increase in recent years.

2. Healthcare Sector Dominance:

The healthcare sector constitutes a significant percentage of reported cybersecurity events. This trend underscores the vulnerability of the healthcare industry to such incidents.

3. Prevalence of Hacking Incidents:

A predominant type of cybersecurity event observed is hacking, which involves unauthorized access with the intent to corrupt or modify data. This highlights the evolving and sophisticated nature of cyber threats.

4. Lack of Cybersecurity Specialists:

Notably, a considerable number of organizations affected by cybersecurity events lack dedicated specialists in cybersecurity. This raises concerns about the preparedness and resilience of these entities against potential threats.

5. Common Targets: Network Servers and Email Systems:

The primary targets of cybersecurity events are network servers and email systems, indicating the sensitivity of these areas to security breaches.

6. Geographical Concentration in the U.S.:

Geographically, cybersecurity events are concentrated in states such as New York, Texas, California, Illinois, and Pennsylvania. This regional pattern may offer insights into the distribution and prevalence of cyber threats.

7. Escalating Impact on Individuals:

The number of affected individuals has seen a continuous rise in recent years, reaching numbers in the millions. This emphasizes the critical importance of prioritizing and enhancing cybersecurity measures to safeguard sensitive information.

These findings collectively underscore the urgency for organizations to strengthen their cybersecurity posture and implement proactive measures to mitigate potential threats.

4 Challenges for integration

While integrating this dataset with other data sources, several challenges may arise, necessitating careful consideration and mitigation strategies:

Data Compatibility:

- **Formats and Standards:** Differences in data formats and standards across datasets can hinder seamless integration. It's essential to ensure that the data adhere to common standards to facilitate effective merging.

Normalization and Standardization:

- **Variable Names and Units:** Inconsistencies in variable names and measurement units may exist among datasets. Normalizing these variations is crucial for accurate integration, requiring a standardized schema for shared variables.

Missing Values:

- **Addressing Missing Data:** The presence of missing values, especially in critical fields, can pose challenges. Imputing missing data or applying appropriate techniques is vital for maintaining data integrity during integration.

Data Quality Issues:

- **Accuracy and Reliability:** Assessing the overall quality of the data is essential. Identify and address issues related to accuracy, reliability, and completeness to ensure the integrated dataset provides meaningful insights.

Inconsistencies:

- **Temporal and Geographical Inconsistencies:** Datasets may have variations in the temporal or geographical granularity, impacting alignment. Addressing these inconsistencies requires careful consideration and potentially harmonizing the data.

Security and Privacy Concerns:

- **Sensitive Information:** Integrating datasets may involve sensitive information. Ensuring compliance with privacy regulations and implementing robust security measures is paramount to protect individual privacy and organizational interests.

Data Governance:

- **Establishing Protocols:** Lack of clear data governance protocols can lead to challenges in managing integrated datasets. Establishing guidelines for data ownership and usage is crucial.

Scalability:

- **Handling Large Datasets:** Integration becomes complex when dealing with large volumes of data. Employing scalable infrastructure and efficient algorithms is necessary for processing and integrating extensive datasets.

Metadata Alignment:

- **Harmonizing Metadata:** Ensuring that metadata, including data dictionaries and descriptions, aligns between datasets is vital for understanding the integrated dataset comprehensively.

Lack of Primary Key:

- **Identification Challenges:** The absence of a standardized primary key or unique identifier across datasets can complicate integration efforts. Establishing a common identifier or exploring alternative matching methods is crucial for linking related records.

Versioning:

- **Dataset Evolution:** Over time, datasets may undergo changes. Managing versioning and tracking dataset evolution is essential for maintaining a clear record of integrated data.

Categorical Variations:

- **Diverse Categories:** Varied categorizations across datasets may impede integration. Establishing a unified categorical framework or employing advanced techniques for handling diverse categories is essential.

Addressing these challenges proactively will contribute to a more successful integration process, allowing for the creation of a comprehensive and cohesive dataset from diverse sources.

Data Integration Challenges

To sum up ...

- Integrating data from different sources can be a tricky task because each source has its way of organizing information. One big challenge is making sure everyone understands the data in the same way. It's like speaking the same language about the data so that everyone uses it consistently. Another challenge is understanding where the data comes from and where it's going. Different places store data differently, and understanding this is like having a map for the data journey. We need to teach our teams about these different places, keep good records, and use tools to help us see how data moves.
- In simpler terms, imagine everyone in a company speaking the same data language, and it's like knowing the map of where the data lives and travels. These challenges need teamwork, clear rules, and handy tools to solve them.
- Dealing with lots of data can be like solving a puzzle. We need to make sure everyone understands the pieces the same way, and we must have a good map to know where each piece belongs. These challenges might seem tough, but there are smart ways to solve them. Things like having clear rules for using data, making sure our systems work well together, and using tools that help us handle lots of data make a big difference. As we move forward, technologies like Estuary Flow show us how we can make these challenges easier, ensuring that dealing with data remains a key part of making good decisions.

5 Recommendations

Data Compatibility:

- **Formats and Standards:** To overcome differences in data formats and standards, establish a data normalization process. Ensure that all datasets adhere to common standards, facilitating seamless integration. Implement data format conversion tools if necessary.

Normalization and Standardization:

- **Variable Names and Units:** Address inconsistencies in variable names and measurement units by implementing a standardized schema. Develop a data dictionary that provides a clear mapping of variable names and units to ensure a harmonized dataset.

Missing Values:

- **Addressing Missing Data:** For handling missing values, employ imputation techniques or utilize statistical methods to impute missing data points. Clearly document the imputation methods applied to maintain transparency in the integration process.

Data Quality Issues:

- **Accuracy and Reliability:** Implement a rigorous data quality assessment process to identify inaccuracies and inconsistencies. Develop data profiling reports and invest in data quality management systems to ensure the reliability of integrated datasets.

Inconsistencies:

- **Temporal and Geographical Inconsistencies:** Address temporal and geographical inconsistencies by standardizing the granularity of data. Harmonize temporal and geographical attributes across datasets to ensure a coherent integration.

Security and Privacy Concerns:

- **Sensitive Information:** Prioritize data security and privacy by implementing encryption, access controls, and anonymization techniques. Conduct regular audits to ensure compliance with privacy regulations and bolster security measures.

Data Governance:

- **Establishing Protocols:** Develop and enforce clear data governance protocols, including guidelines for data ownership, usage, and sharing. Assign data stewards to oversee adherence to these protocols and ensure a well-managed integration process.

Scalability:

- **Handling Large Datasets:** Scale infrastructure and employ efficient algorithms to handle large volumes of data. Utilize distributed computing and cloud-based solutions for scalability, ensuring seamless processing and integration of extensive datasets.

Metadata Alignment:

- **Harmonizing Metadata:** Establish a centralized metadata repository and ensure that metadata aligns between datasets. Implement standardized metadata formats and regularly update the metadata repository to reflect any changes in the integrated datasets.

Lack of Primary Key:

- **Identification Challenges:** Introduce a standardized primary key or unique identifier across datasets. Explore alternative matching methods, such as fuzzy matching or probabilistic matching, to link related records in the absence of a common identifier.

Versioning:

- **Dataset Evolution:** Implement version control mechanisms for datasets. Maintain clear records of dataset versions and changes, allowing for easy tracking of dataset evolution over time.

Categorical Variations:

- **Diverse Categories:** Address varied categorizations by establishing a unified categorical framework. Utilize advanced techniques, such as one-hot encoding or embedding layers, to handle diverse categories and ensure consistency in integration.

Strategies to address challenges

Effectively addressing the challenges of data integration requires collaborative efforts and a commitment to established standards. It is imperative for teams to work cohesively, fostering a shared understanding of data and ensuring consistency in interpretation. Regular meetings among team members play a pivotal role in discussing integration methods, resolving discrepancies, and maintaining alignment across diverse datasets.

Furthermore, the establishment and adherence to common standards are important. Teams must work towards creating a unified language for data, ensuring that variable names, measurement units, and overall data structures conform to standardized schemas. This harmonization facilitates a seamless integration process and contributes to the accuracy and reliability of the integrated dataset.

Staying abreast of modern tools and techniques is equally essential. Team members should actively seek knowledge about the latest advancements in data integration tools, methodologies, and best practices. This continuous learning approach equips teams to leverage cutting-edge solutions, enhancing the efficiency and effectiveness of the data integration process.

Therefore, making data work together smoothly depends on teams working well together, sticking to the same rules, talking to each other a lot, and using new tools when needed. When teams follow these simple ideas, they can solve problems and make data work together in a way that's organized and makes sense.

6 Conclusion

In this report, a comprehensive analysis of the provided data was conducted, aiming to understand its structure, contents, and characteristics. The data summary highlighted key attributes, providing insights into the dataset's format, structure, and basic statistics. The exploratory analysis uncovered notable patterns, trends, and anomalies within the data, shedding light on intriguing findings and peculiarities. As we delved into the challenges for integration, potential difficulties related to data compatibility, normalization, missing values, data quality issues, and inconsistencies were discussed. In light of the identified challenges, recommendations were proposed to address issues and ensure a successful integration process. In summary, this report provides a holistic examination of the data, offering insights, addressing challenges, and presenting actionable recommendations for a seamless integration process.