International Journal of Environmental Sciences ISSN: 2229-7359 Vol. 11 No. 18s, 2025 https://theaspd.com/index.php

# Data Mining For Metadata In Telecom Sector

Bhola Gan Chaudhuri<sup>1</sup>, Shalli Rani<sup>2</sup>

<sup>1,2</sup>Chitkara University Institute of Engineering and Technology, Chitkara University, Punjab, India, bhola\_in@yahoo.com, shallir79@gmail.com

Abstract—In the telecom sector, for metadata mining, IE, clustering, NLP, and complex algorithms such as TF-IDF, Witten-Bell smoothing method, BERT, Genetic Algorithm Cycle are used. Voices and messages gathered by the telecom industry from customers comprise customer reactions, service requests, network upkeep records and others which are unstructured textual data that go through NLP analysis. Since IE is the process of automatically identifying and structuring specific entities such as phone numbers, addresses, service kinds and technical terms from unstructured texts from several sources, IE enriches NLP. For telecom metadata to be analyzed, categorized, or sorted to make them easily discernible, comprehensible, and manageable, clustering is very vital. Another statistical measure called TF-IDF (Term Frequency-Inverse Document Frequency) is applied for analyzing the relevance of a word in a document regarding the importance of that word in a set of documents. The Witten-Bell smoothing formula as a way of estimating the occurrence of new words or events in language modeling in connection with the observed data.

keywords- Data mining, Metadata, Text mining, Witten-Bell smoothing formula, Genetic Algorithm Cycle.

### I. INTRODUCTION

"Data about data" is what metadata basically is. It offers details about the features, qualities, or traits of a certain piece of data making it simpler to handle comprehend and utilize. In the telecom field, the metadata is meant the extra information connected to telecom information. All what telecommunication does the calls the messages, the network traffic, the subscriber's information, and many other aspects of its activities, are fully captured in the metadata which also afford context and structure meanings [1]. Data mining can thus be defined as extracting of more relevant and useful inferences and trends with regards to telecoms from large volumes of data within the context of the telecom industry. Telecom firms can obtain such data from sources like CDRs network logs interaction with customers tariffs and social networking sites. By means of developing data mining over these features one can discover hidden patterns, trends, and relationships that are useful in increasing satisfaction among the customers performance of the networking, organizational processes of the corporations, and even useful in decisionmaking processes. In the context of the telecom industry, it is the assessment of the textual data that the telecom industry produces in the form of NL texts and then making decisions based on outcomes of such data gained through text mining [2]. They are commonly bundled up in telecommunications organizations production and numerous other areas such as networks, customer service call centers social media posts technical documentation et cetera Text mining therefore helps numerous business entities in the telecommunication industry via offering an imperative tool in decision-making processes that focuses on increasing customer satisfaction, efficiency and gaining crucial information.

In the context of telecom metadata is defined as the application of data mining approach that aims at utilizing the metadata of various telecom data and extract from these data set informative interesting and important characteristics. Telecom data such as the CDRs and other data included in the Telecom data like network logs or the subscriber details can be pulled out and analyzed Telecom data like call detail records, network log data, and subscriber data or any data pertaining to telecom. Telecom business itself can partially understand some trends and periodicities in the call numbers, traffic densities, customers, and services by practicing data mining [3]. For instance, metadata mining could be used in telecom operations that need to partition their number of consumers by the frequency when it comes to their usage, their age, sex, and more. Thus, by studying metadata linked to the customers information and services being engaged in the current telecom organizations processes different client segments can be safeguarded to offer customized approaches marketing strategies services and tariffs. That is metadata mining is helpful for the telecom operators in a manner that ensures them the actual pictorial view of Resource Utilization and its performance in the entire Network as a whole. From meta-information regarding network traffic and the available network capacity and servers the operators can define where

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

the troubles are in the network improve the efficiency of connection data rate and thus elevate customers QoS. The originality of the approach based on data mining techniques to predict consumers behaviors and service consumption through the assessment of information. Telecom companies can use evaluation data like call patterns, use habits, billing data and many other things to note the clients who are very probable to leave and can then apply certain strategies to retain them [4]. Exploring telecommunication metadata can lead the telecom organizations to track the fraudulent activities like sim card cloning, call spoofing and billing fraud. Every time a call is made there are metadata attached to it the same goes for the location and the billing information by analyzing this data the telecom companies can identify certain anomalies and patterns that suggest fraudulent activities. Service qualities may include call dropouts, network delay and data rate are some of the parameters telecom companies may evaluate by metadata mining. This paper also reveals that operators can organize resources to such network alterations, find service quality problems and make changes to improve customer tendency through assessing data that are associated with network performance indicators. Regarding the promises of compliance, the study suggests that in the telecom business regulatory compliance could be positively achieved with application of data mining techniques. It has been noted in [5] that depending on the call records metadata, the subscriber information, as well as the billing information may be used by the operators to verify the adherence of the privacy regulations and the data protection legislation and other standards.

In telecommunications, small bits of data transmitted before other data are commonly known as metadata. It provides information and organization and informs about various things concerning to the telecom messages and calls and network load and subscribers' data and about infrastructure as well. The specific details in conscious phone call metadata include the number of the caller and recipient, the time taken in the call, the time that the call was made and call type for instance a voice or SMS among others, the geographical origin and destination of the call and lastly indicators of call quality. Metadata about whom the message was sent by and to whom and what the content of the message is, when it was sent, and whether it has been delivered, all regarding text-messaging, multimedia messaging, and instant messaging. Static and dynamic features of network topology, address of IP, ports, protocol specifications, size of packets, transmission rates and other configuration parameters are among the options of metadata of traffic and performance of the networks. Identifying information includes but is not limited to subscriber identifiers such as phone numbers, subscriber number from the SIM card, demographics like the name and the physical address of the subscriber, account status involving active and suspended accounts, service plans and bill information are other parts of the metadata associated with the subscriber's accounts [7]. Identifiers of geographical area such as the LAC which stands for the location area code TAC which stand for the tracking area code GPS coordinates and ID of the cellular tower are subcategories of the metatags concerning the physical position of the mobile devices and of the networks. Some of the possible types of metadata used for accounting/billing are the call rates/tariffs, use log data, billing cycle/timetables, invoice details and payment history [8].

Characteristics information that relates to the telecom's provider including organization options, organization subscribership organization identifiers and accomplishment position of telecom services. Telecommunications devices, including smartphones, tablets, routers and other equipment and accessories are associated with non-content data including device identifiers like IMEI and MAC addresses, device capabilities such as screen size and processing power, firmware versions, usage history and habits among other features. Application access logs, authentication data, security incidents and detailed information about implemented security policies can be considered as meta-data utilized for security and audit [9]. As mentioned above and shown in the role of metadata for telecom data management in the industry in Figure 1 metadata offers the industry with the capability of monitoring, evaluating and enhancing a wide array of telecoms services, framework and processes. In the telecom industry, it helps in the management, analysis and decision making of data and assists in carving the data into conformed and accurate regulatory standard for its clients.

The process of mining a large amount of data enclosed in the telecom field to look and find out the opinion and tendency of the clients is known as data mining in the telecom industry. A telecom company receives large amount of information from various sources like social media, CDRs, network activity report, emails, and billing records. The above data is analyzed using data mining where large amounts of

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

trends associations are discovered and that will help enhance customer satisfaction and network capacity to minimize the overwhelming nature of corporate initiatives and support the function of strategic management. Data mining is one method through which the telecommunication can develop its clients as per their usage level, age, preference, or tendencies. Segmentation is also useful for customer retention operations together with personalized services and management of marketing promotion [10]. Customer pattern means which services are being used; the experience clients have, and past data are some of the ways data mining is used to arrive at the client attrition rate of a company. Thus, through a combination of the models, it allows the companies to avoid churn by recognizing the consumer and implement the retention initiatives as soon as possible. Data mining can be as an effective tool not only in terms of assessing several indicators related to the performance of telecom operators but also in terms of discovering the limiting factors for the expansion of their networks and improving those possibilities. Other services include predicting network outages, getting the maximum coverage for services and offering a required QoS to the subscribers, all these services are enhanced by data mining.

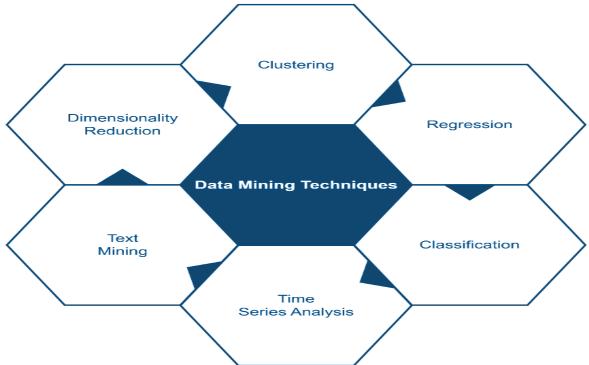


Figure 1: Data Mining Techniques in Telecom Sector

As for the case of call taping, SIM card cloning or billing fraud Ind Ontime's data mining algorithm is applied in the detection of fraud activities. Ultimate goals of telecom businesses include identification of the anomalous conditions and fraud within real-time call and usage and transactional pattern analysis [11]. Data mining is applied in RA where billing data is analyzed to identify billing problems and when the two do not match they are brought to a resolution. Using the analyzed data, it is also possible to define income streams, consider Billing, and identify proper methods for improving the pricing strategies of telecom companies [12]. While analyzing the data of the telecommunication industry, data mining helps to learn about the behavior of customers, competition strategies and trends at the macro-level [13]. By critically evaluating the information that gets posted by consumers of the various social media platforms together with business reviews as well as the trends that prevail within the market, the various telecommunications businesses get to explore new opportunities and also be in a position to scrutinize different threats that are likely to affect the growth of the different businesses before going ahead and making sound business decisions. Other aspects where this can be used include forecasting of network disruptions, system and others are device malfunction and inefficiencies, maintenance needs. For case, there are existing measures that telecom firms can adopt to improve availability and dependability of the network; this is through using predictive maintenance approaches, where momentary data from the network, telecom equipment sensors, logs, and documents, and historical maintenance records are used [14].

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

The telecom industry can also benefit a lot from the contributions made toward data mining for metadata Various areas such as networking, customers and security are some of the areas that data mining would make an easy enhancement.

- 1. The roles played by metadata in the analysis of networks can be implemented through data mining to address potential issues regarding the performance of the networks.
- 2. The following is a research on how data mining in the telecom industry can uncover fraudulent activity. Telecom businesses can block possible risks in an actual time by addressing metadata patterns that relate to suspected behaviors due to the usage of fraud detection algorithms in the company's operations.
- 3. Within the telecom industry, metadata analysis can assist in guaranteeing adherence to industry standards and legal obligations.

Potential drawback of the contribution made to data mining metadata could include; Innovation has led to the stiffness in the telecommunication industry, improvement in the security measurements of the customer, and advancement in the delivery of the services. These I consider as the benefits that can take the industry to another level of competitive advantage in the market.

## II. THEORETICAL BACKGROUND

Data mining hence captures procedures, techniques and technologies used after data collection with the objective of identifying hidden tendencies in the acquired databases. It is possible to classify these techniques into several general groups and every type of method is applied to a particular data mining task or aim [15]. This unearths the realization that classification is one of the strategies in a supervised learning technique, whereby the data have to be mapped and later classified into already-defined classes although based on features [16]. The detailed classification algorithms are known to be decision tree, random forest, SVM, naive Bayes and KNN. The next method of classifying the dataset is regression; this is another supervised learning technique that is used in outcome predictions that are numeric in nature. Regression algorithms are developed in such a manner where it is aimed at seeking a relation between the input and the corresponding output of value. As far as the approaches which can be applied to solve the proposed problem, it is possible to highlight polynomial, ridge and linear regression [17]. Clustering it a technique under unsupervised learning that groups related data into segments or clusters based on the attributes or features that it may have [16]. Gaussian Mixtures Models, Density-Based Spatial Clustering of Applications with Noise are some of its prominent approaches and there are two broad methodologies for clustering namely K-Means and Hierarchical clustering. Implemented strategies that are used in transforming the number of variables or the features of a dataset while preserving crucial information. Other methods that can be used to perform the dimensionality reduction are PCA, t-SNE analysis as well as LDA.

Text mining should be used where data is in textual format because it structured, since it is capable of offering insights and trends. It also results in the generation of activities that include classification of texts, subjective analysis, topical segregation, identification of personalities and places, and summarization of texts. Text mining involves the use of computational techniques for analysis and is generally connected to Artificial Intelligence, and more particularly, to the fields of machine learning and Natural Language Processing. In essence, time series analysis can simply be defined as a technique of trying to look for a model of variation in data that has been collected at different intervals of time. In the case of time series, therefore, applicable time series forecast methods that commonly deployed include, for example. There are the ones mentioned above recurrent neural networks and their types: Long short-term memory and Exponential smoothing methods and the Autoregressive Integrated Moving Average [19].

Basically, metadata is a significant benefit in applying text mining since it offers more values to organizations for utilizing the social data, transforming better decision makers, having a quicker focused on the process of information searching and analyzing. Further extending the metadata for the data can complement this context by using the more sophisticated techniques for the text mining of the raw data, such as topics, themes or sentiment that is associated with the data [20]. Metadata content includes text content description, labels, and note from the data that were obtained from processing the data. It can also help with figuring-out how to make the aforementioned textual information better for making decisions and equip us with a better understanding of the data in question.

The current technique for this text categorization in the data mining is through employing only the genetic or the algorithms of the support vector machines however the process does not yield the results that are

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

desirable. Therefore, one would need to employ higher level of design in stages that would provide all the necessary outcomes. At present, each of the categories namely clustering, information extraction and natural language processing has been used only in one version [21].

Metadata can be augmented with text mining since it is an efficient technique of finding information from large of text documents. This can involve identifying keywords, entities, or relation in text description making a broader metadata capture more thorough and valuable for further analysis [4].

Mining textual metadata will improve the search and the retrieval of textual data by extracting and storing relevant keywords through which the user can easily retrieve data or documents in response to a particular keyword search. Keyword extraction can also be achieved on the metadata in order to enhance the search engine accuracy and relevance of the results [22]. The document clustering or similarity analysis can also be done to metadata for better result of searching. It signified that more abstract patterns, trends or outliers that are difficult to identify by simple analysis of metadata properties can be easily detected when using text mining on metadata. Insights and discoveries, organizational learning may be of great value by considering textual descriptions or annotations as a way of analyzing textual data to look for patterns or connections that may not be apparent otherwise [23].

Text mining techniques can be leveraged to infer metadata attribute semantics. Organizations can make sense of the data using the meaning or context of metadata attributes inferred from textual descriptions or labels and properly use such data [24].

IE, clustering algorithms and NLP are very important in the telecommunication industry since they aid in the extraction of insight refine consumer experiences related to products optimize network operations, which in turn have very significant implications for corporate decisions [25]. Wiring methods belonging to the telecom industries in the text mining domain have been applied as follows:

A. Natural Language Processing: Call center transcripts, emails, chat logs and social media posts are just a few of the channels in which consumer interactions are analyzed using NLP [26]. By applying sentiment analysis, named entity recognition, and topic modeling, NLP helps telecom companies understand customer sentiments, identify common issues, and extract actionable insights for improving customer service [27].

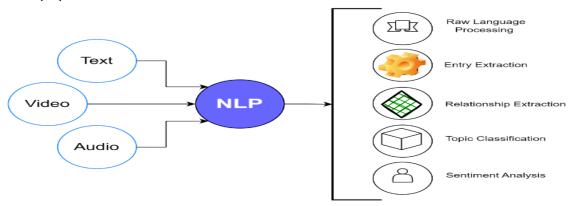


Figure 2: NLP for Data Mining

Figure 3: Flow of Information Extraction for Text Mining NLP enables telecom companies to transcribe voice calls into text allowing for further analysis and processing. This facilitates tasks such as sentiment

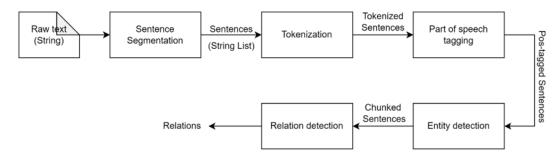


Figure 3: Flow of Information Extraction for Text Mining

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

analysis, call categorization and trend identification helping companies understand customer needs and preferences more effectively. Telecom companies often deal with large volumes of textual data in the form of reports, articles, and customer feedback. NLP techniques for text summarization help in condensing this information into concise summaries, enabling decision-makers to quickly grasp key insights and trends [28]. TF-IDF is a fundamental technique used for various tasks, including text preprocessing, feature extraction, and document similarity computation.

$$\mathbf{TF} = \frac{\text{No. of times a word "X"appears in Documents}}{\text{No. of times a words present in a Document}}$$

$$\mathbf{IDF} = \log \left( \frac{\text{No. of documents present in a corpus}}{\text{No. of Documents where word "X" has appeared}} \right)$$

$$TF(t,d) = \frac{\int_{t,d}}{\sum_{t' \in d} \int_{t',d}}$$

$$IDF(t, D) = log \frac{N}{|\{d \in D: t \in d\}|}$$

$$TF \cdot IDF(t, d, D) = tf(t, d) \cdot idf(t, D)$$

Using natural language processing in data mining can help organizations get a competitive advantage in today's data-driven environment by providing vital insights, improving decision-making, as shown in Figure 2.

B. Information Extraction: IE techniques are employed to pull out customer name, location, phone numbers and product names among other things from natural text. This information is useful for operations such as customer profiling, custom marketing, and bespoken service delivery [29]. Internet business organizations utilize IE as a model for extracting significant events and information from different textual data including the news, microblogs, and telecommunications logs. By closely monitoring competitive activity, market trends and regulatory advancements, businesses can stay informed and adaptable to changes in the marketplace. Internet Explorer makes it easier to identify connections between different pieces of text. For example, identifying the relationship between customers and the products or services they are discussing in their feedback or identifying associations between network performance metrics and specific geographical locations.

Business can effectively transform the raw text into structured data, as indicated in Fig 3, making it possible to apply a host of functions including automated textual processing, decision making and data mining. Typically, each of the stages of the flow helps to read the text data more effectively and then extract the important information and knowledge from the raw text data.

In perhaps the most relevant use of Witten-Bell smoothing in text mining, it calculates the probability of occurrence of a certain phrase or a pattern in a certain context. For example, when extracting entities from text it can apply Witten Bell's method to predict how likely it is to encounter an entity given the words surrounding it [30].

The Witten-Bell smoothing formula calculates the smoothed probability  $P_{WD}(w_i|w_{i-1})$  of an event  $w_i$  occurring given its context  $w_{i-1}$ . Its calculated as:

$$P_{WD}(w_i|w_{i-1}) = \frac{C_{(w_{i-1}w_i)+\delta.V}}{C_{(w_{i-1})+\delta.V}}$$

- $C_{(w_{i-1}w_i)}$  is the count of the specific event  $(w_{i-1}w_i)$  occurring together.
- $C_{(w_{i-1})}$  is the count of the context  $w_{i-1}$ .
- V is the number of unique events in the vocabulary.
- $\delta$  is a smoothing parameter that is often set to 11.
- C. Clustering for Text Mining: A customer segment or a cluster is created based on the similarity of the customers' preferences, their mode of interaction and behavior towards the services. From the above

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

context, marketers in the telecom industries are in a better place to come up with solution to different segment's needs through classification of their consumers into various clusters. As a result, when the network performance data is clutered, the telecom operators will be able to know the clusters of the network nodes or regions that will be similar. This assists in preventing underperformance in aspect that are valued most by the user, identify areas that require network attention and enhance the utilization of networks [30]. The telecom providers employ clustering algorithms to arrive at trends, which may represent such fraud as identify theft, toll fraud and cloning of SIM cards. To curb fraudulent activities, clients may use fake alarm bells or align various forms of irregularities or disparate called patterns which would otherwise incur costs while at the same time, keep the consumer's trust in them.

To telecommunication companies, the textual data obtained from the news, social media, or counterparts' reports concerning the relevant industry, the preferred service aspects, and competitors' positioning could be extracted by subjecting the textual data to cluster analysis. It is used to bring changes in the company's products and services; the marketing strategies and major decisions in a bid to enhance competitiveness in

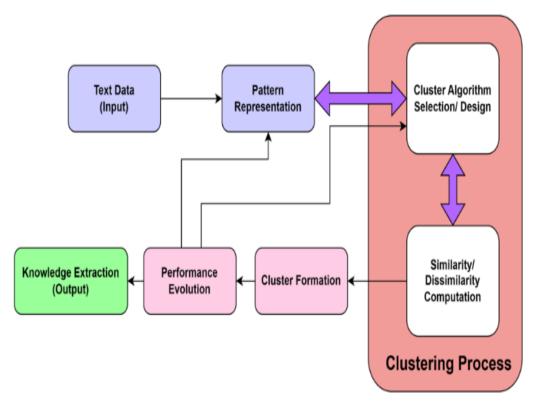


Figure 4: Clustering flow process for Data Mining

the market [31].

The use of highly strategic flows in this process of clustering in data mining will allow the organizations to gain the strategic tools required to correctly dissect the available data and identify the meaningful patterns of the featured data structures as demonstrated in the figure 4. All of them support the overall goal of managerial application of clustering as the way to search for the buried knowledge in the data.

Google created the trained BERT natural language processing model. By capturing the reciprocal interaction between words, it is intended to comprehend the context of words in a phrase. In text mining jobs, BERT can also be used for text clustering. Yes, the K-means clustering technique may be used to group or cluster data points based on their similarity. The goal of the clustering method is to minimize an objective function to divide N data with D dimensions into D clusters. The minimized objective function for a D-dimensional data set  $\{x_1, x_2, ..., x_N\}$  can be seen in the following equation:

$$J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} ||x_n - \mu k||^2$$

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

The membership value of data  $x_n$  in cluster K is represented by the value of  $r_{nk} \in \{0, 1\}$ . The total squares of the distances between each  $\mu k$  centroid and each  $x_n$  data point represents the objective function J.  $r_{nk}$  and  $\mu k$  must be determined to be at the proper values to reduce J. The following equations can be used to iteratively assign values of  $\{r_{nk}\}$  and  $\{\mu k\}$ , respectively, to minimize the objective function J.

$$r_{nk} = \begin{cases} ^{1}\text{,} K = \arg\min K ||x_n - \mu k||^2} \\ 0, & \text{Others} \end{cases}$$

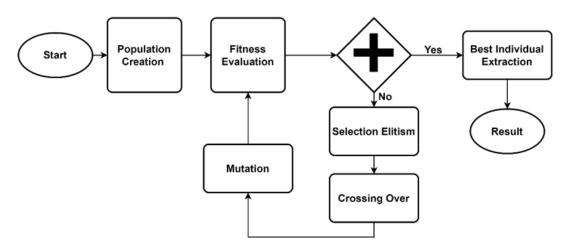


Figure 5: Genetic Algorithm Cycle

$$\mu_{nk} = \frac{\sum_{n=1}^{N} r_{nk} x_n}{\sum_{n=1}^{N} r_{nk}}$$

Telecom companies can achieve business objectives in a rapidly changing industry landscape by driving operational efficiencies, improving customer experiences, and making data-driven decisions by utilizing text mining techniques that leverage NLP, IE, and clustering to unlock valuable insights from their textual data assets [32].

D. Text Mining by Combining techniques of Genetic Algorithm: Genetic algorithm is an optimization technique which is named from the selection process derived from a evolutionary technique. Among the text mining activities that can be mentioned to highly benefit from the proper GA approach, it is possible to list the following ones: tasks in model tuning, parameter optimization and feature selection. So, the objective of designing a feature extraction algorithm for text mining is to extract a subset of features from a text database relevant for a certain job and at the same time guarantees maximum performance for the intended job, such as classification or clustering. The genetic algorithm cycle is explained and depicted in the figure 5 below.

This is how such an algorithm may be defined, essentially, the first step in a genetic process is a setup of the goal function. The fitness function that we came up with to explore the space of solutions is defined as:

$$F(P) = \sum_c F(c) = \sum_c \sum_d W_{c,d} = \sum_c \sum_d sf_{c,d} \times ids_c \ = \ \sum_c \sum_d sf_{c,d} \times \log \frac{N}{ds}_c \ ;$$

Where,  $F(c) \rightarrow$  fitness of c.

Wc,d  $\rightarrow$  Normalized information of c within d.

Sfc,d  $\rightarrow$  Frequency of c within d.

idsc  $\rightarrow$  Inverse frequency of c.

 $dsc \rightarrow Document number containing c.$ 

A similarity function is used to compare the keywords vector of the query and the documents. Every document that matches a preset threshold value is retrieved. The widely recognized cosine measure is the similarity function that is most frequently used:

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

$$sim (q, d_j) = \frac{\sum_{i=1}^{n} W_{i,j} \times W_{i,q}}{\sqrt{\sum_{i=1}^{n} W^2_{i,j}} \times \sqrt{\sum_{i=1}^{n} W^2_{i,q}}};$$

Where, Wi,j  $\rightarrow$  Unit information i in document dj.

Wi, $q \rightarrow$  Unit information of i in query q.

 $n \rightarrow Number of the query q.$ 

Text mining by Genetic Algorithm provides a powerful and adaptive approach for extracting knowledge from text data, offering flexibility and efficiency in discovering patterns and insights from textual information.

- E. Evaluation of Performance: To obtain a thorough picture of text mining systems' performance, it is important to consider a variety of measures. Precision, recall, accuracy, and F-measure metrics are commonly utilized for evaluating text mining systems for tasks that involve text categorization, sentiment analysis, information retrieval, and named entity identification. The following describes the computation of these metrics and their importance:
- a) Precision: The measurements of a model can be done through the percentage ratio of positive predictions that the model makes.

$$Precision = \frac{True \ Positives}{True \ Positives + False \ Positives}$$

b) Recall (Sensitivity): Recall measures the model's ability to provide accurate results for each individual positive class sample that it samples. It is computed by using the formula that is true positive divided by TP plus FN where TP is the true positive; FN is the false negative.

$$Recall = \frac{True \ Positives}{True \ Positives + False \ Negatives}$$

c) F-Measure (F1 Score): F-Measure which provides a balance between recall and precision is computed as the arithmetic mean of the two values divided by the product of the two values. If, memory and precision are on differing levels, the opposing element is very advantageous in document processing.

F1 Score = 
$$2 \times \frac{\text{Precision * Recall}}{\text{Precision + Recall}}$$

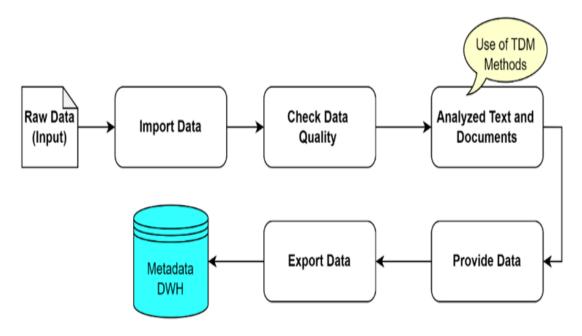


Figure 6: Base-line Model

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

d) Accuracy: Accuracy takes into consideration true positives as well as true negatives that is the overall accuracy of the pursued model in terms of predictions. It is determined with the help of the formula reflecting the ratio of successfully anticipated cases to the number of instances.

Accurecy = 
$$\frac{\text{True Positives} + \text{False Negatives}}{\text{Total no. of Instances}}$$

Since memory and precision are often two sides of the same coin, it means that improvements to one of the aspects can be followed by drawbacks to the other. This it is by noting that the F-Measure is effective when it is used in between the two measures. For unbalanced datasets, accuracy could be misleading because by merely predicting the majority class, a high accuracy could be scored. In such cases F-Measure, accuracy and recall provide additional measures on the merit of the model [33].

#### III. PROPOSED METHODOLOGY

Genetic algorithms, sometimes known as GAs, are optimization methods derived from natural selection and genetics. They are employed in the approximate solution of search and optimization issues. Large volumes of text data are sent to telecom companies from various sources, including social media, network logs, and customer care exchanges. Text data can be accurately classified into categories such as customer sentiment, service concerns, or network problems by using genetic algorithms to optimize the parameters

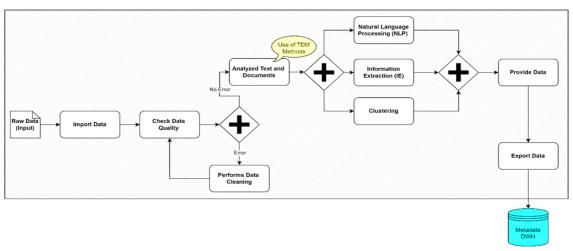


Figure 7: Proposed Methodology

of text classification models, such as SVMs or neural networks.

It is facing challenges in text categorization using genetic algorithms and SVMs alone and considering a more advanced approach. Combining techniques from multiple categories—clustering, information extraction, and natural language processing is a wise strategy to improve the outcomes.

Businesses can effectively use text mining techniques to extract insightful information from metadata, enhance data stewardship procedures, and guarantee the quality and integrity of their data assets by putting in place an effective architecture for text mining of metadata that makes use of NLP, IE, and Clustering techniques. Long-term success in data management and governance is further ensured by ongoing monitoring and improvement methods, which allow businesses to adjust to changing data landscapes and legal requirements [33].

So, here's an architecture that can handle all these requirements:

The Data Ingestion Layer is responsible for gathering information from different sources in the telecom network. This includes things like billing systems, network devices, CDRs and client interactions. To get this information, we use connectors and APIs to extract metadata from various data repositories and streaming sources. For real-time data intake, use technologies like Apache Kafka or Apache NiFi; for batch processing, use frameworks like Apache Spark or Apache Flink to handle massive amounts of metadata. Utilize a distributed, scalable storage system that can manage both structured and unstructured data to

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

hold ingested metadata. For storing raw metadata, make use of data lakes or distributed file systems like Amazon S3 or Hadoop Distributed File System. To store processed and aggregated metadata, employ columnar databases or data warehouses such as Amazon Redshift or Apache Hive.

Before diving into analysis, it's important to prepare the data by doing some preprocessing, transformation, and enrichment on the metadata that has been collected. This step ensures that the metadata is clean, consistent, and ready for analysis. To make sure everything is up to par, we clean up and standardize the metadata to ensure its quality and uniformity.

.In order to make the information even more valuable for analysis, we employ various methods such as entity resolution, normalization, and feature engineering. These techniques help us enhance the metadata and extract meaningful insights from it. For the purposes of classification, regression, clustering, anomaly detection and pattern recognition, apply machine learning methods. Use text mining tools to examine unstructured metadata, like network logs, customer reviews, and call transcripts. Utilize tools for statistical analysis and visualization to investigate trends and patterns in metadata.

Create standards, procedures, and policies for metadata governance to guarantee data stewardship, compliance, and quality. To categorize, tag, and document metadata attributes, use metadata management tools. Establish ownership, access, and data lineage guidelines to monitor and regulate metadata at every

	Average Baseline Result		Average Processed Result	
Collection	Precision	Recall	Precision	Recall
EV-DB	9.2%	92.3%	17.2%	88.2%
CSO-DB	1.8%	90.1%	3.6%	84.6%
PSS-DB	2.3%	98.2%	5.4%	82.7%
PSO-DB	1.4%	88.3%	12.3%	52.5%
PSM-DB	1.7%	97.0%	15.7%	53.1%

Table 1. Text mining performance Results

stage of its existence. Integrate with auditability and policy enforcement tools and data governance frameworks.

To enable data stewards, analysts, and business users to engage with the metadata and obtain insights, provide user-friendly dashboards and interfaces. Create individualized dashboards, reports, and visualizations to display KPIs and the outcomes of the metadata analysis. Provide self-service analytics features so that customers may create ad hoc queries and examine metadata. Create a fault-tolerant, horizontally scalable infrastructure that can manage growing amounts of metadata and user requests. For workload isolation and resource optimization, make use of containerization technologies and distributed computing frameworks. To maximize query performance and minimize latency, put data partitioning techniques and caching mechanisms into practice.

Put strong security measures in place to safeguard private and confidential data as well as sensitive metadata. To prevent unwanted access to metadata, enforce authentication procedures, encryption, and access controls. Assure adherence to industry norms and regulations, including CCPA, HIPAA, and GDPR.

Figure 6 shows Base-line Model architecture and figure 7 shows the proposed methodology, telecom companies can effectively mine metadata for Data Stewardship and Quality, enabling them to improve data governance, enhance service quality, and drive business value.

## IV. RESULT & PERFORMANCE EVALUATION

Text mining evaluation is the technique of ascertaining the effectiveness of text mining methods or algorithms, which can be applied to classify text documents, identity named entities, determine topics in documents, perform sentiment analysis, and more. There are a variety of possibilities regarding the kinds of assessment measures and tools that could be applied to compare the efficiency of different techniques of text mining. Accuracy is defined as the ability to achieve the precise number of occurrences in a given class as there are in the actual data out of all events regarded as belonging to the said class. The amount of

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

samples in a class that has been correctly recognized out of all the instances of that class is called recall. Table 1 contains measures that should not change with balanced data columns.

A type of the task, properties of the data and goals of the study specify metrics and strategies for the evaluation. Performance comparison chart of text mining is presented in the figure 8. Thus, the

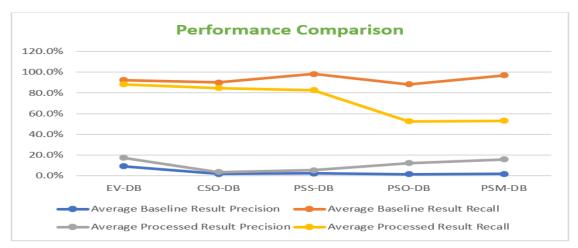


Figure 8: Text mining performance Comparison chart

identification of the indicators used in the assessment of text mining performance is crucial to the effectiveness of the evaluation.

#### V. CONCLUSION

In conclusion, the case of metadata data mining shows that it has a lot of potential for the telecom industry. Thus, it can be useful for strategic decisions, consumer experience, and operation by offering insights and opportunities. Telecom firms can obtain actionable intelligence in the following ways by examining metadata linked to their data, such as call detail records, network logs, customer interactions and service consumption patterns:

- 1. Network Optimization.
- 2. Customer Experience Enhancement.
- 3. Fraud Security and Detection.
- 4. Analytics that predict.
- 5. Effectiveness of Operations.
- 6. Market Intelligence.
- 7. Regulatory Compliance.
- 8. Service Innovation.

Thus, telecom businesses have many opportunities to maintain competitive advantages in the conditions of a constantly evolving industrial environment, develop operational activities, enhance the quality of the customer's interactions, and use metadata mining to derive valuable information. Through complex evaluation techniques and instruments, the telecommunication firms are capable, by using the metadata, to open a new vista of prospects and creativity.

Considering the fact that telecom networks expand rapidly as well as an increased amount of generated data, the future of data mining for the metadata among telecom networks, devices, and applications looks favorable. A large number of data generated through M2M communication and IoT devices will be available as organizations with the Internet of Things (IoT) services evolve. This metadata can be subjected to data mining techniques in order to extract new knowledge for variety of purposes such as industrial process control, smart car and smart city.

#### REFERENCES

1. K. A. Vidhya, G. Aghila, "Text mining process, techniques and tools: an overview", International Journal of Information Technology and Knowledge Management, vol. 2, no. 2, pp. 613-622, 2010.

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

- 2. S.A. Hannan, J. Ahmed, N. Ahmed, R. A. Thakur, "Data Mining and Natural Language Processing Methods for Extracting Opinions from Customer Reviews", International Journal of Computational Intelligence and Information Security, pp.52-58, 2012, ISSN: 1837-7823
- 3. Kaur, Jasleen, J. kumar R. Saini. "A study of text classification natural language processing algorithms for Indian languages", VNSGU J Sci Technol 4.1, pp. 162-167, 2015
- 4. B. Sharma, D. Koundal, "Cattle health monitoring system using wireless sensor network: A survey from innovation perspective", IET Wireless Sensor Systems, vol.8, no.4, pp.143-151, 2018, doi:10.1049/iet-wss.2017.0060
- 5. N. I. Widiastuti, "Deep learning-now and next in text mining and natural language processing", In IOP Conference Series: Materials Science and Engineering, vol.407, no.1, pp.012114, 2018, doi:10.1088/1757-899X/407/1/012114
- 6. Widiastuti, N.I., "Convolution neural network for text mining and natural language processing", In IOP Conference Series: Materials Science and Engineering, vol. 662, no. 5, pp. 052010, 2019, doi:10.1088/1757-899X/662/5/052010
- $7. \quad O.\ Azeroual,\ "Text\ and\ data\ quality\ mining\ in\ CRIS",\ Information,\ vol. 10,\ no. 12,\ pp. 374,\ 2019,\ doi: 10.3390/info 10120374$
- 8. Z. Kong, Y. Cui, Z. Xia, and H. Lv, "Convolution and long short-term memory hybrid deep neural networks for remaining useful life prognostics", Appl. Sci., vol. 9, no. 19, pp. 4156, 2019, doi:10.3390/app9194156.
- 9. S.R. Nayak, J. Mishra, G. Palai, "Analysing roughness of surface through fractal dimension: A review", Image and Vision Computing, vol. 89, no. 21, pp. 34-59, 2019, doi:10.1016/j.imavis.2019.06.015
- 10. M. Bach, A.Werner, and M. Palt, "The proposal of under sampling method for learning from imbalanced datasets", Proc. Comput. Sci., vol. 159, pp. 125–134, 2019, doi: 10.1016/j.procs.2019.09.167.
- 11. K. Verma, S. Bhardwaj, R. Arya, M. S. Islam, M. Bhushan, A. Kumar, P. Samant, "Latest tools for data mining and machine learning", International Journal of Innovative Technology and Exploring Engineering, vol.8, no.9, pp.18-23, 2019, doi:10.35940/ijitee.I1003.0789S19
- 12. H. Kaur, H. S. Pannu, and A. K. Malhi, "A systematic review on imbalanced data challenges in machine learning: Applications and solutions", ACM Comput. Surv., vol. 52, no. 4, pp. 1-36, 2019, doi:10.1145/3343440.
- 13. K. P. Sinaga and M. S. Yang, "Unsupervised K-means clustering algorithm," IEEE Access, vol. 8, pp. 80716–80727, 2020, doi: 10.1109/ACCESS.2020.2988796.
- 14. A. S. Imran, S. M. Daudpota, Z. Kastrati, and R. Batra, "Cross-cultural polarity and emotion detection using sentiment analysis and deep learning on COVID-19 related tweets", IEEE Access, vol. 8, pp. 181074–181090, 2020, doi: 10.1109/ACCESS.2020.3027350.
- 15. P. Kumar, R. Bhatnagar, K. Gaur, and A. Bhatnagar, "Classification of imbalanced data: Review of methods and applications", IOP Conf. Ser., Mater. Sci. Eng., vol. 1099, no. 1, 2021, doi:10.1088/1757-899X/1099/1/012077.
- 16. D. Gupta, S. Wadhwa and S. Rani, "On the Role of Named Data Networking for IoT Content Distribution," International Conference on Communication and Electronics Systems (ICCES), Coimbatre, India, pp. 544-549, 2021, doi: 10.1109/ICCES51350.2021.9488946.
- 17. Tanasescu, L.G. Vines, A. Bologa, A.R. Vaida, "Big Data ETL Process and Its Impact on Text Mining Analysis for Employees", Reviews. Appl. Sci., vol. 12, pp. 7509, 2022, doi:10.3390/app12157509
- 18. K. Thakur, V. Kumar, "Application of text mining techniques on scholarly research articles: Methods and tools", New Review of Academic Librarianship, vol. 28, no. 3, pp.279-302, 2022, doi:10.1080/13614533.2021.1918190
- 19. D. Papakyriakou, I. Barbounakis, "Data mining methods: a review", International journal of computer application, vol. 183, no. 48, pp. 5-19. 2022, doi: 10.5120/ijca2022921884.
- 20. S. Raj, M. Paliwal, "Higher education dashboard implementation using data mining and data warehouse: a review paper", International journal of innovative research in computer science & technology, vol. 10, no. 1, pp. 107-111, 2022, doi: 10.55524ijircst.2022.10.1.19.
- 21. A. Field, C. Y. Park, A. Theophilo, J. Watson-Daniels, and Y. Tsvetkov, "An analysis of emotions and the prominence of positivity in #BlackLives-Matter tweets", Proc. Nat. Acad. Sci. USA, vol. 119, no. 35, Aug. 2022, doi: 10.1073/pnas.2205767119. 22. Q. X. Ng, C. E. Yau, Y. L. Lim, L. K. T. Wong, and T. M. Liew, "Public sentiment on the global outbreak of monkeypox: An unsupervised machine learning analysis of 352,182 Twitter posts", Public Health, vol. 213, pp. 1-4, 2022, doi: 10.1016/j.puhe.2022.09.008.
- 23. A. Kumari, N. Bohra, P. Sangwan, D. Sheoran, S. Kumar, "Data Quality Issues and Metadata Repository of Data Warehouse", An International Interdisciplinary Journal, vol. 01, no.01, pp. 40-51, 2023
- 24. E. Ash, S. Hansen, "Text algorithms in economics", Annual Review of Economics, vol. 15, pp. 659-688, 2023, doi:10.1146/annurev-economics-082222-074352
- 25. R. Olusegun, T. Oladunni, H. Audu, Y. A. O. Houkpati, S. Bengesi, "Text mining and emotion classification on monkeypox Twitter dataset: A deep learning-natural language processing (NLP) approach", IEEE Access, vol. 11, pp. 49882-49894, 2023, doi:10.1109/ACCESS.2023.3277868
- 26. B. G. Chaudhuri, S. Rani, "Future's Backbone Network Monitoring with Metadata in Data Warehouse for Telecom Industry", IEEE, International Conference on Computational Intelligence and Sustainable Engineering Solutions (CISES), pp. 949-954, 2023, doi:10.1109/CISES58720.2023.10183556
- 27. B. G. Chaudhuri, S. Rani, "Managing Metadata in Data Warehouse for Data Quality and Data Stewardship in Telecom Industry-A Compact Survey", IEEE, International Conference on Computing, Communication, and Intelligent Systems (ICCCIS), pp. 353-361. 2023,doi:10.1109/ICCCIS60361.2023.10425001
- 28. S. A. Mohammad, S. Peter, "From simulation to dissemination: automation of data and metadata management", IOP conference series: Earth and environmental science, Bucharest, Romania, vol. 1136, no. 1, pp. 12-18, 2023, doi:10.1088/1755-1315/1136/1/012006.

ISSN: 2229-7359 Vol. 11 No. 18s, 2025

https://theaspd.com/index.php

- 29. S. Digdo, H. A. Alam, R. Nirwantonoa, "Literature study of stunting supplementation in Indonesian utilizing text mining approach", Procedia computer science, Jakarta 11480, Indonesia, vol. 216, pp. 722–729, 2023, doi: 10.1016/j.procs.2022.12.189.
- 30. S. C. Emmadi, P. Agrawal, S. Samudrala, V. Shimpi and M. Natu, "Theory meets practice approach for Event Correlations," IEEE international conference on big data, Sorrento, Italy, vol. 23, no. 7, pp. 1918-1921, 2023, doi: 10.1109/BigData59044.2023.10386269.
- 31. A. Kaur, M. Usama, N. Majid, N. Maarop, "Literature review on metadata governance", Open international journal of informatics. vol. 11, no. 1, pp. 114-120, 2023, doi:10.11113/oiji2023.11n1.235.
- 32. X. Chai, S. Xu, S. Li, J. Zhao, "The Process and Algorithm Analysis of Text Mining System Based on Artificial Intelligence", Procedia Computer Science, vol.228, pp.574-581, 2023, doi.10.1016/j.procs.2023.11.066.
- 33. J. Myllylahti, "A patented method for active data warehousing", 2024, ISBN 978-951-29-9640-7.