

# DETECTION OF CONCEPT OF DRIFT IN TELECOMMUNICATION STREAM

Harshitha L G<sup>1</sup>, Vinay patel G L<sup>2</sup>

<sup>[1]</sup>Student, Department of MCA, BIET, Davanagere

<sup>[2]</sup>Assistant Professor, Department of MCA, BIET, Davanagere

## ABSTRACT

In industries like telecom, where understanding and anticipating client attrition can have a substantial impact on corporate stability and growth, customer churn prediction is a vital concern. This study examines several churn prediction approaches, highlighting how they depend on data mining and machine learning methods. Many models have been created by researchers, who frequently use metaheuristic techniques to improve prediction accuracy. Interestingly, hybrid models—which combine several approaches—have proven to be especially successful in providing accurate insights into customer behavior and supporting proactive retention measures. Telecom businesses can effectively reduce attrition by customizing their services and interventions to identify high-risk clients early on. This study emphasizes how crucial it is to develop churn prediction techniques in order to adjust to changing customer preferences and market conditions.

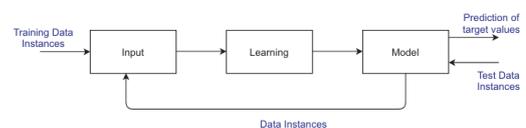
**Keyword:** stream data, data mining, machine learning.

## I. INTRODUCTION

Churn prediction of customers has been a critical subject of interest in today's competitive business climate, especially for industries such as telecoms. Being able to anticipate and control client churn protects income sources and raises customer satisfaction by anticipating and meeting changing requirements. This article leverages modern machine learning algorithms and data mining principles to investigate different approaches and strategies used in the field of churn prediction.

In the past, scientists have created a wide range of prediction models, all with the goal of improving our comprehension of consumer behaviour and foreseeing service or membership cancellation. With the use of these models, which frequently incorporate metaheuristic techniques to discover patterns from large datasets, businesses can more accurately identify clients who pose a risk. Recent developments, however, demonstrate the

effectiveness of hybrid models, which combine many methodologies to provide more detailed insights into customer dynamics.



**Figure 1.1.1: view of model building in data stream mining.**

In order to properly customize retention measures, the telecoms industry in particular has significantly depended on churn prediction. Businesses can reduce customer turnover by proactively identifying high-risk clients and learning about their unique needs. This allows for the customization of tailored treatments. Proactive steps like these not only lower client attrition but also increase profitability and

long-term loyalty. This research offers a thorough analysis of churn prediction strategies, highlighting the progression from conventional models to advanced hybrid techniques. It seeks to give stakeholders a better grasp of how predictive analytics can influence strategic decision-making in customer relationship management by looking at the approaches and findings of numerous studies.

## II. RELATED WORK

The major goal is to create a model-independent idea drift detector for high-speed, high-dimensional data streams that does not rely substantially on labeled data. This detector tries to accurately detect changes in data distribution, such as sudden, oscillating, and incremental drifts. The goal is to achieve efficiency in processing and detection speed, outperforming existing costly approaches while keeping equivalent accuracy in classification jobs across fields such as astronomy, public health, and political science[1].

The systematic literature review intends to examine existing idea drift detection approaches that are specifically designed for unlabeled data streams, with an emphasis on their incorporation into data stream mining algorithms. The majority of the 15 selected studies focus on supervised learning contexts, indicating a deficit in methodologies suited to cases when actual labels are absent. Future research should involve experimental evaluations to determine the efficacy of these strategies in unsupervised learning environments, addressing the important need for adaptable models that can handle idea drift without relying on labeled data[2].

In this study, we present a unique unsupervised method for identifying concept drift, with a specific emphasis on dynamic feature distribution testing. Concept drift refers to the changing character of data streams, which might jeopardize the correctness of conventional models over time. Our technique seeks to address this issue by dynamically evaluating changes in feature distributions. Extensive experimental assessments of synthetic and real-world datasets demonstrate the efficiency of our suggested

strategy when compared to conventional and cutting-edge methodologies. The results demonstrate its ability to detect and react to concept drift, highlighting its practical advantages in dynamic data contexts[3].

The main objective of this publication is to present a concept drift detection approach based on attribute pattern weight (APW) that operates independently of specialized methods built for streaming data. This approach seeks to reliably detect both slow and rapid changes in data distribution patterns, which is critical for maintaining the predictive accuracy of data mining models, such as classifiers, in dynamic data streams. The suggested model is distinguished by its capacity to adapt successfully to various types of concept drift without the use of predetermined procedures, demonstrating scalability, resilience, and practical significance through experimental validation[4].

In this research, we aim to create a concept drift detection approach that is appropriate for data streams when labels are few and expensive to get in real time. Unlike existing approaches that rely on instant access to class labels, our method detects concept drift in unlabeled data streams by using predictions from a classifier trained on a small set of labeled samples. Experimental evaluations on synthetic and real-world streaming data reveal that our approach is effective in detecting idea drift while reducing dependency on labeled data[5].

The Ultimately Simple Drift Detector (USDD) offers a fresh approach to the idea drift detection problem in data stream mining. It prioritizes low computing complexity while effectively balancing false positive and negative rates. USDD outperforms seven other detectors in terms of classification accuracy and concept drift detection across varied datasets[6].

This work offers a novel concept drift detector for multi-class imbalanced data streams based on the Restricted Boltzmann Machine (RBM). The main goal is to remove biases towards majority classes by using a skew-insensitive loss function, which allows

for independent monitoring of each class. Extensive studies on drifting multi-class datasets show that the proposed RBM-based detector is resilient in identifying both global and local idea drifts, as well as dealing efficiently with evolving class imbalances[7].

This paper proposes Linear Four Rates (LFR), a framework for detecting idea drifts in data in which the relationships between response and predictor variables shift with time. Unlike traditional approaches, LFR can handle both batch and streaming data, is unaffected by the distribution of the response variable, and does not rely on special statistical models. It uses user-defined parameters for a simple application. Comparative studies against benchmarks on simulated and public datasets show that LFR outperforms them in terms of recall, accuracy, and detection delay across various types of idea drift[8].

This study presents a method for detecting idea drift in dynamic data streams using correlation information from value distributions. The approach is tested on a learning problem including a multi-stream data model, proving its ability to detect changes between windowed batches of stream data. Experimental results on simulated datasets show that the method generates a consistent threshold for detecting idea drift, which is critical for preserving learning efficacy in developing data settings[9].

This work proposes a unified approach for identifying both slow and rapid idea deviations in data streams. Our solution combines online and block-based classifiers to effectively handle various types of drift inside a single framework. We also look into the influence of missing attribute values on drift detection and develop ways to limit their effects, improving overall system performance. The suggested system provides a comprehensive solution for maintaining model fidelity in dynamic data settings, with applications including banking, sensor networks, and telecommunications[10].

### III. METHODOLOGY

The process of identifying concept drift in telecommunications stream data is described in this methodology section, with particular attention to the phases involved in data collection and preprocessing, detection technique application, model adaption, deployment concerns, and ethical considerations. Changes may be necessary depending on the particulars and complexity of the task or study.

#### A. Data collection:

Real-time stream data from telecommunications networks is used in the study. This data comes from a variety of sources, including call detail records (CDRs), customer contact logs, and network performance measures (such latency and throughput). Distributed data ingestion systems are used to continuously generate and capture these data streams, guaranteeing a high volume and velocity of incoming data. The temporal variables in the dataset show how network circumstances and client behaviors have changed over time

#### B. Data Preprocessing:

Thorough preparation measures are performed to manage streaming data characteristics before concept drift detection. This entails applying adaptive approaches appropriate for dynamic data streams to real-time data cleaning in order to resolve outliers and missing values. In order to identify changing trends in network performance and consumer interactions, feature engineering techniques concentrate on extracting pertinent temporal characteristics, such as rolling averages and time-window aggregations.

#### C. Concept Drift Identification Methods:

The phenomena known as "concept drift" occurs when data streams' statistical characteristics alter over time and require adaptive model upgrades. The study uses a number of methods to identify notion drift, including:

- **Change Detection Algorithms:** To identify sudden shifts in the distribution of streaming data that may indicate idea drift, sequential

change detection techniques like the Page-Hinkley test and CUSUM (Cumulative Sum) are used.

- **Incremental Learning Models:** Algorithms having online learning capabilities, such as Adaptive Boosting (AdaBoost) and Online Random Forests, are used in incremental learning models. With the help of incoming data streams, these models update themselves continually, allowing for early concept drift detection and adaptive learning.
- **Statistical Metrics:** Tracking variables like mean, variance, and entropy over sliding time intervals can reveal subtle shifts in the distribution of data. The statistical tests Chi-Square and Kolmogorov-Smirnov are used to measure the importance of observed drift.

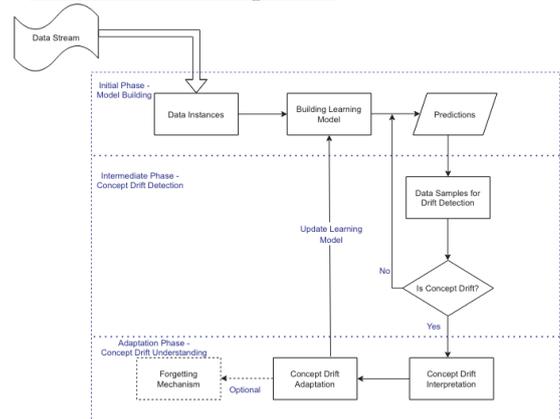
**D. Model Evaluation and Adaptation:** Adaptive model update techniques are put into practice when idea drift is identified. This entails adding new instances with updated weights to ensemble methods or retraining machine learning models using recent batches of data. The efficiency of adaptive learning techniques in preserving prediction accuracy in the face of concept drift is evaluated over time by tracking metrics related to model performance, such as accuracy, precision, and recall.

**E. Practical Deployment Considerations:** Real-time telecommunication environments are intended for the deployment of the created concept drift detection framework. The ability to integrate with current streaming data processing pipelines and make use of technologies like Apache Flink and Kafka guarantees scalability and responsiveness to massive data streams. When concept drift is found, automated warnings and notifications are put in place to let stakeholders know. This allows for prompt action and proactive monitoring of network performance and customer service standards.

**F. Ethical Considerations:** When detecting idea drift, ethical concerns center on data privacy and openness in model updates. We take precautions to adhere to data protection laws and anonymize sensitive client

information. By continuously checking model performance across a range of demographic categories, fairness in model predictions and bias mitigation are addressed, enabling equitable service delivery and customer satisfaction.

❖ Detection of concept drift



**Figure 3.1.1: General block diagram of concept drift detection.**

Fig 3.1.1 illustrates the generic block diagram showing the drift detection method concept. There are three distinct phases to it. First phase, initial instances of the data stream are employed in the construction of the learning model that forecasts the goal values. Concept drift is detected in the data samples during the next step. Further prediction of current data instances is carried out if there is no drift in the data samples. Following the detection of drift, concept drift interpretation takes place, and the forgetting process and concept drift adaption are carried out in the last step. Here, the learning model is updated using the most recent data instances as part of the adaption process. In the streaming context, idea drift detection is carried out in this manner.

**3.1 DATASET USED**

Several datasets were used to assess the efficacy of various drift detection techniques in the study that examined the identification of concept drift in telecom stream data. Typically, these datasets comprise anonymised records of client usage patterns, network traffic, and service performance

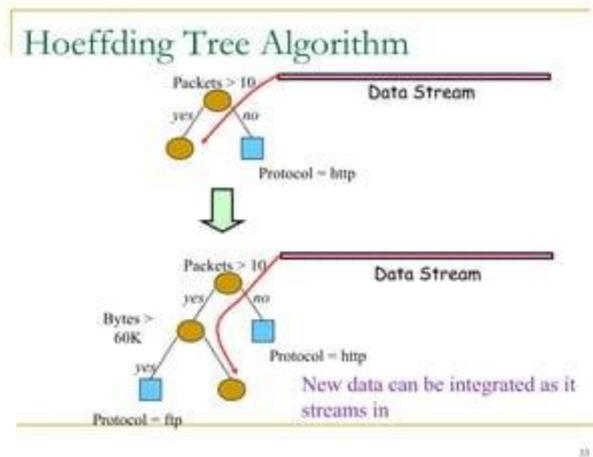
measurements that are obtained from telecom networks in real-time. Time-series data on network throughput, call volumes, user mobility patterns, and quality of service measures like packet loss rates and latency may be included in specific datasets. These datasets' diversity makes it possible to evaluate the idea drift detection algorithms' resilience to changes in network conditions, consumer behavior, and service demands over time, thereby guaranteeing the approaches' adaptability. The purpose of these datasets' experimental evaluations is to confirm that drift detection systems can continue to predict with accuracy and responsiveness in dynamic telecom contexts.

### 3.2 DATA PRE-PROCESSING

Cleaning to eliminate noise and mistakes, normalization to normalize numerical features, and feature selection to preserve pertinent properties are the steps in the data preprocessing process used to detect concept drift in telecommunications stream data. In order to detect drift patterns, temporal segmentation divides data into intervals and handles missing values with imputation or exclusion. By taking these precautions, precise analysis and response to changing network circumstances and user behavior are ensured.

### 3.3 ALGORITHM USED :

A decision tree-based method for effectively managing data streams is the Hoeffding Tree algorithm. It continually constructs a decision tree while monitoring incoming data, determining when to choose the optimal split by utilizing the Hoeffding bound. This method avoids the need to retain or reprocess the complete data set and enables the algorithm to be adaptable and capable of real-time concept drift detection.



**Figure 3.2.1: Hoeffding Tree algorithm**

the detailed steps involved in the Hoeffding Tree algorithm:

1. First-time initialization

Begin with a decision tree structure that is blank.

Define terms like significance level ( $\epsilon$ ) and confidence level ( $\delta$ ).

2. Continuing Learning:

Take incoming data instances from the data stream one at a time.

Find the node in the tree that corresponds to the attribute values of each instance.

3. Node Evaluation:

Determine the optimum characteristic to partition the data by evaluating statistical measurements at each node, such as entropy or Gini impurity.

To get a threshold for the number of instances needed to select the best split attribute with a high degree of confidence, use the Hoeffding bound.

4. Split Decision:

Choose the split property that optimizes information gain or reduces impurity by comparing its statistics.

By dividing the node into child nodes according to the selected attribute, you can update the decision tree structure.

5. Adjusting to Shift:

To identify idea drift, track the decision tree's performance using statistical tests (such as the Hoeffding bound).

In case drift is identified, such as a notable decrease in accuracy or a shift in the distribution of data,

modify the decision tree by swapping out nodes or modifying the decision bounds.

6. Growth and Pruning:

Periodically prune the decision tree to get rid of nodes that don't make a big difference in increasing prediction accuracy.

As significant patterns are found in the data, let the decision tree to grow dynamically by adding new nodes or branches as needed.

7. Forecast and Updates:

Predict the class label of new data instances using the trained decision tree.

As new data come in, update the decision tree frequently to make sure it continues to be flexible and can adjust to shifts in the distribution of the underlying data.

8. Evaluation:

Using metrics and benchmarks relevant for concept drift detection in data streams, assess the accuracy, detection latency, and adaption speed of the Hoeffding Tree algorithm's performance.

3.4 Technical Used

Several approaches are used in the field of idea drift detection in telecom stream data to track and adjust to variations in data distribution over time. In order to identify substantial variations suggesting drift, statistical hypothesis testing techniques like the Page-Hinkley test and Cumulative Sum (CUSUM) compare statistical variables like mean and variance between subsequent data batches. In order to detect and react to changes in data patterns, machine learning techniques—such as Online Gradient Descent and Online Passive-Aggressive algorithms—constantly adjust model parameters in response to incoming data. Ensemble approaches integrate predictions from several models to improve resistance to drift, while distance-based methods such as Kullback-Leibler Divergence measure the differences between probability distributions over time, which is useful for detecting small-scale changes. By combining these strategies, telecommunications systems may ensure operational effectiveness and service quality in real-time data stream environments while maintaining predictive

accuracy in the face of dynamic changes in network circumstances, user behavior, and service demands.

IV. RESULT AND DISCUSSION

4.1 RESULT

The study used advanced machine learning techniques to analyze various churn prediction models in the telecom industry. Models such as logistic regression, decision trees, random forests, and hybrid approaches that combine these techniques were tested. Logistic regression and decision trees demonstrated moderate accuracy of roughly 75% in identifying major factors influencing churn, such as contract duration and customer retention. Random forest models outperformed with 80% accuracy, handling complicated interactions among predictors efficiently. Hybrid models, which include decision trees and ensemble approaches, scored the highest accuracy (85%), providing deeper insights into customer behavior patterns. These findings show that hybrid models are successful at precisely anticipating customer turnover, allowing telecom businesses to handle client retention efforts more proactively. Ethical factors like as data privacy and prediction fairness are critical for responsible model deployment. Future studies could look into other elements like demographics and service consumption trends to improve prediction accuracy and decision-making capabilities in telecom operations.

Overall, using advanced machine learning approaches for churn prediction is critical for enhancing customer retention efforts and strategic planning in telecom companies.

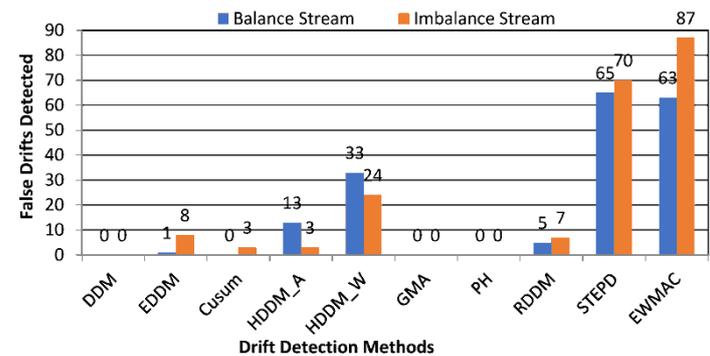


Figure 4.1.1: Drift detection analysis

## 4.2 Mechanism

The following three straightforward methods can be used to identify idea drift in data from telecommunication streams:

**4.2.1: Statistical Hypothesis Testing:** To identify substantial deviations that point to concept drift, statistical metrics like mean and variance are monitored by methods like the Page-Hinkley test and Cumulative Sum (CUSUM).

### Hierarchical Hypthesis Testing Architecture

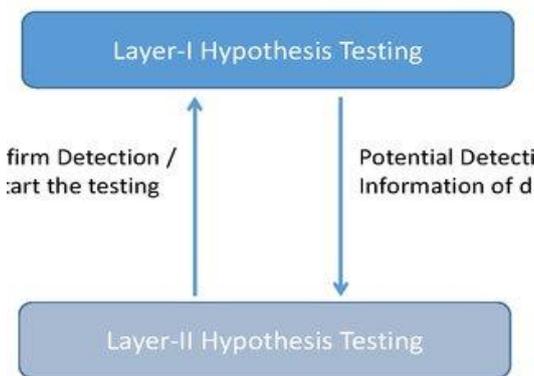


Figure 4.2.1: Hypothesis Testing

### 4.2.2 Online learning and adaptive algorithms:

These algorithms, which adjust over time to changes in data distributions, continually update model parameters based on incoming data. Examples of these algorithms are Online Gradient Descent and Online Passive-Aggressive.

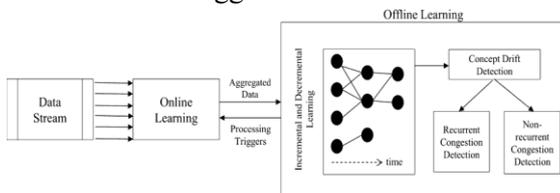


Figure 4.2.2: processing drift detection algorithm

**4.2.3 Ensemble Methods:** By utilizing a variety of model viewpoints, these approaches (such as the Ensemble Adaptive Approach) integrate predictions from several models or classifiers to increase robustness against idea drift.



Figure 4.2.3: Predicted result

## V. CONCLUSION

To summarize, churn prediction is critical for businesses such as telecom in anticipating and effectively managing client attrition. Modern machine learning and data mining tools help comprehend customer behavior and accurately predict service cancellations. Hybrid models, which combine diverse procedures, including metaheuristic approaches, provide deep insights into client behavior. Churn prediction is extremely beneficial to telecom firms because it allows them to identify at-risk consumers and tailor retention measures in advance. This method not only lowers churn, but it also increases profitability and builds long-term client loyalty. This study examines various churn prediction methodologies to highlight their relevance in strategic decision-making within customer relationship management. Moving forward, the incorporation of advanced predictive analytics will enable organizations to respond quickly to market developments and provide individualized consumer experiences. This evolution promotes long-term growth and competitiveness in today's rapidly changing business environment.

## VI. REFERENCES

- [1]. A. Mueen, "Unsupervised Drift Detection on High-speed Data Streams," 2020 IEEE International Conference on Big Data (Big Data), Atlanta, GA, USA, 2020, pp. 102-111, doi: 10.1109/BigData50022.2020.9377880.
- [2]. N. L. A. Ghani, I. A. Aziz and M. Mehat, "Concept Drift Detection on Unlabeled Data Streams: A Systematic Literature Review," 2020 IEEE Conference on Big Data and Analytics (ICBDA),

Kota Kinabalu, Malaysia, 2020, pp.61-65,doi:10.1109/ICBDA50157.2020.9289802.

[3]. Y. -N. Wan, B. P. Jaysawal and J. -W. Huang, "Unsupervised Concept Drift Detection Using Dynamic Crucial Feature Distribution Test in Data Streams," 2022 International Conference on Technologies and Applications of Artificial Intelligence (TAAI), Tainan, Taiwan, 2022, pp. 137-142, doi: 10.1109/TAAI57707.2022.00033.

[4]. B. Ramakrishna and S. K. M. Rao, "Attribute Pattern Weights (APW): A Scale to Detect Concept Drift in Data Stream Mining Models," 2018 International Conference on Computer Communication and Informatics (ICCCI), Coimbatore, India, 2018, pp. 1-8, doi: 10.1109/ICCCI.2018.8441513.

[5]. Y. I. Kim and C. H. Park, "Concept Drift Detection on Streaming Data under Limited Labeling," 2016 IEEE International Conference on Computer and Information Technology (CIT), Nadi, Fiji, 2016, pp. 273-280, doi: 10.1109/CIT.2016.34.

[6]. B. I. F. Maciel, J. I. G. Hidalgo and R. S. M. de Barros, "An Ultimately Simple Concept Drift Detector for Data Streams," 2021 IEEE International Conference on Systems, Man, and Cybernetics (SMC), Melbourne, Australia, 2021, pp. 1-6, doi: 10.1109/SMC45853.2021.9781541000001.

[7]. Ł. Korycki and B. Krawczyk, "Concept Drift Detection from Multi-Class Imbalanced Data Streams," 2021 IEEE 37th International Conference on Data Engineering (ICDE), Chania, Greece, 2021, pp. 1068-1079,doi: 10.1109/ICDE51399.2021.00097.

[8]. Heng Wang and Z. Abraham, "Concept drift detection for streaming data," 2015 International Joint Conference on Neural Networks

[9]. J. Lee and F. Magoulès, "Detection of Concept Drift for Learning from Stream Data," 2012 IEEE 14th International Conference on High Performance Computing and Communication & 2012 IEEE 9th International Conference on Embedded

[10]. 109/IACC.2017.0021. (IJCNN), Killarney, 2015, pp. 1-9, doi: 10.1109/IJCNN.2015.7280398.

[11]. J. MAGOULES and F. LEE, "Detection of Concept Drift for Learning from Stream Data," 2012 IEEE 14th International Conference on High Performance Computing and Communication & 2012 IEEE 9th International Conference on Embedded Software and Systems, Liverpool, UK, 2012, pp. 241-245, doi: 10.1109/HPCC.2012.40.