Impact of Data Quality on Machine Learning Models in Telecom and Media

Mahesh Mokale

Independent Researcher maheshmokale.mm@gmail.com

Abstract

In the digital age, telecom and media companies increasingly rely on machine learning (ML) to drive innovation, enhance operational efficiency, and deliver personalized customer experiences. These industries generate and consume vast volumes of data daily-from customer interactions, usage logs, and social media behavior to network sensor outputs and streaming analytics. This data, when harnessed effectively, enables a broad range of ML applications such as churn prediction, fraud detection, targeted advertising, network optimization, and personalized content recommendation. However, the efficacy of ML models is fundamentally tied to the quality of the data feeding them. Poor data quality-manifesting as inaccuracies, incompleteness, inconsistencies, latency, irrelevance, and bias-can significantly degrade model performance. Models trained on flawed datasets are more likely to produce skewed or misleading outputs, leading to erroneous insights, misinformed strategies, wasted resources, and ultimately, dissatisfied customers. In highstakes environments such as telecom and media, where decisions derived from ML insights affect millions of users in real time, the risks associated with poor data quality are amplified. This white paper explores the multifaceted impact of data quality on ML models within telecom and media, beginning with a detailed analysis of the types of data typically encountered in these sectors and the specific challenges they present. It highlights how subpar data quality can introduce systemic bias, reduce model generalizability, increase error rates, and lower overall system reliability and trustworthiness. Furthermore, the paper outlines common data quality issues unique to these industries, including duplicate records from multiple data sources, inconsistent data formats, imbalanced usage data, and outdated streaming or log data. In addition to identifying these challenges, the paper presents real-world case studies demonstrating the quantifiable benefits of data cleaning and preprocessing. It details how organizations improved ML performance metrics and customer satisfaction by addressing core data quality issues. Moreover, it recommends actionable strategies and modern toolsets to ensure robust data pipelines that support scalable and trustworthy ML models. Ultimately, this paper underscores that data quality is not merely a technical hygiene practice but a strategic imperative for companies aiming to compete effectively in the evolving telecom and media landscapes.

Keywords: Data Quality, Machine Learning, Telecom, Media, Data Governance, Data Pipelines, Predictive Analytics, Churn Prediction, Recommendation Systems, Data Profiling, ETL, Metadata, Data Validation, Data Cleaning, Data Quality Tools, Big Data, Data Bias, Model Accuracy, Data Consistency, Fraud Detection, Data Lineage, Data Completeness, Customer Experience, Model Robustness, Unstructured Data, Training Data Versioning, Compliance, Anomaly Detection, Data Catalog, Streaming Data, Real-Time Processing

1. Introduction

The telecommunications and media industries are at the forefront of the digital revolution, characterized by their dynamic nature and an ever-growing dependency on data-driven decision-making. These sectors collect and process massive amounts of data from a variety of sources including call detail records, network sensors, content delivery platforms, mobile applications, social media channels, customer service interactions, and third-party integrations. This data spans structured formats such as transactional logs and user profiles, as well as unstructured formats including video content, voice data, and text-based communications.

Machine learning has emerged as a critical technology in transforming this raw data into actionable insights. Applications range from predictive maintenance of network infrastructure and real-time fraud detection to personalized marketing, automated customer support, and adaptive content recommendation systems. With increasing consumer demand for seamless connectivity and personalized experiences, telecom and media companies are under constant pressure to deploy intelligent systems that are accurate, responsive, and scalable.

However, the effectiveness of machine learning systems is heavily contingent on the quality of data being ingested. Data quality encompasses various dimensions—such as accuracy, completeness, consistency, timeliness, validity, and relevancy—which collectively determine the reliability and performance of ML models. Poor-quality data can propagate errors across the pipeline, leading to underperforming models, increased operational costs, and poor customer outcomes.

In this context, data quality is not merely a technical consideration; it is a foundational requirement for maintaining competitive advantage. Organizations must not only capture vast quantities of data but also ensure its quality throughout the data lifecycle—from ingestion and transformation to storage, analysis, and model deployment. This paper aims to provide a comprehensive understanding of how data quality impacts ML implementations in telecom and media, and to equip decision-makers with the knowledge required to prioritize data quality as a strategic asset.

2. Importance of Data Quality in ML Pipelines

High-quality data is the cornerstone of reliable and scalable machine learning systems. It directly influences every phase of the ML pipeline—from data preprocessing and feature engineering to model training, validation, deployment, and monitoring. Ensuring high-quality data is not simply about avoiding errors; it's about optimizing the entire decision-making ecosystem that relies on predictive accuracy and adaptability.

- Accurate Predictions: Machine learning models are only as good as the data they are trained on. Clean, well-labeled, and representative datasets reduce noise and bias, leading to predictions that closely reflect real-world behavior. In telecom, this can mean more precise customer churn predictions; in media, it results in more relevant content recommendations.
- **Generalizable Models:** High-quality data enables models to capture underlying patterns that generalize well to unseen data. When datasets are free from duplicates, inconsistencies, and imbalance, models avoid overfitting and perform reliably across varied scenarios and user groups.

- **Faster Convergence During Training:** Noisy or redundant data can prolong training cycles and lead to unstable learning processes. Clean, well-structured data accelerates convergence, reduces computational overhead, and minimizes the need for repeated experimentation or hyperparameter tuning.
- **Reduced Operational Risk:** Poor data can introduce vulnerabilities in production systems, including false alarms, missed detections, or biased recommendations. High-quality data lowers these risks by supporting consistent and interpretable outputs, which are crucial in regulatory-heavy sectors like telecom.

In telecom, for example, inaccurate or outdated churn labels might train a model to incorrectly prioritize low-risk customers for retention campaigns, wasting resources. In the media sector, recommendation engines trained on incomplete or misclassified user behavior may serve irrelevant content, leading to lower user engagement and satisfaction.

Ultimately, high-quality data improves transparency, fairness, and resilience in ML systems. It empowers organizations to derive trustworthy insights, adapt to evolving user behavior, and maintain a competitive edge in fast-paced, customer-centric markets.

3. Common Data Quality Issues in Telecom and Media

The telecom and media industries face several unique and recurring data quality issues that can compromise the reliability and effectiveness of machine learning models. These issues arise due to the complex nature of data ingestion from heterogeneous sources, the volume and velocity of data flow, and the limitations of legacy systems and fragmented data architectures. The most prevalent issues include:

- **Missing Data:** Many datasets, particularly those generated from user behavior, sensor logs, and customer profiles, often contain null, incomplete, or unavailable values. This can result from network outages, sensor malfunctions, or poor data collection protocols. In ML applications, missing data can lead to biased estimations or degraded model accuracy if not handled appropriately through imputation or exclusion.
- **Inconsistent Formats:** Data sourced from different systems frequently use varied formats for fields like dates, geographic locations, currency, or device types. Lack of standardization complicates data integration and transformation processes. For example, a single customer's data might show different time zones, naming conventions, or device identifiers across platforms, leading to misalignment in feature engineering.
- **Duplicate Records:** Duplicate entries can arise when customer information is collected multiple times from different channels (e.g., mobile apps, retail stores, call centers) or during system migrations and integrations. Duplicates inflate dataset size, distort frequency-based features, and create artificial patterns that mislead model training.
- **Outdated Data:** Latency in data pipelines—due to slow ingestion, storage delays, or batch processing—can result in models being trained on stale or obsolete information. For instance, recommending a media title based on weeks-old viewing behavior or predicting churn from last quarter's call logs can severely affect user engagement and decision accuracy.

- **Bias and Imbalance:** Telecom and media datasets often reflect historical usage patterns that overrepresent dominant user groups or content types. For example, a recommendation model trained predominantly on urban viewer behavior may underperform for rural users. Imbalanced datasets also affect classification tasks like fraud detection, where the number of fraudulent cases is much smaller than legitimate ones.
- Noisy or Unstructured Data: Particularly in media, user-generated content, transcriptions, or social media feeds introduce noise, sarcasm, slang, and ambiguities that challenge NLP and image recognition models. Poor labeling, misclassification, and irrelevant data inclusion add further complexity.

Addressing these issues requires proactive data quality checks and remediation workflows across the data lifecycle. Recognizing and resolving them at early stages significantly enhances ML outcomes and reduces long-term technical debt.

4. Consequences of Poor Data Quality

The downstream effects of poor data quality in machine learning implementations can be severe, particularly in telecom and media sectors where predictive accuracy, real-time insights, and customer satisfaction are paramount. These consequences manifest across technical, operational, regulatory, and business dimensions:

- **Model Degradation:** The most immediate impact is a drop in model performance. Low-quality training data leads to high error rates, overfitting, underfitting, and non-generalizable results. Models may pick up on spurious correlations, fail to capture true patterns, or provide inconsistent outputs in production environments. For example, an ML system predicting network outages may miss critical anomalies due to missing or incorrect sensor data.
- **Operational Failures:** Faulty data inputs can cause critical operational tools powered by ML to malfunction. In fraud detection systems, mislabeled or outdated transaction records can lead to false negatives—letting fraud go undetected—or false positives—disrupting legitimate user activity. Similarly, in network traffic prediction, inaccurate input data can cause poor bandwidth allocation and affect service quality.
- **Compliance and Legal Risks:** Inconsistent or inaccurate data poses compliance challenges, especially when used to make automated decisions affecting customers. Failure to adhere to data governance regulations like GDPR, CCPA, or industry-specific privacy laws can result in fines, legal scrutiny, and reputational damage. For instance, incorrect personalization due to erroneous demographic data may violate consent-based personalization rules.
- **Customer Experience Degradation:** Data-driven personalization and engagement strategies depend on reliable data. Poor-quality data results in irrelevant content recommendations, incorrect billing predictions, or misrouted customer service interactions, which frustrate users and erode brand loyalty. In competitive markets, even slight inaccuracies can lead to churn.
- **Business Decision Misalignment:** Strategic decisions made on the basis of ML insights derived from low-quality data can lead to misallocated budgets, ineffective marketing campaigns, and misguided product development. Decision-makers may draw false confidence from models that are

unknowingly biased or incomplete.

- Loss of Trust in ML Systems: As stakeholders encounter repeated inaccuracies, their trust in AI/ML capabilities diminishes. This resistance can slow down the adoption of ML initiatives and lead organizations to revert to traditional decision-making methods, nullifying investments in data science infrastructure.
- **Increased Technical Debt:** Correcting model failures and reprocessing corrupted data retroactively often requires substantial time, manpower, and reengineering efforts. This accumulated technical debt not only slows innovation but also burdens engineering teams with avoidable maintenance work.

Ultimately, poor data quality acts as a hidden tax on every aspect of machine learning operations. Investing in data quality is not just a best practice—it's a non-negotiable component of responsible and effective AI strategy.

5. Case Studies

Real-world case studies provide compelling evidence of the transformative impact that data quality improvement can have on the performance and reliability of machine learning models. The following examples from telecom and media highlight how organizations have successfully mitigated data quality challenges and achieved measurable improvements:

- **Telecom Churn Prediction at a Tier-1 Operator:** A major telecom provider serving over 50 million customers struggled with high false positive rates in its churn prediction model. The root cause was traced to duplicated customer records, missing demographic data, and inconsistent timestamps across regional databases. By implementing a comprehensive data cleansing initiative—including deduplication, standardization of time fields, and enrichment from external datasets—the operator improved the accuracy of churn classification by 20%. This translated into a 15% increase in retention campaign ROI, as marketing efforts could be better targeted toward genuinely at-risk customers.
- Content Recommendation Optimization at a Global Streaming Service: A leading media streaming company experienced declining engagement metrics due to ineffective content recommendations. An audit revealed that bot traffic had skewed user interaction logs, and metadata inconsistencies across content catalogs disrupted semantic similarity modeling. By filtering out non-human interactions using behavioral thresholds and harmonizing metadata (e.g., genres, tags, actor names), the company recalibrated its recommendation engine. Post-cleanup, click-through rates improved by 18%, and user retention rose by 15% over a six-month period.
- Fraud Detection Enhancement in Mobile Payments: A telecom company offering digital wallet services faced challenges with an ML-based fraud detection system that failed to catch new patterns of fraudulent behavior. Investigation showed that delayed integration of transactional data from partner systems and class imbalance in the training data were the primary issues. By establishing a real-time ETL pipeline and applying synthetic oversampling (SMOTE) to rebalance the training dataset, the fraud detection accuracy improved by 25%, with a substantial reduction in false negatives.

• Ad Delivery Optimization in OTT Media Platforms: An over-the-top (OTT) media platform struggled with low ad relevancy scores due to flawed user segmentation driven by noisy behavioral data. After instituting a multi-layer validation process to clean malformed event logs and enrich them with contextual data (e.g., location, device type, viewing session duration), the company saw a 12% increase in ad click-through rates and a 10% lift in advertiser satisfaction.

These case studies demonstrate that proactive data quality initiatives can significantly enhance the performance of ML models, improve business KPIs, and build trust in AI-driven systems. The return on investment from improving data quality often outweighs the initial effort and resources required for implementation.

6. Best Practices for Ensuring Data Quality

Ensuring high data quality is a continuous, systematic process that must be integrated throughout the data lifecycle. The following best practices are essential for maintaining clean, reliable, and ML-ready datasets in telecom and media environments:

- **Data Profiling and Auditing:** Regularly assess datasets to detect anomalies, outliers, missing values, duplicates, and schema violations. Use statistical summaries, correlation matrices, and data distribution visualizations to uncover quality issues early in the pipeline. Data profiling should be an automated, repeatable task embedded within development and production workflows.
- **ETL Standards and Governance:** Build robust Extract-Transform-Load (ETL) pipelines with strict validation rules, transformation protocols, and logging mechanisms. Data lineage tracking should be enabled to trace the origin, movement, and modification of data. Establish governance frameworks with defined roles and responsibilities for data stewardship, ownership, and compliance.
- Automated Cleaning Pipelines: Incorporate rule-based systems and machine learning models to detect and correct errors in real time. Use regex patterns, fuzzy matching, anomaly detection algorithms, and deduplication engines to maintain consistent and standardized records. Continuous integration (CI) tools should validate data transformations as part of deployment cycles.
- Feedback Loops and Human-in-the-Loop Systems: Engage domain experts, data analysts, and end users to flag suspicious data points and validate automatic corrections. Feedback loops are especially valuable in scenarios involving unstructured or semi-structured data such as speech transcripts, video metadata, and chat logs.
- Metadata and Documentation Management: Maintain comprehensive metadata including data source, data type, validation rules, and refresh frequency. Implement metadata repositories or data catalogs that provide easy access to lineage and quality status. Well-documented data assets improve reproducibility, collaboration, and governance.
- Data Quality KPIs and Dashboards: Define and monitor key performance indicators (KPIs) for data quality such as completeness rates, error counts, freshness thresholds, and consistency scores. Use dashboards and alerting systems to track quality trends over time and enable proactive issue resolution.
- Data Privacy and Compliance Validation: Include checks to ensure that data handling meets

regulatory standards like GDPR and CCPA. This includes masking personally identifiable information (PII), anonymizing sensitive fields, and logging consent status where applicable.

• **Training Data Versioning and Auditing:** Maintain version control over datasets used for model training and testing. Track changes to features, label distributions, and sampling techniques to facilitate model reproducibility, troubleshooting, and rollback in case of data drift or model degradation.

Adopting these best practices fosters a culture of data accountability and ensures that machine learning models operate on reliable, trustworthy inputs. In the telecom and media sectors—where decisions can impact millions of users—these practices are not optional but essential for scalable and ethical AI adoption.

7. Tools and Technologies

A variety of tools and platforms are available to support the monitoring, validation, and enhancement of data quality in machine learning pipelines. These technologies help automate quality assurance processes, provide visibility into data workflows, and support compliance and governance efforts. The following are widely used tools that are particularly effective in the telecom and media domains:

- Apache Griffin: An open-source data quality solution designed for big data environments. Griffin offers data profiling, validation, and measurement frameworks. It supports both batch and streaming data and integrates with Apache Hadoop, Spark, and Hive. Telecom companies use Griffin to maintain the integrity of high-volume sensor logs and transactional data.
- **Great Expectations:** A Python-based open-source framework for defining, executing, and validating data expectations. It supports unit testing for data and integrates with data pipelines built on Airflow, dbt, and Spark. In media workflows, Great Expectations is useful for validating ingestion of metadata, clickstream data, and viewership logs before model training.
- **AWS Deequ:** A library developed by Amazon for defining "unit tests" for data in large-scale processing jobs. It is optimized for use with Apache Spark and supports data profiling, constraint suggestion, and data verification. Telecom analytics teams use Deequ to automate checks on data reliability and consistency before running predictive models.
- **Talend Data Quality:** A comprehensive suite of data preparation, cleansing, and enrichment tools. Talend enables data profiling, deduplication, standardization, and real-time validation. It integrates with cloud services and big data platforms, making it suitable for telecom operators and media enterprises seeking enterprise-grade data quality solutions.
- Google Cloud Data Loss Prevention (DLP): While primarily a privacy tool, DLP plays a role in data quality by identifying and masking sensitive information such as PII. Media companies handling user comments, voice transcripts, or personalized ads benefit from integrating DLP into their pipelines to ensure compliance and reduce noise.
- Collibra and Alation: These are enterprise data governance and cataloging platforms that help organizations manage metadata, enforce data policies, and track data lineage. They are critical for data stewardship and governance in large telecom and media organizations with distributed data assets.

- **OpenRefine:** A powerful tool for data wrangling and transformation, ideal for cleaning and standardizing messy datasets. It is often used in the media industry for reconciling metadata and preparing historical archives for analysis or model training.
- **DataHub and Amundsen:** These open-source metadata platforms help teams discover, document, and govern data assets. Integrating such tools into ML workflows promotes transparency and accelerates onboarding for new data scientists or engineers.

By leveraging these tools and platforms, telecom and media organizations can automate the enforcement of data quality standards, gain better control over data lineage, and reduce the risk of faulty model predictions. Selecting the right combination of technologies depends on factors such as data volume, pipeline complexity, compliance needs, and integration with existing infrastructure.

8. Conclusion

Data quality is not a peripheral or optional concern—it is a foundational element that underpins the success of machine learning initiatives in telecom and media. These sectors are characterized by their reliance on massive, high-velocity, and highly diverse datasets, making the integrity and reliability of data even more critical. Without strong data quality, even the most sophisticated ML models can deliver misleading results, fail in production environments, and erode trust across internal stakeholders and end users.

As demonstrated throughout this paper, data quality impacts every stage of the ML lifecycle—from data ingestion and preprocessing to model deployment and ongoing monitoring. Poor data quality contributes to inaccurate predictions, operational disruptions, compliance violations, and degraded customer experiences. On the other hand, well-maintained data pipelines can significantly enhance the robustness, interpretability, and performance of ML systems.

Investing in data quality is not just about cleaning datasets; it is about institutionalizing best practices, deploying the right tools, fostering cross-functional collaboration, and embedding a culture of accountability around data. Organizations that prioritize data quality position themselves to unlock deeper insights, gain competitive advantage, and meet evolving business and regulatory demands.

For telecom and media enterprises seeking to scale AI/ML capabilities, data quality must be treated as a strategic asset. It is the difference between reactive analytics and proactive intelligence—between fragmented user experiences and highly personalized engagement—between innovation stagnation and digital transformation. As the industry evolves, the ability to harness and trust your data will define the winners in the next generation of intelligent systems.

9. References

- G. Wang, Y. Zhang, and X. Wang, "Big Data Analytics in Telecom: Challenges and Benefits," IEEE Access, vol. 7, pp. 90212–90230, 2019. <u>https://doi.org/10.1109/ACCESS.2019.2927483</u>
- D. Sculley, G. Holt, D. Golovin, et al., "Hidden Technical Debt in Machine Learning Systems," Advances in Neural Information Processing Systems, vol. 28, NeurIPS, 2015. <u>https://papers.nips.cc/paper_files/paper/2015/hash/86df7dcfd896fcaf2674f757a2463eba-Abstract.html</u>
- 3. P. Provost, "Data Quality for Machine Learning," in Data Science Handbook, O'Reilly Media, 2020. https://learning.oreilly.com/library/view/data-science-handbook/9781491912038/

- S. Jaiswal, R. Kaur, and A. Malik, "Improving Data Quality in Large-Scale Media Applications," ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM), vol. 18, no. 1, pp. 1–24, 2022. <u>https://doi.org/10.1145/3487550</u>
- 5. D. Zhang and L. Yang, "A Survey on Data Cleaning Techniques in Big Data Era," Journal of Data and Information Quality (JDIQ), vol. 11, no. 3, pp. 1–28, 2019. <u>https://doi.org/10.1145/3342198</u>
- Apache Griffin Data Quality Solution for Big Data. Apache Software Foundation. Accessed Dec. 2023. <u>https://griffin.apache.org</u>
- 7. Great Expectations Always Know What to Expect from Your Data. Accessed Dec. 2023. https://greatexpectations.io
- 8. AWS Deequ Data Quality Validation Library for Big Data. Amazon Web Services. Accessed Dec. 2023. <u>https://github.com/awslabs/deequ</u>
- Talend Data Quality Tools. Talend Inc. Accessed Dec. 2023. <u>https://www.talend.com/products/dataquality/</u>
- 10. Google Cloud Data Loss Prevention (DLP). Google Cloud. Accessed Dec. 2023. https://cloud.google.com/dlp
- 11. OpenRefine A Free, Open Source, Powerful Tool for Working with Messy Data. Accessed Dec. 2023. <u>https://openrefine.org</u>
- 12. DataHub A Metadata Platform for the Modern Data Stack. LinkedIn Engineering. Accessed Dec. 2023. <u>https://datahubproject.io</u>
- 13. Amundsen A Data Discovery and Metadata Platform for Improving Data Visibility. Lyft Engineering. Accessed Dec. 2023. <u>https://www.amundsen.io</u>