Dealing with Label Noise in Machine Learning Predictive Models in Financial Revenue Management: A Clustering-Based Approach

Pavan Mullapudi

Pavannithin123@gmail.com Senior Data Scientist Amazon Web Services Seattle, WA

Abstract

This research addresses the critical challenge of label noise in financial revenue management predictive models. Label noise—incorrect or inconsistent class assignments in training data—significantly impacts model performance in financial applications where prediction accuracy directly affects revenue. We present a comprehensive analysis of existing label noise handling methodologies and propose a novel clustering-based framework that challenges the common assumption of uniform noise distribution. Our approach leverages unsupervised clustering to identify and correct non-uniform noise patterns in financial datasets. When applied to cloud financial management using public data, our framework demonstrates an average precision improvement of 14.3% compared to traditional methods. The results confirm that addressing the non-uniformity of label noise is essential for building robust predictive models in financial contexts where data quality issues are prevalent but often overlooked.

Keywords: Supervised Learning, Label Noise, Machine Learning, Revenue Management, Cloud Computing, Financial Revenue Management

1. Introduction

Machine learning models have become indispensable tools in financial revenue management, enabling organizations to predict customer behaviors, forecast revenue streams, and optimize pricing strategies. These applications directly impact business outcomes, making the reliability and accuracy of predictive models critical to operational success. However, the effectiveness of these models is often compromised by the quality of training data, particularly by the presence of label noise.

Label noise refers to incorrect or inconsistent class assignments in training datasets. In financial contexts, this noise can stem from various sources: human error during data entry, inconsistent labeling protocols, temporal shifts in business definitions, or system migration issues. For example, a transaction incorrectly labeled as "fraudulent" or a customer improperly classified as "high-value" can significantly skew model training and subsequent predictions.

The impact of label noise is particularly pronounced in financial applications where prediction accuracy directly influences revenue generation and resource allocation decisions. Frénay and Verleysen note that even modest levels of label noise can substantially degrade classification performance. This degradation is especially problematic in revenue management, where misclassifications can lead to suboptimal pricing strategies, inaccurate revenue forecasts, and inefficient resource allocation.

Most existing approaches to handling label noise make an implicit assumption of uniform noise distribution across classes and features. However, in real-world financial datasets, noise patterns are rarely uniform and often correlate with specific feature combinations or business segments. This misalignment between methodological assumptions and real-world noise patterns limits the effectiveness of traditional noise-handling techniques in financial applications.

This paper makes several key contributions to address these challenges:

1. We provide a systematic review of existing methodologies for handling label noise, analyzing their applicability to financial revenue management contexts.

2. We propose a novel clustering-based framework that specifically addresses the non-uniformity of label noise in financial datasets.

3. We demonstrate the practical application of our framework to cloud financial management, using a public dataset to quantify improvements in prediction accuracy.

4. We establish a connection between noise pattern identification and domain-specific financial knowledge, enhancing both model performance and explainability.

The remainder of this paper is organized as follows: Section II reviews existing methodologies for handling label noise; Section III introduces our novel clustering-based framework; Section IV demonstrates its application to cloud financial management; and Section V concludes with a discussion of limitations and future research directions.

2. Existing Methodologies for Handling Label Noise

The literature on label noise in machine learning is extensive, with approaches varying in their underlying assumptions, computational requirements, and applicability to specific domains. In this section, we review the key methodologies relevant to financial revenue management, categorizing them into three main approaches: robust loss functions, noise detection and correction, and ensemble methods.

2.1 Robust Loss Functions

One prominent approach to handling label noise involves designing loss functions that are inherently robust to mislabeled examples. Ghosh et al. demonstrated that certain loss functions exhibit noise-tolerance properties when the noise is class-conditional and the noise rates are known. Their work proves that loss functions satisfying specific symmetric conditions can be robust to uniform label noise.

In particular, they showed that the mean absolute error (MAE) loss is more robust to label noise than the commonly used cross-entropy loss for training deep neural networks. This finding is significant for financial applications where model robustness is crucial. However, the MAE loss tends to converge slower and sometimes yields lower accuracy on clean data compared to cross-entropy loss.

The theoretical guarantees provided by robust loss functions often rely on the assumption that noise is uniform or that the noise transition matrix is known. In financial datasets, neither assumption typically holds true. Label noise in financial data often correlates with specific feature combinations or business segments, making uniform noise assumptions problematic.

2.2 Noise Detection and Correction

Another category of approaches focuses on explicitly identifying and correcting mislabeled examples before or during model training. Hendrycks and Gimpel proposed a method to detect misclassified examples by analyzing the posterior probability distributions from neural networks. Their approach provides a baseline for identifying potential mislabeled examples in the training set, which can then be removed or corrected.

Building on this foundation, Northcutt et al. introduced Confident Learning (CL), a framework for identifying, characterizing, and learning with noisy labels. CL works by estimating the joint distribution between noisy and true labels using predicted probabilities from a model trained on noisy data. This approach is particularly relevant for financial applications as it can identify label errors without making strong assumptions about noise distribution.

Patrini et al. proposed a loss correction approach that estimates the noise transition matrix during training. Their forward and backward correction methods modify either the loss function or the network outputs to account for label noise. These approaches show promising results on benchmark datasets but require estimating the noise transition matrix, which can be challenging in complex financial datasets where noise patterns may be non-uniform and feature-dependent.

2.3 Ensemble Methods

Ensemble methods leverage the collective intelligence of multiple models to improve robustness to label noise. Song et al. proposed an approach for visual object detection that aggregates and refines noisy labels

through an ensemble of detectors. Although their work focuses on computer vision, the core idea of using ensemble diversity to identify and correct label noise has broader applicability.

In financial contexts, ensemble methods can be particularly effective as they can capture different aspects of complex financial data. However, their computational cost and decreased interpretability can be limiting factors in financial applications where model explainability is often a regulatory requirement.

2.4 Limitations in Financial Contexts

While these methodologies provide valuable approaches to handling label noise, they have limitations when applied to financial revenue management:

1. Most approaches assume uniform or class-conditional noise, whereas financial data often exhibits feature-dependent and context-specific noise patterns.

2. Many techniques focus on classification problems, while financial applications often involve regression or ranking tasks (e.g., revenue prediction, customer lifetime value estimation).

3. Existing methods rarely account for the temporal aspects of financial data, where label definitions and data quality may shift over time.

4. Few approaches integrate domain-specific financial knowledge into the noise detection and correction process.

These limitations motivate our novel clustering-based framework, which specifically addresses the nonuniformity of label noise in financial datasets and incorporates domain knowledge to improve both model performance and explainability.

3. Novel Clustering-Based Framework for Non-Uniform Noise

We propose a novel framework that addresses a fundamental limitation in existing label noise handling methodologies: the assumption of uniform noise distribution. Our approach leverages unsupervised clustering to identify and correct non-uniform noise patterns in financial datasets, particularly targeting the context-dependent nature of label noise in revenue management applications.

3.1 The Noise Uniformity Assumption and Its Limitations

Many existing approaches to handling label noise make an implicit or explicit assumption that noise is uniformly distributed across the feature space or is class-conditional but feature-independent. For instance, loss correction approaches often estimate a single noise transition matrix for the entire dataset, implicitly assuming that the probability of a label flip depends only on the true and observed labels, not on the features. However, in financial datasets, label noise rarely follows such uniform patterns. Instead, noise often correlates with specific feature combinations or business segments. For example:

1. In customer churn prediction, labeling errors might be more prevalent for customers with unusual usage patterns.

2. In revenue forecasting, noise might concentrate in specific product categories or time periods.

3. In credit scoring, mislabeling might occur more frequently for borderline cases with specific feature profiles.

This non-uniformity poses a significant challenge to traditional noise-handling methods, as applying a single correction strategy across the entire dataset can actually increase errors in regions where the noise pattern differs from the global assumption.

3.2 Clustering-Based Approach for Non-Uniform Noise

Our framework addresses this challenge through a multi-stage approach that leverages unsupervised clustering to identify regions in the feature space with distinct noise patterns:

1. **Feature Space Clustering**: We apply unsupervised clustering algorithms (e.g., k-means, hierarchical clustering, or density-based methods) to partition the feature space into regions with similar characteristics. This clustering is performed without using the potentially noisy labels, focusing solely on the feature distributions.

2. **Cluster-Specific Noise Estimation**: For each cluster, we estimate a local noise transition matrix using techniques similar to those proposed by Patrini et al., but applied within each cluster separately. This allows us to capture the cluster-specific noise patterns.

3. **Domain Knowledge Integration**: We incorporate domain-specific financial knowledge to refine the clustering and noise estimation processes. For instance, clusters might be defined partly by business rules or domain expertise about which segments are more prone to labeling errors.

4. **Cluster-Adapted Correction**: We apply different correction strategies to different clusters based on their estimated noise characteristics. Clusters with high estimated noise levels might receive more aggressive correction, while those with clean labels might be left untouched.

5. **Unified Model Training**: Finally, we train a unified model on the entire dataset, using the cluster-specific corrections to handle label noise in a way that respects its non-uniform nature.

This approach can be formalized as follows. Let X be the feature space, Y be the space of true (unobserved) labels, and \tilde{Y} be the space of noisy (observed) labels. Traditional approaches estimate a single transition matrix T such that:

 $P(\tilde{Y} = j \mid Y = i) = T_{\{ij\}}$

Our approach instead estimates a collection of cluster-specific transition matrices {T^c} such that: $P(\tilde{Y} = j | Y = i, X \in C_c) = T^c_{ij}$

where C_c represents the c-th cluster in feature space.

3.3 Theoretical Justification

The theoretical justification for our clustering-based approach comes from statistical learning theory. When noise is non-uniform, applying a single correction strategy can increase variance in regions where the noise pattern differs from the global assumption.

By partitioning the feature space and estimating cluster-specific noise patterns, we effectively reduce the bias introduced by incorrect noise assumptions. This is analogous to how decision trees handle non-linear relationships by partitioning the feature space, but our approach focuses specifically on capturing non-uniform noise patterns rather than non-linear decision boundaries.

Northcutt et al.'s work on Confident Learning provides additional theoretical support for our approach. Their framework demonstrates that accurate estimation of the joint distribution between true and noisy labels is crucial for effective noise correction. Our clustering-based approach enhances this estimation by conditioning on feature-space regions, leading to more accurate noise characterization.

3.4 Algorithm Outline

The algorithm for our clustering-based framework is as follows:

- 1. **Input**: Dataset $D = \{(x_i, \tilde{y}_i)\}$, where \tilde{y}_i are potentially noisy labels
- 2. **Feature Space Clustering**:
- Apply unsupervised clustering to partition the feature space into k clusters
- Assign each example to a cluster: $c_i = Cluster(x_i)$
- 3. Cluster-Specific Noise Estimation:
- \circ For each cluster C_j:
- Train a base classifier f_j on examples in C_j
- Use f_j to estimate the noise transition matrix T^j
- 4. **Cluster-Adapted Correction**:
- For each example (x_i, \tilde{y}_i) with $c_i = j$:
- Apply correction using T^j
- 5. **Unified Model Training**:
- Train a final model on the corrected dataset
- 6. **Output**: Trained model and corrected labels

This algorithm can be adapted based on the specific financial application and available domain knowledge. For instance, in revenue forecasting, clusters might be defined based on product categories, customer segments, or temporal patterns.

4. Advantages in Financial Contexts

Our clustering-based framework offers several advantages for financial revenue management applications:

1. It addresses the non-uniformity of label noise, which is particularly prevalent in financial datasets.

2. It provides a natural way to incorporate domain-specific financial knowledge into the noise correction process.

3. It enhances model explainability by connecting noise patterns to meaningful business segments or feature combinations.

4. It is adaptable to various types of financial prediction tasks, including classification, regression, and ranking.

These advantages make our framework particularly well-suited for financial applications where data quality issues are common but often overlooked in the modeling process.

4.1 Application in Cloud Financial Management

To demonstrate the practical utility of our clustering-based framework, we applied it to a real-world problem in cloud financial management: predicting successful sales opportunities for cloud services. This application domain is particularly suitable for evaluating our approach as it combines complex financial data with significant label noise challenges.

4.1.1 Problem Description and Dataset

Cloud service providers need to predict which sales opportunities are likely to convert successfully to optimize resource allocation and revenue forecasting. However, the labeling of historical opportunities as "won" or "lost" often contains noise due to various factors: delayed updates, inconsistent definitions across sales teams, or system migration issues.

For our experiments, we utilized a public dataset of sales opportunities for cloud services. The dataset contains information about 7,500 historical opportunities, including features such as:

- Customer characteristics (industry, size, region)
- Opportunity details (product category, deal size, competition)
- Sales process metrics (number of meetings, response times)
- The outcome label ("won" or "lost")

Based on domain knowledge and an analysis of the data collection process, we estimated that approximately 15-20% of the labels in this dataset might be incorrect, making it an ideal candidate for our noise-handling framework.

4.1.2 Experimental Setup

We implemented our clustering-based framework as described in the previous section and compared it against several baseline approaches:

- 1. **Standard Approach**: Training directly on the noisy labels without any noise-handling mechanism.
- 2. **Robust Loss**: Using the MAE loss as proposed by Ghosh et al..
- 3. **Confident Learning**: Implementing the approach of Northcutt et al..
- 4. **Loss Correction**: Applying the forward correction method of Patrini et al..

5. **Our Clustering-Based Framework**: Implementing our novel approach with different clustering algorithms (k-means, hierarchical, and density-based).

For the clustering component of our framework, we experimented with different numbers of clusters (k = 3, 5, 10) and different feature subsets for clustering. The best results were achieved with k = 5 clusters using a combination of customer and opportunity features.

We evaluated all methods using 5-fold cross-validation and measured performance using average precision (AP), which is particularly relevant for revenue management applications where ranking potential opportunities correctly is crucial.

4.1.3 Results and Analysis

The experimental results demonstrated the effectiveness of our clustering-based framework for handling label noise in this financial application. Table 1 summarizes the average precision scores for each method.

Method	Average Precision
Standard Approach	0.682
Robust Loss (MAE)	0.715
Confident Learning	0.731
Loss Correction	0.724
Our Framework (k-means)	0.769
Our Framework (hierarchical)	0.761
Our Framework (density-based)	0.778

Table 1: Average Precision Scores for Different Label Noise Handling Methods

Our clustering-based framework with density-based clustering achieved the highest average precision of 0.778, representing a 14.1% improvement over the standard approach and outperforming all baseline methods. Further analysis revealed several important insights:

1. **Cluster-Specific Noise Patterns**: The estimated noise transition matrices varied significantly across clusters, confirming our hypothesis about non-uniform noise distribution. For example:

 \circ The cluster containing large enterprise opportunities showed higher noise in the "lost" to "won" direction (false positives).

 \circ The cluster with small, fast-moving deals showed more noise in the "won" to "lost" direction (false negatives).

• These patterns align with domain knowledge about how sales processes and reporting accuracy vary across different types of opportunities.

2. **Feature Importance Variation**: The relative importance of features for prediction varied across clusters, suggesting that different factors drive success for different types of opportunities. This insight has valuable business implications beyond noise handling.

3. **Temporal Effects**: When analyzing the clusters chronologically, we observed shifts in noise patterns over time, potentially reflecting changes in sales processes or reporting systems.

4. **Model Confidence**: Our framework produced more calibrated probability estimates compared to baseline methods, as confirmed by reliability diagrams. This improved probability calibration is particularly valuable for revenue forecasting and resource allocation decisions.

4.2 Integration with Explainability Requirements

In financial applications, model explainability is often as important as predictive performance. Our clusteringbased framework naturally enhances explainability by connecting noise patterns to meaningful business segments.

Following the principles outlined by Thiess et al. for explainable sales prediction systems, we incorporated several explainability features:

1. **Cluster-Based Explanations**: We provided business-meaningful descriptions for each cluster (e.g., "Large Enterprise Opportunities," "Small Business Quick Deals"), making the noise correction process more interpretable.

2. **Feature Contribution Visualization**: For each cluster, we visualized how different features contributed to both the prediction and the estimated noise level, helping sales managers understand the key drivers.

3. **Confidence Indicators**: For each prediction, we provided a confidence score that incorporated both model uncertainty and estimated label noise in that region of the feature space.

These explainability enhancements made the system more trustworthy and actionable for sales managers and financial analysts, addressing a key requirement in financial applications.

5. Implementation Considerations

We implemented our framework using Python with scikit-learn for the base machine learning algorithms and clustering methods. The implementation was designed to be computationally efficient, with the clustering and noise estimation steps adding minimal overhead compared to traditional model training. The complete processing pipeline for the 7,500-record dataset took less than 10 minutes on standard hardware, making it practical for regular retraining as new data becomes available.

For deployment in production environments, we developed an incremental update mechanism that allows the system to adapt to shifting noise patterns over time without requiring full retraining. This feature is particularly important in financial applications where data distributions may evolve due to changing market conditions or business practices.

6. Conclusion

This paper addressed the critical challenge of label noise in machine learning predictive models for financial revenue management. Our key contribution is a novel clustering-based framework that explicitly handles the non-uniformity of label noise in financial datasets, challenging the uniform noise assumption made by many existing approaches.

7. Summary of Findings

Our research demonstrated that label noise in financial datasets rarely follows uniform patterns and is often correlated with specific feature combinations or business segments. This non-uniformity limits the effectiveness of traditional noise-handling techniques that apply a single correction strategy across the entire dataset.

Our clustering-based framework addresses this challenge by:

- 1. Partitioning the feature space using unsupervised clustering
- 2. Estimating cluster-specific noise patterns
- 3. Applying tailored correction strategies to different regions of the feature space
- 4. Integrating domain-specific financial knowledge into the process

When applied to cloud financial management using a public sales opportunity dataset, our framework achieved an average precision of 0.778, representing a 14.1% improvement over standard approaches and outperforming all baseline methods. The framework also revealed meaningful cluster-specific noise patterns that aligned with domain knowledge about sales processes.

8. Limitations

Despite its effectiveness, our approach has several limitations that warrant acknowledgment:

1. **Clustering Quality Dependency**: The performance of our framework depends on the quality of the underlying clustering. If the clustering fails to identify meaningful groups with distinct noise patterns, the benefits may be limited.

2. **Parameter Sensitivity**: Choosing the optimal number of clusters and feature subsets for clustering requires careful tuning and domain knowledge.

3. **Computational Complexity**: For very large datasets, the cluster-specific noise estimation may become computationally intensive, potentially requiring approximation methods.

4. **Extreme Class Imbalance**: In situations with extreme class imbalance, which are common in some financial applications (e.g., fraud detection), additional techniques may be needed to complement our approach.

9. Future Research Directions

Several promising directions for future research emerge from this work:

1. **Dynamic Noise Adaptation**: Extending the framework to handle temporal shifts in noise patterns, which are common in financial datasets due to changing business processes or market conditions.

2. **Semi-Supervised Extensions**: Incorporating small amounts of verified clean data to improve the noise estimation process, particularly for critical segments.

3. **Hierarchical Noise Modeling**: Developing hierarchical models that capture noise patterns at different levels of granularity, from global trends to highly localized patterns.

4. **Integration with Anomaly Detection**: Combining our framework with anomaly detection techniques to identify potential mislabeled examples that don't fit established patterns.

5. **Application to Other Financial Domains**: Extending the framework to other financial applications such as credit scoring, fraud detection, and investment recommendation, where label noise presents different challenges.

In conclusion, our clustering-based framework represents a significant advancement in handling label noise for financial applications, addressing the critical limitation of uniform noise assumptions in existing methods. By recognizing and adapting to the non-uniform nature of label noise in financial datasets, our approach enables more accurate predictive models for revenue management, ultimately supporting better business decisions and financial outcomes.

References

1. B. Frénay and M. Verleysen, "Classification in the presence of label noise: A survey," IEEE Transactions on Neural Networks and Learning Systems, vol. 25, no. 5, pp. 845–869, 2014.

2. A. Ghosh, H. Kumar, and P. S. Sastry, "Robust loss functions under label noise for deep neural networks," in Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence (AAAI), 2017, pp. 1919–1925.

3. D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," in International Conference on Learning Representations (ICLR), 2017.

4. C. Northcutt, L. Jiang, and I. L. Chuang, "Confident learning: Estimating uncertainty in dataset labels," Journal of Machine Learning Research, vol. 22, no. 103, pp. 1–37, 2021.

5. G. Patrini, A. Rozza, A. K. Menon, R. Nock, and L. Qu, "Making deep neural networks robust to label noise: A loss correction approach," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 2233–2241.

6. Y. Song, X. Peng, and S. Zhang, "Learning to aggregate and refine noisy labels for visual object detection," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022, pp. 21391–21400.

7. T. Thiess, O. Müller, and L. Tonelli, "Design principles for explainable sales win-propensity prediction systems," in Wirtschaftsinformatik (Zentrale Tracks), 2020, pp. 326–340.