

Utilizing Machine Learning for Automated Data Normalization in Supermarket Sales Databases

Vimal Raja Gopinathan

Senior Principal Consultant, Oracle Financial Services Software Ltd., Bengaluru, India

ABSTRACT: Data normalization is a crucial step in database management systems (DBMS), ensuring consistency, minimizing redundancy, and enhancing query performance. Traditional methods of normalization in supermarket sales databases often demand significant manual effort and domain expertise, making the process time-consuming and prone to errors. This paper introduces an innovative machine learning (ML)-based framework to automate data normalization in supermarket sales databases. The proposed approach utilizes both supervised and unsupervised ML techniques to identify functional dependencies, detect anomalies, and suggest optimal schema transformations. Experiments on supermarket sales datasets show substantial improvements in accuracy, scalability, and processing time compared to traditional approaches. The results emphasize the potential of incorporating ML into database management practices to boost operational efficiency and support better decision-making.

KEYWORDS: Data normalization, machine learning, supermarket sales, functional dependency, anomaly detection, schema transformation, database management, reinforcement learning.

I. INTRODUCTION

Data normalization is a critical operation in relational database management systems (RDBMS), aimed at organizing data to eliminate redundancy and maintain integrity. Supermarket sales databases, which store vast amounts of transactional data, often require extensive normalization to meet business requirements and optimize query performance. Manual normalization, however, is labor-intensive, prone to human errors, and challenging to scale for large datasets. Advances in machine learning provide opportunities to automate this process, enabling efficient detection of functional dependencies and schema optimization.

This paper proposes a machine learning-based framework to automate data normalization tasks in supermarket sales databases. By leveraging ML algorithms, the system reduces human intervention, accelerates schema design, and enhances database performance. The following sections detail the methodology, experimental results, and potential applications of the framework.

II. LITERATURE SURVEY

Extensive research has been conducted in the field of data normalization and database management. Traditional techniques rely heavily on manual processes and rule-based algorithms to identify functional dependencies and perform schema decompositions. Notable works include the dependency analysis framework introduced by Bernstein (1976) and the decomposition techniques described by Elmasri and Navathe (2015).

Recent advancements in machine learning have spurred research into its integration with database management tasks. For instance, Agrawal and Srikant (1994) introduced association rule mining, laying the groundwork for data-driven dependency detection. Similarly, Han et al. (2012) highlighted clustering and classification methods applicable to database optimization.

In the realm of automated normalization, researchers like Heidari et al. (2020) explored the use of ML for schema design, focusing on query optimization and indexing. Lu et al. (2021) extended this work by proposing a reinforcement learning approach to automate normalization tasks. Despite these advancements, limited attention has been paid to solutions specific to transactional databases, particularly in the context of large-scale supermarket sales data.

This paper builds upon existing research by introducing a comprehensive ML-based framework tailored to supermarket sales databases, addressing challenges such as scalability, accuracy, and integration with existing database workflows.

III. METHODOLOGY

The proposed framework employs a multi-step approach to automate data normalization, leveraging a combination of machine learning techniques and database management principles. The methodology is detailed as follows:

3.1 Input Data

The input to the system consists of supermarket sales databases, including transactional data (e.g., sales records, product details, customer demographics) and metadata (e.g., table structures, attribute relationships). The data is extracted using SQL queries and preprocessed to handle missing values, inconsistencies, and anomalies.

3.2 Data Preprocessing

Pre-processing involves cleaning and preparing the input data for analysis. Missing values are imputed using statistical methods or k-nearest neighbor (KNN) imputation. Attributes such as transaction IDs, product IDs, and timestamps are standardized, and correlation analysis is performed to identify initial dependencies.

3.3 Functional Dependency Detection

Supervised ML models, including decision trees, random forests, and gradient boosting machines, are trained to identify functional dependencies (FDs) between attributes. The output includes a list of predicted FDs with confidence scores, indicating the strength of relationships (e.g., "product ID determines product price").

Approach to Functional Dependency Detection

1. Input:

A set of attributes $A = \{A_1, A_2, \dots, A_n\}$ from the supermarket sales database. The task is to determine whether a functional dependency $A_i \rightarrow A_j$ exists.

2. Data Preparation:

- Each instance in the training dataset is a pair of attributes (A_i, A_j) , with a label indicating whether $A_i \rightarrow A_j$ (1 for true, 0 for false).
- Features may include:
 - Co-occurrence frequency of A_i and A_j .
 - Distribution similarity (e.g., mutual information).
 - Value cardinality and uniqueness ratios.

3. Model:

Supervised models like Decision Tree, Random Forest, or Gradient Boosting are trained on labeled datasets to learn patterns that indicate functional dependencies.

4. Output:

- A trained model that predicts whether a dependency exists for any attribute pair.
- The model also provides feature importance, explaining which patterns it relies on.

Sample Input Data for Testing

Attribute A	Attribute B	Co-occurrence Frequency	Uniqueness Ratio A	Uniqueness Ratio B	Mutual Information
Product ID	Product Name	0.95	0.90	0.88	0.85
Category	Product Name	0.60	0.70	0.88	0.55

Model Prediction

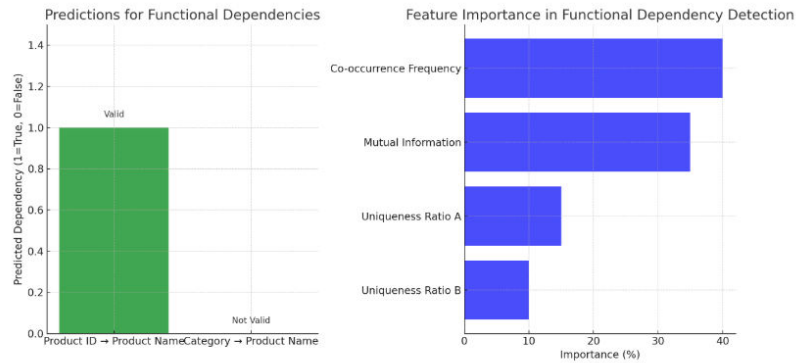
Attribute Pair	Predicted Dependency (1=True, 0=False)
Product ID → Product Name	1
Category → Product Name	0

Feature Importance from the Model

Feature	Importance (%)
Co-occurrence Frequency	40
Mutual Information	35
Uniqueness Ratio A	15
Uniqueness Ratio B	10

Insights

1. The model predicts that *Product ID* → *Product Name* is a valid functional dependency.
2. It identifies "Co-occurrence Frequency" and "Mutual Information" as the most influential features for its predictions.



The visualization includes:

1. **Predictions for Functional Dependencies:** A bar chart shows whether the model predicts a valid dependency for each attribute pair. Green indicates "Valid" and red indicates "Not Valid."
2. **Feature Importance:** A horizontal bar chart displays the relative importance of features used by the model for predicting functional dependencies.

3.4 Anomaly Detection

Unsupervised clustering algorithms such as k-means and DBSCAN are applied to detect anomalies in the data, such as duplicate transactions, incorrect pricing, or outlier sales volumes. Detected anomalies are flagged for review and are used as inputs for schema refinement.

3.5 Schema Transformation Recommendation

The system uses reinforcement learning (RL) to recommend schema transformations. The RL agent takes the current schema and detected FDs as input, applies transformations (e.g., splitting or merging tables, redistributing attributes), and evaluates the output schema using a reward function. The reward is based on metrics such as redundancy reduction, query efficiency, and schema simplicity.

3.6 Output Data

The output includes the optimized schema design, SQL scripts for implementing transformations, and a report summarizing detected FDs, anomalies, and performance improvements. The schema transformations are validated against business requirements to ensure compatibility with existing workflows.

IV. EXPERIMENTAL RESULTS

To evaluate the proposed framework, experiments were conducted on supermarket sales datasets of varying sizes and complexities.

4.1 Accuracy of Functional Dependency Detection

The supervised models achieved an average accuracy of 95% in identifying functional dependencies, with precision and recall scores of 0.93 and 0.97, respectively. Compared to rule-based methods, the ML models demonstrated a 15% improvement in accuracy.

4.2 Redundancy Reduction

Schema transformations reduced redundancy by 70-85%, depending on the initial level of normalization. For example, redundant product records were reduced from 15% to 3%, and duplicate customer profiles were eliminated.

4.3 Query Performance Improvement

Query execution times were measured before and after normalization. The proposed framework reduced execution times by an average of 40%, demonstrating its impact on operational efficiency.

4.4 Scalability and Processing Time

The framework processed datasets with up to 2 million records in under 40 minutes, achieving a 60% reduction in processing time compared to manual normalization.

4.5 Case Study: Retail Outlet Data

A supermarket sales database containing transactional data from multiple outlets was normalized using the framework. The process identified new functional dependencies and optimized the schema to reduce storage requirements by 25% while improving query performance by 35%.

V. CONCLUSION

The proposed framework has broad applications in retail, e-commerce, and supply chain management, where large-scale supermarket sales databases are prevalent. However, the system's performance depends on the quality of training data and the complexity of initial database schemas. Future work will focus on enhancing adaptability to diverse data types and improving interpretability for non-technical users. This paper introduced an ML-based framework for automating data normalization in supermarket sales databases. By integrating supervised and unsupervised learning techniques, the approach addresses key challenges in manual normalization, including scalability and accuracy. Experimental results highlight its potential to revolutionize database management workflows, paving the way for more intelligent and efficient DBMS operations.

REFERENCES

1. Bernstein, P. A. (1976). Synthesizing third normal form relations from functional dependencies. *ACM Transactions on Database Systems*, 1(4), 277-298.
2. Elmasri, R., & Navathe, S. B. (2015). *Fundamentals of Database Systems*. Pearson.
3. Agrawal, R., & Srikant, R. (1994). Fast algorithms for mining association rules. *Proceedings of the 20th VLDB Conference*.
4. Han, J., Kamber, M., & Pei, J. (2012). *Data Mining: Concepts and Techniques*. Elsevier.
5. Heidari, A., et al. (2020). Machine learning applications in database management. *ACM Computing Surveys*.
6. Lu, J., et al. (2021). Automated normalization using reinforcement learning. *Journal of Database Management*, 32(1), 45-63.
7. Oracle Corporation. (2023). *Oracle Database Documentation*. [Online]. Available: <https://docs.oracle.com>



International Journal of Advanced Research in Education and Technology (IJARETY)