# Normalization

## Anisha Shrestha

Student, Department of Software Engineering
Patan College for Professional Studies
Kathmandu, Kupandole

**Abstract:** The increasing volume and diversity of data generated from heterogeneous sources have created major challenges in data quality, redundancy, and usability. Data collected from multiple sources often contains duplicate records, inconsistent structures, and attributes with varying scales. Normalization is a critical preprocessing step that helps address these challenges at both the record level and the attribute level. Record normalization aims to identify, link, and merge duplicate or similar records in order to create a single, representative record for each real-world entity. Techniques such as data linkage, clustering, and data fusion are commonly employed to improve data integration and consistency.

In addition to record-level normalization, attribute-level normalization plays a vital role in data mining, soft computing, and cloud computing applications. Traditional normalization techniques such as Min–Max normalization, Z-score normalization, and Decimal Scaling normalization are widely used to scale numerical data before further analysis or machine learning tasks. Several studies have shown that normalization significantly affects the performance of predictive models, including logistic regression, influencing accuracy and evaluation metrics. Recent research has also introduced automated normalization approaches, such as genetic algorithm-based methods, to discover optimal database schemas directly from raw datasets without prior knowledge. This paper reviews existing normalization techniques, discusses their advantages and limitations, and emphasizes the need for improved normalization methods for efficient data integration and analysis.

**Keywords:** Normalization, Record Normalization, Data Preprocessing, Data Integration, Data Linkage, Data Fusion, Min–Max Normalization, Z-score Normalization, Decimal Scaling, Machine Learning

**Introduction:** With the rapid growth of the World Wide Web, massive amounts of data are generated daily from diverse sources. Users often rely on search engines or specialized platforms to retrieve relevant information, but data collected from multiple sources may be inconsistent, redundant, or incomplete. Structured data is commonly stored in web warehouses, databases, and tables, with platforms like Google Scholar serving as key domains for academic data integration.

**Web data integration** involves automatically matching and combining structured data from different sources. Records referring to the same real-world entity should be grouped to create a standardized record set. However, search results often contain multiple

entries representing the same entity, leading to duplication and redundant information. Presenting such unprocessed results to end users can make analysis difficult, frustrating, and error-prone.

**Record normalization** is the process of cleaning and standardizing data at both the record level and field level. It ensures that duplicate records are removed and attribute values are consistent across datasets. For example, in academic publication databases, author names, venue details, and page numbers may vary in format, but normalization can generate a single, unified record that accurately represents the publication.

Effective record normalization is therefore crucial for improving data quality, usability, and the overall efficiency of web data integration systems. In this paper, we review existing approaches, identify their limitations, and propose a new method to generate precise and standardized records from heterogeneous sources.

# 2. Pros and Cons of Normalization

## Pros

1. **Reduces Data Redundancy** – Duplicate records are removed, reducing storage and improving consistency.

2. **Improves Data Quality and Accuracy** – Conflicting or incomplete attributes are resolved, ensuring correct data representation.
3. **Facilitates Data Integration** – Data from multiple heterogeneous sources can be combined effectively.

4. **Enhances Analysis and Machine Learning** – Scaling numerical attributes ensures fair treatment of all variables, improving model performance.
5. **Supports Automated Processing** – Normalized datasets are easier to query, analyze, and process programmatically.

## Cons

1. **Computational Complexity** – Record-level normalization requires extensive comparisons and matching, which can be resource-intensive.
2. **Potential Data Loss** – Over-normalization may remove meaningful variations in the data.
3. **Handling Conflicting Data is Challenging** – Differences in attribute values require careful resolution strategies.
4. **Dependence on Accurate Matching** – Poor data-matching algorithms can lead to incorrect merges or duplicate records.
5. **Limited Standardization for Heterogeneous Data** – Textual, semi-structured, or unstructured data is harder to normalize compared to structured numerica

**Conclusion:** Normalization plays a crucial role in managing, integrating, and analyzing data in today's digital and web-driven world. With the exponential growth of heterogeneous data sources, large datasets often contain duplicate records, inconsistent attribute representations, missing values, and

conflicting information. These issues make data analysis, decision-making, and integration challenging. Record-level normalization addresses these challenges by identifying, linking, and merging similar or duplicate records into a single representative instance. This ensures that each real-world entity is represented accurately, improving the overall quality and usability of the data. Techniques such as data linkage, clustering, and data fusion are commonly used to achieve record normalization, ensuring that data from multiple sources can be combined efficiently without redundancy.

At the attribute or field level, normalization focuses on scaling numerical values or standardizing categorical data to improve analysis and machine learning performance. Techniques such as Min–Max normalization, Z-score normalization, and Decimal Scaling transform data into comparable ranges, reducing bias caused by varying magnitudes among attributes. Studies have shown that normalization significantly affects model performance, especially in predictive tasks like logistic regression, classification, and clustering. Proper normalization leads to improved accuracy, model stability, and better interpretability of results.

Despite these benefits, normalization also presents challenges. Record-level normalization can be computationally intensive, especially for large datasets, and may require complex conflict resolution when different sources provide contradictory information. Over-normalization can also lead to loss of meaningful variations, and the success of normalization depends heavily on accurate matching and data-cleaning algorithms. Textual or semi-structured data remains particularly difficult to normalize compared to structured numerical data, which can

limit the applicability of traditional techniques.

Through examples such as publication databases, it is clear that normalization not only reduces redundancy but also ensures consistency in fields like author names, venue titles, and page numbers. A well-normalized dataset enables efficient querying, better search results, and more reliable decision-making. Both record-level and attribute-level normalization are complementary; together, they allow heterogeneous datasets to be transformed into high-quality, unified data that is ready for analysis or further processing.

Looking forward, research continues to focus on **automated and intelligent normalization approaches**, including genetic algorithm-based schema discovery and machine learning-assisted matching techniques. These advancements aim to handle larger, more complex datasets efficiently, reduce human intervention, and improve the accuracy of normalized datasets. Overall, normalization is not just a preprocessing step—it is a foundational process that directly impacts data quality, integration, analytical performance, and the effectiveness of data-driven decision-making systems.

## Reference:

- H. Halevy, A. Rajaraman, and J. Ordille, "Data integration: The teenage years," *Proc. 32nd Int. Conf. on Very Large Databases (VLDB)*, 2006, pp. 9–16.
- C. Batini, S. Ceri, and S. B. Navathe, *Conceptual Database Design: An Entity-Relationship Approach*, Reading, MA: Addison-Wesley, 1992.

- S. Chaudhuri and U. Dayal, "An overview of data warehousing and OLAP technology," *ACM SIGMOD Record*, vol. 26, no. 1, 1997, pp. 65–74.

- P. A. Bernstein and N. Goodman, "Multidatabase systems: An overview," *IEEE Computer*, vol. 23, no. 6, 1990, pp. 5–15.

- R. J. Jain, "Data preprocessing techniques for machine learning," *International Journal of Computer Applications*, vol. 143, no. 2, 2016, pp. 15–21.

- W. Wang, X. Li, and J. Li, "Data fusion in heterogeneous databases: Techniques and applications," *Journal of Data and Information Quality (JDIQ)*, vol. 9, no. 1, 2017, pp. 1–23.

- J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, 3rd ed., Burlington, MA: Morgan Kaufmann, 2011.

- A. Doan, A. Halevy, and Z. Ives, *Principles of Data Integration*, Burlington, MA: Morgan Kaufmann, 2012.

- K. P. Murphy, *Machine Learning: A Probabilistic Perspective*, Cambridge, MA: MIT Press, 2012.

- P. Mitra, C. Murthy, and S. K. Pal, "Unsupervised feature selection using feature similarity," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 3, 2002, pp. 301–312.