# Data Preprocessing for Machine Learning

**ARYA BRIJITH**

International Center for AI and Cyber Security Research and Innovations (CCRI), Asia University, Taiwan

(e-mail: arya.brijithk@gmail.com).

**ABSTRACT** A crucial step in the data analysis process is preprocessing, which involves converting raw data into a format that computers and machine learning algorithms can understand. This important phase has a big impact on the precision and efficiency of machine learning models. The importance of data preparation is emphasized as this study explores the many forms of data used in machine learning. Preprocessing guarantees that the data used for modeling are of good quality by resolving problems like noisy, redundant, and missing data. The essay also looks at practical applications to highlight the practical advantages of data preprocessing.

**KEYWORDS** data preprocessing, methods, machine learning, normalization.

## I.    INTRODUCTION

Data preprocessing is one of the key phases carried out during the data analysis process. In this phase the raw data is converted into a format that computers and machine learning algorithms can understand and evaluate.[1] It usually takes a considerable amount of time and effort to perform this task. The output of this step, is highly reliant upon by all the machine learning algorithms.

Machine learning algorithms automatically extracts knowledge from data that can be read by machines. Unfortunately, the quality of the data they use to work on determines how successful they are most of the time. Hence, by increasing the quality of the datasets used the accuracy of the model increases.
To overcome issues, such as noisy data, redundancy data, missing data values, etc., pre-processing is required.[2]

This paper addresses the types of data used in machine learning algorithms also about the importance of data preprocessing and the various methods to enhance the reader's knowledge in terms of the kind of datasets to use while creating a model and how to overcome noises and issues associated with the same. Furthermore, we will discuss about the impact of data preprocessing with some examples.

## II.    TYPES OF DATA

When using any machine learning technique, we work with a variety of data formats.

Data can be classified into the following types/:

- **Structured-** Structured data is set out in a clear way, typically displayed in the form of rows and columns. It follows an ordered pattern. There is a specific location for each piece of information. It is simple to query and examine structured data. It works particularly well for SQL (Structured Query Language) operations on conventional databases. These formats are very well-organized and simple to use. Typically, retrieving and modifying this sort of data is quick and easy. On the negative side, hierarchical data structures or complicated relationships may be difficult for structured data to handle. Furthermore, it might need to be more versatile for storing unstructured or semi-structured data. SQL databases, Excel spreadsheets, and files in CSV (Comma-Separated Values) format are three common examples of structured data.

- **Semi-structured-** Unlike structured data, which adheres to a strict framework, semi-structured data is more flexible. It does, however, include certain organizational components, such as tags or keys, which serve to offer some structure. Data that is semi-structured enables a more adaptable schema. As a result, different parts of the dataset may have different data structures. It's especially helpful when processing data whose structures change or fluctuate. Additionally, it works well for

organizing hierarchical or layered data. Compared to structured data, querying semi-structured data might be more difficult. Extracting relevant information often requires further processing and parsing procedures. These are the most widely used formats for transmitting data in a manner that can be both machine and human-readable. JavaScript Object Notation (JSON) and Extensible Markup Language (XML) are a few examples of semi-structured data formats.

- **Unstructured-** The lack of a predetermined, exact structure or format is what defines unstructured data. It frequently arrives in a format that is simple for people to read but difficult for robots to comprehend without specific methods. Data that is unstructured lacks a predetermined schema. This indicates that the data is not set up in a way that makes it simple for machines to understand it. It could include plain-language text, pictures, videos, music, or other media. These kinds of data frequently contain important information that might not be simple to capture using approaches for structured or semi-structured data. Structured or semi-structured data may struggle to capture significant information that is frequently included in unstructured data. It is an important source of information since it may include extensive multimedia content. It might be difficult to analyze unstructured data using conventional database techniques. To extract useful insights from unstructured data, specialized approaches are often needed, such as natural language processing for text or computer vision for images. Unstructured data examples include text files (such as those in Word or PDF format), pictures, audio files, videos, postings on social media, emails, and more[5-12].
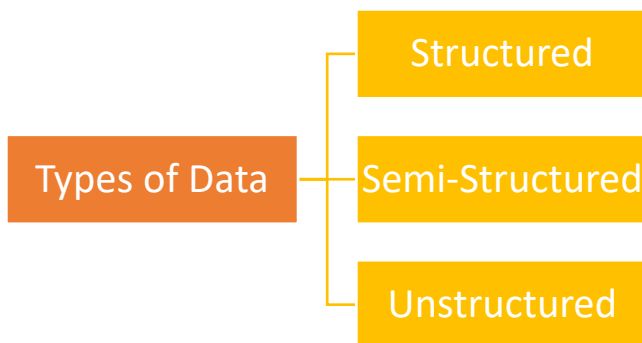


**Table 1: Types of data**

## III.	IMPORTANCE OF DATA PREPROCESSING

Real-world data is seldom without issues. Corrupt data may be the result of issues such as sensor failure, data transmission issues, or incorrect data input, many of which may not have been recognized at the time the data was collected.[4] We know that, datasets is the crucial factor for designing a model. The efficiency of the model decreases if data is not preprocessed. It deals with problems with data quality. Hence, it is crucial to perform data preprocessing.

Preprocessing of the data makes it easier to spot and manage mistakes, outliers, and missing information. This results in clearer, accurate data, which in turn produces more precise analysis and modeling. Machine learning models can perform significantly better with proper preprocessing. It makes sure that the data input into the model is organized in a way that the algorithm can learn from. Models with higher interpretability can result from well-preprocessed data. When the data is organized and clear, it is easier to understand the links between the characteristics. Preprocessing properly can help eliminate noise and unimportant input, which might result in more generalized models that are less prone to overfitting.

The appropriate handling of various data kinds is made possible by data preparation. Datasets from the real world frequently contain a combination of organized, semi-structured, and unstructured information. We can transform this variety of data into a consistent format that machine learning algorithms can analyze well by using preprocessing techniques like normalization, one-hot encoding, and text tokenization. This adaptability is essential in contemporary data science, where a variety of data sources is common.

Overall, data preparation is a critical stage in the workflow for data analysis and machine learning. It makes certain that the analysis's data is of a high caliber, correctly structured, and prepared for machine learning techniques. Preprocessing is essential for gaining insightful information and creating precise

predictive models because it addresses difficulties with data quality, enhances model performance, manages a variety of data formats, and makes feature engineering possible.

## IV. DATA PREPROCESSING METHODS

Here are a few techniques for data preprocessing:

- **Normalization-** Any raw data must go through preprocessing before being subjected to any sort of categorization or identification. The raw data must often be processed since they are unusable in the actual world due to noise. Hence, we perform normalization. There are mainly two kinds of normalization, namely, z-score normalization and min-max normalization. The process of rescaling numerical characteristics to have a mean of 0 and a standard deviation of 1 is known as "Z-score normalization." Each feature is subjected to this modification separately. Assuming that X is a random variable, it is normalized by taking the desired mean value and subtracting it from the initial value, then dividing the result by the standard deviation.[3] The formula is given below.

$$Z = \frac{x - \mu}{\sigma}$$

(z is the standardized value.
x is the original value of the feature.
μ is the mean of the feature.
σ is the standard deviation of the feature.)

Min-Max normalization reduces numerical characteristics to a certain range, often [0, 1] or [-1, 1]. The values are linearly scaled to do this. It is a technique for feature scaling in which numerical feature values are transformed within a predetermined range, often between 0 and 1. The processed data may be represented by mapping the minimum value of X to 0 and the greatest value of Y to 1.[3] The formula is mentioned below.

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)} \times (\max new - \min new) + \min new$$

- **Feature extraction-** In machine learning and data analysis, dimensionality reduction is a technique used to minimize the number of variables in a dataset while maintaining the most pertinent data. When working with high-dimensional data, where the ratio of features to samples is high, it is especially helpful. To convert a set of potentially correlated characteristics into a set of linearly uncorrelated features, principal component analysis (PCA) seeks to discover a rotation. Principal components are the columns employed in this orthogonal transformation. This approach is also intended for matrices with few properties.[5]

## V. CONCLUSION

Data preparation is recognized as a key stage in the fields of data analysis and machine learning. It serves as a gateway, transforming unprocessed data into a form that computers can understand and creating the groundwork for accurate and trustworthy analysis. Preprocessing improves data quality, which increases the accuracy of models and allows them to gain more insight from the data they analyze. Preprocessing is a crucial step in creating reliable and trustworthy machine learning models since it addresses difficulties like noisy or missing data. This article highlights the crucial role that data preprocessing plays in the process of

developing reliable prediction models in the field of machine learning.

## References

1. Ramírez-Gallego, S., Krawczyk, B., García, S., Woźniak, M., & Herrera, F. (2017). A survey on data preprocessing for data stream mining: Current status and future directions. Neurocomputing, 239, 39-57. https://doi.org/10.1016/j.neucom.2017.01.078

2. Kotsiantis, S. B., Kanellopoulos, D., & Pintelas, P. E. (2006). Data preprocessing for supervised leaning. International journal of computer science, 1(2), 111-117.

3. Rahman, A. (2019). Statistics-based data preprocessing methods and machine learning algorithms for big data analysis. International Journal of Artificial Intelligence, 17(2), 44-65.

4. Famili, A., Shen, W. M., Weber, R., & Simoudis, E. (1997). Data preprocessing and intelligent data analysis. Intelligent data analysis, 1(1), 3-23.

5. García, S., Ramírez-Gallego, S., Luengo, J. et al. Big data preprocessing: methods and prospects. Big Data Anal 1, 9 (2016). https://doi.org/10.1186/s41044-016-0014-0

6. Jain, A. K., & Gupta, B. B. (2022). A survey of phishing attack techniques, defence mechanisms and open research challenges. Enterprise Information Systems, 16(4), 527-565.

7. Gupta, S., & Gupta, B. B. (2015, May). PHP-sensor: a prototype method to discover workflow violation and XSS vulnerabilities in PHP web applications. In Proceedings of the 12th ACM international conference on computing frontiers (pp. 1-8).

8. Negi, P., Mishra, A., & Gupta, B. B. (2013). Enhanced CBF packet filtering method to detect DDoS attack in cloud computing environment. arXiv preprint arXiv:1304.7073.

9. Gupta, B. B., Misra, M., & Joshi, R. C. (2012). An ISP level solution to combat DDoS attacks using combined statistical based approach. arXiv preprint arXiv:1203.2400.

10. Chopra, M., Singh, S. K., Gupta, A., Aggarwal, K., Gupta, B. B., & Colace, F. (2022). Analysis & prognosis of sustainable development goals using big data-based approach during COVID-19 pandemic. Sustainable Technology and Entrepreneurship, 1(2), 100012.

11. Mishra, A., Gupta, N., & Gupta, B. B. (2023). Defensive mechanism against DDoS attack based on feature selection and multi-classifier algorithms. Telecommunication Systems, 82(2), 229-244.

12. Chai, Y., Qiu, J., Yin, L., Zhang, L., Gupta, B. B., & Tian, Z. (2022). From data and model levels: Improve the performance of few-shot malware classification. IEEE Transactions on Network and Service Management, 19(4), 4248-4261.