



# Emerging Trends in Data Structures for Modern Machine Learning Applications

Ikvinderpal Singh<sup>1,\*</sup>, Sapandeep Kaur Dhillon<sup>2</sup>

## Abstract

*In the realm of machine learning, data structures play a pivotal role in facilitating efficient data manipulation, storage, and retrieval, thereby significantly impacting the performance and scalability of machine learning algorithms. In recent years, the field of machine learning has witnessed the emergence of novel data structures tailored to address scalability and efficiency challenges inherent in handling large-scale and high-dimensional data. This study provides a look at the data preprocessing, data cleaning and feature extraction in machine learning using data structures. This study explores three key emerging trends: tensor representations, sparse data structures, and graph-based neural networks, highlighting their potential to revolutionize modern machine learning applications. Furthermore, we discussed a case study of the Leveraging Emerging Data Structures in Healthcare Analytics.*

**Keywords:** Machine learning, data structures, tensor representations, sparse data structures, graph-based neural networks

## INTRODUCTION

Machine learning algorithms rely heavily on efficient data representation and manipulation to extract meaningful patterns and make accurate predictions. Data structures serve as the foundational building blocks for organizing and managing data in a structured manner. Data structures play a fundamental role in machine learning, serving as the building blocks for organizing, storing, and manipulating data throughout the entire machine learning pipeline [1, 2]. From the initial preprocessing and cleaning of raw data to the final evaluation and interpretation of model predictions, data structures enable efficient data representation and manipulation. Commonly used data structures such as arrays, matrices, and tensors facilitate the storage and processing of structured data, while more specialized structures like trees, graphs, and hash tables are employed for tasks requiring hierarchical or associative representations. These data structures not only enable algorithms to efficiently access and operate on data but also facilitate feature engineering, model training, and evaluation by providing organized and

optimized data storage formats [3]. Additionally, advancements in data structure design, such as sparse matrices for handling high-dimensional and sparse data, and graph-based data structures for modeling complex relationships, continue to drive innovation and improve the scalability and efficiency of machine learning algorithms. In essence, data structures form the backbone of machine learning, enabling algorithms to leverage data effectively and make accurate predictions or decisions across a wide range of applications.

### \*Author for Correspondence

Ikvinderpal Singh  
E-mail: ips\_sikand@yahoo.com

<sup>1</sup>Assistant Professor, Department of Computer Science & Applications, Trai Shatabdi Guru Gobind Singh Khalsa College, Amritsar, Punjab, India

<sup>2</sup>Assistant Professor, Department of Computer Science, Guru Nanak Dev University, Amritsar, Punjab, India

Received Date: February 16, 2024

Accepted Date: February 19, 2024

Published Date: February 21, 2024

**Citation:** Ikvinderpal Singh, Sapandeep Kaur Dhillon. Emerging Trends in Data Structures for Modern Machine Learning Applications. International Journal of Data Structure Studies. 2024; 2(1): 1–7p.

## DATA PREPROCESSING AND CLEANING

Data preprocessing and cleaning are essential steps in the machine learning pipeline, aimed at

preparing raw data for model training and evaluation [4]. Data structures play a crucial role in facilitating efficient preprocessing and cleaning operations by providing organized and optimized storage formats for the data. Arrays, lists, and hash tables are commonly used data structures for storing and processing raw data during these initial stages. Arrays and matrices are particularly useful for representing tabular data, where rows correspond to samples and columns represent features. These data structures enable operations such as filtering, imputation of missing values, and normalization to be performed efficiently. Linked lists and hash tables are valuable for handling unstructured data, such as text or categorical variables, allowing for efficient storage and retrieval of information during preprocessing tasks like tokenization and encoding. Moreover, specialized data structures like trees and graphs can be employed for more complex data cleaning tasks, such as outlier detection and anomaly removal, where hierarchical or relational representations are required. By leveraging appropriate data structures, machine learning practitioners can streamline the data preprocessing and cleaning process, ensuring that the data is properly formatted and free from inconsistencies or errors before feeding it into the learning algorithm. This not only elevates the caliber of the training data but also boosts the overall effectiveness and dependability of the machine learning model.

### **FEATURE ENGINEERING**

Feature engineering is a crucial step in the machine learning pipeline, aimed at extracting informative and discriminative features from raw data to improve the performance of learning algorithms [5]. Data structures play a pivotal role in facilitating efficient feature engineering operations by providing organized and optimized formats for representing and manipulating data. Arrays, matrices, and tensors are commonly used data structures for storing and processing feature data, enabling operations such as transformation, aggregation, and combination to be performed efficiently. For example, arrays and matrices are suitable for representing tabular data, where each row corresponds to a sample and each column represents a feature. This allows for easy manipulation of feature values, such as scaling, normalization, or encoding categorical variables. Graph-based data structures are valuable for feature engineering in scenarios where data exhibits complex relationships or dependencies, such as social networks or molecular structures. By representing data as a graph, features can be extracted based on graph topology, connectivity, or node attributes, enabling the creation of informative feature representations for machine learning tasks [6]. Additionally, specialized data structures like trees and hash tables can be employed for feature extraction from structured or unstructured data sources, facilitating operations such as text tokenization, image segmentation, or time-series decomposition. Leveraging appropriate data structures for feature engineering not only streamlines the process of extracting relevant features from raw data but also enhances the effectiveness and interpretability of machine learning models by capturing the underlying structure and patterns present in the data.

### **TRENDS IN DATA STRUCTURES FOR MACHINE LEARNING**

Recent advancements in data structure design have led to the development of specialized data structures tailored specifically for machine learning tasks. We explore emerging trends such as tensor representations, sparse data structures, and graph-based neural networks, highlighting their potential to address scalability and efficiency challenges in modern machine learning applications.

### **TENSOR REPRESENTATIONS**

Tensor representations have gained traction in machine learning due to their ability to efficiently handle multi-dimensional data, such as images, videos, and time-series data. Tensors generalize matrices to higher dimensions, enabling the representation of complex relationships and interactions among features. By leveraging tensor representations, researchers can design more expressive and scalable models capable of capturing intricate patterns in large-scale datasets [7, 8].

Tensor decomposition techniques [6], such as Tucker decomposition and tensor factorization, facilitate dimensionality reduction and feature extraction in tensor data, leading to improved model efficiency and interpretability [9]. Moreover, tensor-based operations can be efficiently parallelized

across multiple processing units, enabling accelerated computations on modern hardware architectures like GPUs and TPUs.

### **Representation**

Tensors are arrays with multiple dimensions, extending the concept of matrices to higher-dimensional structures. Each dimension in a tensor corresponds to a different mode or feature, allowing for the representation of complex relationships and interactions among data elements.

### **Role**

Tensor representations play a crucial role in capturing and processing multi-dimensional data structures commonly encountered in machine learning tasks, such as images, videos, time-series data, and multi-modal data.

### **Working Principle**

Tensor operations involve manipulation of multi-dimensional arrays, including element-wise operations, matrix multiplications, and tensor contractions.

These operations aid in accomplishing diverse tasks, including extracting features, reducing dimensionality, and training models.

### **Importance**

By using tensor representations, machine learning models can effectively handle complex data structures and capture intricate relationships among features, enabling more expressive and scalable modeling capabilities.

### **Applications**

Tensor representations are widely used in tasks such as image classification, natural language processing, recommender systems, and sensor data analysis, where the input data exhibits multi-dimensional characteristics.

### **Example**

In image classification, a convolutional neural network (CNN) processes input images as multi-dimensional tensors. Convolutional layers perform convolutions on input tensors to extract spatial features, while pooling layers aggregate features to reduce dimensionality. Finally, fully connected layers process flattened tensors to perform classification.

## **SPARSE DATA STRUCTURES**

Sparse data structures offer an efficient solution for handling datasets with a large number of features or attributes, where the majority of entries are zero or negligible. In many real-world applications, such as text mining, genomics, and recommender systems, data sparsity poses a significant challenge to traditional dense representations, leading to increased memory consumption and computational overhead.

Sparse data structures, such as compressed sparse row (CSR) and compressed sparse column (CSC) matrices, exploit the inherent sparsity of data to store and manipulate only non-zero entries, thereby reducing memory footprint and speeding up computation [10, 11]. Additionally, specialized data structures like hash maps and inverted indices enable efficient retrieval and manipulation of sparse data, facilitating fast feature lookup and model inference.

### **Representation**

Sparse data structures store data in a compressed format, focusing only on non-zero or significant entries. This approach reduces memory consumption and accelerates computations for datasets with a large number of zero or negligible values.

**Role**

Sparse data structures provide an efficient means of representing and processing datasets with a large number of zero or negligible values, reducing memory consumption and computational overhead.

**Importance**

Sparse data structures enable machine learning algorithms to handle high-dimensional data efficiently, especially in applications where the data is inherently sparse, such as text data, genomic data, and high-dimensional feature spaces.

**Working Principle**

Sparse matrices are commonly expressed through formats such as compressed sparse row (CSR) or compressed sparse column (CSC), where solely the non-zero entries, along with their corresponding row and column indices, are stored. This enables efficient storage, manipulation, and multiplication of sparse matrices.

**Applications**

Sparse data structures find applications in tasks such as text mining, document classification, collaborative filtering, and feature engineering, where the input data exhibits sparsity characteristics.

**Example**

In natural language processing, document-term matrices representing word occurrences in documents are often sparse due to the large vocabulary size and sparse nature of text data. Sparse matrices allow efficient storage and computation of term-document statistics for tasks like document classification and clustering.

**GRAPH-BASED NEURAL NETWORKS**

Graph-based neural networks (GNNs) have emerged as a powerful framework for modeling structured data with complex relationships and dependencies, such as social networks, knowledge graphs, and molecular structures. In contrast to conventional neural networks that function on grid-like data structures, Graph Neural Networks (GNNs) directly process data organized in graph structures, enabling greater flexibility and expressive modeling capabilities.

GNNs leverage message passing algorithms to propagate information across nodes and edges in a graph, enabling the integration of local and global context into node representations [12, 13]. GNNs are particularly apt for tasks like node classification, link prediction, and graph generation due to this characteristic. Moreover, recent advancements in graph attention mechanisms and graph convolutional networks have further improved the scalability and efficiency of GNNs, enabling their application to large-scale graph datasets.

**Representation**

Graph-based neural networks (GNNs) function directly on data organized in graph structures, wherein nodes denote entities and edges signify relationships or interactions among entities. GNNs disseminate information across nodes and edges to acquire knowledge about node representations.

**Role**

The purpose of Graph-based neural networks (GNNs) is to model and analyze structured data depicted as graphs, with nodes representing entities and edges denoting relationships or interactions among these entities.

**Working Principle**

GNNs utilize message passing algorithms to gather information from neighboring nodes and iteratively revise node representations [14]. This iterative process empowers GNNs to grasp both local and global graph structures, facilitating tasks like node classification, link prediction, and graph generation.

### **Importance**

GNNs enable machine learning models to capture complex dependencies and interactions among entities in graph-structured data, making them well-suited for tasks involving relational data, social networks, knowledge graphs, and molecular structures.

### **Applications**

GNNs are applied in various tasks including node classification, link prediction, graph generation, recommendation systems, and social network analysis, particularly when the data demonstrates a structure resembling that of a graph.

### **Example**

Within social network analysis, GNNs have the capability to depict relationships among users in the form of a graph, with users represented as nodes and connections or interactions represented by edges. By propagating information across the social graph, GNNs can predict user attributes, identify communities, and detect anomalous behavior.

## **CASE STUDY: LEVERAGING EMERGING DATA STRUCTURES IN HEALTHCARE ANALYTICS**

In the domain of healthcare analytics, the analysis of patient data often involves handling large-scale, multi-dimensional datasets with complex relationships and sparsity. In this case study, we explore how emerging data structures, tensor representations, sparse data structures, and graph-based neural networks, can be leveraged to address challenges in healthcare analytics.

### **Case Scenario: Tensor Representations**

A hospital collects multi-modal patient data, including electronic health records (EHRs), medical images (e.g., MRI, CT scans), and time-series physiological measurements (e.g., heart rate, blood pressure). Analyzing such heterogeneous data requires a unified representation capable of capturing the complex relationships among different data modalities.

#### ***Applications of Tensor Representations***

By representing patient data as a multi-dimensional tensor, where each dimension corresponds to a different data modality (e.g., patient demographics, medical history, image pixels, physiological signals), healthcare practitioners can leverage tensor operations to perform tasks such as:

- Feature extraction: Extracting relevant features from EHRs, medical images, and physiological signals using tensor decomposition techniques (e.g., Tucker decomposition).
- Disease prediction: Building predictive models to diagnose diseases or predict patient outcomes based on multi-modal patient data represented as tensors.
- Treatment recommendation: Recommending personalized treatment plans by analyzing the interactions between different data modalities captured in the tensor representation.

### **Case Scenario: Sparse Data Structures**

A healthcare organization maintains a large database of patient electronic health records (EHRs), where each record contains a wide range of clinical variables. However, many of these variables are sparse, with most entries being zero or missing, posing challenges for traditional dense data representations.

#### ***Applications of Sparse Data Structures***

By employing sparse data structures, such as compressed sparse row (CSR) matrices, healthcare analysts can efficiently handle sparse EHR data by:

- Storing and manipulating patient EHRs in a compressed format, reducing memory consumption and speeding up computations.

- Performing tasks such as patient similarity analysis, disease clustering, and anomaly detection using sparse matrix operations.
- Integrating sparse EHR data with other data sources (e.g., genomic data, medical imaging) to enrich patient profiles and improve predictive modeling accuracy.

### **Case Scenario: Graph-Based Neural Networks**

A healthcare network aims to improve patient care coordination and resource allocation by analyzing the referral patterns among healthcare providers, hospitals, and medical facilities. The representation of the network can be in the form of a graph, where healthcare entities (such as providers and hospitals) are depicted as nodes, and edges symbolize referral relationships.

### **Applications of Graph-Based Neural Networks**

By employing graph-based neural networks (GNNs), healthcare administrators can:

- Model the referral network as a graph and learn embeddings for healthcare entities using GNNs.
- Identify key influencers, communities, and referral patterns within the healthcare network.
- Optimize patient referral pathways, allocate resources effectively, and improve patient outcomes based on insights derived from GNN analysis.

### **CONCLUSION**

These emerging data structures in machine learning operate on different principles tailored to address specific challenges associated with large-scale and high-dimensional data. Tensor representations enable efficient handling of multi-dimensional data, sparse data structures optimize storage and computation for datasets with many zero values, and graph-based neural networks capture complex relationships in structured data. By understanding the working principles of these data structures, researchers and practitioners can effectively leverage them to design more scalable, efficient, and expressive machine learning models for various real-world applications. In healthcare analytics, leveraging emerging data structures such as tensor representations, sparse data structures, and graph-based neural networks offers significant opportunities to extract actionable insights from complex, multi-modal, and interconnected patient data. Through the implementation of these sophisticated methods, healthcare institutions can improve patient care, optimize the allocation of resources, and stimulate innovations in the delivery and management of healthcare.

### **REFERENCES**

1. UNB. (2023). Datasets: Research. Canadian Institute for Cybersecurity. [Online]. Available from: <https://www.unb.ca/cic/datasets/index.html>
2. UNB. (2019). DDoS 2019 Datasets: Research: Canadian Institute for Cybersecurity. [Online]. Available from: <https://www.unb.ca/cic/datasets/ddos-2019.html>
3. Boukerche A, Wang J. Machine learning-based traffic prediction models for intelligent transportation systems. *Comput Netw.* 2020 Nov 9; 181: 107530.
4. Ardabili SF, Mosavi A, Ghamisi P, Ferdinand F, Varkonyi-Koczy AR, Reuter U, Rabczuk T, Atkinson PM. Covid-19 outbreak prediction with machine learning. *Algorithms.* 2020 Oct 1; 13(10): 249.
5. Chollet F. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2017; 1251–1258.
6. Sai Bhargav K, Rajendra V. (2022). Data Structures for Machine Learning. [Online]. Engpaper. Available from: <https://www.engpaper.com/data-structures-for-machine-learning.htm>
7. Koniusz P, Wang L, Cherian A. Tensor representations for action recognition. *IEEE Trans Pattern Anal Mach Intell.* 2021 Aug 24; 44(2): 648–65.
8. Hua C, Rabusseau G, Tang J. High-order pooling for graph neural networks with tensor decomposition. *Adv Neural Inf Process Syst.* 2022 Dec 6; 35: 6021–33.
9. Kuang L, Hao F, Yang LT, Lin M, Luo C, Min G. A tensor-based approach for big data representation and dimensionality reduction. *IEEE Trans Emerg Topics Comput.* 2014 Jun 12; 2(3): 280–91.

10. Harris CR, Millman KJ, Van Der Walt SJ, Gommers R, Virtanen P, Cournapeau D, Wieser E, Taylor J, Berg S, Smith NJ, Kern R. Array programming with NumPy. *Nature*. 2020 Sep 17; 585(7825): 357–62.
11. Niemeyer J, Rottensteiner F, Soergel U. Contextual classification of lidar data and building object detection in urban areas. *ISPRS J Photogramm Remote Sens*. 2014 Jan 1; 87: 152–65.
12. Sadeghi M, Babaie-Zadeh M, Jutten C. Dictionary learning for sparse representation: A novel approach. *IEEE Signal Process Lett*. 2013 Oct 9; 20(12): 1195–8.
13. Allamanis M, Brockschmidt M, Khademi M. Learning to represent programs with graphs. *arXiv preprint arXiv:1711.00740*. 2017 Nov 1.
14. Dettmers T, Minervini P, Stenetorp P, Riedel S. Convolutional 2d knowledge graph embeddings. In *Proceedings of the AAAI conference on artificial intelligence*. 2018 Apr 25; 32(1): 1811–1818.