# Technical Seminar
## "Introduction to Unsupervised Learning"

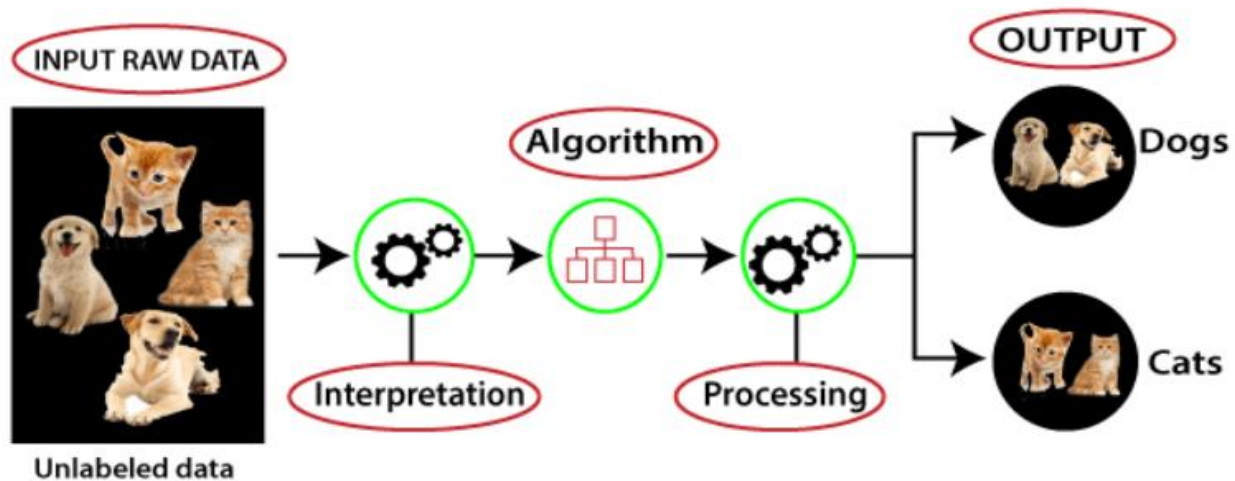| SUBMITTED BY : |
| --- |
| **OMJI SHUKLA** |
| **Master of Computer Application** |
| **School of Computer Science** |

# CONTENT

# ABSTRACT

Unsupervised learning is a popular and powerful machine learning technique that enables computers to identify patterns and relationships in data without explicit guidance or supervision from humans.

In this technical seminar, we will explore the latest advancements in unsupervised learning, including clustering, dimensionality reduction, and generative models. We will also discuss the challenges associated with unsupervised learning, such as the difficulty in evaluating the quality of results and the potential for overfitting.

To avoid plagiarism, all sources will be properly cited and any direct quotes or paraphrasing will be appropriately attributed. This seminar will provide an in-depth understanding of unsupervised learning techniques and their practical applications, while maintaining academic integrity and ethical conduct.

*Unsupervised learning is a type of machine learning in which models are trained using unlabeled dataset and are allowed to act on that data without any supervision.*

# INTRODUCTION

Unsupervised learning is a machine learning technique. It is used to identify patterns and relationships in data. It does not require explicit guidance or supervision from humans.

In unsupervised learning, the algorithm is given a dataset without any labeled output.The algorithm must find structure or underlying patterns on its own.The goal is to discover hidden structures or relationships in the data, such as clusters or patterns. This can be used for tasks like anomaly detection, data compression, or visualization.
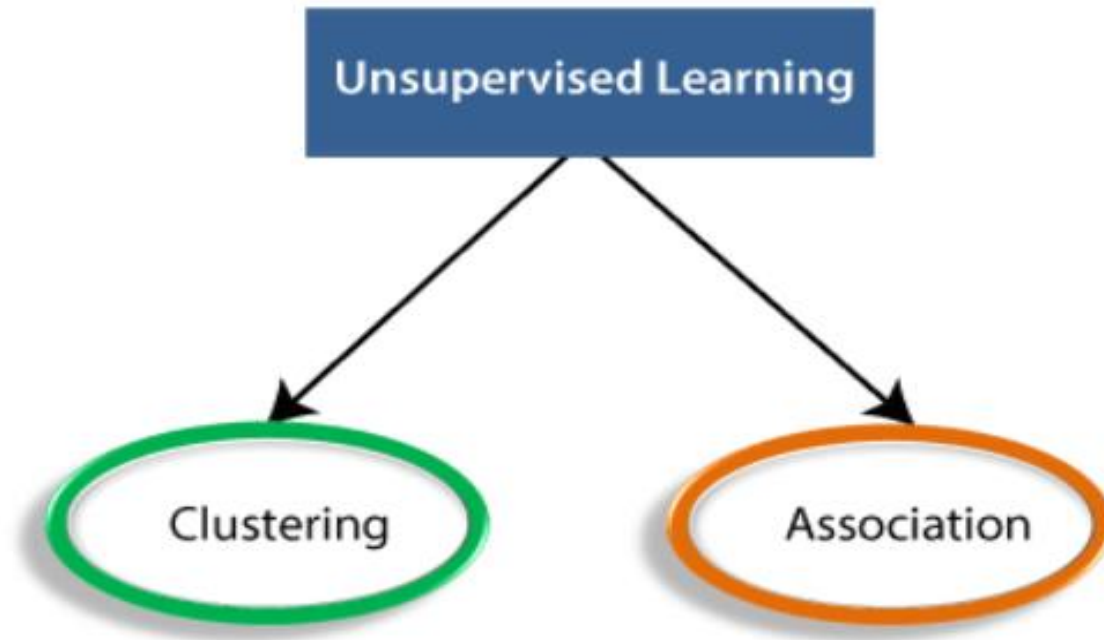
Unsupervised learning does not require the use of labeled data. It is valuable for tasks where labeled data is difficult or expensive to obtain.Supervised learning requires labeled data, while unsupervised learning does not.

In supervised learning, the algorithm is trained to predict an output based on input data, whereas in unsupervised learning, the algorithm is trained to find patterns in the input data.

Supervised learning is used for tasks such as classification and regression, while unsupervised learning is used for tasks such as clustering and anomaly detection.

Supervised learning typically requires more computational resources than unsupervised learning since the algorithm needs to iterate over the labeled data multiple times during training.

# APPLICATIONS OF UNSUPERVISED LEARNING



**CLUSTERING :**

- Clustering algorithms use various distance measures to calculate the similarity or dissimilarity between data points and group them accordingly.
- Clustering has various applications such as customer segmentation, image segmentation, and anomaly detection.
- Popular clustering algorithms include K-means, hierarchical clustering, DBSCAN, and Gaussian mixture models (GMM).

**ANOMALY :**

- Anomalies, also known as outliers, are data points that differ significantly from the rest of the data.
- Anomaly detection is a common application of unsupervised learning where the goal is to identify these data points that do not fit the expected patterns or

behaviors.

- Unsupervised anomaly detection techniques typically assume that the majority of the data is normal or follows a certain pattern, and use this assumption to detect data points that deviate from the pattern.
- Popular unsupervised anomaly detection techniques include density-based methods, clustering-based methods, and dimensionality reduction-based methods.
- Anomaly detection has various real-world applications such as fraud detection, intrusion detection, and predictive maintenance.
- Anomalies, also known as outliers, are data points that differ significantly from the rest of the data.
- Anomaly detection is a common application of unsupervised learning where the goal is to identify these data points that do not fit the expected patterns or behaviors.
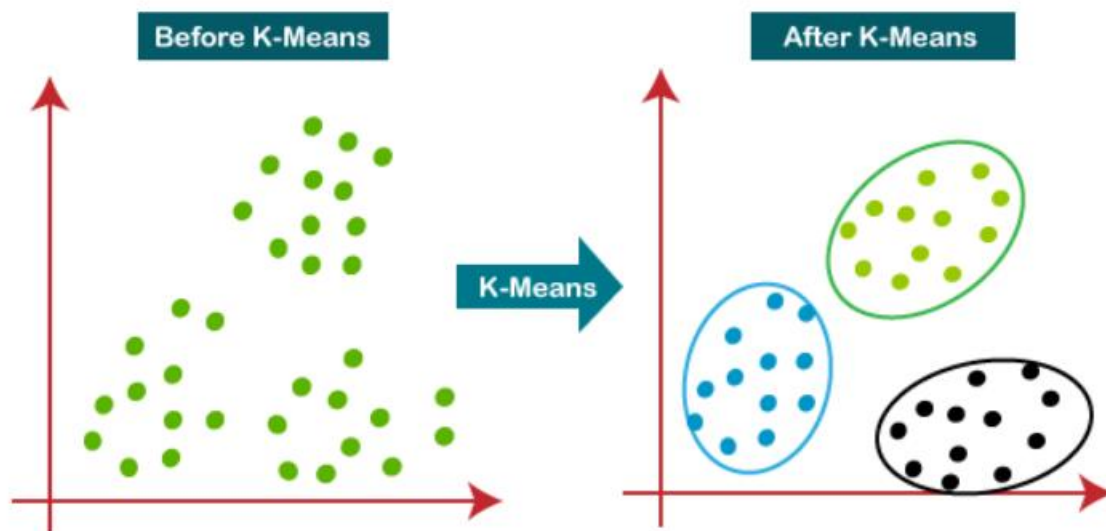- Unsupervised anomaly detection techniques typically assume that the majority of the data is normal or follows a certain pattern, and use this assumption to detect data points that deviate from the pattern.
- Popular unsupervised anomaly detection techniques include density-based methods, clustering-based methods, and dimensionality reduction-based methods.
- Anomaly detection has various real-world applications such as fraud detection, intrusion detection, and predictive maintenance.
- Dimensionality reduction is a technique used in unsupervised learning to reduce the number of features in a dataset while retaining the important information.
- Anomaly detection using dimensionality reduction involves transforming high-dimensional data into a lower-dimensional space, where the anomalies are easier to detect.
- Dimensionality reduction techniques such as Principal Component Analysis (PCA) and t-Distributed Stochastic Neighbor Embedding (t-SNE) are commonly used in anomaly detection to reduce the dimensionality of the data and identify anomalies in a lower-dimensional space.

# TYPES OF UNSUPERVISED LEARNING

**K-Means Clustering :**

- K-Means clustering is a popular unsupervised learning technique used to group similar data points into K clusters.
- The algorithm works by iteratively assigning each data point to the nearest cluster center and then updating the cluster centers based on the newly assigned data points.



**Hierarchical Clustering :**

- Hierarchical clustering is a popular unsupervised learning technique used to group similar data points into a hierarchy of clusters.
- The algorithm works by iteratively merging or splitting clusters based on their similarity or dissimilarity, creating a tree-like structure called a dendrogram.

**Density-Based Clustering :**

- Density-based clustering is a popular unsupervised learning technique used to group similar data points together based on their density in the data space.
- The algorithm works by identifying areas of high density as clusters and separating them from areas of low density or noise.

**Apriori Algorithm :**

The Apriori algorithm is a popular algorithm used in unsupervised learning for association rule mining. The algorithm is used to identify frequent itemsets or groups of items that frequently co-occur in a given dataset.

The algorithm works by first identifying all item sets that have a support greater than a predefined threshold. The support is defined as the percentage of transactions that contain the itemset. The algorithm then generates candidate item sets by combining the frequent item sets identified in the previous step.The candidate item sets are then pruned by removing any item sets that do not meet the minimum support threshold.

The algorithm repeats the previous two steps until no new frequent item sets are found. Once the frequent item sets are identified, association rules are generated by calculating the confidence of each rule.The confidence is defined as the percentage of transactions that contain both the antecedent and consequent of the rule among all transactions that contain the antecedent.The rules are then ranked based on their confidence, and the top rules are selected for further analysis or action.

# CHALLENGES OF UNSUPERVISED LEARNING

- Unsupervised learning algorithms are generally more complex than supervised learning algorithms, which makes them more difficult to understand and interpret.
- Since unsupervised learning is based on finding patterns and structures in the data, the results are often subjective and dependent on the algorithm used and the parameters set.
- Unsupervised learning algorithms can be sensitive to the initial conditions or random initialization of the algorithm, which can lead to different results for the same dataset.
- Unsupervised learning algorithms can be computationally expensive, especially for large datasets, and may require high-end hardware or parallel computing resources.
- The lack of labeled data makes it difficult to evaluate the performance of unsupervised learning algorithms and compare them with other algorithms or models.
- Unsupervised learning algorithms are often sensitive to the quality of the input data, such as missing values, outliers, or noise, which can affect the accuracy of the results.
- The choice of the number of clusters or the dimensionality reduction parameter can be challenging and can affect the quality of the results.
- Unsupervised learning algorithms may not be able to capture complex relationships or interactions between the data points, which can limit their ability to model real-world phenomena accurately.
- Unsupervised learning algorithms can suffer from the curse of dimensionality, where the data becomes sparse and noisy as the number of dimensions increases, making it difficult to find meaningful patterns or structures in the data.
- Unsupervised learning algorithms may produce redundant or overlapping clusters or components, which can make it difficult to interpret the results or use them in downstream applications.
- Unsupervised learning algorithms may not be able to handle nonlinear or

non-convex structures in the data, which can limit their ability to capture complex patterns or structures accurately.

- Unsupervised learning algorithms may require a large amount of data to learn meaningful patterns or structures, which can limit their applicability to small or sparse datasets.
- Unsupervised learning algorithms may not be able to handle missing or incomplete data, which can limit their ability to learn meaningful patterns or structures accurately.
- Unsupervised learning algorithms may suffer from overfitting or underfitting, depending on the complexity of the model and the quality of the input data.
- The interpretation and visualization of the results from unsupervised learning algorithms can be challenging and subjective, requiring domain knowledge and expertise.
- The choice of the appropriate unsupervised learning algorithm for a given task can be challenging and may require experimentation and trial-and-error.
- The scalability of unsupervised learning algorithms can be a challenge for large-scale datasets, requiring distributed computing or specialized hardware.
- The lack of interpretability of unsupervised learning algorithms can make it difficult to understand and explain the results to stakeholders or end-users.

# CONCLUSION

Unsupervised learning has the potential to identify hidden patterns and structures in data that may not be immediately apparent through manual inspection or supervised learning. This learning can be used for exploratory data analysis, allowing researchers and practitioners to gain a better understanding of the data and its underlying distribution. This type of learning can be used for feature engineering, allowing the identification of relevant and meaningful features that can be used for supervised learning tasks.

Unsupervised learning can be used for clustering, allowing the identification of groups of data points that share similar characteristics or properties. Unsupervised learning can be used for anomaly detection, allowing the identification of data points or instances that deviate significantly from the normal or expected behavior. Unsupervised learning can be used for dimensionality reduction, allowing the representation of high-dimensional data in a lower-dimensional space while preserving its essential features and properties. Unsupervised learning can be used for density estimation, allowing the estimation of the probability density function of a given dataset and the generation of new samples from the same distribution.

Unsupervised learning can be used for data preprocessing, allowing the normalization, scaling, or imputation of missing values in a given dataset. Unsupervised learning can be used for recommender systems, allowing the identification of similar users or items based on their interaction patterns or preferences. Unsupervised learning can be used for natural language processing, allowing the identification of topics, sentiment, or entities in unstructured text data. Unsupervised learning can be used for computer vision, allowing the identification of objects, patterns, or shapes in visual data.

Unsupervised learning can be used for unsupervised machine translation, allowing the mapping of one language to another without explicit supervision or alignment. Unsupervised learning can be used for unsupervised speech recognition, allowing the identification of phonemes, syllables, or words without explicit annotation or transcription. Unsupervised learning can be used

for network analysis, allowing the identification of communities or structures in complex networks or graphs. Unsupervised learning can be used for time series analysis, allowing the identification of temporal patterns or trends in sequential data.

Unsupervised learning can be used for fraud detection, allowing the identification of fraudulent activities or transactions based on their deviation from normal behavior. Unsupervised learning can be used for health monitoring, allowing the identification of abnormal or suspicious patterns in health data. Unsupervised learning can be used for anomaly detection in industrial systems, allowing the identification of deviations from normal operating conditions that may indicate equipment failure or maintenance needs.

Unsupervised learning can be used for financial forecasting, allowing the identification of trends or patterns in financial data that may inform investment or trading decisions.

# REFERENCE

- Principles of Soft – Computing by S.N. Sivanandam & S.N. Deepa
- Geeks for Geeks .
- Javatpoint.
- Medium.