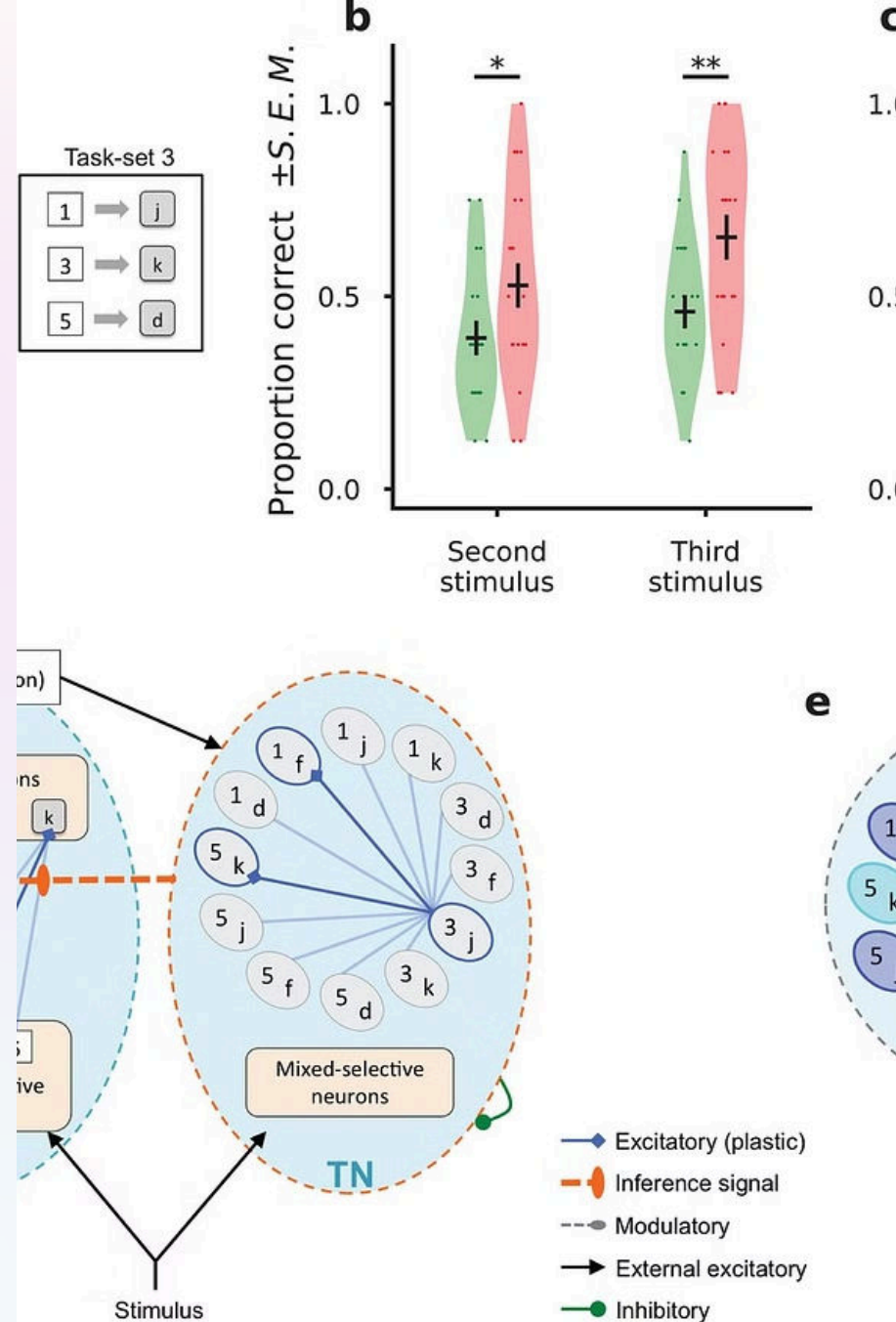# Introduction to Unsupervised Learning

Unsupervised learning is a type of machine learning used to draw inferences from datasets consisting of input data without labeled responses. It explores patterns, relationships, and structures within the data without guidance. This approach enables the discovery of hidden insights and valuable information.

**MA** **by Mvurya Mgala**

# What is unsupervised learning?

### Data Exploration

Unsupervised learning involves exploring data without specific targets or labels.

### Pattern Recognition

It focuses on identifying patterns and relationships within the dataset.

### Cluster Analysis

One key aspect is grouping data points into clusters based on similarities.

# Importance of Unsupervised Learning in Machine Learning

- **Identification of patterns:** Unsupervised learning helps in identifying hidden patterns within data, which can provide valuable insights for decision-making.

- **Data exploration:** It allows for exploration of the data structure, revealing associations and dependencies that are not immediately evident.

- **Feature extraction:** This technique aids in extracting important features from the data, leading to improved model performance in supervised learning tasks.

# Clustering algorithms in unsupervised learning

**1** — **K-means clustering**

An iterative algorithm that partitions data into k distinct clusters based on their attributes.

**2** — **Hierarchical clustering**

Technique that creates a tree of clusters to represent the arrangement of data.

**3** — **Dimensionality reduction**

A method to reduce the number of random variables in large datasets while retaining crucial information.

# Overview of Clustering Algorithms

### K-means Clustering

Partitions data into k clusters based on centroids.

Useful for large datasets and continuous variables.

May produce different results with different initializations.

### Hierarchical Clustering

Clusters data based on the hierarchy of partitions.

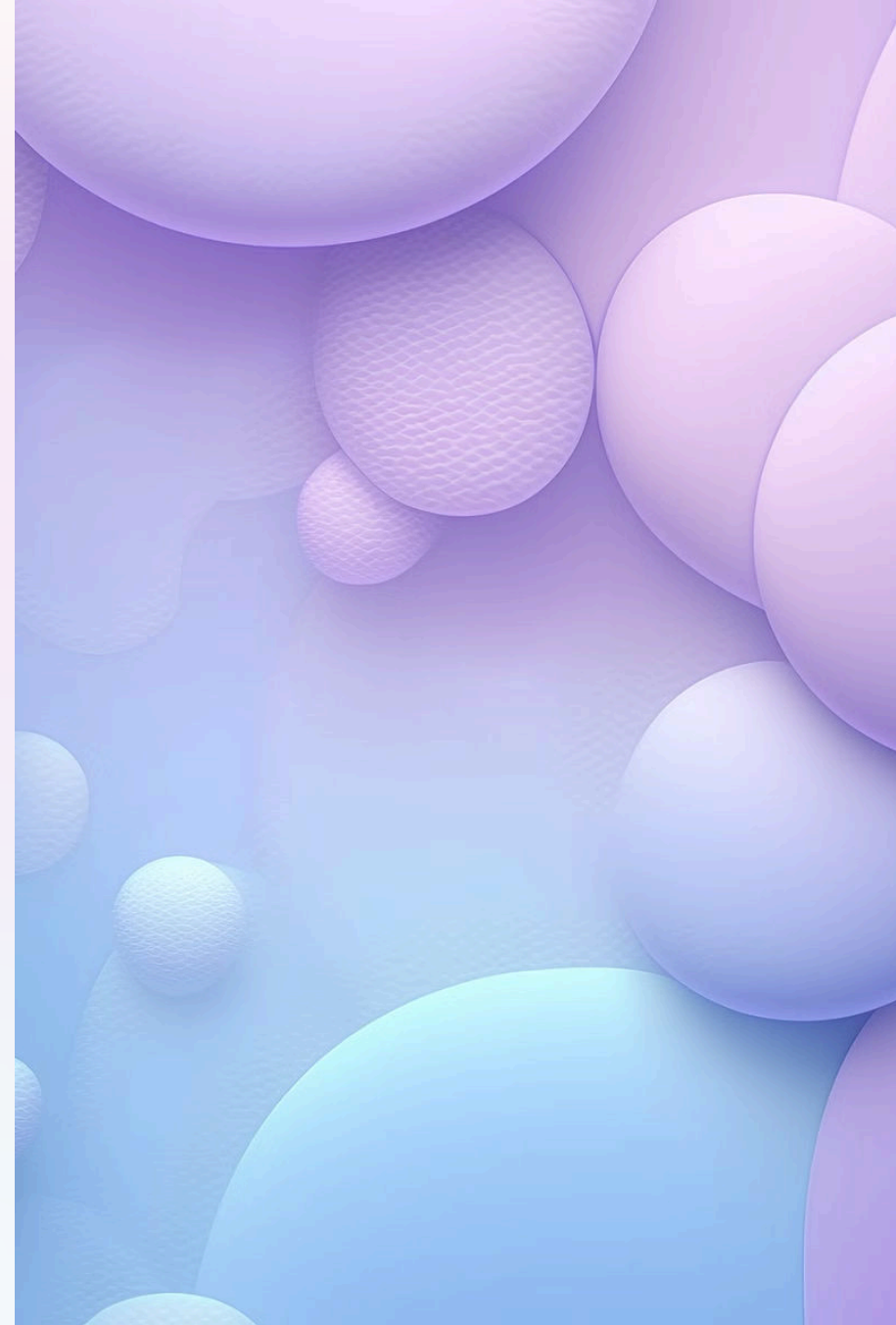Can be agglomerative (bottom-up) or divisive (top-down).

Works well for small to medium-sized datasets.

# K-means clustering algorithm

The K-means clustering algorithm is a fundamental unsupervised learning method used to partition data into clusters based on similarities. It aims to minimize the variance within each cluster and maximize the variance between clusters.

This iterative algorithm assigns data points to the nearest cluster center and updates the center based on the mean of the assigned points. It is widely used in various fields such as image segmentation and market segmentation.

For a visual representation, please refer to **this link**.

# How does K-means clustering work?

### Centroids

Initial centroids are randomly selected within the data space.

### Distance Calculation

Each point is assigned to the nearest centroid based on distance.
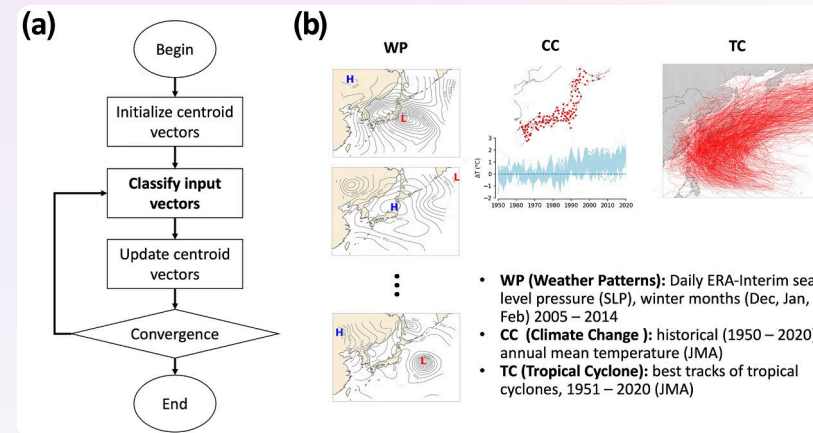
### Cluster Formation

Clusters are formed as points are reassigned to centroids.

# Advantages and Limitations of K-means Clustering

K-means clustering is computationally efficient and easy to implement.

However, it requires the number of clusters to be specified in advance.



(a) Flowchart: Begin → Initialize centroid vectors → **Classify input vectors** → Update centroid vectors → Convergence → End

(b) WP, CC, TC

- **WP (Weather Patterns):** Daily ERA-Interim sea level pressure (SLP), winter months (Dec, Jan, Feb) 2005 – 2014
- **CC (Climate Change ):** historical (1950 – 2020) annual mean temperature (JMA)
- **TC (Tropical Cyclone):** best tracks of tropical cyclones, 1951 – 2020 (JMA)
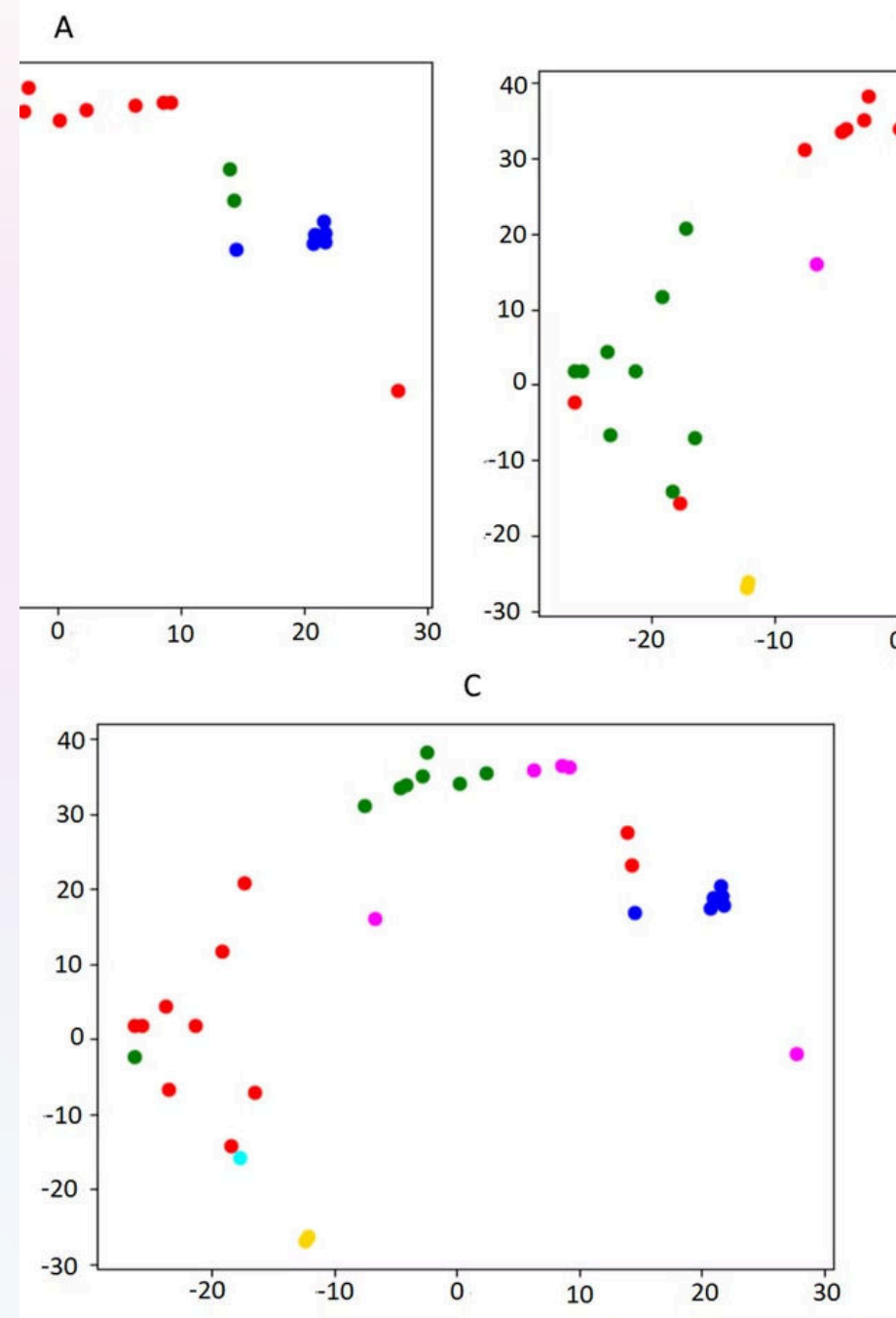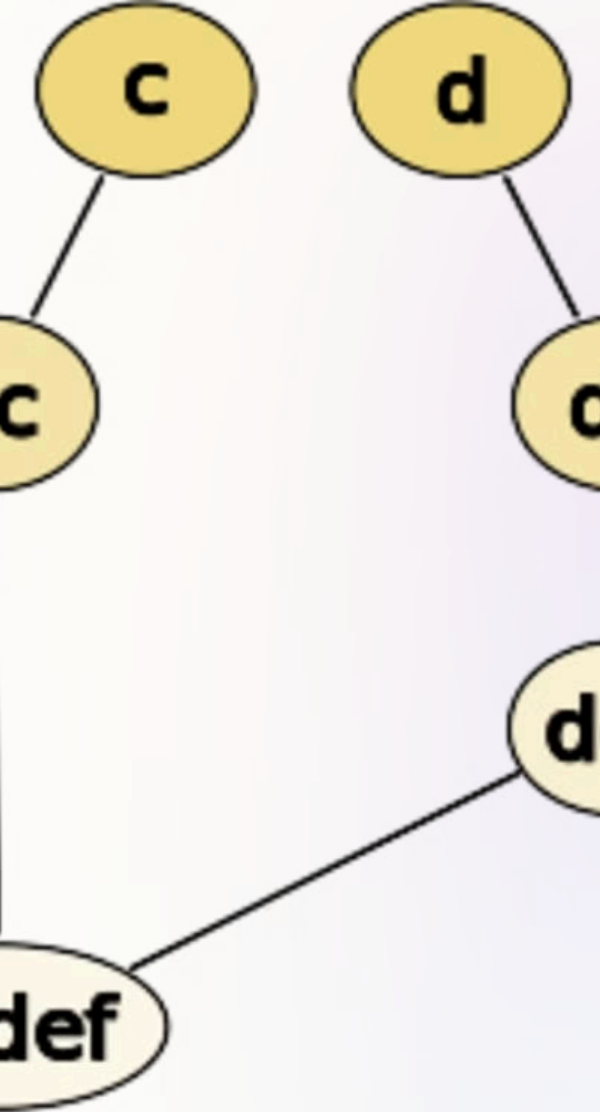
# Hierarchical Clustering Algorithm

Hierarchical clustering is a method of cluster analysis which seeks to build a hierarchy of clusters. It does not require the specification of the number of clusters and visualizes the clustering process through a dendrogram.

It can be implemented using different linkage methods such as single, complete, average, and ward. The process involves recursively merging or splitting clusters.

# How does hierarchical clustering work?

**Step 1**

**1** Start with each data point as a single cluster.

**Step 2**

**2** Merge the closest clusters, creating a new cluster.

**Step 3**

**3** Repeat until all points belong to a single cluster.

# Types of Hierarchical Clustering

**1**

### Agglomerative Clustering

Bottom-up approach, merging data points into clusters

**2**

### Divisive Clustering

Top-down approach, dividing clusters into smaller ones

# Advantages and limitations of hierarchical clustering

### Advantages

Hierarchical clustering can reveal the overall hierarchy of clusters and subclusters, providing a clear understanding of data structure.
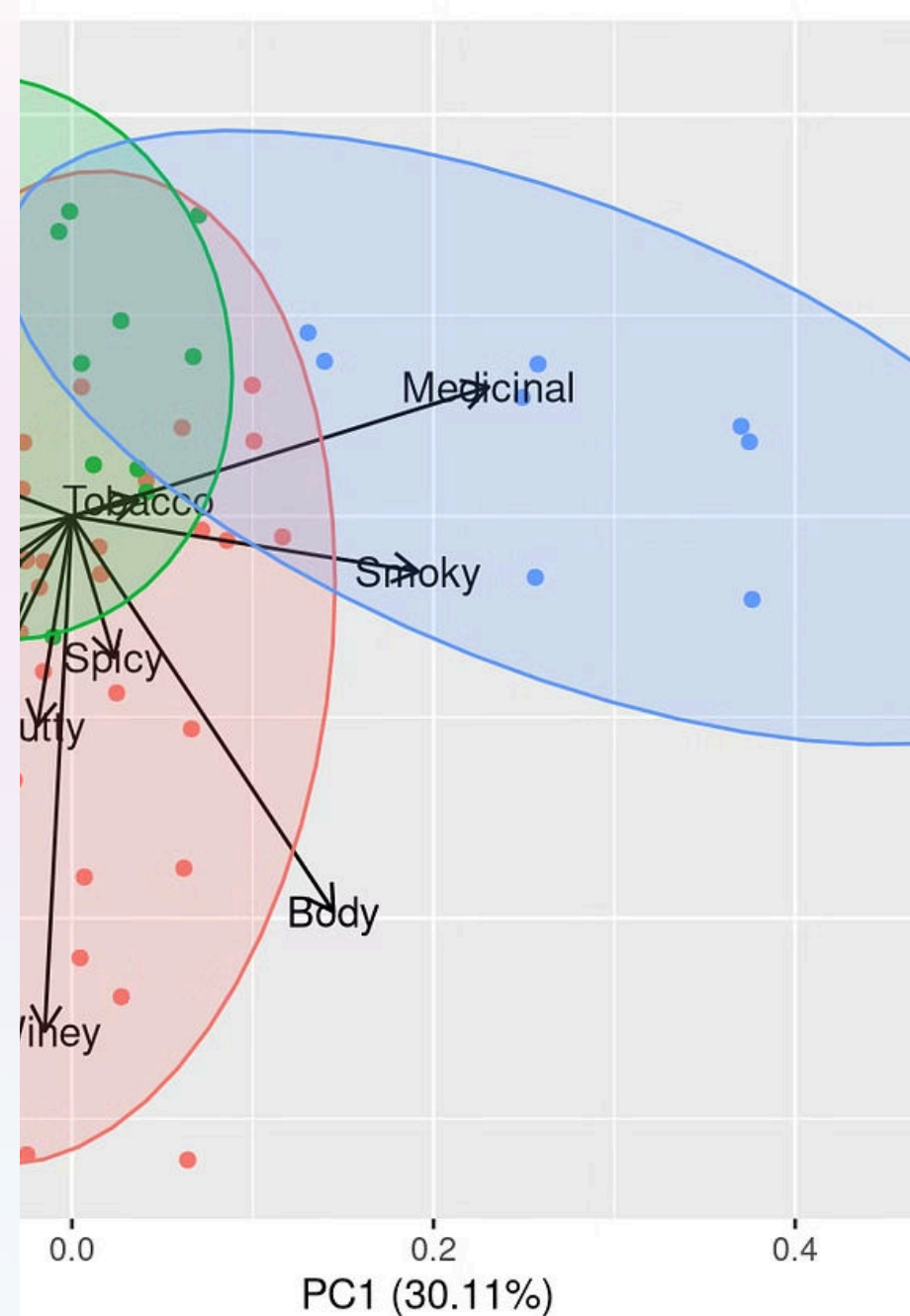
### Limitations

It can be computationally expensive and sensitive to outliers, impacting its scalability in large datasets.

# Dimensionality Reduction in Unsupervised Learning

Dimensionality reduction techniques are used to reduce the number of features in a dataset while preserving its key characteristics. By transforming high-dimensional data into a lower-dimensional space, algorithms like PCA and t-SNE facilitate easier visualization and analysis, leading to improved model performance.

One example is the ability to identify meaningful patterns and relationships within data that would be challenging to discern in its original high-dimensional form.

# Overview of Dimensionality Reduction Techniques

| Technique | Description | Applications | Advantages | Limitations |
| --- | --- | --- | --- | --- |

# Principal Component Analysis (PCA)

### How does PCA work?

PCA identifies the directions of maximum variance in a high-dimensional dataset and projects it onto a new coordinate system.

### Applications of PCA in machine learning

PCA is used for data compression, visualization, and noise reduction in image processing and feature extraction in pattern recognition.

### Advantages and limitations of PCA

PCA reduces data dimensionality but may lose some information. It's best for linearly separable data but may not work well for non-linear relationships.

# How does PCA work?

Principal Component Analysis (PCA) is a statistical technique used to simplify and interpret complex datasets. It works by transforming the data into a new coordinate system, where the greatest variance lies along the first axis, the second greatest variance along the second axis, and so on.

This transformation allows the data to be visualized in a lower-dimensional space while retaining the most important information. It helps in identifying patterns, reducing noise, and improving the speed and efficacy of machine learning algorithms.

# Applications of PCA in Machine Learning

### Data Visualization

PCA helps visualize high-dimensional data for better understanding and insights.

### Feature Selection

It assists in identifying the most significant features for model development.

### Noise Reduction

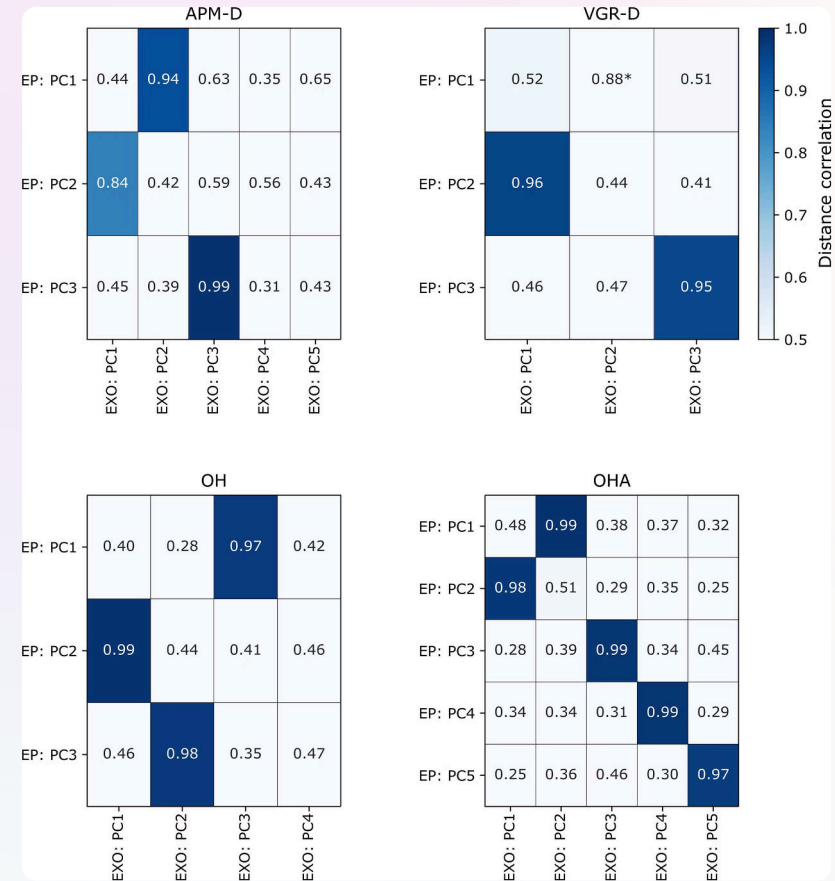PCA reduces noise in data, enhancing model accuracy and performance.

# Advantages and Limitations of PCA

Principal Component Analysis (PCA) helps in reducing the dimensionality of data and identifying important features.

It is widely used for data visualization and feature extraction in machine learning.

However, PCA assumes linear relationships, which may not be suitable for non-linear data.

It can also be sensitive to outliers and noise in the data, impacting its performance.

# t-SNE (t-Distributed Stochastic Neighbor Embedding)

**1**

### Non-linear Dimensionality Reduction

t-SNE is effective in capturing non-linear structures in high-dimensional data.
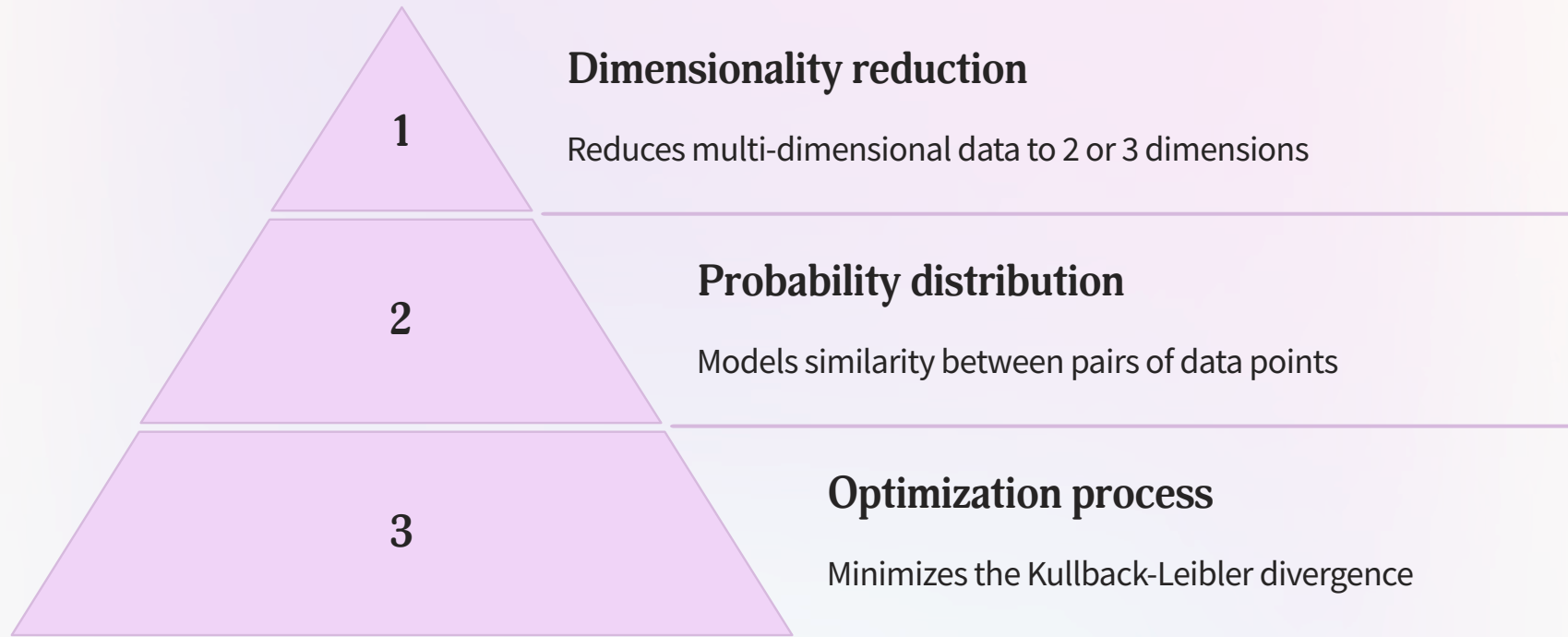
**2**

### Cluster Separation

It's useful for visualizing clusters and identifying patterns in complex datasets.

**3**

### Preserving Proximity

t-SNE preserves the local structure of data points, maintaining their relationships.

# How does t-SNE work?

**Dimensionality reduction**

Reduces multi-dimensional data to 2 or 3 dimensions

**Probability distribution**

Models similarity between pairs of data points

**Optimization process**

Minimizes the Kullback-Leibler divergence

1

2

3

# Applications of t-SNE in Machine Learning

**1** **Data Visualization**

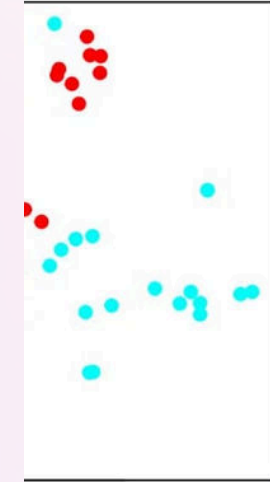t-SNE is used to visualize high-dimensional data in lower dimensions for easier interpretation.

**2** **Feature Extraction**

t-SNE helps in identifying important features by visualizing their relationships in reduced dimensions.
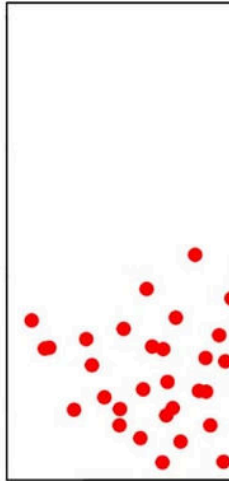
**3** **Clustering Validation**

t-SNE aids in validating the results of clustering algorithms by visualizing the clusters in 2D or 3D space.
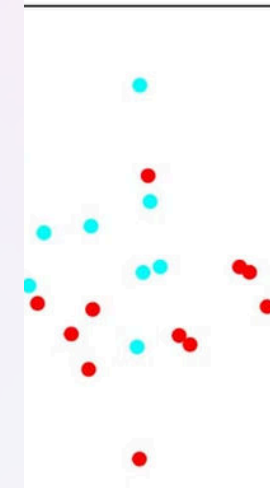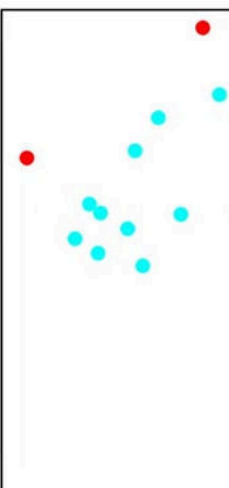
# Advantages and limitations of t-SNE

| t-SNE Advantages | t-SNE Limitations |
| --- | --- |
| Preserves local structure in data visualization. | Computationally expensive for large datasets. |
| Effective in revealing clusters and outliers. | Interpretability of the results can be challenging. |
| Retains non-linear patterns in high-dimensional data. | Sensitive to different parameter settings. |

# Comparison between PCA and t-SNE

## Principal Component Analysis (PCA)

A linear dimensionality reduction technique that seeks to maximize variance.

Great for handling large datasets and identifying underlying patterns.

May struggle with non-linear relationships in the data.

Often used for exploratory data analysis and feature extraction.

## t-SNE (t-Distributed Stochastic Neighbor Embedding)

A non-linear dimensionality reduction technique that focuses on local relationships.

Effective in capturing complex structures and preserving clusters.

Computationally expensive and sensitive to hyperparameters.

Commonly used for visualizing high-dimensional data in lower dimensions.

# Use cases of clustering algorithms in real-world scenarios

## Customer Segmentation

Businesses use clustering algorithms to categorize customers based on purchasing behavior.

## Image Clustering

Media organizations use clustering to organize and classify large image databases efficiently.
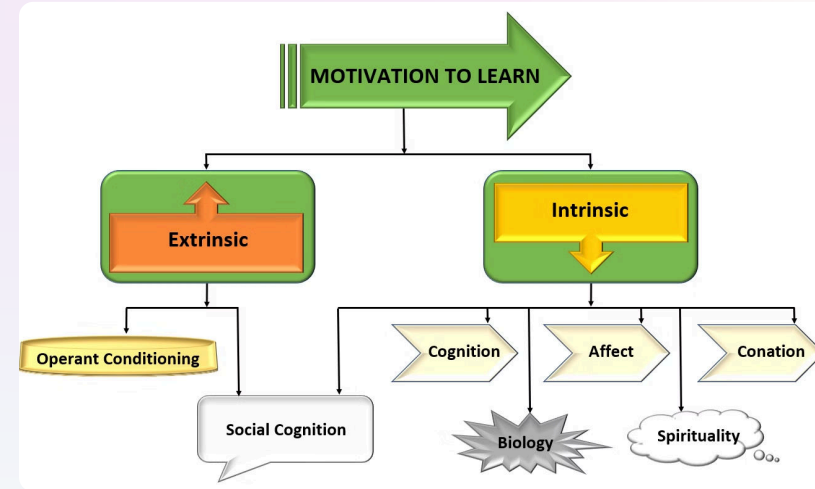
## Anomaly Detection

Clustering algorithms help identify unusual patterns in network data for security purposes.

# Challenges and Considerations in Unsupervised Learning

One of the challenges in unsupervised learning is the difficulty in evaluating the results.

Another consideration is the potential for overfitting due to the absence of labels.

Additionally, interpreting and explaining the outcomes of unsupervised learning models can be complex.

# Conclusion and key takeaways

After exploring clustering algorithms and dimensionality reduction in unsupervised learning, it's clear that they play a crucial role in organizing unstructured data and uncovering patterns. Understanding the advantages and limitations of these techniques is essential for their effective application in real-world scenarios.