

## UNSUPERVISED LEARNING

Unsupervised learning is a foundational approach in data science that focuses on discovering hidden patterns and structures in unlabeled datasets. This document explores the mathematical foundations, techniques, and practical applications of unsupervised learning methods, including clustering, dimensionality reduction, and anomaly detection. The methodologies discussed range from classical algorithms like K-Means and DBSCAN to advanced approaches such as Eigen decomposition and 2D Discrete Cosine Transform (DCT). The document emphasizes real-world applications, challenges, and strategies for optimizing unsupervised models. By integrating mathematical rigor with practical insights, this resource provides a comprehensive guide for implementing unsupervised learning techniques across diverse domains.

*This document is an extension of the research and lecture notes completed at Johns Hopkins University, Whiting School of Engineering, Engineering for Professionals, Artificial Intelligence Master's Program, Computer Science Master's Program and Data Science Master's Program.*

# Contents

<b>1</b>	<b>Introduction to Unsupervised Analysis</b>	<b>1</b>
1.1	Data Preprocessing . . . . .	1
1.2	Advanced Considerations in Unsupervised Analysis . . . . .	3
1.3	Summary . . . . .	3
<b>2</b>	<b>Types of Unsupervised Learning Problems</b>	<b>4</b>
2.1	Clustering . . . . .	4
2.2	Dimensionality Reduction . . . . .	4
2.3	Anomaly Detection . . . . .	5
2.4	Association Rule Learning . . . . .	5
2.5	Imbalanced Data in Unsupervised Learning . . . . .	6
2.6	Summary of Unsupervised Learning Problems . . . . .	6
<b>3</b>	<b>Mathematical Foundations of Unsupervised Learning</b>	<b>7</b>
3.1	Clustering Objectives and Formulations . . . . .	7
3.2	Dimensionality Reduction and Representation Learning . . . . .	8
3.3	Anomaly Detection and Density Estimation . . . . .	8
3.4	Interpretability in Unsupervised Models . . . . .	9
3.5	Applications of Mathematical Foundations . . . . .	9
3.6	Summary of Mathematical Foundations in Unsupervised Learning . . . . .	9
<b>4</b>	<b>K-Means Clustering for Text-Based Data</b>	<b>11</b>
4.1	Mathematical Formulations . . . . .	11
4.2	K-Means Algorithm . . . . .	11
4.3	Implementation Details . . . . .	12
4.4	Advantages and Limitations . . . . .	12
4.5	Analysis of the K-Means Algorithm . . . . .	13
4.6	Correctness Proof . . . . .	13
4.7	Summary . . . . .	14
<b>5</b>	<b>Eigen Decomposition for Variance Expansion and Factor Analysis</b>	<b>15</b>
5.1	Mathematical Formulations . . . . .	15
5.2	Eigen Decomposition Algorithm . . . . .	15
5.3	Implementation Details . . . . .	15
5.4	Advantages and Limitations . . . . .	16
5.5	Analysis of the Eigen Decomposition for Variance Expansion and Factor Analysis Algorithm . . . . .	17
5.6	Correctness Proof . . . . .	18
5.7	Summary Eigen Decomposition for Variance Expansion and Factor Analysis . . . .	18
<b>6</b>	<b>2D Discrete Cosine Transform (2D-DCT) and Inverse 2D-DCT for Unsuper- vised Learning</b>	<b>19</b>
6.1	Mathematical Formulations . . . . .	19
6.2	2D-DCT and 2D-IDCT Algorithms . . . . .	19
6.3	Implementation Details . . . . .	21

6.4	Advantages and Limitations . . . . .	21
6.5	Summary - 2D Discrete Cosine Transform (2D-DCT) and Inverse 2D-DCT for Un-supervised Learning . . . . .	24
<b>7</b>	<b>K-Means Algorithm for Anomaly Detection</b>	<b>25</b>
7.1	Mathematical Formulations . . . . .	25
7.2	K-Means Algorithm for Anomaly Detection Algorithm . . . . .	25
7.3	Implementation Details . . . . .	27
7.4	Advantages and Limitations . . . . .	27
7.5	Analysis of the K-Means Algorithm for Anomaly Detection . . . . .	28
7.6	Time Complexity Analysis . . . . .	28
7.7	Correctness Proof . . . . .	29
7.8	Summary - K-Means Algorithm for Anomaly Detection . . . . .	29
<b>8</b>	<b>DBSCAN (Density-Based Spatial Clustering of Applications with Noise) for Anomaly Detection</b>	<b>30</b>
8.1	Mathematical Formulations . . . . .	30
8.2	DBSCAN for Anomaly Detection Algorithm . . . . .	30
8.3	Implementation Details . . . . .	30
8.4	Advantages and Limitations . . . . .	32
8.5	Time Complexity Analysis . . . . .	32
8.6	Analysis of the DBSCAN Algorithm for Anomaly Detection . . . . .	32
8.7	Correctness Proof . . . . .	33
8.8	Summary - DBSCAN for Anomaly Detection . . . . .	33
<b>9</b>	<b>One-Class SVM for Anomaly Detection</b>	<b>34</b>
9.1	Mathematical Formulations . . . . .	34
9.2	One-Class SVM for Anomaly Detection Algorithm . . . . .	34
9.3	Implementation Details . . . . .	34
9.4	Advantages and Limitations . . . . .	35
9.5	Time Complexity Analysis . . . . .	36
9.6	Correctness Proof . . . . .	36
9.7	Summary - One-Class SVM for Anomaly Detection Algorithm . . . . .	36
<b>10</b>	<b>Evaluation Metrics for Unsupervised Learning</b>	<b>37</b>
10.1	Clustering Evaluation Metrics . . . . .	37
10.2	Dimensionality Reduction Metrics . . . . .	38
10.3	Anomaly Detection Metrics . . . . .	38
10.4	Time Complexity of Evaluation Metrics . . . . .	38
10.5	Summary of Evaluation Metrics for Unsupervised Learning . . . . .	39
<b>11</b>	<b>Recent Advances in Unsupervised Learning</b>	<b>40</b>
11.1	Clustering in Unsupervised Learning . . . . .	40
11.2	Dimensionality Reduction and Representation Learning . . . . .	41
11.3	Anomaly Detection in Unsupervised Learning . . . . .	43
11.4	Interpretability in Unsupervised Learning . . . . .	44

11.5 Applications and Future Directions . . . . .	44
11.6 Summary of Recent Advances in Unsupervised Learning . . . . .	46
<b>12 Unsupervised Learning Methods in scikit-learn</b>	<b>48</b>
12.1 Overview of Scikit-learn . . . . .	48
12.2 Common Steps for Unsupervised Learning with Scikit-learn . . . . .	48
12.3 Supported Unsupervised Learning Methods . . . . .	49
12.3.1 Clustering Algorithms . . . . .	49
12.3.2 Dimensionality Reduction Methods . . . . .	49
12.3.3 Anomaly Detection Methods . . . . .	50
12.3.4 Interpretability Techniques . . . . .	50
12.4 Evaluation of Unsupervised Models . . . . .	50
12.5 Applications of Scikit-learn in Unsupervised Learning . . . . .	51
12.6 Conclusion . . . . .	51
<b>13 Summary</b>	<b>52</b>
<b>14 Module Questions</b>	<b>53</b>
<b>15 Recommended Kaggle Datasets for Unsupervised Learning</b>	<b>54</b>
15.1 Clustering Methods . . . . .	54
15.2 Dimensionality Reduction . . . . .	54
15.3 Anomaly Detection . . . . .	54
15.4 Interpretability in Unsupervised Learning . . . . .	54
15.5 Deep Learning for Unsupervised Tasks . . . . .	55
<b>16 Group Exercise</b>	<b>56</b>

# 1 Introduction to Unsupervised Analysis

Definition of unsupervised learning and its focus on discovering patterns and structures in unlabeled data.

Real-world applications (e.g., customer segmentation, anomaly detection, dimensionality reduction).

Overview of unsupervised analysis workflow:

- Data preprocessing.
- Pattern discovery.
- Model evaluation and interpretation.

## Applications of Unsupervised Learning

Unsupervised learning methods are extensively used to uncover hidden structures and relationships within data, enabling exploratory data analysis and pattern recognition [19]. Techniques such as clustering, dimensionality reduction, and anomaly detection empower researchers to identify trends, group similar entities, and detect rare events without relying on labeled data. Additionally, unsupervised learning serves as a foundation for semi-supervised and transfer learning, expanding its utility in complex data analysis tasks.

## What Are the Benefits?

Unsupervised learning offers several advantages. Firstly, it enables **exploratory analysis** by discovering hidden patterns and clusters in data, making it particularly valuable for tasks like customer segmentation [23]. Secondly, it is **data-efficient** as it operates without labeled data, reducing the dependency on expensive or time-consuming annotations. Thirdly, unsupervised learning provides **dimensionality reduction**, allowing the representation of high-dimensional data in a compact and interpretable form. Lastly, it acts as a **preprocessing step** for supervised learning tasks by extracting meaningful features, improving downstream model performance.

## What Are the Pitfalls?

Despite its strengths, unsupervised learning has several challenges. One major issue is **lack of ground truth**, which makes evaluating model performance difficult [18]. Another concern is **scalability**; certain algorithms, such as hierarchical clustering, may struggle with large datasets due to their high computational complexity [36]. Additionally, unsupervised learning methods are **sensitive to hyperparameters**, such as the number of clusters in K-Means or the density threshold in DBSCAN, which can significantly affect outcomes. Finally, **interpretability** remains a challenge, as understanding and explaining the discovered patterns often require domain expertise and advanced techniques.

### 1.1 Data Preprocessing

Effective data preprocessing is critical for robust unsupervised analysis. Real-world datasets often contain imperfections such as missing values, outliers, and irrelevant features, which can adversely

impact the results [15]. Preprocessing transforms raw data into a structured format suitable for analysis, enhancing both accuracy and interpretability.

This section covers essential preprocessing techniques:

- **Feature Scaling and Normalization:** Ensures consistency in data ranges for algorithms sensitive to magnitudes.
- **Handling Missing Data:** Strategies to manage incomplete datasets effectively.
- **Dimensionality Reduction:** Reducing feature dimensions to improve computational efficiency and visualization.

By addressing these issues, unsupervised learning models can operate on clean, well-prepared data, leading to more reliable and meaningful outcomes [23].

## Feature Scaling and Normalization

Feature scaling ensures that all data dimensions contribute equally to the analysis by standardizing their ranges [15]. Common methods include **Standardization**, which transforms features to have a mean of zero and a standard deviation of one, and **Min-Max Scaling**, which normalizes values to a fixed range, typically  $[0, 1]$ . These techniques are particularly crucial for distance-based algorithms like K-Means and DBSCAN.

## Handling Missing Data

Missing data can distort patterns and reduce the effectiveness of unsupervised learning algorithms [26]. Basic imputation methods, such as replacing missing values with the mean or median, provide a simple solution. Advanced techniques, such as **Multiple Imputation by Chained Equations (MICE)**, generate multiple plausible datasets, ensuring robust inferences by accounting for uncertainty [7].

## Dimensionality Reduction

Dimensionality reduction simplifies datasets by projecting them into lower-dimensional spaces, retaining critical structures while discarding noise. Techniques like **Principal Component Analysis (PCA)** and **t-Distributed Stochastic Neighbor Embedding (t-SNE)** improve visualization and computational efficiency, particularly for high-dimensional datasets [24].

## Practical Implementation Considerations

Implementing unsupervised learning models effectively requires attention to various practical aspects, including algorithm choice, parameter tuning, and scalability [15]. Strategies like pipeline creation streamline the process by chaining preprocessing steps and model training into a unified framework, ensuring consistent results across multiple stages of analysis. Additionally, distributed computing platforms, such as Apache Spark, enable scalable solutions for large datasets by parallelizing computations [11].

## 1.2 Advanced Considerations in Unsupervised Analysis

Advanced considerations in unsupervised learning focus on model evaluation, visualization, and interpretability to ensure actionable insights. Techniques such as clustering validation metrics, dimensionality reduction for visualization, and feature attribution methods enhance understanding and usability. These considerations are essential for refining models and addressing real-world challenges in exploratory data analysis.

### Model Evaluation and Validation

Evaluating unsupervised models requires specialized techniques due to the absence of labeled data. Key approaches include:

- **Silhouette Score:** Measures cluster cohesion and separation.
- **Explained Variance:** Assesses the quality of dimensionality reduction.
- **Reconstruction Error:** Quantifies anomalies in reconstruction-based models [19].

### Interpretation and Explainability

Interpreting unsupervised models is critical for understanding discovered patterns and making informed decisions. Techniques like **SHAP** (SHapley Additive exPlanations) and visualization methods, such as t-SNE, enable domain experts to uncover meaningful insights [29, 35]. Cluster explanation using prototypes and decision rules further enhances interpretability by summarizing patterns in an intuitive format.

## 1.3 Summary

Unsupervised analysis is a cornerstone of data exploration in data science, enabling the discovery of hidden patterns, relationships, and anomalies within unlabeled data. Its advantages include exploratory capabilities, data efficiency, and dimensionality reduction, making it invaluable for applications like customer segmentation, anomaly detection, and visualization. Challenges such as scalability, sensitivity to hyperparameters, and interpretability highlight the need for robust pre-processing, thoughtful algorithm selection, and advanced tools for evaluation and explanation. By mastering these techniques, practitioners can leverage unsupervised learning to extract actionable insights and drive data-driven decision-making in diverse domains.

## 2 Types of Unsupervised Learning Problems

Unsupervised learning involves discovering patterns, structures, or relationships within unlabeled data. Depending on the nature of the task and data, unsupervised learning problems can be categorized into several types. This lecture explores these categories, highlighting their unique characteristics, challenges, and applications.

### 2.1 Clustering

Clustering aims to group data points into clusters based on similarity, where points within the same cluster are more similar to each other than to those in other clusters. Examples include:

- **Customer segmentation:** Grouping customers based on purchasing behavior.
- **Document clustering:** Organizing documents into topics or categories.

#### Key Techniques:

- **K-Means:** Partitions the data into  $k$  clusters by minimizing intra-cluster variance.
- **DBSCAN:** Identifies clusters based on density and treats low-density regions as noise.
- **Hierarchical Clustering:** Builds a tree of clusters (dendrogram) for nested clustering.

#### Evaluation Metrics:

- **Silhouette Score:** Measures how similar points are to their cluster compared to other clusters.
- **Calinski-Harabasz Index:** Evaluates cluster compactness and separation.

Clustering provides insights into the underlying structure of data, enabling exploratory analysis and decision-making.

### 2.2 Dimensionality Reduction

Dimensionality reduction transforms high-dimensional data into a lower-dimensional space while retaining meaningful structures and relationships. Examples include:

- **Image compression:** Reducing the dimensionality of image data for efficient storage.
- **Visualization:** Projecting high-dimensional data into 2D or 3D for interpretability.

#### Key Techniques:

- **Principal Component Analysis (PCA):** Projects data onto orthogonal axes that maximize variance.
- **t-SNE:** Preserves local similarity for effective visualization of complex data.
- **Autoencoders:** Neural network-based methods that learn compact data representations.

#### Evaluation Metrics:

- **Explained Variance:** Quantifies the proportion of data variance captured by reduced dimensions.



- **Reconstruction Error:** Measures how accurately the original data is reconstructed from the reduced representation.

Dimensionality reduction enhances computational efficiency and aids in data exploration and visualization.

## 2.3 Anomaly Detection

Anomaly detection identifies rare or unusual data points that deviate significantly from the majority. Examples include:

- **Fraud detection:** Identifying fraudulent transactions in financial data.
- **Network intrusion detection:** Detecting irregular patterns in network traffic.

### Key Techniques:

- **Density-Based Methods:** Identify anomalies as points in low-density regions (e.g., DBSCAN).
- **Boundary-Based Methods:** Use models like One-Class SVM to define decision boundaries separating normal and anomalous points.
- **Reconstruction-Based Methods:** Autoencoders detect anomalies by measuring reconstruction error.

### Evaluation Metrics:

- **Precision-Recall Curves:** Evaluate performance in highly imbalanced datasets.
- **Reconstruction Error:** High reconstruction error often indicates anomalies.

Anomaly detection is critical in domains requiring high sensitivity to rare but significant events.

## 2.4 Association Rule Learning

Association rule learning discovers relationships between variables in a dataset, often used for market basket analysis. Examples include:

- **Retail:** Identifying products frequently purchased together.
- **Web usage mining:** Discovering patterns in user behavior on websites.

### Key Techniques:

- **Apriori Algorithm:** Identifies frequent itemsets and generates association rules.
- **FP-Growth:** Efficiently finds frequent patterns without candidate generation.

### Evaluation Metrics:

- **Support:** Measures how often an itemset appears in the dataset.
- **Confidence:** Measures the likelihood of a consequent itemset given an antecedent itemset.
- **Lift:** Evaluates the strength of an association rule compared to random co-occurrence.

Association rule learning provides actionable insights for recommendation systems and strategic decision-making.

## 2.5 Imbalanced Data in Unsupervised Learning

Imbalanced data occurs when certain groups or patterns are underrepresented. Examples include:

- **Fraud detection:** Few fraudulent transactions compared to non-fraudulent ones.
- **Rare event detection:** Identifying failures in a manufacturing process.

### Challenges:

- Clustering algorithms may overfit to dominant patterns, ignoring minority groups.
- Dimensionality reduction and anomaly detection may fail to adequately capture rare patterns.

### Solutions:

- **Oversampling:** Generate synthetic data points to balance groups.
- **Algorithmic Adjustments:** Incorporate density-based techniques that are less sensitive to imbalance.

Addressing imbalanced data is essential for unbiased analysis and effective model performance in real-world settings.

## 2.6 Summary of Unsupervised Learning Problems

Unsupervised learning encompasses a diverse range of problems, including clustering, dimensionality reduction, anomaly detection, and association rule learning. These methods are instrumental in discovering patterns, summarizing data, and identifying rare events. By understanding the unique characteristics and challenges of each type, practitioners can design and implement robust unsupervised models tailored to specific data and application requirements.

### 3 Mathematical Foundations of Unsupervised Learning

Mathematics underpins unsupervised learning, providing the formal framework for clustering, dimensionality reduction, anomaly detection, and interpretability. By leveraging probability theory, linear algebra, optimization techniques, and graph theory, unsupervised learning methods uncover hidden structures and patterns in unlabeled data. This section explores the mathematical principles that drive key unsupervised learning techniques, including clustering objectives, low-dimensional representations, density estimation, and feature attributions.

#### 3.1 Clustering Objectives and Formulations

Clustering aims to partition data into groups where points in the same group are more similar to each other than to those in other groups. Common mathematical formulations include:

##### K-Means Clustering

K-Means minimizes the within-cluster sum of squares (WCSS):

$$\text{WCSS} = \sum_{k=1}^K \sum_{\mathbf{x}_i \in S_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2,$$

where:

- $K$  is the number of clusters.
- $S_k$  is the set of points in cluster  $k$ .
- $\boldsymbol{\mu}_k$  is the centroid of cluster  $k$ .

K-Means uses an iterative optimization approach, alternating between assigning points to the nearest centroid and updating centroids to the mean of assigned points.

##### DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

DBSCAN defines clusters as density-connected regions and labels low-density points as noise. A core point is defined by:

$$|\text{Neighborhood}_\epsilon(\mathbf{x}_i)| \geq \text{minPts},$$

where  $\text{Neighborhood}_\epsilon(\mathbf{x}_i)$  is the set of points within a radius  $\epsilon$  of  $\mathbf{x}_i$ . Clusters are formed by density-connected core points.

##### Spectral Clustering

Spectral clustering leverages the eigenvalues of the Laplacian matrix  $L = D - W$ , where  $D$  is the degree matrix and  $W$  is the adjacency matrix of the graph. The normalized Laplacian  $L_{\text{norm}}$  is used to partition data by solving:

$$L_{\text{norm}} \mathbf{u}_i = \lambda_i \mathbf{u}_i,$$

where  $\mathbf{u}_i$  are the eigenvectors, and  $\lambda_i$  are the eigenvalues.

### 3.2 Dimensionality Reduction and Representation Learning

Dimensionality reduction transforms high-dimensional data into a lower-dimensional space while preserving critical structures. Mathematical methods include:

#### Principal Component Analysis (PCA)

PCA identifies directions (principal components) that maximize variance. The principal components  $\mathbf{w}_j$  are the eigenvectors of the covariance matrix  $\Sigma$ , satisfying:

$$\Sigma \mathbf{w}_j = \lambda_j \mathbf{w}_j,$$

where  $\lambda_j$  are the eigenvalues.

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

t-SNE minimizes the Kullback-Leibler (KL) divergence between high-dimensional and low-dimensional probability distributions:

$$\text{KL}(P\|Q) = \sum_{i \neq j} P_{ij} \log \frac{P_{ij}}{Q_{ij}},$$

where  $P_{ij}$  is the pairwise similarity in the original space, and  $Q_{ij}$  is the similarity in the reduced space.

#### Autoencoders

Autoencoders are neural networks that encode input data into a latent space and reconstruct it:

$$\min_{\theta} \|\mathbf{x} - f_{\text{dec}}(f_{\text{enc}}(\mathbf{x}; \theta_{\text{enc}}); \theta_{\text{dec}})\|^2,$$

where  $f_{\text{enc}}$  and  $f_{\text{dec}}$  are the encoder and decoder functions, respectively, and  $\theta = \{\theta_{\text{enc}}, \theta_{\text{dec}}\}$  are the parameters.

### 3.3 Anomaly Detection and Density Estimation

Anomaly detection identifies instances that deviate from normal patterns using density-based, boundary-based, or reconstruction-based methods.

#### One-Class SVM

One-Class SVM constructs a hyperplane in a transformed feature space that separates normal data from the origin:

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu N} \sum_{i=1}^N \xi_i,$$

subject to  $(\mathbf{w} \cdot \phi(\mathbf{x}_i)) \geq 1 - \xi_i$ ,  $\xi_i \geq 0$ , where  $\nu$  controls the fraction of outliers.

## Isolation Forest

Isolation Forest isolates anomalies by recursively partitioning the data. Anomalies are points that require fewer splits to isolate:

$$\text{Path Length}(\mathbf{x}) \approx \log_2(n),$$

where  $n$  is the number of samples in a partition.

## 3.4 Interpretability in Unsupervised Models

Interpretability techniques ensure unsupervised models are accessible and trustworthy.

### Feature Contribution Analysis

SHAP (SHapley Additive exPlanations) attributes the contribution of each feature to clustering or anomaly detection outputs:

$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} [v(S \cup \{i\}) - v(S)],$$

where  $v(S)$  is the value of subset  $S$ , and  $N$  is the set of all features.

### Visualization with t-SNE and UMAP

t-SNE and UMAP provide interpretable low-dimensional visualizations, helping identify clusters and anomalies.

## 3.5 Applications of Mathematical Foundations

Mathematical foundations of unsupervised learning drive applications across domains:

- **Clustering:** Customer segmentation, social network analysis, and document categorization.
- **Dimensionality Reduction:** Image compression, gene expression analysis, and visualization of high-dimensional datasets.
- **Anomaly Detection:** Fraud detection, cybersecurity, and equipment failure prediction.
- **Interpretability:** Explaining clustering results in marketing analytics or anomaly detection in healthcare.

## 3.6 Summary of Mathematical Foundations in Unsupervised Learning

This section explored the mathematical principles foundational to unsupervised learning, covering key areas such as clustering objectives, dimensionality reduction, anomaly detection, and interpretability. Clustering techniques, including formulations for K-Means, DBSCAN, and spectral clustering, were highlighted for their ability to uncover patterns in data. Dimensionality reduction methods like PCA, t-SNE, autoencoders, and UMAP were discussed as tools for efficient representation learning and visualization. Anomaly detection approaches, including density estimation, boundary-based methods, and reconstruction error, demonstrated their utility in identifying rare and unusual data points. Additionally, interpretability techniques, such as feature contribution

analysis and visualization, were emphasized for enhancing the transparency and trustworthiness of unsupervised models. By mastering these mathematical foundations, practitioners are equipped to design, evaluate, and deploy unsupervised learning models effectively, addressing real-world challenges in data exploration, anomaly detection, and representation learning.

## 4 K-Means Clustering for Text-Based Data

In this section, we explore the K-Means clustering algorithm, a fundamental unsupervised learning technique for partitioning text-based data into  $k$  distinct clusters. K-Means operates by iteratively assigning data points to the nearest cluster centroids and then recalculating these centroids based on the current assignments. The goal is to minimize the within-cluster sum of squares (WCSS), ensuring that data points within a cluster are as similar as possible. This algorithm is particularly effective for text data when the text is represented as numerical vectors, such as TF-IDF or Word2Vec embeddings. K-Means has wide applications in document clustering, topic modeling, and text-based market segmentation. Despite its simplicity and efficiency, K-Means is sensitive to the initial choice of centroids and assumes that clusters are spherical. In the following subsections, we delve into the mathematical foundations, implementation details, and the strengths and limitations of K-Means for text-based data clustering [21, 33].

### 4.1 Mathematical Formulations

Given a dataset of  $N$  observations, represented as  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ , the objective of K-Means is to partition  $\mathbf{X}$  into  $k$  clusters, each represented by a centroid  $\mathbf{c}_j$  for  $j = 1, 2, \dots, k$ . The goal is to minimize the within-cluster sum of squares (WCSS), defined as:

$$\text{WCSS} = \sum_{i=1}^N \min_{j \in \{1, \dots, k\}} \|\mathbf{x}_i - \mathbf{c}_j\|^2$$

where  $\|\mathbf{x}_i - \mathbf{c}_j\|^2$  is the squared Euclidean distance between a data point  $\mathbf{x}_i$  and the centroid  $\mathbf{c}_j$ .

### 4.2 K-Means Algorithm

---

**Algorithm 1** K-Means Clustering for Text-Based Data

---

```
1: function KMEANS( $\mathbf{X}, k, \text{max\_iterations}$ )
2:   Initialize  $k$  cluster centroids  $\mathbf{C} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k]$  randomly from  $\mathbf{X}$ 
3:   Initialize assignments  $\mathbf{A} = [a_1, a_2, \dots, a_N]$ , where  $a_i$  is the cluster assignment for  $\mathbf{x}_i$ 
4:   for iteration = 1 to max_iterations do
5:                                      $\triangleright$  Assignment Step
6:     for each document vector  $\mathbf{x}_i \in \mathbf{X}$  do
7:        $a_i \leftarrow \arg \min_{j \in \{1, \dots, k\}} \text{Distance}(\mathbf{x}_i, \mathbf{c}_j)$ 
8:     end for
9:                                      $\triangleright$  Update Step
10:    for  $j = 1$  to  $k$  do
11:       $\mathbf{c}_j \leftarrow \frac{1}{|S_j|} \sum_{\mathbf{x}_i \in S_j} \mathbf{x}_i$ , where  $S_j = \{\mathbf{x}_i \mid a_i = j\}$ 
12:    end for
13:                                      $\triangleright$  Check for Convergence
14:    if centroids  $\mathbf{C}$  have not changed significantly then
15:      break
16:    end if
17:  end for
18:  return  $\mathbf{A}, \mathbf{C}$ 
19: end function
```

---

The K-Means algorithm can be described as follows:

1. Initialize  $k$  centroids randomly from the dataset  $\mathbf{X}$ .
2. Assign each data point to the nearest centroid.
3. Update each centroid by computing the mean of the data points assigned to it.
4. Repeat steps 2 and 3 until convergence (i.e., centroids do not change significantly).

### 4.3 Implementation Details

#### Data Preprocessing:

- Convert text data to numerical vectors using techniques like Term Frequency-Inverse Document Frequency (TF-IDF) or Word2Vec.
- Normalize the feature vectors to ensure that each feature contributes equally to the distance calculations.

#### Model Processing:

- Initialize  $k$  centroids randomly from the data.
- Use a distance metric such as the Euclidean distance or cosine similarity for clustering.

#### Evaluation:

- Use the Elbow Method to determine the optimal number of clusters  $k$ .
- Calculate the Silhouette Score to measure the quality of the clusters.

### 4.4 Advantages and Limitations

#### Advantages:

- Simplicity: Easy to understand and implement.
- Scalability: Efficient for large datasets.
- Versatility: Applicable to a wide range of domains, including text data, image processing, and market segmentation [32].

#### Limitations:

- Sensitivity to Initialization: Results may vary depending on the initial choice of centroids.
- Fixed Number of Clusters: Requires specifying  $k$  beforehand.
- Assumes Spherical Clusters: Performs poorly when clusters have irregular shapes or different densities.



## 4.5 Analysis of the K-Means Algorithm

The time complexity of the K-Means algorithm can be broken down as follows:

- Initialization: Randomly initializing  $k$  centroids takes  $O(k)$ .
- Assignment Step: For each of the  $N$  data points, the distance to each of the  $k$  centroids needs to be calculated. Assuming the data points have  $d$  dimensions, this step takes  $O(N \times k \times d)$  per iteration.
- Update Step: For each of the  $k$  clusters, the centroid is recalculated by averaging the points assigned to that cluster. In the worst case, all  $N$  points could belong to a single cluster, leading to a complexity of  $O(N \times d)$  for updating all centroids.

Given that the algorithm runs for a maximum of  $I$  iterations, the total time complexity is:

$$O(I \times N \times k \times d)$$

In practice,  $I$  (the number of iterations) is relatively small, making K-Means computationally efficient for moderate-sized datasets [28].

K-Means is a popular clustering algorithm due to its simplicity and efficiency. It works well when:

- The data clusters are *spherical* and *well-separated*.
- The number of clusters  $k$  is known or can be estimated using techniques like the *Elbow Method* or the *Silhouette Score*.

However, K-Means has some limitations:

- Sensitivity to Initialization: Poor initialization of centroids can lead to suboptimal clustering.
- Fixed Number of Clusters: Requires specifying  $k$  in advance.
- Non-Convex Clusters: Performs poorly on datasets with non-spherical or overlapping clusters.

### Summary - K-Means for Unsupervised Learning

K-Means is a cornerstone algorithm for clustering tasks in unsupervised learning. It is easy to implement, computationally efficient, and widely used for tasks such as market segmentation, image compression, and document clustering. Despite its limitations, it remains a go-to method when the dataset's structure aligns with its assumptions [22].

## 4.6 Correctness Proof

### Correctness of the K-Means Algorithm

The K-Means algorithm is guaranteed to converge to a local minimum of the objective function, which is the within-cluster sum of squares (WCSS). This is because each iteration of the algorithm strictly decreases the objective function until no further improvement is possible.

## Theorem

The K-Means algorithm converges to a local minimum of the within-cluster sum of squares (WCSS) in a finite number of iterations.

## Proof

The K-Means algorithm consists of two key steps: assignment and update.

1. Assignment Step: Each data point is assigned to the nearest centroid. This step decreases the WCSS since each point is moved to the closest cluster, minimizing the distance contribution to the WCSS.
2. Update Step: Centroids are updated by recalculating the mean of the points assigned to each cluster. The mean minimizes the sum of squared distances within the cluster.

Since there is a finite number of ways to partition  $N$  data points into  $k$  clusters, and the WCSS decreases or remains the same with each iteration, the algorithm must converge in a finite number of iterations.

Therefore, the K-Means algorithm converges to a local minimum of the WCSS [28].

## Summary for Correctness Proof

The K-Means algorithm is guaranteed to converge due to the monotonically decreasing nature of the objective function (WCSS) and the finite number of possible cluster assignments. However, it converges to a *local minimum*, not necessarily the global minimum. This highlights the importance of good initialization strategies, such as K-Means++ [6].

## 4.7 Summary

K-Means clustering is a widely used unsupervised learning algorithm for partitioning text-based data into  $k$  distinct clusters by minimizing the within-cluster sum of squares (WCSS). It works iteratively by assigning each data point to the nearest cluster centroid and then updating the centroids based on the current assignments. The algorithm is particularly effective when text data is represented as numerical vectors like TF-IDF or Word2Vec embeddings. K-Means is computationally efficient and finds applications in document clustering, topic modeling, and market segmentation. However, it is sensitive to the initial choice of centroids, requires the number of clusters  $k$  to be predefined, and assumes spherical cluster shapes. Despite these limitations, K-Means remains a fundamental tool in unsupervised learning due to its simplicity and effectiveness [21, 33].

## 5 Eigen Decomposition for Variance Expansion and Factor Analysis

Eigen decomposition is a fundamental technique in linear algebra that plays a crucial role in unsupervised data analysis, particularly in dimensionality reduction and exploratory factor analysis. By decomposing the covariance matrix of a dataset, eigen decomposition identifies principal components—directions in which the data’s variance is maximized. These principal components are obtained as eigenvectors corresponding to the largest eigenvalues of the covariance matrix. The algorithm consists of two key steps: variance expansion, where data is projected onto these principal components, and factor analysis, where the components are rotated (e.g., using Varimax rotation) to enhance interpretability. This approach is widely used in Principal Component Analysis (PCA) and factor analysis to reveal the latent structure within high-dimensional datasets, facilitating tasks such as visualization, noise reduction, and feature extraction [24, 8].

### 5.1 Mathematical Formulations

Given a dataset  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ , the covariance matrix  $\mathbf{\Sigma}$  is defined as:

$$\mathbf{\Sigma} = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$$

The eigen decomposition of the covariance matrix  $\mathbf{\Sigma}$  is:

$$\mathbf{\Sigma} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$$

where: -  $\mathbf{V}$  is the matrix of eigenvectors. -  $\mathbf{\Lambda}$  is the diagonal matrix of eigenvalues.

The data  $\mathbf{X}$  can be projected onto the principal components to expand variance:

$$\mathbf{Z} = \mathbf{X} \mathbf{V}$$

For Factor Analysis, a rotation matrix  $\mathbf{R}$  (e.g., Varimax) is applied to improve interpretability:

$$\mathbf{Z}_{\text{rot}} = \mathbf{Z} \mathbf{R}, \quad \mathbf{V}_{\text{rot}} = \mathbf{V} \mathbf{R}$$

### 5.2 Eigen Decomposition Algorithm

### 5.3 Implementation Details

**Data Preprocessing:**

- Centering the Data: Subtract the mean of each feature to ensure the data has zero mean.
- Standardization: Scale the features to have unit variance if the scales of the features differ significantly.

---

**Algorithm 2** Eigen Decomposition for Variance Expansion and Factor Analysis

---

```
1: function EIGENDECOMPOSITION( $\mathbf{X}$ )
2:                                     ▷ Step 1: Compute the covariance matrix
3:    $\mathbf{\Sigma} \leftarrow \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$ 
4:                                     ▷ Step 2: Perform eigen decomposition
5:   Compute eigenvectors  $\mathbf{V}$  and eigenvalues  $\mathbf{\Lambda}$  of  $\mathbf{\Sigma}$  such that  $\mathbf{\Sigma} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$ 
6:                                     ▷ Step 3: Map the data to expand variance
7:    $\mathbf{Z} \leftarrow \mathbf{X} \mathbf{V}$ 
8:   return  $\mathbf{Z}, \mathbf{V}, \mathbf{\Lambda}$ 
9: end function
10: function FACTORANALYSIS( $\mathbf{Z}, \mathbf{V}$ )
11:                                     ▷ Step 4: Perform rotation for factor analysis
12:   Choose a rotation matrix  $\mathbf{R}$  (e.g., Varimax rotation)
13:    $\mathbf{Z}_{\text{rot}} \leftarrow \mathbf{Z} \mathbf{R}$ 
14:    $\mathbf{V}_{\text{rot}} \leftarrow \mathbf{V} \mathbf{R}$ 
15:   return  $\mathbf{Z}_{\text{rot}}, \mathbf{V}_{\text{rot}}$ 
16: end function
```

---

**Model Processing:**

- Compute Covariance Matrix: Calculate the covariance matrix  $\mathbf{\Sigma}$  of the centered data.
- Eigen Decomposition: Decompose  $\mathbf{\Sigma}$  into eigenvectors and eigenvalues.
- Project Data: Map the data onto the eigenvectors to obtain the principal components.

**Evaluation:**

- Variance Explained: Compute the proportion of variance explained by each principal component.
- Rotation: Apply rotation (e.g., Varimax) to make the factors more interpretable.

## 5.4 Advantages and Limitations

**Advantages:**

- Dimensionality Reduction: Reduces the number of features while retaining most of the variance [24].
- Interpretability: Factor rotation improves the interpretability of principal components.
- Efficiency: Eigen decomposition is computationally efficient for moderate-sized datasets.

**Limitations:**

- Assumes Linearity: Assumes that the principal components capture linear relationships in the data.
- Sensitive to Noise: Eigen decomposition can be affected by noise and outliers [8].
- Computational Complexity: Can be computationally expensive for very large datasets.

## 5.5 Analysis of the Eigen Decomposition for Variance Expansion and Factor Analysis Algorithm

The Eigen Decomposition for Variance Expansion and Factor Analysis\* algorithm is a two-step process for analyzing data in an unsupervised manner. It leverages linear algebra techniques to transform and interpret high-dimensional data. The algorithm consists of:

1. Variance Expansion: By performing eigen decomposition on the covariance matrix of the data, the data is projected onto the principal components (eigenvectors) that maximize variance.
2. Factor Analysis: The principal components are rotated (e.g., using Varimax rotation) to improve interpretability by aligning factors with the data's inherent structure [24].

This algorithm is particularly useful in dimensionality reduction, data visualization, and identifying latent factors in datasets.

### Time Complexity Analysis

The time complexity of the algorithm can be analyzed in two parts:

#### 1. Computing the Covariance Matrix:

- Given a dataset  $\mathbf{X}$  with  $N$  data points and  $d$  dimensions, calculating the covariance matrix  $\mathbf{\Sigma} = \frac{1}{N-1}\mathbf{X}^T\mathbf{X}$  takes  $O(N \times d^2)$ .

#### 2. Eigen Decomposition:

- Performing eigen decomposition on a  $d \times d$  covariance matrix  $\mathbf{\Sigma}$  has a time complexity of  $O(d^3)$  [16].

#### 3. Mapping Data to Eigenvectors:

- Projecting the data onto the eigenvectors  $\mathbf{Z} = \mathbf{XV}$  takes  $O(N \times d^2)$ .

#### 4. Factor Analysis (Rotation):

- Rotating the principal components using a rotation matrix  $\mathbf{R}$  (e.g., Varimax) takes  $O(d^2)$ .

Total Time Complexity: The overall time complexity is:

$$O(N \times d^2 + d^3)$$

For large datasets where  $N \gg d$ , the complexity is dominated by  $O(N \times d^2)$ .

### Summary - Eigen Decomposition for Variance Expansion and Factor Analysis

Eigen decomposition is a powerful technique for unsupervised learning tasks, enabling the transformation of data into a set of principal components that capture the maximum variance. The algorithm involves computing the covariance matrix, performing eigen decomposition, and optionally rotating the principal components for better interpretability. This approach is widely used in

dimensionality reduction, data visualization, and exploratory factor analysis. Despite its computational cost, eigen decomposition remains a cornerstone of unsupervised data analysis due to its ability to reveal the underlying structure of complex datasets [8].

## 5.6 Correctness Proof

### Correctness of Eigen Decomposition

#### Theorem

The eigen decomposition of a symmetric covariance matrix  $\Sigma$  always exists, and the resulting eigenvectors form an orthonormal basis.

#### Proof

The covariance matrix  $\Sigma$  is symmetric by definition ( $\Sigma = \Sigma^T$ ). According to the **Spectral Theorem**, any real symmetric matrix can be diagonalized by an orthogonal matrix. Therefore, there exists a matrix of eigenvectors  $\mathbf{V}$  and a diagonal matrix of eigenvalues  $\Lambda$  such that:

$$\Sigma = \mathbf{V}\Lambda\mathbf{V}^T$$

The eigenvectors in  $\mathbf{V}$  are orthonormal, meaning  $\mathbf{V}^T\mathbf{V} = \mathbf{I}$ , where  $\mathbf{I}$  is the identity matrix. This ensures that the transformation  $\mathbf{Z} = \mathbf{XV}$  preserves the variance structure of the data [24].

Thus, the eigen decomposition of the covariance matrix  $\Sigma$  always exists and produces valid principal components.

### Summary for Correctness Proof

The correctness of eigen decomposition for variance expansion is guaranteed by the Spectral Theorem for symmetric matrices. The decomposition yields orthonormal eigenvectors and corresponding eigenvalues that capture the variance structure of the data. This ensures that the transformation into principal components is both mathematically sound and interpretable. Factor analysis via rotation further enhances the interpretability of these components by aligning them with the underlying data structure [16].

## 5.7 Summary Eigen Decomposition for Variance Expansion and Factor Analysis

Eigen decomposition is a powerful linear algebra technique used in unsupervised learning to analyze the structure of datasets by transforming them into principal components that capture the maximum variance. This process involves decomposing the covariance matrix into eigenvectors and eigenvalues, projecting the data onto these eigenvectors, and optionally applying rotations such as Varimax to improve interpretability. The algorithm is essential for tasks like Principal Component Analysis (PCA) and Factor Analysis, facilitating dimensionality reduction, feature extraction, and data visualization. Despite its computational complexity, eigen decomposition provides valuable insights into the latent structure of high-dimensional data, making it a cornerstone technique in exploratory data analysis [24, 8].

## 6 2D Discrete Cosine Transform (2D-DCT) and Inverse 2D-DCT for Unsupervised Learning

The 2D Discrete Cosine Transform (2D-DCT) is a fundamental tool in signal and image processing, widely used for transforming spatial-domain data into the frequency domain. It decomposes 2D datasets, such as images, into frequency components that capture patterns in different orientations—vertical, horizontal, and diagonal. This ability to analyze data in terms of its frequency content makes 2D-DCT particularly valuable for unsupervised learning tasks like feature extraction, clustering, and anomaly detection. The Inverse 2D Discrete Cosine Transform (2D-IDCT) complements this process by reconstructing the original spatial data from its frequency-domain representation. Together, these transformations enable efficient data compression, noise reduction, and pattern recognition, while preserving the integrity of the original information [3, 17, 40].

### 6.1 Mathematical Formulations

Given a 2D dataset  $\mathbf{X}$  of size  $M \times N$ , the 2D Discrete Cosine Transform of  $\mathbf{X}$  produces a matrix  $\mathbf{C}$  of DCT coefficients, defined as:

$$C_{u,v} = \alpha_u \alpha_v \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} \mathbf{X}_{x,y} \cos\left(\frac{\pi(2x+1)u}{2M}\right) \cos\left(\frac{\pi(2y+1)v}{2N}\right)$$

where:

$$\alpha_u = \begin{cases} \frac{1}{\sqrt{M}} & \text{if } u = 0 \\ \sqrt{\frac{2}{M}} & \text{if } u \neq 0 \end{cases}, \quad \alpha_v = \begin{cases} \frac{1}{\sqrt{N}} & \text{if } v = 0 \\ \sqrt{\frac{2}{N}} & \text{if } v \neq 0 \end{cases}$$

The Inverse 2D Discrete Cosine Transform reconstructs  $\mathbf{X}$  from  $\mathbf{C}$  as:

$$\mathbf{X}_{x,y} = \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} \alpha_u \alpha_v C_{u,v} \cos\left(\frac{\pi(2x+1)u}{2M}\right) \cos\left(\frac{\pi(2y+1)v}{2N}\right)$$

### 6.2 2D-DCT and 2D-IDCT Algorithms

#### Discrete Cosine Transform (2D-DCT) Algorithm:

The 2D-DCT computes the frequency coefficients of a 2D dataset by applying the DCT formula to each pair of frequency indices  $(u, v)$ . The resulting coefficients capture the data's frequency patterns in different directions:

- Vertical Coefficients: Capture low-frequency vertical patterns.
- Horizontal Coefficients: Capture low-frequency horizontal patterns.
- Diagonal Coefficients: Capture high-frequency diagonal patterns.

---

**Algorithm 3** 2D Discrete Cosine Transform (2D-DCT)

---

```
1: function DCT2D(X)  
2:                                     ▷ Step 1: Initialize dimensions  
3:    $M, N \leftarrow$  number of rows and columns of X  
4:   Initialize C as a zero matrix of size  $M \times N$   
5:                                     ▷ Step 2: Compute 2D-DCT coefficients  
6:   for  $u = 0$  to  $M - 1$  do  
7:     for  $v = 0$  to  $N - 1$  do  
8:        $C_{u,v} \leftarrow \alpha_u \alpha_v \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} \mathbf{X}_{x,y} \cos\left(\frac{\pi(2x+1)u}{2M}\right) \cos\left(\frac{\pi(2y+1)v}{2N}\right)$   
9:     end for  
10:  end for  
11:  return C =  $[C_{u,v}]$   
12: end function  
13:                                     ▷ Scaling factors  
14:  $\alpha_u = \begin{cases} \frac{1}{\sqrt{M}} & \text{if } u = 0 \\ \sqrt{\frac{2}{M}} & \text{if } u \neq 0 \end{cases}$   
15:  $\alpha_v = \begin{cases} \frac{1}{\sqrt{N}} & \text{if } v = 0 \\ \sqrt{\frac{2}{N}} & \text{if } v \neq 0 \end{cases}$ 
```

---

---

**Algorithm 4** Inverse 2D Discrete Cosine Transform (2D-IDCT)

---

```
1: function IDCT2D(C)  
2:                                     ▷ Step 1: Initialize dimensions  
3:    $M, N \leftarrow$  number of rows and columns of C  
4:   Initialize X as a zero matrix of size  $M \times N$   
5:                                     ▷ Step 2: Compute 2D-IDCT  
6:   for  $x = 0$  to  $M - 1$  do  
7:     for  $y = 0$  to  $N - 1$  do  
8:        $\mathbf{X}_{x,y} \leftarrow \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} \alpha_u \alpha_v C_{u,v} \cos\left(\frac{\pi(2x+1)u}{2M}\right) \cos\left(\frac{\pi(2y+1)v}{2N}\right)$   
9:     end for  
10:  end for  
11:  return X =  $[\mathbf{X}_{x,y}]$   
12: end function  
13:                                     ▷ Scaling factors  
14:  $\alpha_u = \begin{cases} \frac{1}{\sqrt{M}} & \text{if } u = 0 \\ \sqrt{\frac{2}{M}} & \text{if } u \neq 0 \end{cases}$   
15:  $\alpha_v = \begin{cases} \frac{1}{\sqrt{N}} & \text{if } v = 0 \\ \sqrt{\frac{2}{N}} & \text{if } v \neq 0 \end{cases}$ 
```

---



## Inverse Discrete Cosine Transform (2D-IDCT) Algorithm:

The 2D-IDCT reconstructs the original dataset from its DCT coefficients. By summing the contributions of all frequency components, the 2D-IDCT ensures that the transformation is reversible, accurately restoring the original data.

### 6.3 Implementation Details

#### Data Preprocessing:

- Normalization: Normalize the data to a standard range (e.g.,  $[0, 1]$ ) to stabilize the transformation.
- Centering: Subtract the mean to center the data, ensuring it has zero mean.

#### Model Processing:

- Compute 2D-DCT: Apply the DCT formula to transform the data into frequency coefficients.
- Classify Coefficients:
  - ◊ Vertical: Low-frequency coefficients along the vertical direction ( $u$ ).
  - ◊ Horizontal: Low-frequency coefficients along the horizontal direction ( $v$ ).
  - ◊ Diagonal: High-frequency coefficients representing fine details.
- Compute 2D-IDCT: Reconstruct the original data from the DCT coefficients.

#### Evaluation:

- Reconstruction Error: Measure the difference between the original and reconstructed data.
- Energy Distribution: Analyze the distribution of energy across the coefficients to identify dominant patterns.

### 6.4 Advantages and Limitations

#### Advantages:

- Efficient Representation: Captures most of the data's energy in a few low-frequency coefficients [3].
- Feature Extraction: Useful for identifying patterns in image and signal data.
- Compression: Reduces data dimensionality, making it ideal for image compression (e.g., JPEG).
- Reversible: The 2D-IDCT ensures lossless reconstruction.

#### Limitations:

- Computational Complexity: Standard implementation has  $O(M^2N^2)$  time complexity; optimized versions (e.g., Fast DCT) reduce this to  $O(MN \log(MN))$ .

- Sensitive to Noise: High-frequency coefficients can be affected by noise [40].
- Assumes Stationarity: Assumes that the data’s frequency characteristics are stationary.

### **Analysis of the 2D Discrete Cosine Transform (2D-DCT) and Inverse 2D-DCT**

The 2D Discrete Cosine Transform (2D-DCT) is a fundamental transformation used in image and signal processing to convert spatial-domain data into the frequency domain. The 2D-DCT decomposes a matrix into frequency components, capturing vertical, horizontal, and diagonal patterns. This decomposition is crucial in unsupervised learning tasks, such as feature extraction, clustering, and anomaly detection [3, 17].

The Inverse 2D Discrete Cosine Transform (2D-IDCT) reconstructs the original matrix from its DCT coefficients. Together, the 2D-DCT and 2D-IDCT ensure reversible transformations, enabling compression and reconstruction applications, such as in JPEG image compression [40].

### **Time Complexity Analysis**

#### **2D Discrete Cosine Transform (2D-DCT):**

The 2D-DCT involves calculating  $M \times N$  coefficients, where each coefficient requires a double summation over  $M$  rows and  $N$  columns. Therefore, the time complexity of the 2D-DCT is:

$$O(M^2N^2)$$

#### **Inverse 2D Discrete Cosine Transform (2D-IDCT):**

Similarly, the 2D-IDCT reconstructs each of the  $M \times N$  data points by summing over  $M$  rows and  $N$  columns. Hence, the time complexity of the 2D-IDCT is also:

$$O(M^2N^2)$$

### **Optimized Versions:**

Fast algorithms, such as the Fast Discrete Cosine Transform (FDCT), reduce the complexity to  $O(MN \log(MN))$ , similar to the Fast Fourier Transform (FFT) [40].

### **Analysis of the 2D-DCT and 2D-IDCT for Identifying Coefficients**

The 2D-DCT transforms data into a frequency representation, allowing the identification of specific types of patterns:

- Vertical Coefficients: Represent low-frequency components in the vertical direction, capturing smooth variations and vertical edges.
- Horizontal Coefficients: Represent low-frequency components in the horizontal direction, capturing smooth variations and horizontal edges.
- Diagonal Coefficients: Represent high-frequency components, capturing fine details, textures, and diagonal edges.

By analyzing these coefficients, unsupervised learning models can detect patterns, anomalies, and features for tasks such as clustering, segmentation, and compression.

### Summary - 2D-DCT and 2D-IDCT for Identifying Coefficients

The 2D Discrete Cosine Transform (2D-DCT) decomposes a 2D matrix into frequency components, making it useful for identifying vertical, horizontal, and diagonal patterns. The Inverse 2D Discrete Cosine Transform (2D-IDCT) reconstructs the original data from these frequency coefficients. These algorithms are essential for tasks such as image compression, feature extraction, and unsupervised analysis. While computationally expensive in their standard form, optimized versions like the Fast DCT provide efficiency gains [3, 40].

### Correctness Proof

#### Correctness of the 2D-DCT and 2D-IDCT

##### Theorem

The Inverse 2D Discrete Cosine Transform (2D-IDCT) accurately reconstructs the original data from the 2D-DCT coefficients.

##### Proof

The 2D-DCT of a matrix  $\mathbf{X}$  of size  $M \times N$  produces a matrix of coefficients  $\mathbf{C}$ , where:

$$C_{u,v} = \alpha_u \alpha_v \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} \mathbf{X}_{x,y} \cos\left(\frac{\pi(2x+1)u}{2M}\right) \cos\left(\frac{\pi(2y+1)v}{2N}\right)$$

The 2D-IDCT reconstructs  $\mathbf{X}$  as:

$$\mathbf{X}_{x,y} = \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} \alpha_u \alpha_v C_{u,v} \cos\left(\frac{\pi(2x+1)u}{2M}\right) \cos\left(\frac{\pi(2y+1)v}{2N}\right)$$

By substituting  $C_{u,v}$  into the IDCT expression and applying the orthogonality of the cosine functions, it can be shown that the summation reduces to  $\mathbf{X}_{x,y}$ . The scaling factors  $\alpha_u$  and  $\alpha_v$  ensure normalization, preserving the original values [40].

Thus, the 2D-IDCT perfectly reconstructs the original matrix from its 2D-DCT coefficients.

### Summary for Correctness Proof

The correctness of the 2D-DCT and 2D-IDCT is guaranteed by the orthogonality of the cosine basis functions. This ensures that the transformation to the frequency domain (DCT) and the reconstruction to the spatial domain (IDCT) are lossless, preserving all information. The scaling factors normalize the coefficients, ensuring accurate reconstruction [3, 40].

## 6.5 Summary - 2D Discrete Cosine Transform (2D-DCT) and Inverse 2D-DCT for Unsupervised Learning

The 2D Discrete Cosine Transform (2D-DCT) and its inverse (2D-IDCT) are essential algorithms for analyzing 2D data by decomposing it into frequency components. The 2D-DCT identifies vertical, horizontal, and diagonal patterns by transforming data from the spatial domain to the frequency domain. The 2D-IDCT accurately reconstructs the original data, ensuring a reversible transformation. These algorithms play a crucial role in unsupervised learning tasks such as feature extraction, image compression, and anomaly detection. Despite the computational complexity of the standard implementations, optimized versions like the Fast DCT improve efficiency. The robustness and versatility of 2D-DCT and 2D-IDCT make them indispensable tools in data science and signal processing [3, 40].

## 7 K-Means Algorithm for Anomaly Detection

The K-Means Algorithm for Anomaly Detection is an effective unsupervised learning technique for identifying unusual patterns or outliers in a dataset. By clustering data points into a predefined number of groups,  $k$ , the algorithm assigns each point to the nearest cluster centroid based on a distance metric, such as the Euclidean distance. Anomalies are detected as data points that lie significantly far from their respective centroids, exceeding a specified threshold distance. This approach is particularly useful in various applications, including fraud detection, network security, and quality control, where detecting deviations from normal behavior is crucial. K-Means provides a simple yet powerful method for anomaly detection by leveraging the natural groupings within the data [32, 22].

### 7.1 Mathematical Formulations

Let the dataset be  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ , where each  $\mathbf{x}_i \in \mathbb{R}^d$  is a  $d$ -dimensional vector.

1. Cluster Centroids:

The  $k$  centroids  $\mathbf{C} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k]$  are updated as:

$$\mathbf{c}_j = \frac{1}{|S_j|} \sum_{\mathbf{x}_i \in S_j} \mathbf{x}_i,$$

where  $S_j$  is the set of points assigned to cluster  $j$ .

2. Distance Calculation:

The distance between a data point  $\mathbf{x}_i$  and its assigned centroid  $\mathbf{c}_{a_i}$  (using the Euclidean distance) is given by:

$$d_i = \|\mathbf{x}_i - \mathbf{c}_{a_i}\|_2 = \sqrt{\sum_{l=1}^d (x_{i,l} - c_{a_i,l})^2}.$$

3. Anomaly Detection:

A point  $\mathbf{x}_i$  is identified as an anomaly if:

$$d_i > \text{threshold}.$$

### 7.2 K-Means Algorithm for Anomaly Detection Algorithm

The K-Means algorithm clusters the data into  $k$  clusters and flags anomalies based on the distance of each data point from its assigned cluster centroid. The steps are as follows:

1. Initialization: Randomly initialize  $k$  centroids.
2. Assignment: Assign each data point to the nearest centroid.
3. Update: Recalculate the centroids as the mean of the assigned points.
4. Convergence Check: Stop if centroids do not change significantly.

---

**Algorithm 5** K-Means Algorithm for Anomaly Detection

---

```
1: function KMEANSANOMALYDETECTION( $\mathbf{X}, k, \text{threshold}, \text{max\_iterations}$ )
2:                                      $\triangleright$  Step 1: Initialize parameters
3:    $N \leftarrow$  number of data points in  $\mathbf{X}$ 
4:   Initialize  $k$  centroids  $\mathbf{C} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k]$  randomly from  $\mathbf{X}$ 
5:   Initialize assignments  $\mathbf{A} = [a_1, a_2, \dots, a_N]$  where  $a_i$  is the cluster for  $\mathbf{x}_i$ 
6:                                      $\triangleright$  Step 2: Iterate until convergence or max iterations
7:   for iteration = 1 to max_iterations do
8:                                      $\triangleright$  Assignment Step
9:     for each data point  $\mathbf{x}_i \in \mathbf{X}$  do
10:       $a_i \leftarrow \arg \min_{j \in \{1, \dots, k\}} \text{Distance}(\mathbf{x}_i, \mathbf{c}_j)$ 
11:    end for
12:                                      $\triangleright$  Update Step
13:    for  $j = 1$  to  $k$  do
14:       $\mathbf{c}_j \leftarrow \frac{1}{|S_j|} \sum_{\mathbf{x}_i \in S_j} \mathbf{x}_i$ , where  $S_j = \{\mathbf{x}_i \mid a_i = j\}$ 
15:    end for
16:                                      $\triangleright$  Check for Convergence
17:    if centroids  $\mathbf{C}$  have not changed significantly then
18:      break
19:    end if
20:  end for
21:                                      $\triangleright$  Step 3: Calculate distances to centroids
22:  for each data point  $\mathbf{x}_i \in \mathbf{X}$  do
23:     $d_i \leftarrow \text{Distance}(\mathbf{x}_i, \mathbf{c}_{a_i})$ 
24:  end for
25:                                      $\triangleright$  Step 4: Identify anomalies
26:  Anomalies  $\leftarrow \{\mathbf{x}_i \mid d_i > \text{threshold}\}$ 
27:  return  $\mathbf{A}, \mathbf{C}, \text{Anomalies}$ 
28: end function
```

---

5. Anomaly Detection: Calculate the distance of each point to its centroid and flag points with distances exceeding a threshold.

### 7.3 Implementation Details

#### Data Preprocessing:

- Normalization: Scale features to a standard range (e.g.,  $[0, 1]$ ) or use z-score normalization.
- Handling Missing Values: Impute or remove missing values to ensure data completeness.
- Feature Selection: Choose relevant features to improve clustering performance.

#### Model Processing:

- Initialization: Select  $k$  centroids randomly from the data.
- Distance Metric: Use Euclidean distance for measuring the distance between points and centroids.
- Iteration Limit: Set a maximum number of iterations (e.g., 100).

#### Evaluation:

- Distance Threshold: Choose a threshold based on the distribution of distances.
- Visualization: Plot the clusters and anomalies for visual inspection.
- Metrics: Use metrics like the Silhouette Score to evaluate cluster quality.

### 7.4 Advantages and Limitations

#### Advantages:

- Simplicity: Easy to understand and implement.
- Efficiency: Linear time complexity  $O(N \times k \times d \times I)$ , where  $I$  is the number of iterations.
- Interpretability: Anomalies are identified based on clear distance criteria.

#### Limitations:

- Fixed Number of Clusters: Requires specifying  $k$  beforehand, which may not be optimal for all datasets.
- Sensitivity to Initialization: Different initial centroids can lead to different results.
- Assumption of Spherical Clusters: Performs poorly if clusters have non-spherical shapes or varying densities.

## 7.5 Analysis of the K-Means Algorithm for Anomaly Detection

The K-Means algorithm for anomaly detection works by clustering data points into  $k$  distinct groups and identifying anomalies as points that are significantly distant from their assigned cluster centroids. The algorithm can be broken down into the following key phases:

1. Initialization: Randomly initialize  $k$  cluster centroids from the dataset  $\mathbf{X}$ .
2. Assignment Step: Assign each data point to the nearest centroid based on a distance metric (typically the Euclidean distance).
3. Update Step: Recalculate the centroids as the mean of the points assigned to each cluster.
4. Convergence Check: Stop the iteration when centroids do not change significantly or after a set number of iterations.
5. Anomaly Detection: Calculate the distance of each point to its assigned centroid and flag points whose distances exceed a specified threshold as anomalies.

This method leverages the assumption that anomalies are located far from the dense regions of normal data points.

## 7.6 Time Complexity Analysis

The time complexity of the K-Means algorithm for anomaly detection can be analyzed as follows:

1. Initialization:

Randomly selecting  $k$  initial centroids takes  $O(k)$  time.

2. Assignment Step:

For each of the  $N$  data points, the distance to each of the  $k$  centroids is calculated. Assuming the data points have  $d$  dimensions, this step takes:

$$O(N \times k \times d)$$

per iteration.

3. Update Step:

Recomputing the centroids by averaging the points assigned to each cluster takes  $O(N \times d)$  for each cluster, leading to:

$$O(k \times N \times d)$$

per iteration.

4. Convergence:

The algorithm typically converges within  $I$  iterations, where  $I$  is the maximum number of iterations.

||

5. Distance Calculation:



Calculating the distance of each data point to its assigned centroid for anomaly detection takes:

$$O(N \times d)$$

Total Time Complexity:

Putting it all together, the overall time complexity is:

$$O(I \times N \times k \times d)$$

where  $I$  is the number of iterations,  $N$  is the number of data points,  $k$  is the number of clusters, and  $d$  is the dimensionality of the data.

## 7.7 Correctness Proof

### Theorem:

The K-Means algorithm for anomaly detection converges to a local minimum of the sum of squared distances to the centroids (within-cluster sum of squares, WCSS).

### Proof:

1. Assignment Step: Each data point is assigned to the cluster whose centroid minimizes the distance to that point. This step reduces the within-cluster sum of squares (WCSS).
2. Update Step: The centroids are updated to the mean of the points assigned to each cluster. The mean minimizes the sum of squared distances within each cluster.
3. Monotonic Decrease: The WCSS decreases with each assignment and update step. Since the WCSS is bounded below by zero, the algorithm must converge to a local minimum.
4. Anomaly Detection: Once the centroids converge, anomalies are identified based on the distances to these centroids. The points with distances exceeding the threshold are flagged as anomalies.

Thus, the algorithm is guaranteed to converge to a local minimum of the WCSS, and anomalies are correctly identified based on the distance threshold [32, 22].

## 7.8 Summary - K-Means Algorithm for Anomaly Detection

The K-Means algorithm for anomaly detection is a straightforward and efficient method for identifying anomalies in unsupervised learning. It clusters the data into  $k$  groups and detects anomalies as data points that are significantly distant from their assigned cluster centroids. The algorithm's strengths lie in its simplicity, scalability, and clear interpretability. However, it assumes spherical clusters and requires specifying the number of clusters  $k$  beforehand. Despite these limitations, it remains a popular choice for anomaly detection in various applications, such as fraud detection, network security, and quality control [32, 22].

## 8 DBSCAN (Density-Based Spatial Clustering of Applications with Noise) for Anomaly Detection

The DBSCAN (Density-Based Spatial Clustering of Applications with Noise) algorithm is a density-based clustering method used for identifying clusters of arbitrary shapes and detecting anomalies (noise). Unlike partition-based methods such as K-Means, DBSCAN does not require specifying the number of clusters in advance. Instead, it groups points based on the density of their neighborhood defined by a radius  $\epsilon$  and a minimum number of points  $\text{minPts}$ . Points that do not belong to any dense region are labeled as anomalies. This makes DBSCAN particularly useful for unsupervised learning tasks where identifying irregularities is crucial, such as fraud detection, network security, and outlier detection in regression analysis [12].

### 8.1 Mathematical Formulations

Given a dataset  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ , where each  $\mathbf{x}_i \in \mathbb{R}^d$  is a  $d$ -dimensional point, the key concepts in DBSCAN are defined as follows:

1. Neighborhood: The  $\epsilon$ -neighborhood of a point  $\mathbf{x}_i$  is defined as:

$$N_\epsilon(\mathbf{x}_i) = \{\mathbf{x}_j \mid \text{Distance}(\mathbf{x}_i, \mathbf{x}_j) \leq \epsilon\}$$

2. Core Point: A point  $\mathbf{x}_i$  is a core point if:

$$|N_\epsilon(\mathbf{x}_i)| \geq \text{minPts}$$

3. Density-Reachable: A point  $\mathbf{x}_j$  is density-reachable from  $\mathbf{x}_i$  if there exists a chain of core points  $\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_m}$  such that:

$$\mathbf{x}_i = \mathbf{x}_{i_1}, \mathbf{x}_j = \mathbf{x}_{i_m}, \quad \text{and} \quad \mathbf{x}_{i_k} \in N_\epsilon(\mathbf{x}_{i_{k-1}})$$

4. Anomalies: Points that are not part of any cluster are labeled as anomalies (noise).

### 8.2 DBSCAN for Anomaly Detection Algorithm

The DBSCAN algorithm clusters the dataset based on the density of data points and identifies anomalies as points that do not meet the density criteria. The steps of Algorithm 6 are as follows:

### 8.3 Implementation Details

#### Data Preprocessing:

- Normalization: Scale the data to ensure consistent distance measurements.
- Distance Metric: Typically, **\*\*Euclidean distance\*\*** is used, but other metrics like Manhattan or cosine distance can be applied.

#### Model Processing:

- Parameter Tuning: Carefully choose  $\epsilon$  and  $\text{minPts}$  to optimize performance.
- Cluster Expansion: Expand clusters by recursively exploring neighbors.

---

**Algorithm 6** DBSCAN Algorithm for Anomaly Detection

---

```
1: function DBSCANANOMALYDETECTION( $\mathbf{X}, \epsilon, \text{minPts}$ )  
2:                                      $\triangleright$  Step 1: Initialize parameters  
3:    $N \leftarrow$  number of data points in  $\mathbf{X}$   
4:   Initialize cluster label array  $\mathbf{L} = [-1, -1, \dots, -1]$  (length  $N$ )  
5:   Initialize cluster counter  $c \leftarrow 0$   
6:   Initialize an empty set of anomalies  $\text{Anomalies} \leftarrow \emptyset$   
7:                                      $\triangleright$  Step 2: Iterate through all data points  
8:   for each data point  $\mathbf{x}_i \in \mathbf{X}$  do  
9:     if  $\mathbf{L}[i] \neq -1$  then  
10:      continue                                      $\triangleright$  Skip if already processed  
11:    end if  
12:                                      $\triangleright$  Step 3: Find neighbors within  $\epsilon$  distance  
13:     $\text{Neighbors} \leftarrow \{\mathbf{x}_j \mid \text{Distance}(\mathbf{x}_i, \mathbf{x}_j) \leq \epsilon\}$   
14:    if  $|\text{Neighbors}| < \text{minPts}$  then  
15:       $\mathbf{L}[i] \leftarrow 0$                                       $\triangleright$  Mark as noise (potential anomaly)  
16:       $\text{Anomalies} \leftarrow \text{Anomalies} \cup \{\mathbf{x}_i\}$   
17:      continue  
18:    end if  
19:                                      $\triangleright$  Step 4: Expand cluster  
20:     $c \leftarrow c + 1$                                       $\triangleright$  Create new cluster  
21:     $\mathbf{L}[i] \leftarrow c$   
22:    Initialize SeedSet  $\leftarrow \text{Neighbors}$   
23:    while SeedSet is not empty do  
24:      Remove a point  $\mathbf{x}_k$  from SeedSet  
25:      if  $\mathbf{L}[k] = 0$  then  
26:         $\mathbf{L}[k] \leftarrow c$                                       $\triangleright$  Previously marked as noise, now part of cluster  
27:      end if  
28:      if  $\mathbf{L}[k] \neq -1$  then  
29:        continue                                      $\triangleright$  Skip if already processed  
30:      end if  
31:       $\mathbf{L}[k] \leftarrow c$   
32:                                      $\triangleright$  Find neighbors of  $\mathbf{x}_k$   
33:       $\text{NewNeighbors} \leftarrow \{\mathbf{x}_j \mid \text{Distance}(\mathbf{x}_k, \mathbf{x}_j) \leq \epsilon\}$   
34:      if  $|\text{NewNeighbors}| \geq \text{minPts}$  then  
35:        SeedSet  $\leftarrow \text{SeedSet} \cup \text{NewNeighbors}$   
36:      end if  
37:    end while  
38:  end for  
39:  return  $\mathbf{L}, \text{Anomalies}$   
40: end function
```

---

## Evaluation:

- Silhouette Score: Measure cluster quality.
- Visualization: Plot clusters and anomalies to assess results.

## 8.4 Advantages and Limitations

### Advantages:

- No Need for  $k$ : Unlike K-Means, DBSCAN does not require specifying the number of clusters.
- Identifies Arbitrary Shapes: Can detect clusters of arbitrary shapes and sizes.
- Robust to Noise: Effectively identifies anomalies as noise points.

### Limitations:

- Parameter Sensitivity: Performance depends heavily on  $\epsilon$  and minPts.
- Scalability: Time complexity is  $O(N^2)$  for large datasets.

## 8.5 Time Complexity Analysis

The time complexity of the DBSCAN algorithm can be analyzed in the following steps:

1. Neighborhood Search: For each data point, finding its  $\epsilon$ -neighborhood requires computing the distance to all other points. The complexity for this operation is  $O(N)$ , where  $N$  is the number of data points.
2. Iteration Over All Points: Since the neighborhood search is performed for all  $N$  data points, the overall complexity of this step is  $O(N^2)$ .
3. Optimized Search: Using spatial indexing structures like KD-Trees or Ball Trees can reduce the neighborhood search complexity to  $O(N \log N)$ .

Total Time Complexity:

$$O(N^2) \quad (\text{without spatial indexing})$$

$$O(N \log N) \quad (\text{with spatial indexing})$$

In practice, DBSCAN with spatial indexing is efficient for datasets with a moderate number of dimensions [12].

## 8.6 Analysis of the DBSCAN Algorithm for Anomaly Detection

DBSCAN clusters data points based on density and flags anomalies as points that do not belong to any dense region. The main phases of the algorithm are:

1. Neighborhood Identification: For each point, identify neighbors within a radius  $\epsilon$ .
2. Cluster Formation: If a point has at least minPts neighbors, it becomes a core point, and a new cluster is formed.
3. Cluster Expansion: Expand the cluster by iteratively including density-reachable points.

4. Anomaly Detection: Points that do not meet the density criteria are labeled as anomalies (noise).

DBSCAN is effective for detecting anomalies in datasets where clusters have arbitrary shapes and varying densities. However, its performance depends on the choice of  $\epsilon$  and minPts.

## 8.7 Correctness Proof

### Theorem:

The DBSCAN algorithm correctly identifies clusters and anomalies based on density reachability.

### Proof:

1. Core Points: A point  $\mathbf{x}_i$  is classified as a core point if its  $\epsilon$ -neighborhood contains at least minPts points. This ensures that all core points lie in dense regions.
2. Density-Reachable Points: If  $\mathbf{x}_j$  is density-reachable from  $\mathbf{x}_i$ , it belongs to the same cluster as  $\mathbf{x}_i$ . The recursive expansion ensures that all points in a dense region are connected.
3. Anomalies: Points that are not density-reachable from any core point are labeled as anomalies. These points lie in low-density regions and do not meet the clustering criteria.

Since the algorithm iterates through all points and expands clusters based on density reachability, it guarantees that dense regions are correctly identified as clusters, and outliers are labeled as anomalies.

## Summary for Correctness Proof

The correctness of DBSCAN is based on the definitions of core points, density-reachable points, and noise. The algorithm accurately clusters dense regions and identifies anomalies by checking whether points meet the density criteria. This ensures reliable anomaly detection in datasets with complex structures [43].

## 8.8 Summary - DBSCAN for Anomaly Detection

DBSCAN is a powerful unsupervised learning algorithm for clustering and anomaly detection. It identifies clusters based on density and flags points in low-density regions as anomalies. DBSCAN is effective for datasets with clusters of arbitrary shapes and varying densities. Its advantages include not requiring the number of clusters in advance and robustness to noise. However, it is sensitive to the choice of  $\epsilon$  and minPts and can be computationally expensive for large datasets. Despite these limitations, DBSCAN remains a widely-used method for anomaly detection in data science [12, 43].

## 9 One-Class SVM for Anomaly Detection

The One-Class Support Vector Machine (One-Class SVM) is an unsupervised learning algorithm designed for anomaly detection by learning the boundary of normal data points. It fits within the context of regression analysis by determining whether a new observation fits the learned distribution of the training data. Unlike standard SVMs, which separate two classes, the One-Class SVM aims to capture the region where most of the data lies and classify points outside this region as anomalies [42]. This makes it highly suitable for tasks where anomalies (outliers) need to be identified without labeled examples.

### 9.1 Mathematical Formulations

Given a dataset of  $N$  observations represented as  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ , where  $\mathbf{x}_i \in \mathbb{R}^d$ , the One-Class SVM finds a hyperplane that separates the data points from the origin by solving the following optimization problem:

$$\min_{\mathbf{w}, \rho, \xi_i} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu N} \sum_{i=1}^N \xi_i - \rho,$$

subject to:

$$(\mathbf{w} \cdot \phi(\mathbf{x}_i)) \geq \rho - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, N.$$

Where: -  $\mathbf{w}$  is the weight vector defining the hyperplane. -  $\phi(\mathbf{x}_i)$  is a feature mapping function (depends on the chosen kernel). -  $\rho$  is the offset. -  $\xi_i$  are slack variables to allow for some points to be outside the boundary. -  $\nu$  controls the proportion of anomalies.

The decision function  $f(\mathbf{x})$  is given by:

$$f(\mathbf{x}) = \text{sign}((\mathbf{w} \cdot \phi(\mathbf{x})) - \rho).$$

A data point  $\mathbf{x}_i$  is classified as an anomaly if  $f(\mathbf{x}_i) < 0$ .

### 9.2 One-Class SVM for Anomaly Detection Algorithm

### 9.3 Implementation Details

#### Data Preprocessing:

- Normalization: Scale the data to have zero mean and unit variance using z-score normalization.
- Feature Selection: Select relevant features to improve model performance.
- Handling Missing Values: Impute or remove missing values to ensure data completeness.

---

**Algorithm 7** One-Class SVM for Anomaly Detection

---

```
1: function ONECLASSSVMANOMALYDETECTION( $\mathbf{X}$ ,  $\nu$ , kernel)
2:                                     ▷ Step 1: Initialize parameters
3:    $N \leftarrow$  number of data points in  $\mathbf{X}$ 
4:    $\nu \leftarrow$  Anomaly proportion (e.g., 0.05 for 5% anomalies)
5:   kernel  $\leftarrow$  Kernel function (e.g., 'RBF', 'linear', 'polynomial')
6:                                     ▷ Step 2: Train the One-Class SVM model
7:   Train the One-Class SVM model OCSVM with parameter  $\nu$  and kernel function on dataset  $\mathbf{X}$ 
8:                                     ▷ Step 3: Calculate decision function scores
9:   for each data point  $\mathbf{x}_i \in \mathbf{X}$  do
10:     $s_i \leftarrow$  OCSVM.DecisionFunction( $\mathbf{x}_i$ )
11:   end for
12:                                     ▷ Step 4: Identify anomalies
13:   Anomalies  $\leftarrow \{\mathbf{x}_i \mid s_i < 0\}$ 
14:   return OCSVM, Anomalies
15: end function
```

---

### Model Processing:

- Kernel Choice: Common kernels include:
  - ◊ Radial Basis Function (RBF): Handles non-linear boundaries.
  - ◊ Linear Kernel: Suitable for linearly separable data.
  - ◊ Polynomial Kernel: Captures more complex relationships.
- Training: Fit the One-Class SVM model using the specified  $\nu$  and kernel.

### Evaluation:

- Decision Function: Evaluate the decision function scores for each data point.
- Visualization: Plot the data points and mark anomalies for visual inspection.
- Performance Metrics: Use precision, recall, and F1-score if ground truth labels are available.

## 9.4 Advantages and Limitations

### Advantages:

- Unsupervised: No need for labeled data.
- Flexible: Works well with different kernel functions for linear and non-linear data.
- Robust: Can detect anomalies even in high-dimensional datasets.

### Limitations:

- Sensitive to Parameters: Choice of  $\nu$  and kernel significantly affects performance.
- Computational Complexity: Training can be slow for large datasets.
- High False Positive Rate: May flag normal points as anomalies if the boundary is too tight [42].

## 9.5 Time Complexity Analysis

The time complexity of the One-Class SVM for anomaly detection is primarily determined by the training phase:

- Training Phase: The complexity is  $O(N^2)$  for training on  $N$  data points. For large datasets, this can be computationally expensive.
- Prediction Phase: Evaluating the decision function for each data point takes  $O(N \cdot d)$ , where  $d$  is the feature dimension.

## 9.6 Correctness Proof

**Theorem:**

The One-Class SVM algorithm correctly identifies anomalies by finding a decision boundary that encloses the majority of the data points.

**Proof:**

1. Optimization: The One-Class SVM solves an optimization problem that minimizes the boundary's complexity while enclosing most data points.
2. Decision Boundary: The decision function  $s_i = (\mathbf{w} \cdot \phi(\mathbf{x}_i)) - \rho$  ensures that data points within the boundary have positive scores, while anomalies lie outside the boundary with negative scores.
3. Guaranteed Support: By controlling the parameter  $\nu$ , the algorithm guarantees that a proportion  $\nu$  of the data points will be identified as anomalies [42].

## 9.7 Summary - One-Class SVM for Anomaly Detection Algorithm

The One-Class SVM for Anomaly Detection is a powerful unsupervised learning method for detecting anomalies by learning a boundary that encapsulates the majority of the data. It leverages kernel methods to handle non-linear data and can identify anomalies without requiring labeled data. Despite its computational complexity and sensitivity to parameter tuning, it remains a widely used technique for anomaly detection in high-dimensional datasets [42].



## 10 Evaluation Metrics for Unsupervised Learning

The evaluation of unsupervised learning models is critical to ensure their effectiveness and reliability across various tasks. Unlike supervised learning, where labeled data is available for performance comparison, unsupervised methods rely on indirect metrics to quantify their success. This lecture explores evaluation metrics tailored for unsupervised tasks, focusing on clustering, dimensionality reduction, and anomaly detection.

**Purpose of Evaluation Metrics** Evaluation metrics for unsupervised learning provide insights into the quality of clusters, the preservation of variance in data, and the detection of anomalies. These metrics ensure that models such as K-Means clustering, Eigen Decomposition, 2D-DCT, DBSCAN, and One-Class SVM achieve their intended objectives. Metrics like Silhouette Score, Explained Variance, and Precision-Recall curves for anomalies guide model selection and optimization [8].

**Challenges in Unsupervised Learning** In the absence of ground truth labels, unsupervised learning evaluation often depends on intrinsic measures (e.g., cluster cohesion) or extrinsic measures (e.g., comparison to known results). This lecture outlines metrics applicable to unsupervised tasks, including clustering for text-based data, dimensionality reduction [3], and anomaly detection [42].

### 10.1 Clustering Evaluation Metrics

**Silhouette Score** The Silhouette Score evaluates the quality of clustering by considering both intra-cluster cohesion and inter-cluster separation. For each data point  $\mathbf{x}_i$ , the silhouette score is defined as:

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)},$$

where:

- $a_i$  is the average distance between  $\mathbf{x}_i$  and other points in its cluster.
- $b_i$  is the average distance between  $\mathbf{x}_i$  and points in the nearest neighboring cluster.

A high Silhouette Score (close to 1) indicates well-separated and cohesive clusters, making it ideal for evaluating K-Means clustering for text-based data [8].

**Elbow Method for K** The Elbow Method determines the optimal number of clusters  $k$  in K-Means by analyzing the within-cluster sum of squares (WCSS):

$$\text{WCSS} = \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{c}_{a_i}\|^2,$$

where  $\mathbf{c}_{a_i}$  is the centroid of the cluster to which  $\mathbf{x}_i$  is assigned. The "elbow point" is where adding more clusters yields diminishing returns in WCSS reduction [8].

**Calinski-Harabasz Index** The Calinski-Harabasz Index measures the ratio of between-cluster variance to within-cluster variance, favoring well-separated and compact clusters:

$$\text{CH} = \frac{\text{trace}(\mathbf{B})}{\text{trace}(\mathbf{W})} \cdot \frac{N - k}{k - 1},$$

where  $\mathbf{B}$  and  $\mathbf{W}$  are the between-cluster and within-cluster scatter matrices, respectively [8].

## 10.2 Dimensionality Reduction Metrics

**Explained Variance** For Eigen Decomposition and 2D-DCT, explained variance quantifies how much variance is retained after transformation. Given a dataset  $\mathbf{X}$ , the explained variance for the  $j$ -th component is:

$$\text{Explained Variance Ratio}_j = \frac{\lambda_j}{\sum_{i=1}^d \lambda_i},$$

where  $\lambda_j$  is the eigenvalue corresponding to the  $j$ -th component. High explained variance indicates effective dimensionality reduction [3].

**Reconstruction Error** Reconstruction error measures the loss of information during dimensionality reduction and is particularly relevant for the Inverse 2D-DCT:

$$\text{Reconstruction Error} = \|\mathbf{X} - \mathbf{X}_{\text{reconstructed}}\|_F^2,$$

where  $\|\cdot\|_F$  denotes the Frobenius norm [3].

## 10.3 Anomaly Detection Metrics

**Precision-Recall Curve** For anomaly detection algorithms such as K-Means, DBSCAN, and One-Class SVM, the Precision-Recall (PR) curve evaluates the trade-off between precision and recall. These metrics are defined as:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}},$$

where TP, FP, and FN represent true positives, false positives, and false negatives, respectively. The PR curve emphasizes performance on the minority class (anomalies) [42].

**Area Under PR Curve (PR-AUC)** The PR-AUC summarizes the overall performance of anomaly detection methods, with higher values indicating better precision-recall trade-offs [42].

**Cluster-Based Metrics** For DBSCAN, anomaly detection relies on identifying noise points (cluster label 0) as anomalies. Metrics such as the number of noise points and their average distance to clusters provide insights into the effectiveness of the detection [bishop2006pattern].

## 10.4 Time Complexity of Evaluation Metrics

- Silhouette Score:  $O(N^2)$  for pairwise distance calculations, where  $N$  is the number of data points.
- Explained Variance:  $O(d^3)$  for Eigen Decomposition on a  $d$ -dimensional dataset.
- Precision-Recall Curve:  $O(N)$  for evaluating predictions on  $N$  data points.
- Reconstruction Error:  $O(Nd)$  for calculating the Frobenius norm of reconstruction loss.

## 10.5 Summary of Evaluation Metrics for Unsupervised Learning

This section highlighted key evaluation metrics tailored for unsupervised methods. Clustering metrics, such as the Silhouette Score and Calinski-Harabasz Index, provide valuable insights into cluster quality, particularly for algorithms like K-Means and DBSCAN. Dimensionality reduction metrics, including explained variance and reconstruction error, assess the effectiveness of techniques like Eigen Decomposition and 2D-DCT. For anomaly detection, metrics such as Precision-Recall curves and PR-AUC are instrumental in evaluating the performance of methods like One-Class SVM. These metrics collectively ensure that unsupervised models are robust, interpretable, and aligned with specific task objectives.

## 11 Recent Advances in Unsupervised Learning

Unsupervised learning is a critical area of machine learning that focuses on uncovering hidden patterns and structures in unlabeled data. Over the past decade, advancements in unsupervised learning have broadened its applications across diverse domains, including natural language processing (NLP), computer vision, and anomaly detection. These developments have addressed key challenges such as scalability, interpretability, and the handling of complex, high-dimensional data.

This section explores the latest innovations in clustering, dimensionality reduction, anomaly detection, and interpretability techniques, highlighting their role in solving real-world problems. Enhanced clustering algorithms, such as K-Means++ and DBSCAN, have improved the ability to identify patterns in text and image datasets. Similarly, advanced dimensionality reduction methods like Eigen Decomposition and 2D-DCT provide interpretable low-dimensional representations of high-dimensional data. Meanwhile, robust anomaly detection algorithms, including One-Class SVM and autoencoder-based methods, enable precise identification of rare events across large and complex datasets. Finally, interpretability techniques, such as SHAP and visualization tools like UMAP, ensure that the outputs of unsupervised models are accessible and trustworthy for domain experts and practitioners. These advances collectively represent a significant step forward in making unsupervised learning scalable, transparent, and impactful.

### 11.1 Clustering in Unsupervised Learning

Clustering is a foundational technique in unsupervised learning, aiming to group data points into clusters based on their similarity. Advances in clustering algorithms have expanded their applicability to high-dimensional, dynamic, and complex datasets, making clustering a critical tool in fields like Natural Language Processing, image analysis, and bioinformatics.

#### K-Means Variants

The standard K-Means algorithm has been refined to address its limitations, such as sensitivity to initialization and poor performance with non-spherical clusters. For example, K-Means++ improves initialization by selecting initial centroids probabilistically, leading to better convergence and higher-quality clusters [5]. Techniques such as Mini-Batch K-Means enhance scalability by processing data in smaller batches, making clustering feasible for large datasets. Clustering evaluation metrics like Silhouette Score and Calinski-Harabasz Index provide quantitative assessments of cluster compactness and separation, ensuring robust evaluation of clustering performance [41, 9].

#### DBSCAN for Complex Structures

Density-based clustering algorithms, such as DBSCAN, overcome the limitations of K-Means by identifying clusters of arbitrary shapes and varying densities. DBSCAN groups points into density-connected regions and identifies noise points as outliers, making it highly effective for non-spherical data distributions [12]. This capability is particularly valuable in applications such as text-based clustering, where data often exhibit irregular structures. Advanced variants like HDBSCAN improve upon DBSCAN by automatically determining the optimal value for parameters like  $\epsilon$ , enhancing usability and robustness [10].

## Spectral Clustering for Graph-Based Data

Spectral clustering leverages the eigenvalues of similarity matrices to partition data into clusters. Unlike K-Means, spectral clustering can handle non-convex clusters and is particularly useful for graph-based data, where relationships between data points are critical. This method finds applications in community detection in social networks and clustering in bioinformatics datasets [44].

## Gaussian Mixture Models (GMMs)

GMMs model clusters as a mixture of Gaussian distributions, providing a probabilistic approach to clustering. Unlike K-Means, which assigns each point to a single cluster, GMMs compute the probability of a data point belonging to multiple clusters. This soft clustering approach is particularly beneficial for datasets with overlapping clusters, such as in image segmentation and market segmentation [8].

## Hierarchical Clustering

Hierarchical clustering builds a tree-like structure (dendrogram) to represent data clusters at various levels of granularity. Agglomerative methods start with individual data points and merge them into larger clusters, while divisive methods begin with all data points in a single cluster and iteratively split them. This approach is useful for datasets where a nested hierarchy of clusters is meaningful, such as taxonomy generation in biology [37].

## Subspace Clustering for High-Dimensional Data

Subspace clustering addresses the challenges of clustering in high-dimensional datasets by identifying clusters in lower-dimensional subspaces. Algorithms like CLIQUE and PROCLUS effectively discover meaningful patterns in datasets with many irrelevant dimensions, making them valuable in gene expression analysis and recommendation systems [1, 2].

## Deep Clustering

Recent advances in deep learning have enabled deep clustering, which integrates feature learning and clustering into a unified framework. Methods like Deep Embedded Clustering (DEC) jointly optimize a neural network for feature extraction and clustering, improving performance on complex datasets, such as images and text [45]. These approaches leverage the representational power of neural networks to uncover hierarchical and non-linear patterns in data.

## 11.2 Dimensionality Reduction and Representation Learning

Dimensionality reduction techniques aim to project high-dimensional data into a lower-dimensional space while preserving essential structures and relationships. These methods enhance interpretability, computational efficiency, and the ability to handle high-dimensional datasets in unsupervised learning.

## **Eigen Decomposition and Factor Analysis**

Eigen Decomposition and Factor Analysis are classical techniques that extract latent structures by analyzing the covariance matrix of the data. Eigenvectors and eigenvalues provide principal directions of variance, enabling variance expansion and latent factor discovery. These methods are essential in exploratory analysis and dimensionality reduction tasks across domains like social science and bioinformatics [24].

## **2D Discrete Cosine Transform (2D-DCT)**

The 2D-DCT is a powerful method that transforms spatial data, such as images, into the frequency domain, capturing vertical, horizontal, and diagonal patterns. It reduces data redundancy, making it a cornerstone in applications like image compression (e.g., JPEG) and feature extraction for pattern recognition. Advances in 2D-DCT and its inverse (2D-IDCT) have improved representation learning for image and video datasets [3].

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

t-SNE is a non-linear dimensionality reduction technique that visualizes high-dimensional data by preserving local similarities in a low-dimensional space. It is widely used in clustering and exploratory data analysis, particularly for visualizing complex datasets like text embeddings and biological data. t-SNE minimizes the divergence between probability distributions of data points in high- and low-dimensional spaces, creating interpretable visual representations [31].

## **Principal Component Analysis (PCA)**

PCA is a linear dimensionality reduction method that identifies principal components by maximizing variance along orthogonal directions. It is widely used for preprocessing, noise reduction, and feature selection. PCA is computationally efficient and applicable to diverse datasets, making it a foundational tool in dimensionality reduction [39].

## **Autoencoders**

Autoencoders are neural network-based representation learning models that encode data into a compressed latent space and then reconstruct the input data. Variants like denoising autoencoders and variational autoencoders (VAEs) enhance robustness and probabilistic modeling capabilities. Autoencoders are extensively used in image generation, anomaly detection, and data compression tasks [25].

## **Uniform Manifold Approximation and Projection (UMAP)**

UMAP is a state-of-the-art technique for non-linear dimensionality reduction that preserves both local and global structures. It constructs a high-dimensional graph of the data, then optimizes its low-dimensional representation. UMAP is faster and more scalable than t-SNE, making it suitable for large datasets in fields like genomics, NLP, and computer vision [35].

### 11.3 Anomaly Detection in Unsupervised Learning

Anomaly detection focuses on identifying rare or unusual instances that deviate from the majority of the data. Recent advancements have improved the precision, scalability, and applicability of anomaly detection methods across various domains.

#### Density-Based Methods

Density-based methods like DBSCAN and its variants detect anomalies by identifying low-density regions as outliers. These methods are particularly effective in applications like fraud detection, network intrusion detection, and environmental monitoring, where anomalies are sparse and irregularly distributed. The algorithm clusters points based on density connectivity, and points that do not belong to any cluster are labeled as anomalies [12].

#### One-Class SVM

The One-Class SVM formulates anomaly detection as a one-class classification problem, where the majority of the data forms a single class, and anomalies are identified as deviations from this class. By leveraging kernel functions such as RBF and polynomial kernels, the One-Class SVM models complex data distributions and detects anomalies in high-dimensional spaces. It is widely used in tasks like text-based anomaly detection and network security [42].

#### Autoencoder-Based Methods

Autoencoders are neural network models that learn compact representations of data through encoding and decoding. For anomaly detection, autoencoders are trained to reconstruct normal data with minimal error. When anomalies are fed into the model, the reconstruction error is typically high, enabling their identification. Variants like Variational Autoencoders (VAEs) and Deep Autoencoders enhance robustness and scalability, making them ideal for detecting anomalies in images, videos, and time-series data [25].

#### Isolation Forest

The Isolation Forest algorithm isolates anomalies by recursively partitioning the data space. Anomalies, being rare and different, are easier to isolate, requiring fewer splits on average. The method's linear time complexity and small memory footprint make it well-suited for large-scale datasets in domains like financial fraud detection and manufacturing defect analysis [27].

#### Time-Series Anomaly Detection

For sequential data, methods like Long Short-Term Memory (LSTM) networks and Seasonal Hybrid Extreme Studentized Deviate (S-H-ESD) models are used. LSTM networks learn temporal dependencies and can detect anomalies based on deviations from expected patterns. S-H-ESD is a statistical method that identifies anomalies in time-series data by accounting for seasonality and extreme deviations. These methods are widely applied in system monitoring, predictive maintenance, and healthcare analytics [34, 20].

## 11.4 Interpretability in Unsupervised Learning

Interpretability techniques are crucial for understanding the outputs of unsupervised learning models. Recent advancements emphasize making unsupervised methods more transparent and accessible, enabling domain experts and practitioners to trust and effectively utilize these models.

### Feature Contribution Analysis

Techniques like SHAP (SHapley Additive exPlanations) provide a framework for quantifying the contribution of each feature to cluster assignment or anomaly detection. By attributing predictions to individual features, SHAP offers consistent and additive explanations, making unsupervised models more transparent [30]. For instance, in anomaly detection using One-Class SVM or clustering with K-Means, SHAP can highlight the key features driving the model's decisions.

### Visualization of Reduced Dimensions

Dimensionality reduction techniques, such as Eigen Decomposition and 2D-DCT, inherently support interpretability by reducing data complexity and enabling intuitive visualization of key features. Visualizing data in two or three dimensions allows practitioners to identify clusters, anomalies, and patterns that may otherwise remain hidden in high-dimensional spaces [24, 3].

### Cluster Explanation with Prototypes and Exemplars

Prototype-based techniques select representative data points from clusters to provide interpretable summaries. Methods like K-Medoids and exemplar-based clustering use central or typical instances to represent clusters, making the results understandable for domain experts [38]. These techniques are particularly useful in applications like customer segmentation or medical diagnosis.

### Rule-Based Explanations for Clustering

Rule-based methods generate human-readable rules to explain cluster structures. For example, decision tree-based approaches can extract if-then rules to describe clusters formed by K-Means or DBSCAN. These methods provide interpretable insights into cluster boundaries and features, facilitating their application in fields like retail analytics and finance [4].

### Interactive Visualization Tools

Modern visualization tools, such as t-SNE (t-distributed Stochastic Neighbor Embedding) and UMAP (Uniform Manifold Approximation and Projection), enable dynamic exploration of high-dimensional data. These tools are often combined with clustering techniques to visually analyze the relationships between clusters and identify anomalies. Interactive platforms, such as Plotly and Tableau, further enhance interpretability by allowing users to drill down into specific clusters or anomalies [35].

## 11.5 Applications and Future Directions

Advances in unsupervised learning have significantly broadened its applicability across diverse fields. These methods are now integral to solving complex challenges in data science and AI.



## Applications Across Domains

- **Natural Language Processing (NLP):** Clustering methods like K-Means are widely applied in NLP for grouping similar documents, enabling tasks such as topic modeling, document classification, and information retrieval. Anomaly detection techniques, such as One-Class SVM, are employed to identify rare patterns in text corpora, such as unusual linguistic structures in sentiment analysis or errors in machine translation. These methods streamline unstructured textual data processing, a critical need in modern AI systems.
- **Healthcare:** Unsupervised learning is transformative in healthcare, where clustering algorithms group patients with similar genetic or lifestyle profiles, aiding personalized medicine. Dimensionality reduction techniques like Principal Component Analysis (PCA) simplify complex datasets, such as electronic health records or genomic data, enhancing the analysis of disease progression and treatment outcomes. These applications underpin AI-driven diagnostics and precision medicine.
- **Cybersecurity:** Anomaly detection algorithms enhance cybersecurity by identifying irregularities in network traffic that may indicate cyberattacks or insider threats. These models enable adaptive threat detection systems by automatically flagging suspicious activities without requiring labeled training data, thus forming a cornerstone of modern AI-powered cybersecurity systems.
- **Astronomy and Scientific Research:** Clustering and dimensionality reduction are integral to analyzing high-dimensional astronomical datasets. For instance, clustering algorithms categorize celestial objects, such as stars and galaxies, based on their spectral properties, while dimensionality reduction methods extract meaningful patterns for efficient analysis. These tools facilitate AI-driven discoveries in vast datasets, accelerating advancements in scientific exploration.

## Applications of Dimensionality Reduction

- **Image Analysis:** Techniques like the 2D Discrete Cosine Transform (2D-DCT) reduce the dimensionality of image data, capturing essential features for tasks such as compression, facial recognition, and object detection. These methods are fundamental in AI-powered vision systems, making image processing more efficient and accurate.
- **Marketing and Customer Segmentation:** Dimensionality reduction simplifies the analysis of customer datasets, enabling businesses to identify key features driving purchasing behavior and engagement. This improves clustering accuracy in customer segmentation and supports targeted marketing campaigns using AI-driven analytics.
- **Supply Chain and Logistics:** Dimensionality reduction optimizes supply chain systems by summarizing data on suppliers, inventory levels, and delivery routes. Reduced feature sets help detect patterns, streamline operations, and enhance decision-making in dynamic environments, enabling robust AI-powered logistics solutions.

## Applications of Anomaly Detection

- **Fraud Detection:** Anomaly detection models isolate unusual patterns in financial datasets to identify fraudulent activities, such as credit card fraud or insider trading. These models

enhance real-time fraud detection systems by operating unsupervised across diverse datasets, ensuring adaptability and scalability in financial AI applications.

- **Cybersecurity:** Detecting anomalies in network traffic helps identify potential malware, phishing attempts, or insider threats. Unsupervised anomaly detection methods, such as DBSCAN and One-Class SVM, allow adaptive and scalable AI systems to protect sensitive infrastructure.
- **Manufacturing and Quality Control:** Anomaly detection identifies defects in manufacturing processes by analyzing sensor data or product quality metrics. These models ensure high standards of production and maintenance by flagging deviations from expected behavior in real-time.

## Applications of Interpretability Techniques

- **Healthcare:** Feature contribution analysis using SHAP helps explain why certain patients are grouped together or identified as anomalies. These insights enable practitioners to understand the drivers of AI-based diagnoses and tailor treatment plans effectively.
- **Retail and Marketing:** Interpretability tools, such as rule-based clustering explanations, help businesses understand customer segments and refine their marketing strategies. These tools enhance trust in AI-driven systems by providing clear insights into customer behavior.
- **Environmental Science:** Visualization of reduced dimensions allows researchers to uncover patterns in climate data, such as identifying regions with similar weather trends or detecting anomalies in temperature records. These insights support AI-driven decision-making in sustainability and resource management.

Future work in unsupervised learning focuses on addressing several key challenges. Researchers aim to develop scalable algorithms capable of handling large datasets efficiently, which is crucial given the increasing volume and complexity of modern data. Efforts are also directed toward enhancing the interpretability of unsupervised models, enabling better understanding and trust in their outcomes, particularly in high-stakes applications. Additionally, integrating domain knowledge into unsupervised learning frameworks is a priority, as it can guide the learning process, improve model performance, and ensure the results are more relevant and actionable for specific contexts.

## 11.6 Summary of Recent Advances in Unsupervised Learning

Recent progress in unsupervised learning has substantially enhanced its effectiveness and versatility, enabling its application to increasingly complex datasets across various domains. Clustering algorithms, such as K-Means++, DBSCAN, and deep clustering, provide robust solutions for pattern discovery in structured and unstructured data. Dimensionality reduction techniques, including Eigen Decomposition, PCA, and UMAP, enable efficient processing and visualization of high-dimensional datasets, ensuring interpretability and scalability. Advanced anomaly detection methods, such as One-Class SVM, Isolation Forest, and autoencoder-based approaches, offer reliable tools for identifying outliers and irregularities in domains like cybersecurity, healthcare, and fraud detection.

Additionally, interpretability techniques have addressed the “black box” nature of unsupervised models, providing transparency through feature contribution analysis, rule-based explanations, and

dynamic visualization tools. These methods ensure that unsupervised learning outputs are not only accurate but also actionable for practitioners in fields ranging from marketing to environmental science. Future work in unsupervised learning aims to develop scalable algorithms, integrate domain knowledge, and further enhance interpretability, solidifying its role as a cornerstone of modern machine learning.

## 12 Unsupervised Learning Methods in scikit-learn

Scikit-learn is a versatile Python library for machine learning, offering robust tools for unsupervised learning alongside its capabilities for supervised tasks. Unsupervised learning methods in Scikit-learn enable the discovery of hidden structures, patterns, or anomalies in data without requiring labeled outputs. This lecture highlights Scikit-learn's functionality for clustering, dimensionality reduction, anomaly detection, and interpretability, demonstrating their implementation and evaluation.

### 12.1 Overview of Scikit-learn

Scikit-learn is built on core Python libraries like NumPy, SciPy, and matplotlib, ensuring compatibility with the broader data science ecosystem. Key features of Scikit-learn include:

- A consistent API for unsupervised learning models.
- Preprocessing utilities for feature scaling, normalization, and missing value handling.
- Support for hyperparameter tuning with grid search and randomized search.
- Comprehensive tools for evaluating clustering, dimensionality reduction, and anomaly detection.

### 12.2 Common Steps for Unsupervised Learning with Scikit-learn

A typical workflow for unsupervised learning in Scikit-learn involves the following steps:

**Step 1: Data Preparation** Load the dataset and preprocess it:

```
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

**Step 2: Model Selection and Training** Choose an unsupervised learning algorithm and fit it to the data:

```
from sklearn.cluster import KMeans

model = KMeans(n_clusters=3, random_state=42)
model.fit(X_scaled)
```

**Step 3: Evaluation** Evaluate the results using metrics like Silhouette Score or explained variance:

```
from sklearn.metrics import silhouette_score

score = silhouette_score(X_scaled, model.labels_)
print(f"Silhouette Score: {score}")
```

## 12.3 Supported Unsupervised Learning Methods

Scikit-learn supports a variety of unsupervised learning methods:

### 12.3.1 Clustering Algorithms

**K-Means Clustering** K-Means groups data into clusters based on distance from centroids:

```
from sklearn.cluster import KMeans

model = KMeans(n_clusters=5, random_state=42)
model.fit(X_scaled)
```

**DBSCAN (Density-Based Spatial Clustering of Applications with Noise)** DBSCAN identifies clusters based on density and detects outliers as noise:

```
from sklearn.cluster import DBSCAN

model = DBSCAN(eps=0.5, min_samples=5)
model.fit(X_scaled)
```

**Spectral Clustering** Spectral clustering is effective for graph-based and non-convex data structures:

```
from sklearn.cluster import SpectralClustering

model = SpectralClustering(n_clusters=3, random_state=42)
model.fit(X_scaled)
```

### 12.3.2 Dimensionality Reduction Methods

**Principal Component Analysis (PCA)** PCA projects high-dimensional data into a lower-dimensional space:

```
from sklearn.decomposition import PCA

pca = PCA(n_components=2)
X_reduced = pca.fit_transform(X_scaled)
```

**t-Distributed Stochastic Neighbor Embedding (t-SNE)** t-SNE visualizes high-dimensional data by preserving local similarities:

```
from sklearn.manifold import TSNE

tsne = TSNE(n_components=2, random_state=42)
X_embedded = tsne.fit_transform(X_scaled)
```

**Uniform Manifold Approximation and Projection (UMAP)** UMAP is a fast, scalable dimensionality reduction technique:

```
import umap

umap_model = umap.UMAP(n_components=2)
X_umap = umap_model.fit_transform(X_scaled)
```

### 12.3.3 Anomaly Detection Methods

**One-Class SVM** One-Class SVM identifies anomalies based on deviations from a learned boundary:

```
from sklearn.svm import OneClassSVM

model = OneClassSVM(kernel='rbf', gamma=0.1, nu=0.05)
model.fit(X_scaled)
```

**Isolation Forest** Isolation Forest isolates anomalies by partitioning the data:

```
from sklearn.ensemble import IsolationForest

model = IsolationForest(random_state=42)
model.fit(X_scaled)
```

**DBSCAN for Outlier Detection** DBSCAN identifies anomalies as noise points:

```
outliers = X[model.labels_ == -1]
```

### 12.3.4 Interpretability Techniques

**Feature Importance with SHAP** SHAP explains feature contributions for clustering or anomaly detection:

```
import shap

explainer = shap.Explainer(model.predict, X_scaled)
shap_values = explainer(X_scaled)
shap.summary_plot(shap_values, X_scaled)
```

**Visualization Tools** Use PCA, t-SNE, or UMAP for visualizing clusters or anomalies:

```
import matplotlib.pyplot as plt

plt.scatter(X_reduced[:, 0], X_reduced[:, 1], c=model.labels_)
plt.show()
```

## 12.4 Evaluation of Unsupervised Models

Scikit-learn provides tools for evaluating unsupervised learning models:

**Silhouette Score** Measures cluster quality by comparing intra-cluster distance to inter-cluster distance:

```
from sklearn.metrics import silhouette_score

score = silhouette_score(X_scaled, model.labels_)
print(f"Silhouette Score: {score}")
```

**Explained Variance (PCA)** Indicates how much variance is captured by each principal component:

```
print(f"Explained Variance Ratio: {pca.explained_variance_ratio_}")
```

**Reconstruction Error (Autoencoders)** Evaluates anomaly detection models based on reconstruction loss:

```
reconstruction_error = np.mean((X_scaled - X_reconstructed)**2, axis=1)
```

## 12.5 Applications of Scikit-learn in Unsupervised Learning

Scikit-learn's unsupervised learning capabilities are widely applied across domains:

- **NLP:** Document clustering and topic modeling.
- **Healthcare:** Patient segmentation and anomaly detection in medical data.
- **Cybersecurity:** Detecting network intrusions and malware.
- **Finance:** Fraud detection and risk assessment.
- **Retail:** Customer segmentation and product recommendation.

## 12.6 Conclusion

Scikit-learn provides a comprehensive framework for implementing unsupervised learning methods. With a wide array of algorithms and evaluation tools, Scikit-learn empowers practitioners to efficiently discover patterns, reduce dimensionality, and detect anomalies in diverse datasets. Mastering these tools allows students and professionals to address real-world challenges in data science and AI.

## 13 Summary

Unsupervised learning techniques enable the discovery of underlying patterns, structures, and anomalies in unlabeled datasets, serving as critical tools for data exploration and representation. This document delves into various unsupervised learning approaches, including clustering algorithms such as K-Means and DBSCAN, which group data based on similarity, and dimensionality reduction methods like PCA and t-SNE, which transform high-dimensional data into interpretable forms. Additionally, anomaly detection techniques, including One-Class SVM and Isolation Forest, are discussed for identifying rare or unusual events.

Mathematical foundations, such as eigen decomposition and density estimation, underpin these techniques, providing frameworks for clustering objectives, representation learning, and model interpretability. The document also highlights advanced topics like feature attribution using SHAP and visualization with t-SNE for explaining patterns.

Real-world applications span across domains, from customer segmentation and fraud detection to image compression and network security. Practical considerations, including data preprocessing, hyperparameter tuning, and scalability, are addressed to ensure robust and efficient model deployment. By mastering these methodologies, practitioners can leverage unsupervised learning to extract actionable insights and drive data-driven decision-making in complex datasets.



## 14 Module Questions

### Introduction to Unsupervised Learning

- **Question:** What is the primary focus of unsupervised learning?
- **Question:** What are some real-world applications of unsupervised learning?
- **Question:** Why is preprocessing important in unsupervised learning workflows?

### Clustering Techniques

- **Question:** What are the key objectives of clustering in unsupervised learning?
- **Question:** How does DBSCAN differ from K-Means?
- **Question:** What are common metrics for evaluating clustering performance?

### Dimensionality Reduction

- **Question:** Why is dimensionality reduction important in unsupervised learning?
- **Question:** What is the difference between PCA and t-SNE?
- **Question:** How is the explained variance used in dimensionality reduction?

### Anomaly Detection

- **Question:** What are the primary techniques for anomaly detection in unsupervised learning?
- **Question:** How is reconstruction error used in anomaly detection?
- **Question:** Why are precision-recall curves important for evaluating anomaly detection models?

### Evaluation and Interpretation

- **Question:** What challenges exist in evaluating unsupervised models?
- **Question:** How does SHAP improve the interpretability of unsupervised learning models?
- **Question:** Why is visualization critical in unsupervised learning?

## 15 Recommended Kaggle Datasets for Unsupervised Learning

Below are Kaggle datasets associated with various unsupervised learning techniques, along with a brief explanation of their use.

### 15.1 Clustering Methods

**Dataset:** Mall Customers Dataset (<https://www.kaggle.com/vjchoudhary7/customer-segmentation-tut>)

**Use:** This dataset contains demographic and spending data of mall customers. It is ideal for clustering tasks to segment customers into groups based on shopping behavior.

**Dataset:** Wholesale Customers Dataset (<https://www.kaggle.com/uciml/wholesale-customers-data-se>)

**Use:** This dataset provides wholesale customer sales data and is suitable for clustering techniques to identify customer segments and optimize marketing strategies.

### 15.2 Dimensionality Reduction

**Dataset:** World Happiness Report Dataset (<https://www.kaggle.com/unsdsn/world-happiness>)

**Use:** This dataset includes happiness scores and related metrics. Dimensionality reduction methods like PCA or t-SNE can help visualize and analyze global happiness trends.

**Dataset:** Google Landmark Dataset (<https://www.kaggle.com/google/google-landmarks-dataset>)

**Use:** This dataset contains images of landmarks. Techniques like t-SNE or UMAP can reduce the dimensionality of image feature vectors for clustering or visualization.

### 15.3 Anomaly Detection

**Dataset:** Credit Card Fraud Detection Dataset (<https://www.kaggle.com/mlg-ulb/creditcardfraud>)

**Use:** This dataset provides anonymized credit card transaction data. Anomaly detection techniques can be used to identify fraudulent transactions.

**Dataset:** Network Intrusion Detection Dataset (NSL-KDD) (<https://www.kaggle.com/sampadab17/network-intrusion-detection>)

**Use:** This dataset contains network traffic data for intrusion detection. Anomaly detection methods can identify irregular patterns indicative of cyberattacks.

### 15.4 Interpretability in Unsupervised Learning

**Dataset:** HR Analytics: Employee Attrition Dataset (<https://www.kaggle.com/arashnic/hr-analytics-job-change-of-data-scientists>)

**Use:** This dataset includes employee demographics and performance metrics. Clustering and visualization techniques like SHAP or t-SNE can be applied to understand patterns in employee attrition.

**Dataset:** Retail Customer Churn Dataset (<https://www.kaggle.com/shubhendra1985/retail-customer-churn>)

**Use:** This dataset helps analyze customer behavior and churn patterns. Interpretability tools can provide insights into underlying factors influencing churn.

## 15.5 Deep Learning for Unsupervised Tasks

**Dataset:** Fashion MNIST (<https://www.kaggle.com/zalando-research/fashionmnist>)

**Use:** This dataset contains images of clothing items. Autoencoders and deep clustering methods can be applied for feature learning and grouping similar images.

**Dataset:** MNIST Handwritten Digits (<https://www.kaggle.com/c/digit-recognizer>)

**Use:** This dataset includes handwritten digit images. Dimensionality reduction techniques like UMAP or t-SNE can visualize the data, while autoencoders can learn compact representations.

## 16 Group Exercise

The objective of this exercise is to apply the concepts of unsupervised learning, including clustering, dimensionality reduction, anomaly detection, and interpretability, to a real-world dataset from Kaggle. The exercise involves data exploration, application of unsupervised learning techniques, and interpretation of results to reinforce practical understanding.

### Dataset

The dataset for this exercise is the **Mall Customers Dataset** from Kaggle (<https://www.kaggle.com/vjchoudhary7/customer-segmentation-tutorial-in-python>).

### Task Description

Participants will work in groups to complete the following tasks:

#### 1. Data Exploration and Preprocessing

1. Load the dataset and examine its structure, including features such as customer demographics and spending scores.
2. Handle missing values using appropriate imputation techniques.
3. Perform exploratory data analysis (EDA) to understand distributions, correlations, and patterns within the data.
4. Normalize numerical features for compatibility with clustering algorithms.

#### 2. Clustering Analysis

1. Apply **K-Means Clustering** to segment customers based on demographic and spending features.
2. Use the **Elbow Method** or **Silhouette Score** to determine the optimal number of clusters.
3. Visualize the clusters using dimensionality reduction techniques such as **Principal Component Analysis (PCA)** or **t-SNE**.
4. Interpret the characteristics of each cluster and suggest potential marketing strategies for different customer segments.

#### 3. Dimensionality Reduction

1. Apply **Principal Component Analysis (PCA)** to reduce the dimensionality of the dataset while retaining significant variance.
2. Visualize the data in the reduced-dimensional space and discuss the effectiveness of the dimensionality reduction.
3. Compare the clustering results on the original dataset versus the PCA-transformed dataset.

## 4. Anomaly Detection

1. Apply an **Isolation Forest** or **DBSCAN** to identify outliers or anomalies in customer spending behavior.
2. Visualize the anomalies and discuss their potential implications for the business.

## 5. Interpretability Techniques

1. Use **SHAP** (SHapley Additive exPlanations) or **Local Interpretable Model-agnostic Explanations (LIME)** to interpret the clustering results and feature contributions for customer segmentation.
2. Discuss how interpretability techniques enhance understanding and trust in unsupervised models.

## Deliverables

Each group must submit:

1. A Jupyter Notebook with all code, visualizations, and analysis.
2. A summary report (2-3 pages) detailing the steps taken, results, and key insights.
3. A brief presentation (5 minutes) highlighting the findings and learnings from the exercise.

## Evaluation Criteria

- **Completeness:** Were all tasks completed as described in the instructions?
- **Accuracy:** Are the results accurate and appropriately interpreted?
- **Clarity:** Is the code well-documented and easy to follow? Is the report concise and clear?
- **Creativity:** Were any additional insights or innovative approaches explored?

## Group Discussion Points

At the end of the exercise, groups will discuss:

1. Challenges faced during data preprocessing and clustering.
2. Insights gained from dimensionality reduction and anomaly detection techniques.
3. The role of interpretability techniques in understanding clustering results.

## Resources

- Kaggle Dataset: <https://www.kaggle.com/vjchoudhary7/customer-segmentation-tutorial-in-py>
- Scikit-learn Documentation: <https://scikit-learn.org>
- SHAP and LIME Documentation: <https://shap.readthedocs.io>, <https://github.com/marcotcr/lime>

- Matplotlib and Seaborn for Visualization: <https://matplotlib.org>, <https://seaborn.pydata.org>

## References

- [1] Charu C Aggarwal et al. “Automatic subspace clustering of high dimensional data for data mining applications”. In: *Proceedings of the 1998 ACM SIGMOD International Conference on Management of Data*. 1998, pp. 94–105.
- [2] Charu C. Aggarwal et al. “PROCLUS: A fast clustering algorithm for high-dimensional data”. In: *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*. 2000, pp. 61–72.
- [3] Nasir Ahmed, T Natarajan, and K R Rao. “Discrete cosine transform”. In: *IEEE Transactions on Computers* 23.1 (1974), pp. 90–93.
- [4] Artur Andrzejak, Dejan Silva, and Luciano Lopes. “Decision trees and clustering for anomaly detection in performance data”. In: *International Symposium on Integrated Network Management (IM)*. 2004, pp. 240–246.
- [5] David Arthur and Sergei Vassilvitskii. “k-means++: The advantages of careful seeding”. In: (2007), pp. 1027–1035.
- [6] David Arthur and Sergei Vassilvitskii. “k-means++: the advantages of careful seeding”. In: *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms* (2007), pp. 1027–1035.
- [7] Melissa J Azur et al. “Multiple Imputation by Chained Equations: What is it and how does it work?” In: *International Journal of Methods in Psychiatric Research* 20.1 (2011), pp. 40–49.
- [8] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006. URL: <https://www.microsoft.com/en-us/research/uploads/prod/2006/01/Bishop-Pattern-Recognition-and-Machine-Learning-2006.pdf>.
- [9] Tadeusz Calinski and Jerzy Harabasz. “A dendrite method for cluster analysis”. In: *Communications in Statistics* 3.1 (1974), pp. 1–27.
- [10] Ricardo JG Campello, Davoud Moulavi, and Jörg Sander. “Density-based clustering based on hierarchical density estimates”. In: *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. 2013, pp. 160–172.
- [11] Jeffrey Dean and Sanjay Ghemawat. “MapReduce: Simplified Data Processing on Large Clusters”. In: *Communications of the ACM*. Vol. 51. 1. ACM, 2008, pp. 107–113.
- [12] Martin Ester et al. “A density-based algorithm for discovering clusters in large spatial databases with noise”. In: *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*. 1996, pp. 226–231.
- [13] Hugging Face. *Hugging Face Contribution Guide*. Accessed: 2024-10-17. 2024. URL: <https://huggingface.co/docs/transformers/main/en/contributing>.
- [14] Hugging Face. *Hugging Face Developer Guide*. Accessed: 2024-10-17. 2024. URL: <https://huggingface.co/docs/transformers/main/en/developers>.
- [15] Aurélien Géron. *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O’Reilly Media, 2019.
- [16] Gene H Golub and Charles F Van Loan. *Matrix Computations*. Johns Hopkins University Press, 2013.
- [17] Rafael C Gonzalez and Richard E Woods. *Digital Image Processing*. Prentice Hall, 2008.
- [18] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. URL: <https://doi.org/10.1117/12.2664346>.
- [19] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. 2nd. Springer, 2009.
- [20] Joshua Hochenbaum, Oscar Vallis, and Ashish Kejariwal. “Automatic anomaly detection in the cloud via statistical learning”. In: *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)* (2017), pp. 485–493.
- [21] A. K. Jain, M. N. Murty, and P. J. Flynn. “Data clustering: A review”. In: *ACM Computing Surveys* 31 (2 1999), pp. 264–323.
- [22] Anil K Jain. “Data clustering: 50 years beyond K-Means”. In: *Pattern Recognition Letters* 31.8 (2010), pp. 651–666.

- [23] G. James et al. *An Introduction to Statistical Learning*. Springer, 2013.
- [24] Ian T Jolliffe. *Principal Component Analysis*. Springer, 2002.
- [25] Diederik P. Kingma and Max Welling. “Auto-encoding variational Bayes”. In: *International Conference on Learning Representations (ICLR)*. 2013.
- [26] Roderick JA Little and Donald B Rubin. *Statistical Analysis with Missing Data*. John Wiley & Sons, 2002.
- [27] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. “Isolation forest”. In: *Eighth IEEE International Conference on Data Mining*. IEEE. 2008, pp. 413–422.
- [28] Stuart P Lloyd. “Least squares quantization in PCM”. In: *IEEE Transactions on Information Theory* 28.2 (1982), pp. 129–137.
- [29] Scott M Lundberg and Su-In Lee. “A Unified Approach to Interpreting Model Predictions”. In: (2017), pp. 4765–4774.
- [30] Scott M Lundberg and Su-In Lee. “A unified approach to interpreting model predictions”. In: *Advances in Neural Information Processing Systems* 30 (2017), pp. 4765–4774.
- [31] Laurens van der Maaten and Geoffrey Hinton. “Visualizing data using t-SNE”. In: *Journal of Machine Learning Research* 9 (2008), pp. 2579–2605.
- [32] J MacQueen. “Some methods for classification and analysis of multivariate observations”. In: *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*. University of California Press. 1967, pp. 281–297.
- [33] James MacQueen. “Some methods for classification and analysis of multivariate observations”. In: *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*. University of California Press. 1967, pp. 281–297.
- [34] Pankaj Malhotra et al. “Long short term memory networks for anomaly detection in time series”. In: *Proceedings of the European Symposium on Artificial Neural Networks (ESANN)* (2015).
- [35] Leland McInnes, John Healy, and James Melville. “UMAP: Uniform manifold approximation and projection for dimension reduction”. In: *arXiv preprint arXiv:1802.03426* (2018).
- [36] Douglas C Montgomery, Elizabeth A Peck, and G Geoffrey Vining. *Introduction to Linear Regression Analysis*. 5th ed. John Wiley & Sons, 2012.
- [37] Fionn Murtagh and Pedro Contreras. “Algorithms for hierarchical clustering: an overview”. In: *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 2.1 (2012), pp. 86–97.
- [38] Haesun Park and Chiho Jun. “A simple and fast algorithm for K-medoids clustering”. In: *Expert Systems with Applications* 36.2 (2009), pp. 3336–3341.
- [39] Karl Pearson. “On lines and planes of closest fit to systems of points in space”. In: *Philosophical Magazine* 2.11 (1901), pp. 559–572.
- [40] K R Rao. *Discrete Cosine Transform: Algorithms, Advantages, Applications*. Academic Press, 1990.
- [41] Peter J Rousseeuw. “Silhouettes: a graphical aid to the interpretation and validation of cluster analysis”. In: *Journal of Computational and Applied Mathematics* 20 (1987), pp. 53–65.
- [42] Bernhard Schölkopf et al. “Estimating the Support of a High-Dimensional Distribution”. In: *Neural Computation* 13.7 (2001), pp. 1443–1471.
- [43] Erich Schubert et al. “DBSCAN revisited, revisited: Why and how you should (still) use DBSCAN”. In: *ACM Transactions on Database Systems (TODS)* 42.3 (2017), pp. 1–21.
- [44] Jianbo Shi and Jitendra Malik. “Normalized cuts and image segmentation”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2000, pp. 731–737.
- [45] Junyuan Xie, Ross Girshick, and Ali Farhadi. “Unsupervised deep embedding for clustering analysis”. In: *International Conference on Machine Learning*. 2016, pp. 478–487.