

MACHINE LEARNING ALGORITHMS

Page | 1

Below are comprehensive notes on several popular machine learning algorithms for clustering, density estimation, dimensionality reduction, and visualization. These notes cover the concepts, methodologies, strengths, limitations, and examples for:

- **K-means Clustering and Hierarchical Clustering**
- **Gaussian Mixture Models (GMMs)**
- **Principal Component Analysis (PCA)**
- **t-SNE (t-Distributed Stochastic Neighbor Embedding)**

1. K-means Clustering and Hierarchical Clustering

K-means Clustering

- **Definition:** A popular partitioning method that divides a dataset into k distinct, non-overlapping clusters based on feature similarity.
- **How It Works:**
 1. **Initialization:** Randomly selects k centroids (cluster centers) from the data.
 2. **Assignment Step:** Each data point is assigned to the nearest centroid, typically using Euclidean distance.
 3. **Update Step:** New centroids are computed as the mean of all points assigned to each cluster.
 4. **Iteration:** Steps 2 and 3 repeat until convergence (i.e., cluster assignments no longer change significantly or a maximum number of iterations is reached).
- **Advantages:**
 - Simple and efficient with computational complexity roughly $O(n \times k \times I)$ (where n is the number of data points, k is the number of clusters, and I is the number of iterations).



- Scalable to large datasets.
- **Limitations:**
 - Requires the number of clusters k to be predefined.
 - Sensitive to initial centroid placement, which may result in different clusterings on different runs.
 - Assumes spherical clusters of similar size; struggles with clusters of varying shapes and densities.
- **Use Cases:**
 - Market segmentation, image compression, anomaly detection, and organizing large datasets into meaningful groups.

Hierarchical Clustering

- **Definition:** A clustering approach that builds a hierarchy of clusters, often represented in a dendrogram, without requiring a predefined number of clusters.
- **Types:**
 - **Agglomerative (Bottom-Up):**
 - Starts with each data point as its own cluster and iteratively merges the closest pairs of clusters.
 - **Divisive (Top-Down):**
 - Starts with all data points in a single cluster and recursively splits the cluster into smaller groups.
- **How It Works (Agglomerative Example):**
 1. **Initialization:** Each data point is treated as an individual cluster.
 2. **Distance Calculation:** Compute a distance (or similarity) matrix between clusters.
 3. **Merge Clusters:** The two closest clusters are merged into one.
 4. **Update Distance Matrix:** Recalculate distances between the new cluster and existing clusters using linkage criteria (single, complete, or average linkage, etc.).

5. **Iteration:** Repeat merging until a stopping criterion (e.g., desired number of clusters or a distance threshold) is met.

- **Advantages:**

- Does not require a preset number of clusters.
- The dendrogram provides a visual summary of clustering structure and can reveal cluster hierarchies.

- **Limitations:**

- Computationally intensive for very large datasets.
- Sensitive to noise and outliers.
- Choice of linkage criteria impacts the resulting clusters.

- **Use Cases:**

- Creating taxonomies in biology, organizing documents, and exploratory data analysis when the number of clusters is not known a priori.

2. Gaussian Mixture Models (GMMs)

- **Definition:** GMMs are probabilistic models that assume the data is generated from a mixture of several Gaussian distributions (components), each with its own mean and covariance.

- **How It Works:**

1. **Assumption:** Data is modeled as a weighted sum of k Gaussian distributions.

2. **Expectation-Maximization (EM) Algorithm:**

- **E-step (Expectation):** Compute the probability that each data point belongs to each Gaussian component (posterior probabilities).
- **M-step (Maximization):** Update the parameters (means, covariances, and mixture weights) to maximize the likelihood of the observed data.



3. **Iteration:** Alternate between the E-step and M-step until convergence (i.e., changes in the likelihood, parameters, or posteriors are minimal).

- **Advantages:**

- More flexible than K-means; can model clusters with different shapes and sizes.
- Provides a soft clustering assignment where each point is given a probability of belonging to each cluster.

- **Limitations:**

- May converge to local optima; initialization and choice of k are important.
- Computation can be expensive, especially for high-dimensional data.
- Assumes the underlying distributions are Gaussian, which may not fit all datasets.

- **Use Cases:**

- Density estimation, clustering in complex datasets, anomaly detection, image segmentation, and modeling in finance or bioinformatics.

3. Principal Component Analysis (PCA)

- **Definition:** A widely used dimensionality reduction technique that transforms a high-dimensional dataset into a lower-dimensional space while preserving as much variance as possible.
- **How It Works:**
 1. **Standardization:** Subtract the mean from each feature, often scaling to unit variance.
 2. **Covariance Matrix Computation:** Calculate the covariance matrix of the standardized data.
 3. **Eigen Decomposition:** Compute eigenvalues and eigenvectors of the covariance matrix.
 4. **Principal Components:** The eigenvectors corresponding to the largest eigenvalues represent the directions of maximum variance.

5. **Projection:** Project the original data onto the space defined by these principal components to reduce dimensionality.

- **Advantages:**

- Reduces the dimensionality of data, making it easier to visualize and process.
- Helps remove noise and redundancy from features.
- Improves computational efficiency in subsequent analysis.

- **Limitations:**

- Assumes linear relationships between features.
- The new principal components can be hard to interpret as they are linear combinations of original features.
- Sensitive to outliers if not preprocessed.

- **Use Cases:**

- Data compression, visualization of high-dimensional data, feature extraction for supervised learning, and noise reduction in datasets.

4. t-SNE (t-Distributed Stochastic Neighbor Embedding)

- **Definition:** A nonlinear dimensionality reduction technique primarily used for visualizing high-dimensional datasets in 2 or 3 dimensions. It emphasizes maintaining the local structure of the data.

- **How It Works:**

1. **Pairwise Similarities in High Dimensions:**

- Computes pairwise similarities between data points using conditional probabilities. Similar points obtain higher probabilities.

2. **Mapping to Low-Dimensional Space:**

- The algorithm defines a similar probability distribution for points in the low-dimensional space.

3. **Kullback-Leibler Divergence:**



- Minimizes the divergence between the high-dimensional and low-dimensional probability distributions to find an optimal embedding.

4. Optimization:

- Iterative gradient descent is used to adjust the low-dimensional representation to best reflect the original similarities.

- **Advantages:**

- Particularly effective for visualizing clusters and complex patterns in the data.
- Captures local relationships, making it great for seeing subclusters and groupings.

- **Limitations:**

- Computationally intensive on large datasets.
- The results can differ significantly with different initializations or hyperparameter choices (e.g., perplexity).
- Primarily a visualization tool rather than a robust analytical method for further processing.

- **Use Cases:**

- Visualizing high-dimensional data in fields such as genetics, image processing, and natural language processing, allowing human interpretation of complex clustering or manifold structure.

Summary

- **Clustering (K-means & Hierarchical):**

- **K-means:** Fast, partitioning method; best for spherical clusters when k is known.
- **Hierarchical:** Builds a tree of clusters; useful for exploratory analysis without a prior k .

- **Density Estimation (GMMs):**

- Models data as a mixture of Gaussian distributions, providing a probabilistic, flexible clustering solution.
- **Dimensionality Reduction (PCA & t-SNE):**
 - **PCA:** Linear method that reduces dimensions by capturing maximum variance; good for preprocessing.
 - **t-SNE:** Nonlinear method that preserves local structure; ideal for visualizing high-dimensional data in a human-interpretable form.

These methods, individually or in combination, are widely used in practice for exploratory data analysis, pattern recognition, and preparing data for downstream tasks in machine learning and artificial intelligence projects.