

Unsupervised Learning: Clustering and Autoencoders

Unsupervised learning is a fundamental branch of machine learning where the goal is to uncover hidden patterns or structures in data without labeled outputs. One of the most prominent techniques in unsupervised learning is clustering, which involves grouping similar data points together.

Clustering aims to form groups (clusters) such that data points within the same cluster are more similar to each other than to those in other clusters. This has wide applications in document organization, recommendation systems, genomics, finance, and data summarization.

There are three main clustering techniques: partitional, hierarchical, and density-based. Partitional clustering, such as k-means, divides data into non-overlapping subsets. Hierarchical clustering builds a tree of clusters either by merging (agglomerative) or splitting (divisive). Density-based methods like DBSCAN identify clusters as dense regions separated by sparse areas, handling noise and irregular shapes effectively.

Clusters can be categorized by structure: well-separated clusters have clear boundaries; prototype-based clusters revolve around a central point (centroid or medoid); contiguity-based clusters rely on local proximity; and density-based clusters are defined by high-density regions.

The k-means algorithm is a popular partitional method. It starts by selecting K initial centroids, assigns each point to the nearest centroid, and updates centroids iteratively. The objective function minimized is the Sum of Squared Errors (SSE), which measures the compactness of clusters. However, k-means is sensitive to the initial choice of centroids. Poor initialization can lead to suboptimal clustering. K-means++ addresses this by selecting initial centroids with a probability proportional to their distance from existing centroids, improving convergence and accuracy.

Despite its simplicity, k-means has limitations. It struggles with clusters of varying sizes, densities, non-convex shapes, and is sensitive to outliers. These issues can lead to incorrect clustering results.

Hierarchical clustering offers an alternative by constructing a dendrogram - a tree-like diagram

showing nested clusters. Agglomerative clustering starts with individual points and merges them, while divisive clustering begins with one cluster and splits it. The choice of inter-cluster similarity measure affects the clustering outcome. Single link (minimum distance) can handle non-elliptical shapes but is sensitive to noise. Complete link (maximum distance) is less sensitive to outliers but favors globular clusters. Average linkage and Ward's method (which minimizes variance) offer balanced approaches.

Hierarchical clustering does not require specifying the number of clusters in advance and can reveal meaningful taxonomies. However, it is computationally intensive and sensitive to the choice of distance metric.

Autoencoders represent another form of unsupervised learning. These neural networks learn to compress data into a lower-dimensional representation (encoding) and reconstruct it (decoding). They are useful for dimensionality reduction, anomaly detection, and feature learning. Deep autoencoders with multiple hidden layers can capture complex structures in data, often outperforming traditional methods like PCA.

In conclusion, unsupervised learning techniques like clustering and autoencoders are powerful tools for discovering structure in unlabeled data. Understanding their mechanisms, strengths, and limitations is essential for effective application in real-world problems.