

Okay, here is a structured study note based on the provided PDF document, including both English and Chinese explanations:

Unsupervised Learning: Structure Analysis

无监督学习：结构分析

I. Introduction

一、简介

- **Definition of Unsupervised Learning (无监督学习的定义):**

- Unsupervised learning does not elevate any attribute to the category of label: all the attributes are treated equally. The essence of unsupervised learning is encapsulated in the simple question **where is my data?** (无监督学习不会将任何属性提升到标签的类别：所有属性都被平等对待。无监督学习的本质可以用一个简单的问题来概括：我的数据在哪里？)
- The attribute space is infinite and mostly empty, and the answer to this question will be a **model** that will allow us to identify the regions where we could expect to find samples. (属性空间是无限的，而且大部分是空的，这个问题的答案将是一个模型，它将允许我们识别可以预期找到样本的区域。)

- **Supervised vs. Unsupervised Learning (监督学习与无监督学习):**

- **Supervised learning (监督学习):** Designates one attribute as a **label (标签)** and treats the rest as **predictors (预测因子)**. It builds a model that estimates the value of the label based on the value of the predictors. The model mirrors the underlying structure of the dataset, and the notion of quality is defined as a function of the discrepancy between the true and estimated values of the label. (指定一个属性作为标签，并将其余属性视为预测变量。它建立一个基于预测变量的值来估计标签值的模型。该模型反映了数据集的基本结构，质量的概念被定义为标签的真实值和估计值之间差异的函数。)
- **Unsupervised learning (无监督学习):** Does not have a designated label. It aims to find patterns and structures within the data without prior knowledge of what those patterns might be. (没有指定的标签。它的目标是在没有事先知道这些模式可能是什么的情况下，在数据中找到模式和结构。)

II. Main Approaches of Unsupervised Learning

二、无监督学习的主要方法

- **Density estimation (密度估计):** Creates models that allow us to quantify the probability of finding a sample within a region of the attribute space (**probability density 概率密度**). (创建模型，使我们能够量化在属性空间的某个区域内找到样本的概率（概率密度）。)
- **Structure analysis (结构分析):** Creates models that identify regions within the attribute space (**cluster analysis 聚类分析**) or directions (**component analysis 成分分析**) with a high density of samples. (创建模型，识别属性空间内样本密度高的区域（聚类分析）或方向（成分分析）。)

III. Applications of Unsupervised Learning

三、无监督学习的应用

- **Summarizing data (数据总结):** Provide summaries of a population in the form of **prototypes (原型)** samples. (以原型样本的形式提供总体的摘要。)
- **Discovery (探索发现):** Discover hidden **structure (结构)**, generate new knowledge (e.g., genetics-based migration studies), or change the way we represent our data (e.g., compression). (发现隐藏的结构，产生新的知识（如基于遗传学的迁徙研究），或改变我们表示数据的方式（如压缩）。)
- **Quantitative applications (定量应用):** Build **class densities (类别密度)** that describe the probability of finding a sample from a given class in a region; identify **anomalies (异常值)**. (建立描述在某一区域找到某一类样本的概率的类别密度；识别异常值。)

IV. Cluster Analysis

四、聚类分析

- **Definition (定义):** Clustering is a family of unsupervised learning algorithms that describe the structure of a dataset as groups, or clusters, of **similar (相似的)** samples. (聚类是一系列无监督学习算法，将数据集的结构描述为一组组相似样本或聚类。)
- **Similarity (相似度):** A notion of **similarity (相似度)** is needed to partition a dataset into clusters. (将数据集划分为聚类需要一个相似性的概念。)
- **Proximity as Similarity (基于距离的相似度):** Clusters can be defined as groups of samples that are **close (相近)** to one another. **Proximity (距离)** is used as the notion of similarity. (聚类可以定义为彼此接近的样本组。距离被用作相似性的概念。)
- **Squared Distance (平方距离):** Given two samples x_i and x_j consisting of P attributes, $x_{i,1}, \dots, x_{i,P}$ and $x_{j,1}, \dots, x_{j,P}$, the **squared distance (平方距离)** $d_{i,j}$ is defined as: $d_{i,j} = (x_{i,1} - x_{j,1})^2 + \dots + (x_{i,P} - x_{j,P})^2$. (给定两个由 P 个属性组成的样本 x_i 和 x_j , $x_{i,1}, \dots, x_{i,P}$ 和 $x_{j,1}, \dots, x_{j,P}$, 平方距离 $d_{i,j}$ 定义为: $d_{i,j} = (x_{i,1} - x_{j,1})^2 + \dots + (x_{i,P} - x_{j,P})^2$ 。)
- **Quality Metric (质量指标):**
 - **Intra-cluster sample scatter (类内样本离散度):** The sum of the square distances between samples in the same cluster. (同一聚类中样本之间的平方距离之和。)
 - **Inter-cluster sample scatter (类间样本离散度):** The sum of the distances between samples in different clusters. (不同聚类中样本之间的距离之和。)
 - **Goal (目标):** The best clustering arrangement has the **lowest intra-cluster (最小的类内)** sample scatter and **highest inter-cluster (最大的类间)** sample scatter. (最佳聚类排列具有最小的类内样本离散度和最大的类间样本离散度。)

V. K-means Clustering

五、 K-均值聚类

- **Concept (概念):** Partitions a dataset into K clusters represented by their **mean (均值)**. (将数据集划分为由其均值表示的 K 个聚类。)
- **Prototypes (原型):** Obtained as the **centre (or mean) (中心 (或均值))** of each cluster. (作为每个聚类的中心 (或均值) 获得。)
- **Process (过程):**
 1. Prototypes are obtained as the centre (or mean) of each cluster.
 2. Samples are re-assigned to the cluster with the closest prototype. (样本被重新分配到具有最接近原型的聚类。)
- **Local Optimum (局部最优):** The final solution is a **local optimum (局部最优)**, not necessarily the global one. (最终解决方案是局部最优，而不一定是全局最优。)
- **Determining the Number of Clusters (确定聚类数量):**
 - **Elbow method (肘部法则):** The true number of clusters can be identified by observing the value of K beyond which the improvement slows down. (可以通过观察 K 值来确定真实的聚类数量，超过这个值，改进的速度就会减慢。)

VI. Density-based Clustering (DBSCAN)

六、 基于密度的聚类 (DBSCAN)

- **Concept (概念):** Clusters are defined as groups of **connected (连接的)** samples, suitable for **non-convex (非凸)** clusters. (聚类被定义为连接样本的组，适用于非凸聚类。)
- **Definitions (定义):**
 - **Radius (半径) r**
 - **Threshold (阈值) t**
 - **Core (核心点):** Its density is equal or higher than the threshold t . (其密度等于或高于阈值 t 。)
 - **Border (边界点):** Its density is lower than the threshold t , but contains a core sample within its neighborhood. (其密度低于阈值 t ，但在其邻域内包含一个核心样本。)
 - **Outlier (离群点):** Any other sample. (任何其他样本。)
- **Process (过程):**
 1. Identify core, border and outlier samples.
 2. Pair of core samples that are within each other's neighborhood are connected. Connected core samples form the **backbone (主干)** of a cluster. (彼此位于邻域内的核心样本对被连接起来。连接的核心样本形成一个聚类的主干。)
 3. Border samples are assigned to the cluster that has more core samples in the neighborhood of the border sample.
 4. Outlier samples are not assigned to any cluster.

VII. Hierarchical Clustering

七、 层次聚类

- **Concept (概念):** Progressively building clustering arrangements at **different levels (不同层次)**. (在不同层次逐步建立聚类排列。)
- **Dendrogram (树状图):** The representation of the relationship between clusters at different levels is called a **dendrogram (树状图)**. (在不同层次上表示聚类之间关系的图被称为树状图。)
- **Strategies (策略):**
 - **Divisive (分裂式)** or top-down: Splits clusters starting from the top of the dendrogram and stops at the bottom level. (从树状图的顶部分割聚类，并在底层停止。)
 - **Agglomerative (凝聚式)** or bottom-up: Merges two clusters, starting from the bottom until we reach the top level. (从底部开始合并两个聚类，直到到达顶层。)
- **Merging/Splitting Criteria (合并/分裂标准):**
 - **Single linkage (单一连接):** Uses the distance between the two closest samples from two clusters. (使用两个聚类中最近的两个样本之间的距离。)
 - **Complete linkage (完全连接):** Uses the distance between the two further samples from each pair of clusters. (使用每对聚类中两个较远样本之间的距离。)
 - **Group average (组平均):** Uses the average distance between samples in two clusters. (使用两个聚类中样本之间的平均距离。)

VIII. Summary

八、 总结

- Unsupervised learning answers the question "where is my data?" in the attribute space. (无监督学习回答了“我的数据在哪里”在属性空间中的问题。)
- The answer is a mathematical/computer model that can indicate sample locations (clustering) or probability of finding a sample in a region (density estimation). (答案是一个数学/计算机模型，可以指示样本位置 (聚类) 或在某个区域找到样本的概率 (密度估计)。)
- The absence of a target label makes defining a quality metric less obvious. (由于没有目标标签，定义质量指标就不那么明显了。)
- **K-means** is prototype-based and produces spherical clusters. (K-均值是基于原型的，并产生球形聚类。)

- **DBSCAN** is density-based, suitable for non-convex scenarios, and does not require pre-defining the number of clusters. (DBSCAN 是基于密度的，适用于非凸情况，不需要预先定义聚类的数量。)
- **Hierarchical clustering** explores structure at multiple levels. (层次聚类在多个层次上探索结构。)
- Clustering quality should be evaluated based on the specific goals of the application. (应根据应用的具体目标评估聚类质量。)

This comprehensive summary provides a clear and organized understanding of the concepts presented in the PDF document. Remember to refer back to the original document for visual aids and more detailed examples. Let me know if you have any other questions.