



Unsupervised Learning

Madson L. D. Dias

Huawei / IFCE

March 26, 2021

Agenda

① Unsupervised Learning

- Introduction

- Applications

- Supervised Learning vs. Unsupervised Learning

② Clustering

- Formal Definition

- Types of clustering

- Clustering using k -means

③ Dimensionality reduction

- Principal Component Analysis

Agenda

① Unsupervised Learning

- Introduction

- Applications

- Supervised Learning vs. Unsupervised Learning

② Clustering

- Formal Definition

- Types of clustering

- Clustering using k -means

③ Dimensionality reduction

- Principal Component Analysis

Agenda

① Unsupervised Learning

- Introduction

- Applications

- Supervised Learning vs. Unsupervised Learning

② Clustering

- Formal Definition

- Types of clustering

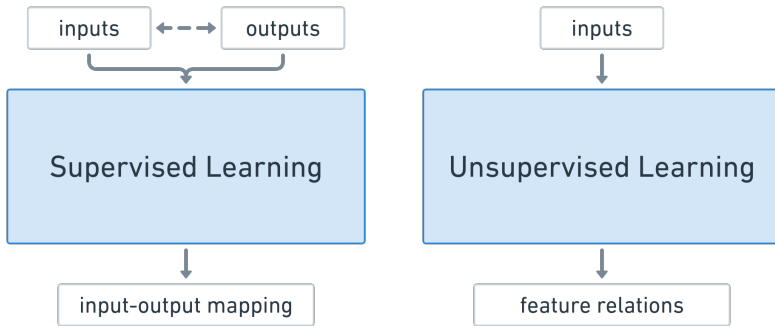
- Clustering using k -means

③ Dimensionality reduction

- Principal Component Analysis

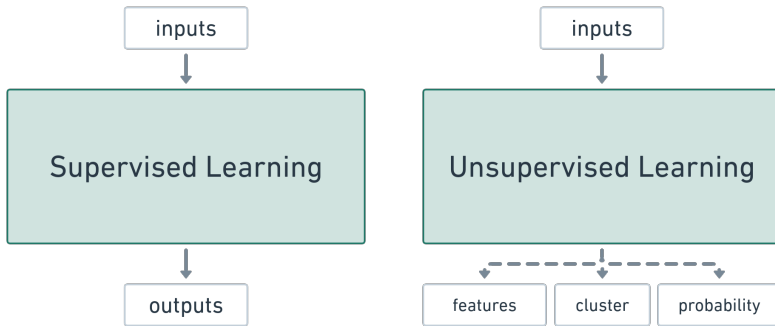
Unsupervised Learning

- A type of machine learning algorithm used to draw inferences from data sets consisting of **input data without labels**
- Task of inferring a function to describe **hidden structure** from unlabeled data.



Unsupervised Learning

- A type of machine learning algorithm used to draw inferences from data sets consisting of **input data without labels**
- Task of inferring a function to describe **hidden structure** from unlabeled data.



Agenda

① Unsupervised Learning

Introduction

Applications

Supervised Learning vs. Unsupervised Learning

② Clustering

Formal Definition

Types of clustering

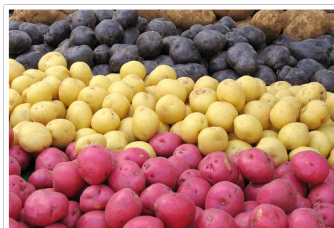
Clustering using k -means

③ Dimensionality reduction

Principal Component Analysis

Applications of Unsupervised Learning

Clustering



- Process of **grouping data into different clusters** or groups.
- Same group → similar and different groups → dissimilar
- Used for analyzing and grouping data without labels

Applications of Unsupervised Learning

Dimensionality reduction



- Process of **reducing** the number of **variables**
- Simplify the data **without losing** too much information
- This method is also called feature extraction

Applications of Unsupervised Learning

Anomaly detection



- Identification of **rare observations**, which brings suspicions
- The model is **trained** with a **lot of normal** instances
- Also called as Novelty Detection and Outlier Detection

Agenda

① Unsupervised Learning

Introduction

Applications

Supervised Learning vs. Unsupervised Learning

② Clustering

Formal Definition

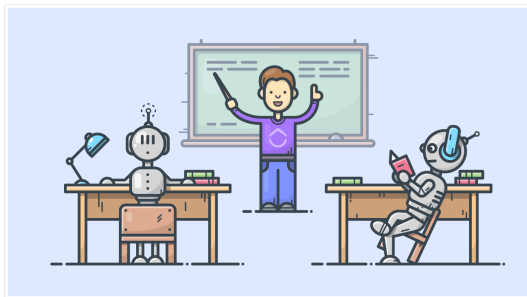
Types of clustering

Clustering using k -means

③ Dimensionality reduction

Principal Component Analysis

Supervised Learning vs. Unsupervised Learning



Parameter	Supervised Learning	Unsupervised Learning
Dataset	Labelled Dataset	Unlabelled Dataset
Learning	Guided learning	Guided by some metric
Complexity	Simpler method	Computationally complex
Accuracy	More Accurate	Less Accurate

Agenda

① Unsupervised Learning

Introduction

Applications

Supervised Learning vs. Unsupervised Learning

② Clustering

Formal Definition

Types of clustering

Clustering using k -means

③ Dimensionality reduction

Principal Component Analysis

Agenda

① Unsupervised Learning

Introduction

Applications

Supervised Learning vs. Unsupervised Learning

② Clustering

Formal Definition

Types of clustering

Clustering using k -means

③ Dimensionality reduction

Principal Component Analysis

Formal definition of Clustering

- Given a dataset $\mathcal{D} = \{\mathbf{x}_n \in \mathbb{R}^D\}_{n=1}^N$, where \mathbf{x} is a vector with D dimensions and N is the number of elements of the data set¹. The task of clustering is to find a function $g(\mathbf{x})$ in that $g : \mathbb{R}^D \rightarrow \mathbb{N}^K$, where $\mathcal{C} \in \mathbb{N}^K$ is a set of sets $\mathcal{C} = \{\mathcal{C}_k\}_{k=1}^{K \leq N}$ containing clusters $\mathcal{C}_k = \{\mathbf{x}_n : g(\mathbf{x}_n) = k\}$ such that

$$\mathcal{C}_k \neq \emptyset, \forall k \quad (1)$$

$$\bigcup_{k=1}^K \mathcal{C}_k = \mathcal{D} \quad (2)$$

$$\mathcal{C}_k \cap \mathcal{C}_j = \emptyset, \forall k, j \text{ where } k \neq j \quad (3)$$

¹The \mathbb{R}^D space is called “data space” or “input space”.

Agenda

① Unsupervised Learning

Introduction

Applications

Supervised Learning vs. Unsupervised Learning

② Clustering

Formal Definition

Types of clustering

Clustering using k -means

③ Dimensionality reduction

Principal Component Analysis

Types of clustering algorithms



- Clustering itself can be categorized into two types
 - Hard Clustering
 - Soft Clustering
- We can also categorize clustering methods by their methodology and learning algorithms

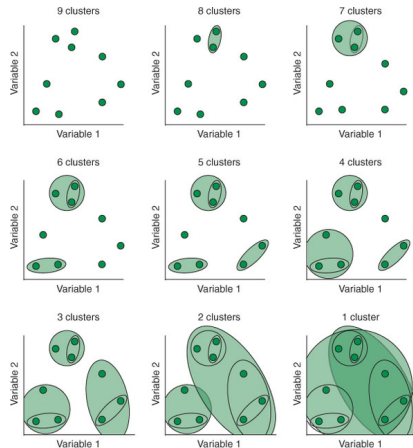
Density-based clustering



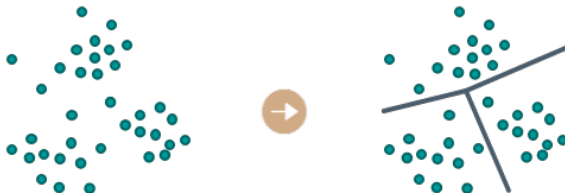
- In this method, the clusters are created based upon the **density of the data points** which are represented in the data space.
- Regions that become **dense** due to the huge number of data points residing in that region are considered as **clusters**
- Points in the **sparse** region are considered as noise or **outliers**

Hierarchical Clustering

- Divides the clusters based on the **distance metrics**
- **Agglomerative** or **Divisive**
- Create a **distance matrix** of all the existing clusters and perform the **linkage** between the clusters depending on the **criteria** of the linkage
 - Single Linkage
 - Complete Linkage
 - Average Linkage



Partitioning Clustering



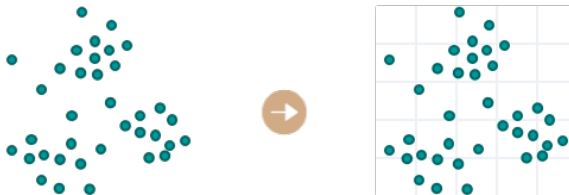
- The clusters are **partitioned** based upon the characteristics of the data points
- We need to specify the **number of clusters** to be created
- Iterative process to reassign the data points between clusters based upon the **distance**

Fuzzy clustering



- In fuzzy clustering, the assignment of the data points in any of the clusters is not decisive
- Here, one data point can belong to **more than one cluster**
- It provides the outcome as the **probability** of the data point belonging to each of the clusters.

Grid-based clustering



- The data set is represented into a **grid structure** which comprises of grids (also called cells).
- After partitioning the data sets into cells, it computes the **density** of the cells which helps in identifying the clusters.
- This makes it appropriate for dealing with **big data sets**

Agenda

① Unsupervised Learning

Introduction

Applications

Supervised Learning vs. Unsupervised Learning

② Clustering

Formal Definition

Types of clustering

Clustering using k -means

③ Dimensionality reduction

Principal Component Analysis

k -means clustering

- One of the most popular choices for analysts to create clusters.
- We need to specify the number of clusters to be created for this clustering method.

Definition

- Given a dataset $\mathcal{D} = \{\mathbf{x}_n \in \mathcal{R}^D\}_{n=1}^N$, k -means clustering aims to partition the N observations into $K(\leq N)$ sets so as to minimize the *within-cluster* sum of squares (i.e. variance). Formally, the objective is to find:

$$\mathcal{C}^* = \arg \min_{\mathcal{C}} \sum_{k=1}^K \sum_{\mathbf{x} \in \mathcal{C}_k} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \quad (4)$$

where $\boldsymbol{\mu}_k$ is the centroid of the cluster \mathcal{C}_k .

k -means clustering

Algorithm

1. Specify the number K of clusters to assign
2. Randomly initialize K centroids
3. Repeat until the centroid positions do not change
 - 3.1. Assign each point to its closest centroid
 - 3.2. Compute the new centroid (mean) of each cluster

*“Talk is cheap.
Show me the code.”*

- Linus Torvalds

<https://github.com/omadson/vds>

Agenda

① Unsupervised Learning

Introduction

Applications

Supervised Learning vs. Unsupervised Learning

② Clustering

Formal Definition

Types of clustering

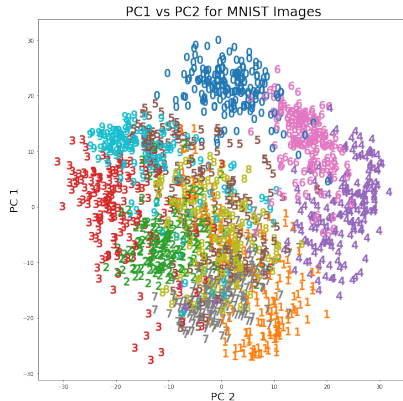
Clustering using k -means

③ Dimensionality reduction

Principal Component Analysis

Dimensionality reduction

- Transformation of data from a high-dimensional space into a **low-dimensional space**
- Ideally close to its **intrinsic dimension**.
- Methods
 - Feature selection
 - Feature projection



Agenda

① Unsupervised Learning

Introduction

Applications

Supervised Learning vs. Unsupervised Learning

② Clustering

Formal Definition

Types of clustering

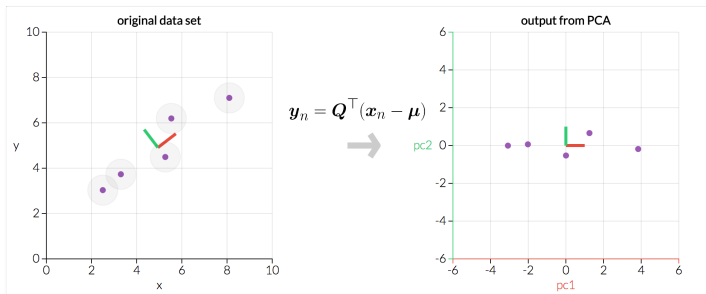
Clustering using k -means

③ Dimensionality reduction

Principal Component Analysis

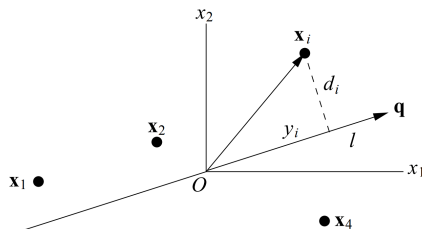
Principal Component Analysis

- PCA is often used to **reduce the dimensionality** of large data sets, by transforming a large set of variables into a smaller one that still contains most of the information in the large set.



Principal Component Analysis

- Consider a set of points $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$ in a D -dimensional space, such that their mean $\boldsymbol{\mu} = \mathbf{0}$, i.e., centroid is at the origin.



- The PCA wants to find a line l through the origin that maximizes the projections x'_n of the points \mathbf{x}_n on l .
- Let \mathbf{q} denote the unit vector along line l .

Principal Component Analysis

- The projection x'_n of \mathbf{x}_n on l is $x'_n = \mathbf{x}_n^\top \mathbf{q}$.
- The mean squared projection is the variance V over all points

$$V = \frac{1}{N} \sum_{n=1}^N x_n'^2 \quad (5)$$

$$= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n^\top \mathbf{q})^2 \quad (6)$$

$$= \frac{1}{N} \sum_{n=1}^N (\mathbf{q}^\top \mathbf{x}_n)(\mathbf{x}_n^\top \mathbf{q}) \quad (7)$$

$$= \mathbf{q}^\top \left[\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right] \mathbf{q} \quad (8)$$

Principal Component Analysis

- The middle factor is the covariance matrix \mathbf{C} of the data points

$$\mathbf{C} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \quad (9)$$

- We want to find a unit vector \mathbf{q} that maximizes the variance V

$$\mathbf{q}^* = \arg \max_{\mathbf{q}} \mathbf{q}^\top \mathbf{C} \mathbf{q}, \quad (10)$$

- Put $\|\mathbf{q}\| = 1$ constraint to avoid overflow
- The constraint optimization problem

$$\text{maximize } V = \mathbf{q}^\top \mathbf{C} \mathbf{q} \quad \text{subject to} \quad \|\mathbf{q}\| = 1. \quad (11)$$

Principal Component Analysis

Lagrange multiplier method

- Consider this problem

$$\text{maximize } f(\mathbf{x}) \text{ subject to } g(\mathbf{x}) = c. \quad (12)$$

- Lagrange multiplier method introduces a Lagrange multiplier λ to combine $f(\mathbf{x})$ and $g(\mathbf{x})$ as

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda(g(\mathbf{x}) - c) \quad (13)$$

- Then, we can solve using

$$\frac{\partial L(\mathbf{x}, \lambda)}{\partial \mathbf{x}} = 0, \quad \frac{\partial L(\mathbf{x}, \lambda)}{\partial \lambda} = 0 \quad (14)$$

Principal Component Analysis

- Use Lagrange multiplier method, combine V and the constraint
- Consider this problem

$$\text{maximize } L = \mathbf{q}^\top \mathbf{C} \mathbf{q} - \lambda(\mathbf{q}^\top \mathbf{q} - 1) \quad (15)$$

- Now, we differentiate L with respect to \mathbf{q} and λ :

$$\frac{\partial L}{\partial \mathbf{q}} = 2\mathbf{q}^\top \mathbf{C} - 2\lambda \mathbf{q}^\top = 0 \quad (16)$$

$$\frac{\partial L}{\partial \lambda} = \mathbf{q}^\top \mathbf{q} - 1 = 0 \quad (17)$$

- Eq. (16) gives

$$\mathbf{q}^\top \mathbf{C} = \lambda \mathbf{q}^\top \iff \mathbf{C} \mathbf{q} = \lambda \mathbf{q} \quad (18)$$

- This is called an *eigenvector equation*

Principal Component Analysis

General PCA

- PCA transforms \mathbf{x}_n into a new vector \mathbf{y}_n through \mathbf{Q} as follows:

$$\mathbf{x}'_n = \mathbf{Q}^\top (\mathbf{x}_n - \boldsymbol{\mu}) = \sum_{m=1}^M (\mathbf{x}_n - \boldsymbol{\mu})^\top \mathbf{q}_m \mathbf{q}_m \quad (19)$$

- Each component of \mathbf{x}'_n is

$$x'_{nm} = (\mathbf{x}_n - \boldsymbol{\mu})^\top \mathbf{q}_m \quad (20)$$

- This is the projection of $\mathbf{x}_n - \boldsymbol{\mu}$ on \mathbf{q}_m .

*“Talk is cheap.
Show me the code.”*

- Linus Torvalds

<https://github.com/omadson/vds>

Bye-Bye!

Thank you!