Machine Learning Unsupervised Learning

FAST DISCOVERING THE FUTURE

(日) (문) (문) (문) (문)

- ✓ Ingredients of Machine Learning
- ✓ Classification Basics, Basic Linear Classifier
- ✓ K-Nearest Neighbours and Naive Bayes Classifier
- ✓ Linear and Quadratic Discriminant Analysis
- ✓ Support Vector Machines (SVM)
- Decision Trees
- ✓ Ensemble Methods (Bagging, Weighted Voting, Stacking)
- ✓ Regression Methods
- ✓ Evaluation and Scoring of Classifiers
- ✓ Ensemble Methods (Boosting)

- Concepts of Cluster Analysis
- Hierarchical Clustering
- K-means Clustering
- K-medoids Clustering
- DBSCAN
- Dissimilarity Measure Selection

Machine Learning Map



• Given data objects



Image: Image:

< ∃ ►

- Given data objects
- Find a grouping (clustering) such that the objects are:



- Given data objects
- Find a grouping (clustering) such that the objects are:
 - similar (related) to the objects in the same group



- Given data objects
- Find a grouping (clustering) such that the objects are:
 - similar (related) to the objects in the same group
 - dissimilar (unrelated) from objects in other groups



Intuition building

- Intuition building
- Hypothesis generation

- Intuition building
- Hypothesis generation
- Discover structures and patterns in high-dimensional data

- Intuition building
- Hypothesis generation
- Discover structures and patterns in high-dimensional data
- Summarizing / compressing large data

• Suppose that we need to put you in some groups for projects

- Suppose that we need to put you in some groups for projects
- How to define those groups?
 - Work experience
 - Age
 - Education
 - Preferences

∃ >

- Suppose that we need to put you in some groups for projects
- How to define those groups?
 - Work experience
 - Age
 - Education
 - Preferences
- Which way is correct? Depends on the goal!

• Google news uses the clustering algorithm to categories the news related to the same topic

∃ >

- Google news uses the clustering algorithm to categories the news related to the same topic
- Segmentation of the people according to the items purchased in e-commerce applications

- Google news uses the clustering algorithm to categories the news related to the same topic
- Segmentation of the people according to the items purchased in e-commerce applications
- Segmentation of the customers of the fashion company

- Google news uses the clustering algorithm to categories the news related to the same topic
- Segmentation of the people according to the items purchased in e-commerce applications
- Segmentation of the customers of the fashion company
- Take a collection of 1,000,000 different genes, and find a way to automatically group these genes into groups that are somehow similar or related by different variables

- Google news uses the clustering algorithm to categories the news related to the same topic
- Segmentation of the people according to the items purchased in e-commerce applications
- Segmentation of the customers of the fashion company
- Take a collection of 1,000,000 different genes, and find a way to automatically group these genes into groups that are somehow similar or related by different variables
- Find stocks which have common price fluctuations

- Google news uses the clustering algorithm to categories the news related to the same topic
- Segmentation of the people according to the items purchased in e-commerce applications
- Segmentation of the customers of the fashion company
- Take a collection of 1,000,000 different genes, and find a way to automatically group these genes into groups that are somehow similar or related by different variables
- Find stocks which have common price fluctuations
- Image segmentation

- Google news uses the clustering algorithm to categories the news related to the same topic
- Segmentation of the people according to the items purchased in e-commerce applications
- Segmentation of the customers of the fashion company
- Take a collection of 1,000,000 different genes, and find a way to automatically group these genes into groups that are somehow similar or related by different variables
- Find stocks which have common price fluctuations
- Image segmentation
- Search result grouping

- Google news uses the clustering algorithm to categories the news related to the same topic
- Segmentation of the people according to the items purchased in e-commerce applications
- Segmentation of the customers of the fashion company
- Take a collection of 1,000,000 different genes, and find a way to automatically group these genes into groups that are somehow similar or related by different variables
- Find stocks which have common price fluctuations
- Image segmentation
- Search result grouping
- Social network analysis

- Google news uses the clustering algorithm to categories the news related to the same topic
- Segmentation of the people according to the items purchased in e-commerce applications
- Segmentation of the customers of the fashion company
- Take a collection of 1,000,000 different genes, and find a way to automatically group these genes into groups that are somehow similar or related by different variables
- Find stocks which have common price fluctuations
- Image segmentation
- Search result grouping
- Social network analysis
- Anomaly detection

Example: Supermarket data

Pampers Customer 1 Task Cluster customers into segments by their typical shopping baskets **Customer 3 Customer 4**

Example: Supermarket data

Pampers Customer 1 Task Cluster customers into segments by their typical shopping baskets **Customer 3** Motivation C Launch a discount campaign addressing a particular segment of customers

- Fuzzy (Soft) clustering: Each object belongs to each cluster with some weight (the weight can be zero)
- Non-fuzzy (Hard) clustering each object belongs to exactly one cluster



Clustering Types

- Partitional clustering finds a fixed number of clusters
- **Hierarchical clustering** creates a series of clusterings contained in each other (nested clusters)



• Hierarchical clustering is usually visualized as a dendrogram (tree)

- Hierarchical clustering is usually visualized as a dendrogram (tree)
- Each subtree corresponds to a cluster

- Hierarchical clustering is usually visualized as a dendrogram (tree)
- Each subtree corresponds to a cluster
- Height of branching shows distance between the objects

- Hierarchical clustering is usually visualized as a dendrogram (tree)
- Each subtree corresponds to a cluster
- Height of branching shows distance between the objects
- There are two strategies: Agglomerative and Divisive

- Hierarchical clustering is usually visualized as a dendrogram (tree)
- Each subtree corresponds to a cluster
- Height of branching shows distance between the objects
- There are two strategies: Agglomerative and Divisive
 - **Agglomerative** clustering (Bottom-up): Start with a single-object clusters (singletons) and recursively merge them into larger clusters.

- Hierarchical clustering is usually visualized as a dendrogram (tree)
- Each subtree corresponds to a cluster
- Height of branching shows distance between the objects
- There are two strategies: Agglomerative and Divisive
 - **Agglomerative** clustering (Bottom-up): Start with a single-object clusters (singletons) and recursively merge them into larger clusters.
 - **Divisive** clustering (Top down): Start with a cluster containing all data points and recursively divide it into smaller clusters.

Hierarchical Clustering

• Hierarchical clustering is usually visualized as a dendrogram (tree)



Hierarchical Clustering

• Hierarchical clustering is usually visualized as a dendrogram (tree)


Algorithm for Agglomerative Hierarchical Clustering: Join the two closest objects



Join the two closest objects (for example, closest with respect to Euclidean distance)



Keep joining the closest pairs



Keep joining the closest pairs



Keep joining the closest pairs



- 4 ∃ ▶

Keep joining the closest pairs





Keep joining the closest pairs



After $10 \ {\rm steps} \ {\rm we} \ {\rm have} \ 4 \ {\rm clusters} \ {\rm left}$





After $10 \ {\rm steps} \ {\rm we} \ {\rm have} \ 4 \ {\rm clusters} \ {\rm left}$



Several ways to measure distance between clusters:

• Single linkage (a.k.a. nearest neighbour technique)

After $10 \ {\rm steps} \ {\rm we} \ {\rm have} \ 4 \ {\rm clusters} \ {\rm left}$



Several ways to measure distance between clusters:

- Single linkage (a.k.a. nearest neighbour technique)
- Complete linkage (a.k.a furthest neighbour technique)

After $10 \ {\rm steps} \ {\rm we} \ {\rm have} \ 4 \ {\rm clusters} \ {\rm left}$



Several ways to measure distance between clusters:

- Single linkage (a.k.a. nearest neighbour technique)
- Complete linkage (a.k.a furthest neighbour technique)
- Average linkage
 - -Weighted
 - -Unweighted

In this example and at this stage we have the same result as in our partitional clustering example



In this example and at this stage we have the same result as in our partitional clustering example





Image: Image:

In the final step the two remaining clusters are joined into a single cluster



 Single Linkage - distance is computed between the two MOST similar parts of clusters (two closest points).

- **Single** Linkage distance is computed between the two MOST similar parts of clusters (two closest points).
 - Suffers from *chaining* meaning that clusters can be too spread out, and not compact enough

- **Single** Linkage distance is computed between the two MOST similar parts of clusters (two closest points).
 - Suffers from *chaining* meaning that clusters can be too spread out, and not compact enough
 - Clusters can violate the "compactness" property that all observations within each cluster tend to be similar to one another.

- **Single** Linkage distance is computed between the two MOST similar parts of clusters (two closest points).
 - Suffers from *chaining* meaning that clusters can be too spread out, and not compact enough
 - Clusters can violate the "compactness" property that all observations within each cluster tend to be similar to one another.
- **Complete** Linkage distance is computed between the two LEAST similar parts of clusters (two most distant points).

- Single Linkage distance is computed between the two MOST similar parts of clusters (two closest points).
 - Suffers from *chaining* meaning that clusters can be too spread out, and not compact enough
 - Clusters can violate the "compactness" property that all observations within each cluster tend to be similar to one another.
- **Complete** Linkage distance is computed between the two LEAST similar parts of clusters (two most distant points).
 - Avoids chaining, but suffers from *crowding*, meaning that clusters are compact, but not far enough apart

- Single Linkage distance is computed between the two MOST similar parts of clusters (two closest points).
 - Suffers from *chaining* meaning that clusters can be too spread out, and not compact enough
 - Clusters can violate the "compactness" property that all observations within each cluster tend to be similar to one another.
- **Complete** Linkage distance is computed between the two LEAST similar parts of clusters (two most distant points).
 - Avoids chaining, but suffers from *crowding*, meaning that clusters are compact, but not far enough apart
 - Observations assigned to a cluster can be much closer to members of other clusters than they are to some members of their own cluster

- Single Linkage distance is computed between the two MOST similar parts of clusters (two closest points).
 - Suffers from *chaining* meaning that clusters can be too spread out, and not compact enough
 - Clusters can violate the "compactness" property that all observations within each cluster tend to be similar to one another.
- **Complete** Linkage distance is computed between the two LEAST similar parts of clusters (two most distant points).
 - Avoids chaining, but suffers from *crowding*, meaning that clusters are compact, but not far enough apart
 - Observations assigned to a cluster can be much closer to members of other clusters than they are to some members of their own cluster
- Average Linkage distance is computed between clusters' centroids. This is a balanced approach: clusters tend to be relatively compact and relatively far apart.

Pro: The clusters can be visualized and interpreted with dendrograms

- Pro: The clusters can be visualized and interpreted with dendrograms
- **Pro:** The number of clusters is not pre-defined and can be chosen after inspecting the dendrogram

- Pro: The clusters can be visualized and interpreted with dendrograms
- **Pro:** The number of clusters is not pre-defined and can be chosen after inspecting the dendrogram
- Pro: In some problems, hierarchy of clusters is preferred over flat clusters

- Pro: The clusters can be visualized and interpreted with dendrograms
- **Pro:** The number of clusters is not pre-defined and can be chosen after inspecting the dendrogram
- Pro: In some problems, hierarchy of clusters is preferred over flat clusters
- **Con:** Hierarchical clustering can't handle big data well, because of its quadratic time complexity

1. Choose K, the number of potential clusters



- 1. Choose K, the number of potential clusters
- 2. Initialise cluster centers randomly within the data



- 1. Choose K, the number of potential clusters
- 2. Initialise cluster centers randomly within the data
- 3. Instances are clustered to the nearest (Euclidean distance) cluster centre



- 1. Choose K, the number of potential clusters
- 2. Initialise cluster centers randomly within the data
- 3. Instances are clustered to the nearest (Euclidean distance) cluster centre



- 1. Choose K, the number of potential clusters
- 2. Initialise cluster centers randomly within the data
- 3. Instances are clustered to the nearest (Euclidean distance) cluster centre
- 4. Centroids of each of the K clusters become new cluster centers



- 1. Choose K, the number of potential clusters
- 2. Initialise cluster centers randomly within the data
- 3. Instances are clustered to the nearest (Euclidean distance) cluster centre
- 4. Centroids of each of the K clusters become new cluster centers



- 1. Choose K, the number of potential clusters
- 2. Initialise cluster centers randomly within the data
- 3. Instances are clustered to the nearest (Euclidean distance) cluster centre
- 4. Centroids of each of the K clusters become new cluster centers
- 5. Steps 3 and 4 are repeated until convergence



- 1. Choose K, the number of potential clusters
- 2. Initialise cluster centers randomly within the data
- 3. Instances are clustered to the nearest (Euclidean distance) cluster centre
- 4. Centroids of each of the K clusters become new cluster centers
- 5. Steps 3 and 4 are repeated until convergence



FAST Foundation

21 Dec 2020 26 / 36

Hierarchical vs K-means



- 4 ⊒ →

How to choose K in K-means?

• Elbow method: increase K until it does not help to describe data better
How to choose K in K-means?

- Elbow method: increase K until it does not help to describe data better
- We are interested in finding K such that the sum of within-group Euclidean distances is smaller

$$J = \sum_{j=1}^{K} \sum_{i=1}^{n} \|\mathbf{x}_{i}^{(j)} - \mathbf{c}_{j}\|^{2},$$

where c_j is the centroid (mean) of the j^{th} cluster



Summary of K-means

Pro: Simple, easy to implement

- Pro: Simple, easy to implement
- Pro: Easy to interpret the clustering results;

- Pro: Simple, easy to implement
- Pro: Easy to interpret the clustering results;
- Pro: Fast and efficient in terms of computational cost

- Pro: Simple, easy to implement
- Pro: Easy to interpret the clustering results;
- Pro: Fast and efficient in terms of computational cost
- **Con:** The number of clusters (K) needs to be defined in advance

- Pro: Simple, easy to implement
- Pro: Easy to interpret the clustering results;
- Pro: Fast and efficient in terms of computational cost
- **Con:** The number of clusters (K) needs to be defined in advance
- **Con:** The dissimilarity measure is fixed to Euclidean distance and features should be quantitative (numeric)

- Pro: Simple, easy to implement
- Pro: Easy to interpret the clustering results;
- Pro: Fast and efficient in terms of computational cost
- **Con:** The number of clusters (K) needs to be defined in advance
- **Con:** The dissimilarity measure is fixed to Euclidean distance and features should be quantitative (numeric)
- **Con:** Squared Euclidean distance places the highest influence on the largest distances, therefore this approach is sensitive to outliers in the data

- Pro: Simple, easy to implement
- Pro: Easy to interpret the clustering results;
- Pro: Fast and efficient in terms of computational cost
- **Con:** The number of clusters (K) needs to be defined in advance
- **Con:** The dissimilarity measure is fixed to Euclidean distance and features should be quantitative (numeric)
- **Con:** Squared Euclidean distance places the highest influence on the largest distances, therefore this approach is sensitive to outliers in the data
- **Con:** In K Means clustering, the results produced by running the algorithm multiple times might differ because of the random initialization of the centroids. While results are reproducible in Hierarchical clustering.

K-medoids (Partitioning Around Medoids) clustering

1. Choose K, the number of potential clusters

- 1. Choose K, the number of potential clusters
- 2. Initialise cluster medoids (central points) randomly within the data

- 1. Choose K, the number of potential clusters
- 2. Initialise cluster medoids (central points) randomly within the data
- 3. Instances are clustered to the nearest cluster medoid according to a predefined dissimilarity measure

- 1. Choose K, the number of potential clusters
- 2. Initialise cluster medoids (central points) randomly within the data
- 3. Instances are clustered to the nearest cluster medoid according to a predefined dissimilarity measure

- 1. Choose K, the number of potential clusters
- 2. Initialise cluster medoids (central points) randomly within the data
- 3. Instances are clustered to the nearest cluster medoid according to a predefined dissimilarity measure
- 4. Medoids of each of the K clusters are updated, taking the ones that are closer to all other points in the cluster

- 1. Choose K, the number of potential clusters
- 2. Initialise cluster medoids (central points) randomly within the data
- 3. Instances are clustered to the nearest cluster medoid according to a predefined dissimilarity measure
- 4. Medoids of each of the K clusters are updated, taking the ones that are closer to all other points in the cluster

- 1. Choose K, the number of potential clusters
- 2. Initialise cluster medoids (central points) randomly within the data
- 3. Instances are clustered to the nearest cluster medoid according to a predefined dissimilarity measure
- 4. Medoids of each of the K clusters are updated, taking the ones that are closer to all other points in the cluster
- 5. Steps 3 and 4 are repeated until convergence







FAST Foundation

21 Dec 2020 31 / 36































FAST Foundation

21 Dec 2020 31 / 36







FAST Foundation

21 Dec 2020 31 / 36
Density-Based Spatial Clustering of Applications with Noise (DBSCAN)



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



• A key component in cluster analysis is the notion of similarity (or dissimilarity) between the individual data objects being clustered

- A key component in cluster analysis is the notion of similarity (or dissimilarity) between the individual data objects being clustered
- How can we measure distances (dissimilarities)?

- A key component in cluster analysis is the notion of similarity (or dissimilarity) between the individual data objects being clustered
- How can we measure distances (dissimilarities)?
- Choosing a good distance measure is critically important for clustering

- A key component in cluster analysis is the notion of similarity (or dissimilarity) between the individual data objects being clustered
- How can we measure distances (dissimilarities)?
- Choosing a good distance measure is critically important for clustering
- Distance measure formalizes how to describe which objects are "close" to each other, and which are not

- A key component in cluster analysis is the notion of similarity (or dissimilarity) between the individual data objects being clustered
- How can we measure distances (dissimilarities)?
- Choosing a good distance measure is critically important for clustering
- Distance measure formalizes how to describe which objects are "close" to each other, and which are not
- On the same dataset the clustering result can be very different if changing the distance measure









1. Distance measure is usually required to be a metric:

- 1. Distance measure is usually required to be a metric:
 - $d(x,y) \ge 0, d(x,y) = 0$ iff x = y (positive definite)

- 1. Distance measure is usually required to be a metric:
 - $d(x,y) \ge 0, \, d(x,y) = 0 \text{ iff } x = y \text{ (positive definite)}$
 - d(x,y) = d(y,x) (symmetry)

- 1. Distance measure is usually required to be a metric:
 - $d(x,y) \ge 0, \, d(x,y) = 0 \text{ iff } x = y \text{ (positive definite)}$
 - d(x,y) = d(y,x) (symmetry)
 - $d(x,z) \leq d(x,y) + d(y,z)$ (triangle inequality)

- 1. Distance measure is usually required to be a metric:
 - $d(x,y) \ge 0, \, d(x,y) = 0 \text{ iff } x = y \text{ (positive definite)}$
 - d(x,y) = d(y,x) (symmetry)
 - $d(x,z) \leq d(x,y) + d(y,z)$ (triangle inequality)
- 2. Sometimes triangle inequality is dropped (semi-metric) and we refer to the measure as **dissimilarity** measure

1. Decide which features should be included into distance calculation

- 1. Decide which features should be included into distance calculation
- 2. Make all features into numeric by encoding each categorical feature as multiple binary features (one-hot encoding)

- 1. Decide which features should be included into distance calculation
- 2. Make all features into numeric by encoding each categorical feature as multiple binary features (one-hot encoding)
- 3. Rescale all features so that the difference of size 1 means roughly the same in each feature

- 1. Decide which features should be included into distance calculation
- 2. Make all features into numeric by encoding each categorical feature as multiple binary features (one-hot encoding)
- 3. Rescale all features so that the difference of size 1 means roughly the same in each feature
- 4. Think about some well-known measures (euclidean, manhattan etc.), see which makes more sense to use or design your own dissimilarity measure

- ✓ Concepts of Cluster Analysis
- ✓ Hierarchical Clustering
- ✓ K-means Clustering
- ✓ K-medoids Clustering
- DBSCAN
- ✓ Dissimilariy Measure Selection