

Introduction to Machine Learning

Introduction to Clustering

Andres Mendez-Vazquez

July 31, 2018

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Supervised Learning vs. Unsupervised Learning

Supervised learning:

- Discover patterns in the data that relate data attributes with a target (class) attribute.

Unsupervised learning:

The data have no target attribute.

Supervised Learning vs. Unsupervised Learning

Supervised learning:

- Discover patterns in the data that relate data attributes with a target (class) attribute.

Unsupervised learning:

The data have no target attribute.

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- **Clustering**
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Clustering

Clustering

It is a technique for finding similarity groups in data, called clusters.

Clustering

Clustering

It is a technique for finding similarity groups in data, called clusters.

Called

An unsupervised learning task as no class values denoting an a priori grouping of the data instances are given, which is the case in supervised learning.

Clustering

Clustering

It is a technique for finding similarity groups in data, called clusters.

Called

An unsupervised learning task as no class values denoting an a priori grouping of the data instances are given, which is the case in supervised learning.

Due to historical reasons

Clustering is often considered synonymous with unsupervised learning.

• In fact, association rule mining is also unsupervised.

Clustering

Clustering

It is a technique for finding similarity groups in data, called clusters.

Called

An unsupervised learning task as no class values denoting an a priori grouping of the data instances are given, which is the case in supervised learning.

Due to historical reasons

Clustering is often considered synonymous with unsupervised learning.

- In fact, association rule mining is also unsupervised.

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- **Pattern Recognition**
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Pattern Recognition

Definition

Search for structure in data

Pattern Recognition

Definition

Search for structure in data

Elements of Numerical Pattern Recognition

1 Process Description

- ▶ Feature Nomination, Test Data, Design Data

2 Feature Analysis

- ▶ Preprocessing, Extraction, Selection, ...

3 Cluster Analysis

- ▶ Labeling, Validity, ...

4 Classifier Design

- ▶ Classification, Estimation, Prediction, Control, ...

Pattern Recognition

Definition

Search for structure in data

Elements of Numerical Pattern Recognition

- 1 Process Description
 - ▶ Feature Nomination, Test Data, Design Data
- 2 Feature Analysis
 - ▶ Preprocessing, Extraction, Selection, ...
- 3 Cluster Analysis
 - ▶ Labeling, Validity, ...
- 4 Classifier Design
 - ▶ Classification, Estimation, Prediction, Control, ...

Pattern Recognition

Definition

Search for structure in data

Elements of Numerical Pattern Recognition

- 1 Process Description
 - ▶ Feature Nomination, Test Data, Design Data
- 2 Feature Analysis
 - ▶ Preprocessing, Extraction, Selection, ...
- 3 Cluster Analysis
 - ▶ Labeling, Validity, ...
- 4 Classifier Design
 - ▶ Classification, Estimation, Prediction, Control, ...

Pattern Recognition

Definition

Search for structure in data

Elements of Numerical Pattern Recognition

- 1 Process Description
 - ▶ Feature Nomination, Test Data, Design Data
- 2 Feature Analysis
 - ▶ Preprocessing, Extraction, Selection, ...
- 3 Cluster Analysis
 - ▶ Labeling, Validity, ...
- 4 Classifier Design
 - ▶ Classification, Estimation, Prediction, Control, ...

Pattern Recognition

Definition

Search for structure in data

Elements of Numerical Pattern Recognition

- 1 Process Description
 - ▶ Feature Nomination, Test Data, Design Data
- 2 Feature Analysis
 - ▶ Preprocessing, Extraction, Selection, ...
- 3 **Cluster Analysis**
 - ▶ Labeling, Validity, ...
- 4 Classifier Design
 - ▶ Classification, Estimation, Prediction, Control, ...

Pattern Recognition

Definition

Search for structure in data

Elements of Numerical Pattern Recognition

- 1 Process Description
 - ▶ Feature Nomination, Test Data, Design Data
- 2 Feature Analysis
 - ▶ Preprocessing, Extraction, Selection, ...
- 3 **Cluster Analysis**
 - ▶ Labeling, Validity, ...
- 4 Classifier Design
 - ▶ Classification, Estimation, Prediction, Control, ...

Pattern Recognition

Definition

Search for structure in data

Elements of Numerical Pattern Recognition

- 1 Process Description
 - ▶ Feature Nomination, Test Data, Design Data
- 2 Feature Analysis
 - ▶ Preprocessing, Extraction, Selection, ...
- 3 **Cluster Analysis**
 - ▶ Labeling, Validity, ...
- 4 Classifier Design
 - ▶ Classification, Estimation, Prediction, Control, ...

Pattern Recognition

Definition

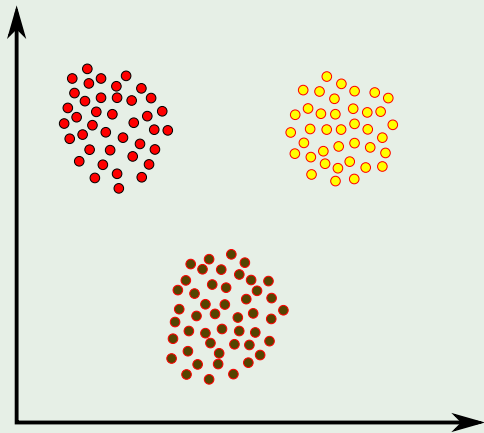
Search for structure in data

Elements of Numerical Pattern Recognition

- 1 Process Description
 - ▶ Feature Nomination, Test Data, Design Data
- 2 Feature Analysis
 - ▶ Preprocessing, Extraction, Selection, ...
- 3 **Cluster Analysis**
 - ▶ Labeling, Validity, ...
- 4 Classifier Design
 - ▶ Classification, Estimation, Prediction, Control, ...

An illustration

The data set has three natural groups of data points, i.e., 3 natural clusters.



Examples

Example 1

Groups people of similar sizes together to make “small”, “medium” and “large” T-Shirts.

Example 2

In marketing, segment customers according to their similarities.

Examples

Example 1

Groups people of similar sizes together to make “small”, “medium” and “large” T-Shirts.

Example 2

In marketing, segment customers according to their similarities.

How we create this classes?

For this, we use the following concept

Clustering!!!

Basically,

We want to "reveal" the organization of patterns into "sensible" clusters (groups).

Intuitively,

Clustering is one of the most primitive mental activities of humans, used to handle the huge amount of information they receive every day.

How we create this classes?

For this, we use the following concept

Clustering!!!

Basically

We want to “reveal” the organization of patterns into “sensible” clusters (groups).

Clustering is one of the most primitive mental activities of humans, used to handle the huge amount of information they receive every day.

How we create this classes?

For this, we use the following concept

Clustering!!!

Basically

We want to “reveal” the organization of patterns into “sensible” clusters (groups).

Actually

Clustering is one of the most primitive mental activities of humans, used to handle the huge amount of information they receive every day.

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- **Why Clustering?**
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

What is clustering for?

Data Reduction

We can use clustering to reduce the amount of samples in each group to reduce the processing of information.

- Example:

- ▶ In data transmission, a representative for each cluster is defined.
- ▶ Then, instead of transmitting the data samples, we transmit a code number corresponding to the representative.
- ▶ Thus, data compression is achieved.

What is clustering for?

Data Reduction

We can use clustering to reduce the amount of samples in each group to reduce the processing of information.

- Example:
 - ▶ In data transmission, a representative for each cluster is defined.
 - ▶ Then, instead of transmitting the data samples, we transmit a code number corresponding to the representative.
 - ▶ Thus, data compression is achieved.

Hypothesis Generation

In this case we apply cluster analysis to a data set in order to infer some hypotheses concerning the nature of the data.

What is clustering for?

Data Reduction

We can use clustering to reduce the amount of samples in each group to reduce the processing of information.

- Example:
 - ▶ In data transmission, a representative for each cluster is defined.
 - ▶ Then, instead of transmitting the data samples, we transmit a code number corresponding to the representative.
 - ▶ Thus, data compression is achieved.

Hypothesis Generation

In this case we apply cluster analysis to a data set in order to infer some hypotheses concerning the nature of the data.

What is clustering for?

Data Reduction

We can use clustering to reduce the amount of samples in each group to reduce the processing of information.

- Example:
 - ▶ In data transmission, a representative for each cluster is defined.
 - ▶ Then, instead of transmitting the data samples, we transmit a code number corresponding to the representative.
 - ▶ Thus, data compression is achieved.

Hypothesis Generation

In this case we apply cluster analysis to a data set in order to infer some hypotheses concerning the nature of the data.

What is clustering for?

Data Reduction

We can use clustering to reduce the amount of samples in each group to reduce the processing of information.

- Example:
 - ▶ In data transmission, a representative for each cluster is defined.
 - ▶ Then, instead of transmitting the data samples, we transmit a code number corresponding to the representative.
 - ▶ Thus, data compression is achieved.

Hypothesis Generation

In this case we apply cluster analysis to a data set in order to infer some hypotheses concerning the nature of the data.

What is clustering for?

Hypothesis testing

In this context, cluster analysis is used for the verification of the validity of a specific hypothesis.

- Example:

- ▶ “Big Companies Invest Abroad”

What is clustering for?

Hypothesis testing

In this context, cluster analysis is used for the verification of the validity of a specific hypothesis.

- Example:

- ▶ “Big Companies Invest Abroad”

Prediction based on clusters

First: We apply cluster analysis to the available data set.

Second: The resulting clusters are characterized based on the characteristics of the patterns by which they are formed.

Third: If we are given an unknown pattern, we can determine the cluster to which it is more likely to belong.

What is clustering for?

Hypothesis testing

In this context, cluster analysis is used for the verification of the validity of a specific hypothesis.

- Example:
 - ▶ “Big Companies Invest Abroad”

Diagnostic-based clustering

First: We apply cluster analysis to the available data set.

Second: The resulting clusters are characterized based on the characteristics of the patterns by which they are formed.

Third: If we are given an unknown pattern, we can determine the cluster to which it is more likely to belong.

What is clustering for?

Hypothesis testing

In this context, cluster analysis is used for the verification of the validity of a specific hypothesis.

- Example:
 - ▶ “Big Companies Invest Abroad”

Prediction based on groups

First: We apply cluster analysis to the available data set.

Second: The resulting clusters are characterized based on the characteristics of the patterns by which they are formed.

Third: If we are given an unknown pattern, we can determine the cluster to which it is more likely to belong.

What is clustering for?

Hypothesis testing

In this context, cluster analysis is used for the verification of the validity of a specific hypothesis.

- Example:
 - ▶ “Big Companies Invest Abroad”

Prediction based on groups

First: We apply cluster analysis to the available data set.

Second: The resulting clusters are characterized based on the characteristics of the patterns by which they are formed.

Third: If we are given an unknown pattern, we can determine the cluster to which it is more likely to belong.

What is clustering for?

Hypothesis testing

In this context, cluster analysis is used for the verification of the validity of a specific hypothesis.

- Example:
 - ▶ “Big Companies Invest Abroad”

Prediction based on groups

- First:** We apply cluster analysis to the available data set.
- Second:** The resulting clusters are characterized based on the characteristics of the patterns by which they are formed.
- Third:** If we are given an unknown pattern, we can determine the cluster to which it is more likely to belong.

What is clustering for?

Hypothesis testing

In this context, cluster analysis is used for the verification of the validity of a specific hypothesis.

- Example:
 - ▶ “Big Companies Invest Abroad”

Prediction based on groups

- First:** We apply cluster analysis to the available data set.
- Second:** The resulting clusters are characterized based on the characteristics of the patterns by which they are formed.
- Third:** If we are given an unknown pattern, we can determine the cluster to which it is more likely to belong.

Aspects of clustering

A clustering algorithm - They are Many!!!

- Partition clustering.
- Hierarchical clustering.
- etc.

Aspects of clustering

A clustering algorithm - They are Many!!!

- Partition clustering.
- Hierarchical clustering.
- etc.

Based on a function

A distance (similarity, or dissimilarity) function.

Aspects of clustering

A clustering algorithm - They are Many!!!

- Partition clustering.
- Hierarchical clustering.
- etc.

Based on a function

A distance (similarity, or dissimilarity) function.

Clustering quality

- Inter-clusters distance \rightarrow maximized.
- Intra-clusters distance \rightarrow minimized.

Aspects of clustering

A clustering algorithm - They are Many!!!

- Partition clustering.
- Hierarchical clustering.
- etc.

Based in a function

A distance (similarity, or dissimilarity) function.

Improving quality

- Inter-clusters distance \rightarrow maximized.
- Intra-clusters distance \rightarrow minimized.

Aspects of clustering

A clustering algorithm - They are Many!!!

- Partition clustering.
- Hierarchical clustering.
- etc.

Based in a function

A distance (similarity, or dissimilarity) function.

Clustering quality

- Inter-clusters distance \rightarrow maximized.
- Intra-clusters distance \rightarrow minimized.

Aspects of clustering

A clustering algorithm - They are Many!!!

- Partition clustering.
- Hierarchical clustering.
- etc.

Based in a function

A distance (similarity, or dissimilarity) function.

Clustering quality

- Inter-clusters distance \rightarrow maximized.
- Intra-clusters distance \rightarrow minimized.

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- **Two Important Models of Clustering**

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Introduction

Observation

Clusters are considered as groups containing data objects that are similar to each other.

When they are not

We assume there is some way to measure that difference!!

Introduction

Observation

Clusters are considered as groups containing data objects that are similar to each other.

When they are not

We assume there is some way to measure that difference!!!

Therefore

It is possible to give a mathematical definition of two important types of clustering

- Partitional Clustering
- Hierarchical Clustering

Therefore

It is possible to give a mathematical definition of two important types of clustering

- Partitional Clustering
- Hierarchical Clustering

Given

Given a set of input patterns $X = \{x_1, x_2, \dots, x_N\}$, where $x_j = (x_{1j}, x_{2j}, \dots, x_{dj})^T \in \mathbb{R}^d$, with each measure x_{ij} called a feature (attribute, dimension, or variable).

Therefore

It is possible to give a mathematical definition of two important types of clustering

- Partitional Clustering
- Hierarchical Clustering

Given

Given a set of input patterns $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, where $\mathbf{x}_j = (x_{1j}, x_{2j}, \dots, x_{dj})^T \in \mathbb{R}^d$, with each measure x_{ij} called a feature (attribute, dimension, or variable).

Hard Partitional Clustering

Definition

Hard partitional clustering attempts to seek a K -partition of X , $C = \{C_1, C_2, \dots, C_K\}$ with $K \leq N$ such that

- $C_i \neq \emptyset$ for $i = 1, \dots, K$.
- $\bigcup_{i=1}^K C_i = X$.
- $C_i \cap C_j = \emptyset$ for all $i, j = 1, \dots, K$ and $i \neq j$.

Hard Partitional Clustering

Definition

Hard partitional clustering attempts to seek a K -partition of X , $C = \{C_1, C_2, \dots, C_K\}$ with $K \leq N$ such that

- 1 $C_i \neq \emptyset$ for $i = 1, \dots, K$.
- 2 $\bigcup_{i=1}^K C_i = X$.
- 3 $C_i \cap C_j = \emptyset$ for all $i, j = 1, \dots, K$ and $i \neq j$.

The third property can be relaxed

- To obtain smoother versions of Hard Clustering

Hard Partitional Clustering

Definition

Hard partitional clustering attempts to seek a K -partition of \mathbf{X} , $\mathbf{C} = \{C_1, C_2, \dots, C_K\}$ with $K \leq N$ such that

- 1 $C_i \neq \emptyset$ for $i = 1, \dots, K$.
- 2 $\cup_{i=1}^K C_i = \mathbf{X}$.
- 3 $C_i \cap C_j = \emptyset$ for all $i, j = 1, \dots, K$ and $i \neq j$.

The third property can be relaxed

- To obtain smoother versions of Hard Clustering

Hard Partitional Clustering

Definition

Hard partitional clustering attempts to seek a K -partition of \mathbf{X} , $\mathbf{C} = \{C_1, C_2, \dots, C_K\}$ with $K \leq N$ such that

- 1 $C_i \neq \emptyset$ for $i = 1, \dots, K$.
- 2 $\cup_{i=1}^K C_i = \mathbf{X}$.
- 3 $C_i \cap C_j = \emptyset$ for all $i, j = 1, \dots, K$ and $i \neq j$.

• The third property can be relaxed

• To obtain smoother versions of Hard Clustering

Hard Partitional Clustering

Definition

Hard partitional clustering attempts to seek a K -partition of \mathbf{X} , $\mathbf{C} = \{C_1, C_2, \dots, C_K\}$ with $K \leq N$ such that

- 1 $C_i \neq \emptyset$ for $i = 1, \dots, K$.
- 2 $\cup_{i=1}^K C_i = \mathbf{X}$.
- 3 $C_i \cap C_j = \emptyset$ for all $i, j = 1, \dots, K$ and $i \neq j$.

The third property can be relaxed

- To obtain smoother versions of Hard Clustering

For example, Fuzzy and Possibilistic clustering

Given the membership $A_i(\mathbf{x}_j) \in [0, 1]$

We can have $\sum_{i=1}^K A_i(\mathbf{x}_j) = 1, \forall j$ and $\sum_{j=1}^N A_i(\mathbf{x}_j) < N, \forall i$.

Hierarchical Clustering

Definition

Hierarchical clustering attempts to construct a tree-like, nested structure partition of \mathbf{X} , $\mathbf{H} = \{H_1, \dots, H_Q\}$ with $Q \leq N$ such that

- If $C_i \in H_m$ and $C_j \in H_l$ and $m > l$ imply $C_i \subset C_j$ or $C_i \cap C_j = \emptyset$ for all $i, j = 1, \dots, K$, $i \neq j$ and $m, l = 1, \dots, Q$.

Hierarchical Clustering

Definition

Hierarchical clustering attempts to construct a tree-like, nested structure partition of \mathbf{X} , $\mathbf{H} = \{H_1, \dots, H_Q\}$ with $Q \leq N$ such that

- If $C_i \in H_m$ and $C_j \in H_l$ and $m > l$ imply $C_i \subset C_j$ or $C_i \cap C_j = \emptyset$ for all $i, j = 1, \dots, K$, $i \neq j$ and $m, l = 1, \dots, Q$.

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- **Types of Features**
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Features

Usually

A data object is described by a set of features or variables, usually represented as a multidimensional vector.

Thus

We tend to collect all that information as a pattern matrix of dimensionality $N \times d$.

When

A feature can be classified as continuous, discrete, or binary.

Features

Usually

A data object is described by a set of features or variables, usually represented as a multidimensional vector.

Thus

We tend to collect all that information as a pattern matrix of dimensionality $N \times d$.

When

A feature can be classified as continuous, discrete, or binary.

Features

Usually

A data object is described by a set of features or variables, usually represented as a multidimensional vector.

Thus

We tend to collect all that information as a pattern matrix of dimensionality $N \times d$.

Then

A feature can be classified as continuous, discrete, or binary.

Thus

Continuous

A continuous feature takes values from an uncountably infinite range set.

Discrete

Discrete features only have finite, or at most, a countably infinite number of values.

Binary

Binary or dichotomous features are a special case of discrete features when they have exactly two values.

Thus

Continuous

A continuous feature takes values from an uncountably infinite range set.

Discrete

Discrete features only have finite, or at most, a countably infinite number of values.

Binary

Binary or dichotomous features are a special case of discrete features when they have exactly two values.

Thus

Continuous

A continuous feature takes values from an uncountably infinite range set.

Discrete

Discrete features only have finite, or at most, a countably infinite number of values.

Binary

Binary or dichotomous features are a special case of discrete features when they have exactly two values.

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- **Measurement Levels**

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Measurement Level

Meaning

Another property of features is the measurement level, which reflects the relative significance of numbers.

We have four from lowest to highest:

Nominal, ordinal, interval, and ratio.

Measurement Level

Meaning

Another property of features is the measurement level, which reflects the relative significance of numbers.

We have four from lowest to highest

Nominal, ordinal, interval, and ratio.

Measurement Level

Nominal

- Features at this level are represented with labels, states, or names.
- The numbers are meaningless in any mathematical sense and no mathematical calculation can be made.
- For example, peristalsis \in {hypermotile, normal, hypomotile, absent}.

Measurement Level

Nominal

- Features at this level are represented with labels, states, or names.
- The numbers are meaningless in any mathematical sense and no mathematical calculation can be made.
- For example, peristalsis \in {hypermotile, normal, hypomotile, absent}.

Ordinal

- Features at this level are also names, but with a certain order implied.
- However, the difference between the values is again meaningless.
- For example, abdominal distension \in {slight, moderate, high}.

Measurement Level

Nominal

- Features at this level are represented with labels, states, or names.
- The numbers are meaningless in any mathematical sense and no mathematical calculation can be made.
- For example, peristalsis \in {hypermotile, normal, hypomotile, absent}.

Ordinal

- Features at this level are also names, but with a certain order implied.
- However, the difference between the values is again meaningless.
- For example, abdominal distension \in {slight, moderate, high}.

Measurement Level

Nominal

- Features at this level are represented with labels, states, or names.
- The numbers are meaningless in any mathematical sense and no mathematical calculation can be made.
- For example, peristalsis \in {hypermotile, normal, hypomotile, absent}.

Ordinal

- Features at this level are also names, but with a certain order implied.
- However, the difference between the values is again meaningless.
- For example, abdominal distension \in {slight, moderate, high}.

Measurement Level

Nominal

- Features at this level are represented with labels, states, or names.
- The numbers are meaningless in any mathematical sense and no mathematical calculation can be made.
- For example, peristalsis \in {hypermotile, normal, hypomotile, absent}.

Ordinal

- Features at this level are also names, but with a certain order implied.
- However, the difference between the values is again meaningless.
- For example, abdominal distension \in {slight, moderate, high}

Measurement Level

Nominal

- Features at this level are represented with labels, states, or names.
- The numbers are meaningless in any mathematical sense and no mathematical calculation can be made.
- For example, peristalsis \in {hypermotile, normal, hypomotile, absent}.

Ordinal

- Features at this level are also names, but with a certain order implied.
- However, the difference between the values is again meaningless.
- For example, abdominal distension \in {slight, moderate, high}.

Measurement Level

Interval

- Features at this level offer a meaningful interpretation of the difference between two values.
- However, there exists no true zero and the ratio between two values is meaningless.
- For example, the concept cold can be seen as an interval.

Measurement Level

Interval

- Features at this level offer a meaningful interpretation of the difference between two values.
- However, there exists no true zero and the ratio between two values is meaningless.
- For example, the concept cold can be seen as an interval.

Ratio

- Features at this level possess all the properties of the above levels.
- There is an absolute zero.
- Thus, we have the meaning of ratio!!!

Measurement Level

Interval

- Features at this level offer a meaningful interpretation of the difference between two values.
- However, there exists no true zero and the ratio between two values is meaningless.
- For example, the concept cold can be seen as an interval.

Ratio

- Features at this level possess all the properties of the above levels.
- There is an absolute zero.
- Thus, we have the meaning of ratio!!!

Measurement Level

Interval

- Features at this level offer a meaningful interpretation of the difference between two values.
- However, there exists no true zero and the ratio between two values is meaningless.
- For example, the concept cold can be seen as an interval.

Ratio

- Features at this level possess all the properties of the above levels.
- There is an absolute zero.
- Thus, we have the meaning of ratio!!!

Measurement Level

Interval

- Features at this level offer a meaningful interpretation of the difference between two values.
- However, there exists no true zero and the ratio between two values is meaningless.
- For example, the concept cold can be seen as an interval.

Ratio

- Features at this level possess all the properties of the above levels.
- There is an absolute zero.

● Thus, we have the meaning of ratio!!!

Measurement Level

Interval

- Features at this level offer a meaningful interpretation of the difference between two values.
- However, there exists no true zero and the ratio between two values is meaningless.
- For example, the concept cold can be seen as an interval.

Ratio

- Features at this level possess all the properties of the above levels.
- There is an absolute zero.
- Thus, we have the meaning of ratio!!!

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- **Similarity Measures**
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Similarity Measures

Definition

A Similarity Measure s on X is a function

$$s : X \times X \rightarrow \mathbb{R} \quad (1)$$

Similarity Measures

Definition

A Similarity Measure s on X is a function

$$s : X \times X \rightarrow \mathbb{R} \quad (1)$$

Such that

- 1 $s(\mathbf{x}, \mathbf{y}) \leq s_0$.
- 2 $s(\mathbf{x}, \mathbf{x}) = s_0$ for all $\mathbf{x} \in X$.
- 3 $s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in X$.
- 4 If $s(\mathbf{x}, \mathbf{y}) = s_0 \iff \mathbf{x} = \mathbf{y}$.
- 5 $s(\mathbf{x}, \mathbf{y}) s(\mathbf{y}, \mathbf{z}) \leq s(\mathbf{x}, \mathbf{z}) [s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{z})]$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$.

Similarity Measures

Definition

A Similarity Measure s on X is a function

$$s : X \times X \rightarrow \mathbb{R} \quad (1)$$

Such that

- 1 $s(\mathbf{x}, \mathbf{y}) \leq s_0$.
- 2 $s(\mathbf{x}, \mathbf{x}) = s_0$ for all $\mathbf{x} \in X$.
- 3 $s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in X$.
- 4 If $s(\mathbf{x}, \mathbf{y}) = s_0 \iff \mathbf{x} = \mathbf{y}$.
- 5 $s(\mathbf{x}, \mathbf{y}) s(\mathbf{y}, \mathbf{z}) \leq s(\mathbf{x}, \mathbf{z}) [s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{z})]$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$.

Similarity Measures

Definition

A Similarity Measure s on X is a function

$$s : X \times X \rightarrow \mathbb{R} \quad (1)$$

Such that

- 1 $s(\mathbf{x}, \mathbf{y}) \leq s_0$.
- 2 $s(\mathbf{x}, \mathbf{x}) = s_0$ for all $\mathbf{x} \in X$.
- 3 $s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in X$.

4 If $s(\mathbf{x}, \mathbf{y}) = s_0 \iff \mathbf{x} = \mathbf{y}$.

5 $s(\mathbf{x}, \mathbf{y})s(\mathbf{y}, \mathbf{z}) \leq s(\mathbf{x}, \mathbf{z}) [s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{z})]$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$.

Similarity Measures

Definition

A Similarity Measure s on X is a function

$$s : X \times X \rightarrow \mathbb{R} \quad (1)$$

Such that

- 1 $s(\mathbf{x}, \mathbf{y}) \leq s_0$.
 - 2 $s(\mathbf{x}, \mathbf{x}) = s_0$ for all $\mathbf{x} \in X$.
 - 3 $s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in X$.
 - 4 If $s(\mathbf{x}, \mathbf{y}) = s_0 \iff \mathbf{x} = \mathbf{y}$.
- 5 $s(\mathbf{x}, \mathbf{y})s(\mathbf{y}, \mathbf{z}) \leq s(\mathbf{x}, \mathbf{z}) [s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{z})]$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$.

Similarity Measures

Definition

A Similarity Measure s on X is a function

$$s : X \times X \rightarrow \mathbb{R} \quad (1)$$

Such that

- 1 $s(\mathbf{x}, \mathbf{y}) \leq s_0$.
- 2 $s(\mathbf{x}, \mathbf{x}) = s_0$ for all $\mathbf{x} \in X$.
- 3 $s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in X$.
- 4 If $s(\mathbf{x}, \mathbf{y}) = s_0 \iff \mathbf{x} = \mathbf{y}$.
- 5 $s(\mathbf{x}, \mathbf{y}) s(\mathbf{y}, \mathbf{z}) \leq s(\mathbf{x}, \mathbf{z}) [s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{z})]$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$.

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- **Dissimilarity Measures**

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Dissimilarity Measures

Defintion

A Dissimilarity Measure d on X is a function

$$d : X \times X \rightarrow \mathbb{R} \quad (2)$$

Dissimilarity Measures

Defintion

A Dissimilarity Measure d on X is a function

$$d : X \times X \rightarrow \mathbb{R} \quad (2)$$

Such that

- 1 $d(\mathbf{x}, \mathbf{y}) \geq d_0$ for all $\mathbf{x}, \mathbf{y} \in X$.
- 2 $d(\mathbf{x}, \mathbf{x}) = d_0$ for all $\mathbf{x} \in X$.
- 3 $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in X$.
- 4 If $d(\mathbf{x}, \mathbf{y}) = d_0 \iff \mathbf{x} = \mathbf{y}$.
- 5 $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$.

Dissimilarity Measures

Defintion

A Dissimilarity Measure d on X is a function

$$d : X \times X \rightarrow \mathbb{R} \quad (2)$$

Such that

- 1 $d(\mathbf{x}, \mathbf{y}) \geq d_0$ for all $\mathbf{x}, \mathbf{y} \in X$.
- 2 $d(\mathbf{x}, \mathbf{x}) = d_0$ for all $\mathbf{x} \in X$.
- 3 $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in X$.
- 4 If $d(\mathbf{x}, \mathbf{y}) = d_0 \iff \mathbf{x} = \mathbf{y}$.
- 5 $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$.

Dissimilarity Measures

Definition

A Dissimilarity Measure d on X is a function

$$d : X \times X \rightarrow \mathbb{R} \quad (2)$$

Such that

- 1 $d(\mathbf{x}, \mathbf{y}) \geq d_0$ for all $\mathbf{x}, \mathbf{y} \in X$.
- 2 $d(\mathbf{x}, \mathbf{x}) = d_0$ for all $\mathbf{x} \in X$.
- 3 $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in X$.

4 If $d(\mathbf{x}, \mathbf{y}) = d_0 \iff \mathbf{x} = \mathbf{y}$.

5 $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$.

Dissimilarity Measures

Definition

A Dissimilarity Measure d on X is a function

$$d : X \times X \rightarrow \mathbb{R} \quad (2)$$

Such that

- 1 $d(\mathbf{x}, \mathbf{y}) \geq d_0$ for all $\mathbf{x}, \mathbf{y} \in X$.
 - 2 $d(\mathbf{x}, \mathbf{x}) = d_0$ for all $\mathbf{x} \in X$.
 - 3 $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in X$.
 - 4 If $d(\mathbf{x}, \mathbf{y}) = d_0 \iff \mathbf{x} = \mathbf{y}$.
- 5 $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$.

Dissimilarity Measures

Definition

A Dissimilarity Measure d on X is a function

$$d : X \times X \rightarrow \mathbb{R} \quad (2)$$

Such that

- 1 $d(\mathbf{x}, \mathbf{y}) \geq d_0$ for all $\mathbf{x}, \mathbf{y} \in X$.
- 2 $d(\mathbf{x}, \mathbf{x}) = d_0$ for all $\mathbf{x} \in X$.
- 3 $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in X$.
- 4 If $d(\mathbf{x}, \mathbf{y}) = d_0 \iff \mathbf{x} = \mathbf{y}$.
- 5 $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$.

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- **Real-Valued Vectors**
 - Dissimilarity Measures
 - Similarity Measures
- **Discrete-Valued Vectors**
 - Dissimilarity Measures
 - Similarity Measures
- **Between Sets**

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- **Real-Valued Vectors**
 - **Dissimilarity Measures**
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Dissimilarity Measures

The Hamming distance for $\{0, 1\}$

The vectors here are collections of zeros and ones. Thus,

$$d_H(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^d (x_i - y_i)^2 \quad (3)$$

Euclidean distance

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d |x_i - y_i|^2 \right)^{\frac{1}{2}} \quad (4)$$

Dissimilarity Measures

The Hamming distance for $\{0, 1\}$

The vectors here are collections of zeros and ones. Thus,

$$d_H(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^d (x_i - y_i)^2 \quad (3)$$

Euclidean distance

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d |x_i - y_i|^{\frac{1}{2}} \right)^2 \quad (4)$$

Dissimilarity Measures

The weighted l_p metric DMs

$$d_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d w_i |x_i - y_i|^p \right)^{\frac{1}{p}} \quad (5)$$

Where

Where x_i, y_i are the i th coordinates of x and y , $i = 1, \dots, d$ and $w_i \geq 0$ is the i th weight coefficient.

Important

A well-known representative of the latter category of measures is the Euclidean distance.

Dissimilarity Measures

The weighted l_p metric DMs

$$d_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d w_i |x_i - y_i|^p \right)^{\frac{1}{p}} \quad (5)$$

Where

Where x_i, y_i are the i th coordinates of x and y , $i = 1, \dots, d$ and $w_i \geq 0$ is the i th weight coefficient.

Important

A well-known representative of the latter category of measures is the Euclidean distance.

Dissimilarity Measures

The weighted l_p metric DMs

$$d_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d w_i |x_i - y_i|^p \right)^{\frac{1}{p}} \quad (5)$$

Where

Where x_i, y_i are the i th coordinates of x and y , $i = 1, \dots, d$ and $w_i \geq 0$ is the i th weight coefficient.

Important

A well-known representative of the latter category of measures is the Euclidean distance.

Dissimilarity Measures

The weighted l_2 metric DM can be further generalized as follows

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T B (\mathbf{x} - \mathbf{y})} \quad (6)$$

Where B is a symmetric, positive definite matrix

A special case

This includes the Mahalanobis distance as a special case

Dissimilarity Measures

The weighted l_2 metric DM can be further generalized as follows

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T B (\mathbf{x} - \mathbf{y})} \quad (6)$$

Where B is a symmetric, positive definite matrix

A special case

This includes the Mahalanobis distance as a special case

Dissimilarity Measures

Special l_p metric DMs that are also encountered in practice are

$$d_1(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^d w_i |x_i - y_i| \quad (\text{Weighted Manhattan Norm})$$

$$d_\infty(\mathbf{x}, \mathbf{y}) = \max_{1 \leq i \leq d} w_i |x_i - y_i| \quad (\text{Weighted } l_\infty \text{ Norm})$$

Dissimilarity Measures

Special l_p metric DMs that are also encountered in practice are

$$d_1(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^d w_i |x_i - y_i| \quad (\text{Weighted Manhattan Norm})$$

$$d_\infty(\mathbf{x}, \mathbf{y}) = \max_{1 \leq i \leq d} w_i |x_i - y_i| \quad (\text{Weighted } l_\infty \text{ Norm})$$

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- **Real-Valued Vectors**
 - Dissimilarity Measures
 - **Similarity Measures**
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Similarity Measures

The inner product

$$s_{inner}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^d x_i y_i \quad (7)$$

Closely related to the inner product is the cosine similarity measure

$$s_{cosine}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \quad (8)$$

Where $\|\mathbf{x}\| = \sqrt{\sum_{i=1}^d x_i^2}$ and $\|\mathbf{y}\| = \sqrt{\sum_{i=1}^d y_i^2}$, the length of the vectors!!!

Similarity Measures

The inner product

$$s_{inner}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^d x_i y_i \quad (7)$$

Closely related to the inner product is the cosine similarity measure

$$s_{cosine}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \quad (8)$$

Where $\|\mathbf{x}\| = \sqrt{\sum_{i=1}^d x_i^2}$ and $\|\mathbf{y}\| = \sqrt{\sum_{i=1}^d y_i^2}$, the length of the vectors!!!

Similarity Measures

Pearson's correlation coefficient

$$r_{\text{Pearson}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \quad (9)$$

Similarity Measures

Pearson's correlation coefficient

$$r_{\text{Pearson}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \quad (9)$$

Where $\mathbf{x} = (x_1 - \bar{x}, \dots, x_d - \bar{x})^T$ and $\mathbf{y} = (y_1 - \bar{y}, \dots, y_d - \bar{y})^T$

$$\bar{x} = \frac{1}{d} \sum_{i=1}^d x_i \quad \text{and} \quad \bar{y} = \frac{1}{d} \sum_{i=1}^d y_i \quad (10)$$

Similarity Measures

Pearson's correlation coefficient

$$r_{\text{Pearson}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \quad (9)$$

Where $\mathbf{x} = (x_1 - \bar{x}, \dots, x_d - \bar{x})^T$ and $\mathbf{y} = (y_1 - \bar{y}, \dots, y_d - \bar{y})^T$

$$\bar{x} = \frac{1}{d} \sum_{i=1}^d x_i \quad \text{and} \quad \bar{y} = \frac{1}{d} \sum_{i=1}^d y_i \quad (10)$$

Properties

- It is clear that $r_{\text{Pearson}}(\mathbf{x}, \mathbf{y})$ takes values between -1 and +1.

• The difference between $s_{\text{similarity}}$ and r_{Pearson} does not depend directly on \mathbf{x} and \mathbf{y} but on their corresponding difference vectors.

Similarity Measures

Pearson's correlation coefficient

$$r_{\text{Pearson}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \quad (9)$$

Where $\mathbf{x} = (x_1 - \bar{x}, \dots, x_d - \bar{x})^T$ and $\mathbf{y} = (y_1 - \bar{y}, \dots, y_d - \bar{y})^T$

$$\bar{x} = \frac{1}{d} \sum_{i=1}^d x_i \quad \text{and} \quad \bar{y} = \frac{1}{d} \sum_{i=1}^d y_i \quad (10)$$

Properties

- It is clear that $r_{\text{Pearson}}(\mathbf{x}, \mathbf{y})$ takes values between -1 and +1.
- The difference between s_{inner} and r_{Pearson} does not depend directly on \mathbf{x} and \mathbf{y} but on their corresponding difference vectors.

Similarity Measures

The Tanimoto measure or distance

$$s_T(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \mathbf{x}^T \mathbf{y}} \quad (11)$$

Something Notable

$$s_T(\mathbf{x}, \mathbf{y}) = \frac{1}{1 + \frac{(\mathbf{x}-\mathbf{y})^T(\mathbf{x}-\mathbf{y})}{\mathbf{x}^T \mathbf{y}}} \quad (12)$$

Meaning

The Tanimoto measure is inversely proportional to the squared Euclidean distance.

Similarity Measures

The Tanimoto measure or distance

$$s_T(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \mathbf{x}^T \mathbf{y}} \quad (11)$$

Something Notable

$$s_T(\mathbf{x}, \mathbf{y}) = \frac{1}{1 + \frac{(\mathbf{x}-\mathbf{y})^T(\mathbf{x}-\mathbf{y})}{\mathbf{x}^T \mathbf{y}}} \quad (12)$$

Meaning

The Tanimoto measure is inversely proportional to the squared Euclidean distance.

Similarity Measures

The Tanimoto measure or distance

$$s_T(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \mathbf{x}^T \mathbf{y}} \quad (11)$$

Something Notable

$$s_T(\mathbf{x}, \mathbf{y}) = \frac{1}{1 + \frac{(\mathbf{x}-\mathbf{y})^T(\mathbf{x}-\mathbf{y})}{\mathbf{x}^T \mathbf{y}}} \quad (12)$$

Meaning

The Tanimoto measure is inversely proportional to the squared Euclidean distance.

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- **Discrete-Valued Vectors**
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

We consider now vectors whose coordinates belong to a finite set

Given $F = \{0, 1, 2, \dots, k - 1\}$ where k is a positive integer

There are exactly k^d vectors $\mathbf{x} \in F^d$

Now consider $\mathbf{x}, \mathbf{y} \in F^d$

With $A(\mathbf{x}, \mathbf{y}) = [a_{ij}]$, $i, j = 0, 1, \dots, k - 1$ be a $k \times k$ matrix

where

The element a_{ij} is the number of places where the first vector has the i symbol and the corresponding element of the second vector has the j symbol.

We consider now vectors whose coordinates belong to a finite set

Given $F = \{0, 1, 2, \dots, k - 1\}$ where k is a positive integer

There are exactly k^d vectors $\mathbf{x} \in F^d$

Now consider $\mathbf{x}, \mathbf{y} \in F^d$

With $A(\mathbf{x}, \mathbf{y}) = [a_{ij}]$, $i, j = 0, 1, \dots, k - 1$ be a $k \times k$ matrix

The element a_{ij} is the number of places where the first vector has the i symbol and the corresponding element of the second vector has the j symbol.

We consider now vectors whose coordinates belong to a finite set

Given $F = \{0, 1, 2, \dots, k - 1\}$ where k is a positive integer

There are exactly k^d vectors $\mathbf{x} \in F^d$

Now consider $\mathbf{x}, \mathbf{y} \in F^d$

With $A(\mathbf{x}, \mathbf{y}) = [a_{ij}]$, $i, j = 0, 1, \dots, k - 1$ be a $k \times k$ matrix

Where

The element a_{ij} is the number of places where the first vector has the i symbol and the corresponding element of the second vector has the j symbol.

Using Indicator functions

We can think of each element a_{ij}

$$a_{ij} = \sum_{r=1}^d I [(i, j) == (x_r, y_r)] \quad (13)$$

Thanks to Ricardo Llamas and Gibran Felix!!! Class Summer 2015.

Contingency Matrix

Thus

This matrix is called a **contingency table**.

For example if $v = 0$ and $k = 1$

With $\mathbf{x} = [0, 1, 2, 1, 2, 1]^T$ and $\mathbf{y} = [1, 0, 2, 1, 0, 1]^T$

Then

$$A(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix} \quad (14)$$

Contingency Matrix

Thus

This matrix is called a **contingency table**.

For example if $d = 6$ and $k = 3$

With $\mathbf{x} = [0, 1, 2, 1, 2, 1]^T$ and $\mathbf{y} = [1, 0, 2, 1, 0, 1]^T$

Then

$$A(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix} \quad (14)$$

Contingency Matrix

Thus

This matrix is called a **contingency table**.

For example if $d = 6$ and $k = 3$

With $\mathbf{x} = [0, 1, 2, 1, 2, 1]^T$ and $\mathbf{y} = [1, 0, 2, 1, 0, 1]^T$

Then

$$A(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix} \quad (14)$$

Thus, we have that

Easy to verify that

$$\sum_{i=0}^{k-1} \sum_{j=0}^{k-1} a_{ij} = d \quad (15)$$

Something Notable

Most of the proximity measures between two discrete-valued vectors may be expressed as combinations of elements of matrix $A(x, y)$.

Thus, we have that

Easy to verify that

$$\sum_{i=0}^{k-1} \sum_{j=0}^{k-1} a_{ij} = d \quad (15)$$

Something Notable

Most of the proximity measures between two discrete-valued vectors may be expressed as combinations of elements of matrix $A(\mathbf{x}, \mathbf{y})$.

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Between Sets

Dissimilarity Measures

The Hamming distance

It is defined as the number of places where two vectors differ:

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^{k-1} \sum_{j=0, j \neq i}^{k-1} a_{ij} \quad (16)$$

The summation of all the off-diagonal elements of A , which indicate the positions where \mathbf{x} and \mathbf{y} differ.

Special case: $k=2$, these vectors are binary valued

The vectors $\mathbf{x} \in F^d$ are binary valued and the Hamming distance becomes

$$\begin{aligned} d_H(\mathbf{x}, \mathbf{y}) &= \sum_{i=1}^d (x_i + y_i - 2x_i y_i) \\ &= \sum_{i=1}^d (x_i - y_i)^2 \end{aligned}$$

Dissimilarity Measures

The Hamming distance

It is defined as the number of places where two vectors differ:

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^{k-1} \sum_{j=0, j \neq i}^{k-1} a_{ij} \quad (16)$$

The summation of all the off-diagonal elements of A , which indicate the positions where \mathbf{x} and \mathbf{y} differ.

Special case $k = 2$, thus vectors are binary valued

The vectors $\mathbf{x} \in F^d$ are binary valued and the Hamming distance becomes

$$\begin{aligned} d_H(\mathbf{x}, \mathbf{y}) &= \sum_{i=1}^d (x_i + y_i - 2x_i y_i) \\ &= \sum_{i=1}^d (x_i - y_i)^2 \end{aligned}$$

Dissimilarity Measures

The l_1 distance

It is defined as in the case of the continuous-valued vectors,

$$d_1(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^d w_i |x_i - y_i| \quad (17)$$

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- **Discrete-Valued Vectors**
 - Dissimilarity Measures
 - **Similarity Measures**
- Between Sets

Similarity Measures

The Tanimoto measure

Given two sets X and Y , we have that

$$s_T(X, Y) = \frac{n_{X \cap Y}}{n_X + n_Y - n_{X \cap Y}} \quad (18)$$

where $n_X = |X|$, $n_Y = |Y|$, $n_{X \cap Y} = |X \cap Y|$.

Similarity Measures

The Tanimoto measure

Given two sets X and Y , we have that

$$s_T(X, Y) = \frac{n_{X \cap Y}}{n_X + n_Y - n_{X \cap Y}} \quad (18)$$

where $n_X = |X|$, $n_Y = |Y|$, $n_{X \cap Y} = |X \cap Y|$.

This measure is defined for discrete-valued vectors.

The measure takes into account all pairs of corresponding coordinates, except those whose corresponding coordinates (x_i, y_i) are both 0.

Similarity Measures

The Tanimoto measure

Given two sets X and Y , we have that

$$s_T(X, Y) = \frac{n_{X \cap Y}}{n_X + n_Y - n_{X \cap Y}} \quad (18)$$

where $n_X = |X|$, $n_Y = |Y|$, $n_{X \cap Y} = |X \cap Y|$.

Thus, if x, y are discrete-valued vectors

The measure takes into account all pairs of corresponding coordinates, except those whose corresponding coordinates (x_i, y_i) are both 0.

Thus

We now define

$$n_x = \sum_{i=1}^{k-1} \sum_{j=0}^{k-1} a_{ij} \text{ and } n_y = \sum_{i=0}^{k-1} \sum_{j=1}^{k-1} a_{ij}$$

in other words

n_x and n_y denotes the number of the nonzero coordinate of x and y .

Thus, we have

$$s_T(x, y) = \frac{\sum_{i=1}^{k-1} \sum_{j=i}^{k-1} a_{ij}}{n_x + n_y - \sum_{i=1}^{k-1} \sum_{j=i}^{k-1} a_{ij}} \quad (19)$$

Thus

We now define

$$n_x = \sum_{i=1}^{k-1} \sum_{j=0}^{k-1} a_{ij} \text{ and } n_y = \sum_{i=0}^{k-1} \sum_{j=1}^{k-1} a_{ij}$$

In other words

n_x and n_y denotes the number of the nonzero coordinate of x and y .

Thus, we have

$$s_T(x, y) = \frac{\sum_{i=1}^{k-1} \sum_{j=1}^{k-1} a_{ij}}{n_x + n_y - \sum_{i=1}^{k-1} \sum_{j=1}^{k-1} a_{ij}} \quad (19)$$

Thus

We now define

$$n_x = \sum_{i=1}^{k-1} \sum_{j=0}^{k-1} a_{ij} \text{ and } n_y = \sum_{i=0}^{k-1} \sum_{j=1}^{k-1} a_{ij}$$

In other words

n_x and n_y denotes the number of the nonzero coordinate of \mathbf{x} and \mathbf{y} .

Thus, we have

$$s_T(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i=1}^{k-1} \sum_{j=i}^{k-1} a_{ij}}{n_X + n_Y - \sum_{i=1}^{k-1} \sum_{j=i}^{k-1} a_{ij}} \quad (19)$$

Outline

1 Supervised Learning vs. Unsupervised Learning

- Supervised vs Unsupervised
- Clustering
- Pattern Recognition
- Why Clustering?
- Two Important Models of Clustering

2 Features

- Types of Features
- Measurement Levels

3 Similarity and Dissimilarity Measures

- Similarity Measures
- Dissimilarity Measures

4 Proximity Measures between Two Points

- Real-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- Discrete-Valued Vectors
 - Dissimilarity Measures
 - Similarity Measures
- **Between Sets**

Between Sets

Jaccard Similarity

Given two sets X and Y , we have that

Between Sets

Jaccard Similarity

Given two sets X and Y , we have that

$$J(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} \quad (20)$$

Although there are more stuff to look for

Please Refer to

- 1 Theodoridis' Book Chapter 11.
- 2 Rui Xu and Don Wunsch. 2009. *Clustering*. Wiley-IEEE Press.