

# Data Mining - Appunti

Francesco Lorenzoni

Febrero 2025



# Contents

<b>I Introduction to Data Mining</b>	<b>7</b>
<b>1 Introduction</b>	<b>13</b>
1.1 Definitions . . . . .	13
1.1.1 Knowledge Discovery Loop . . . . .	13
1.1.2 KDD Process . . . . .	13
1.1.3 Data Mining Process . . . . .	14
1.2 Data Understanding . . . . .	16
1.2.1 Data Quality . . . . .	17
1.2.1.1 Histograms . . . . .	17
1.2.1.2 Statistics notions . . . . .	17
1.2.1.3 Box-Plot . . . . .	17
1.3 Data Understanding - Lab . . . . .	17
1.3.1 Data Collections . . . . .	18
1.3.2 Data Types . . . . .	18
1.3.2.1 Tabular . . . . .	18
1.3.2.2 Transaction . . . . .	18
1.3.2.3 Graph . . . . .	19
1.3.2.4 Sequential . . . . .	20
1.3.2.5 Spatial . . . . .	20
1.3.2.6 Attribute types . . . . .	20
1.3.2.7 Values types . . . . .	20
1.3.3 Data Syntax and Semantics . . . . .	20
1.4 Data Cleaning . . . . .	21
1.4.1 Handling Duplicates . . . . .	22
1.4.1.1 Duplicate Features . . . . .	22
1.4.2 Handling Missing Values . . . . .	22
1.4.3 Outliers . . . . .	23
1.4.3.1 Flower Example . . . . .	23
1.5 Data Preparation . . . . .	23
1.5.1 Aggregation . . . . .	24
1.5.2 Reduction . . . . .	24
1.5.2.1 Sampling . . . . .	24
1.5.2.2 Dimensionality Reduction . . . . .	24
1.5.2.3 Feature Subset Selection . . . . .	26
<b>2 Data Representation</b>	<b>27</b>
2.1 Principal Component Analysis (PCA) . . . . .	27
2.1.1 Observations . . . . .	28
2.2 t-SNE . . . . .	28
2.2.1 Similarity phase . . . . .	29
2.2.2 Embedding phase . . . . .	29
2.2.3 Optimization . . . . .	29
2.3 UMAP . . . . .	30
<b>3 Data Cleaning</b>	<b>31</b>
3.1 Anomalous Values . . . . .	31
3.1.1 Discretization . . . . .	31
3.1.1.1 Natural Binning . . . . .	32
3.1.1.2 Equal Frequency Binning . . . . .	32
3.1.1.3 How many bins? . . . . .	32

3.2 Supervised discretization . . . . .	32
3.2.1 Entropy-based Discretization . . . . .	32
3.3 Binarization . . . . .	32
3.3.1 Attribute Transformation . . . . .	33
3.3.1.1 Normalization . . . . .	33
3.3.1.2 Transformation functions . . . . .	33
<b>4 Cluster analysis</b>	<b>35</b>
4.1 Definitions . . . . .	35
4.2 Types of Clustering . . . . .	36
4.3 Similarity . . . . .	38
4.3.1 Similarity and Dissimilarity for Different Attribute Types . . . . .	38
4.3.2 Euclidean Distance . . . . .	38
4.3.2.1 Minkowski Distance . . . . .	38
4.3.3 Binary Similarity . . . . .	39
4.3.4 Cosine Similarity . . . . .	39
4.3.5 Correlation . . . . .	39
4.3.6 Entropy . . . . .	40
4.4 K-Means . . . . .	40
4.4.1 Evaluating K-Means clusters . . . . .	41
4.4.2 Limitations of K-Means . . . . .	41
4.4.2.1 Empty Clusters . . . . .	42
4.4.3 Workflow . . . . .	42
4.4.4 Choosing the number of clusters K . . . . .	42
4.5 Hierarchical Clustering . . . . .	43
4.5.1 Agglomerative vs Divisive . . . . .	43
4.5.1.1 Updating Proximity Matrix . . . . .	43
4.5.2 Divisive Hierarchical Clustering . . . . .	44
4.5.3 Complexity and Limitations . . . . .	44
4.5.3.1 Limitations . . . . .	44
4.6 Density based - DBSCAN . . . . .	44
4.7 Cluster Validity . . . . .	46
4.8 Towards cluster validation . . . . .	46
4.8.1 Measuring validity through correlation . . . . .	46
4.8.2 Internal measures . . . . .	46
4.8.2.1 SSE - Sum of Squared Error . . . . .	46
4.8.2.2 Cohesion and Separation . . . . .	47
4.8.2.3 Silhouette coefficient . . . . .	48
4.8.3 External measures . . . . .	48
4.8.3.1 Entropy . . . . .	48
4.8.3.2 Purity . . . . .	49
<b>5 Anomaly Detection</b>	<b>51</b>
5.1 Outliers . . . . .	51
5.2 Outlier Detection Algorithms . . . . .	51
5.2.1 Distributions . . . . .	51
5.2.1.1 Grubbs test . . . . .	52
5.2.2 Thresholding . . . . .	53
5.2.3 Manifold . . . . .	53
5.2.3.1 Grading neighbors connectivity . . . . .	54
5.2.4 Reach . . . . .	54
5.2.4.1 Reach ratio factor . . . . .	55
5.2.5 Concentration . . . . .	56
5.2.6 Neighborhoods . . . . .	56
<b>6 K-Means</b>	<b>59</b>
6.1 Bisecting K-Means . . . . .	59
6.2 X-Means . . . . .	59
6.2.0.1 Bayesian Information Criterion . . . . .	59
6.3 Mixture Models and the EM Algorithm . . . . .	60
<b>7 Association Analysis</b>	<b>63</b>
7.1 Basic Concepts . . . . .	63
7.1.1 Frequent Itemset . . . . .	63

7.2	Apriori Algorithm . . . . .	64
7.2.1	Closed Itemsets . . . . .	64
7.2.2	Maximal Itemsets . . . . .	65
7.3	Confidence . . . . .	65
7.3.1	Drawbacks . . . . .	66
7.3.1.1	What rules do we want . . . . .	66
7.4	Other criteria . . . . .	66
7.4.1	Lift . . . . .	66
7.4.2	Interest . . . . .	67
7.4.3	Other measures . . . . .	67
7.5	Non-binary Attributes . . . . .	67
7.5.1	Categorical attributes . . . . .	67
7.5.2	Continuous attributes . . . . .	68
7.6	Association Rule Mining . . . . .	68
7.6.1	Rule lists . . . . .	68
7.6.2	Branch and bound algorithms . . . . .	68
7.7	FP Tree Growth . . . . .	69
<b>8</b>	<b>Sequential Pattern Mining</b>	<b>71</b>
8.1	Definitions . . . . .	71
8.1.1	Exercises . . . . .	71
8.1.1.1	Exercise 1 . . . . .	71
8.1.1.2	Exercise 2 . . . . .	72
8.2	Towards an Algorithm . . . . .	72
8.2.1	GSP - Generalized Sequential Pattern . . . . .	72
8.2.1.1	Candidate Generation in GSP . . . . .	73
8.2.1.2	Candidate Generation - Merging Procedure . . . . .	73
8.2.1.3	Merging Examples . . . . .	73
8.2.1.4	Candidate Pruning . . . . .	74
8.3	Timing Constraints . . . . .	74
8.3.1	Contiguous Subsequences . . . . .	74
<b>9</b>	<b>Supervised Machine Learning</b>	<b>75</b>
9.1	Experience or not . . . . .	75
9.2	Improper models . . . . .	75
9.3	Searching for models . . . . .	76
9.3.1	Data Partitioning . . . . .	77
9.3.1.1	Partitioning the dataset . . . . .	77
9.4	Performance evaluation . . . . .	78
9.4.1	Confusion Matrix . . . . .	79
<b>10</b>	<b>Decision Trees</b>	<b>81</b>
10.1	Classification . . . . .	81
10.2	Classification with Decision Trees . . . . .	81
10.2.1	Hunt's Algorithm . . . . .	81
10.2.1.1	Example: Building a Decision Tree Step-by-Step . . . . .	82
10.2.1.2	Splitting Strategies . . . . .	83
10.2.1.3	Continuous attributes . . . . .	83
10.2.1.4	Choosing the best split . . . . .	84
10.2.2	Gini Index . . . . .	84
10.2.2.1	Gini Index Examples . . . . .	84
10.2.2.2	Gini for a collection of nodes . . . . .	85
10.2.2.2.1	Entropy . . . . .	85
10.2.2.3	Comparing measures . . . . .	85
10.3	Decision Trees Wrap Up . . . . .	85
10.4	Model selection . . . . .	87
10.4.1	Evaluating Trees . . . . .	87
10.4.1.1	Estimating Statistical Bounds . . . . .	87
10.4.1.2	Address overfitting . . . . .	88
10.4.1.3	Observations . . . . .	88
10.4.2	Test conditions . . . . .	88
10.5	Model Selection . . . . .	88
10.5.1	Metrics for Performance evaluation . . . . .	88

10.6 Methods for Performance evaluated . . . . .	89
10.7 Methods for Model Comparison . . . . .	89
10.7.1 ROC - Receiver Operating Characteristic . . . . .	89
10.7.2 Significance Testing . . . . .	89
10.7.2.1 Confidence Interval for accuracy . . . . .	89
<b>11 Rule-based Classification</b>	<b>91</b>
11.1 Introduction to Rule-based Classifiers . . . . .	91
11.1.1 Example of Rule-based Classifier . . . . .	91
11.1.2 Application of Rule-Based Classifier . . . . .	91
11.2 Rule Coverage and Accuracy . . . . .	91
11.2.1 How does a Rule-based Classifier Work? . . . . .	91
11.3 Characteristics of Rule Sets . . . . .	91
11.3.1 Strategy 1: Mutually Exclusive and Exhaustive Rules . . . . .	91
11.3.2 Strategy 2: Non-mutually Exclusive and Non-exhaustive Rules . . . . .	92
11.3.3 Ordered Rule Set . . . . .	92
11.3.4 Rule Ordering Schemes . . . . .	92
11.4 Building Classification Rules . . . . .	92
11.5 Direct Method: Sequential Covering . . . . .	92
11.5.1 Learn-One-Rule Function . . . . .	92
11.5.2 Rule Growing Strategies . . . . .	92
11.6 Rule Evaluation for Growing Rules . . . . .	92
11.6.1 Based on Rule Coverage . . . . .	92
11.6.2 Based on Support Count: FOIL's Information Gain . . . . .	93
11.6.3 Based on Statistical Test: Likelihood Ratio Statistic . . . . .	93
11.7 Direct Method: RIPPER . . . . .	93
11.7.1 For 2-class Problem . . . . .	93
11.7.2 For Multi-class Problem . . . . .	93
11.7.3 Growing a Rule in RIPPER . . . . .	93
11.7.4 Building a Rule Set in RIPPER . . . . .	93
11.8 Minimum Description Length (MDL) . . . . .	94
11.9 Indirect Method: C4.5rules . . . . .	94
11.9.1 Algorithm . . . . .	94
11.9.2 Pessimistic Error Estimate . . . . .	94
11.9.3 Class Ordering in C4.5rules . . . . .	94
11.10 Advantages of Rule-Based Classifiers . . . . .	94
11.11 Example: C4.5 vs C4.5rules vs RIPPER . . . . .	94
11.11.1 C4.5rules . . . . .	94
11.11.2 RIPPER . . . . .	94
11.11.3 Performance Comparison . . . . .	94

## Part I

# Introduction to Data Mining



---

<b>1</b>	<b>Introduction</b>	<b>13</b>
1.1	Definitions . . . . .	13
1.1.1	Knowledge Discovery Loop . . . . .	13
1.1.2	KDD Process . . . . .	13
1.1.3	Data Mining Process . . . . .	14
1.2	Data Understanding . . . . .	16
1.2.1	Data Quality . . . . .	17
1.3	Data Understanding - Lab . . . . .	17
1.3.1	Data Collections . . . . .	18
1.3.2	Data Types . . . . .	18
1.3.3	Data Syntax and Semantics . . . . .	20
1.4	Data Cleaning . . . . .	21
1.4.1	Handling Duplicates . . . . .	22
1.4.2	Handling Missing Values . . . . .	22
1.4.3	Outliers . . . . .	23
1.5	Data Preparation . . . . .	23
1.5.1	Aggregation . . . . .	24
1.5.2	Reduction . . . . .	24
<b>2</b>	<b>Data Representation</b>	<b>27</b>
2.1	Principal Component Analysis (PCA) . . . . .	27
2.1.1	Observations . . . . .	28
2.2	t-SNE . . . . .	28
2.2.1	Similarity phase . . . . .	29
2.2.2	Embedding phase . . . . .	29
2.2.3	Optimization . . . . .	29
2.3	UMAP . . . . .	30
<b>3</b>	<b>Data Cleaning</b>	<b>31</b>
3.1	Anomalous Values . . . . .	31
3.1.1	Discretization . . . . .	31
3.2	Supervised discretization . . . . .	32
3.2.1	Entropy-based Discretization . . . . .	32
3.3	Binarization . . . . .	32
3.3.1	Attribute Transformation . . . . .	33
<b>4</b>	<b>Cluster analysis</b>	<b>35</b>
4.1	Definitions . . . . .	35
4.2	Types of Clustering . . . . .	36
4.3	Similarity . . . . .	38
4.3.1	Similarity and Dissimilarity for Different Attribute Types . . . . .	38
4.3.2	Euclidean Distance . . . . .	38
4.3.3	Binary Similarity . . . . .	39
4.3.4	Cosine Similarity . . . . .	39
4.3.5	Correlation . . . . .	39
4.3.6	Entropy . . . . .	40
4.4	K-Means . . . . .	40
4.4.1	Evaluating K-Means clusters . . . . .	41
4.4.2	Limitations of K-Means . . . . .	41
4.4.3	Workflow . . . . .	42
4.4.4	Choosing the number of clusters K . . . . .	42
4.5	Hierarchical Clustering . . . . .	43
4.5.1	Agglomerative vs Divisive . . . . .	43
4.5.2	Divisive Hierarchical Clustering . . . . .	44
4.5.3	Complexity and Limitations . . . . .	44
4.6	Density based - DBSCAN . . . . .	44
4.7	Cluster Validity . . . . .	46
4.8	Towards cluster validation . . . . .	46
4.8.1	Measuring validity through correlation . . . . .	46
4.8.2	Internal measures . . . . .	46
4.8.3	External measures . . . . .	48

<b>5 Anomaly Detection</b>	<b>51</b>
5.1 Outliers . . . . .	51
5.2 Outlier Detection Algorithms . . . . .	51
5.2.1 Distributions . . . . .	51
5.2.2 Thresholding . . . . .	53
5.2.3 Manifold . . . . .	53
5.2.4 Reach . . . . .	54
5.2.5 Concentration . . . . .	56
5.2.6 Neighborhoods . . . . .	56
<b>6 K-Means</b>	<b>59</b>
6.1 Bisecting K-Means . . . . .	59
6.2 X-Means . . . . .	59
6.3 Mixture Models and the EM Algorithm . . . . .	60
<b>7 Association Analysis</b>	<b>63</b>
7.1 Basic Concepts . . . . .	63
7.1.1 Frequent Itemset . . . . .	63
7.2 Apriori Algorithm . . . . .	64
7.2.1 Closed Itemsets . . . . .	64
7.2.2 Maximal Itemsets . . . . .	65
7.3 Confidence . . . . .	65
7.3.1 Drawbacks . . . . .	66
7.4 Other criteria . . . . .	66
7.4.1 Lift . . . . .	66
7.4.2 Interest . . . . .	67
7.4.3 Other measures . . . . .	67
7.5 Non-binary Attributes . . . . .	67
7.5.1 Categorical attributes . . . . .	67
7.5.2 Continuous attributes . . . . .	68
7.6 Association Rule Mining . . . . .	68
7.6.1 Rule lists . . . . .	68
7.6.2 Branch and bound algorithms . . . . .	68
7.7 FP Tree Growth . . . . .	69
<b>8 Sequential Pattern Mining</b>	<b>71</b>
8.1 Definitions . . . . .	71
8.1.1 Exercises . . . . .	71
8.2 Towards an Algorithm . . . . .	72
8.2.1 GSP - Generalized Sequential Pattern . . . . .	72
8.3 Timing Constraints . . . . .	74
8.3.1 Contiguous Subsequences . . . . .	74
<b>9 Supervised Machine Learning</b>	<b>75</b>
9.1 Experience or not . . . . .	75
9.2 Improper models . . . . .	75
9.3 Searching for models . . . . .	76
9.3.1 Data Partitioning . . . . .	77
9.4 Performance evaluation . . . . .	78
9.4.1 Confusion Matrix . . . . .	79
<b>10 Decision Trees</b>	<b>81</b>
10.1 Classification . . . . .	81
10.2 Classification with Decision Trees . . . . .	81
10.2.1 Hunt's Algorithm . . . . .	81
10.2.2 Gini Index . . . . .	84
10.3 Decision Trees Wrap Up . . . . .	85
10.4 Model selection . . . . .	87
10.4.1 Evaluating Trees . . . . .	87
10.4.2 Test conditions . . . . .	88
10.5 Model Selection . . . . .	88
10.5.1 Metrics for Performance evaluation . . . . .	88
10.6 Methods for Performance evaluated . . . . .	89
10.7 Methods for Model Comparison . . . . .	89

10.7.1 ROC - Receiver Operating Characteristic . . . . .	89
10.7.2 Significance Testing . . . . .	89
<b>11 Rule-based Classification</b>	<b>91</b>
11.1 Introduction to Rule-based Classifiers . . . . .	91
11.1.1 Example of Rule-based Classifier . . . . .	91
11.1.2 Application of Rule-Based Classifier . . . . .	91
11.2 Rule Coverage and Accuracy . . . . .	91
11.2.1 How does a Rule-based Classifier Work? . . . . .	91
11.3 Characteristics of Rule Sets . . . . .	91
11.3.1 Strategy 1: Mutually Exclusive and Exhaustive Rules . . . . .	91
11.3.2 Strategy 2: Non-mutually Exclusive and Non-exhaustive Rules . . . . .	92
11.3.3 Ordered Rule Set . . . . .	92
11.3.4 Rule Ordering Schemes . . . . .	92
11.4 Building Classification Rules . . . . .	92
11.5 Direct Method: Sequential Covering . . . . .	92
11.5.1 Learn-One-Rule Function . . . . .	92
11.5.2 Rule Growing Strategies . . . . .	92
11.6 Rule Evaluation for Growing Rules . . . . .	92
11.6.1 Based on Rule Coverage . . . . .	92
11.6.2 Based on Support Count: FOIL's Information Gain . . . . .	93
11.6.3 Based on Statistical Test: Likelihood Ratio Statistic . . . . .	93
11.7 Direct Method: RIPPER . . . . .	93
11.7.1 For 2-class Problem . . . . .	93
11.7.2 For Multi-class Problem . . . . .	93
11.7.3 Growing a Rule in RIPPER . . . . .	93
11.7.4 Building a Rule Set in RIPPER . . . . .	93
11.8 Minimum Description Length (MDL) . . . . .	94
11.9 Indirect Method: C4.5rules . . . . .	94
11.9.1 Algorithm . . . . .	94
11.9.2 Pessimistic Error Estimate . . . . .	94
11.9.3 Class Ordering in C4.5rules . . . . .	94
11.10 Advantages of Rule-Based Classifiers . . . . .	94
11.11 Example: C4.5 vs C4.5rules vs RIPPER . . . . .	94
11.11.1 C4.5rules . . . . .	94
11.11.2 RIPPER . . . . .	94
11.11.3 Performance Comparison . . . . .	94

---

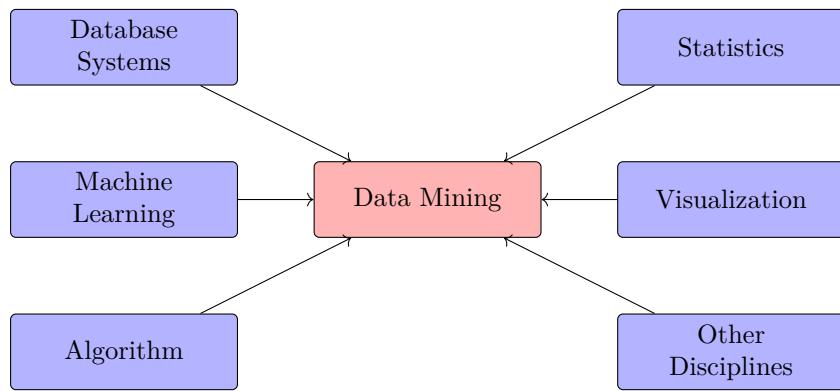


# Chapter 1

## Introduction

**Definition 1.1 (Data Mining)** is the use of efficient techniques for the analysis of very large collections of data and the extraction of useful and possibly unexpected patterns in data (hidden knowledge). The goal is to extract (human-readable) knowledge and insight from raw data.

- ◊ Knowledge implies we are often not just trying to solve a task
- ◊ Insight implies that we should infer non-obvious knowledge
- ◊ Human-readable implies that knowledge should be (when possible) understood by humans: focus on interpretability!
- ◊ Raw data implies we'll need to clean it



### 1.1 Definitions

#### 1.1.1 Knowledge Discovery Loop

Large collections tend to be heterogeneous in source, domain, language and refinement. The first step is to store the data, which however does not assess its heterogeneity. Data cleaning and integration tackle this problem, so that we get integrated sources, homogenous language, and data cleared of noise and outliers.

To look for insight on the data we have to answer questions on the data as a stakeholder. We may see patterns and ask ourselves their nature. Pattern extraction and validation lead to possible insight. Insight may lead to noticing that some data missing may be useful, and we may want to collect it, going back at previous steps.



Figure 1.1: Knowledge Discovery Loop  
This essentially summarizes the KDD process.

#### 1.1.2 KDD Process

The KDD process consists of the following steps:

1. **Data Cleaning:** Remove noise and inconsistent data.

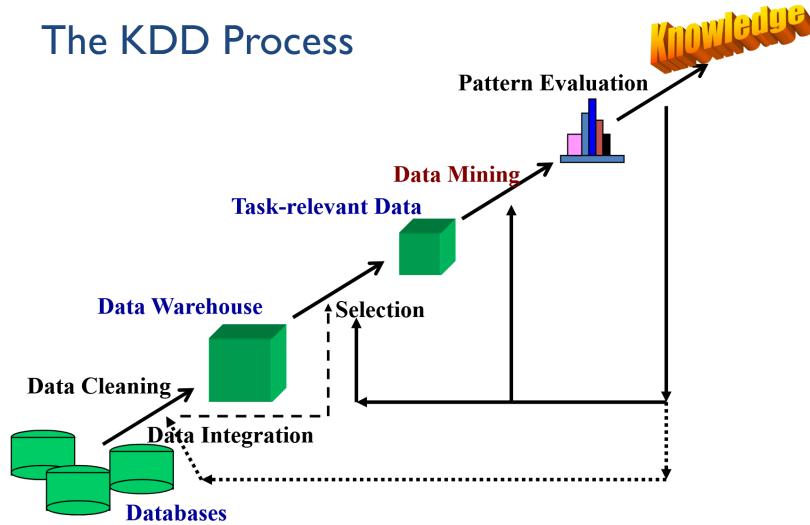


Figure 1.2: KDD Process

## 2. Data Integration:

Combine multiple data sources.

Involves the process of data understanding, data cleaning, merging data coming from multiple sources and transforming them to load them into a **Data Warehouse**.

**Data Warehouse** is a database targeted to answer specific business questions

## 3. Data Selection:

Select relevant data for analysis.

## 4. Data Transformation:

Transform data into suitable formats for mining (summary, aggregation, etc.).

## 5. Data Mining:

Apply algorithms to extract patterns.

◊ *Prediction Methods*

    Use some variables to predict unknown or future values of other variables.

◊ *Description Methods*

    Find human-interpretable patterns that describe the data.

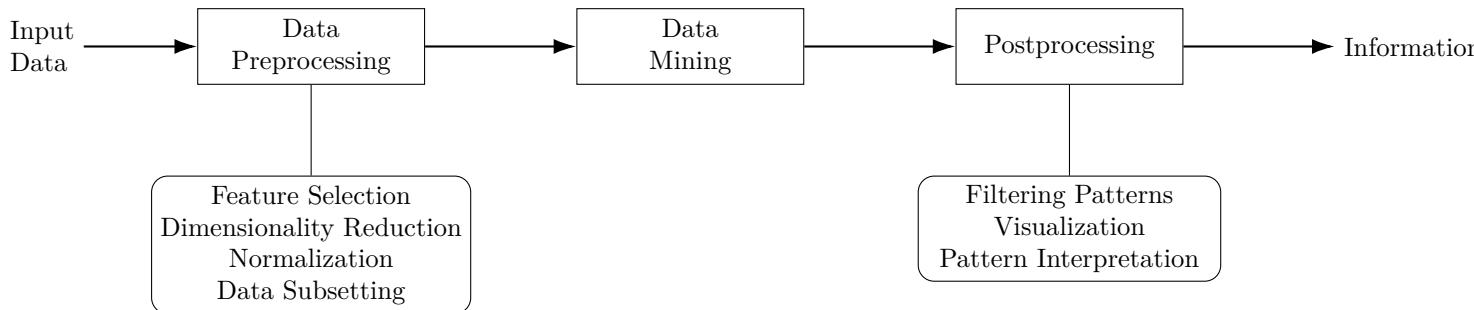
## 6. Pattern Evaluation:

Identify truly interesting patterns.

## 7. Knowledge Presentation:

Present the mined knowledge in an understandable way.

### 1.1.3 Data Mining Process



**Definition 1.2 (Primary Data)** *Original data that has been collected for a specific purpose.*

*Primary data is not altered by humans*

**Definition 1.3 (Secondary Data)** *Data that has been already collected and made available for other purposes.*

*Secondary data may be obtained from many sources*

**Definition 1.4 (Association rule discovery)** *Given a set of records each of which contain some number of items from a given collection.*

*Produce dependency rules which will predict occurrence of an item based on occurrences of other items.*

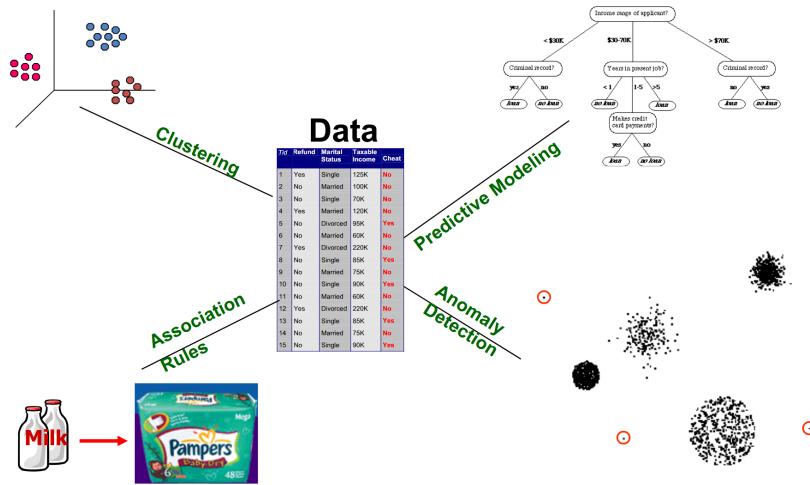


Figure 1.3: Data Mining methods

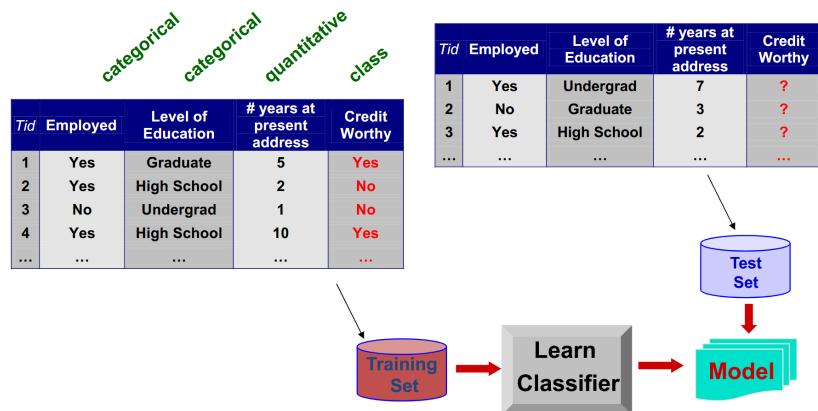


Figure 1.4: Classification Process

<u>Association Use Cases</u>
◊ <b>Market-basket analysis</b> Rules are used for sales promotion, shelf management, and inventory management
◊ <b>Telecommunication alarm diagnosis</b> Rules are used to find combination of alarms that occur together frequently in the same time period
◊ <b>Medical Informatics</b> Rules are used to find combination of patient symptoms and test results associated with certain diseases

## 1.2 Data Understanding

**Definition 1.5 (Data)** *Data is a collection of data objects and their attributes.*

*An attribute is a property or characteristic of an object. A collection of attributes describe an object (record).*

If data objects have the same fixed set of numeric attributes, then the data objects can be thought of as points in a multi-dimensional space, where each dimension represents a distinct attribute.

Such data set can be represented by an  $m \times n$  matrix, where there are  $m$  rows, one for each object, and  $n$  columns, one for each attribute.

Data Types:

- ◊ Document data
- ◊ Transaction data
- ◊ Graph data
- ◊ Ordered data
  - Spatial data
  - Temporal data

The type of the attribute depends on the following properties:

- ◊ Distinctness:  $=\neq$
- ◊ Order:  $<>$
- ◊ Differences are meaningful:  $+-$
- ◊ Ratios are meaningful:  $*/$

Attribute types:

- ◊ Nominal/Categorical: attribute values in a finite domain (*distinctness*)
- ◊ Binary: special case of nominal with two values
- ◊ Ordinal: attribute values have a total ordering (*distinctness* and *order*)
- ◊ Numeric: quantity (integer or real-valued) (*distinctness*, *order*, *differences*)
- ◊ Ratio-Scaled: we can speak of values as being an order of magnitude larger than the unit of measurement (*all 4 properties*)  
length, counts, elapsed time (A baseball game lasting 3 hours is 50% longer than a game lasting 2 hours)
- ◊ Discrete/Continuous: attribute values are discrete (finite or countably infinite) or continuous (real-valued).

Attribute Type	Description	Examples	Operations
Nominal	Nominal attribute values only distinguish. $(=, \neq)$	zip codes, employee ID numbers, eye color, sex: {male, female}	mode, entropy, contingency correlation, $\chi^2$ test
Ordinal	Ordinal attribute values also order objects. $(<, >)$	hardness of minerals, {good, better, best}, grades, street numbers	median, percentiles, rank correlation, run tests, sign tests
Interval	For interval attributes, differences between values are meaningful. $(+, -)$	calendar dates, temperature in Celsius or Fahrenheit	mean, standard deviation, Pearson's correlation, <i>t</i> and <i>F</i> tests
Ratio	For ratio variables, both differences and ratios are meaningful. $(*, /)$	temperature in Kelvin, monetary quantities, counts, age, mass, length, current	geometric mean, harmonic mean, percent variation

Table 1.1: Attribute types, examples and operations

### 1.2.1 Data Quality

Examples of data quality problems:

- ◊ Wrong data
- ◊ Duplicate data
- ◊ Noise and outliers
- ◊ Missing values

In order to know our data and discover quality issues we need use descriptive statistics for getting a global picture and summarize properties of data and compare such statistics with the expected behaviour. All around we can exploit visualization techniques that can help in detecting general or unusual patterns and trends, as well as outliers.

#### 1.2.1.1 Histograms

A histogram shows the frequency distribution for a numerical attribute. The range of the numerical attribute is discretized into a fixed number of intervals (**bins**).

The number of bins according to Sturges' rule is:

$$k = \lceil \log_2 n + 1 \rceil$$

where  $n$  is the number of records in the data set. Sturges' rule is suitable for data from normal distributions and from data sets of moderate size.

#### 1.2.1.2 Statistics notions

Notorious Mean/Median/Mode...The degree in which data tend to spread is called the *dispersion*, or **variance** of the data.

The most common measures for data dispersion are **range** (The distance between the largest and the smallest values), **standard deviation**, the **five-number summary** (based on *quartiles*), and the **inter-quartile range**.

$$\text{variance}(x) = \sigma^2 = s_x^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2$$

Standard deviation  $\sigma$  is the square root of variance  $\sigma^2$ .

Because of outliers, other measures are often used:

- ◊ absolute average deviation (AAD)

$$\text{AAD}(x) = \frac{1}{m} \sum_{i=1}^m |x_i - \bar{x}|$$

- ◊ median average deviation (MAD)

$$\text{MAD}(x) = \text{median}(|x_1 - \bar{x}|, \dots, |x_m - \bar{x}|)$$

#### 1.2.1.3 Box-Plot

## 1.3 Data Understanding - Lab

Data comes from diverse sources, and generally is not tailor-made for some downstream task. We need to start from basics:

- ◊ What features are available?
- ◊ What are they measuring, exactly?
- ◊ What properties do they have?
- ◊ What are their relations?
- ◊ Are there outliers?
- ◊ ...

Data can be of different nature which may co-occur:

- ◊ **Temporal**: the data describes events over time
- ◊ **Sequential**: the data spans some ordering
- ◊ **Relational**: the data describes event in between instances
- ◊ **Spatial**: the data describes space
- ◊ **Independent**: instances in data are independent observations

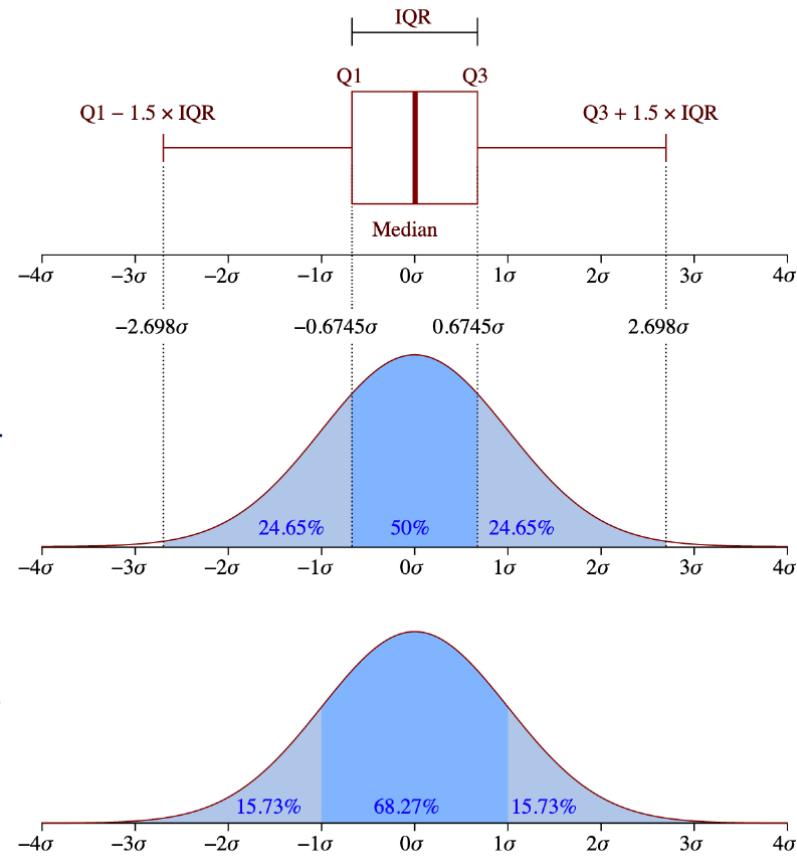


Figure 1.5: Box-Plot

### 1.3.1 Data Collections

We refer to single instances in the collections as objects/records/instances, which are described by attributes.

<b>Id</b>	<b>Age</b>	<b>Income</b>	<b>Marital</b>	<b>Loan</b>
0	30	2.5k	Married	Yes
1	24	1.4k	Single	No
...	...	...	...	...

Table 1.2: Grant Data

- ◊ Attributes: `Id` , `Age` , `Income` , `Marital`,`Loan` grant
- ◊ Records: 0, 30, 2.5k, Married, Yes , 1, 24,1.4k, Single, No

### 1.3.2 Data Types

#### 1.3.2.1 Tabular

When records are independent, and described by the same finite set of features, they are often represented in a tabular form: the data matrix. Each row is a record, each dimension is an attribute.

Records on the rows, attributes on the columns.

<b>Id</b>	<b>Age</b>	<b>Bike used</b>	<b>Length</b>	<b>Duration</b>	<b>Date</b>	<b>Cyclist</b>
0	28	Colnago VRS4	152.4	3:43:12	15-5-2025	Alessandro Covi
1	40	Cervelo RS5	72.4	2:55:01	4-3-2024	Gianni Affino

Table 1.3: Cyclist Data

#### 1.3.2.2 Transaction

A feature contains a (multi)set of items.

PurchaseId	Cart	Bought on
0	Bread, Milk	17:12-15-5-2025
1	Notebook, Pens, Bread, Basil	8:04-4-3-2024

Table 1.4: Transaction Data

Records on the rows, attributes on the columns.

### 1.3.2.3 Graph

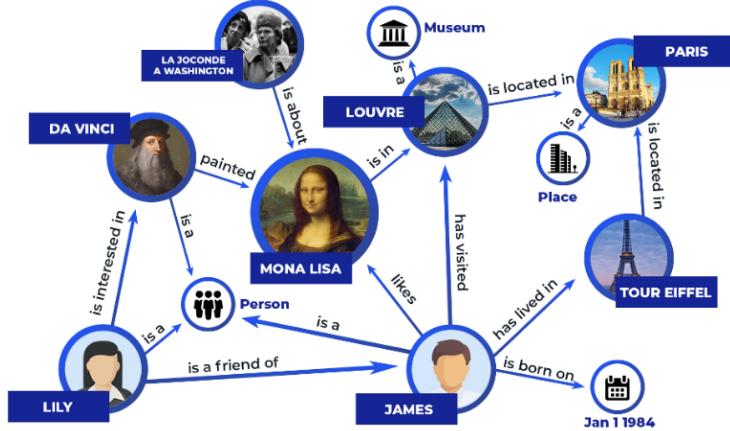


Figure 1.6: Graph Data

Data is linked, either on records or features. Records are nodes in a graph, attributes can vary wildly across records.

### 1.3.2.4 Sequential

Records are sequences (of variable length): attributes are indexed (order or time).

Image on the right lacks two images ☺

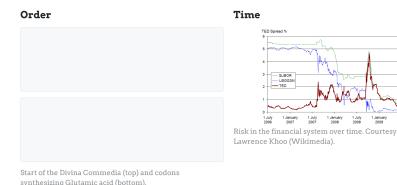


Figure 1.7: Sequential Data

### 1.3.2.5 Spatial

Records are associated with locations in space: attributes can include coordinates, regions, etc.

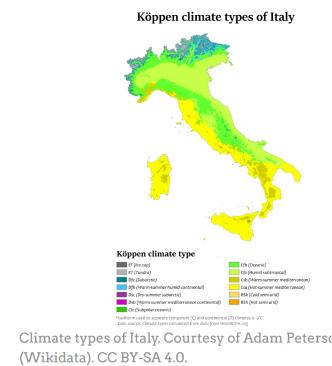


Figure 1.8: Spatial Data

### 1.3.2.6 Attribute types

Type	Description	Example
Numerical	Values have a total ordering, and represent some numerical quantity	Age, dates
Ordinal	Values have a total ordering, and represent some quantity	Dress size, Cup size
Binary	Values are one of two categories: no ordering	Boolean values
Categorical	Values of one of multiple categories: no ordering	Country, Job

Table 1.5: Types of Data

### 1.3.2.7 Values types

Values can be either:

- ◊ **Discrete** Defined in a finite or countably finite domain, e.g., country, job, cup size. Note: ordinal values may be discrete too!
- ◊ **Continuous** Defined in a continuous and infinite domain, e.g., distance.

### 1.3.3 Data Syntax and Semantics

Given the categorization of the records and attributes of your data, we can study its general behavior. We leverage some basic statistical tools, first of all by drawing the empirical distribution of the attributes.



Figure 1.9: Data Syntax and Semantics

Useful statistics for data semantics

- ◊ **Expected value**

$$\mathbb{E}[X] = \sum_{x \in \text{dom}(X)} \Pr(X = x) x$$

A statistic representative of the value of an attribute, weighing values and their probability

- ◊ **Variance**

$$\sigma^2(X) = \mathbb{E} \left[ \sum_{x \in \text{dom}(X)} (x - \mathbb{E}[X])^2 \right]$$

Distance from the expected value of all records: the data spread

- ◊ **Quantiles**

$$q^p = x \text{ s.t. } \Pr(X \leq x) = q^p$$

Inflection points defining values for a threshold, e.g., if the 99-th percentile is , then we

- ◊ **Interquantile range**

$$q^{75} - q^{25}$$

Distance between quantiles: how spread are inflection points?

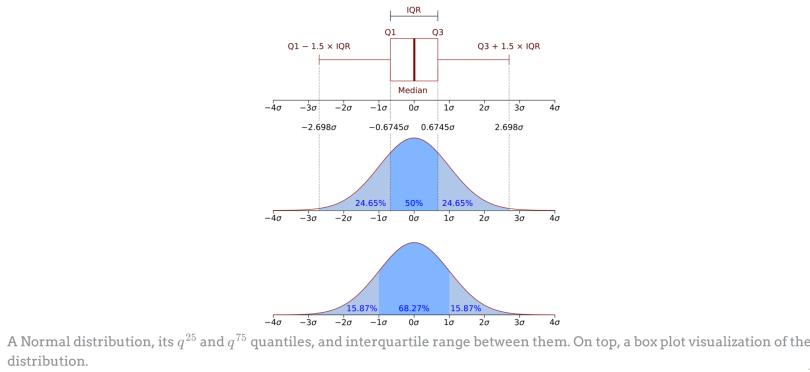


Figure 1.7: Statistics Graph

Statistical summary of the distribution are typically accompanied by visual and semantic one.  
Erroneous or weird values to be cleaned later may already pop up in these basic steps. Outlier values typically skew statistics. Variance is often replaced by absolute/median average deviation

## 1.4 Data Cleaning

There are some concepts to be aware of when dealing with data quality, hence data cleaning.

**Data accuracy** is the degree to which data correctly describes the "real world" object or event being described.

- ◊ Syntactic: values outside domain, e.g., Eataly in Country
- ◊ Semantic: values in domain, but semantically wrong, e.g., age is 3, and weight is 82kg

**Completeness** is the degree to which all required data is known.

Some attributes are not collected, or are collected partially, e.g., temperature was not recorded by the sensor.

**Biased gathering** is the degree to which data may be over/under-representative, e.g., the bank may only provide data about successful loan applicants.

**Timeliness** is the degree to which data is up to date.

Remember: *garbage in, garbage out!*<sup>1</sup> In a task-agnostic view, we are interested in addressing the above by tackling:

- ◊ **Duplicates:** skews the data distribution

<sup>1</sup>i.e. if you have garbage data, you'll get garbage results

- ◊ **Missing values:** give false/partial information
- ◊ **Noise:** uninformative of the data
- ◊ **Poor accuracy:** gives wrong data
- ◊ **Outliers:** skews the data distribution and models of the data

### 1.4.1 Handling Duplicates

Remove them... when appropriate! Not all duplicates are garbage, it depends on what insight you can gather from it.

#### Case A

You have data on registration to your website, with several duplicate e-mails. Insights:

- ◊ The “Sign in” button is hard to find
- ◊ The “Sign in” button is less visible than the “Sign up” button
- ◊ Your site is so anonymous people forget they signed up already

#### 1.4.1.1 Duplicate Features

Duplicate features may be more tricky. Features convey similar, although not equal, information to others.  
Examples:

- ◊ Resting heart rate and heart rate under continuous high effort
- ◊ Education level and reading skills
- ◊ Rent and available bank deposit

These pairs of features are not per se one duplicate of the other, but are strongly related: when one grows, so does the other, and when one goes down, so does the other.

Linear (and rank) relationships between two features  $X, Y$  can be quantified with their correlation. Correlation ranges in  $[-1, 1]$ , from perfectly negative to perfectly positive correlation.

Given two lists of values  $x^{i^n}, y^{i^n}$  we can compute two main correlation types.

- ◊ **Pearson correlation**

$$\rho_P^{X,Y} = \frac{\mathbb{E} [(x^i - \mathbb{E}[X])(y^i - \mathbb{E}[Y])]}{\sigma_X \sigma_Y}$$

Measures linear correlation between two numerical features and their values.

- ◊ **Spearman correlation**

$$\rho_S = \rho_P^{\text{rank}(X), \text{rank}(Y)}$$

Measures monotonic correlation between two ordinal or numerical features.

### 1.4.2 Handling Missing Values

Data may be missing for any number of reasons (at random or not at random).

- ◊ A record has a large and/or significant set of missing attributes
- ◊ An attribute has a large percentage of missing values

We have two choices: **dropping** or **imputing**.

#### Dropping

If a record has a large and/or significant set of missing attributes, or an attribute has a large percentage of missing values, we can drop the record/attribute.

- ◊ High percentage of missing values
- ◊ Missing values in critical attributes, e.g., a patient in cardiology has no heart rate data

#### Imputing

Imputing means replacing the missing value with a “best guess” value.

If a record has a small set of missing attributes, or an attribute has a small percentage of missing values, we can impute the missing values. We have to create a model to predict the missing value.

- ◊ Low percentage of missing values
- ◊ Reasonably good understanding of the attribute semantics/distribution

- ◊ Presence of related attributes

### 1.4.3 Outliers

Quantiles and distributions inform us on what values may be outlier. They are typically dropped, and unlike missing values, almost never imputed. We'll tackle algorithms later in the course.

#### 1.4.3.1 Flower Example

There is a dataset with 5 attributes: sepal length, sepal width, petal length, petal width, and species (type).

Sepal L.	Sepal W.	Petal L.	Petal W.	Type
5.1	3.5	1.4	0.2	Setosa
7.0	3.2	4.7	1.4	Versicolor
...	...	...	...	...

Table 1.6: Flower Dataset Example



The three Iris types in the dataset.

Figure 1.8: Flower Data plotted

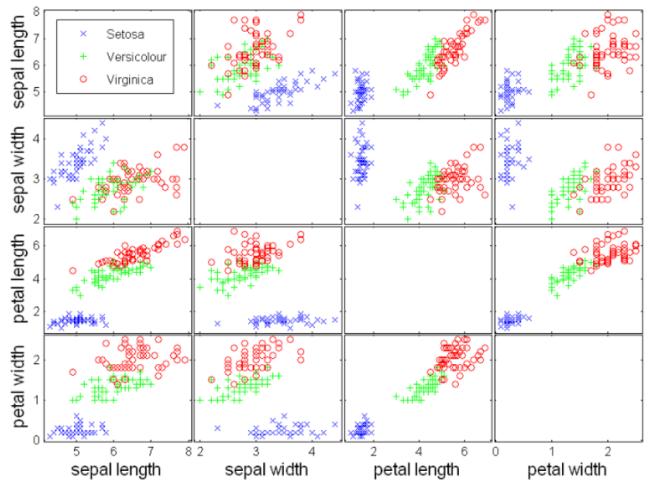
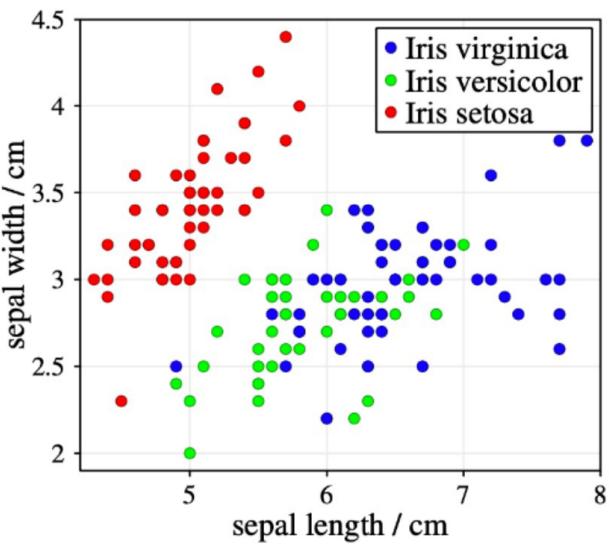


Figure 1.9: Scatter plot of sepal length and width, and scatter matrix: scatter plots of all pairs of attributes in the Iris dataset.

Plot bivariate (or trivariate) data, eyeing data correlation and outliers.

## 1.5 Data Preparation

We will delve into the following techniques of data preparation:

- ◊ Aggregation
- ◊ Data Reduction: Sampling
- ◊ Dimensionality Reduction

- ◊ Feature subset selection
- ◊ Feature creation
- ◊ Discretization and Binarization
- ◊ Attribute Transformation

### 1.5.1 Aggregation

Aggregation is the process of combining two or more attributes (or objects) into a single attribute (or object).

*Purpose*

- ◊ Data reduction
  - Reduce the number of attributes or objects
- ◊ Change of scale
  - Cities aggregated into regions, states, countries, etc.
  - Days aggregated into weeks, months, or years
- ◊ More “stable” data
  - Aggregated data tends to have less variability

### 1.5.2 Reduction

Reduction is simply reducing the amount of data. We may reduce the number of **records** by sampling or clustering, or the number of **attributes** (*columns*) by selecting a subset of them, or by creating a new —smaller— set of attributes from the old one.

#### 1.5.2.1 Sampling

Sampling is the main technique employed for data reduction.

It is often used for both the preliminary investigation of the data and the final data analysis.

Sampling is typically used in data mining because processing the entire set of data of interest is too expensive or time consuming.

The key principle for effective sampling is the following:

- ◊ Using a sample will work almost as well as using the entire data set, if the sample is representative
- ◊ A sample is representative if it has approximately the same properties (of interest) as the original set of data
- ◊ **Simple Random Sampling**
  - There is an *equal probability* of selecting any particular item
  - Sampling **without replacement**
    - As each item is selected, it is removed from the population
  - Sampling **with replacement**
    - Objects are not removed from the population as they are selected for the sample.
    - In sampling with replacement, the same object can be picked up more than once
- **Stratified sampling**
  - Split the data into several partitions; then draw random samples from each partition
  - Approximation of the percentage of each class
  - Suitable for distribution with peaks: each peak is a **layer**

#### 1.5.2.2 Dimensionality Reduction

This consists in reducing the number of attributes (or features) in the data. We want a selection of a subset of attributes that is as small as possible and sufficient for the data analysis.

- ◊ removing (more or less) irrelevant features
  - Contain no information that is useful for the data mining task at hand
  - Example: students' ID is often irrelevant to the task of predicting students' GPA
- ◊ removing redundant features
  - Duplicate much or all of the information contained in one or more other attributes
  - Example: purchase price of a product and the amount of sales tax paid

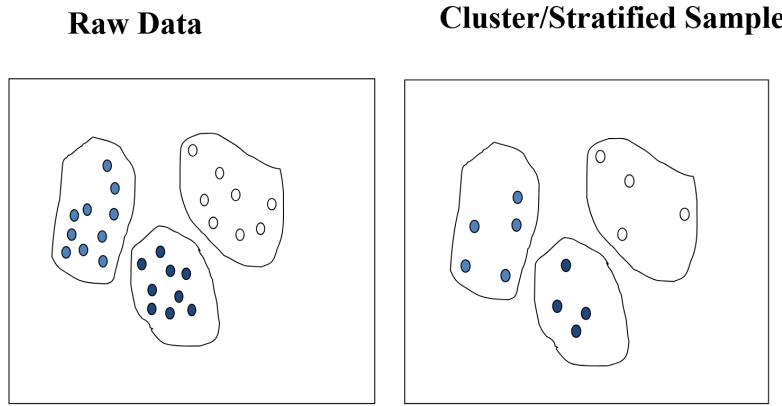


Figure 1.10: Stratified Sampling

### Curse of Dimensionality

When dimensionality increases, data becomes **increasingly sparse** in the space that it occupies.

Definitions of density and distance between points, which are critical for clustering and outlier detection, become less meaningful.

This phenomenon is known as the **curse of dimensionality**.

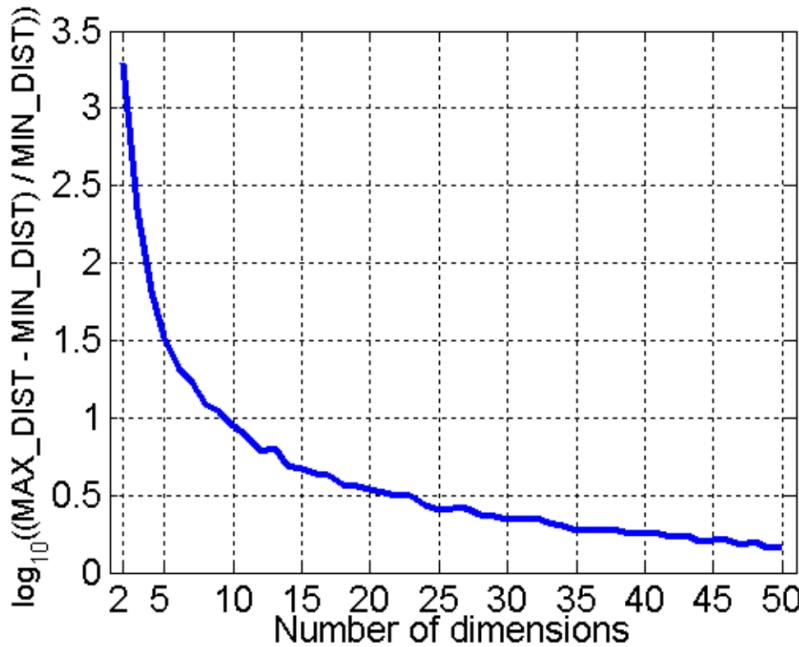


Figure 1.11:  $\log_{10}((\text{MAX\_DIST} - \text{MIN\_DIST}) / \text{MIN\_DIST})$  decreases as the dimensionality increases, meaning that the difference between the farthest and nearest neighbor distances becomes less significant

Purposes of dimensionality reduction include:

- ◊ Avoid curse of dimensionality
- ◊ Reduce amount of time and memory required by data mining algorithms
- ◊ Allow data to be more easily visualized
- ◊ May help to eliminate irrelevant features or reduce noise

Techniques to do so include:

- ◊ Principal Components Analysis (PCA)
- ◊ Singular Value Decomposition
- ◊ Others: supervised and non-linear techniques

### 1.5.2.3 Feature Subset Selection

Feature subset selection consists in selecting a subset of the original features. The goal is to find a minimal subset of features that is as good as the entire set of features for the data analysis task at hand.

For removing irrelevant features, it is needed a **performance measure** indicating how well a feature or subset of features performs w.r.t. the considered data analysis task.

For removing **redundant features**, either a *performance measure* for subsets of features or a *correlation measure* is needed.

#### Filter Methods

- ◊ Selection after analyzing the **significance** and **correlation** with other attributes
- ◊ Selection is independent of any data mining task
- ◊ The operation is a pre-processing

#### Wrapper Methods

- ◊ Selecting the top-ranked features using as reference a DM task
- ◊ Incremental Selection of the “best” attributes  
“Best” = with respect to a specific measure of statistical significance (e.g.: information gain)

#### Embedded Methods

- ◊ Selection as part of the data mining algorithm
- ◊ During the operation of the DM algorithm, the algorithm itself decides which attributes to use and which to ignore (e.g. Decision tree)

### Feature Selection Techniques

- ◊ **Selecting the top-ranked features:** Choose the features with the best evaluation when single features are evaluated.
- ◊ **Selecting the top-ranked subset:** Choose the subset of features with the best performance. This requires exhaustive search and is impossible for larger numbers of features. (For 20 features there are already more than one million possible subsets.)
- ◊ **Forward selection:** Start with the empty set of features and add features one by one. In each step, add the feature that yields the best improvement of the performance.
- ◊ **Backward elimination:** Start with the full set of features and remove features one by one. In each step, remove the feature that yields the least decrease in performance.

# Chapter 2

## Data Representation

- ◊ By **correlation**

I want to represent data according to the correlation of the dataset  
Algorithm: PCA

- ◊ By **neighborhood**

I want to represent the data so that similar instances are similar  
Algorithm: t-SNE

- ◊ By **manifold**

I want to represent the data so that its manifold is preserved  
Algorithm: UMAP

### 2.1 Principal Component Analysis (PCA)

PCA is a linear dimensionality reduction technique that transforms the data to a new coordinate system such that the greatest variance by any projection of the data comes to lie on the first coordinate (called the first principal component), the second greatest variance on the second coordinate, and so on.

Essentially, PCA exploits spectral decomposition of the whole dataset to find a new basis for the data.

Data can often be correlated, and linear dependencies can exist among variables, e.g.,

- ◊ Rent is linearly dependent on salary and food expenses
- ◊ Bank deposit is linearly dependent on salary and work
- ◊ Cardio is linearly dependent on  $VO_2max$

Vectors are  $m$ -dimensional elements in a field, and enjoy both addition and multiplication by scalar.

Composing these two, we can generate an infinite number of vectors: this is a **vector space**, and is defined by the basis vectors involved in the composition.

A matrix  $A$  defines a space...and thus a linear transformation!  $Av$  linearly combines the columns of with coefficients given by  $v$ .

Eigenpairs  $(\lambda, v)$  of a square matrix  $A$  are defined by the equation  $Av = \lambda v$ .

The eigenvectors  $v_1, \dots, v_m$  of a matrix  $A$  define the stretching of the space, and their eigenvalues  $\lambda_1 > \dots > \lambda_m$  define the stretching factor.

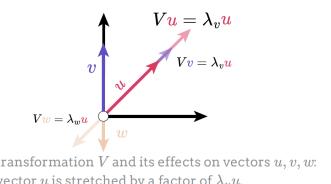
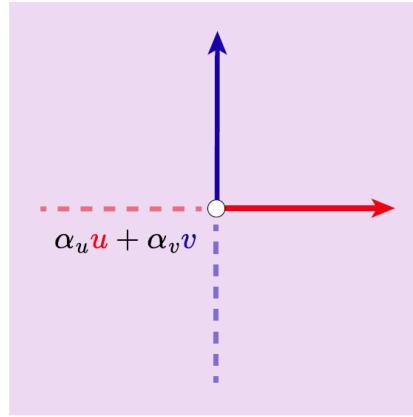


Figure 2.2: Eigenvectors of a matrix

PCA projects some data  $X$  to  $\hat{X}$  through a linear transformation  $A$ :  $AX = \hat{X}$ .

Fun fact #1: for a mean-centered  $\bar{X}$ , the slope is directly proportional to the covariance!



Two vectors  $u, v$  (in red and blue), and the plane spanned by all their linear combinations  $\alpha_u u + \alpha_v v$  (in purple).

Figure 2.1: Vector space spanned by two vectors

$$\bar{\Sigma} = \begin{bmatrix} \sigma_{\bar{X}^1}^2 & \cdots & \text{cov}(\bar{X}^1, \bar{X}^n) \\ \cdots & \cdots & \cdots \\ & & \sigma_{\bar{X}^n}^2 \end{bmatrix}$$

There is some pretty complicated linear algebra behind PCA, but the main steps are the following:

1. Mean-center your data  $X$  to get  $\bar{X}$
2. Compute its eigenvectors matrix  $\bar{V}$
3. Transpose  $V$  to obtain the transformation  $V^T$
4. Project the data:  $\bar{X}$  through  $V^T \bar{X}$ , obtaining the PCA-transformed data  $\hat{X}$

**In simple terms:** PCA looks at your data and finds the “most important directions” - imagine you have a cloud of points and you want to find the best line that captures the main trend. PCA finds not just one line, but multiple directions ordered by importance. It then rotates your data so that the first dimension captures the most variation, the second dimension captures the second most variation. This allows you to keep only the first two dimensions while retaining most of the information, effectively reducing the complexity of your data while preserving its essential structure.

### 2.1.1 Observations

- ◊ PCA redefines data by removing collinearity: if your data has low covariance, the transformation will have minimal effect.
  - ◊ PCA performs a linear transformation to tackle linear relationships between variables. Nonlinear relationships are not influenced.
  - ◊ PCA does not work very well for high complexity data.
- Uses*
- ◊ **Feature selection:** high covariance of a feature may indicate disposability.
  - ◊ **Dimensionality reduction:** trimming columns of lets us reduce the dimension of the resulting data.
  - ◊ **Clustering preprocessing:** correlated features inflate object similarity

## 2.2 t-SNE

t-SNE (t-distributed Stochastic Neighbor Embedding) is a nonlinear dimensionality reduction technique particularly well suited for embedding high-dimensional data for visualization in a low-dimensional space of two or three dimensions. It works by modeling the data as a distribution of points in a high-dimensional space and then finding a lower-dimensional representation that preserves the pairwise similarities between points.

t-SNE focuses on data clusters rather than subspace representation, and again maps the original data  $X$  to a representation  $\hat{X}$ .

t-SNE tackles this problem in two phases:

1. **Similarity phase** In the original space  $\mathcal{X}$ , how similar is  $x_i$  to  $x_j$ ?

How similar is  $x_i$  to  $x_j$ ? Even better, what is the probability that  $x_i$  is a neighbor of  $x_j$ ?

2. **Embedding phase** In the mapped space  $\hat{\mathcal{X}}$ , how similar is  $\hat{x}_i$  to  $\hat{x}_j$ ?

### 2.2.1 Similarity phase

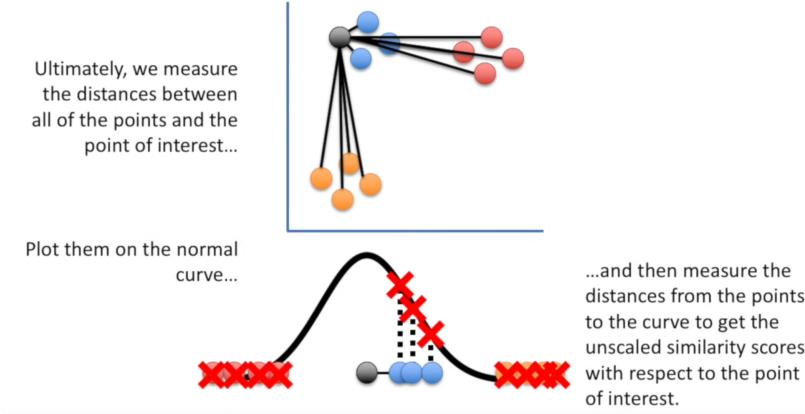


Figure 2.2: t-SNE plotting distance on an X axis and then projecting it on a normal distribution curve, to get the probability of being a neighbor

The similarity phase computes the similarity between points in the original high-dimensional space. This is typically done by converting the Euclidean distances between points into conditional probabilities that represent similarities. The probability that point  $x_j$  is a neighbor of point  $x_i$  is given by a Gaussian distribution centered at  $x_i$ . The variance of this Gaussian is controlled by a parameter called **perplexity**, which can be thought of as a smooth measure of the effective number of neighbors.

This yields a neighboring matrix  $P$  where each entry  $p_{ij}$  represents the probability that point  $x_j$  is a neighbor of point  $x_i$  in the original high-dimensional space.

The conditional probability is computed as:

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2/2\sigma_i^2)}$$

where  $\sigma_i$  is the variance of the Gaussian centered at point  $x_i$ . The perplexity parameter determines  $\sigma_i$  through a binary search to match the desired effective number of neighbors.

### 2.2.2 Embedding phase

In the embedding phase, t-SNE defines a similar probability distribution over the points in the low-dimensional map. However, instead of using a Gaussian distribution, it uses a Student's t-distribution with one degree of freedom (heavy-tailed distribution) to avoid the “crowding problem” where moderate distant points are forced to be too far apart in the low-dimensional representation.

The probability in the low-dimensional space is:

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}$$

where  $y_i$  and  $y_j$  are the low-dimensional counterparts of  $x_i$  and  $x_j$ .

### 2.2.3 Optimization

t-SNE minimizes the Kullback-Leibler divergence between the probability distributions  $P$  (high-dimensional) and  $Q$  (low-dimensional):

$$C = \sum_i KL(P_i || Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}$$

The algorithm uses gradient descent to find the low-dimensional representation  $Y$  that minimizes this cost function, effectively preserving the neighborhood structure of the original high-dimensional data.

#### Key advantages:

- ◊ Excellent for visualization of high-dimensional data
- ◊ Preserves local neighborhood structure
- ◊ Can reveal clusters and patterns not visible in linear methods

#### Key limitations:

- ◊ Computationally expensive (quadratic in the number of points)
- ◊ Non-deterministic (different runs can give different results)
- ◊ Sensitive to hyperparameters, especially perplexity
- ◊ Not suitable for embedding new data points (no explicit mapping function)

## 2.3 UMAP

UMAP (Uniform Manifold Approximation and Projection) is a nonlinear dimensionality reduction technique that is particularly effective for visualizing high-dimensional data in a low-dimensional space. It is based on manifold learning and topological data analysis, aiming to preserve both local and global structure of the data.

The computed distances induce a connectivity graph, and thus an adjacency matrix  $A$ , its edges measuring distances among instances. After turning distances into probabilities, UMAP optimizes a distance on  $A$ , to make it so that all and only the edges on the original manifold also appear in the transformed manifold with the same magnitude.

For the set of edges  $E$ , UMAP minimizes

$$-\sum_{e \in E} \left( \underbrace{\Pr(e; X) \log(\Pr(e; Z))}_{\text{existing edges}} + \underbrace{(1 - \Pr(e; Z)) \log(1 - \Pr(e; X))}_{\text{non-existing edges}} \right),$$

where  $\Pr(e; X), \Pr(e; Z)$  indicate the probability of edge  $e$  in the original and transformed representation, respectively.

# Chapter 3

## Data Cleaning

### Key Points

- ◊ How to handle **anomalous values**
- ◊ How to handle **outliers**
- ◊ **Data Transformations**

### 3.1 Anomalous Values

- ◊ **Missing** values - NULL, ?
- ◊ **Unknown** Values - Values without a real meaning
- ◊ **Not Valid** Values - Values not significant

We can handle anomalous values with various techniques. First of all we can **eliminate** the records with anomalous values.

We could also **substitute** anomalous values, knowing, however, that it could influence the original distribution of numerical values.

To mitigate this effect, we can use **mean/median/mode** to substitute missing values, or estimate missing values using the **probability distribution** of existing values.

We could also **segment data** and apply the above for every segment where there happen to be missing values, hence using mean/mode/median or the probability distribution of the segment.

Finally, we could use *predictive models* (**classification/regression**) to estimate missing values, using the other attributes as predictors.

#### 3.1.1 Discretization

*Discretization is the process of converting a **continuous** attribute into an **ordinal** attribute.*

- ◊ A potentially infinite number of values are mapped into a small number of categories
- ◊ Discretization is commonly used in classification
- ◊ Many classification algorithms work best if both the independent and dependent variables have only a few values

### Unsupervised Discretization

- ◊ No label for instances
- ◊ The number of classes is unknown

Techniques of binning:

- ◊ **Natural** binning - Intervals with the same width
- ◊ **Equal Frequency** binning - Intervals with the same frequency
- ◊ **Statistical** binning - Use statistical information (Mean, variance, Quartile)

### 3.1.1.1 Natural Binning

This is a fairly simple approach: we sort the values, subdividing the range of values into  $k$  parts with the same size.

$$\sigma = \frac{x_{max} - x_{min}}{k}$$

Then, element  $x_j$  is assigned to bin —i.e. belong to the class—  $i$  if:

$$x_j \in [x_{min} + i\sigma, x_{min} + (i + 1)\sigma)$$

Note however that this approach is sensitive to outliers, which can skew the range of values, generating very unbalanced distributions.

### 3.1.1.2 Equal Frequency Binning

Sort and count the elements, definition of  $k$  intervals of  $f$ , where, having  $N$  elements:

$$f = \frac{N}{k}$$

The element  $x_j$  is assigned to bin  $i$  if:

$$i \times f \leq j < (i + 1) \times f$$

This is not always suitable for highlighting interesting correlations.

### 3.1.1.3 How many bins?

In both cases, the number of bins  $k$  is a parameter to be set. A rule of thumb is to set  $k$  according to Sturges' optimal number of classes  $C$  for  $N$  elements:

$$C = 1 + \frac{10}{3} \log_{10}(N)$$

where  $N$  is the number of elements. The optimal width of the classes depends on the variance and the number of data points:

$$h = \frac{3.5 \cdot s}{\sqrt{N}}$$

## 3.2 Supervised discretization

### 3.2.1 Entropy-based Discretization

Minimizes the entropy wrt a label, with the goal of maximizing the purity of intervals (information gain).

## 3.3 Binarization

Binarization is the process of converting a continuous attribute into a binary attribute.

This can be useful in various scenarios, such as when we want to simplify the model or when we need to handle categorical variables. It is common to apply this for **association** analysis.

There are several techniques for binarization of continuous attributes:

- ◊ **Thresholding** - Assigning values above a certain threshold to one class and values below to another.
- ◊ **One-Hot Encoding** - Creating binary columns for each category in a categorical variable.
- ◊ **Binarization with Decision Trees** - Using decision tree algorithms to find optimal splits for binarization.

There are also techniques for binarization of categorical attributes:

Table 3.1: Conversion of a categorical attribute to three binary attributes.

Categorical Value	Integer Value	$x_1$	$x_2$	$x_3$
awful	0	0	0	0
poor	1	0	0	1
OK	2	0	1	0
good	3	0	1	1
great	4	1	0	0

This however leads to highlighting associations that are not really there, for example, in the table above,  $x_2$  and  $x_3$  may look correlated, but, in fact, they are not.

Table 3.2: Conversion of a categorical attribute to five asymmetric binary attributes.

Categorical Value	Integer Value	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$
awful	0	1	0	0	0	0
poor	1	0	1	0	0	0
OK	2	0	0	1	0	0
good	3	0	0	0	1	0
great	4	0	0	0	0	1

### 3.3.1 Attribute Transformation

An attribute transform is a function that maps the entire set of values of a given attribute to a new set of replacement values such that each old value can be identified with one of the new values.

- ◊ Simple functions:  $x^k$ ,  $\log(x)$ ,  $e^x$ ,  $|x|$
- ◊ **Normalization** Refers to various techniques to adjust to differences among attributes in terms of frequency of occurrence, mean, variance, range
- ◊ Take out unwanted, common signal, e.g., seasonality
- ◊ In statistics, **standardization** refers to subtracting off the means and dividing by the standard deviation

The transformation  $Y = T(X)$  should:

- ◊ preserve the relevant information of  $X$
- ◊ eliminates at least one of the problems of  $X$
- ◊ is more useful of  $X$

#### 3.3.1.1 Normalization

- ◊ Min-Max Normalization

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new\_max}_A - \text{new\_min}_A) + \text{new\_min}_A$$

- ◊ z-score Normalization

$$v' = \frac{v - \mu}{\sigma} = \frac{v - \text{mean}_A}{\text{stddev}_A}$$

- ◊ Decimal Scaling

$$v' = \frac{v}{10^j} \quad \text{where } j \text{ is the smallest integer such that } \max(|v'|) < 1$$

#### 3.3.1.2 Transformation functions

- ◊ Exponential Transformation

$$T_p(x) = \begin{cases} ax^p + b & (p \neq 0) \\ c \log x + d & (p = 0) \end{cases}$$

where  $a, b, c, d, p \in \mathbb{R}$

- Preserve the order
- Preserve some basic statistics
- They are continuous functions
- They are derivable
- They are specified by simple functions

◊ Logarithmic Transformation - Stabilizing Variance

$$T(x) = c \log x + d$$

- Applicable to positive values
- Makes homogenous the variance in log-normal distributions  
E.g.: normalize seasonal peaks

◊ ??? - Stabilizing Variance

$$T(x) = ax^p + b$$

- Square-root Transformation
  - $p = 1/c$ ,  $c$  integer number
  - To make homogenous the variance of particular distributions e.g., Poisson Distribution
- Reciprocal Transformation
  - $p < 0$
  - Suitable for analyzing time series, when the variance increases too much wrt the mean

# Chapter 4

## Cluster analysis

Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups.

The aim of clustering is to ease data **understanding** and **summarization**, i.e. reduce the size of data sets.

The following are *NOT* cluster analysis:

- ◊ Simple segmentation
  - Dividing students into different registration groups alphabetically, by last name
- ◊ Results of a query
  - Groupings are a result of an external specification
  - Clustering is a grouping of objects based on the data
- ◊ Supervised classification
  - Have class label information
- ◊ Association Analysis
  - Local vs. global connections

These ain't clustering essentially because there is no **similarity** measure involved, which instead is the core of clustering.

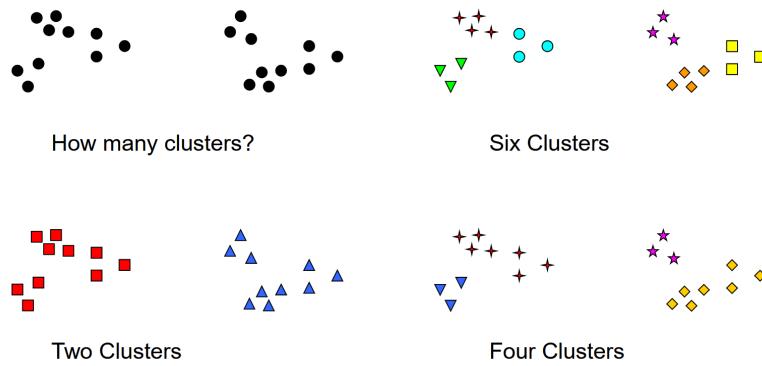


Figure 4.1: “Cluster” may be an ambiguous term. How to define a cluster? How big is a cluster? How many clusters are there?

### 4.1 Definitions

A **clustering** is a set of clusters. There is an important distinction between hierarchical and partitional sets of clusters.

- ◊ *Partitional Clustering*  
A division of data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset.
- ◊ *Hierarchical Clustering*

A set of nested clusters organized as a hierarchical tree.

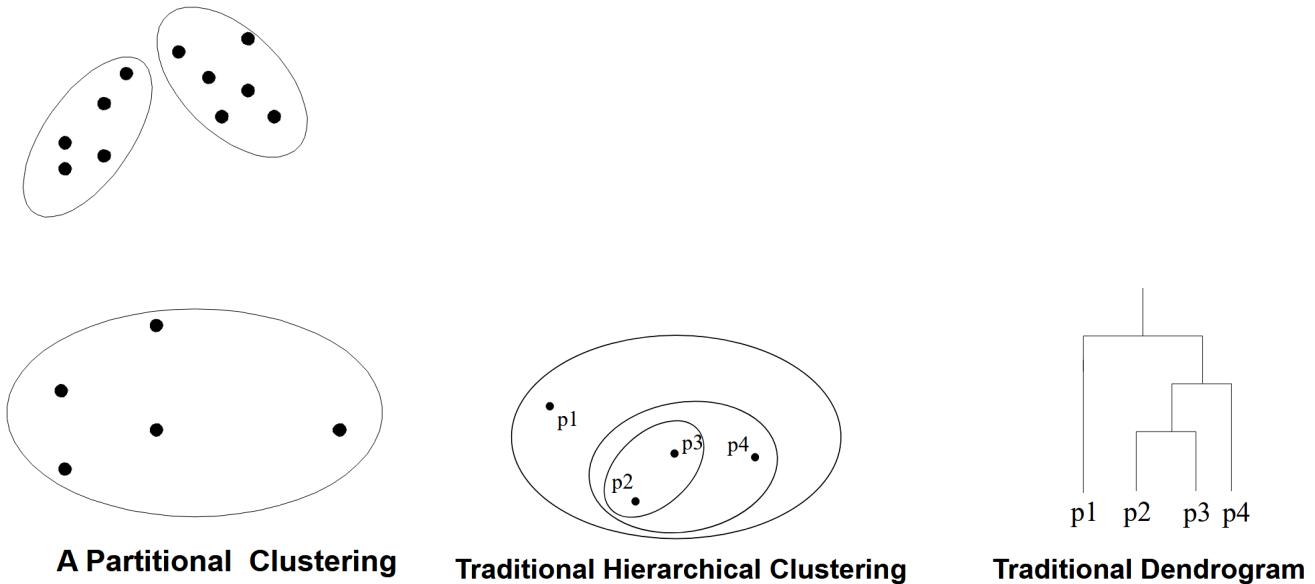


Figure 4.2: Partitional vs Hierarchical clustering

There are other distinctions among clusterings:

- ◊ Exclusive versus non-exclusive
  - In non-exclusive clusterings, points may belong to multiple clusters.
  - Can represent multiple classes or ‘border’ points
- ◊ Fuzzy versus non-fuzzy
  - In fuzzy clustering, a point belongs to every cluster with some weight between 0 and 1
  - Weights must sum to 1
  - Probabilistic clustering has similar characteristics
- ◊ Partial versus complete
  - In some cases, we only want to cluster some of the data
- ◊ Heterogeneous versus homogeneous
  - Clusters of widely different sizes, shapes, and densities

## 4.2 Types of Clustering

- ◊ Well-separated clusters A cluster is a set of points such that any point in a cluster is closer (or more similar) to every other point in the cluster than to any point not in the cluster
- ◊ Center-based clusters

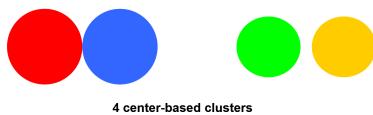


Figure 4.3: Center-based clusters

- ◊ Contiguous clusters (Nearest neighbor or Transitive)
  - Each point is closer to at least one point in its cluster than to any point in another cluster.
  - Graph based clustering
  - This approach can have trouble when noise is present since a small bridge of points can merge two distinct clusters
- ◊ Density-based clusters

- A cluster is a set of objects such that an object in a cluster is closer (more similar) to the “center” of a cluster, than to the center of any other cluster
- The center of a cluster is often a centroid, the average of all the points in the cluster, or a medoid, the most “representative” point of a cluster



Figure 4.3: Contiguity-based clusters

- A cluster is a dense region of points, which is separated by low-density regions, from other regions of high density.
  - Used when the clusters are irregular or intertwined, and when noise and outliers are present.
  - Points which are not classified as part of any cluster are considered noise. In the figure they may be the gray background points.
- ◊ Property or Conceptual
  - ◊ Described by an Objective Function
    - Finds clusters that minimize or maximize an objective function.
    - Enumerate all possible ways of dividing the points into clusters and evaluate the ‘goodness’ of each potential set of clusters by using the given objective function.
    - NP Hard)
    - Can have global or local objectives.
      - Hierarchical clustering algorithms typically have local objectives
      - Partitional algorithms typically have global objectives



Figure 4.4: Density-based clusters

## 4.3 Similarity

- ◊ Similarity
  - Numerical measure of how alike two data objects are.
  - Is higher when objects are more alike.
  - Often falls in the range [0,1]
- ◊ Dissimilarity
  - Numerical measure of how different are two data objects
  - Lower when objects are more alike
  - Minimum dissimilarity is often 0
  - Upper limit varies
- ◊ Proximity refers to a similarity or dissimilarity

### 4.3.1 Similarity and Dissimilarity for Different Attribute Types

Table 4.1: Similarity and dissimilarity for simple attributes

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$	$s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$
Ordinal	$d = \frac{ p-q }{n-1}$ (values mapped to integers 0 to $n - 1$ , where $n$ is the number of values)	$s = 1 - \frac{ p-q }{n-1}$
Interval or Ratio	$d =  p - q $	$s = -d, s = \frac{1}{1+d}$ or $s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

$p$  and  $q$  are the attribute values for two data objects.

### 4.3.2 Euclidean Distance

The most common way to measure distance between two data objects is the Euclidean distance:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

where  $n$  is the number of dimensions (attributes) and  $x_k$  and  $y_k$  are, respectively, the  $k^{th}$  attributes (components) or data objects  $\mathbf{x}$  and  $\mathbf{y}$ . Standardization is necessary, if scales differ.

- ◊ Standardization is necessary, if scales differ.

#### 4.3.2.1 Minkowski Distance

Minkowski Distance is a generalization of Euclidean Distance.  $r$  is a parameter that defines the type of distance,  $n$  is the number of dimensions (attributes) and  $x_k$  and  $y_k$  are, respectively, the  $k^{th}$  attributes (components) or data objects  $\mathbf{x}$  and  $\mathbf{y}$ .

$$d(\mathbf{x}, \mathbf{y}) = \left( \sum_{k=1}^n |x_k - y_k|^r \right)^{\frac{1}{r}}$$

- ◊ For  $r = 1$ , it is the Manhattan distance (or city block/Hamming distance)
- ◊ For  $r = 2$ , it is the Euclidean distance
- ◊ For  $r \rightarrow \infty$ , it is the supremum distance (or Chebyshev distance)

### Metrics and Similarities

1.  $d(x, y) \geq 0$  for all  $x$  and  $y$  and  $d(x, y) = 0$  only if  $x = y$ . (Positive definiteness)
2.  $d(x, y) = d(y, x)$  for all  $x$  and  $y$ . (Symmetry)
3.  $d(x, z) \leq d(x, y) + d(y, z)$  for all points  $x$ ,  $y$ , and  $z$  (Triangle Inequality).

A *distance*  $d$  is a *metric* if it satisfies the three conditions above.

1.  $s(x, y) = 1$  (or maximum similarity) only if  $x = y$ .
2.  $s(x, y) = s(y, x)$  for all  $x$  and  $y$ . (Symmetry)

These instead are properties of a *similarity* measure.

#### 4.3.3 Binary Similarity

Computing similarity among objects described by binary attributes is slightly different from numerical attributes.

- ◊ Simple Matching Coefficient (SMC)

$$SMC = \frac{\# \text{matches}}{\# \text{attributes}} = \frac{f_{11} + f_{00}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

- ◊ Jaccard Coefficient

$$J = \frac{\# \text{11matches}}{\# \text{non-zero attribute values}} = \frac{f_{11}}{f_{11} + f_{10} + f_{01}}$$

$$p = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0$$

$$q = 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 1$$

$$f_{01} = 2 \ (\# \text{attributes where } p \text{ was 0 and } q \text{ was 1})$$

$$f_{10} = 1 \ (\# \text{attributes where } p \text{ was 1 and } q \text{ was 0})$$

$$f_{00} = 7 \ (\# \text{attributes where } p \text{ was 0 and } q \text{ was 0})$$

$$f_{11} = 0 \ (\# \text{attributes where } p \text{ was 1 and } q \text{ was 1})$$

$$SMC = \frac{0 + 7}{10} = 0.7$$

$$J = \frac{0}{0 + 1 + 2} = 0$$

#### 4.3.4 Cosine Similarity

Used often in text mining, where each dimension corresponds to a term (word) and the value of the dimension is the frequency of the term in the document.

If  $d_1$  and  $d_2$  are two document vectors, then

$$\cos(d_1, d_2) = (d_1 \cdot d_2) / \|d_1\| \|d_2\|$$

where  $\cdot$  indicates vector dot product and  $\|d\|$  is the length of vector  $d$ .

Example:

$$d_1 = 3 \ 2 \ 0 \ 5 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0$$

$$d_2 = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 2$$

$$d_1 \cdot d_2 = 3 * 1 + 2 * 0 + 0 * 0 + 5 * 0 + 0 * 0 + 0 * 0 + 0 * 0 + 2 * 1 + 0 * 0 + 0 * 2 = 5$$

$$\|d_1\| = \sqrt{3^2 + 2^2 + 0^2 + 5^2 + 0^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2} = \sqrt{42}$$

$$\|d_2\| = \sqrt{1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 0^2 + 2^2} = \sqrt{6}$$

$$\cos(d_1, d_2) = \frac{5}{\sqrt{42} * \sqrt{6}} = \frac{5}{\sqrt{252}} \approx 0.315$$

#### 4.3.5 Correlation

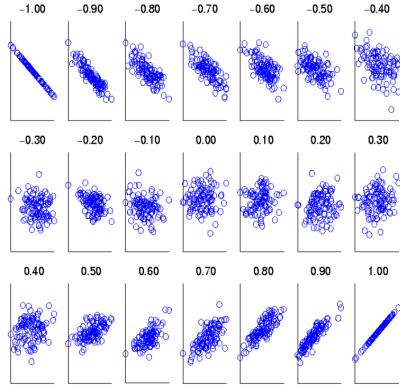


Figure 4.4: Scatter plots showing the similarity from -1 to 1

$$\text{corr}(\mathbf{x}, \mathbf{y}) = \frac{\text{covariance}(\mathbf{x}, \mathbf{y})}{\text{stddev}(\mathbf{x}) * \text{stddev}(\mathbf{y})} = \frac{\sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})}{\sqrt{\sum_{k=1}^n (x_k - \bar{x})^2} \sqrt{\sum_{k=1}^n (y_k - \bar{y})^2}}$$

### 4.3.6 Entropy

Information relates to possible outcomes of an event such as transmission of a message, flip of a coin, or measurement of a piece of data. The amount of information is inversely related to the probability of an event. Entropy is the commonly used measure of information content. For a discrete random variable  $X$  with  $n$  possible values  $x_1, x_2, \dots, x_n$  each having probability  $p_1, p_2, \dots, p_n$ , the entropy  $H(X)$  is defined as:

$$H(X) = - \sum_{i=1}^n p_i \log_2(p_i)$$

Entropy is measured in bits and is  $0 \leq H(X) \leq \log_2(n)$ . Thus, entropy is a measure of how many bits it takes to represent an observation of  $X$  on average.

The information one variable provides about another is called mutual information. The mutual information  $I(X; Y)$  of two discrete random variables  $X$  and  $Y$  is defined as:

$$I(X, Y) = H(X) + H(Y) - H(X, Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(x, y)$$

where  $H(X, Y)$  is the joint entropy of  $X$  and  $Y$ , and  $p(x, y)$  is the probability that  $x$  and  $y$  occur together.

## 4.4 K-Means

- ◊ *Partitional* clustering approach
  - ◊ The number of clusters  $K$ , must be specified
  - ◊ Each cluster is associated with a **centroid** (center point), typically being the mean of the points in the cluster
  - ◊ Each point is assigned to the cluster with the closest centroid. This may be measured as Euclidean distance, cosine similarity, correlation, etc.
- Using these measures makes K-Means to converge in the first few iterations.
- ◊ Complexity is  $O(nKId)$ , where  $n$  is the number of data points,  $K$  is the number of clusters,  $I$  is the number of iterations, and  $d$  is the number of dimensions (attributes)
  - ◊ The basic algorithm is very simple

---

#### Algorithm 1 K-Means Clustering Algorithm

- 1: Select  $K$  points as the initial centroids.
  - 2: **repeat**
  - 3:     Form  $K$  clusters by assigning all points to the closest centroid.
  - 4:     Recompute the centroid of each cluster.
  - 5: **until** The centroids don't change
- 

Correlation measures the linear relationship between objects (binary or continuous). To compute correlation, we standardize data objects,  $p$  and  $q$ , and then take their dot product (covariance/standard deviation).

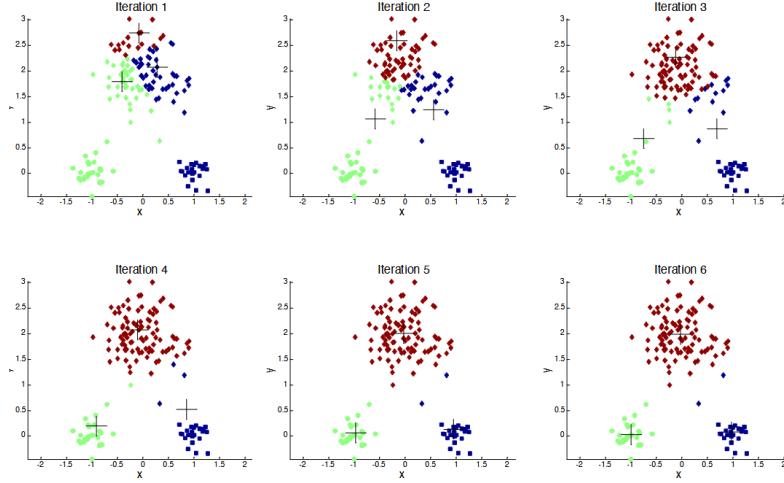


Figure 4.5: K-Means Clustering Iterations

#### 4.4.1 Evaluating K-Means clusters

The most common measure to evaluate the quality of K-Means clusters is the **Sum of Squared Errors** (SSE), also known as **Inertia**. It quantifies how tightly the data points in a cluster are grouped around their centroid. The SSE is calculated as follows:

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} \|x - \mu_i\|^2$$

- ◊  $K$  is the number of clusters,
- ◊  $C_i$  is the set of points in cluster  $i$ ,
- ◊  $x$  is a data point in cluster  $C_i$ ,
- ◊  $\mu_i$  is the centroid of cluster  $C_i$ ,
- ◊  $\|x - \mu_i\|^2$  is the squared Euclidean distance between point  $x$  and centroid  $\mu_i$ .

#### 4.4.2 Limitations of K-Means

K-means has problems when clusters are of differing

- ◊ Clusters have different sizes
- ◊ Clusters have different densities
- ◊ Clusters have Non-globular shapes
- ◊ The data contains **outliers**.

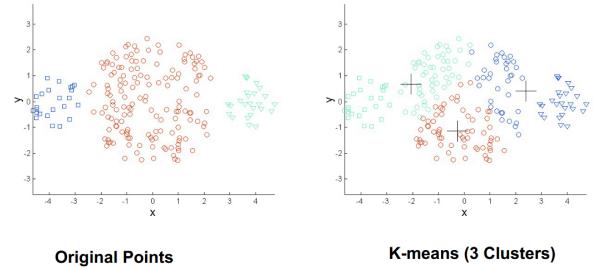


Figure 4.6: K-Means with different cluster sizes  
A solution may be to use more clusters, hence increasing  $K$ , and then merging them.

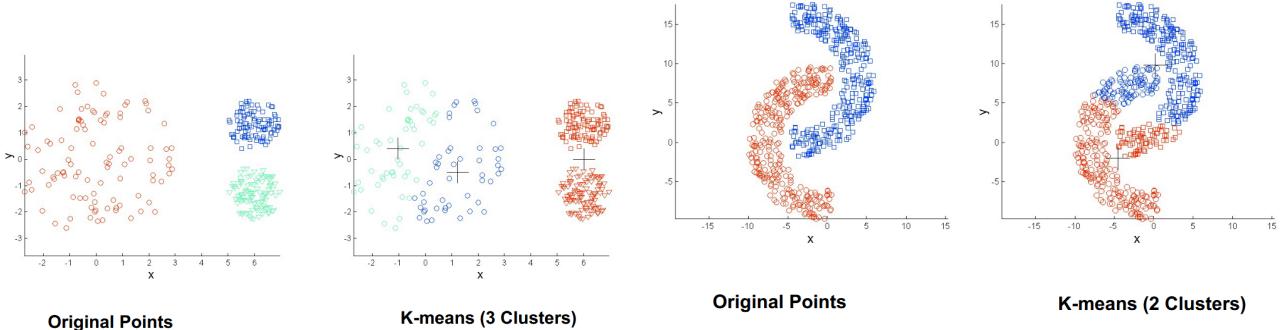


Figure 4.6: K-Means with different cluster densities and non-globular shapes

#### 4.4.2.1 Empty Clusters

K-Means sometimes produces empty clusters. This happens when no points are assigned to a cluster during the assignment step. This can occur if the initial centroids are poorly chosen or if the data distribution is such that certain centroids end up being too far from any data points.

Solutions include the following strategies, which may be iterated until no empty clusters remain:

- ◊ Choose a point and assign it to the cluster
  - Choose the point that contributes most to SSE
  - Choose a point from the cluster with the highest SSE

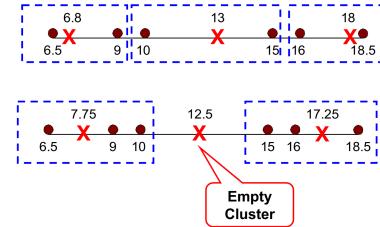


Figure 4.7: K-Means with Empty Clusters

#### 4.4.3 Workflow

1. Data Preprocessing
    - ◊ Handle missing values
    - ◊ Normalize/standardize data
  2. Choose the number of clusters  $K$
  3. Post processing and validation
    - ◊ Eliminate small clusters that may represent outliers
    - ◊ Split ‘loose’ clusters, i.e., clusters with relatively high SSE
    - ◊ Merge clusters that are ‘close’ and that have relatively low SSE
    - ◊ Can use these steps during the clustering process
- ISODATA

#### 4.4.4 Choosing the number of clusters K

If there are  $K$  “real” clusters then the chance of selecting one centroid from each cluster is small, especially if  $K$  is large.

A solution is to run K-Means multiple times with different initial centroids and choose the best result, but probability of getting a good result is still small.

A different approach is to use **hierarchical clustering** to determine initial centroids for K-Means.

Another option is to **incrementally** update centroids as points are assigned to clusters, instead of waiting until all points are assigned. This is more expensive and introduces order dependence.

We can define a goodness measure of a cluster  $c$  (hence of the set of Clusters  $C$ ) as the SSE from the cluster centroid:

$$SSE_C(c, s) = \sum_{i=1}^n (d_i, s_c)^2 \quad (4.1)$$

$$G(C, s) = \sum_{c \in C} SSE_C(c, s) \quad (4.2)$$

where  $s_c$  is the centroid of cluster  $c$ ,  $d_i$  is a point in cluster  $c$ , and  $C$  is the set of clusters.

Re-assignment of points to clusters and re-computation of centroids is guaranteed to *not increase* (or *monotonically decreases*)  $G(C, s)$ , hence K-Means converges to a local minimum.

At any step we have some value for  $G(C, s)$ ,

1. Fix  $s$ , optimize  $C \rightarrow$  Assignment step: assign  $d_i$  to the closest centroid  $\Rightarrow G(C', s) \leq G(C, s)$
2. Fix  $C'$ , optimize  $s \rightarrow$  Update step: recompute centroids  $\Rightarrow G(C', s') \leq G(C', s) \leq G(C, s)$

In this way the new cost results smaller than the original one, leading to convergence to a local minimum.

## 4.5 Hierarchical Clustering

Hierarchical clustering produces a set of nested clusters organized as a hierarchical tree can be visualized as a dendrogram, which is a tree-like diagram that records the sequences of merges or splits.

We do not have to assume any particular number of clusters, any desired number of clusters can be obtained by “cutting” the dendrogram at the proper level, each possibly corresponding to meaningful taxonomies.

### 4.5.1 Agglomerative vs Divisive

There are two types of hierarchical clustering, both exploiting a similarity or distance matrix:

- ◊ **Agglomerative** (bottom-up)

Start with the points as individual clusters, and, at each step, merge the closest pair of clusters until only one cluster (or  $k$  clusters) remain.

- ◊ **Divisive** (top-down)

Start with all points in one cluster and, at each step, split the cluster into smaller clusters until each point is its own cluster (or  $k$  clusters are formed).

#### Algorithm 2 Agglomerative Hierarchical Clustering Algorithm

- 1: Compute the proximity matrix
- 2: Let each data point be a cluster
- 3: **repeat**
- 4:     Merge the two closest clusters
- 5:     Update the proximity matrix
- 6: **until** only a single cluster remains

**Proximity** matrix is a square matrix that contains the distances (or similarities) between each pair of data points. “Proximity” is just a general term that can refer to either distance or similarity, depending on the context.

#### 4.5.1.1 Updating Proximity Matrix

When two clusters  $c_i$  and  $c_j$  are merged into a new cluster  $c_k$ , we need to update the proximity matrix to reflect the distances between the new cluster  $c_k$  and all other existing clusters. There are several methods to do this, each defining the distance between clusters differently:

- ◊ **Single Linkage** (or Nearest Neighbor) The distance between two clusters is defined as the minimum distance between any single point in the first cluster and any single point in the second cluster.

$$d(c_i, c_j) = \min_{x \in c_i, y \in c_j} d(x, y)$$

- ◊ **Complete Linkage** (or Farthest Neighbor) The distance between two clusters is defined as the maximum distance between any single point in the first cluster and any single point in the second cluster.

$$d(c_i, c_j) = \max_{x \in c_i, y \in c_j} d(x, y)$$

- ◊ **Average Linkage** The distance between two clusters is defined as the average distance between all pairs of points, one from each cluster.

$$d(c_i, c_j) = \frac{1}{|c_i||c_j|} \sum_{x \in c_i} \sum_{y \in c_j} d(x, y)$$

Note that **Single**, **Complete** and **Average** Linkage have all the same complexity!  $O(n^3)$ , where  $n$  is the number of data points.

- ◊ **Centroid Linkage** The distance between two clusters is defined as the distance between their centroids.

$$d(c_i, c_j) = d(\mu_i, \mu_j)$$

where  $\mu_i$  and  $\mu_j$  are the centroids of clusters  $c_i$  and  $c_j$ , respectively.

- ◊ **Ward's Method** This method minimizes the total within-cluster variance. When two clusters are merged, the increase in the total within-cluster variance is calculated, and the pair of clusters that results in the smallest increase is merged.

$$d(c_i, c_j) = \frac{|c_i||c_j|}{|c_i| + |c_j|} \|\mu_i - \mu_j\|^2$$

where  $|c_i|$  and  $|c_j|$  are the sizes of clusters  $c_i$  and  $c_j$ , and  $\mu_i$  and  $\mu_j$  are their centroids.

### 4.5.2 Divisive Hierarchical Clustering

First of all, build MST (Minimum Spanning Tree) of the data points, then iteratively remove the longest edge in the MST to split the data into clusters.

Start with a single cluster which consists of any point, and at each step, look for the closest pair of points  $p, q$  such that  $p$  is in the current tree while  $q$  is not. Add  $q$  to the tree and the edge  $(p, q)$  to the MST. Repeat until all points are in the tree.

---

#### Algorithm 3 Divisive Hierarchical Clustering Algorithm

---

- 1: Compute a MST for the proximity graph
  - 2: **repeat**
  - 3:     Create a new cluster by breaking the link corresponding to the longest edge in the MST (smallest proximity)
  - 4: **until** only singleton clusters remain
- 

### 4.5.3 Complexity and Limitations

- ◊ Time Complexity -  $O(n^3)$  for naive implementations of hierarchical clustering
  - $O(n^2)$  to compute the proximity matrix
  - $O(n)$  to find the closest pair of clusters
  - $O(n)$  to update the proximity matrix
  - Repeated  $n$  times
- More efficient algorithms can achieve  $O(n^2 \log n)$ .
- ◊ Space Complexity -  $O(n^2)$  for storing the proximity matrix

#### 4.5.3.1 Limitations

Once a decision is made to combine two clusters, it cannot be undone.

No global objective function is directly minimized.

Different schemes have problems with one or more of the following:

- ◊ Sensitivity to noise and outliers
- ◊ Difficulty handling clusters of different sizes and non- globular shapes
- ◊ Breaking large clusters

## 4.6 Density based - DBSCAN

**Definition 4.1 (Density)** *Density* is the number of points within a specified radius (*EPS* or  $\epsilon$ ).

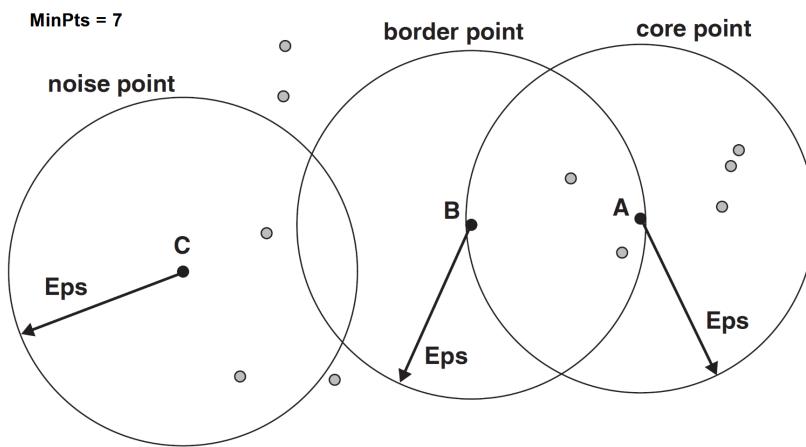


Figure 4.7: Core, Border, and Noise Points in DBSCAN

A point is a **core point** if it has at least *MINPTS* points within distance *EPS* (including itself). These are points which are at the interior of a cluster.

A point is a **border point** if it has fewer than *MINPTS* points within distance *EPS*, but is in the neighborhood of a core point. These are points which are on the edge of a cluster.

A point is a **noise point** (or outlier) if it is neither a core point nor a border point. These are points which are not part of any cluster.

- DBSCAN**
1. Label points as core, border, or noise based on thresholds  $R$  (radius of neighborhood) and `min_pts` (min number of neighbors)
  2. Connect core points that are within  $R$  of each other, hence are neighbors, putting them in the same cluster
  3. Associate border points to the nearest core point, hence to the same cluster as the core point, and remove noise.

**Algorithm 4** DBSCAN Algorithm

```

1: current_cluster_label  $\leftarrow 0$ 
2: for all core points do
3:   if the core point has no cluster label then
4:     current_cluster_label  $\leftarrow$  current_cluster_label + 1
5:     Label the current core point with cluster label current_cluster_label
6:   for all points in the  $Eps$ -neighborhood, except  $i^{th}$  the point itself do
7:     if the point does not have a cluster label then
8:       Label the point with cluster label current_cluster_label

```

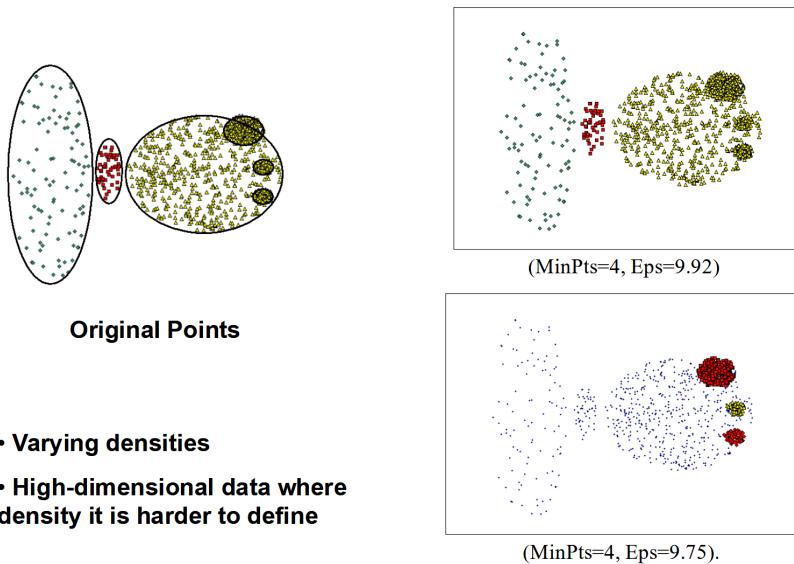


Figure 4.8: DBSCAN Not behaving well in given situations

**Determining Eps and MinPts**

To determine EPS and MINPTS, we can use the k-distance graph.

The idea is that for points in a cluster, their  $k^{th}$  nearest neighbors are at roughly the same distance. Noise points have the  $k^{th}$  nearest neighbor at distance.

So, plot sorted distance of every point to its  $k^{th}$  nearest neighbor.

## 4.7 Cluster Validity

For supervised classification we have a variety of measures to evaluate how good our model is. such as accuracy, precision, recall, etc...

*“ For cluster analysis, the analogous question is how to evaluate the “goodness” of the resulting clusters? ” —*

“Clusters are in the eye of the beholder”, but we still want to evaluate them to avoid finding patterns in noise, compare clustering algorithms, or to compare clusters or sets of clusters.

## 4.8 Towards cluster validation

1. Determining the **clustering tendency** of a set of data, i.e., distinguishing whether non-random structure actually exists in the data.
2. Comparing the results of a cluster analysis to externally known results, e.g., to externally given class labels.
3. Evaluating how well the results of a cluster analysis fit the data without reference to external information.
4. Comparing the results of two different sets of cluster analyses to determine which is better
5. Determining the “correct” number of clusters.

Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following three types:

- ◊ **External Index:** Used to measure the extent to which cluster labels match externally supplied class labels.  
Entropy
- ◊ **Internal Index:** Used to measure the goodness of a clustering structure without respect to external information.  
Sum of Squared Error (SSE)
- ◊ **Relative Index:** Used to compare two different clusterings or clusters.  
Often an external or internal index is used for this function, e.g., SSE or entropy

Sometimes these are referred to as *criteria* instead of *indices*; however, sometimes criterion is the general strategy and index is the numerical measure that implements the criterion.

### 4.8.1 Measuring validity through correlation

Two matrices:

- ◊ Proximity matrix
- ◊ Ideal similarity matrix
  - One row and one column for each data point
  - An entry is 1 if the associated pair of points belong to the same cluster
  - An entry is 0 if the associated pair of points belongs to different clusters

Compute the correlation between the two matrices; note that since the matrices are symmetric, only the correlation between  $n(n - 1)/2$  entries needs to be calculated.

### 4.8.2 Internal measures

#### 4.8.2.1 SSE - Sum of Squared Error

SSE (Sum of Squared Error) is a common measure of cluster cohesion, defined as:

$$SSE = \sum_{i=1}^k \sum_{x \in C_i} dist^2(m_i, x)$$

where  $m_i$  is the centroid of cluster  $C_i$ .

This is a measure of how tightly the data points in a cluster are grouped together. Lower values of SSE indicate more cohesive clusters. It can be a nice measure which does not need external information, and may also be used to determine the appropriate number of clusters by plotting the SSE for different values of  $k$  and looking for an “elbow” in the graph.

- Correlation of ideal similarity and proximity matrices for the K-means clusterings of the following two data sets.

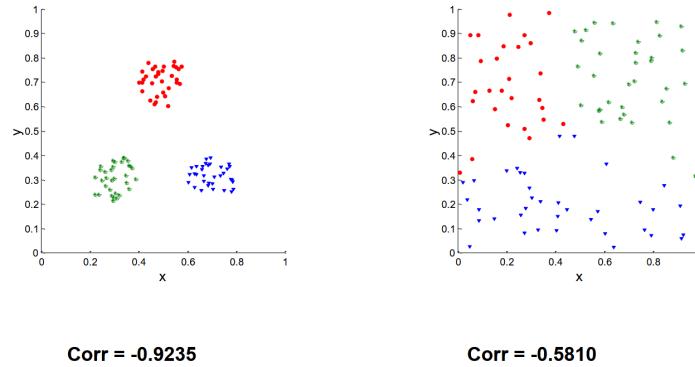


Figure 4.9: Correlation of two sets. Actually, in real world data it's very rare to observe a  $-0.9235$  correlation score

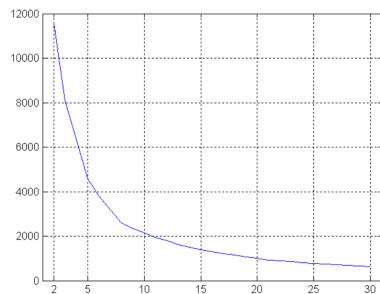


Figure 4.10: SSE of cluster found using K-means with various values of  $k$

#### 4.8.2.2 Cohesion and Separation

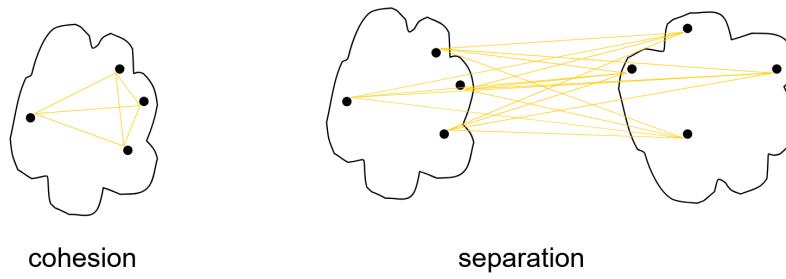


Figure 4.11: Cohesion and separation

A proximity graph based approach can also be used for cohesion and separation.

- ◊ Cluster **cohesion** is the sum of the weight of all links within a cluster, it measures how closely related are objects in a cluster. (SSE is a measure of cohesion)

$$SSE = WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

- ◊ Cluster **separation** is the sum of the weights between nodes in the cluster and nodes outside the cluster. Measures how distinct or well-separated a cluster is from other clusters. (Squared error between clusters is a measure of separation)

$$BSS = \sum_i |C_i|(m - m_i)^2 \quad |C_i| \text{ size of cluster } C_i$$

Note that  $BSS + WSS = TSS$  where  $TSS$  is the total sum of squares, and is constant and independent of the clustering.

#### 4.8.2.3 Silhouette coefficient

Silhouette coefficient combines ideas of both cohesion and separation, but for individual points, as well as clusters and clusterings.

Consider a point  $i$  in cluster  $A$ :

- ◊ Calculate  $a = \text{average distance of } i \text{ to the points in its cluster}$
- ◊ Calculate  $b = \min(\text{average distance of } i \text{ to points in another cluster})$
- ◊ The silhouette coefficient for a point is then given by

$$s = \frac{b - a}{\max(a, b)}$$

- ◊ Typically between 0 and 1. The closer to 1 the better.
- ◊ Can calculate the average silhouette coefficient for a cluster or a clustering.

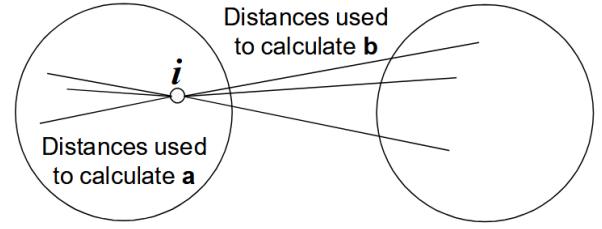


Figure 4.12: Silhouette coefficient for  $i$

#### 4.8.3 External measures

External measures compare the clustering results to an externally known class labels. Examples may be **entropy** and **purity**.

Table 4.2: K-means Clustering Results for LA Document Data Set

Cluster	Entertainment	Financial	Foreign	Metro	National	Sports	Entropy	Purity
1	3	5	40	506	96	27	1.2270	0.7474
2	4	7	280	29	39	2	1.1472	0.7756
3	1	1	1	7	4	671	0.1813	0.9796
4	10	162	3	119	73	2	1.7487	0.4390
5	331	22	5	70	13	23	1.3976	0.7134
6	5	358	12	212	48	13	1.5523	0.5525
<b>Total</b>	<b>354</b>	<b>555</b>	<b>341</b>	<b>943</b>	<b>273</b>	<b>738</b>	<b>1.1450</b>	<b>0.7203</b>

##### 4.8.3.1 Entropy

Entropy measures how the various clusters are distributed with respect to the class labels.

For each cluster, the class distribution of the data is calculated first, i.e., for cluster  $j$  we compute  $p_{ij}$ , the “probability” that a member of cluster  $j$  belongs to class  $i$  as follows:  $p_{ij} = m_{ij}/m_j$ , where  $m_j$  is the number of values in cluster  $j$  and  $m_{ij}$  is the number of values of class  $i$  in cluster  $j$ .

Then using this class distribution, the entropy of *each* cluster  $j$  is calculated using the standard formula:

$$e_j = - \sum_{i=1}^L p_{ij} \log_2 p_{ij}$$

where  $L$  is the number of classes.

The total entropy for a set of clusters is calculated as the sum of the entropies of each cluster weighted by the size of each cluster, i.e.:

$$e = \sum_{i=1}^K \frac{m_i}{m} e_j$$

where  $m_j$  is the size of cluster  $j$ ,  $K$  is the number of clusters, and  $m$  is the total number of data points.

#### 4.8.3.2 Purity

Using the terminology derived for entropy, the purity of cluster  $j$ , is given by:

$$purity_j = \max_i p_{ij}$$

and the overall purity of a clustering by:

$$purity = \sum_{i=1}^K \frac{m_i}{m} purity_j$$



# Chapter 5

## Anomaly Detection

We may refer to anomalies in data also as *outliers*.

### 5.1 Outliers

- ◊ Inherently **fuzzy** - An instance has a degree of outlierness, which we can threshold to decide whether an instance is an outlier or not.
- ◊ **Data-dependent** - Outlier are exceptions to the data. But outliers themselves define the data...?
- ◊ **Not noise** - Noise is random, outliers are exceptional.
- ◊ **Mono/multi-dimensional** - An outlier can be so on one just one dimension, or on multiple.

Outliers are something that is either **unusual** or **extreme**, or both.

Outliers are, by nature, defined in terms of other instances. Whatever approach we use to detect them, we should take into account that they influence it as well.

### 5.2 Outlier Detection Algorithms

Algorithms used to detect outliers usually involve two key steps:

1. **Grading** - Define a grading function  $\tilde{o}$  that assigns to each instance  $x$  a degree of outlierness/anomaly  $\tilde{o}(x)$ .
2. **Thresholding** - Decide on a threshold  $\hat{o}$  such that instances with  $\tilde{o}(x) > \hat{o}$  are considered outliers.

We can categorize outlier detection algorithms depending on the *axis* on which they operate:

- ◊ **Locality** - is the outlier detection performed in a local neighborhood (local) or considering the entire dataset (global)?
- ◊ **Sensitivity** - is the outlier detection sensitive to the presence of other outliers (sensitive) or not (robust)?
- ◊ **Interpretability** - can we interpret/explain why an instance is considered an outlier (interpretable) or not (black-box)?

#### 5.2.1 Distributions

Locality	Global/Local
Sensitive	Robust/Sensitive
Interpretable	Black-box/Interpretable

In case of a normal distribution  $\mathcal{N}(\mu, \sigma)$ , we can define outliers as instances that are more than  $k$  standard deviations away from the mean  $\mu$  (see Figure 5.2).

This approach is **global**, **robust** and **interpretable**.

Locality	Global
Sensitive	Outliers influence the distribution but may be removed by Grubb's test
Interpretable	Black-box, no clear explanation for outlierness, simply, there are not many similar instances.

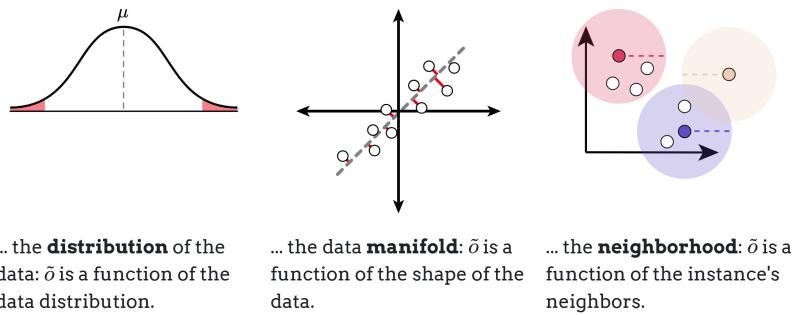
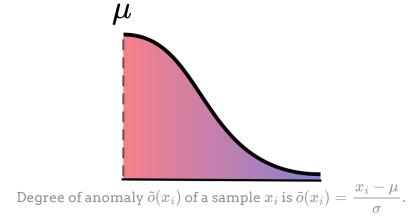


Figure 5.1: Distribution, manifold, neighborhood

The degree of anomaly  $\tilde{o}(x)$  can be defined as:

$$\text{z-score} = \tilde{o}(x) = \frac{|x - \mu|}{\sigma}$$

Figure 5.2: For a normal distribution  $\mathcal{N}(\mu, \sigma)$ , we can define outliers as instances that are more than  $k$  standard deviations away from the mean  $\mu$ .

### 5.2.1.1 Grubbs test

*z-scores* generate sample-dependent outlier degrees  $\tilde{o}(x_1), \tilde{o}(x_2), \dots, \tilde{o}(x_n)$ , but does not tackle the **+1 problem**.

The **+1 problem** is the problem of deciding whether the most extreme instance in a dataset is an outlier or not. Grubb's test iterates over detected outliers, removing one layer of outliers at a time, until no more outliers are found.

---

#### Algorithm 5 Grubbs's test for outliers

---

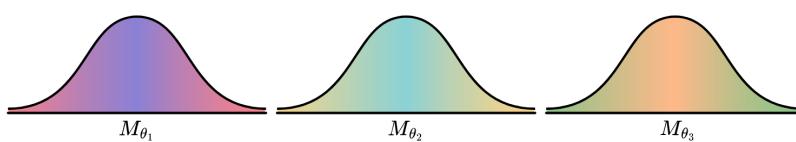
- 1: Find current outlier set  $\hat{X}$
  - 2: If  $\hat{X} = \emptyset$ , stop
  - 3:  $X = X \setminus \hat{X}$  - remove outliers from dataset
  - 4: Go to step 1
- 

Data may vary *locally*: subsets of the data each follow a different distribution.

Assumption: there exists a partition of the data, each block distributed according to a Normal distribution.

Thus we could use multiple models to model the data, and define outliers as instances that are outliers with respect to their local distribution.

One of  $k$  models  $M_{\theta_0}, M_{\theta_1}, \dots, M_{\theta_{k-1}}$  is sampled, each with a sampling probability  $m_i$ . Different distributions sample different regions of the density.

Figure 5.2: A mixture of Normals  $M_{\theta_0}, M_{\theta_1}, M_{\theta_2}$ ; each sampled with probability  $m_0, m_1, m_2$ , respectively.

Locality	Local
Sensitive	Outliers influence the distribution, may be unstable, but may be removed by Grubb's test (?)
Interpretable	<b>Black-box</b> , no clear explanation for outliers, simply, there are not many similar instances.

### 5.2.2 Thresholding

#### Grubb's test

Choosing a threshold  $\hat{o}$  is arbitrary, but there are algorithms such as Grubb's test which define their own thresholding mechanism. In Grubb's test, the threshold is defined as:

$$\hat{o} = n \frac{\sum(x_i - \bar{X})^4}{(\sum(x_i - \bar{X})^2)^2}$$

*z-scores* assume a Normal distribution, but often this is not the case. Yet, we can still identify tails of a distribution, and in turn, anomalies.

- ◊ Markov inequality - for a variable (distribution)  $X$  with positive values and threshold  $\beta$ , it holds

$$P(X \geq \beta) \leq \frac{\mathbb{E}(X)}{\beta}$$

Thus, given an estimate of the variable's expected value, we can retrieve the inverse of an image of its cumulative distribution

- ◊ Chebyshev inequality - for a variable (distribution)  $X$  with mean  $\mu$  ( $= \mathbb{E}(X)$ ) and standard deviation  $\sigma$ , and threshold  $\beta$ , it holds

$$P(|X - \mathbb{E}(X)| > \beta) \leq \frac{\sigma^2}{\beta^2}$$

That is, the probability of deviation from the mean is inversely proportional to the deviation

### 5.2.3 Manifold

Distributional approaches define the density, but do not describe the data itself.  $\tilde{o}$  is defined in terms of the manifold: does the given instance lie in the manifold? Just like the distributional approach, we must assume the manifold family. To preserve the interpretability of our results, we initially stick to linear manifolds.

We can define the anomaly degree  $\tilde{o}(x)$  as the distance of  $x$  from the manifold.

We can use PCA to find the linear manifold that best fits the data.

PCA finds the directions of maximum variance in the data, and uses them as a new basis for the data.

A matrix  $A$  spans a linear space, thus every vector  $b$  in its spanned space is defined as a linear combination of  $A$ :  $b = Ax$ . For non fullrank matrices  $A$ , such a solution  $x$  may not exist. Thus, we need to project on the data manifold.

#### Least Squares

Least Squares is a method to find the best fitting solution to an overdetermined system of equations  $Ax = b$  (more equations than unknowns).

The best fitting solution is the one that minimizes the residual sum of squares (RSS):

$$RSS = \|Ax - b\|_2^2$$

Least squares assumes a linear manifold, and squared norm as distance metric.

The instability of least squares is due to the data collinearity. A possible solution: de-correlate the data! PCA does exactly that.

Some mathy examples are displayed in the lecture slides...

Manifold-based algorithms are as flexible as the defined manifold. Like with mixture models, neighbor-based approaches reintroduce locality: outliers are defined in function of their neighbors:

Locality	Global
Sensitive	Strongly influenced by outliers
Interpretable	Partial: which instances have lower degrees? What even is a “low” degree?

Table 5.1: Least squares

- ◊ **Connectivity** - An outlier is defined in terms of the connectivity to its neighbors
- ◊ **Concentration** - An outlier is defined in terms of its neighbor concentration

Each instance has a posting list of neighbors, from the closest to the farthest: the lower the aggregated position in other lists, the higher the connectivity degree.

- ◊ Posting position defines connectivity: it is not density
- ◊ Connectivity is asymmetric: I may be your closest instance, you may not be mine

### 5.2.3.1 Grading neighbors connectivity

Posting matrices are often used as a base on which to measure different indices of connectivity, e.g.,

- ◊ **hub** - instance  $x_i$  is at least the  $t^{th}$  neighbor of at least  $k$  instances.

Definition used by ODIN: given a posting matrix  $A$ ,  $x_i$  is a *hub* if it appears at least  $k$  times in the first  $t$  columns of  $A$ . Hence,  $x_i$  is an outlier if the opposite is true:

$$\tilde{o}(x_i) = \begin{cases} 1 & \text{if } |i| i \in A_{\neq i, \leq t} | < k \\ 0 & \text{otherwise} \end{cases}$$

- ◊ **popularity** - instance  $x_i$  is on average the  $t^{th}$  neighbor of at least  $k$  instances

Given a posting matrix  $A$ ,  $x_i$  is an outlier if, on average, is not less than the  $t - th$  neighbor of other instances:

$$\hat{o}(x_i) = \frac{\sum_{l=0, l \neq i}^{n-1} \sum_{j=0}^{n-1} \mathbb{1}\{a_{l,j} = x_i\} l}{n-1} > t.$$

where  $\mathbb{1}\{a_{l,j} = x_i\}$  is an indicator function that equals 1 if  $a_{l,j} = x_i$  and 0 otherwise, representing the position in the posting list.

- ◊ **ostracism** - instance  $x_i$  is at worst  $t^{th}$  neighbor of other  $k$  instances

$$A = \begin{bmatrix} 0 & 2.28 & 0.16 & 0.21 \\ 2.21 & 0 & 1.21 & 3.91 \\ 0.16 & 1.21 & 0 & 0.76 \\ 0.21 & 3.91 & 0.76 & 0 \end{bmatrix}$$

Connectivity and concentration can be approximated through similar structures: we go from *postings* matrix to distance matrix! To ease notation, we use a row-sorted distance matrix  $A_\gamma$ , so that row  $i$  holds increasing distances from instance  $x_i$ .

$$A_\gamma = \begin{bmatrix} 0.16 & 0.21 & 2.28 \\ 1.21 & 2.28 & 3.91 \\ 0.16 & 0.76 & 1.21 \\ 0.21 & 0.76 & 3.91 \end{bmatrix}$$

A distance matrix  $A$  (top), and its row-sorted version  $A_\gamma$  (bottom). First column of 0s trimmed from  $A_\gamma$ .

### 5.2.4 Reach

An instance  $x$  has reach  $\gamma^k(x)$  if the  $k - th$  nearest neighbor is at distance  $\gamma^k$ , and average reach  $\bar{\gamma}^k(x)$  if the average of  $\{\gamma^1, \dots, \gamma^k\}$  is  $\bar{\gamma}^k(x)$ .

Our row-sorted distance matrix  $A_\gamma$  is the *reach* matrix of the data! Indeed,  $A_\gamma$  defines both reach and average reach.

The reach of instance  $x_i$  is encoded in the row-sorted distance matrix  $A_\gamma$  as follows:

$$A_\gamma = \begin{bmatrix} \gamma^1(x_1) & \gamma^2(x_1) & \gamma^3(x_1) \\ \gamma^1(x_2) & \gamma^2(x_2) & \gamma^3(x_2) \\ \gamma^1(x_3) & \gamma^2(x_3) & \gamma^3(x_3) \end{bmatrix}$$

where  $\gamma^j(x_i)$  denotes the  $j$ -th nearest neighbor of instance  $x_i$ .

For example, consider the following transformation:

$$A_\gamma \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ 0 & \frac{1}{2} & \frac{1}{3} \\ 0 & 0 & \frac{1}{3} \end{bmatrix} = \begin{bmatrix} \bar{\gamma}^1(x_1) & \bar{\gamma}^2(x_1) & \bar{\gamma}^3(x_1) \\ \bar{\gamma}^1(x_2) & \bar{\gamma}^2(x_2) & \bar{\gamma}^3(x_2) \\ \bar{\gamma}^1(x_3) & \bar{\gamma}^2(x_3) & \bar{\gamma}^3(x_3) \end{bmatrix}$$

$A_\gamma$  explicitly encodes reach ( $A_\gamma$  itself) and average reach.

#### 5.2.4.1 Reach ratio factor

Assumption: Inliers have lower reach than their neighbors. We formalize this in a reach (ratio factor):

$$\tilde{o}_{i,j}^k = \frac{\bar{\gamma}^k(x_i)}{\bar{\gamma}^k(x_j)}$$

which 1 is for pairs  $x_i, x_j$  with equal k-neighbors concentration, and  $> 1$  for instances with different concentrations,  $x_i$  laying in a sparser area of the space.

**Local outlier factor** generalizes outlier factor by averaging the outlier factor over the neighbors of an instance:

$$\tilde{o}(x_i) = \sum_{x_j \in \text{neigh}(x_i)} \tilde{o}_{i,j}^k$$

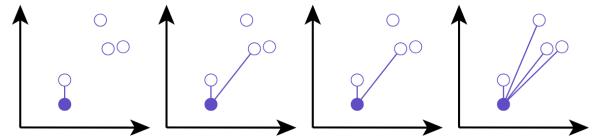


Figure 5.3: Neighbors at different  $k$  values.

Local outlier factor respects the posting matrix, as it creates *clusters* of neighbors.

**Connectivity outlier factor** is based on the idea that outliers are poorly connected to their neighbors. It is defined as:

$$\tilde{o}(x_i) = \sum_{x_j \in \text{connect\_neigh}(x_i)} \tilde{o}_{i,j}^k$$

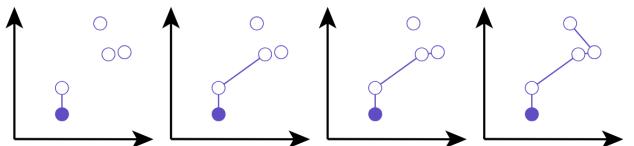


Figure 5.4: Neighbors with different  $k$  values.

Connectivity Outlier Factor does not respects the posting matrix. Rather, it creates *chains* of neighbors.

**k-NN** outlier factor (kOF) replaces the average reach at  $k$  (denoted with  $\bar{\gamma}^k$ ) with the maximum reach at  $k$  (denoted with  $\hat{\gamma}^k$ ):

$$\tilde{o}(x_i) = \gamma^k(x_i)$$

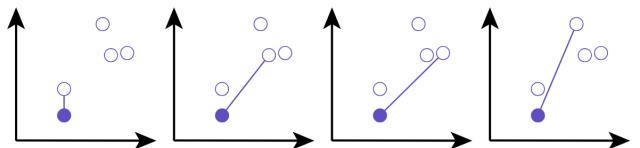


Figure 5.5: Neighbors at different  $k$  values

### 5.2.5 Concentration

We compute concentration on a two-radii approach:

- ◊ **concentration radius**  $\varepsilon$ : determines the hyper-spheres  $B(x_i, \varepsilon)$  estimating concentration  $c^\varepsilon(x_i)$  of  $x_i$  within a radius  $\varepsilon$
- ◊ **neighborhood radius**  $\delta$ : proportional to  $\varepsilon$ , determines the neighborhood  $B_i$  of  $x_i$  as the instances laying within  $B(x_i, \delta)$

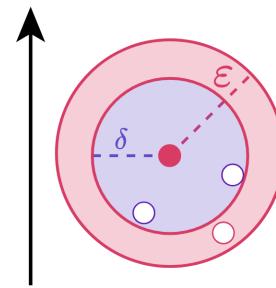


Figure 5.6: The two radii  $\varepsilon, \delta$ : the former is used to estimate *concentration*, the latter to choose which neighbors to compare concentration against. **Note:**  $\delta$  may also be larger than  $\varepsilon$ !

Locality	Local
<b>Sensitive</b>	Choice of neighborhood, connectivity parameter
<b>Interpretable</b>	Partial: can inspect what instances lead to different reaches

Table 5.2: Grading connectivity factors

### 5.2.6 Neighborhoods

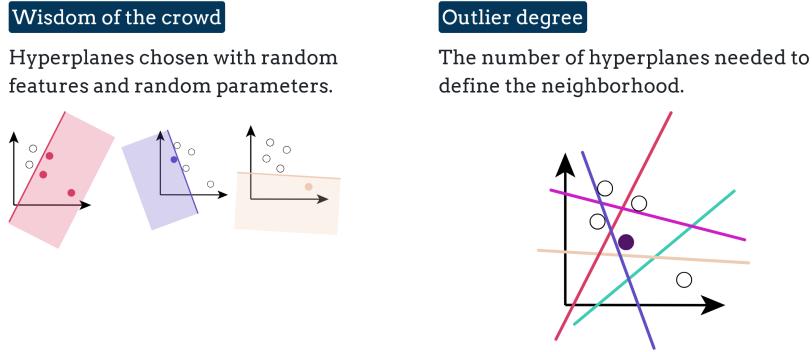


Figure 5.3: Faster and easier to compute hyperplane-based neighborhoods.

An **isolation tree**  $t$  is a random tree which randomly partitions the space into a set of blocks.

- ◊ Splits are sampled randomly
- ◊ Tree grows up to a predefined height, or until all leaves contain one instance

Outlier degree  $\tilde{o}^t(x_i) = \frac{\text{path}(x_i, t)}{c}$ , where  $c$  defines the average path length in the tree.

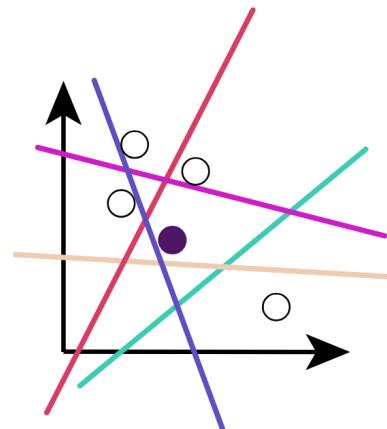


Figure 5.4: Isolation tree

An isolation forest  $T$  is comprised of several isolation trees, further sampling the hyperplane space. Outlier degree  $\tilde{o}(x_i) = 2 - \frac{\sum_{t \in T} \text{path}(x_i, t)}{|T|c}$ , where  $|T|$  is the number of trees in the forest.

	<b>Global and local</b>
<b>Sensitive</b>	Dataset noise can be interpreted as outlier
<b>Interpretable</b>	One of the most interpretable. Splits are induced by the tree, if the tree is univariate.

Table 5.3: Grading Isolation Forests



# Chapter 6

## K-Means

### 6.1 Bisection K-Means

Instead of partitioning the data set into  $K$  clusters in each iteration, bisection k-means algorithm splits one cluster into two sub clusters at each bisection step (by using k-means) until  $K$  clusters are obtained.

---

**Algorithm 6** Bisection K-means algorithm

---

- 1: Initialize the list of clusters to contain the cluster consisting of all points.
  - 2: **repeat**
  - 3:     Remove a cluster from the list of clusters.
  - 4:     {Perform several “trial” bisections of the chosen cluster.}
  - 5:     **for**  $i = 1$  to  $number\_of\_trials$  **do**
  - 6:         Bisect the selected cluster using basic K-means.
  - 7:     Select the two clusters from the bisection with the lowest total SSE.
  - 8:     Add these two clusters to the list of clusters.
  - 9: **until** Until the list of clusters contains  $K$  clusters.
- 

The algorithm is exhaustive terminating at singleton clusters (unless  $K$  is known)

- ◊ Note that Terminating at singleton clusters
  - Is time consuming
  - Singleton clusters are meaningless
  - Intermediate clusters are more likely to correspond to real classes
  - No criterion for stopping bisections before singleton clusters are reached

The resulting clusters can be refined by using their centroids as the initial centroids for the basic K-means.

### 6.2 X-Means

**X-Means** clustering algorithm is an extended K-Means which tries to automatically determine the number of clusters based on BIC scores.

The X-Means goes into action after each run of K-Means, making local decisions about which subset of the current centroids should split in order to better fit the data.

The splitting decision is done by computing the *Bayesian Information Criterion* (BIC).

#### 6.2.0.1 Bayesian Information Criterion

- ◊ A strategy to stop the Bisecting algorithm when meaningful clusters are reached to avoid over-splitting
- ◊ Using BIC as splitting criterion of a cluster in order to decide whether a cluster should split or no
- ◊ BIC measures the improvement of the cluster structure between a cluster and its two children clusters.
- ◊ Compute the BIC score of:
  - A cluster

- Two children clusters
- ◊ BIC approximates the probability that the  $M_j$  is describing the real clusters in the data

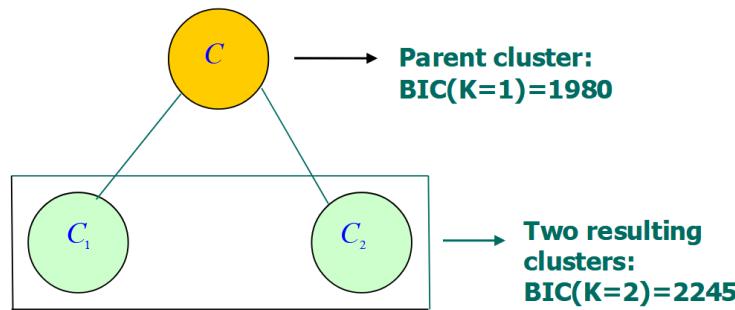


Figure 6.1: BIC splitting

The BIC score of the parent cluster is less than BIC score of the generated cluster structure  $\Rightarrow$  we accept the bisection

Forward search for the appropriate value of  $k$  in a given range  $[r_1, r_{max}]$ ; we recursively split each cluster and use BIC score to decide if we should keep each split.

1. Run K-means with  $k=r_1$
  2. Improve structure
  3. If  $k > r_{max}$  Stop and return the best-scoring model
- ◊ Use local BIC score to decide on keeping a split
  - ◊ Use global BIC score to decide which  $K$  to output at the end

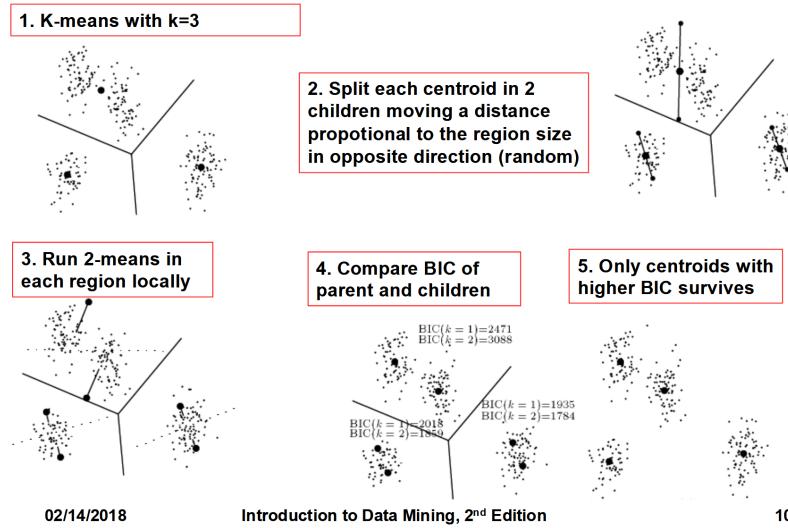


Figure 6.2: X-Means process

## 6.3 Mixture Models and the EM Algorithm

A **mixture model** is a probabilistic model for representing the presence of subpopulations within an overall population, without requiring that an observed data set should identify the subpopulation to which an individual observation belongs. The model assumes that the data is generated from a mixture of several distributions, each representing a different subpopulation.

**Algorithm 7** EM algorithm

---

- 1: Select an initial set of model parameters.
  - 2: (As with K-means, this can be done randomly or in a variety of ways.)
  - 3: **repeat**
  - 4:     **Expectation Step** For each object, calculate the probability that each object belongs to each distribution, i.e., calculate  $\text{prob}(\text{distribution}_j | \mathbf{x}_i, \Theta)$ .
  - 5:     **Maximization Step** Given the probabilities from the expectation step, find the new estimates of the parameters that maximize the expected likelihood.
  - 6:     **until** The parameters do not change.
  - 7: (Alternatively, stop if the change in the parameters is below a specified threshold.)
-



# Chapter 7

## Association Analysis

**Association Rule Mining** refers to, given a set of transactions, finding rules that will predict the occurrence of an item based on the occurrences of other items in the transaction

### 7.1 Basic Concepts

#### 7.1.1 Frequent Itemset

An **itemset** is a collection of one or more items. An itemset with  $k$  items is called a  $k$ -itemset.

An itemset is **frequent** if its **support** (the fraction of transactions that contain the itemset) is greater than or equal to a user-specified minimum support threshold.

Transaction ID	Items Purchased
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Coke, Beer, Diaper
4	Bread, Milk, Beer, Diaper
5	Bread, Milk, Diaper, Coke

Table 7.1: Example of transactions

- ◊ Itemset - *Milk, Bread, Diaper*
- ◊ Support count ( $\sigma$ ) - Frequency of occurrence of an itemset -  $\sigma(Milk, Bread, Diaper) = 2$
- ◊ Support - Fraction of transactions that contain an itemset -  $s(Milk, Bread, Diaper) = \frac{2}{5}$
- ◊ An itemset whose support is greater than or equal to a *minsup* threshold is called a frequent itemset.
- ◊ Association Rule - An implication expression of the form  $X \Rightarrow Y$ , where  $X$  and  $Y$  are itemsets - *Milk, Diaper*  $\Rightarrow$  *Beer*
- ◊ Rule Evaluation Metrics -
  - Support ( $s$ ) - The proportion of transactions that contain the itemset .
  - Confidence ( $c$ ) - The proportion of the transactions that contain  $X$  which also contain  $Y$ .

Given a set of transactions T, the goal of association rule mining is to find all rules having

- ◊ support  $\geq$  *minsup* threshold
- ◊ confidence  $\geq$  *minconf* threshold

Brute-force approach:

- ◊ List all possible association rules
- ◊ Compute the support and confidence for each rule
- ◊ Prune rules that fail the *minsup* and *minconf* thresholds

Computationally prohibitive!

Two-step approach:

- ◊ Frequent Itemset Generation - Generate all itemsets whose support  $\geq \text{minsup}$   
This is still computationally expensive
- ◊ Rule Generation - Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

## 7.2 Apriori Algorithm

The purpose of the Apriori algorithm is to find all frequent itemsets in a transaction database, which means identifying itemsets that satisfy the support constraint (threshold), i.e. finding all itemsets whose support is greater than or equal to a user-specified minimum support threshold.

The Apriori principle holds due to the following property of support of the support measure:

**Definition 7.1 (Anti-monotone Property)** *Anti-monotone property of support is formulated as:*

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

This states that the support of an itemset never exceeds the support of its subsets.

---

### Algorithm 8 Apriori Algorithm

---

- ```

1:  $k \leftarrow 1$ 
2:  $F_1 \leftarrow \{\text{frequent 1-itemsets}\}$                                 ▷ Scan database and count support of each item
3: repeat
4:    $L_{k+1} \leftarrow \text{Candidate Generation from } F_k$           ▷ Generate candidate  $(k + 1)$ -itemsets
5:    $L_{k+1} \leftarrow \text{Candidate Pruning of } L_{k+1}$           ▷ Prune candidates with infrequent  $k$ -subsets
6:   Scan transaction database to count support of each candidate in  $L_{k+1}$ 
7:    $F_{k+1} \leftarrow \text{candidates in } L_{k+1} \text{ with support } \geq \text{minsup}$ 
8:    $k \leftarrow k + 1$ 
9: until  $F_k = \emptyset$ 
10: return  $\bigcup_k F_k$   ▷ Return all frequent itemsets

```
- 

**Definition 7.2 (Closed Itemset)** An itemset  $X$  is **closed** if none of its immediate supersets has the same support as the itemset  $X$ .

$X$  is not closed if at least one of its immediate supersets has support count as  $X$ .

| TID | Items        |
|-----|--------------|
| 1   | {A, B}       |
| 2   | {B, C, D}    |
| 3   | {A, B, C, D} |
| 4   | {A, B, D}    |
| 5   | {A, B, C, D} |

Table 7.2: Example Transaction Database

| Itemset   | Support |
|-----------|---------|
| {A}       | 4       |
| {B}       | 5       |
| {C}       | 3       |
| {D}       | 4       |
| {A,B}     | 4       |
| {A,C}     | 2       |
| {A,D}     | 3       |
| {B,C}     | 3       |
| {B,D}     | 4       |
| {C,D}     | 3       |
| {A,B,C}   | 2       |
| {A,B,D}   | 3       |
| {A,C,D}   | 2       |
| {B,C,D}   | 2       |
| {A,B,C,D} | 2       |

Table 7.2: Frequent Itemsets and Their Support

### 7.2.1 Closed Itemsets

From the example above, the **closed itemsets** are:

- ◊ {B} with support 5 - closed because its immediate superset {A,B} has support  $4 \neq 5$
- ◊ {A,B} with support 4 - closed because its immediate supersets {A,B,C} and {A,B,D} have support 2 and 3 respectively, both  $\neq 4$

- ◊  $\{B,D\}$  with support 4 - closed because its immediate superset  $\{A,B,D\}$  has support  $3 \neq 4$  and  $\{B,C,D\}$  has support  $2 \neq 4$
- ◊  $\{A,B,D\}$  with support 3 - closed because its immediate superset  $\{A,B,C,D\}$  has support  $2 \neq 3$
- ◊  $\{C,D\}$  with support 3 - closed because its immediate supersets have different support
- ◊  $\{A,B,C,D\}$  with support 2 - closed because it has no supersets

Examples of **non-closed itemsets**:

- ◊  $\{A\}$  with support 4 is **not closed** because its immediate superset  $\{A,B\}$  has the same support 4
- ◊  $\{A,C\}$  with support 2 is **not closed** because its immediate superset  $\{A,B,C\}$  has the same support 2
- ◊  $\{A,B,C\}$  with support 2 is **not closed** because its immediate superset  $\{A,B,C,D\}$  has the same support 2

Closed itemsets are important because they provide a compact representation of all frequent itemsets while preserving complete support information.

### 7.2.2 Maximal Itemsets

**Definition 7.3 (Maximal Itemset)** An itemset  $X$  is **maximal** if none of its immediate supersets is frequent.  $X$  is not maximal if at least one of its immediate supersets is frequent.

Assuming a minimum support threshold of 2, the **maximal frequent itemsets** from the example are:

- ◊  $\{B\}$  with support 5 - maximal if we consider only singleton itemsets, but **not maximal** overall because its supersets  $\{A,B\}$ ,  $\{B,C\}$ ,  $\{B,D\}$  are frequent
- ◊  $\{A,B,D\}$  with support 3 - maximal because its only superset  $\{A,B,C,D\}$  has support 2, which is still frequent, so **not maximal**
- ◊  $\{B,D\}$  with support 4 - **not maximal** because its superset  $\{A,B,D\}$  is frequent
- ◊  $\{C,D\}$  with support 3 - **not maximal** because its supersets  $\{A,C,D\}$ ,  $\{B,C,D\}$ , and  $\{A,B,C,D\}$  are frequent (support  $\geq 2$ )
- ◊  $\{A,B,C,D\}$  with support 2 - **maximal** because it has no supersets and it is frequent

Therefore, with  $\text{minsup} = 2$ , the only **maximal frequent itemset** is:

- ◊  $\{A,B,C,D\}$  with support 2

Every maximal frequent itemset is also closed, but not every closed itemset is maximal. Maximal itemsets provide an even more compact representation than closed itemsets, but they only preserve information about which itemsets are frequent, not their exact support counts.

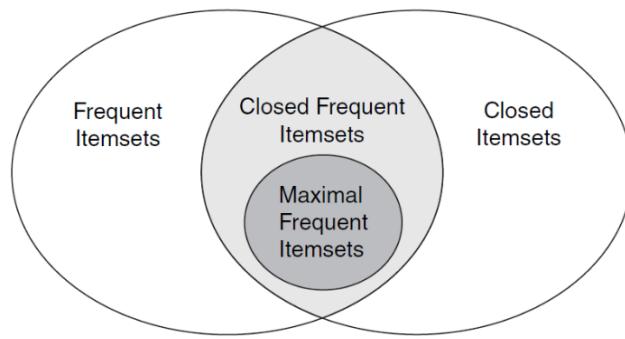


Figure 7.1: Relationships between frequent, closed, closed frequent and maximal frequent itemsets

## 7.3 Confidence

The **confidence** of an association rule  $X \Rightarrow Y$  is a measure of the reliability of the rule. It is defined as the conditional probability that a transaction contains the itemset  $Y$  given that it contains the itemset  $X$ . Mathematically, confidence is expressed as:

$$c(X \Rightarrow Y) = \frac{s(X \cup Y)}{s(X)}$$

where:

- ◊  $s(X \cup Y)$  is the support of the itemset that contains both  $X$  and  $Y$ .

- ◊  $s(X)$  is the support of the itemset  $X$ .

The confidence value ranges from 0 to 1, where a higher confidence indicates a stronger association between the itemsets  $X$  and  $Y$ . For example, a confidence of 0.8 means that 80% of the transactions that contain  $X$  also contain  $Y$ .

### 7.3.1 Drawbacks

Confidence has some drawbacks:

- ◊ It does not consider the overall frequency of the consequent itemset  $Y$  in the dataset.
- ◊ It can be misleading in cases where the consequent itemset  $Y$  is very common in the dataset, leading to high confidence values even when there is no real association between  $X$  and  $Y$ .

To address these issues, other measures such as **lift** and **conviction** are often used in conjunction with confidence to provide a more comprehensive evaluation of association rules.

#### 7.3.1.1 What rules do we want

- ◊  $\text{Confidence}(X \Rightarrow Y)$  should be sufficiently high
  - To ensure that people who buy  $X$  will more likely buy  $Y$  than not buy  $Y$
- ◊  $\text{Confidence}(X \Rightarrow Y) > \text{support}(Y)$ 
  - Otherwise, rule will be misleading because having item  $X$  actually reduces the chance of having item  $Y$  in the same transaction
- ◊ Is there any measure that capture this constraint?
  - Answer: Yes. There are many of them.

## 7.4 Other criteria

$\text{confidence}(X \Rightarrow Y) = \text{support}(Y)$  is equivalent to:

$$\begin{aligned} P(Y|X) &= P(Y) \\ P(X, Y) &= P(X)P(Y) \quad (\text{independence}) \end{aligned}$$

$$\begin{aligned} P(X, Y) &< P(X)P(Y) \quad (\text{negative correlation}) \\ P(X, Y) &> P(X)P(Y) \quad (\text{positive correlation}) \end{aligned}$$

### 7.4.1 Lift

The **lift** of an association rule  $X \Rightarrow Y$  is a measure of how much more likely the occurrence of itemset  $Y$  is when itemset  $X$  is present, compared to when  $Y$  occurs independently of  $X$ . It is defined as:

$$\text{lift}(X \Rightarrow Y) = \frac{c(X \Rightarrow Y)}{s(Y)} = \frac{s(X \cup Y)}{s(X) \times s(Y)}$$

where:

- ◊  $c(X \Rightarrow Y)$  is the confidence of the rule.
- ◊  $s(Y)$  is the support of the itemset  $Y$ .

In the slides it is defined as:

$$\text{lift}(X \Rightarrow Y) = \frac{P(X, Y)}{P(X) \times P(Y)}$$

A lift value greater than 1 indicates a positive association between  $X$  and  $Y$ , meaning that the presence of  $X$  increases the likelihood of  $Y$ . A lift value less than 1 indicates a negative association, while a lift value equal to 1 suggests that  $X$  and  $Y$  are independent.

Lift is particularly useful for identifying interesting rules that may not be apparent from confidence alone, as it accounts for the overall frequency of the consequent itemset  $Y$  in the dataset.

### 7.4.2 Interest

The **interest** of an itemset  $X$  is a measure of how much the actual support of  $X$  deviates from what would be expected if the items in  $X$  were independent. It is defined as:

$$\text{interest}(X) = s(X) - \prod_{i \in X} s(\{i\})$$

where:

- ◊  $s(X)$  is the support of the itemset  $X$ .
- ◊  $\prod_{i \in X} s(\{i\})$  is the product of the supports of the individual items in  $X$ .

In the slides it is defined as:

$$\text{interest}(X) = \frac{P(X, Y)}{P(X), P(Y)}$$

A positive interest value indicates that the items in  $X$  co-occur more frequently than would be expected under independence, suggesting a positive association. A negative interest value indicates that the items co-occur less frequently than expected, suggesting a negative association. An interest value of zero suggests that the items are independent.

### 7.4.3 Other measures

In the slide two other measures are defined:

$$\begin{aligned} PS &= P(X, Y) - P(X)P(Y) \\ \sigma\text{-coefficient} &= \frac{P(X, Y) - P(X)P(Y)}{\sqrt{P(X)P(Y)(1 - P(X))(1 - P(Y))}} \end{aligned}$$

These measures also aim to capture the degree of association between itemsets, taking into account their individual supports and the expected co-occurrence under independence.

## 7.5 Non-binary Attributes



| Gender | ... | Age | Annual Income | No of hours spent online per week | No of email accounts | Privacy Concern |
|--------|-----|-----|---------------|-----------------------------------|----------------------|-----------------|
| Female | ... | 26  | 90K           | 20                                | 4                    | Yes             |
| Male   | ... | 51  | 135K          | 10                                | 2                    | No              |
| Male   | ... | 29  | 80K           | 10                                | 3                    | Yes             |
| Female | ... | 45  | 120K          | 15                                | 3                    | Yes             |
| Female | ... | 31  | 95K           | 20                                | 5                    | Yes             |
| Male   | ... | 25  | 55K           | 25                                | 5                    | Yes             |
| Male   | ... | 37  | 100K          | 10                                | 1                    | No              |
| Male   | ... | 41  | 65K           | 8                                 | 2                    | No              |
| Female | ... | 26  | 85K           | 12                                | 1                    | No              |
| ...    | ... | ... | ...           | ...                               | ...                  | ...             |

| Male | Female | ... | Age < 13 | Age ∈ [13, 21] | Age ∈ [21, 30] | ... | Privacy = Yes | Privacy = No |
|------|--------|-----|----------|----------------|----------------|-----|---------------|--------------|
| 0    | 1      | ... | 0        | 0              | 1              | ... | 1             | 0            |
| 1    | 0      | ... | 0        | 0              | 0              | ... | 0             | 1            |
| 1    | 0      | ... | 0        | 0              | 1              | ... | 1             | 0            |
| 0    | 1      | ... | 0        | 0              | 0              | ... | 1             | 0            |
| 0    | 1      | ... | 0        | 0              | 0              | ... | 1             | 0            |
| 1    | 0      | ... | 0        | 0              | 1              | ... | 1             | 0            |
| 1    | 0      | ... | 0        | 0              | 0              | ... | 0             | 1            |
| 0    | 1      | ... | 0        | 0              | 1              | ... | 0             | 1            |
| ...  | ...    | ... | ...      | ...            | ...            | ... | ...           | ...          |

Figure 7.2: Handling non-binary attributes by changing the actual columns

### 7.5.1 Categorical attributes

Categorical attributes are variables that can take on a limited, fixed number of possible values, representing distinct categories or groups. These attributes are often used in data mining and machine learning tasks to classify or group data points based on their characteristics.

We have to apply association analysis to non-asymmetric binary attributes, so we have to formulate rules like:

$$\text{Gender} = \text{Male}, \text{Age} \in [21, 30] \Rightarrow \text{Noofhoursonline} \geq 10$$

Some attributes can have many possible values, with many of these values having very low support. To deal with this, we can use **attribute generalization**, which involves replacing specific attribute values with more general categories based on a predefined **hierarchy** or **taxonomy**. For example, instead of using specific ages, we can group them into age ranges like [0-10), [10-20), [20-30), etc. This helps to reduce the number of distinct values and increases the support for each category, making it easier to identify meaningful patterns in the data.

In some other cases we also may have a hierarchy of values for an attribute. For example, for the attribute “Location”, we may have a hierarchy like City → State → Country. In such cases, we can use the hierarchy to generalize the attribute values and find association rules at different levels of granularity.

### 7.5.2 Continuous attributes

Continuous attributes are variables that can take on an infinite number of values within a given range. These attributes are often used in data mining and machine learning tasks to represent measurements or quantities that can vary continuously.

To handle such values we have various methods:

- ◊ Discretization-based
- ◊ Statistics-based
- ◊ Non-discretization (such as MINAPRIORI)

// TODO skipped support counting with hash tree

## 7.6 Association Rule Mining

Association Rule Mining refers to, given a set of transactions, finding rules that will predict the occurrence of an item based on the occurrences of other items in the transaction. Formally, an association rule is an implication of the form  $X \Rightarrow Y$ , where  $X$  and  $Y$  are disjoint itemsets. The goal is to discover interesting relationships between items in large datasets.

### 7.6.1 Rule lists

A rule list is a sorted set of rules: starting from rule 1, if rule  $i$  does not support the given instance, move to rule  $i + 1$ .

```
(1) if age in (23, 26) and priors in (2, 3) then recidivous
(2) if age in (18, 20) then recidivous
(3) if sex is male and age in (21, 22) then recidivous
(4) if priors > 3 then recidivous
(5) else not_recidivous
```

Support, confidence, and other measures still hold, but we need to reconsider support: in a rule list, the first rule to support an instance is the one predicting. To adjust, we treat support as an indicator variable, evaluating to 1 for the *first* rule of a given list to apply, and otherwise: 0.

$$supp_A(r, x) = \begin{cases} 1 & \text{if } r \in A \text{ is the first rule to satisfy instance } x \\ 0 & \text{otherwise} \end{cases}$$

| age | priors | sex | $r_1$ | $r_2$ | $r_3$ | $r_4$ |
|-----|--------|-----|-------|-------|-------|-------|
| 24  | 4      | m   | 0     | 0     | 0     | 1     |
| 20  | 1      | f   | 0     | 1     | 0     | 0     |
| 20  | 5      | m   | 0     | 1     | 0     | 0     |

Table 7.3: Support indicator variable for a rule list

### 7.6.2 Branch and bound algorithms

Family of optimization algorithms that defines a space of solutions to search, according to a given objective function. Exploring the whole space is unfeasible, hence branch and bound algorithms define:

- ◊ A set of feasible solutions to explore from a starting set: the branches of the solution tree to explore
- ◊ A bound for the objective: it allows to prune branches, thus reducing the search space

A search tree has an exponentially large number of states! For a tree of order  $k$  and depth  $d$ ,  $\mathcal{O}(k^d)$  states! So we do not need to go very deep in the tree, just enough to find a good solution.

Exploration populates a queue of states to consider: the larger the queue, the larger the computational cost. A search algorithm simply

- ◊ Inserts states in the queue
- ◊ Pops states from the queue

## 7.7 FP Tree Growth

The FP-tree contains a compressed representation of the transaction database.

A trie (prefix-tree) data structure is used

Each transaction is a path in the tree (paths may overlap). Once the FP-tree is constructed the recursive, divide-and-conquer FP-Growth algorithm is used to enumerate all frequent itemsets.

Since transactions are sets of items, we need to transform them into ordered sequences so that we can have prefixes. We need to impose an order to the items. Initially, we assume a lexicographic order.

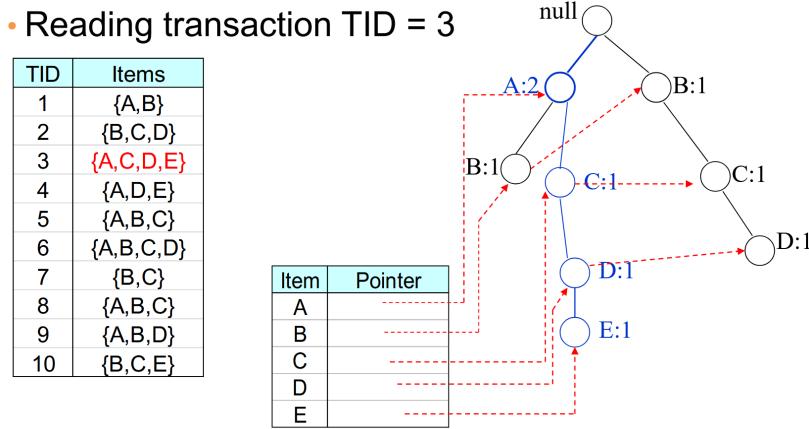


Figure 7.3: Each transaction is a path in the trie. There is a header table to point to the first occurrence of each item in the tree.

The size of the tree depends on the compressibility of the data

- ◊ *Extreme case:* All transactions are the same, the FP-tree is a single branch
- ◊ *Extreme case:* All transactions are different the size of the tree is the same as that of the database (bigger actually since we need additional pointers)

The size of the tree also depends on the ordering of the items. We could order the items according to their frequency from larger to smaller, but we need to do an extra pass over the dataset to count frequencies.

---

### Algorithm 9 FP-Growth Algorithm

---

- 1: **for** each suffix  $X$  **do**
  - 2:   **Phase 1:** Construct the prefix tree for  $X$  and compute the support using the header table and the pointers
  - 3:   **if**  $X$  is frequent **then**
  - 4:     **Phase 2:** Construct the *conditional* FP-tree for  $X$ :
  - 5:       1. Recompute support
  - 6:       2. Prune infrequent items
  - 7:       3. Prune leaves and recurse
-



# Chapter 8

## Sequential Pattern Mining

Sequential pattern mining is the task of discovering statistically relevant patterns between data examples where the values are delivered in a sequence. A sequence is an ordered list of events. Each event is a set of items that occur together. A sequential pattern is a sequence that occurs frequently in a sequence database.

Consider the following example of a sequence of different transactions by a customer at an online store:

`< {Digital Camera,iPad} {memory card} {headphone,iPad cover} >`

This sequence indicates that the customer first bought a digital camera and an iPad, then later bought a memory card, and finally bought a headphone and an iPad cover together.

The knowledge we can extract here is that probably the customer didn't realize he needed the memory card when he bought the digital camera, so he bought it later. Also, after buying the iPad, he probably realized he needed accessories for it, so he bought them together.

Databases of transactions usually have a temporal information which *Sequential patterns* can exploit.

### 8.1 Definitions

A sequence is an ordered list of elements (transactions/itemsets). An element is a set of items that occur at the same time, and it is also associated with a timestamp.

$$s = \langle e_1 e_2 e_3 \rangle$$

$$e_i = \{i_1, i_2, \dots, i_k\}$$

Length of a sequence,  $|s|$ , is given by the number of elements of the sequence.

A k-sequence is a sequence that contains k events (items).

A sequence  $s_a = \langle a_1 a_2 \dots a_n \rangle$  is a subsequence of another sequence  $s_b = \langle b_1 b_2 \dots b_m \rangle$  (denoted as  $s_a \subseteq s_b$ ) if there exist integers  $1 \leq i_1 < i_2 < \dots < i_n \leq m$  such that  $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$ .

| Data sequence                                | Subsequence                      | T/F |
|----------------------------------------------|----------------------------------|-----|
| $\langle \{2, 4\} \{3, 5, 6\} \{8\} \rangle$ | $\langle \{2\} \{3, 5\} \rangle$ | T   |
| $\langle \{1, 2\} \{3, 4\} \rangle$          | $\langle \{1\} \{2\} \rangle$    | F   |
| $\langle \{2, 4\} \{2, 4\} \{2, 5\} \rangle$ | $\langle \{2\} \{4\} \rangle$    | T   |

Table 8.1: Examples of subsequence containment

**Definition 8.1 (Subsequence support)** *The support of a subsequence w is the fraction of data sequences in the sequence database S that contain w.*

#### 8.1.1 Exercises

##### 8.1.1.1 Exercise 1

Find instances/occurrences of the following subsequences in the sequence database

$$\begin{aligned}
 & \langle \{C\} \{H\} \{C\} \rangle \\
 & \langle \{A\} \{F\} \rangle \\
 & \langle \{A\} \{A\} \{D\} \rangle \\
 & \langle \{A\} \{A, B\} \{F\} \rangle \\
 \\ 
 & \langle \{A, C\} \{C, D\} \{F, H\} \{A, B\} \{B, C, D\} \{E\} \{A, B, D\} \{F\} \rangle \\
 & t = 0 \ t = 1 \ t = 2 \ t = 3 \ t = 4 \ t = 5 \ t = 6 \ t = 7
 \end{aligned}$$

### 8.1.1.2 Exercise 2

Find instances/occurrences of the following subsequences in the sequence database

$$\langle \{C\} \{H\} \{C\} \rangle \langle \{A\} \{B\} \rangle \langle \{C\} \{C\} \{E\} \rangle \langle \{A\} \{E\} \rangle$$

$$\begin{aligned}
 & \langle \{A, C\} \{C, D, E\} \{F\} \{A, H\} \{B, C, D\} \{E\} \{A, B, D\} \{F\} \rangle \\
 & t = 0 \ t = 1 \ t = 2 \ t = 3 \ t = 4 \ t = 5 \ t = 6 \ t = 7
 \end{aligned}$$

## 8.2 Towards an Algorithm

**Definition 8.2 (Sequential Pattern Mining)** Given a sequence database  $S$  and a minimum support threshold  $\text{min\_sup}$ , the **sequential pattern mining** task is to find all subsequences  $w$  such that  $\text{support}(w) \geq \text{min\_sup}$ .

The most trivial, yet inefficient, way to find all sequential patterns is to generate all possible  $k$ -subsequences for  $k = 1, 2, \dots$  and count their support in the sequence database. This approach is computationally expensive due to the exponential number of possible subsequences.

### 8.2.1 GSP - Generalized Sequential Pattern

Follows the same structure of the Apriori algorithm, i.e. starting from short patterns and finding longer ones at each iteration.

It is based on “Apriori principle” (or “anti-monotonicity of support”) which states that if a sequence is not frequent, none of its super-sequences can be frequent.

$$S_1 \subseteq S_2 \implies \text{support}(S_1) \geq \text{support}(S_2)$$

**Proof:** Any input sequence that contains  $S_2$  will also contain  $S_1$

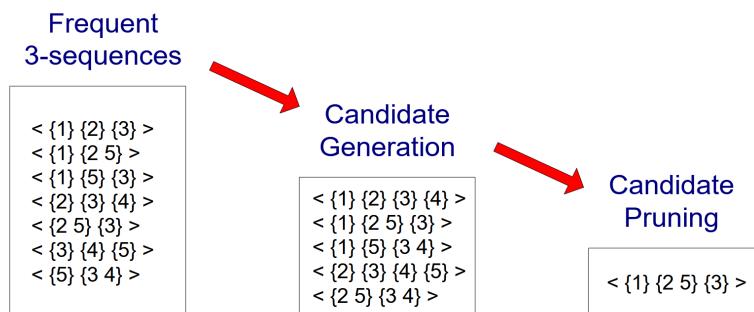


Figure 8.1: GSP example

**Algorithm 10** GSP - Generalized Sequential Pattern

---

```

1: Step 1: Make the first pass over the sequence database  $D$  to yield all the 1-element frequent sequences
2:
3: repeat
4:   Candidate Generation:
5:     Merge pairs of frequent subsequences found in the  $(k - 1)$ -th pass to generate candidate sequences that
       contain  $k$  items
6:
7:   Candidate Pruning:
8:     Prune candidate  $k$ -sequences that contain infrequent  $(k - 1)$ -subsequences
9:
10:  Support Counting:
11:    Make a new pass over the sequence database  $D$  to find the support for these candidate sequences
12:
13:  Candidate Elimination:
14:    Eliminate candidate  $k$ -sequences whose actual support is less than  $min\_sup$ 
15: until no new frequent sequences are found

```

---

**8.2.1.1 Candidate Generation in GSP**

Given  $n$  events:  $i_1, i_2, i_3, \dots, i_n$

◊ Candidate 1-subsequences:

$$\langle\{i_1\}\rangle, \langle\{i_2\}\rangle, \langle\{i_3\}\rangle, \dots, \langle\{i_n\}\rangle$$

◊ Candidate 2-subsequences:

$$\langle\{i_1, i_2\}\rangle, \langle\{i_1, i_3\}\rangle, \dots, \langle\{i_1\}i_1\rangle, \langle\{i_1\}i_2\rangle, \dots, \langle\{i_{n-1}\}i_n\rangle$$

◊ Candidate 3-subsequences:

$$\begin{aligned} &\langle\{i_1, i_2, i_3\}\rangle, \langle\{i_1, i_2, i_4\}\rangle, \dots, \langle\{i_1, i_2\}i_1\rangle, \langle\{i_1, i_2\}i_2\rangle, \dots, \\ &\langle\{i_1\}i_1, i_2\rangle, \langle\{i_1\}i_1, i_3\rangle, \dots, \langle\{i_1\}i_1\rangle i_1, \langle\{i_1\}i_1\rangle i_2, \dots \end{aligned}$$

**Remark:** events within an element are ordered

◊ YES:  $\langle\{i_1, i_2, i_3\}\rangle$     NO:  $\langle\{i_3, i_1, i_2\}\rangle$

**8.2.1.2 Candidate Generation - Merging Procedure**

◊ **Base case ( $k = 2$ ):**

- Merging two frequent 1-sequences  $\langle\{i_1\}\rangle$  and  $\langle\{i_j\}\rangle$  will produce two candidate 2-sequences:  $\langle\{i_1\}i_j\rangle$  and  $\langle\{i_1, i_j\}\rangle$
- Special case:  $i_1$  can be merged with itself:  $\langle\{i_1\}i_1\rangle$

◊ **General case ( $k > 2$ ):**

- A frequent  $(k - 1)$ -sequence  $w_1$  is merged with another frequent  $(k - 1)$ -sequence  $w_2$  to produce a candidate  $k$ -sequence if the subsequence obtained by removing the **first event** in  $w_1$  is the same as the one obtained by removing the **last event** in  $w_2$

- The resulting candidate after merging is given by the sequence  $w_1$  extended with the last event of  $w_2$ .
  - If last two events in  $w_2$  belong to the same element  $\Rightarrow$  last event in  $w_2$  becomes part of the last element in  $w_1$ :

$$\langle\{d\}a\{b\}\rangle + \langle\{a\}b,c\rangle = \langle\{d\}a\{b,c\}\rangle$$

- Otherwise, the last event in  $w_2$  becomes a separate element appended to the end of  $w_1$ :

$$\langle\{a,d\}b\rangle + \langle\{d\}b,c\rangle = \langle\{a,d\}b,c\rangle$$

- Special case: check if  $w_1$  can be merged with itself

- Works when it contains only one event type:  $\langle\{a\}a\rangle + \langle\{a\}a\rangle = \langle\{a\}aa\rangle$

**8.2.1.3 Merging Examples**

◊ Merging the sequences  $w_1 = \langle\{1\}2\ 3\}4\rangle$  and  $w_2 = \langle\{2\}3\}4\ 5\rangle$

- will produce the candidate sequence  $\langle\{1\}\{2\ 3\}\{4\ 5\}\rangle$  because the last two events in  $w_2$  (4 and 5) belong to the same element
- ◊ Merging the sequences  $w_1 = \langle\{1\}\{2\ 3\}\{4\}\rangle$  and  $w_2 = \langle\{2\ 3\}\{4\}\{5\}\rangle$ 
  - will produce the candidate sequence  $\langle\{1\}\{2\ 3\}\{4\}\{5\}\rangle$  because the last two events in  $w_2$  (4 and 5) do not belong to the same element
- ◊ We **do not have to** merge the sequences  $w_1 = \langle\{1\}\{2\ 6\}\{4\}\rangle$  and  $w_2 = \langle\{1\}\{2\}\{4\ 5\}\rangle$  to produce the candidate  $\langle\{1\}\{2\ 6\}\{4\ 5\}\rangle$ 
  - Notice that if the latter is a viable candidate, it will be obtained by merging  $w_1$  with  $\langle\{2\ 6\}\{4\ 5\}\rangle$

#### 8.2.1.4 Candidate Pruning

Candidate pruning follows —again— the Apriori principle: If a k-sequence  $W$  contains a  $(k - 1)$ -subsequence that is not frequent, then  $W$  is not frequent and can be pruned.

##### Method:

- ◊ Enumerate all  $(k - 1)$ - subsequences:
  - $\{a, b\}\{c\}\{d\} \rightarrow \{b\}\{c\}\{d\}, \{a\}\{c\}\{d\}, \{a, b\}\{d\}, \{a, b\}\{c\}$
- ◊ Each subsequence generated by cancelling 1 event in  $W$ 
  - Number of  $(k - 1)$ - subsequences =  $k$
- ◊ Remark: candidates are generated by merging two “mother”  $(k - 1)$ - subsequences that we know to be frequent
  - Correspond to remove the first event or the last one
  - Number of significant  $(k - 1)$ - subsequences to test =  $k - 2$
  - Special cases: at step  $k = 2$  the pruning has no utility, since the only  $(k - 1)$ - subsequences are the “mother” ones

## 8.3 Timing Constraints

In some applications, it is useful to impose timing constraints on the sequential patterns to be mined. Typical timing constraints include:

- ◊  $x_g$  - **max gap**: Each element of the pattern instance must be at most  $x_g$  time after the previous one
- ◊  $n_g$  - **min gap**: Each element of the pattern instance must be at least  $n_g$  time after the previous one
- ◊  $m_s$  - **max span**: The overall duration of the pattern instance must be at most  $m_s$  time

#### 8.3.1 Contiguous Subsequences

**Definition 8.3 (Contiguous Subsequence)**  $s$  is a **contiguous subsequence** of  $w = \langle e_1 \rangle \langle e_2 \rangle \dots \langle e_k \rangle$  if any of the following conditions hold:

1.  $s$  is obtained from  $w$  by deleting an item from either  $e_1$  or  $e_k$  (avoids internal “jumps”)
2.  $s$  is obtained from  $w$  by deleting an item from any element  $e_i$  that contains **at least 2 items** (not interesting for our usage)
3.  $s$  is a contiguous subsequence of  $s'$  and  $s'$  is a contiguous subsequence of  $w$  (recursive definition)

**Examples:** Consider  $s = \langle\{1\}\{2\}\rangle$

- ◊  $s$  is a contiguous subsequence of:
 
$$\langle\{1\}\{2, 3\}\rangle$$

- ◊  $s$  is not a contiguous subsequence of:
 
$$\langle\{1\}\{3\}\{2\}\rangle \quad \text{and} \quad \langle\{2\}\{1\}\{3\}\{2\}\rangle$$

In some domains, we may have only one very long time series, for example:

- ◊ monitoring network traffic events for attacks
- ◊ monitoring telecommunication alarm signals

Goal is to find frequent sequences of events in the time series, so we have to count “instances”, but which ones? This problem is also known as *frequent episode mining*.

# Chapter 9

## Supervised Machine Learning

A machine is said to learn if, when tackling a task, it is able to improve its own performance through experience.

- ◊ Task  $T$ : the problem we are trying to solve, e.g., predict the cancer risk of a patient
- ◊ Experience  $E$ : the experience on the task provided to the model, e.g., some dataset
- ◊ Performance  $P$ : a measure of success, e.g., the success rate in predicting cancer
- ◊ Model  $F$ : a function  $f$  solving the task with a learning algorithm  $A$

| Task                      | Predict                        | Example                                                               |
|---------------------------|--------------------------------|-----------------------------------------------------------------------|
| Binary classification     | one of two discrete labels     | Is the patient at high risk of developing cancer, or not?             |
| Multilabel classification | any of several discrete labels | Of all the possible syndromes, which is the patient going to develop? |
| Multiclass classification | one of several discrete labels | The student is going to major in...?                                  |
| Regression                | a continuous label             | The student's grade is going to be...?                                |

Table 9.1: Types of supervised learning tasks

### 9.1 Experience or not

Models are not developed to aid on known experiences, rather on **unknown** ones. The performance of a model can't be uniquely measured on its performance on the given experience, but rather on novel experiences which the model was not preview to. We want to achieve a low generalization error.

- ◊ **Optimization**  
Maximizes performance on the given experience  $E$ : optimization performance
- ◊ **Machine Learning**  
Maximizes performance on the given, and expected non-given, experience  $E$  and  $\bar{E}$ : generalization performance

If we do not know the non-given experience  $\bar{E}$  how can we ever expect to be effective on it?

**Equal distribution assumption.** It is assumed that the given experience  $E$ , and the non- given experience  $\bar{E}$  are sampled from the same distribution  $Pr^E$ .

Given a learning algorithm  $A$ , can I always expect the performance to transfer?

**No.**

Sampling from the data distribution is a random process, and the resulting models learned end up erring in terms of:

- ◊ Bias: the expected performance decrease w.r.t. the best model
- ◊ Variance: the variance with respect to different samples

Ideally, we want to have learning algorithms with low enough bias and variance.

### 9.2 Improper models

There are two categories of improper models:

- ◊ Underfit. High bias, high variance. Models which have poor performance, regardless of samples. They have not learnt enough!
- High variance, low bias. Models which ought to improve their performance on the given experience  $E$
- ◊ Overfit. Low bias, high variance. Models which have overfit on the given experience, and ought to improve their generalization performance. They have learnt too much!
- Low variance, high bias. Models which ought to improve their generalization performance on the unknown experience  $\bar{E}$

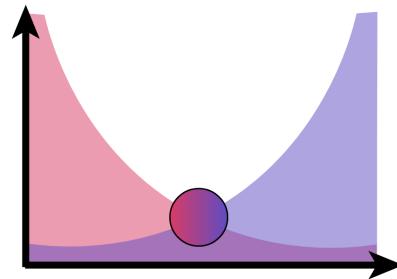


Figure 9.1: Bias (red) and variance (blue) as a decomposition of model performance

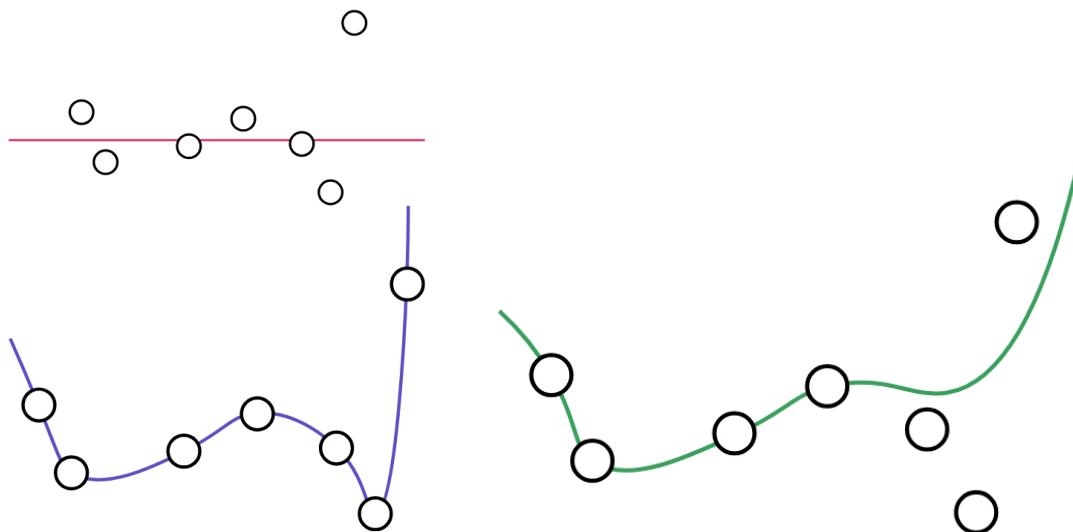


Figure 9.1: Two models approximating a dataset: one model has too low a capacity and underfits the data (in red), while another has too high a capacity and overfits the data (in blue). These are opposed to a well-fitted model (in green) which generalizes well, without being too strict and leaving some space for variance in data.

### 9.3 Searching for models

- ◊ *Tackling the generalization gap.* Data-based strategies to maximize generalization performance
- ◊ *Performance evaluation.* How to measure model performance/error
- ◊ *Parameterization.* Exploring the space of models

We design two phases in the model search:

- ◊ **Model selection**  
A learning phase wherein, among all possible models in a model space  $\mathcal{F}$ , we select a model  $f$
- ◊ **Model Validation**  
A learning phase wherein we estimate the generalization performance (error) of the selected model  $f$

Model validation cannot affect model selection: it only goes one way, from selection to validation. Thus, we need to incorporate in model selection some strategy to avoid under/overfitting

### 9.3.1 Data Partitioning

The standard approach is model agnostic: we can apply this to any learning algorithm or family of models we want. We operate a tripartite partitioning of  $(X, Y)$ :

- ◊ **Training dataset**  $(X^{tr}, Y^{tr})$ : search through the models' space  $\mathcal{F}$
- ◊ **Validation dataset**  $(X^{vl}, Y^{vl})$ : guesstimate the generalization performance of candidate models  $f_1, \dots, f_k$
- ◊ **Test dataset**  $(X^{ts}, Y^{ts})$ : estimate the generalization performance of the selected model  $f_i$

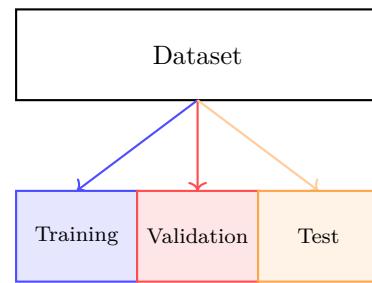


Figure 9.2: A partition of the dataset (in black) in training (blue), validation (red), and test (beige).

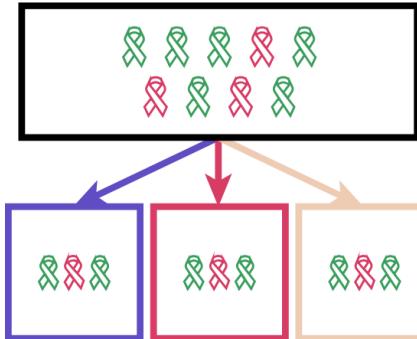
| Task                | Training         | Validation       | Test          |
|---------------------|------------------|------------------|---------------|
| Disease diagnosis   | Biology lectures | Homework         | Exam          |
| Learning a language | Duolingo         | Exchange student | Living abroad |
| Pandemic diffusion  | Black plague     | Ebola            | Covid         |

Table 9.2: Analogy between data partitioning and real-world learning

#### 9.3.1.1 Partitioning the dataset

How to partition the dataset properly?

- ◊ Size. Test and validation set of similar size, training set of much larger size, e.g., a ratio of 4:1. Some learning algorithms are more data-hungry, so this is a starting baseline.
- ◊ Distribution. Ideally, same distribution for all three datasets. Random stratified sampling is used



A partition of the dataset (in black) in training (blue), validation (red), and test (beige). Patients with (red cross) and without (green cross) cancer are split evenly among the blocks.

Figure 9.2: Partitioned data

Hold-out leverages the three-blocks partition train-validation-test.

- ◊ Learn candidate models  $f_1, \dots, f_k$  on the training set
- ◊ Evaluate them on the validation set
- ◊ Estimate the generalization error on the test set

Stretching hold-out, we aim to further increase the size of the validation set. We partition a given set of data, e.g., the training dataset, in two folds:

- ◊ in set 1, block 1 is a training dataset, block 2 is the validation dataset
- ◊ in set 2, block 1 is a training dataset, block 2 is the validation dataset

Now I can learn and guesstimate a model  $f_i$  on both folds!

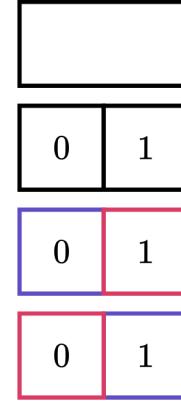
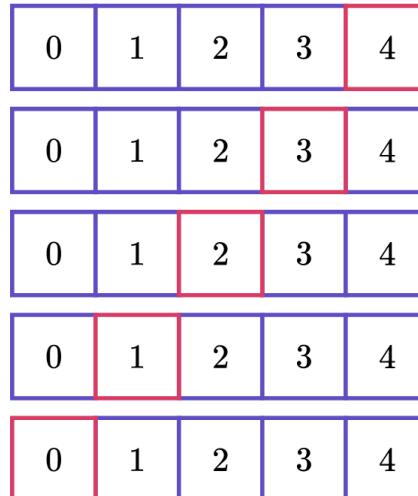


Figure 9.3: A set partitioned in two blocks, and the resulting folds. In each fold, a block acts as validation set (in red), and the other as training set (in blue).



A set partitioned in  $k = 5$  blocks, and the resulting *folds*. In each fold, a block acts as validation set (red), and the other as training set (blue).

Figure 9.3: A set partitioned in  $k$  blocks, and the resulting folds. In each fold, a block acts as validation set (in red), and the other as training set (in blue).

## 9.4 Performance evaluation

With proper partitioning, we are now able to feed experiences (data) to the candidate models  $f_1, \dots, f_k$  we wish to select. How do we evaluate them?

**Classification.** Classification tasks generally aim to measure a Hamming distance ( $\ominus$ ) between the gold labels, and the labels given by the model. Either measured as error or performance.

**Reminder:** we indicate the vector of  $n$  gold labels with  $Y \in y$ , the model of interest  $f$  with  $f$ , its prediction on an instance  $x^i$  with  $f(x^i)$ , and the indicator function with  $\mathbb{1}$ .

| Task               | Measure            | Formulation                                                                     |
|--------------------|--------------------|---------------------------------------------------------------------------------|
| Binary, Multiclass | Accuracy           | $\frac{1}{n} \sum_{i=1}^n \mathbb{1}(Y_i = f(x^i))$ or $1 - \text{Error rate}$  |
| Binary, Multiclass | Error rate         | $\frac{1}{n} \sum_{i=1}^n \mathbb{1}(Y_i \neq f(x^i))$ or $1 - \text{Accuracy}$ |
| Multilabel         | Jaccard similarity | $\frac{1}{n} \sum_{i=1}^n \frac{Y_i \cap f(x^i)}{Y_i \cup f(x^i)}$              |
| Multilabel         | Hamming error      | $\frac{1}{n} \sum_{i=1}^n 1 - (Y_i \ominus f(x^i))$                             |

Table 9.3: Performance measures for classification tasks

#### 9.4.1 Confusion Matrix

In some binary cases, one label  $y$  is for us of interest (positive label), e.g., patients we predict will have cancer, while the other is not (negative label). We can construct a *confusion matrix* out of the predictions  $f(x^i)$  of the model, that we can then leverage to define more performance measures.

|        |          | $Y$      |          |
|--------|----------|----------|----------|
|        |          | positive | negative |
| $f(X)$ | positive | $tp$     | $fp$     |
|        | negative | $fn$     | $tn$     |

Figure 9.4: A confusion matrix: columns defined by the gold labels  $Y$ , and rows defined by the predicted labels  $f(X)$ .

| Task   | Measure      | Formulation                          | Description                                                                            |
|--------|--------------|--------------------------------------|----------------------------------------------------------------------------------------|
| Binary | Precision    | $\frac{tp}{tp + fp}$                 | Of all the positive predictions, how many, in proportion, are correct?                 |
| Binary | Recall       | $\frac{tp}{tp + fn}$                 | Of all the positive instances, how many, in proportion, have been correctly predicted? |
| Binary | $f_1$ -score | $h(\text{precision}, \text{recall})$ | Harmonic mean of precision and recall                                                  |
| Binary | Accuracy     | $\frac{tp + tn}{tp + tn + fp + fn}$  | Accuracy                                                                               |

Table 9.4: Binary classification performance measures from confusion matrix

|        |          | $Y$           |               |
|--------|----------|---------------|---------------|
|        |          | positive      | negative      |
| $f(X)$ | positive | $\omega_{tp}$ | $\omega_{fp}$ |
|        | negative | $\omega_{fn}$ | $\omega_{tn}$ |

Confusion matrices are limited in their weighting, as any entry, e.g., true positives ( $tp$ ), has a unitary weight in all performance measures. Yet, in some cases, some false weigh heavier than others, e.g., diagnosing a false positive cancer is far worse than diagnosing a false negative. Thus, we introduce the cost matrix, holding one weight per each entry in the confusion matrix, weighing it in performance measures.

Figure 9.5: Confusion 2



# Chapter 10

## Decision Trees

### 10.1 Classification

Classification is a data mining task that assigns a class label to a record based on its attribute values.

We start off with **training set** of records, each characterized by a tuple  $(x, y)$ , where  $x$  is the attribute set and  $y$  is the class label

- ◊  $x$ : attribute, predictor, independent variable, input
- ◊  $y$ : class, response, dependent variable, output

The goal of classification is to learn a function (often called “**model**”)  $f : X \rightarrow Y$  that maps each attribute set  $x$  into one of the predefined class labels  $y$ .

The **goal** of classification is to accurately predict the class labels of **unseen records** (i.e., records not in the training set).

A test set is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.

- ◊ **Base classifier:** a classifier built from the training set
  - Decision Tree based Methods
  - Rule-based Methods
  - Nearest-neighbor
  - Neural Networks
  - Deep Learning
  - Naïve Bayes and Bayesian Belief Networks
  - Support Vector Machines
- ◊ **Ensemble classifier:** a classifier that combines multiple base classifiers to improve accuracy
  - Bagging
  - Boosting
  - Random Forests

### 10.2 Classification with Decision Trees

A decision tree is a flowchart-like tree structure, where each internal node denotes a test on an attribute, each branch represents the outcome of the test, and each leaf node holds a class label.

There are various algorithms to build decision trees, such as Hunt’s algorithm (among the earliest ones), ID3, C4.5, CART, etc.

#### 10.2.1 Hunt’s Algorithm

Let  $D_t$  be the set of training records at node  $t$ . The general procedure of Hunt’s algorithm is as follows:

- ◊ If  $D_t$  contains records that belong the same class  $y_t$ , then  $t$  is a leaf node labeled as  $y_t$

- ◇ If  $D_t$  contains records that belong to more than one class, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.

#### 10.2.1.1 Example: Building a Decision Tree Step-by-Step

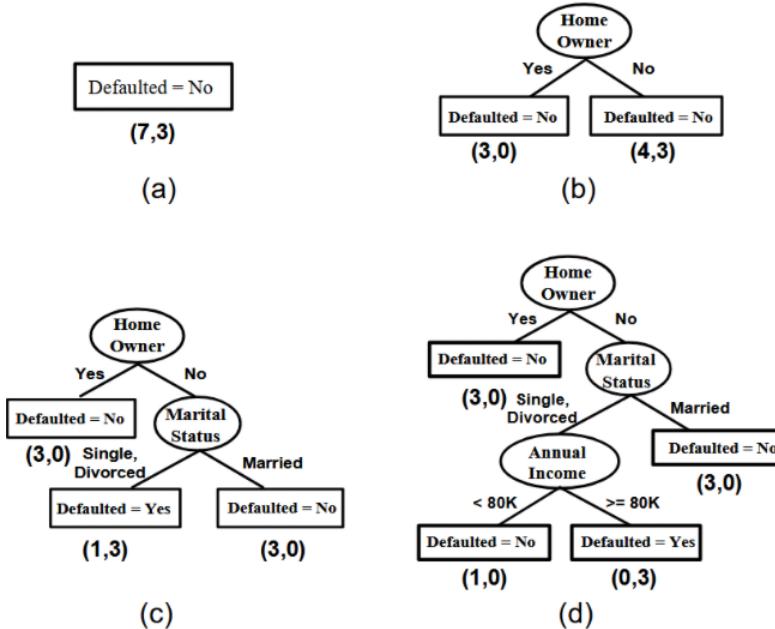


Figure 10.1: Hunt's algorithm steps applied to Table 10.1

| ID | Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|----|------------|----------------|---------------|--------------------|
| 1  | Yes        | Single         | 125K          | No                 |
| 2  | No         | Married        | 100K          | No                 |
| 3  | No         | Single         | 70K           | No                 |
| 4  | Yes        | Married        | 120K          | No                 |
| 5  | No         | Divorced       | 95K           | Yes                |
| 6  | No         | Married        | 60K           | No                 |
| 7  | Yes        | Divorced       | 220K          | No                 |
| 8  | No         | Single         | 85K           | Yes                |
| 9  | No         | Married        | 75K           | No                 |
| 10 | No         | Single         | 90K           | Yes                |

Table 10.1: Example dataset for classification

Using the dataset in Table 10.1, we can illustrate how Hunt's algorithm builds a decision tree iteratively:

**Step (a):** The algorithm starts with all training records at the root node. The initial dataset contains 10 records with 7 instances of “No” and 3 instances of “Yes” for the Defaulted Borrower class. Since the node contains records from both classes, we need to split the data.

**Step (b):** The algorithm selects *Home Owner* as the first splitting attribute. This creates two branches:

- ◇ **Home Owner = Yes:** Contains 3 records, all with class label “Defaulted = No” (3,0). Since all records belong to the same class, this becomes a leaf node.
- ◇ **Home Owner = No:** Contains 7 records with 4 instances of “No” and 3 instances of “Yes” (4,3). Since this node contains mixed classes, further splitting is required.

**Step (c):** The algorithm continues by splitting the *Home Owner = No* branch using *Marital Status* as the splitting attribute:

- ◇ **Home Owner = No, Marital Status = Single or Divorced:** Contains 4 records with 1 instance of “No” and 3 instances of “Yes” (1,3). The majority class is “Yes”, but the node is not pure.
- ◇ **Home Owner = No, Marital Status = Married:** Contains 3 records, all with class label “Defaulted = No” (3,0). This becomes a leaf node.

**Step (d):** The algorithm makes a final split on the remaining impure node (*Home Owner = No, Marital Status = Single or Divorced*) using *Annual Income* with threshold 80K:

- ◊ **Home Owner = No, Marital Status = Single/Divorced, Annual Income < 80K:** Contains 1 record with class “Defaulted = No” (1,0). This becomes a leaf node.
- ◊ **Home Owner = No, Marital Status = Single/Divorced, Annual Income  $\geq 80K$ :** Contains 3 records with class “Defaulted = Yes” (0,3). This becomes a leaf node.

At this point, all leaf nodes contain records from a single class (pure nodes), and the decision tree is complete. The tree can now be used to classify new instances by traversing from the root to a leaf based on the attribute values of the test record.

### 10.2.1.2 Splitting Strategies

When building decision trees, there are different strategies for splitting nodes based on categorical attributes:

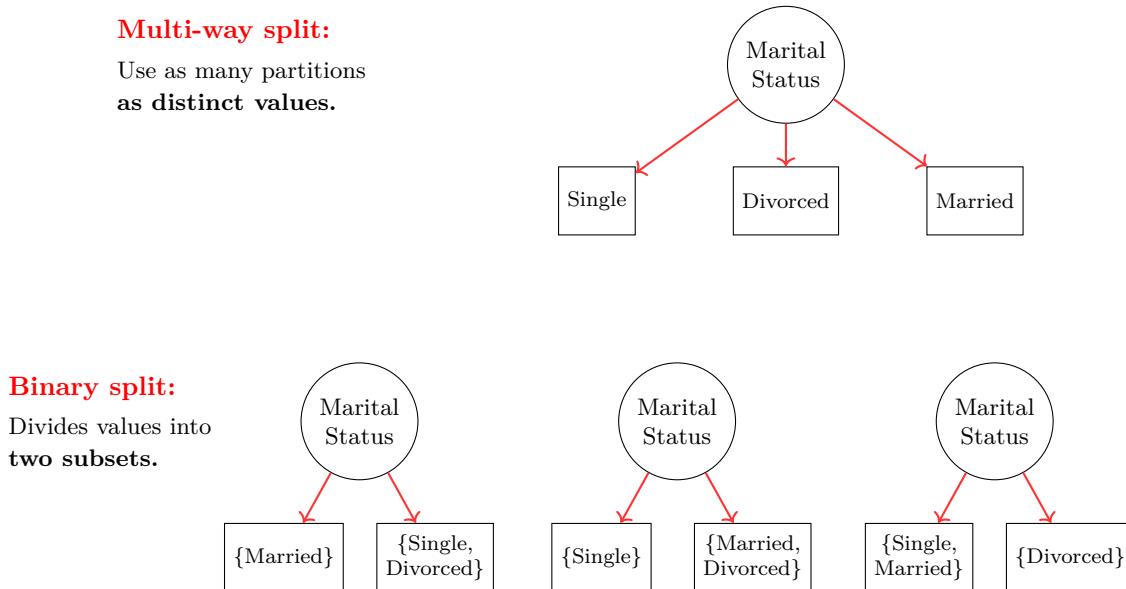


Figure 10.2: Splitting strategies for categorical attributes: Multi-way vs Binary splits

**Multi-way split:** Creates as many child nodes as there are distinct values in the attribute. For example, splitting on Marital Status with three values (Single, Divorced, Married) creates three branches. This approach is intuitive but can lead to data fragmentation, especially when attributes have many distinct values.

**Binary split:** Divides the attribute values into two subsets, creating only two child nodes. There are multiple ways to partition the values (e.g., {Married} vs {Single, Divorced}, or {Single} vs {Married, Divorced}, etc.). Binary splits are preferred in algorithms like CART as they create more balanced trees and are computationally more efficient, though they may require more splits to achieve the same separation.

### 10.2.1.3 Continuous attributes

There are two main approaches to handle continuous attributes in decision trees:

- ◊ **Discretization** to form an ordinal categorical attribute based on a set of **intervals**. Such intervals may be determined using domain knowledge or algorithms like equal-width or equal-frequency binning (percentiles), or clustering.  
The discretization may be static, so performed only once at the beginning of the tree construction, or dynamic, so performed at each node during tree construction.
- ◊ **Binary splits** based on thresholding, e.g.,  $\text{Income} \leq 100K$  vs  $\text{Income} > 100K$ .  
All possible splits should be considered to find the best threshold that optimizes a certain criterion (e.g., information gain, Gini index, etc.). This may result in an intensive computation.

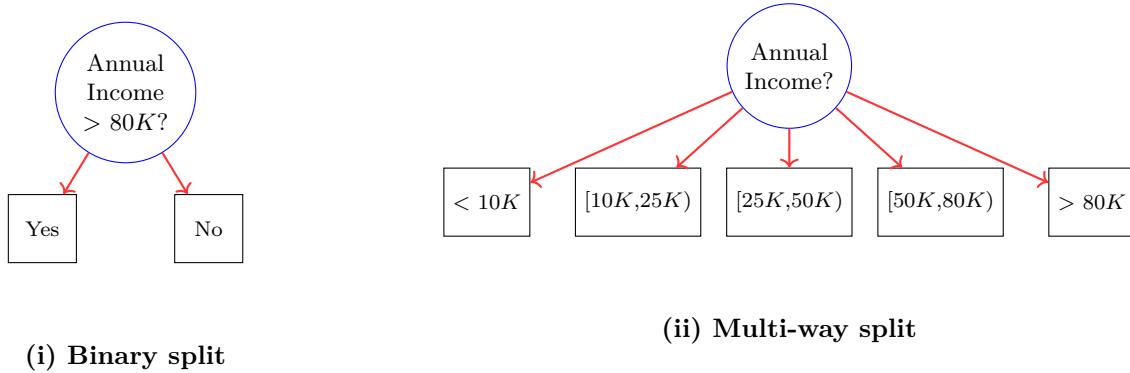


Figure 10.3: Splitting strategies for continuous attributes

#### 10.2.1.4 Choosing the best split

Nodes with purer homogeneous class distributions are preferred. To measure the **impurity** of a node, we can use metrics such as **Gini Index**, **Information Gain(?)**, **Entropy** or **Misclassification error**.

$$\textbf{Gini Index: } Gini(t) = 1 - \sum_{i=1}^c p(i|t)^2$$

$$\textbf{Entropy: } Entropy(t) = - \sum_{i=1}^c p(i|t) \log_2 p(i|t)$$

**Misclassification Error:**  $ME(t) = 1 - \max_i p(i|t)$

Practically, we compute the chosen impurity measure  $P$  before splitting attributes, and then we compute it again ( $M$ ) after the split for each child node.  $M$  is usually a weighted average of the impurity measures of the child nodes.

At this point we may choose the attribute test condition that produces the highest gain  $G = P - M$ , or, equivalently, the lowest impurity measure after splitting ( $M$ ).

### 10.2.2 Gini Index

### 10.2.2.1 Gini Index Examples

The following examples illustrate how the Gini Index measures node impurity:

|           |          |                                                |                   |
|-----------|----------|------------------------------------------------|-------------------|
| <b>C1</b> | <b>0</b> | $P(C1) = 0/6 = 0$                              | $P(C2) = 6/6 = 1$ |
| <b>C2</b> | <b>6</b> | $Gini = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$ |                   |

This represents a **pure node** where all instances belong to class C2. The Gini Index is 0, indicating no impurity.

|           |          |
|-----------|----------|
| <b>C1</b> | <b>1</b> |
| <b>C2</b> | <b>5</b> |

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$Gini = 1 - (1/6)^2 - (5/6)^2 = 0.278$$

This node has low impurity with most instances belonging to C2.

This node has higher impurity with a less uniform distribution. Note that the Gini Index reaches its maximum value of 0.5 (for two classes) when the distribution is perfectly balanced, e.g.,  $P(C1) = P(C2) = 0.5$ .

|    |          |
|----|----------|
| C1 | <b>2</b> |
| C2 | <b>4</b> |

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$Gini = 1 - (2/6)^2 - (4/6)^2 = 0.444$$

### 10.2.2.2 Gini for a collection of nodes

To compute the Gini Index for a collection of nodes after a split, we use a weighted average based on the number of instances in each child node.

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} Gini(i)$$

where:

- ◊  $k$ : number of child nodes after the split
- ◊  $n_i$ : number of instances in child node  $i$
- ◊  $n$ : total number of instances in the parent node before the split
- ◊  $Gini(i)$ : Gini Index of child node  $i$

Gini index is used in decision tree algorithms such as CART, SLIQ, SPRINT

There are some slides concerning optimizations on how the Gini index is computed when dealing with continuous attributes. These optimizations are not included here for brevity. By the way, they probably are already included in some Python library implementing decision trees ☺.

**10.2.2.2.1 Entropy** Entropy is another impurity measure used in decision trees, particularly in the ID3 and C4.5 algorithms. It quantifies the uncertainty or randomness in the class distribution of a node.

It is pretty similar to Gini index, but it tends to be more sensitive to changes in the class distribution. For both of them, the lower the value, the purer the node.

#### Note impurity limitations

Node impurity measures tend to prefer splits that result in large number of partitions, each being small but pure. This may lead to overfitting, as the model captures noise in the training data rather than the underlying patterns.

To mitigate this, we may consider also the number of child nodes created by a split when choosing the best attribute to split on. Or, as it happens in CART, we may use binary splits only, leading to have at most two child nodes per split.

### 10.2.2.3 Comparing measures

The image displays a 2-class problem. We can easily see that there is consistency among the three measures: they all reach their minimum (0) when the node is pure (i.e., all instances belong to one class) and their maximum when the classes are evenly distributed (i.e.,  $P(C1) = P(C2) = 0.5$ ). Furthermore, if a node  $N_1$  has lower entropy than node,  $N_2$ , then the Gini index and error rate of  $N_1$ , will also be lower than that of  $N_2$ .

In some cases Gini may get better, but the misclassification error may stay exactly the same.

## 10.3 Decision Trees Wrap Up

- ◊ Easy to interpret for small-sized trees
- ◊ Accuracy is comparable to other classification techniques for many simple data sets
- ◊ Robust to noise (especially when methods to avoid overfitting are employed)
- ◊ Can easily handle redundant or irrelevant attributes
  - **Redundant** attributes: provide no additional information because they are highly correlated with other attributes
  - **Irrelevant** attributes: provide no useful information for predicting the target class
- ◊ Inexpensive to construct
- ◊ Extremely fast at classifying unknown record
  - Construction cost:  $O(MN\log N)$  having  $M$  attributes and  $N$  records
  - Testing cost:  $O(\log N) = O(w)$  per record, where  $w$  is the depth of the tree
- ◊ Handle Missing Values

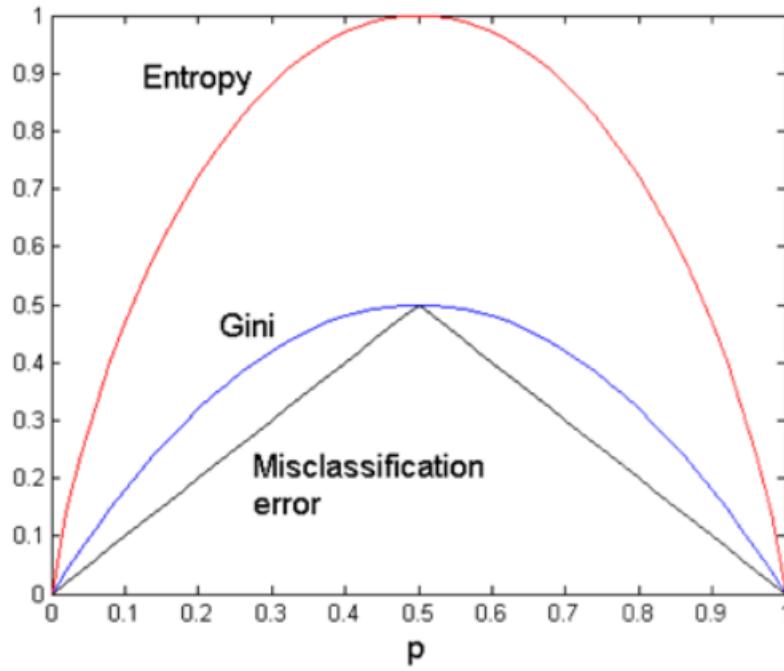


Figure 10.4: Comparing Entropy, Gini and Misclassification error

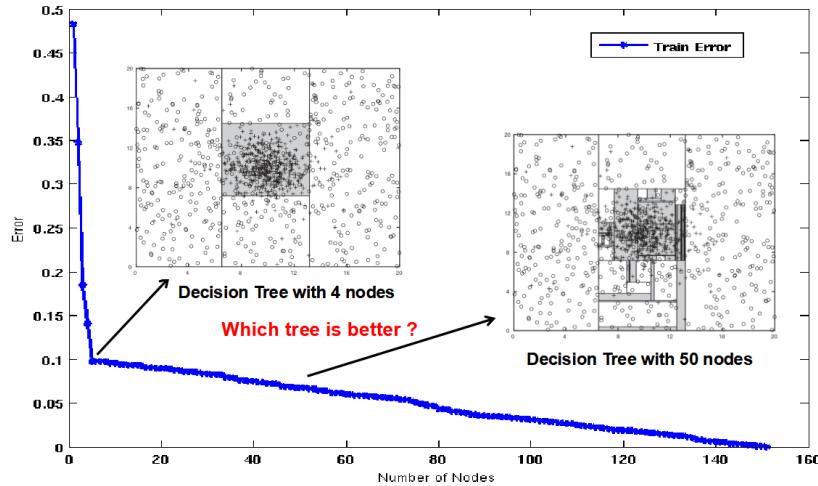


Figure 10.5: The model on the right with 50 nodes probably overfits the data. This may lead to a high error rate on a test set

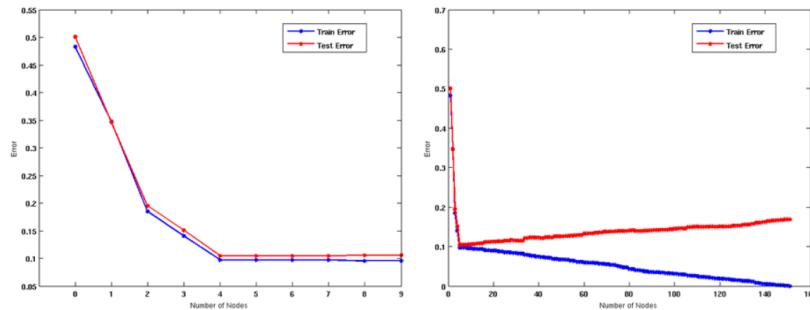


Figure 10.6: Underfitting and overfitting

## 10.4 Model selection

Divide training data into two parts: a **Training set** used for model building and a **Validation set** used for estimating generalization error.

Note that the validation set is not the same as the test set. This is because the test set is used only at the very end to get an unbiased estimate of the generalization error.

This is a nice way to avoid overfitting, but the drawback is that there is less data available for training.

**Definition 10.1 (Occam's Razor)** *Given two models of similar generalization errors, one should prefer the simpler model over the more complex model.*

*This is due to a complex model having a greater chance of being fitted accidentally by errors in data*

Occam's Razor suggests to include a penalty for model complexity in the estimate of generalization error.

### 10.4.1 Evaluating Trees

**Definition 10.2 (Pessimistic Error Estimate of decision tree T)** *Pessimistic Error Estimate of decision tree T with k leaf nodes:*

$$\text{err}_{\text{gen}}(T) = \text{err}(T) + \Omega \times \frac{k}{N_{\text{train}}}$$

- ◊  $\text{err}(T)$ : error rate on all training records
- ◊  $\Omega$ : Relative cost of adding a leaf node
- ◊  $k$ : number of leaf nodes
- ◊  $N_{\text{train}}$ : total number of training record

**Minimum Description Length** (MDL) is another way to implement Occam's Razor.

**Definition 10.3 (Minimum Description Length)** *The best model is the one that minimizes the total description length:*

$$DL(\text{Model}) + DL(\text{Data}|\text{Model})$$

where:

- ◊  $DL(\text{Model})$ : number of bits to describe the model
- ◊  $DL(\text{Data}|\text{Model})$ : number of bits to describe the data when the model is known

#### 10.4.1.1 Estimating Statistical Bounds

We can apply a statistical correction to the training error rate of the model that is indicative of its model complexity.

Doing so, we can obtain a statistical upper bound on the true error rate of the model.

To compute this statistical bound, we need the probability distribution of the training error, which can be either available from data or assumed.

In decision trees, the number of errors committed by a leaf node can be assumed to follow a **binomial distribution**. This is because:

- ◊ Each instance in the leaf is classified as either correct or incorrect (binary outcome)
- ◊ We assume instances are independent
- ◊ The probability of error is constant for all instances in the leaf

Given a leaf node with  $N$  instances and observed error rate  $e$ , we can compute a corrected error estimate  $e'(N, e, \alpha)$  using the binomial confidence interval:

$$e'(N, e, \alpha) = \frac{e + \frac{z_{\alpha/2}^2}{2N} + z_{\alpha/2} \sqrt{\frac{e(1-e)}{N} + \frac{z_{\alpha/2}^2}{4N^2}}}{1 + \frac{z_{\alpha/2}^2}{N}}$$

where  $z_{\alpha/2}$  is the critical value from the standard normal distribution at confidence level  $\alpha$  (e.g.,  $z_{0.125} \approx 1.15$  for  $\alpha = 0.25$ ).

The total generalized error for a tree  $T$  is then:

$$e'(T) = \sum_{\text{leaves } L} N_L \times e'(N_L, e_L, \alpha)$$

This statistical correction penalizes small leaf nodes (which have higher uncertainty) and helps decide whether a split actually improves the model or just fits noise in the training data.

#### 10.4.1.2 Address overfitting

- ◊ Pre-pruning: stop growing the tree earlier
  - Typical stopping conditions for a node are:
    - Stop if all instances belong to the same class
    - Stop if all the attribute values are the same
  - More restrictive conditions:
    - Stop if number of instances is less than some user-specified threshold
    - Stop if expanding the current node does not improve impurity measures (e.g., Gini or information gain).
    - Stop if estimated generalization error falls below certain threshold
- ◊ Post-pruning: grow the full tree, then remove nodes that do not help, following a bottom-up approach
  - If generalization error improves after trimming, replace sub-tree by a leaf node.
  - Class label of leaf node is determined from majority class of instances in the sub-tree
  - Can use MDL for post-pruning

#### 10.4.1.3 Observations

The number of possible decision trees can be very large, many decision tree algorithms employ a heuristic-based approach to guide their search in the vast hypothesis space.

That is splitting the records based on an attribute test that optimizes a certain criterion (e.g., information gain, Gini index, etc.) at each node.

*“But then, how should training records be split at each node?  
And, how should the splitting procedure stop?” —*

We need a method to specify a test condition and a measure to evaluate the quality of a split, as well as a stopping criterion to prevent overfitting, which could be that if all records belong to the same class or if further splitting does not improve the model significantly.

#### 10.4.2 Test conditions

Defining the test conditions to split the attributes depends on the type of attributes, whether they are categorical, binary, nominal, or continuous.

### 10.5 Model Selection

#### 10.5.1 Metrics for Performance evaluation

To evaluate the performance of a classification model we must focus on the predictive capability of a model. Rather than how fast it takes to classify or build models, scalability, etc.

We exploit a **confusion matrix** to evaluate the performance of a classification model, built like the following:

The most widely used metric based on the matrix is the **accuracy**.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Note that accuracy may be misleading when dealing with imbalanced classes. Consider a dataset where 95% of instances belong to class A and only 5% to class B. A model that always predicts class A will have an accuracy of 95%, but it fails to identify any instances of class B.

Another metric is  $C(i|j)$ , the **cost of erroneously classifying** an instance of class  $j$  as class  $i$ .

We have also:

- ◊ **Precision:**

$$\text{Precision} = \frac{TP}{TP + FP}$$

- ◊ **Recall:**

$$\text{Recall} = \frac{TP}{TP + FN}$$

- ◊ **F-measure:**

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- ◊ **Weighted accuracy:**

$$\text{Weighted Accuracy} = \frac{TP + TN}{TP + TN + \beta \times FP + (1 - \beta) \times FN}$$

or

$$\text{Weighted Accuracy} = \frac{w_1 \cdot TP + w_2 \cdot TN}{w_1 \cdot TP + w_2 \cdot TN + w_3 \cdot FP + w_4 \cdot FN}$$

where  $\beta$  is a user-defined parameter that adjusts the weight of false positives and false negatives.

## 10.6 Methods for Performance evaluated

## 10.7 Methods for Model Comparison

To compare two models, we can use statistical tests to determine if the observed differences in performance are significant or due to random chance.

Comparing models is key when tuning hyperparameters, selecting features, or choosing between different algorithms.

### 10.7.1 ROC - Receiver Operating Characteristic

ROC curves plot the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings.

The area under the ROC curve (AUC) provides a single measure of overall model performance.

Every point on the ROC curve represents the performance of each classifier, each having different trade-off between sensitivity and specificity.

### 10.7.2 Significance Testing

Given two models:

- ◊ Model M1:  $accuracy = 85\%$ , tested on 30 instances
- ◊ Model M2:  $accuracy = 75\%$ , tested on 5000 instances

Is M1 really better than M2? Can the difference in performance measure be explained as a result of random fluctuations in the test set?

We can use significance testing to answer these questions.

#### 10.7.2.1 Confidence Interval for accuracy

We can determine a confidence interval for accuracy, in order to estimate the range within which the true accuracy of the model lies, with a certain level of confidence.



# Chapter 11

## Rule-based Classification

### 11.1 Introduction to Rule-based Classifiers

Rule-based classifiers classify records by using a collection of “if...then...” rules. A **rule** has the form  $(\text{Condition}) \rightarrow y$ , where **Condition** is a conjunction of tests on attributes and **y** is the class label.

Examples of classification rules include  $(\text{Blood Type} = \text{Warm}) \wedge (\text{Lay Eggs} = \text{Yes}) \rightarrow \text{Birds}$  and  $(\text{Taxable Income} < 50K) \wedge (\text{Refund} = \text{Yes}) \rightarrow \text{Evade} = \text{No}$ .

#### 11.1.1 Example of Rule-based Classifier

Consider the following rule set:

- ◊ R1:  $(\text{Give Birth} = \text{no}) \wedge (\text{Can Fly} = \text{yes}) \rightarrow \text{Birds}$
- ◊ R2:  $(\text{Give Birth} = \text{no}) \wedge (\text{Live in Water} = \text{yes}) \rightarrow \text{Fishes}$
- ◊ R3:  $(\text{Give Birth} = \text{yes}) \wedge (\text{Blood Type} = \text{warm}) \rightarrow \text{Mammals}$
- ◊ R4:  $(\text{Give Birth} = \text{no}) \wedge (\text{Can Fly} = \text{no}) \rightarrow \text{Reptiles}$
- ◊ R5:  $(\text{Live in Water} = \text{sometimes}) \rightarrow \text{Amphibians}$

#### 11.1.2 Application of Rule-Based Classifier

A rule  $r$  **covers** an instance  $x$  if the attributes of the instance satisfy the condition of the rule. For example, the rule R1 covers a hawk (classifying it as a Bird), while the rule R3 covers the grizzly bear (classifying it as a Mammal).

## 11.2 Rule Coverage and Accuracy

The **coverage** of a rule is the fraction of records that satisfy the antecedent of the rule, while the **accuracy** is the fraction of records that satisfy the antecedent that also satisfy the consequent of the rule. For example, the rule  $(\text{Status} = \text{Single}) \rightarrow \text{No}$  might have a coverage of 40% and an accuracy of 50%.

#### 11.2.1 How does a Rule-based Classifier Work?

When applying a rule-based classifier, different scenarios can occur. A lemur triggers rule R3 and is classified as a mammal. A turtle triggers both R4 and R5, which requires a conflict resolution strategy. A dogfish shark triggers none of the rules, necessitating the use of a default class.

## 11.3 Characteristics of Rule Sets

### 11.3.1 Strategy 1: Mutually Exclusive and Exhaustive Rules

A classifier contains **mutually exclusive rules** if the rules are independent of each other, meaning every record is covered by at most one rule. The classifier has **exhaustive coverage** if it accounts for every possible combination of attribute values, ensuring that each record is covered by at least one rule.

### 11.3.2 Strategy 2: Non-mutually Exclusive and Non-exhaustive Rules

When rules are not mutually exclusive, a record may trigger more than one rule. This conflict can be resolved using either an ordered rule set or an unordered rule set with voting schemes. When rules are not exhaustive, a record may not trigger any rules, and the solution is to use a default class for such cases.

### 11.3.3 Ordered Rule Set

In an ordered rule set, rules are rank ordered according to their priority. Such an ordered rule set is known as a **decision list**. When a test record is presented to the classifier, it is assigned to the class label of the highest ranked rule it has triggered. If none of the rules fired, it is assigned to the default class (typically the majority class).

### 11.3.4 Rule Ordering Schemes

**Rule-based ordering:** Individual rules are ranked based on their quality measured by accuracy, coverage, or size (number of attribute tests in the rule antecedent). The resulting rule set is known as a decision list, where the record  $X$  is classified by the rule with the highest priority and any other rule that satisfies is ignored. Each rule in a decision list implies the negation of the rules that come before it in the list, making rules in a decision list more difficult to interpret.

**Class-based ordering:** Rules that belong to the same class appear together. The classes are sorted in order of decreasing “importance”, such as by decreasing order of prevalence, or alternatively based on the misclassification cost per class. Within each class, the rules are not ordered.

**Unordered rules:** These use a voting schema to resolve conflicts.

## 11.4 Building Classification Rules

There are two main approaches for building classification rules. The **direct method** extracts rules directly from data (examples include RIPPER, CN2, and Holte’s 1R), while the **indirect method** extracts rules from other classification models such as decision trees or neural networks (e.g., C4.5rules).

## 11.5 Direct Method: Sequential Covering

**Algorithm:**

1. Start from an empty rule
2. For each class:
  - i. Grow a rule using the Learn-One-Rule function
  - ii. Remove training records covered by the rule
  - iii. Repeat Step until stopping criterion is met

### 11.5.1 Learn-One-Rule Function

The goal of the Learn-One-Rule function is to extract a classification rule covering many positive records and none (or few) negative ones. Since finding the optimal rule requires high computational time, a greedy strategy is employed by refining an initial rule based on some evaluation measure. Rules are extracted one class at a time, and the criterion for deciding the order of the class to consider depends on class prevalence and misclassification error for a given class.

### 11.5.2 Rule Growing Strategies

Two common strategies exist for growing rules. The **general-to-specific** approach starts with an empty rule  $\{ \}$  and iteratively adds conjuncts like  $Refund = No$ ,  $Status = Single$ ,  $Income > 80K$  until the rule is sufficiently specific. Conversely, the **specific-to-general** approach starts with a specific rule covering a selected record (e.g.,  $Refund = No$ ,  $Status = Single$ ,  $Income = 85K$  with Class=Yes) and then generalizes it by removing conditions.

## 11.6 Rule Evaluation for Growing Rules

### 11.6.1 Based on Rule Coverage

Evaluation measures can be based on the coverage of the rule with respect to records with class  $c$ .

### 11.6.2 Based on Support Count: FOIL's Information Gain

FOIL's Information Gain compares an initial rule  $R_0 : \{ \} \Rightarrow \text{class}$  with a rule after adding a conjunct  $R_1 : \{A\} \Rightarrow \text{class}$ :

$$\text{Gain}(R_0, R_1) = p_1 \times \left[ \log_2 \frac{p_1}{p_1 + n_1} - \log_2 \frac{p_0}{p_0 + n_0} \right]$$

where  $p_0$  and  $n_0$  are the number of positive and negative instances covered by  $R_0$ , respectively, and  $p_1$  and  $n_1$  are the corresponding counts for  $R_1$ . FOIL (First Order Inductive Learner) is an early rule-based learning algorithm.

### 11.6.3 Based on Statistical Test: Likelihood Ratio Statistic

Given a rule, we can compute the Likelihood ratio statistic  $R = 2 \sum_i f_i \log \frac{f_i}{e_i}$ , where  $f_i$  is the number of records covered by the rule and  $e_i$  is the expected frequency of a rule that makes random predictions. A large  $R$  value suggests that the number of correct predictions made by the rule is significantly larger than that expected by random guessing.

#### Example:

Dataset: 60 positive records and 100 negative records

*Rule 1:* covers 50 positive records and 5 negative examples

- ◊ Expected frequency for the positive class is  $e_+ = 55 \times 60/160 = 20.625$
- ◊ Expected frequency for the negative class is  $e_- = 55 \times 100/160 = 34.375$

*Rule 2:* covers 2 positive records and no negative examples

- ◊ Expected frequency for the positive class is  $e_+ = 2 \times 60/160 = 0.75$
- ◊ Expected frequency for the negative class is  $e_- = 2 \times 100/160 = 1.25$

## 11.7 Direct Method: RIPPER

### 11.7.1 For 2-class Problem

For a 2-class problem, RIPPER chooses one of the classes as the positive class and the other as the negative class. It then learns rules for the positive class, while the negative class becomes the default class.

### 11.7.2 For Multi-class Problem

For multi-class problems, RIPPER orders the classes according to increasing class prevalence (the fraction of instances that belong to a particular class). It learns the rule set for the smallest class first, treating the rest as the negative class, then repeats the process with the next smallest class as the positive class.

### 11.7.3 Growing a Rule in RIPPER

RIPPER starts from an empty rule and adds conjuncts as long as they improve FOIL's information gain. It stops when the rule no longer covers negative examples and then immediately prunes the rule using incremental reduced error pruning. The measure for pruning is  $v = (p - n)/(p + n)$ , where  $p$  is the number of positive examples covered by the rule in the validation set and  $n$  is the number of negative examples. The pruning method deletes any final sequence of conditions that maximizes  $v$ .

### 11.7.4 Building a Rule Set in RIPPER

RIPPER uses a sequential covering algorithm that finds the best rule covering the current set of positive examples and eliminates both positive and negative examples covered by the rule. Each time a rule is added to the rule set, the algorithm computes the new description length and stops adding new rules when this description length is  $d$  bits longer than the smallest description length obtained so far.

## 11.8 Minimum Description Length (MDL)

The Minimum Description Length principle is expressed as  $\text{Cost}(\text{Model}, \text{Data}) = \text{Cost}(\text{Data}|\text{Model}) + \alpha \times \text{Cost}(\text{Model})$ , where Cost is the number of bits needed for encoding. The goal is to search for the least costly model. Here,  $\text{Cost}(\text{Data}|\text{Model})$  encodes the misclassification errors, while  $\text{Cost}(\text{Model})$  uses node encoding (number of children) plus splitting condition encoding.

## 11.9 Indirect Method: C4.5rules

### 11.9.1 Algorithm

C4.5rules extracts rules from an unpruned decision tree. For each rule  $r : A \rightarrow y$ , it considers an alternative rule  $r' : A' \rightarrow y$  where  $A'$  is obtained by removing one of the conjuncts in  $A$ . The algorithm compares the pessimistic error rate for  $r$  against all  $r'$ s and prunes the rule if one of the alternative rules has a lower pessimistic error rate. This process repeats until generalization error can no longer be improved. After removing duplicate rules, C4.5rules uses class ordering instead of ordering individual rules. For multi-class problems, rather than assigning a default class to test records not covered by any rule, C4.5rules looks for the rule that most closely matches the record.

### 11.9.2 Pessimistic Error Estimate

The Pessimistic Error Estimate of a rule set  $T$  with  $k$  rules is computed as  $\text{err}(T) + W \times \frac{k}{N_{\text{train}}}$ , where  $\text{err}(T)$  is the error rate on all training records,  $W$  is a trade-off hyper-parameter representing the relative cost of adding a rule,  $k$  is the number of rule nodes, and  $N_{\text{train}}$  is the total number of training records.

### 11.9.3 Class Ordering in C4.5rules

C4.5rules orders subsets of rules rather than individual rules, using class ordering. Each subset is a collection of rules with the same rule consequent (class). The algorithm computes the description length of each subset as  $L(\text{error}) + g \times L(\text{model})$ , where  $g$  is a parameter that takes into account the presence of redundant attributes in a rule set (with a default value of 0.5).

## 11.10 Advantages of Rule-Based Classifiers

Rule-based classifiers have characteristics quite similar to decision trees: they are as highly expressive as decision trees, easy to interpret, have comparable performance, and can handle redundant attributes. Additionally, they are better suited for handling imbalanced classes, though they are harder to handle missing values in the test set.

## 11.11 Example: C4.5 vs C4.5rules vs RIPPER

### 11.11.1 C4.5rules

- ◊  $(\text{Give Birth} = \text{No}, \text{Can Fly} = \text{Yes}) \rightarrow \text{Birds}$
- ◊  $(\text{Give Birth} = \text{No}, \text{Live in Water} = \text{Yes}) \rightarrow \text{Fishes}$
- ◊  $(\text{Give Birth} = \text{Yes}) \rightarrow \text{Mammals}$
- ◊  $(\text{Give Birth} = \text{No}, \text{Can Fly} = \text{No}, \text{Live in Water} = \text{No}) \rightarrow \text{Reptiles}$
- ◊  $() \rightarrow \text{Amphibians}$

### 11.11.2 RIPPER

- ◊  $(\text{Live in Water} = \text{Yes}) \rightarrow \text{Fishes}$
- ◊  $(\text{Have Legs} = \text{No}) \rightarrow \text{Reptiles}$
- ◊  $(\text{Give Birth} = \text{No}, \text{Can Fly} = \text{No}, \text{Live In Water} = \text{No}) \rightarrow \text{Reptiles}$
- ◊  $(\text{Can Fly} = \text{Yes}, \text{Give Birth} = \text{No}) \rightarrow \text{Birds}$
- ◊  $() \rightarrow \text{Mammals}$

### 11.11.3 Performance Comparison

When comparing C4.5, C4.5rules, and RIPPER, the algorithms produce different classification results. The confusion matrices reveal that RIPPER may have different error patterns, particularly in distinguishing between classes like Amphibians, Reptiles, and Mammals. While C4.5 and C4.5rules tend to produce similar results due to their shared origin from decision trees, RIPPER's sequential covering approach leads to a different rule structure and potentially different classification behavior.

