



LINFO2364 Mining Patterns in Data

ISSAMBRE L'HERMITE DUMONT

This summary may not be up-to-date, the newer version is available at this address:
<https://github.com/SimonDesmidt/Syntheses>

Academic year 2025-2026 - Q2



UCLouvain

Contents

1	Introduction	2
1.1	Definitions	2
2	Preprocessing	4
2.1	Useful statistical values	4
2.2	Data plot	5
2.3	Distances and similarities	6
2.4	Preprocessing techniques	6
3	Itemset mining	8
3.1	Definitions	8
3.2	Apriori algorithm	9

Introduction

In our data-driven world, the ability to extract meaningful information from vast datasets is crucial. Understanding all the aspect of this discipline is essential to derive those meaningful information, and this is the goal of this course. First, we need to define some key concepts.

1.1 Definitions

Definition 1.1. A pattern is a recurring structure in a dataset.

Patterns can be simple or complex, relevant or irrelevant. Their advantages is that they are interpretable. When found, relevant patterns can be used to make predictions, to understand the underlying structure of the data, and to make informed decisions.

Definition 1.2. Data mining is the process of discovering interesting patterns, models, and other kinds of knowledge in large data sets.

1.1.1 Type of data

We can mine data out of various types of structure of data:

- **Tabular data:** Data is organized in rows and columns. Example: spreadsheets, databases.
- **Sequences:** Data points are ordered in a sequence. Example: DNA sequences, text data.
- **Graphs, trees, networks:** Data is represented as nodes and edges. Example: social networks, web graphs.

Those structures can be discrete, continuous, enumerable data, etc. Those structures, can be combined to form more complex data types. And they can be highly structured, semi-structured, or unstructured.

Definition 1.3. Highly structured data are relational databases, with uniform record or table-like structures, with a fixed set of well-defined attributes. This is rarely the case in real-world data.

Definition 1.4. Semi-structured data are not as structured as in relational databases, but presents some structure with clearly defined semantic meaning. For example:

- **Transactional dataset:** structured into transactions, but each transaction is an unstructured set of values
- **Sequence data set:** unstructured collection of ordered sequences of values
- **Graphs:** set of nodes connected by a set of edges, with edges labelled given some semantic

Definition 1.5. Unstructured data have no predefined structure or organization. For example: text documents, images, audio files, videos.

Those requires advanced techniques to extract patterns, like deep learning or domain-specific methods. We can also categorize data based on the way they are generated:

- **Stored data:** Data is collected and stored in a finite set.
- **Streamed data:** Data is continuously generated and updated over time. Example: video surveillance, etc.

1.1.2 Types of data mining

Techniques and algorithms may vary depending on the data but also on way we will use mined patterns. We can categorize data mining tasks into two main types:

- **Descriptive data mining:** finds patterns that characterize the properties of the data.
- **Predictive data mining:** finds patterns that can be used to make predictions by induction.

We can identify patterns by multiples ways:

- **Frequent patterns:** patterns that identify a recurring structure in the data.
- **Associations patterns:** patterns that show rules of implications between different attributes
- **Correlations:** patterns that has positive or negative correlation between different attributes.

We can uses patterns for predictions:

- **Classification:** assign new data to classes based on similarities with historic data.
- **Regression:** predict numerical values based on the new data and the historic data.
- **Feature selection:** identify the most relevant features.

Preprocessing

To analyse data, and retrieve meaningful patterns, data often needs to be preprocessed. This step is crucial as raw data is often incomplete, inconsistent, and noisy. Preprocessing improves the quality of the data and the efficiency of the mining algorithms.

2.1 Useful statistical values

Definition 2.1. The mean (or average) of a dataset is the sum of all values divided by the number of values.

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad (2.1)$$

Definition 2.2. The midrange is the average of the largest and smallest values in the set.

Definition 2.3. The variance of a dataset measures is defined like this:

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (2.2)$$

Definition 2.4. The standard deviation of a dataset is the square root of the variance:

$$\sigma = \sqrt{\sigma^2} = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2} \quad (2.3)$$

Definition 2.5. The mode for a set of data is the value that occurs most frequently. When there are multiple values with the same highest frequency, the dataset can be unimodal, bimodal, trimodal or multimodal.

Definition 2.6. Quantiles are points taken at regular intervals of a data distribution, dividing it into essentially equal-sized consecutive sets. The k th q -quantile for a given data distribution is the value x such that at most $\frac{k}{q}$ of the data values are less than x and at most $\frac{q-k}{q}$ of the data values are more than x , where k is an integer such that $0 < k < q$. There are $q - 1$ q -quantiles.

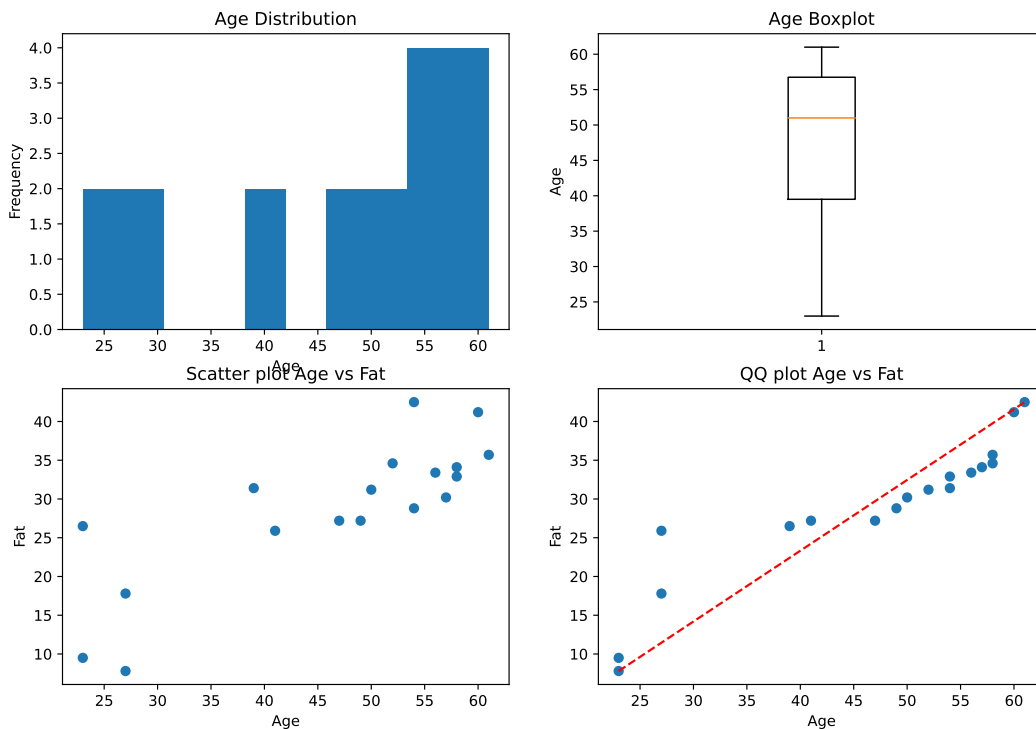
Usually we use quartiles ($q = 4$), where Q_2 is the median, or percentiles ($q = 100$).

2.2 Data plot

Plotting the data can be useful to observe some characteristics of the dataset. we can represent the data under multiples forms:

- **Histogram** is a bar chart representing the count of values present in each bucket. The buckets can be defined for each values, for ranges of values (equal-width) or for buckets containing the same number of values (equal-frequency).
- **Boxplots** are a popular way of visualizing a distribution. Typically, the ends of the box are at the quartiles so that the box length is the interquartile range. The median is marked by a line within the box. Two lines (called whiskers) outside the box extend to the smallest (minimum) and largest (maximum) observations. To highlight outliers, the whiskers are usually extended to extreme low and high observations only if these values are less than 1.5 times the interquartile range beyond the quartiles. Otherwise, the whiskers terminate at the most extreme observations occurring within 1.5 times the interquartile range.
- **Scatter plot**, each pair of values is treated as a pair of coordinates in an algebraic sense and plotted as points in the plane. It helps visually determine whether there is a relationship, pattern, or trend between the two numeric attributes.
- **Quantile-quantile plot**, or q-q plots are graphs that plot the quantiles of one univariate distribution against the corresponding quantiles of another.

Here is an example of the four plots for the same dataset:



2.3 Distances and similarities

To obtain relation between data, we can use distance and similarity measures. Distance measures the dissimilarity between two data points, while similarity measures the closeness between them. The choice of distance or similarity measure depends on the type of data and the specific problem at hand.

Definition 2.7. The norm of a vector x is:

$$||x|| = \sqrt{\sum_{i=1}^n x_i^2} \quad (2.4)$$

Definition 2.8. The Minkowski distance between two vectors x and y is:

$$d(x, y) = \sqrt[p]{\sum_{i=1}^n |x_i - y_i|^p} \quad (2.5)$$

It is a generalization of the Euclidean distance ($p = 2$) and the Manhattan distance ($p = 1$).

Definition 2.9. The supremum distance (Minkowski when $p \rightarrow \infty$) between two vectors x and y is:

$$d(x, y) = \max_i |x_i - y_i| \quad (2.6)$$

So two vectors are similar if their distances is close to 0.

Definition 2.10. The cosine similarity between two vectors x and y is:

$$s(x, y) = \frac{x \cdot y}{||x|| \cdot ||y||} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \cdot \sqrt{\sum_{i=1}^n y_i^2}} \quad (2.7)$$

Here two vectors are similar if their cosine similarity is close to 1, and dissimilar if it is close to 0.

2.4 Preprocessing techniques

First we can smooth data using binning methods. Binning methods smooth a sorted data value by consulting its neighborhood. We can smooth using bins by partitioning the sorted data into a number of equal-frequency or equal-width bins, and then replacing each data value in a bin by the mean, median, or boundary values of the bin. Another type of method to manipulate data is to use normalization. There exists multiple normalization techniques:

- **Vector normalization:** scales a vector x by it's norm to have a unit norm vector:

$$x' = \frac{x}{||x||} \quad (2.8)$$

- **Min-max normalization:** scales the vector linearly to fit in a specific range $[new_min, new_max]$

$$x' = \frac{(x - \min(x))(new_max - new_min)}{\max(x) - \min(x)} + new_min \quad (2.9)$$

- **Z-score normalization:** (or standardization) rescales the data to have a mean of 0 and a standard deviation of 1:

$$x' = \frac{x - \bar{x}}{\sigma} \quad (2.10)$$

- **Normalization by decimal scaling:**

$$x' = \frac{x}{10^j} \quad \text{where } j \text{ is the smallest integer such that } \max(|x'|) < 1 \quad (2.11)$$

j can be computed as $j = \lceil \log_{10}(\max(|x|)) + \epsilon \rceil$ with ϵ a small value to assure $\max(|x'|) < 1$.

For missing values, we can (last 3 methods introduce bias in the data):

- Ignore the tuple: usually done when the class label is missing or when multiple values are missing in the same tuple.
- Fill the missing value manually: time-consuming if lots of missing data
- Use a global constant (Unknown or $-\infty$): mining algorithms might mistakenly think that they form an interesting concept by having this "missing" value in common
- Use a measure of the central tendency of the attribute (e.g., mean or median)
- Use the measure of the central tendency for all samples belonging to the same class
- Use the most probable value to fill: using regression, inference-based tools, or decision tree induction.

We can also reduce the size of the data using sampling:

- **SRSWOR** (Simple random sample without replacement of size s): this sampling is created by drawing samples from D , and every time a sample is drawn, it is not to be placed back into the data set D .
- **SRSWR** (Simple random sample with replacement of size s): this sampling is similar to SRSWOR, except that when a sample is drawn, it can be redrawn again, potentially.
- **Stratified sampling:** If D is divided into mutually disjoint parts called strata, a stratified sample of D is generated by obtaining a sample at each stratum. This helps ensure a representative sample, each stratum need to be defined manually.

Itemset mining

3.1 Definitions

Definition 3.1. An itemset is a set of items and an itemset containing k items is called a k -itemset.

Frequent itemset mining finds associations and correlations between items in large transactional or relational datasets. It is widely used in market basket analysis for example. It is the analysis of the buying habits of customers by finding associations between the different items that customers place in their shopping basket.

Definition 3.2. A transactional dataset is a collection of transactions, where each transaction is a set of items.

With \mathcal{L} the set of possible items, \mathcal{T} the set of possible transactions, a transactional dataset \mathcal{D} can be seen as the function $\mathcal{D} : \mathcal{T} \rightarrow 2^{\mathcal{L}}$. It can also be represented as a binary matrix, where rows represent transactions and columns represent items, mathematically it is like the function $\mathcal{D} : \mathcal{T} \times \mathcal{L} \rightarrow \{0, 1\}$.

Definition 3.3. An itemset $l \subset \mathcal{L}$ covers (or matches) a transaction $t \subset \mathcal{T}$ if every item from l is in t . Consider the match function, with $D(t, l)$ the function representing the transactional dataset:

$$\text{match}(l, t) = \begin{cases} 1 & \text{if } \forall i \in l, D(t, i) = 1 \\ 0 & \text{otherwise} \end{cases} \quad (3.1)$$

Definition 3.4. The occurrence frequency of an itemset (also called frequency, support count, count or absolute support) is the number of transactions that contain the itemset (i.e., the size of the cover)

$$\text{support}_{\text{count}}(l) = \sum_{t \in \mathcal{T}} \text{match}(l, t) \quad (3.2)$$

Definition 3.5. The relative support of an itemset l in a database \mathcal{D} is the percentage of transactions containing the itemset:

$$\text{support}_{\mathcal{D}}(l) = \frac{\text{support}_{\text{count}}(l)}{|\mathcal{T}|} \quad (3.3)$$

Definition 3.6. An association rule represents an implication of the form $X \Rightarrow Y$, where X and Y are itemsets.

Definition 3.7. The rule support is the support of the itemset containing all the elements of the rule $\text{support}_{\text{count}}(X \Rightarrow Y) = \text{support}_{\text{count}}(X \cup Y)$.

We consider an association rule or an itemset as frequent if its support is greater than a user-defined minimum support threshold θ .

Definition 3.8. The rule confidence is the percentage of transactions containing X that contains Y too

$$\text{confidence}(X \Rightarrow Y) = P(Y|X) = \frac{\text{support}_{\text{count}}(X \cup Y)}{\text{support}_{\text{count}}(X)} \quad (3.4)$$

3.2 Apriori algorithm

The key idea behind the Apriori algorithm is that all nonempty subsets of a frequent itemset must also be frequent.

Algorithm 1 Apriori algorithm

```

1: Input: Transactional dataset  $\mathcal{D}$ , minimum support threshold  $\theta$ 
2: Output: Set  $F$  of frequent itemsets
3:  $F = \emptyset$ 
4:  $C_1 = \{\{i\} : i \in \mathcal{L}\}$ 
5:  $L_1 = \{c \in C_1 : \text{support}_{\text{count}}(c) \geq \theta\}$ 
6:  $F = F \cup L_1$ 
7:  $k = 2$ 
8: while  $L_k$  is not empty do
9:    $C_k = \text{apriori\_gen}(L_{k-1})$ 
10:   $L_k = \{c \in C_k : \text{support}_{\text{count}}(c) \geq \theta\}$ 
11:   $F = F \cup L_k$ 
12:   $k = k + 1$ 
13: end while
14: return  $F$ 

```

The $\text{apriori_gen}(L_{k-1})$ function generates candidate k -itemsets from the frequent $(k-1)$ -itemsets L_{k-1} by joining them and pruning those that have infrequent subsets.

Algorithm 2 $\text{apriori_gen}(L_{k-1})$

```

1: Input: Set  $L_{k-1}$  of frequent  $(k-1)$ -itemsets
2: Output: Set  $C_k$  of candidate  $k$ -itemsets
3:  $C_k = \emptyset$ 
4: for each  $l_1 \in L_{k-1}$  do
5:   for each  $l_2 \in L_{k-1}$  do
6:     if  $l_1[k-2] = l_2[k-2]$  and  $l_1[k-1] < l_2[k-1]$  then
7:        $c = l_1 \cup l_2$ 
8:       if all  $(k-1)$ -subsets of  $c$  are in  $L_{k-1}$  then
9:          $C_k = C_k \cup \{c\}$ 
10:      end if
11:    end if
12:  end for
13: end for
14: return  $C_k$ 

```

The number of candidate itemsets is bounded by $\mathcal{O}\left(\binom{|\mathcal{L}|}{k}\right)$.

The two main weaknesses of the Apriori algorithm are:

- Potential generation of huge candidate sets
- Repeated scan of the database and checks for pattern matching to the transactions

The Apriori algorithm can be improved using:

- **Hash-based techniques:** Instead of building L_2 from the join operation, build it similarly to L_1 . When scanning the transactions, generate the 2-itemset subset of the transaction, hash them and map them to corresponding buckets containing counters to be increased. L_2 is the set of the 2-itemsets associated to buckets of a count higher than θ
- **Transaction reduction:** A transaction that does not contain any frequent k -itemsets cannot contain any frequent $(k+1)$ -itemsets. They can be flagged and avoided at next steps.
- **Partitioning:** The dataset can be partitioned in N non-overlapping sub-database. Given θ , the minimum support count threshold for the whole database, $\theta' = \lfloor \frac{\theta}{N} \rfloor$ is the threshold for each partition. The candidate set C_k is the union of all local frequent k -itemset of each partition as any itemset that is potentially frequent with respect to D must occur as a frequent itemset in at least one of the partitions.
- **Sampling:** The frequent itemsets of a sample of the database are computed to serve as C_k . Some degree of accuracy is traded for efficiency, thus a lower value than θ' is used to counterbalanced the effect.