



Fouille de données

▷ Data

Philippe Lenca et Romain Billot

`philippe.lenca@telecom-bretagne.eu`

Telecom Bretagne
2015-2016





Outline

Data

Attributes

Datasets

Characteristics

Bibliographie

- 1 Data description
- 2 Attribute types
- 3 Dataset types
- 4 Dataset characteristics
- 5 Bibliographie



Outline

Data

Attributes

Datasets

Characteristics

Bibliographie

- 1 Data description
- 2 Attribute types
- 3 Dataset types
- 4 Dataset characteristics
- 5 Bibliographie



Readings

Data

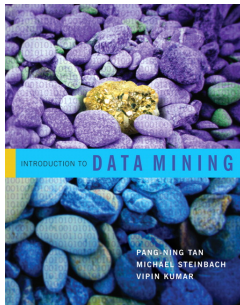
Attributes

Datasets

Characteristics

Bibliographie

Introduction to Data Mining: Pang-Ning Tan, Michael Steinbach, Vipin Kumar [?]





What is a Dataset?

Data

Attributes

Datasets

Characteristics

Bibliographie

A collection of data objects described by attributes

- **Objects:**
the **atomic elements** from a dataset

(examples, records, prototypes, objects, cases, points, samples, instances)

- **Attributes:**
a **property** of an object

(features, variables, fields, characteristics)

	Attributes				
	<i>TID</i>	Marital status	Sex	Income	Age
Objects	id ₁	single	M	100k€	22
	id ₂	single	F	150k€	28
	id ₃	married	F	120k€	32
	id ₄	single	F	250k€	42
	id ₅	divorced	M	95k€	25
	id ₆	single	M	120k€	55
	id ₇	married	F	95k€	33
	id ₈	divorced	M	150k€	47
	id ₉	single	F	100k€	29
	id ₁₀	single	M	80k€	26

Objects × Attributes

↪ Concepts: content inside the data that can be learned.



Outline

Data

Attributes

Datasets

Characteristics

Bibliographie

- 1 Data description
- 2 Attribute types**
- 3 Dataset types
- 4 Dataset characteristics
- 5 Bibliographie



Attribute types

Data

Attributes

Datasets

Characteristics

Bibliographie

Attribute values are numbers or symbols assigned to an attribute

- An attribute can be mapped to different **measurement scale**

Temperature can be measured in Celsius or Fahrenheit

Age can be measured in month or year

Height can be measured in feet or meter

- Different attributes can be mapped to the same scale

ID and age could be integers

Customer satisfaction on price and quality can be mapped on 'bad' 'reasonable' 'good'

↪ But properties of attribute values can be different.
Readings: [?], [?].



Attribute types

Data

Attributes

Datasets

Characteristics

Bibliographie

<i>TID</i>	Marital status	Sex	Income	Age	Relation
id ₁	single	M	100k€	22	▼
id ₂	single	F	150k€	28	▲
id ₃	married	F	120k€	32	▲
id ₄	single	F	250k€	42	—
id ₅	divorced	M	95k€	25	▼
id ₆	single	M	120k€	55	▲
id ₇	married	F	95k€	33	▲
id ₈	divorced	M	150k€	47	▲
id ₉	single	F	100k€	29	▲
id ₁₀	single	M	80k€	26	—

<i>TID</i>	Marital status	Sex	Income	Age	Relation
id ₁	S	1	135k\$	22	▼
id ₂	S	2	202k\$	28	▲
id ₃	M	2	162k\$	32	▲
id ₄	S	2	337k\$	42	—
id ₅	D	1	128k\$	25	▼
id ₆	S	1	162k\$	55	▲
id ₇	M	2	128k\$	33	▲
id ₈	D	1	202k\$	47	▲
id ₉	S	2	135k\$	29	▲
id ₁₀	S	1	108k\$	26	—

↔ Attribute types and scales of measurement have to be properly specified.



Levels of measurement

Data

Attributes

Datasets

Characteristics

Bibliographie

Four main levels of measurement

- nominal
- ordinal
- interval
- ratio



Levels of measurement

Data

Attributes

Datasets

Characteristics

Bibliographie

Nominal scale

Values are names or labels (identifiers):

- no ordering, no distance measure, no relation among values
- only equality tests and set membership can be performed
- categorical variables
- boolean as a special case

Examples

- ID numbers, ZIP codes
- color, marital status, sex

↔ Standard set structure (unordered).



Levels of measurement

Data

Attributes

Datasets

Characteristics

Bibliographie

Ordinal scale

Values with an order:

- no distance measure (no relative size or degree of difference)
- relation w.r.t. the order, addition does not make sense
- with a preference system one can transform an nominal scale toward an ordinal one

Examples

- rankings ('bad' < 'medium' < 'good')
- height ('short' < 'medium' < 'tall')

↔ Totally ordered set.



Levels of measurement

Data

Attributes

Datasets

Characteristics

Bibliographie

Interval scale

Values ordered and measured in fixed and equal units:

- zero point is not defined
- difference of two values makes sense
- sum or product does not make sense

Examples

- temperatures in Celsius or Fahrenheit
- calendar dates

↪ Affine line.



Levels of measurement

Data

Attributes

Datasets

Characteristics

Bibliographie

Ratio scale

Values ordered and measured on a scale with a zero point:

- estimation of the ratio between a magnitude of a continuous quantity and a unit magnitude of the same kind
- ratio quantities are treated as real numbers
- all mathematical operations are allowed

Examples

- temperatures in Kelvin
- distance, time

↔ One-dimensional vector space.



Levels of measurement

Data

Attributes

Datasets

Characteristics

Bibliographie

Brief synthesis

		Nominal	Ordinal	Interval	Ratio
distinctness	$= \neq$	▲	▲	▲	▲
order	$< >$		▲	▲	▲
addition	$+ -$			▲	▲
multiplication	$\times /$				▲

↪ Most schemes consider just two or three levels of measurement: nominal, ordinal and numerical. See also which scales are considered by the software.



Levels of measurement

Data

Attributes

Datasets

Characteristics

Bibliographie

Brief synthesis

Type	Description	Examples	Operations
Nominal	Distinction	ID, color, sex	mode, entropy, contingency correlation, χ^2 test
Ordinal	Order	grades	median, percentiles, rank correlation, run tests, sign tests
Interval	Differences	temperatures in C or F	mean, standard deviation, Pearson's correlation, t and F tests
Ratio	Differences and ratio	age	geometric mean, harmonic mean, percent variation

↪ Levels of measurement & legal operations.



Levels of measurement

Data

Attributes

Datasets

Characteristics

Bibliographie

Brief synthesis

Type	Transformation	Comments
Nominal	One to One	Any permutation of values
Ordinal	Monotonic increasing	$v_{new} = f(v_{old})$
Interval	Positive linear (affine)	$v_{new} = a \times v_{old} + b$
Ratio	Positive similarities	$v_{new} = a \times v_{old}$

↪ Levels of measurement & legal transformations.



Discrete and Continuous Attributes

Data

Attributes

Datasets

Characteristics

Bibliographie

Discrete Attribute

- finite or countably infinite set of values
- often represented as integer variables

Continuous Attribute

- infinite set of values
- real numbers
- in practice it can only be measured and represented using a finite number of digits (floating-point variables)

↔ Under software constraints.



Why Specifying Attribute Types?

Data

Attributes

Datasets

Characteristics

Bibliographie

Consistency, efficiency analysis

- check for valid values, missing values
- make adequate comparisons and legal operations
- express the best possible patterns

↪ Under software constraints.



Attribute roles

Data

Attributes

Datasets

Characteristics

Bibliographie

Main roles

- id
- predictive
- label, class
- descriptive

↔ Under software constraints.



Outline

Data

Attributes

Datasets

Characteristics

Bibliographie

- 1 Data description
- 2 Attribute types
- 3 Dataset types**
- 4 Dataset characteristics
- 5 Bibliographie



Dataset types

Data

Attributes

Datasets

Characteristics

Bibliographie

Three main dataset types

- record e.g.
 - data matrix
 - document data
 - transaction Data
- graph e.g.
 - World Wide Web
 - molecular structures
- ordered e.g.
 - spatial data
 - temporal data
 - sequential data
 - genetic Sequence data



Dataset types

Data

Attributes

Datasets

Characteristics

Bibliographie

Record data: tables, matrix.

A collection of records described with a fixed set of attributes.

<i>TID</i>	Marital status	Sex	Income	Age	Relation
id ₁	single	M	100k€	22	▼
id ₂	single	F	150k€	28	▲
id ₃	married	F	120k€	32	▲
id ₄	single	F	250k€	42	—
id ₅	divorced	M	95k€	25	▼
id ₆	single	M	120k€	55	▲
id ₇	married	F	95k€	33	▲
id ₈	divorced	M	150k€	47	▲
id ₉	single	F	100k€	29	▲
id ₁₀	single	M	80k€	26	—



Dataset types

Data

Attributes

Datasets

Characteristics

Bibliographie

Record data: document data

Each document becomes a 'term' vector:

- each term is a component (attribute) of the vector
- the value of each component is the number of times the corresponding term occurs in the document

<i>ID</i>	processus	décision	données
D_1	2	1	4
D_2	1	2	1
D_3	1	3	2



Dataset types

Data

Attributes

Datasets

Characteristics

Bibliographie

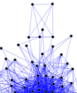
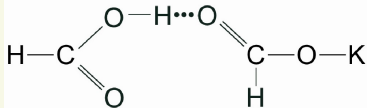
Record data: transaction data

- each record (transaction) involves a set of items

<i>TID</i>	Items
id_1	Bread, Coke, Milk
id_2	Bread, Coke
id_3	Coke, Milk
id_4	Beer, Bread, Diaper, Milk
...	...



Generic graph.





Dataset types

Data

Attributes

Datasets

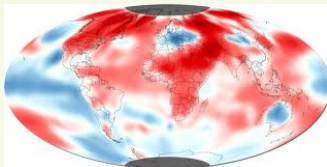
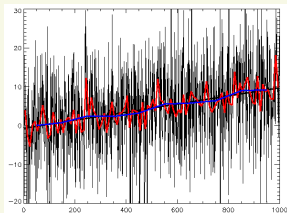
Characteristics

Bibliographie

Ordered data

(CE) (A) (AE)
(BD) (C) (E)
(AB) (D) (CE)

AAATCGGACGCCGGGCTATA
CAACCGTACCCCGGGCTATA
GGTCCGTACCCCGGGCTAAA





Outline

Data

Attributes

Datasets

Characteristics

Bibliographie

- 1 Data description
- 2 Attribute types
- 3 Dataset types
- 4 Dataset characteristics**
- 5 Bibliographie



Dataset characteristics

Data

Attributes

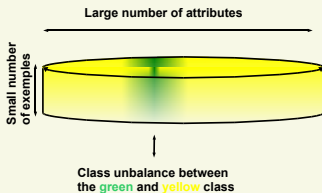
Datasets

Characteristics

Bibliographie

Main characteristics

- **dimensionality**: number of attributes
- **class label unbalance**: prior class probabilities
- **size**: number of samples
- **quality**: missing values, outliers, precision
- **sparsity**: presence of attributes





Dataset characteristics

Data

Attributes

Datasets

Characteristics

Bibliographie

Main troubles

- dimensionality: curse of dimensionality, most existing algorithms only work well on data in **large** quantity, with a **reasonable** number of attributes.
- class label balance: high-unbalanced data (large differences in prior class probabilities) & **bad performance** of classifiers for the **minority class**.
- size: **robustness**, **complexity** issues,
- quality: **quality** of the outputs
- sparsity: presence of attributes, mining of associations among attribute sets only works if they **frequently** co-occur

↔ Challenges to the data mining community, especially when these characteristics are mixed together.



Dataset characteristics

Data

Attributes

Datasets

Characteristics

Bibliographie

Main characteristics

- **dimensionality**: number of attributes
- **class label unbalance**: prior class probabilities
- **size**: number of samples
- **quality**: missing values, outliers, precision
- **sparsity**: presence of attributes

↔ Challenges to the data mining community, especially when these characteristics are mixed together.

A first step of the analysis is to take a 'picture' of the dataset.



Outline

Data

Attributes

Datasets

Characteristics

Bibliographie

- 1 Data description
- 2 Attribute types
- 3 Dataset types
- 4 Dataset characteristics
- 5 Bibliographie**



Bibliographie

Data

Attributes

Datasets

Characteristics

Bibliographie