

# Minería de datos

## U2 Entendimiento de datos

---

Héctor Maravillo

## Section 1

# Knowledge representation

---

# Input

The **input** of a learning methods takes the form of:

- **Concepts.** A **concept description** refers to the entity or idea that is intended to be learned. This should be **intelligible**, meaning that it can be understood and communicable, and **operational**, that is, it can be applied to actual examples.
- **Instances.** Each instance is an **individual**, independent **example** of the concept to be learned.
- **Attributes.** Each instance is characterized by the values of attributes that measure different aspects of the instance.

# Datasets

Typically, a **dataset** in machine learning is represented as a matrix of instances versus attributes, which in database terms is a single relation, or a **flat file**.

Order ID	Order ID	Ship Date	Customer ID	Customer Name	State	Region	Product ID
1	CA-2015-MN179	42322	42325	MN-179351408	Michael Nguyen	Tennessee	Off-PA-2847
2	CA-2015-MY18251	42251	42255	MY-182951402	Muhammed Yedou	Texas	Central US
3	CA-2015-KK15205	42308	42314	KK-152051404	Jamie Kuntz	California	Western US
4	CA-2015-RB1904	42626	42633	RB-190451408	Robert Barrios	Tennessee	Southern US
5	CA-2013-CK1220	41117	41117	CK-122051408	Chloris Kastensold	Florida	Southern US
6	CA-2015-MP1811	42353	42357	MP-181751404	Mike Pellerier	California	Western US
7	CA-2015-AK10638	42277	42283	AK-106301402	Ann Steele	Texas	Central US
8	CA-2013-AS1009	41587	41591	AS-10091402	Adam Shillingsburg	Missouri	Central US
9	CA-2013-AA1071	41542	41548	AA-107051404	Aaron McFarland	Washington	Western US
10	CA-2013-BG1169	41278	41283	BG-116951406	Brookline Gillingham	New York	Eastern US
11	US-2015-IB16000	42257	42263	IB-160001402	Joy Bell	Texas	Central US
12	CA-2013-PS1897	41608	41630	PS-189701408	Paul Stevenson	Kentucky	Southern US
13	US-2012-SB2017	41085	41071	SB-201701402	Sarah Bern	Illinois	Central US
14	CA-2015-AP1091	42251	42255	AP-109151408	Arthur Pritchep	Kentucky	Southern US
15	CA-2015-DI13551	42125	42130	DI-135501404	Don Jones	California	Western US
16	CA-2015-BH1171	42019	42023	BH-117101408	Brosina Hoffman	Mississippi	Southern US
17	US-2014-LB1679	41975	41975	LB-167951404	Laurel Beltran	California	Western US
18	CA-2015-LI1697	42321	42327	LI-169751402	Lindsey Shagati	Texas	Central US
19	CA-2012-TN2104	41191	41191	TN-210401404	Tania Norrell	California	Western US
20	CA-2015-ZC2191	42315	42320	ZC-219101404	Zoschus Carroll	California	Western US
21	CA-2012-AS1024	41192	41196	AS-102401404	Alan Shonely	California	Western US
22	CA-2014-WB2181	41654	41660	WB-218051406	William Brown	Pennsylvania	Eastern US
23	CA-2015-CY1234	42309	42314	CY-123451402	Craig Yedou	Texas	Central US
24	CA-2013-CS1250	41581	41588	CS-122051402	Chris Selenick	Wisconsin	Central US
25	CA-2014-JL15505	41738	41743	JL-155051406	Jeremy Lonsdale	New York	Eastern US
26	CA-2014-QJ1925	41928	41933	QJ-192551408	Quincy Jones	Virginia	Southern US
27	CA-2015-KP1678	42259	42260	KP-167851404	Karen Ferguson	California	Western US

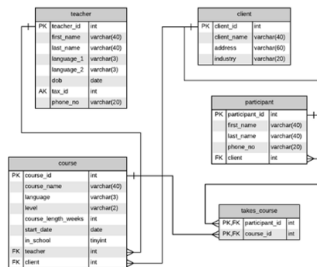
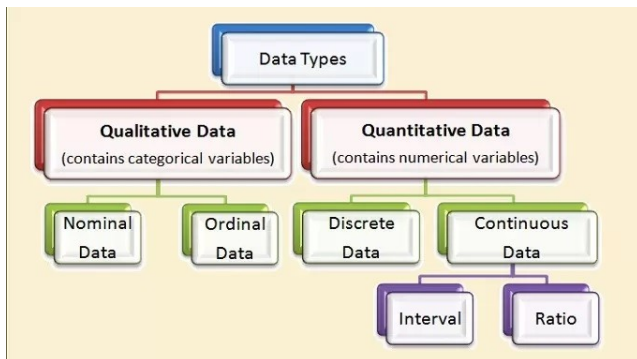


Figure: Flat file vs. relational database

# Attributes

The **value** of an attribute for a particular instance is a **measurement** of the quantity to which the attribute refers.

There is a clear distinction between attributes that are **numerical** (sometimes called **continuous**) and those that are **nominal**.



# Attributes

Statistics text often introduce four *levels of measurements*:

- **Nominal** attributes take on values in a prespecified, finite set of possibilities. The values themselves serve just as **labels** or names.
- **Ordinal** quantities are ones that make it possible to **rank order** the categories. However, although there is a notion of **ordering**, there is no notion of distance.
- **Interval** quantities have values that are not only ordered but also measured in **fixed and equal units**.
- **Ratio** quantities are ones for which the measurement method inherently defines a **zero point**.

# Knowledge representation

In the usual supervised machine learning, the learning algorithm receives as **input** the **background knowledge** and the **learning data**. It then searches its **hypothesis space** and returns the final hypothesis as a result.

It is therefore necessary to **efficiently represent** the background knowledge, the input data and the hypothesis space in order to allow for efficient usage and generation of new knowledge.

# Knowledge representation

Different kinds of knowledge are represented in different ways.

Most frequently, the following types of **knowledge representation** are used:

- Propositional calculus.
- First-order predicate calculus.
- Discriminant and regression functions.
- Probability distribution.



# Propositional calculus

In mathematical logic, **propositional calculus** is a formal deduction system, whose atomic formulae are **propositional variables**.

Propositional calculus is frequently used for representation of:

- Learning examples in classification and regression problems
- Hypotheses in symbolic learning of classification and regression rules.

# Representation with attributes

An **attribute** (feature) is a **variable** with a specified **domain** (a set of allowed values). It represents some property of an example.

# Representation with attributes

Representation with attributes is defined by

- A set of  $l$  attributes

$$A = \{A_0, \dots, A_a\}$$

- For each discrete attribute  $A_i$ : a set of allowable values

$$V_i = \{V_1, \dots, V_m\}$$

- For each continuous attribute  $A_i$ : an interval of allowable values

$$V_i = [\min_i, \max_i]$$

# Representation with attributes

- The dependent (target) variable is by definition represented with  $A_0$ :
  - In **classification** problems  $A_0$  is a discrete attribute (Class  $C$  with values  $\{C_1, \dots, C_{m_0}\}$ ).
  - In **regression** problems  $A_0$  is a continuous attribute.
- Each example is described with an attribute vector

$$t_j = (v_{1,j}, v_{2,j}, \dots, v_{a,j})$$

its value of the dependent variable is given by

$$c_j = v_{0,j}$$

- The set of learning examples is a set of attribute vectors

$$T = \{t_1, \dots, t_n\}$$

# Classification and regression rules

Classification and regression rules are usually made of two parts:

- A conditional part (**antecedent**). Usually, it is a conjunction of propositions  $T_i$ .
- An unconditional part (**consequent**). It includes the proposition  $T_0$ :

$$T_{i_1} \wedge T_{i_2} \wedge \dots \wedge T_{i_l} \Rightarrow T_0$$

**Rule representation** can be generalized by allowing other logical operators in the antecedent. They may include disjunction, negation, equivalence, and others.

# Classification and regression rules

The propositions are formed like

$$T_i = (A_i \in V_i')$$

where  $V_i'$  is

- A subset of allowed values for

$$A_i : V_i^* \subset V_i$$

if the attribute  $A_i$  is discrete.

- A sub-interval of the interval of the allowed values for the attribute

$$A_i : V_i^* \subset [\min_i^*, \max_i^*], \quad \min_i^* \geq \min_i, \quad \max_i^* \leq \max_i$$

- In regression problems, the interval in the consequential proposition is sometimes reduced to a single point

$$\min_0^* = \max_i^*$$

- In classification problems, the consequential proposition may include only one value:

$$|V_0^*| = 1$$

# Association rules

**Association rules** are a generalization of classification rules to problems where not only the target variable (class), but also the associations between attributes are interest.

The consequential part of the rule therefore does not include the proposition about the target variable  $T_0$ , but an arbitrary conjunction of propositions:

$$T_{i_1} \wedge T_{i_2} \wedge \dots \wedge T_{i_l} \Rightarrow T_{j_1} \wedge T_{j_2} \wedge \dots \wedge T_{j_k}$$

In practice,  $k > 2$  is only rarely used.



# Hypothesis space

In symbolic rule learning the hypothesis space  $\mathcal{H}$  is defined as the power set of all possible rule.

The following properties of hypotheses are important:

- **Consistency:** A hypothesis  $H$  is **consistent** if for each rule  $r \in H$  it holds that it correctly predicts the class value for all learning examples covered by the antecedent of the rule.
- **Completeness:** a hypothesis  $H$  is **complete** if every learning example is covered by at least one rule  $r \in H$ .

## Section 2

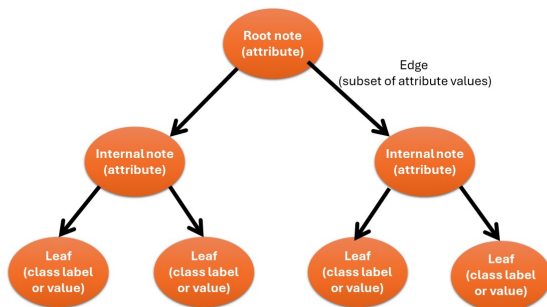
### Decision trees

---

# Decision trees

A **decision tree** is a special case of a set of decision rules.

A decision tree consists of **internal nodes** -corresponding to attributes  
edges- corresponding to subsets of attribute values, and **terminal nodes**  
(leaves) corresponding to class labels.



**Figure:** Schematic representation of a decision or regression tree.

# Decision trees

Each path starting in a tree root (topmost node) and ending in a tree leaf corresponds to a **decision rule**.

The conditions (pairs of attribute names and corresponding sets of attribute values) evaluated in each internal node are conjunctively connected.

Decision tree learning is one of the most widely used and practical methods for **inductive inference**. It is a method for approximating discrete-valued functions that is **robust to noisy data** and capable of learning disjunctive expressions.

# Regression trees

If leaf labels are continuous, such a tree is called a **regression tree**.

Knowledge representation with decision or regression trees can be generalized by evaluating in each node an **arbitrary function of several attributes** instead of a single attribute. For example, in regression trees terminal nodes (leaves) often consist of linear functions of subsets of attributes.

# Appropriate problems for decision tree learning

Decision tree learning is generally best suited to problems with the following characteristics:

- Instances are represented by **attribute-valued pairs**. real-valued attributes as well.
- The target function has **discrete output values**.
- Decision tree learning methods are **robust to errors**: in classification of the training examples and in the attribute values.
- The training data may contain **missing attribute values**.

## ID3 Algorithm

**ID3 Algorithm** is an algorithm for **empirical learning**, developed by Ross Quinlan in 1986. Its main task is constructing a decision tree.

The central choice in the **ID3 algorithm** is selecting which attribute to test at each node in the tree.

To select the attribute that is most useful for classifying examples is used the **information gain**, that measures how well a given attribute separates the training examples according to their target classification.

# Learning from examples

Let  $U$  denote the set of examples. Any subset  $C$  of  $U$  may be considered as a **concept** to be learned.

An example from  $C$  is called a **positive example** for  $C$ , an example from  $\bar{C} = U - C$  is called a **negative example** for  $C$ .

A concept  $C$  is typically represented by a unique value of a variable  $d$ , called a **decision**. Thus each value  $w$  of a decision  $d$  describes a unique concept  $C$ .

The task of learning from examples is to learn the **description**  $D$  of the concept  $C$ .



## Learning from examples

The description  $D$  learned from concept  $C$  may be expressed by a **set of rules** or a **decision tree**. The description  $D$ , characterizing the set  $C$ , is easier to comprehend for humans, and can be used by other computer programs, such as expert systems.

Two important properties of a description  $D$ :

- **Completeness:** For any example  $e$  from concept  $C$  there exists a description  $D$  such that  $D$  describes  $e$ .
- **Consistency:** For every example  $e$  from concept  $C$  there are not two different descriptions  $D$  and  $D'$  such that both  $D$  and  $D'$  describe  $e$ .

# Learning from examples

Say that the original **production rule** is

$$C_1 \wedge C_2 \wedge \dots \wedge C_j \rightarrow A$$

where  $C_1, C_2, \dots, C_j$  are **conditions** and  $A$  is an **action**.

Usually conditions are presented as ordered pairs  $(a, v)$ , where  $a$  is an attribute and  $v$  is its value, and action  $A$  is presented as an ordered pair  $(d, w)$ , where  $d$  is a decision and  $w$  is its value. Thus  $(d, w)$  describes some concept  $C$ .

## Information gain

Let  $a$  be the attribute with values  $a_1, a_2, \dots, a_l$  and let  $d$  be a decision with values  $d_1, d_2, \dots, d_k$ . Then the **information gain** is defined as

$$I(a \rightarrow d) = H(d) - H(d|a)$$

where  $H(d)$  is the **entropy** of decision  $d$ :

$$H(d) = - \sum_{i=1}^k p(d_i) \log p(d_i) \quad (1)$$

$$H(d|a) = \sum_{j=1}^l p(a_j) H(d|a_j) \quad (2)$$

$$= - \sum_{j=1}^l p(a_j) \sum_{i=1}^k p(d_i|a_j) \log p(d_i|a_j) \quad (3)$$

The **entropy**  $H(X)$  of a discrete variable random  $X$  is defined by

$$H(X) = - \sum_{x \in X} p(x) \log p(x)$$

The log is to the base 2 and entropy is expressed in **bits**. In this case, the entropy can be interpreted as the minimum numbers of bits of information needed to encode the classification of an arbitrary member of  $S$ .

Note that entropy is a **functional** of the distribution of  $X$ . It does not depend on the actual values taken by the random variable  $X$ , but only on the probability.

# Entropy

If the base of the logarithm is  $b$ , we denote the entropy as  $H_b(X)$ .

If the base of the logarithm is  $e$ , the entropy is measured in **nats**. Unless other specified, we will take all the logarithms to **base 2**.

Note that

$$H_b(X) = \log_b a H_a(X)$$

For calculations, we use the convention that  $0 \log 0 = 0$ . Adding terms of zero probability does not change the entropy.

## ID3 Algorithm

At this point, the corresponding decision tree, created by **ID3**, has the root, labeled by the attribute  $a$  and outgoing arches from the root, each such arch corresponds to a value  $a_j$  of the attribute.

The set of all examples with the same value  $a_j$  of attribute  $a$  consists of a new set  $S$  of examples.

When all members of  $S$  are members of the same concept  $C$ , they do not need to be further partitioned, so  $S$  is a label of the leaf of the tree. However, when  $S$  contains examples from at least two different concepts, the node of the tree labeled by  $S$  is the root of a new subtree.

An attribute to be a label for the root of this new subtree is selected among remaining attributes again on the basis of the **maximum of information gain** criterion.

# Example

	<i>Attributes</i>			<i>Decision</i>
	Temperature	Headache	Nausea	Flu
1	high	yes	no	yes
2	very_high	yes	yes	yes
3	normal	no	no	no
4	high	yes	yes	yes
5	high	no	yes	no
6	normal	yes	no	no
7	normal	no	yes	no

Figure: Example of the decision table.

# Example

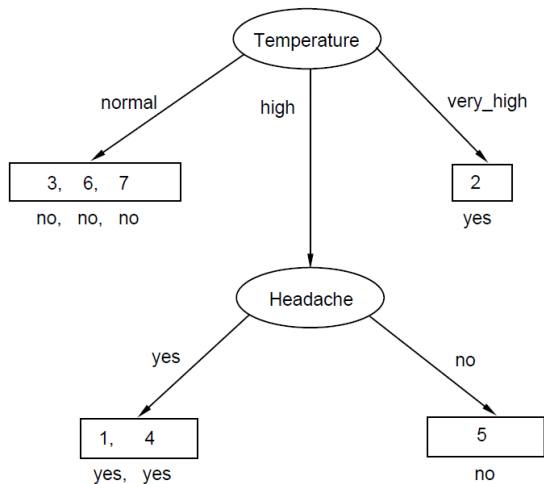


Figure: Final decision tree (Information Gain Version)



## Example

The following rules can be deduced from the above decision tree:

- $(Temperature, normal) \rightarrow (Flu, no)$
- $(Temperature, high) \wedge (Headache, yes) \rightarrow (Flu, yes)$
- $(Temperature, high) \wedge (Headache, no) \rightarrow (Flu, no)$
- $(Temperature, very\ high) \rightarrow (Flu, yes)$

## Gain Ratio Version

The **information gain version** of **ID3**, builds the decision tree with some bias—attributes with **greater** number of values are preferred by the algorithm.

In order to avoid this bias, another version of ID3 has been developed. This version uses another criterion for attribute selection, called gain ratio.

The **gain ratio** is defined as follows

$$\frac{I(a \rightarrow d)}{H(a)}$$

# Example

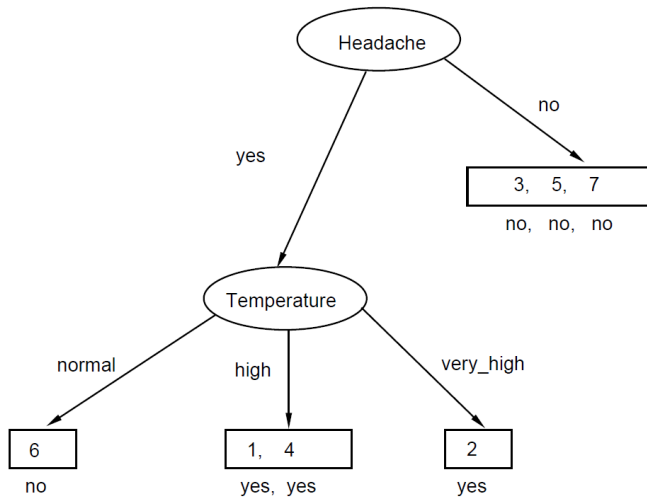


Figure: Final decision tree (Gain Ratio Version)

## Example

The following rules can be deduced from the above decision tree:

$$(Headache, yes) \wedge (Temperature, normal) \rightarrow (Flu, no)$$

$$(Headache, yes) \wedge (Temperature, high) \rightarrow (Flu, yes)$$

$$(Headache, yes) \wedge (Temperature, very\ high) \rightarrow (Flu, yes)$$

$$(Headache, no) \rightarrow (Flu, no)$$

# CART

CART stands for **Classification and Regression Tree**, is a decision tree algorithm applicable to both classification and regression.

CART employs the **Gini Index** for selecting the splitting features.

Intuitively, Gini Index represents the likelihood of two samples we randomly selected from data set belonging to different classes.

Given a feature set  $A$ , we select the feature with the lowest Gini index as the splitting feature.

## Gini Index

If a data set  $T$  contains examples from  $n$  classes, **gini index**,  $gini(T)$  is defined as

$$gini(T) = 1 - \sum_{j=1}^n p_j^2$$

where  $p_j$  is the relative frequency of class  $j$  in  $T$ .

The attribute provides the **smallest**  $gini(T)$  is chosen **to split** the node.

Note that Gini indices equal to 0 indicate that all the data belong to a single category.

# References

- Cover, T. & Thomas, J. *Elements of Information Theory*, 2nd ed., John Wiley & Sons.
- Duda, R., Hart, P. & Stork, D. (2001) *Pattern Classification*. 2nd ed., Wiley.
- Grzynala-Busse, J. (1993). Selected Algorithms of Machine Learning from Examples. *Fundamental Informaticae*, 18, 193–207.
- Kononenko, I. & Kukar, M. *Machine Learning and Data Mining. Introduction to principles and algorithms*, Horwood Publishing, 2007.
- Mitchell, T., *Machine Learning*, McGraw-Hill, 1997.
- Witten, I. & Frank, E. *Data Mining. Practical Machine Learning Tools and Techniques*, 2nd ed., Morgan Kaufmann Publishers, 2005.