

Introduction- to t'ecnicas classi fi cac'i'on



Maria-Amparo Vila
vila@decsai.ugr.es

Investigaci'on group Databases and Intelligent
Systems Informaci'on <https://idbis.ugr.es/>

Department of Sciences of the computation and
artificial intelligence
University of Granada

Introduction

Basic concept

Definition

Classification is the process of learning a function applying a set of attributes X_1, \dots, X_n in another attribute Y . Yes:

- Yes Y It is discrete, boolean, nominal etc. we have **Models classified** themselves
- Yes Y is continuous we **Regression models**

The function is called **learns** también **Model classification** in general

Introduction

B'asicos concept

Seg'one we have the goal of learning:

explanatory models Tambi'en called descriptive. try

mostar c'omo depends Y of $X_{one..N}$: decisi'on spanning trees of classi
fi ers Bayesian models regresi'on, model rules

predictive models . They do not seek both show the dependence

as given an item or_i with values $x_{ij}, j = one..N$ get the value Y_i of the target
variable. Yes Y It is discreet, the class to which it belongs. M'etodos the
nearby M'as neighbor, m'etodos based on neural networks. SVM

Introduction

The general process of classifying cases on

one.- September 1 data values are considered in Y . set of training

items variables X_{one} X_{two} X_N				Y
X_{one}	X_{eleven}	X_{12}	$X_{one\ N}$	Y_{one}
\vdots	\vdots	\vdots	\vdots	\vdots
X_M	$X_{M\ one}$	$X_{M\ two}$	X_{MN}	Y_M

two.- Is constructed (learn) the classified models cases on

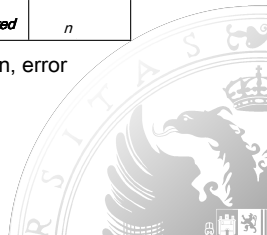
Introduction

The general process of classifying cases on

- 3.- Tested in another dataset **joint test** calculating the values Y_{pred}

items variables X_{one}					
		X_N	Y	Y_{pred}
X_{one}	X_{eleven}	X_{one}	N	Y_{one}
\vdots	\vdots	\vdots	\vdots	\vdots
X_n	$X_{n\ one}$	X_n	N	Y_n
				Y_{pred}	n

- 4.- is evaluated in the model segment one different criteria: precision, error rate, classification, scalability, interpretability, complexity, etc.



Introduction

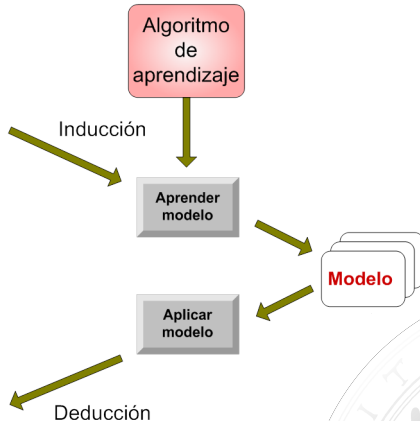
The general process of classifying cars on

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Conjunto de
entrenamiento

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

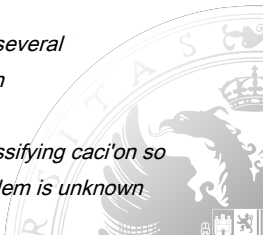
Conjunto de prueba



Caci'on classified by trees Decisi'on

b'asicas ideas

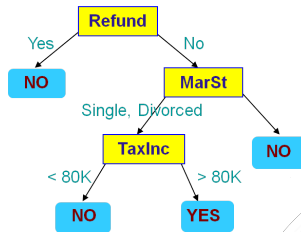
- *The spanning trees of decisi'on try to find a structure for explaining jer'arquica c'omo different parts of the input space corresponding to different values of the attribute object.*
- *The tree has three types of nodes:*
 - *Root node where begins*
 - *internal nodes each of which has an input shaft and two or more output partitions the subspace corresponding to this node*
 - *Leaf node has output shafts and est'a labeled with a target attribute valos*
- *At each node that is not part of sheet input space is divided into several subsets seg' one value of a given attribute, to reach leaf nodes.*
- *In principle the importance of each attribute in the process of classifying caci'on so it is possible that different results are obtained for the same problem is unknown*



Categorization by decision trees

Example: tree from data

Tid	categorico		categorico	continuo	clase
	Refund	Marital Status	Taxable Income	Cheat	
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	



Conjunto de
entrenamiento

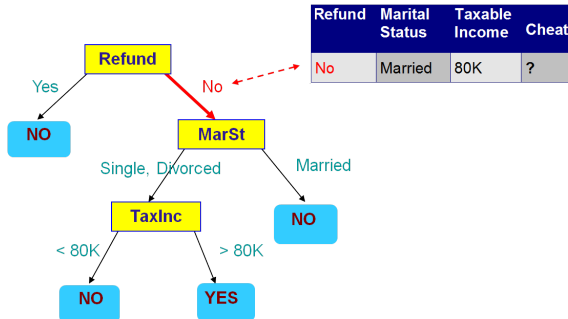


Modelo de clasificación:
Árbol de decisión

15

Cacit'on classified by trees Decisi'on

Example: classifying cacit' on an item



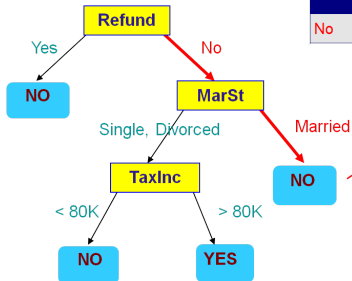
Modelo de clasificaci'':
rbol de decisi''

Cacit'on classified by trees Decisi'on

Example: classifying cacit' on an item

Caso de prueba

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	No



Clase 'No'

Modelo de clasificaci3n:
3rbol de decisi3n

Classification trees Decision

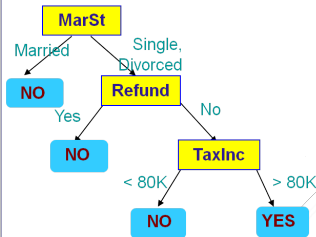
Example: other tree

Tid	Refund	Marital Status	Taxable Income	Cheat
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

Conjunto de
entrenamiento



Modelo de clasificación:
Árbol de decisión



Constructing Decision Trees

basic ideas

- Finding a spanning tree of decision is not trivial. Is two M where M is the number of examples
- A reasonably spanning tree is sought not explain properly training set
- It uses a **Greedy strategy** converting the problem NP problem.
- It is part of a root node and goes ramified each node of the "best" possible way



Constructing Decision Trees

basic ideas

? Algorithm "divide and conquer"

one. We started with all training examples in root
spanning tree of decision.

two. Examples are dividing into function of the attribute
Select to refine the spanning tree on each node.

3. If a node contains examples of solo a class is transformed into leaf

Four. The attributes used to refinement are chosen from function
one heuristic.

5. The form of refined also



Constructing Decision Trees on

Application problems

on the heuristic

When building a decision tree stops?

- When all remaining examples belong to the same class (is a leaf spanning tree with the class label).
- When there are attributes for which a majority class (is a leaf labeled with the most frequent class in the node).
- When we are not classifying data.



Construcci' Trees on Decisi' on

Aplicaci' problems

on the heuristic

¿ Given a non-leaf node, as is partitioned

binary nodes It has no problem

nominal nodes Two options:

- Partitioning all values (partici'on nary)
- Group and converted into binary

ordinal nodes Two options:

- Repartition all values
- Group and converted into binary by fixing a cut point ($\leq v, > v$)



Constructing Decision Trees on

Application problems

on the heuristic

Criteria for Node Partition

Given a non-leaf node, as is partitioned

Numerical nodes We have two options:

- Discretize the attribute and treated as ordinal
- Group and converted into binary by fixing a cut point ($\leq v, > v$)

Utilize different algorithms partition different forms: single binary CART

ID3 only discrete attributes and partition nary attributes for categorical

C4.5 partition nary and binary to continuous, etc.



Constructing Decision Trees on

Application problems

on the heuristic

Given a non-leaf node, ¿Which attribute is chosen to partition?

- We want a tree as shallow as possible. The goal is to reach as leaf nodes as early as possible.
- Partitions need to separate elements of different classes
- Partitions are looking for homogeneous nodes
- We will use measures based on the diversity of classes in each element of the partition. *measures impurity*



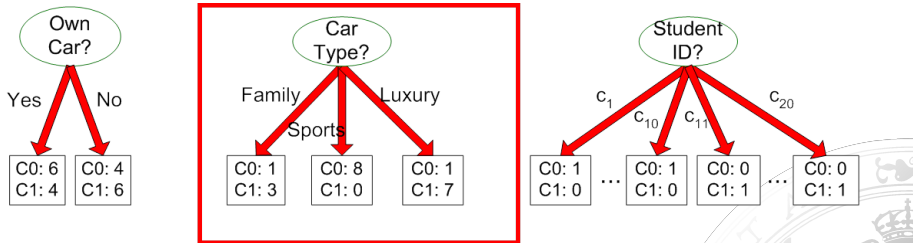
Construisci Decision Trees on Decisi

on

Seleziona i criteri

on attributes

Example



Construisci Trees on Decisi'

on

Selecci' criteria

on attributes

?Selecci'on measures

- Be $p(i/t)$ $i \in \{\text{one, two... } c\}$ the fracci'on of items belonging to the class i that est'a in a given node t , obviously c It is the n' umber of classes
- $p(i/t)$ It is an approximation of the probability of finding an item class i in that partici'on t It represents.
- According to the above idea, the more uniform are the values of $p(i/t)$ It is less desirable t to be selected.
- The worst possible value for $p(i/t)$ $i \in \{\text{one, two... } c\}$ is $(\text{one}/nt, \dots, \text{one}/nt)$ where nt It is the n' umber of element t . The best is $(0 \dots, \text{one}, \dots 0)$



Construisci Trees on Decisi

on

Seleziona i criteri

on attributes

? Seleziona misure basate su Entropy

entropy $Entropy(t) = - \sum_c$

$$p(i/t) \log_{two}(p(i/t))$$

Gini index $Gini = 1 - \sum_c$

$$p(i/t)^2$$

Error classificazione $error = 1 - \max_i (p(i/t))$

Example

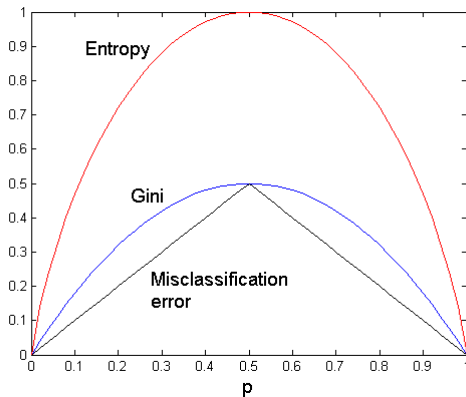
Distribuzione node			Gini	entropy	Error
N_{one}	Class1 0		0	0	0
	Class2 1				
N_{two}	Class1 Class2		0.278	0.650	0.167
	0.2 0.8				
N_{one}	Class1 Class2		0.5	one	0.5
	0.5 0.5				



Construisci Trees on Decision

Seleziona i criteri per la scelta degli attributi

Seleziona le misure basate sull'Entropia, per due classi



Construisci Trees on Decisi

on

Selecc'i criteria

on attributes

? Informaci'on gain

To see c'omo operates a divisi'on compare the measure of a parent node to the child nodes. Sean p a parent node

$v_j, j = \text{one} \dots k$ their children, $N(p)$ n' umber of elements in the node p Y

$N(v_j)$ n' umber of elements v_j . We define:

$$\text{Gain } 4 = I(p) - \sum_{j=\text{one}} \frac{N(v_j)}{N(p)} I(v_j) \text{ where } I(.) \text{ It is one of the measures before fi ned. Sis } I(.) \text{ is the Entropy is called } \text{Informaci'on gain } 4 \text{ info}$$

Proporci'on Gain Gain $\text{Gainratio} = 4 \text{ info}$

$$\text{SplitInfo} = - \sum_{j=\text{one}} \frac{N(v_j)}{N(p)} \log_{\text{two}} \left(\frac{N(v_j)}{N(p)} \right)$$

SplitInfo where

4 info used ID3 and GainRatio in C4.5. CART, SLIQ..utilizan the Gini'indice

Constructing Decision Trees

Selection criteria *on attributes*

Comparison rules division

Information gain Biased towards attributes with many different values.

Proportion gain Partitions tend to prefer slightly balanced (with a partition M as much larger than the other)

Gini index It works worse when there are many kinds and tends to partition favoring one and no similar purity.

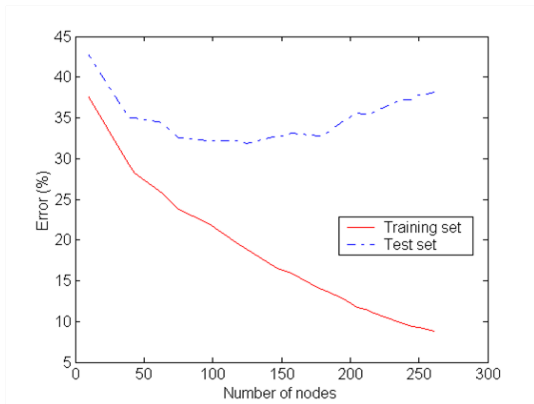
No rule of division is significantly better than the others



Constructing Decision Trees on Decision

Additional questions: overlearning

The greater complexity, models classified based on features fit the training set **overlearning**



Constructing Decision Trees on

Additional questions: overlearning

A solution to overlearning are *Técnicas pruning* which they are developed to simplify the spanning tree. To prune a spanning tree of decision

- one subtree by a leaf node (corresponding to frequent M's class in the subtree) is replaced
- Or, a subtree other subtree contained in the first. There are techniques

previous Poda It is shrinking the spanning tree when it is generating

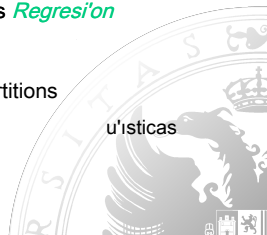
Poda post the spanning tree once generated is reduced For criteria see bibliography a basic



Construcci' Trees on Decisi' on

additional issues

- The classified caci'on / predicci'on by spanning trees of decisi'on is one of the M'as t'ecnicas studied within the classification caci'on.
- There are many variants and extensions of algorithms b'asicos funci'on of: partici'on mechanisms of space, use of t'ecnicas pruning, use of additional mechanisms as the rules of asociaci'on etc.
- The idea has spread to predict m'etodos for continuous attributes *Regresi'on trees*
- They have been extended to consider criteria partici'on fuzzy partitions domain labels establishing LING for attributes discretized *Fuzzy Decison Trees*



Construcci' Trees on Decisi' on

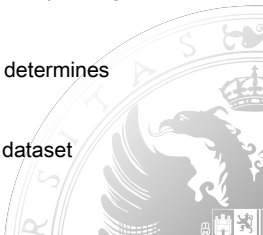
additional issues

Advantages of spanning trees of decisi'on

- Facil interpretaci'on (when peque~ us).
- Quickly to classify new data.
- Precisi'on comparable to other t'ecnicas.

Some algorithms ef fi cient and scalable

- PUBLIC (Rastogi and Shim, VLDB'1998) integrates pruning process spanning tree construcci'on
- RainForest (Gehrke et al., VLDB'1998) separates what algorithm determines scalability
- BOAT (Gehrke et al., PODS'1999) s'olo needs to run 2 times the dataset



Categori'on classified by rules

b'asicas ideas

objective

Classify records using a collection rule "if then"

The form of a rule is Condition - \rightarrow Y Where:

- condition is a conjunction conditions on the value of various attributes, tambi'en antecedent
- Y is the value of third chapter is devoted class consequent

Examples of rules

- Blood Type = hot \wedge Lays eggs = Si - \rightarrow Bird
- Income 6 30 \wedge Devoluci'on = yes - \rightarrow no evader



Caci'on classified by rules

b'asicas ideas

*Given a rule r we say that **it covers** an instance x the dataset if that body satisfies the history of the rule*

Example

R1: (Viviparo = no) \wedge (Puede volar = yes) \rightarrow Pajaro

R2: (Viviparo = no) \wedge (Acuatico = yes) \rightarrow Pez

R3: ((Viviparo = yes) \wedge (Sangre = caliente) \rightarrow Mamífero

R4: ((Viviparo = no) \wedge (Puede volar = no) \rightarrow Reptil

R5: (Acuatico = a veces) \rightarrow Anfibios

First name	Blood	Viviparous	Can fly	Water	Class
Hawk	Hot	do not	yes	do not	bird
Bear	Hot	yes	do not	do not	mammal
Platypus	Hot	do not	do not	Sometimes	mammal
Lemur	Hot	yes	do not	do not	mammal
Turtle	cold	do not	do not	Sometimes	an bio fi

Cac'ion classified by rules

b'asicas ideas

Example

- Halcon is covered by R1
- Platypus and turtle are covered by R5 A rule is:

Coverage Proportion of records that satisfy their background

Precision Proportion of records and meeting records
consistent

Example

Ruler	Cober.	Preci.
R1	1/5	1/5
R2	0	0
R3	2/5	2/5
R4	2/5	0
R5	2/5	1/5



Caci'on classified by rules

b'asicas ideas

A test set is classi fi ed register to register firing rules corresponding to the values of each of them, and scoring his class

Example

Hot Gorri'on	no	yes	no	bird			
--------------	----	-----	----	------	--	--	--

Register shoot Rule R1



Cacì'on classified by rules

b'asicas ideas

Regarding your aplicacì'on a set of rules can be:

Mutually exclusive Yes:

- Each rule can be applied independently
- Any registration est'a covered as much by a rule

Comprehensive Yes:

- There is a rule for any possible attribute values combinacì'on
- Any registration est'a covered by at least one rule



Caci'on classified by rules

b'asicas ideas

Example. set of mutually exclusive and exhaustive rules

r1: (Sangre= fria) \rightarrow No mamífero

r2: (Sangre=caliente) \wedge (Vivíparo = yes) \rightarrow Mamífero

R3: ((Sangre = caliente) \wedge (Viviparo = No) \rightarrow No Mamifero

If a set of rules is mutually exclusive and exhaustive, each record to classify triggers a rule ys'olo one



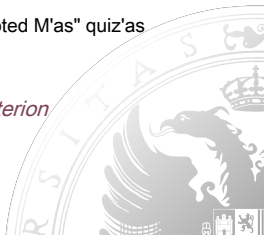
Caci'on classified by rules

b'asicas ideas

¿Qu'e do when a set of rules does not have these properties?

- If not comprehensive malaria, it defines a default class and assigned to it the uncovered records
- If not exclusive, a record may be covered by several contradictory rules.
solutions:
 - the rules are sorted by alg' criteria: coverage, precis'i'on class that define etc., and priority rule applies M'as
 - all the rules for the registration skyrocket and assigned the class "voted M'as" quiz'as weighted by the weight of the rules

The mayor'ia algorithms do not produce unique rules follow a criterion ordenaci'on



Cac'ion classified by rules

Extracci' on rules: general ideas

Given a set of TRAINING WHAT C'omo extract a set of rules?

From a spanning tree of decisi'on Simply describe the spanning tree by a Set of rules.

- They are mutually exclusive
- They are exhaustive
- Informaci'on contain all the spanning tree

direct M'etodos Act' ohn directly on the data, known they are those of M'as m'etodos *sequential coating* .

CN2, RIPPER and its variants etc.

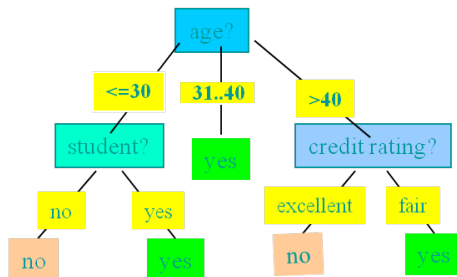


Caci'on classified by rules

Extracci' Rules on using spanning trees of decisi'

on

Example



IF (age<=30) AND (student=no) THEN buys_computer = no

IF (age<=30) AND (student=yes) THEN
buys_computer = yes

IF (30<age<=40) THEN buys_computer = yes

IF (age>40) AND (credit_rating=excellent) THEN
buys_computer = no

IF (age>40) AND (credit_rating=fair) THEN
buys_computer = yes

Cac'ion classified by rules

Extracci' direct on rules: b'asicas Ideas

The process is the b'asico **sequential coating**

one. Start with a set of rules vac'io

two. generate **best rule** covering a particular class

3. TO nadir learned rule to set

Four. Remove the examples set of training covered by
Rule

5. Failure to comply with the **stopping rule** go 2 otherwise stop



Cac'ion classified by rules

Extracci' direct on rules: b'asicas Ideas

To extract the best rule:

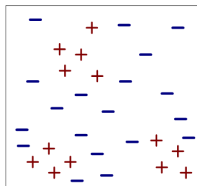
- one.** classes are ordered (no criteria high to low or contrary)
- two.** They are considered positive examples of that class and the negative rest
- 3.** The best rule is covering many positive examples and few negative. evaluaci'on measures used for this purpose.
- Four.** A rule may need to be pruned if it causes problems of overlearning.



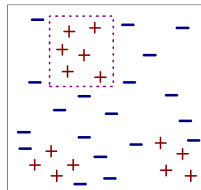
Caci'on classified by rules

Extracci' direct on rules: b'asicas Ideas

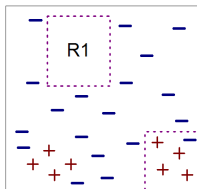
Example rules selecci'on



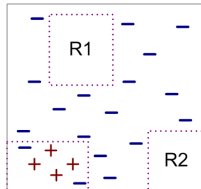
(i) Original Data



(ii) Step 1



(iii) Step 2



(iv) Step 3

Cac'ion classified by rules

Extracci' direct on rules: b'asicas Ideas

some algorithms

- FOIL (Quinlan, Machine Learning, 1990)
- CN2 (Clark and Boswell, EWSL'1991)
- RIPPER (Cohen, ICML'1995)
- PNrul (Joshi and Kumar Agarwal, SIGMOD'2001)



Bayesian classification

basics ideas

- There are problems in which the relation between items and classes have a random component
- Even items (instances) with equal values in the attributes may belong to different classes (diagnostic of seg' diseases
one symptom)
- **basic principle** If you can not ensure qu'e class belongs to an instance, **assign the class that is more likely to belong**



Bayesian classification

basic ideas

- Two problems:

1 Given a particular instance, it is estimated the class to which it belongs? *By Bayes theorem*

2 How are stored / calculate the classes

likely set one possible combination of values
attributes so efficient? *Hypothesis simplified classes*

Hypothesis of independence for discrete attributes

Method leads to "Naive Bayes"

Hypothesis joint standard for continuous attributes

distribution Leads to Analysis

Discriminant and its variants



Bayesian classification

probabilistic Model: Bayes theorem

- Given random events A and C we have:

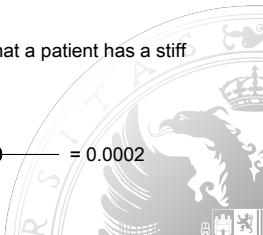
$$P(C|A) = \frac{P(A, C)}{P(C)} \quad P(A|C) = \frac{P(A, C)}{P(C)} \quad P(A) \Rightarrow$$

$$\Rightarrow P(C|A) = \frac{P(A|C) P(C)}{P(A)}$$

- Example of Use*

- It is known that meningitis causes stiff neck in 50% of cases.
- It is known that the probability of having meningitis is 1/50000 and that a patient has a stiff neck 1/20.
- So:

$$P(\text{Men} | \text{Rg}) = \frac{P(\text{Rg} | \text{Men}) P(\text{Men})}{P(\text{Rg})} = \frac{0.5 \times \frac{1}{50000}}{\frac{1}{20}} = 0.0002$$



Bayesian classi fi cation

probabilistico model

- Suppose that both attributes $X_{one}.. X_N$ as the class Y They are random variables.
- Given an instance (item) with attribute values $x_{one}.. x_N$ we want to predict the value of its class Y
- Specifically want to find the value that maximizes expression:

$$Prob (Y = y / X_1 = x_{one}, X_2 = x_{two}, ..., X_n = x_N)$$

to simplify $P (y / x_{one}, x_{two}.., x_N)$

- **issue** :

Can we calculate $P (y / x_{one}, x_{two}.., x_N)$, known ($x_{one}, .. x_N$) from the data?

- **Solution** The use of Bayes' theorem.



Bayesian classificazione

probabilistico model algorithm basico

one. For all $Y \in \text{do my}$ calculate:

$$P(y | x_{\text{one}}, x_{\text{two}}, \dots, x_n) = \frac{P(x_{\text{one}}, x_{\text{two}}, \dots, x_n | y) P(y)}{P(x_{\text{one}}, \dots, x_n)}$$

two. To choose \hat{Y} such that

$$P(\hat{Y} | x_{\text{one}}, x_{\text{two}}, \dots, x_n) = \max_{Y \in \text{do my}} \frac{P(x_{\text{one}}, x_{\text{two}}, \dots, x_n | Y) P(Y)}{P(x_{\text{one}}, \dots, x_n)}$$

3. It really is equivalent to choose \hat{Y} such that

$$P(\hat{Y} | x_{\text{one}}, x_{\text{two}}, \dots, x_n) = \max_{Y \in \text{do my}} P(x_{\text{one}}, x_{\text{two}}, \dots, x_n | Y) P(Y)$$

Four. **issue**. Come estimate $P(x_{\text{one}}, x_{\text{two}}, \dots, x_n | Y)$? In principle it is joint distributions of N random variables, each conditional value of the domain class

Bayesian classification

Naive Bayes classifier

conditional independence between attributes is assumed

X_1, \dots, X_N so that:

$$\forall Y \in \text{Dom}(Y) P(X_1, X_2, \dots, X_N | Y) = P(X_1 | Y) P(X_2 | Y) \dots P(X_N | Y)$$

one. $\forall Y \in \text{Dom}(Y) \forall j \in \{1, \dots, N\}$ It can be estimated $P(X_j | Y)$ using the training set

two. Since a new item value (x_1, \dots, x_N) It is categorized in class z such that:

$$P(z) P(x_1 | z) \dots P(x_N | z) = \max_{Y \in \text{Dom}(Y)} P(Y) P(x_1 | Y) \dots P(x_N | Y)$$



Bayesian classification

Naive Bayes classification: simple example

Tid	Refund	Marital Status	Taxable Income	Evade
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

Chance of classes

$$P(y) = M/Y, P(S) = 10/30, P(N) = 7/10$$

Probability of discrete attributes

$$P(x_j | y) = m_{x_j y} / m_y$$

$$P(\text{Status} = \text{Married} | \text{no}) = 4/7$$

Probability of continuous attributes

Hypothesis normal

$$P(x_j | y) = \frac{1}{\sigma_j \sqrt{2\pi}} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

$$P(\text{Income} = 120 | \text{No}) = 0.0072$$

Bayesian classi fi caci'on

Naive Bayes classi fi er: simple example

Be $X = (\text{Refund} = \text{NO}, \text{Married}, \text{Income} = 120)$

naive Bayes Classifier:

$P(\text{Refund}=\text{Yes}|\text{No}) = 3/7$
 $P(\text{Refund}=\text{No}|\text{No}) = 4/7$
 $P(\text{Refund}=\text{Yes}|\text{Yes}) = 0$
 $P(\text{Refund}=\text{No}|\text{Yes}) = 1$
 $P(\text{Marital Status}=\text{Single}|\text{No}) = 2/7$
 $P(\text{Marital Status}=\text{Divorced}|\text{No})=1/7$
 $P(\text{Marital Status}=\text{Married}|\text{No}) = 4/7$
 $P(\text{Marital Status}=\text{Single}|\text{Yes}) = 2/7$
 $P(\text{Marital Status}=\text{Divorced}|\text{Yes})=1/7$
 $P(\text{Marital Status}=\text{Married}|\text{Yes}) = 0$

For taxable income:

If class=No: sample mean=110
sample variance=2975

If class=Yes: sample mean=90
sample variance=25

$$\begin{aligned}P(X|\text{Class}=\text{No}) &= P(\text{Refund}=\text{No}|\text{Class}=\text{No}) \\&\times P(\text{Married}|\text{Class}=\text{No}) \\&\times P(\text{Income}=120K|\text{Class}=\text{No}) \\&= 4/7 \times 4/7 \times 0.0072 = 0.0024\end{aligned}$$

$$\begin{aligned}P(X|\text{Class}=\text{Yes}) &= P(\text{Refund}=\text{No}|\text{Class}=\text{Yes}) \\&\times P(\text{Married}|\text{Class}=\text{Yes}) \\&\times P(\text{Income}=120K|\text{Class}=\text{Yes}) \\&= 1 \times 0 \times 1.2 \times 10^{-9} = 0\end{aligned}$$

Since $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$

Therefore $P(\text{No}|X) > P(\text{Yes}|X)$
 $\Rightarrow \text{Class} = \text{No}$

Bayesian classification

Naive Bayes classifier: estimation

on the probabilities of discrete attributes

In general $\forall Y \in \text{Dom}(Y) \forall j \in \{1, \dots, N\} P(x_j | Y)$ it is estimated:

$$P(x_j | y) = \gamma + \frac{m_{x_j y}}{\gamma m_{x_j} + m_y}$$

where m_{x_j} is the number of elements $\text{dom}(X_j) \cap \text{dom}(Y)$ it is estimated as:

$$p(y) = \gamma + \frac{m_y}{\gamma m_Y + m}$$

where m_Y is the number of elements having $\text{dom}(Y)$ m the number of elements of the dataset

- constant γ It is called **CORRECTION Laplace**.
- Usually it is taken equal to zero; but to treat cases of non-existent attribute values is taken equal to 1 or 1/2
- It is used when the training set is peque~

do not

Bayesian classification

Naive Bayes classification: estimation

on the probabilities of continuous attributes

- A continuous attribute X_j is considered distributed $N(\mu_{jY}, \sigma_{jY}^2)$ for all class value Y . And its conditional probability is given by $f(x_{jY}) = N(x_{jY} | \mu_{jY}, \sigma_{jY}^2)$

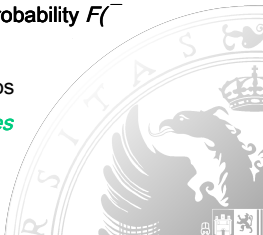
$f(x_{jY})$

- When it has a set of continuous attributes \bar{X} , it can be avoided assuming conditional independence,

$\forall Y, (\bar{X} | Y)$ It is distributed as a Multivariate Normal

$N(\mu_{\bar{X}|Y}, \Sigma_{\bar{X}|Y})$. You can then calculate the joint conditional probability $F(\bar{x} | y)$.

- When we have numerical attributes do not impose the method of *hypothesis independence* and we have *complete Bayesian classification*



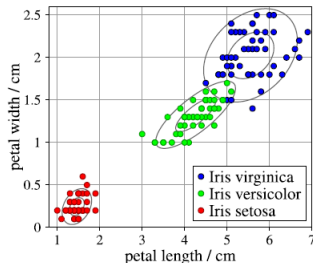
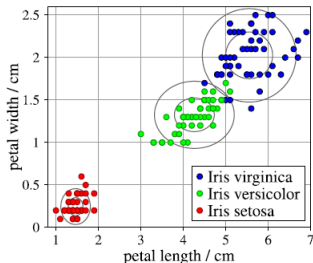
Bayesian classi fi caci'on

Naive Bayes classi fi er: estimaci'

on the probabilities of continuous attributes

Example

Iris type	Iris setosa	Iris versicolor	Iris virginica
Prior probability	0.333	0.333	0.333
Petal length	1.46 ± 0.17	4.26 ± 0.46	5.55 ± 0.55
Petal width	0.24 ± 0.11	1.33 ± 0.20	2.03 ± 0.27



Bayesian classification

Naive Bayes classifier overview

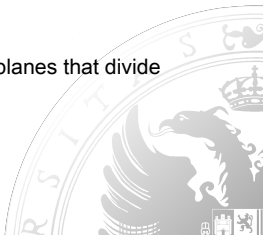
- They are robust against noise isolated points and working on midrange
- Ignoring orders let you manage instances values that are at the stage of ESTIMATION odds.
- They are robust to irrelevant attributes because if an attribute X It has no influence on $YP(X/Y)$ It tends to distribution uniform.
- The conditional independence hypothesis everyone attributes can be very strong.
 - Bayesian networks generalize the model and make this hypothesis flexible M'as



Bayesian classification

Discriminant analysis

- It is a particular case of *classical full Bayes*
- Hypothesis Simplifications:
 - Numerical attributes
 - Distribution regular multivariate. With restrictions:
 - Covariance matrix equal for classes
 - very different half
 - Initially two classes only
- The calculation is based on class optima calculate a set of hyperplanes that divide the space for classes.



Classify instances based on

- The classifiers studied so far, working in two stages:
 - inductively** Learning model classification
 - inferential** Apply the model to the test examples set
- They are "forward classifiers" (eager learners)

Basic Idea

Why not store the entire training set and when you get a test example search, items that "Most will appear" and assign class ?.

They are methods "lazy" Lazy Learners



Classify instances based on

Set of Stored Cases

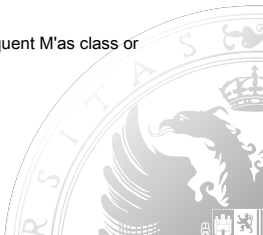
Atr1	AtrN	Class
			A
			B
			B
			C
			A
			C
			B

Unseen Case

Atr1	AtrN

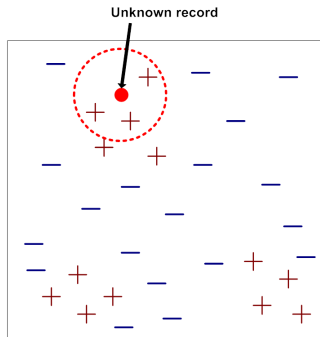
Classify instances based on

- Methods obviously are predictive. They explain nothing.
- Examples:
 - The "memoristic classifiers" (Rote Learners)
 - Stores the training set and only assigns a class to an example when there is a training item that is exactly like it.
 - K-nearest neighbors (K nearest neighbor) (K-NN)
 - Select the k-items that "most resemble" the example and assign the frequent class or "important class"



K-nearest neighbors M'as (K-NN)

b'asicas ideas



requirements

Distance between records $d(.,.)$

The value k fixed

Algorithm

1. Sea x calculate $d(x, e), \forall e \in \text{AND}$

two.- Identify $\{ \text{and } i, i = \text{one} \dots k \}$

k nearest neighbors M'as x

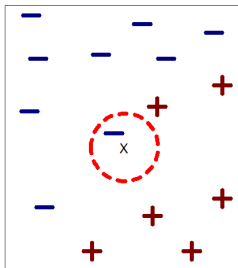
3.- classes $\{ \text{and } i, i = \text{one} \dots k \}$,

get the kind of x

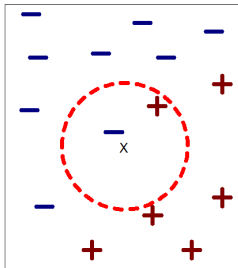
Mayor'ia or weighted by mayor'ia

K-nearest neighbors M'as (K-NN)

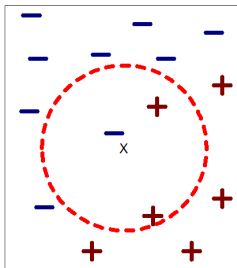
b'asicas ideas



(a) 1-nearest neighbor



(b) 2-nearest neighbor



(c) 3-nearest neighbor

It is important to fix the value of k . Since it can lead to overlearning or error.

K-nearest neighbors M'as (K-NN)

Aplicaci' problems on

The function away .

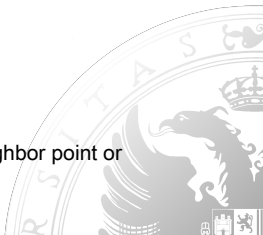
- Est'a's method designed to num'ericos attributes
- You can use distances proposals on the subject grouping
- Should take into account problems of scale

The value of k .

- It is best to get it through validaci'on cross
(*is subsequently ver'a*)

Class mechanism selecci'on .

- Usually choose the majority class
- It can be weighted by the distance from the neighbor point or more functions so this fi sticadas



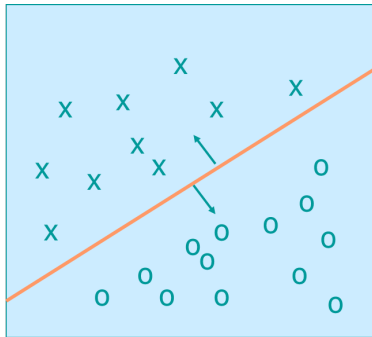
Other predictive models classification

on

Classifiers based on Neural Networks

They have been studied in other subjects

Support Vector Machines SVMs



Other predictive models classifi caci'

on

Support Vector Machines SVMs

Advantage . • **Precisi'**

on high

- Robustness against noise

drawbacks • **Expensive to train. (Little scalable and ef fi cient)**

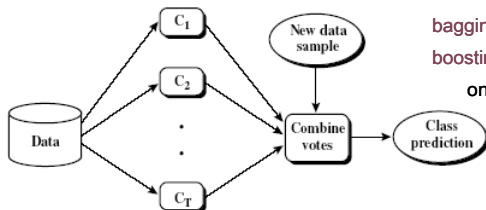
- Di fi cult to interpret (Increase dimension space to separate classes hierplanos)



Other models classified caci'on

"Emsembles"

Combine several models to improve the classi fi er
preci'sion deciding qu'e class is assigned:



bagging Votaci'on by majority

boosting Votaci'on weighted seg'

one classi fi er quality (*AdaBoost*)

Evaluaci'on models classified cac'i'on

b'asicas ideas

Process steps classi fi cac'i'on:

- one.-** September 1 data values are considered in Y . **set of training**
- two.-** Is constructed (learn) the classified models cac'i'on
- 3.-** Tested in another dataset **joint test** caculando the values Y_{pred}
- 4.- is eval' ua in the model seg'** one different criteria: preci'si'on, Error classi fi cac'i'on, scalability, interpretability, complexity, etc.



Evaluation models classified according

basic ideas

- *Evaluation different aspects of modeling*

Metrics They are measures of the "quality" of a process classification.

Models So they used to estimate reliable measures quality.

Comparison Are techniques that compare performance concerning two models classification



Evaluation of models classified according to

basic ideas

- *Criteria for quality measures*

Precision ¿Cómo well classifies the model?

Efficiency needed to build / use the classifier time

Sturdiness Against noise and nulls

scalability You admit large datasets?

interpretability Does it explain the model?

Complexity Tree with many nodes etc.

The first two may be metrics while 's'olo can be measured in some cases.

last

Precision measures ser'an key



Evaluation models classification

Precision measures on

Sean: **Model classifier**, and $x_i = x_{i1}, \dots, x_{in}$ an example of the independent variable of a test set $T = \{x_1, y_1, \dots, x_n, y_n\}$.

$\hat{y}_i = M(x_i)$ the result of applying the process to $M(x_i)$.

We define:

Precision M.

$$Acc = \frac{1}{n} \sum_{i=1}^n I(\hat{y}_i = y_i)$$

where $I(e)$ It equals 1 if e is true, and 0 if e is false.

Error rate M.

$$Er = \frac{1}{n} \sum_{i=1}^n I(\hat{y}_i \neq y_i) = 1 - Acc$$

Evaluation models classification

Precision measures on

It seems that it is best to have a low error rate / high precision but:

If it is excessively adjusted to set s'olo learn this training model. overlearning

The appearance of overlearning is due to various causes, some dependent models; but others can be smoothed.



Evaluation of models classified on

Precision measures on

Some considerations on overlearning

- The greater the complexity of a model classification, Models excessively adjusted to the training set (in spanning trees of decision)
- Another cause of overlearning is the presence of noise points (on models that partition the space: analysis discriminant, SVM
eg)
- By Finally a shortage of points in a class against other overlearning can give problems, since the precision not take into account the number of classes

To avoid this' Last issue must take into account the weight classes precision measures for this purpose

Arrays Confusion



Evaluation of models classified categorical

Precision measures on: confusion matrices on

Given a classification process if we define n_{ij} number of elements assigned to the class i when estimated in class j we have:

Confusion matrix

predicted some	Y_{one}	Y_{two}	...	Y_k	TOT_{pred}
\hat{Y}_{one}	$n_{one one}$	$n_{two one}$...	$n_{k one}$	m_{one}
\vdots	\vdots	\vdots	...	\vdots	\vdots
\hat{Y}_k	$n_{k one}$	$n_{k two}$...	$n_{k k}$	m_k
TOT_{CIER}	n_{one}	n_{two}	...	n_k	n

Confusion matrices allow define measures associated with the classes and global measures weighted

Evaluation of models classified categorical data

Precision measures on: confusion matrices on

They can define:

Precision of a class $\forall i \in \{1 \dots k\}, \text{prec}_i = n_{ii} / m_i$

"Recall" of a class $\forall i \in \{1 \dots k\}, \text{rec}_i = n_{ii} / n_i$

F as a class $\forall i \in \{1 \dots k\}, F_i = 2 n_{ii} / (m_i + n_i)$

Precision and recall the overall remain the same but you can define:

$$\text{Global F} = 1 / \sum_{i=1}^k \frac{1}{F_i}$$



Evaluation of models classified caci' on

Precisi' measures on: matrices confusi' on

Example

Iris data using "sepal length" and "sepal width". Naive Bayes classi fi er. 120 training examples and test data 30.

global data $acc = 0.733$ $er = 0.267$

Confusion Matrix

predicted certain	setosa	versicolor	Virginica	
setosa	10	0	0	10
versicolor	0	7	5	12
virginica	0	3	5	8
	10	10	10	30

Data associated with classes

	Accuracy	Recall	F-measure
setosa	one	one	one
versicolor	0.583	0.7	0.636
virginica	0.625	0.5	0.556

Overall = 0.731 F



Evaluation of models classified categorical on

binary problems: matrices confusion

on

Binary matrices confusion

- Classes are now P (positive), N (negative)
- The confusion matrix is

predicted some	positives	negatives	
positives	TP	FP	$P_{pred} = TP + FP$
negatives	FN	TN	$N_{pred} = FN + TN$
	$P_d = TP + FN$	$N_d = FP + TN$	n

- the measurements are calculated using the previous expressions.
- In the event that there is a great discrepancy between cases you can use the cost matrix

Evaluation of models classified on

binary problems: ROC curves (Receiver Operating Characteristics)

Hypothesis

- A binary problem
- There is a measure $S(.)$ and a threshold ρ such that if $S(x_i) > \rho$ is categorized x_i as a positive case. For example, a classifier Bayes $S(x_i) = \text{Prob}(P | x_i)$ so we consider a positive example c'omo $S(x_i) > 0.8$

The ROC curve is constructed by varying the threshold ρ between its minimum and maximum value and representing:

- On the Y axis the "positive rate certain" $TPR = TP$
- X in the "False Positive Rate" $FPR = FP$

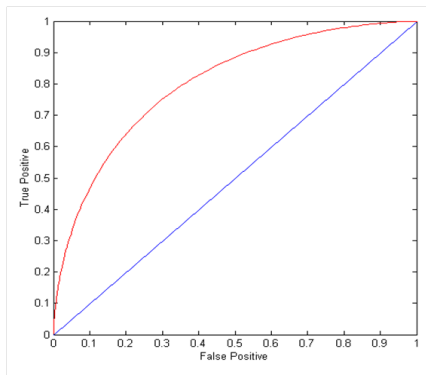
$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$

Evaluation of models classified

on

binary problems: ROC curves



Yes ρ It is m'nimo all positive (1.1) If ρ The same is m'aximo all negative n'

umber of false to true value on the diagonal If the curve above est'a **Good**

If the curve below est'a **Bad**

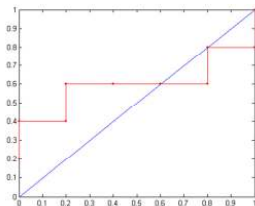
With values close to the curve (0.1) the best area under the curve ≈ 1 perfect

Evaluation of models classified caci'

on

binary problems: ROC curves

Example costruccion



Ejemplo	P(+ E)	Clase
1	0.95	+
2	0.93	+
3	0.87	-
4	0.85	-
5	0.85	-
6	0.85	+
7	0.76	-
8	0.53	+
9	0.43	-
10	0.25	+

Clase	+	-	+	-	+	-	+	-	+	-	+	-
TP	8	4	4	3	3	3	3	3	3	3	3	3
FP	3	3	3	3	3	3	3	3	3	3	3	3
TN	3	3	3	3	3	3	3	3	3	3	3	3
FN	3	3	3	3	3	3	3	3	3	3	3	3
TPR	3	3	3	3	3	3	3	3	3	3	3	3
FPR	3	3	3	3	3	3	3	3	3	3	3	3



M'etodos for evaluaci' on models classified caci'on

issue In a real case, with a given dataset. ¿C'omo organize training and test sets?

- Precisi'on to evaluate a model classi fi caci'on, the whole test must be independent
- Divide the dataset into two. For example, podr'amos reserve 2/3 of the examples available for training and the remaining 1/3 utilizar'amos the test set.
- *issue* ¿C'omo split, qu'e data will training and what to test?
- A first idea *random Selecci'on* . Examples of the test set are drawn
 - By drawing overall uniform
 - Strati fi ed by lot seg' one classes



Methods for evaluation on models classified

Validation on cross

issue

Selection random s'olo once made can be biased

- As solution repeats h Sometimes the process and the average precision ser'a $acc = \text{one} / h \sum$
 h

$$i = \text{one} \text{ } acc_i$$

- Another alternative: *Validation cross*

one. the dataset is divided into h equal parts

two. They are caught $h - \text{one}$ parts and training of test.

3. Leaves of varying the test to repeat process h times. The precision is average.



Methods for evaluation on models classification

Validation on cross

Variants validation cross

Two-fold cross validation . In this case $h = 2$.

Leave-one-out . If we have N examples in the data set, we divide N times, leaving $N-1$ training and test case. N execution process performed classification on.

Validation cross stratify each . Partitions are maintaining the initial proportion of elements in each class



M'etodos for evaluaci' on models classified caci'on

Bootstrapping

- the training set is sampled with replacement. With what examples can be repeated
- If N is sufficiently large sample fi tama~ It contains no N about 63.2% of the examples.
- The data are chosen not part of the set test
- **The process is repeated b times with an accuracy of ϵ_i , $i = \text{one}.. b$.**
- Precisi'on total can be calculated in various ways, the usual M'as is:

$$acc_{Boost} = 0.632 (1 / b) \sum_{i=one}^b \epsilon_i + 0.368 acc_s$$

where acc_s precisi'on is obtained using the set of Total Training



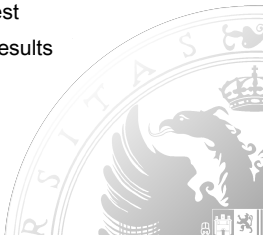
Comparaci' on classi fi ers

b'asicas ideas

Classi fi ers to compare:

Comparaci'on gen'etica .

- one.** a set of problems is chosen
- two.** a quality measure is chosen. Habitually precision.
- 3.** is eval' ohn and performs alg' a statistical test (Mean difference, ANOVA etc.) for comparing results



Comparison of classifiers

basic ideas

Classifiers to compare:

Comparison against particular problem .

- If the classifiers are binary ROC curves can be used to compare your action on a particular problem.
- They can be used, validation cross or bootstrapping to generate experiments and compare results without socks.

Not classifiers easily compare in general, Most likely that some work better than others on specific problems

one class

