
Induction of Rules



Lecturer: JERZY STEFANOWSKI
Institute of Computing Sciences
Poznan University of Technology
Poznan, Poland
Lecture 6 /7
SE Master Course 2008/2009

Źródła

- Wykład częściowo oparty na moim wykładzie szkoleniowym dla COST Action Spring School on Data Mining and MCDA – Troina 2008 oraz wcześniejszych wystąpieniach konferencyjnych.
- Proszę także przeczytać stosowane rozdziały z mojej rozprawy habilitacyjnej – dostępna na mojej stronie WWW.

Outline of this lecture

1. Rule representation
2. Basic algorithms for rule induction – idea of „Sequential covering”
3. MODLEM → exemplary algorithm for inducing a minimal set of rules.
4. Classification strategies
5. Descriptive properties of rules
6. Explore → discovering a richer set of rules
7. Logical relations (ILP) and rule induction
8. Final remarks

Rules - preliminaries

- **Rules** → popular symbolic representation of knowledge derived from data;
 - Natural and easy form of representation → possible inspection by human and their interpretation.
- Standard form of rules
IF *Conditions* THEN *Class*
- Other forms: Class IF Conditions; Conditions → Class

Example: The set of decision rules induced from PlaySport:

if outlook = overcast **then** Play = yes

if temperature = mild **and** humidity = normal **then** Play = yes

if outlook = rainy **and** windy = FALSE **then** Play = yes


if humidity = normal **and** windy = FALSE **then** Play = yes

if outlook = sunny **and** humidity = high **then** Play = no

if outlook = rainy **and** windy = TRUE **then** Play = no

Polish contribution – prof. Ryszard Michalski

- Father of Machine Learning and rule induction



Ryszard S. Michalski
(1937 - 2007)

PRC Chaired Professor of Computational Sciences and Health Informatics
Director of the Center for Discovery Science and Health Informatics

George Mason University

This page has been visited **15691** since January 1, 1999

6/27/06 R.S. Michalski gives a keynote address at the International Conference on Machine Learning, to celebrate the return of the conference to Carnegie-Mell after 26 years since the very first conference was organized there by Carbonell, Michalski and Mitchell

Articles in Mason Gazette:

- 7/31/07 New Center to Help Investigators Discover New Knowledge in Medical Databases
- 3/12/03 University Wins 10th Patent For Machine Learning Invention
- 11/19/02 Spotlight on Research: Grants Support Machine Learning and Inference Research
- 7/27/00 Michalski Receives Prestigious Science Honor

Interests

Research areas:
Machine Learning, Data Mining and Knowledge Discovery, Inductive Databases and Knowledge Scouts, Non-Darwinian Evolutionary Computation and Practical applications of these areas to Bioinformatics, Medicine, User Modeling, Intrusion Detection, and Very Complex System Design.

More about fathers of machine learning

- J. Carbonel, R. Michalski, T. Mitchell (2006)



Rules – more preliminaries

- A set of rules – a disjunctive set of conjunctive rules.
- Also DNF form:
 - *Class IF Cond_1 OR Cond_2 OR ... Cond_m*
- Various types of rules in data mining
 - Decision / classification rules
 - Association rules
 - Logic formulas (ILP)
 - Other → action rules, ...
- **MCDA** → attributes with some additional preferential information and ordinal classes.

Why Decision Rules?

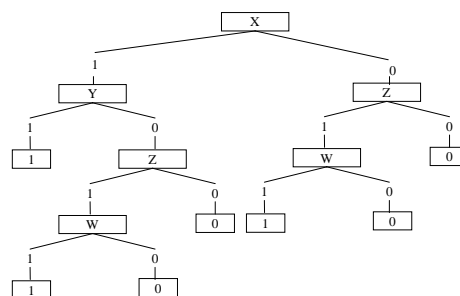
- Decision rules are more compact.
- Decision rules are more understandable and natural for human.
- Better for descriptive perspective in data mining.
- Can be nicely combined with background knowledge and more advanced operations, ...

Example: Let $X \in \{0,1\}$, $Y \in \{0,1\}$,
 $Z \in \{0,1\}$, $W \in \{0,1\}$. The rules are:

if $X=1$ and $Y=1$ **then** 1

if $Z=1$ and $W=1$ **then** 1

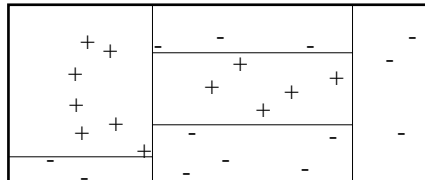
Otherwise 0;



Decision rules vs. decision trees:

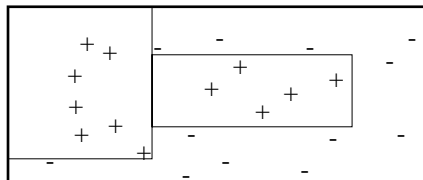
- Trees – splitting the data space (e.g. C4.5)

Decision boundaries of decision trees



- Rules – covering parts of the space (AQ, CN2, LEM)

Decision boundaries of decision rules



Rules – more formal notations

- A rule corresponding to class K_j is represented as

if P then Q

where $P = w_1$ and w_2 and ... and w_m is a condition part and Q is a decision part (object x satisfying P is assigned to class K_j)

- Elementary condition w_i ($a \text{ rel } v$), where $a \in A$ and v is its value (or a set of values) and rel stands for an operator as $=, <, \leq, \geq, >$.
- $[P]$ is a cover of a condition part of a rule \rightarrow a subset of examples satisfying P .
 - *if* ($a_2 = \text{small}$) *and* ($a_3 \leq 2$) *then* ($d = C1$) $\{x_1, x_7\}$

Rules - properties

- $B \rightarrow$ a set of examples from K_j ,
- A rule if P then Q is discriminant in DT iff $[P] = \bigcap [w_i] \subseteq B$,
- otherwise $(P \cap B \neq \emptyset)$ the rule is partly discriminating
 - Rule accuracy (or confidence) $|[P \cap K]|/|[P]|$
- Rule cannot not have a redundant condition part, i.e. there is no other $P^* \subset P$ such that $[P^*] \subseteq B$.
- Rule sets induced from DT
 - Minimal set of rules
 - Other sets of rules (all rules, satisfactory)

An example of rules induced from data table

Minimal set of rules

- if $(a_2 = s) \wedge (a_3 \leq 2)$ then $(d = C1)$ $\{x_1, x_7\}$
- if $(a_2 = n) \wedge (a_4 = c)$ then $(d = C1)$ $\{x_3, x_4\}$
- if $(a_2 = w)$ then $(d = C2)$ $\{x_2, x_6\}$
- if $(a_1 = f) \wedge (a_4 = a)$ then $(d = C2)$ $\{x_5, x_8\}$

Partly discriminating rule:

- if $(a_1 = m)$ then $(d = C1)$ $\{x_1, x_3, x_7 \mid x_6\}$ 3/4

id.	a_1	a_2	a_3	a_4	d
x_1	m	s	1	a	C1
x_2	f	w	1	b	C2
x_3	m	n	3	c	C1
x_4	f	n	2	c	C1
x_5	f	n	2	a	C2
x_6	m	w	2	c	C2
x_7	m	s	2	b	C1
x_8	f	s	3	a	C2

How to learn decision rules?

- Typical algorithms based on the scheme of a sequential covering and heuristically generate a minimal set of rule covering examples:
 - see, e.g., AQ, CN2, LEM, PRISM, MODLEM, Other ideas – PVM, R1 and RIPPER).
- Other approaches to induce „richer” sets of rules:
 - Satisfying some requirements (Explore, BRUTE, or modification of association rules, „Apriori-like”).
 - Based on local „reducts” → boolean reasoning or LDA.
- Specific optimization, eg. genetic approaches.
- Transformations of other representations:
 - Trees → rules.
 - Construction of (fuzzy) rules from ANN.



Covering algorithms

- A strategy for generating a rule set directly from data:
 - for each class in turn find a rule set that covers all examples in it (excluding examples not in the class).
- The main procedure is iteratively repeated for each class.
 - Positive examples from this class vs. negative examples.
- This approach is called a *covering* approach because at each stage a rule is identified that covers some of the instances.
- A sequential approach.
- For a given class it conducts in a stepwise way a general to specific search for the best rules (learn-one-rule) guided by the evaluation measures.

Original covering idea (AQ, Michalski 1969, 86)

for each class K_i **do**

$E_i := P_i \cup N_i$ (P_i positive, N_i negative example)

RuleSet(K_i) := empty

repeat {**find-set-of-rules**}

find-one-rule R covering some positive examples

and no negative ones

add R to RuleSet(K_i)

delete from P_i all pos. ex. covered by R

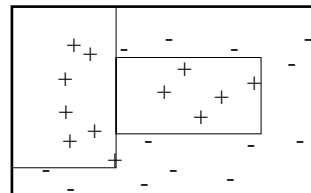
until P_i (set of pos. ex.) = empty

Find one rule:

Choosing a positive example called a seed.

Find a limited set of rules characterizing
the seed → **STAR**.

Choose the best rule according to LEF criteria.



Another variant – CN2 algorithm

- Clark and Niblett 1989; Clark and Boswell 1991
- Combine ideas AQ with TDIDT (search as in AQ, additional evaluation criteria or pruning as for TDIDT).
 - AQ depends on a seed example
 - Basic AQ has difficulties with noise handling
 - Latter solved by rule truncation (pos-pruning)
- Principles:
 - Covering approach (but stopping criteria relaxed).
 - Learning one rule – not so much example-seed driven.
 - Two options:
 - Generating an unordered set of rules (First Class, then conditions).
 - Generating an ordered list of rules (find first the best condition part than determine Class).

General schema of inducing minimal set of rules

- The procedure conducts a general to specific (greedy) search for the best rules (**learn-one-rule**) guided by the evaluation measures.
- At each stage add to the current condition part next elementary tests that optimize possible rule's evaluation (no backtracking).

Procedure Sequential covering (K_j Class; A attributes; E examples, τ - acceptance threshold);

```

begin
   $R := \emptyset;$     {set of induced rules}
   $r := \text{learn-one-rule}(Y_j \text{ Class; } A \text{ attributes; } E \text{ examples})$ 
  while  $\text{evaluate}(r, E) > \tau$  do
    begin
       $R := R \cup r;$ 
       $E := E \setminus [R];$     {remove positive examples covered by  $R$ }
       $r := \text{learn-one-rule}(K_j \text{ Class; } A \text{ attributes; } E \text{ examples});$ 
    end;
  return  $R$ 
end.
```



The contact lenses data



Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	No	Reduced	None
Young	Myope	No	Normal	Soft
Young	Myope	Yes	Reduced	None
Young	Myope	Yes	Normal	Hard
Young	Hypermetrope	No	Reduced	None
Young	Hypermetrope	No	Normal	Soft
Young	Hypermetrope	Yes	Reduced	None
Young	Hypermetrope	Yes	Normal	hard
Pre-presbyopic	Myope	No	Reduced	None
Pre-presbyopic	Myope	No	Normal	Soft
Pre-presbyopic	Myope	Yes	Reduced	None
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetrope	No	Reduced	None
Pre-presbyopic	Hypermetrope	No	Normal	Soft
Pre-presbyopic	Hypermetrope	Yes	Reduced	None
Pre-presbyopic	Hypermetrope	Yes	Normal	None
Presbyopic	Myope	No	Reduced	None
Presbyopic	Myope	No	Normal	None
Presbyopic	Myope	Yes	Reduced	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetrope	No	Reduced	None
Presbyopic	Hypermetrope	No	Normal	Soft
Presbyopic	Hypermetrope	Yes	Reduced	None
Presbyopic	Hypermetrope	Yes	Normal	None

Example: contact lens data 2

- Rule we seek: If ?
then recommendation = hard
- Possible conditions:

Age = Young	2/8
Age = Pre-presbyopic	1/8
Age = Presbyopic	1/8
Spectacle prescription = Myope	3/12
Spectacle prescription = Hypermetrope	1/12
Astigmatism = no	0/12
Astigmatism = yes	4/12
Tear production rate = Reduced	0/12
Tear production rate = Normal	4/12

ACK: slides coming from witten&eibe WEKA

Modified rule and covered data

- Condition part of the rule with the best elementary condition added:

If astigmatism = yes
then recommendation = hard

- Examples covered by condition part:

Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	Yes	Reduced	None
Young	Myope	Yes	Normal	Hard
Young	Hypermetrope	Yes	Reduced	None
Young	Hypermetrope	Yes	Normal	hard
Pre-presbyopic	Myope	Yes	Reduced	None
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetrope	Yes	Reduced	None
Pre-presbyopic	Hypermetrope	Yes	Normal	None
Presbyopic	Myope	Yes	Reduced	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetrope	Yes	Reduced	None
Presbyopic	Hypermetrope	Yes	Normal	None

Further specialization, 2

- Current state:

If astigmatism = yes and ? then recommendation = hard

- Possible conditions:

Age = Young	2/4
Age = Pre-presbyopic	1/4
Age = Presbyopic	1/4
Spectacle prescription = Myope	3/6
Spectacle prescription = Hypermetrope	1/6
Tear production rate = Reduced	0/6
Tear production rate = Normal	4/6

Two conditions in the rule

- The rule with the next best condition added:

If astigmatism = yes and tear production rate = normal then recommendation = hard

- Examples covered by modified rule:

Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	Yes	Normal	Hard
Young	Hypermetrope	Yes	Normal	hard
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetrope	Yes	Normal	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetrope	Yes	Normal	None

Further refinement, 4

- Current state:

```
If astigmatism = yes
    and tear production rate = normal
    and ?
then recommendation = hard
```

- Possible conditions:

Age = Young	2/2
Age = Pre-presbyopic	1/2
Age = Presbyopic	1/2
Spectacle prescription = Myope	3/3
Spectacle prescription = Hypermetrope	1/3

- Tie between the first and the fourth test
 - We choose the one with greater coverage

The result

- Final rule:

```
If astigmatism = yes
    and tear production rate = normal
    and spectacle prescription = myope
then recommendation = hard
```

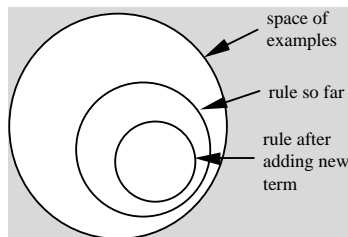
- Second rule for recommending “hard lenses”:
(built from instances not covered by first rule)

```
If age = young and astigmatism = yes
    and tear production rate = normal
then recommendation = hard
```

- These two rules cover all “hard lenses”:
 - Process is repeated with other two classes

A simple covering algorithm

- Generates a rule by adding tests that maximize rule's accuracy
- Similar to situation in decision trees: problem of selecting an attribute to split on
 - But: decision tree inducer maximizes overall purity
- Each new term reduces rule's coverage:



Evaluation of candidates in Learning One Rule

- When is a candidate for a rule R treated as “good”?
 - High accuracy $P(K|R)$;
 - High coverage $|[P]| = n$.
- Possible evaluation functions: $\frac{n_K(R)}{n(R)}$
 - *Relative frequency*:
 - where n_K is the number of correctly classified examples form class K , and n is the number of examples covered by the rule \rightarrow problems with small samples;
 - Laplace estimate:
 - Good for uniform prior distribution of k classes $\frac{n_K(R) + 1}{n(R) + k}$
 - *m-estimate of accuracy*: $(n_K(R) + mp) / (n(R) + m)$,
 - where n_K is the number of correctly classified examples, n is the number of examples covered by the rule, p is the prior probability of the class predicted by the rule, and m is the weight of p (domain dependent – more noise / larger m).

Other evaluation functions of rule R and class K

Assume rule R specialized to rule R'

- Entropy (Information gain and others versions).
- Accuracy gain (increase in expected accuracy)

$$P(K|R') - P(K|R)$$

- Many others
- Also weighted functions, e.g.

$$WAG(R', R) = \frac{n_K(R')}{n_K(R)} \cdot (P(K | R') - P(K | R))$$

$$WIG(R', R) = \frac{n_K(R')}{n_K(R)} \cdot (\log_2(K | R') - \log_2(K | R))$$

MODLEM – Algorithm for rule induction

- MODLEM [Stefanowski 98] generates a minimal set of rules.
- Its extra specificity – handling directly numerical attributes during rule induction; elementary conditions, e.g. $(a \geq v)$, $(a < v)$, $(a \in [v_1, v_2))$ or $(a = v)$.
- Elementary condition evaluated by one of three measures: class entropy, Laplace accuracy or Grzymala 2-LEF.

obj.	a1	a2	a3	a4	D
x1	m	2.0	1	a	C1
x2	f	2.5	1	b	C2
x3	m	1.5	3	c	C1
x4	f	2.3	2	c	C1
x5	f	1.4	2	a	C2
x6	m	3.2	2	c	C2
x7	m	1.9	2	b	C1
x8	f	2.0	3	a	C2

<i>if</i> (a1 = m) <i>and</i> (a2 ≤ 2.6) <i>then</i> (D = C1) {x1, x3, x7}
<i>if</i> (a2 ∈ [1.45, 2.4]) <i>and</i> (a3 ≤ 2) <i>then</i> (D = C1) {x1, x4, x7}
<i>if</i> (a2 ≥ 2.4) <i>then</i> (D = C2) {x2, x6}
<i>if</i> (a1 = f) <i>and</i> (a2 ≤ 2.15) <i>then</i> (D = C2) {x5, x8}

Procedure Modlem

Procedure MODLEM

```

(input  $B$  - a set of positive examples from a given decision concept;
  $criterion$  - an evaluation measure;
output  $T$  - single local covering of  $B$ , treated here as rule condition parts)
begin
   $G := B$ ; {A temporary set of rules covered by generated rules}
   $T := \emptyset$ ;
  while  $G \neq \emptyset$  do {look for rules until some examples remain uncovered}
  begin
     $T := \emptyset$ ; {a candidate for a rule condition part}
     $S := U$ ; {a set of objects currently covered by  $T$ }
    while ( $T = \emptyset$ ) or (not( $[T] \subseteq B$ )) do {stop condition for accepting a rule}
    begin
       $t := \emptyset$ ; {a candidate for an elementary condition}
      for each attribute  $q \in C$  do {looking for the best elementary condition}
      begin
         $new_t := \text{Find\_best\_condition}(q, S)$ ;
        if Better( $new_t, t, criterion$ ) then  $t := new_t$ ;
        {evaluate if a new condition is better than previous one
        according to the chosen evaluation measure}
      end;
       $T := T \cup \{t\}$ ; {add the best condition to the candidate rule}
       $S := S \cap [t]$ ; {focus on examples covered by the candidate}
    end; { while not( $[T] \subseteq B$  ) }
    for each elementary condition  $t \in T$  do
      if  $[T - t] \subseteq B$  then  $T := T - \{t\}$ ; {test a rule minimality}
     $T := T \cup \{T\}$ ; {store a rule}
     $G := B - \bigcup_{T \in \mathcal{T}} [T]$ ; {remove already covered examples}
  end; { while  $G \neq \emptyset$  }
  for each  $T \in T$  do
    if  $\bigcup_{T' \in \mathcal{T} - T} [T'] = B$  then  $T := T - T$  {test minimality of the rule set}
  end {procedure}

```

Set of positive examples

Looking for the best rule

Testing conjunction

Finding the most discriminatory single condition

Extending the conjunction

Testing minimality

Removing covered examples

Find best condition

function Find_best_condition

```

(input  $c$  - given attribute;  $S$  - set of examples; output  $best_I$  - bestcondition)
begin
   $best_I := \emptyset$ ;
  if  $c$  is a numerical attribute then
  begin
     $H :=$  list of sorted values for attribute  $c$  and objects from  $S$ ;
    {  $H(i)$  -  $i$ th unique value in the list }
    for  $i := 1$  to length( $H$ )-1 do
      if object class assignments for  $H(i)$  and  $H(i+1)$  are different then
      begin
         $v := (H(i) + H(i+1))/2$ ;
        create a  $new_I$  as either ( $c < v$ ) or ( $c \geq v$ );
        if Better( $new_I, best_I, criterion$ ) then  $best_I := new_I$ ;
      end
    end
  else { attribute is nominal }
  begin
    for each value  $v$  of attribute  $c$  do
      if Better( $(c = v), best_I, criterion$ ) then  $best_I := (c = v)$ ;
    end
  end {function}.

```

Preparing the sorted value list

Looking for the best cut point between class assignments

Testing each candidate

Return the best evaluated condition

An Example (1)



No.	Age	Job	Period	Income	Purpose	Dec.
1	m	u	0	500	K	r
2	sr	p	2	1400	S	r
3	m	p	4	2600	M	d
4	st	p	16	2300	D	d
5	sr	p	14	1600	M	p
6	m	u	0	700	W	r
7	sr	b	0	600	D	r
8	m	p	3	1400	D	p
9	sr	p	11	1600	W	d
10	st	e	0	1100	D	p
11	m	u	0	1500	D	p
12	m	b	0	1000	M	r
13	sr	p	17	2500	S	p
14	m	b	0	700	D	r
15	st	p	21	5000	S	d
16	m	p	5	3700	M	d
17	m	b	0	800	K	r

Class (Decision = r)

$$E = \{1, 2, 6, 7, 12, 14, 17\}$$

List of candidates

(Age=m) {1,6,12,14,17+; 3,8,11,16-}

(Age=sr) {2,7+; 5,9,13-}

(Job=u) {1,6+; 11-}

(Job=p) {2+, 3,4,8,9,13,15,16-}

(Job=b) {7,12,14,17+; \emptyset }

(Pur=K) {1,17+; \emptyset }

(Pur=S) {2+;13,15-}

{Pur=W} {6+, 9-}

{Pur=D} {7,14+; 4,8,10,11-}

{Pur=M} {12+;5,16-}

An Example (2)

- Numerical attributes: Income

500	600	700	800	1000	1100	1400	1500	1600	2300	2500	2600	3700	5000
1+	7+	6+	17+	12+	10-	2+	11-	9-	4-	13-	3-	10-	15-
		14+				8-		5-					

(Income < 1050) {1,6,7,12,14,17+; \emptyset }

(Income < 1250) {1,6,7,12,14,17+;10-}

(Income < 1450) {1,2,6,7,12,14,17+;8,10-}

Period

(Period < 1) {1,6,7,14,17+;10,11-}

(Period < 2.5) {1,2,6,7,12,14,17+;10,11-}

Example (3) - the minimal set of induced rule

1. if (Income<1050) then (Dec=r) [6]
 2. if (Age=sr) and (Period<2.5) then (Dec=r) [2]
 3. if (Period∈[3.5,12.5)) then (Dec=d) [2]
 4. if (Age=st) and (Job=p) then (Dec=d) [3]
 5. if (Age=m) and (Income∈[1050,2550)) then (Dec=p) [2]
 6. if (Job=e) then (Dec=p) [1]
 7. if (Age=sr) and (Period≥12.5) then (Dec=p) [2]
- For inconsistent data:
 - Approximations of decision classes (rough sets)
 - Rule post-processing (a kind of post-pruning) or extra testing and earlier acceptance of rules.

Mushroom data (UCI Repository)

- Mushroom records drawn from The Audubon Society Field Guide to North American Mushrooms (1981).
- This data set includes descriptions of hypothetical samples corresponding to 23 species of mushrooms in the Agaricus and Lepiota Family. Each species is identified as definitely edible, definitely poisonous, or of unknown edibility.
- Number of examples: 8124.
- Number of attributes: 22 (all nominally valued)
- Missing attribute values: 2480 of them.
- Class Distribution:
 - edible: 4208 (51.8%)
 - poisonous: 3916 (48.2%)

MOLDEM rule set (Implemented in WEKA)

=== Classifier model (full training set) ===

Rule 1.(odor is in: {n, a, l})&(spore-print-color is in: {n, k, b, h, o, u, y, w})&(gill-size = b)
=> (class = e); [3920, 3920, 93.16%, 100%]

Rule 2.(odor is in: {n, a, l})&(spore-print-color is in: {n, h, k, u}) => (class = e); [3488,
3488, 82.89%, 100%]

Rule 3.(gill-spacing = w)&(cap-color is in: {c, n}) => (class = e); [304, 304, 7.22%,
100%]

Rule 4.(spore-print-color = r) => (class = p); [72, 72, 1.84%, 100%]

Rule 5.(stalk-surface-below-ring = y)&(gill-size = n) => (class = p); [40, 40, 1.02%,
100%]

Rule 6.(odor = n)&(gill-size = n)&(bruises? = t) => (class = p); [8, 8, 0.2%, 100%]

Rule 7.(odor is in: {f, s, y, p, c, m}) => (class = p); [3796, 3796, 96.94%, 100%]

Number of rules: 7

Number of conditions: 14

Approaches to Avoiding Overfitting

- **Pre-pruning:** stop learning the decision rules before they reach the point where they perfectly classify the training data
- **Post-pruning:** allow the decision rules to overfit the training data, and then post-prune the rules.

Pre-Pruning

The criteria for stopping learning rules can be:

- **minimum purity** criterion requires a certain percentage of the instances covered by the rule to be positive;
- **significance test** determines if there is a significant difference between the distribution of the instances covered by the rule and the distribution of the instances in the training sets.

Post-Pruning

1. Split instances into *Growing Set* and *Pruning Set*,
2. Learn set *SR* of rules using *Growing Set*,
3. Find the best simplification *BSR* of *SR*.
4. **while** ($\text{Accuracy}(\text{BSR}, \text{Pruning Set}) >$
 $\text{Accuracy}(\text{SR}, \text{Pruning Set})$) **do**
 - 4.1 $\text{SR} = \text{BSR}$;
 - 4.2 Find the best simplification *BSR* of *SR*.
5. **return** *BSR*;

Applying rule set to classify objects

- **Matching** a new object description x to condition parts of rules.
 - Either object's description satisfies all elementary conditions in a rule, or not.
- IF (a1=L) and (a3 ≥ 3) THEN Class +
 $x \rightarrow (a1=L), (a2=s), (a3=7), (a4=1)$
- Two ways of assigning x to class K depending on the set of rules:
 - Unordered set of rules (AQ, CN2, PRISM, LEM)
 - Ordered list of rules (CN2, c4.5rules)

Applying rule set to classify objects

- The rules are ordered into priority decision list!
 - Another way of rule induction – rules are learned by first determining Conditions and then Class (CN2)
- Notice:** mixed sequence of classes K_1, \dots, K in a rule list
- But: ordered** execution when classifying a new instance: rules are sequentially tried and the first rule that 'fires' (covers the example) is used for final decision
- Decision list {R1, R2, R3, ..., D}**: rules R_i are interpreted as **if-then-else** rules
- If no rule fires, then DefaultClass (majority class in input data)

Priority decision list (C4.5 rules)

The screenshot displays the C4.5 software interface. The main window shows a list of rules generated from 16 attributes and 300 training cases. The rules are as follows:

- Rule 1: [99.4%] IF physician fee freeze = n THEN democrat
- Rule 2: [94.7%] IF mx missile = y AND synfuels corporation cutback = y THEN democrat
- Rule 3: [63.0%] IF physician fee freeze = u AND mx missile = n THEN democrat
- Rule 4: [94.0%] IF physician fee freeze = y AND immigration = y THEN republican
- Rule 5: [91.2%] IF physician fee freeze = y AND mx missile = n THEN republican
- Rule 6: [82.0%] IF adoption of the budget resolution = n AND education spending = u THEN republican
- Rule 7: [50.0%] IF physician fee freeze = u AND mx missile = u THEN republican

The default class is 'democrat'. Errors in training set: 11 (3.7%), Errors in test set: 6 (4.4%).

Below the rules, a 'Cross-validation (rules)' window shows the following table:

Ruleset	Size	Errors	Errors (test)
1	5	10 (3.7%)	1 (3.3%)
2	5	10 (3.7%)	1 (3.3%)
3	5	11 (4.1%)	0 (0.0%)
4	4	10 (3.7%)	3 (10.0%)
5	5	9 (3.3%)	2 (6.7%)
6	4	11 (4.1%)	2 (6.7%)
7	5	11 (4.1%)	0 (0.0%)
8	5	10 (3.7%)	1 (3.3%)
9	2	12 (4.4%)	3 (10.0%)
10	3	11 (4.1%)	2 (6.7%)

A 'Confusion matrix (test set)' window shows the following data:

Dep. \ C45	democrat	republican
democrat	18	1
republican	1	11

Specific list of rules - RIPPER (Mushroom data)

The screenshot shows the Weka Explorer interface with the RIPPER classifier applied to Mushroom data. The classifier window displays the following rules:

- [odor = e] => class=p (2160.0/0.0)
- [gill-size = n] and [gill-color = b] => class=p (1152.0/0.0)
- [gill-size = n] and [odor = p] => class=p (256.0/0.0)
- [odor = c] => class=p (192.0/0.0)
- [spore-print-color = l] => class=p (72.0/0.0)
- [stalk-surface-above-ring = k] and [gill-spacing = c] => class=p (68.0/0.0)
- [habitat = l] and [cap-color = w] => class=p (8.0/0.0)
- [stalk-color-above-ring = y] => class=p (8.0/0.0)
- => class=e (4208.0/0.0)

Number of Rules : 9
Time taken to build model: 4.11 seconds

--- Stratified cross-validation ---
=== Summary ===

Metric	Value	Percentage
Correctly Classified Instances	8124	100 %
Incorrectly Classified Instances	0	0 %
Kappa statistic	1	1 %
Mean absolute error	0	
Root mean squared error	0	
Relative absolute error	0	%
Root relative squared error	0	%
Total Number of Instances	8124	

=== Detailed Accuracy By Class ===

TP Rate	FP Rate	Precision	Recall	F-Measure	Class
1	0	1	1	1	e
1	0	1	1	1	p

--- Confusion Matrix ---

a	b	<- classified as
4208	0	a = e
0	3916	b = p

Learning ordered set of rules

- RuleList := empty; $E_{cur} := E$
- **repeat**
 - learn-one-rule R
 - RuleList := RuleList ++ R
 - $E_{cur} := E_{cur} - \{\text{all examples covered by R}\}$
(Not only positive examples !)
- **until** performance(R, E_{cur}) < ThresholdR
- RuleList := sort RuleList by performance(R,E)
- RuleList := RuleList ++ DefaultRule(E_{cur})

CN2 – unordered rule set

```
WinCn2_16.attributes (crx.aex) 490 examples (crx.aex) 30 rules (induced)
Data Rules Cross-validation Trace Output
Unordered Laplacian Unset 5 0.05 10 0
Reading attributes and examples...
490 examples!
Finished reading attribute and example file!
Running CN on current example set...
Finished inducing rules!

-----
| UN-ORDERED RULE LIST |
-----
IF A8 < 10.75
AND A9 = T
AND 5.50 < A11 < 18.50
THEN DECISION = Y [69 0]

IF A15 > 5676.00
THEN DECISION = Y [19 0]

IF A2 > 19.00
AND A4 = U
AND A8 < 11.75
AND A9 = T
AND A14 < 91.00
THEN DECISION = Y [67.50 0]

IF A3 > 1.70
AND A9 = T
AND A15 > 241.50
THEN DECISION = Y [80 0]

IF A6 = X
AND 1.33 < A8 < 7.88
THEN DECISION = Y [11 0]

IF A2 < 26.08
AND A9 = T
AND 20.00 < A14 < 106.00
THEN DECISION = Y [32.50 0]

IF A0 > 12.75
AND A14 < 107.00
THEN DECISION = Y [12 0]

Lister - [c:\User\jurek\students\CichyCN2\Exo\Examples\crx.aex]
Plik Edytuj Opcje Pomoc
**ATTRIBUTE AND EXAMPLE FILE**
R1: B A;
R2: (FLOAT)
R3: (FLOAT)
R4: U V L;
R5: G P GG;
R6: M Q H R CC K C D X I E AA FF J;
R7: V H BB FF J Z U DD N;
R8: (FLOAT)
R9: T F;
R10: T F;
R11: (FLOAT)
R12: F I;
R13: G S P;
R14: (FLOAT)
R15: (FLOAT)
DECISION: V H;

@
@ 30.89 0 U G W V 1.25 T T 1 F G 202 0 V;
A 58.67 4.46 U C Q W 3.04 T T 6 F C 43 560 V;
A 24.50 .5 U C Q H 1.5 T F 0 F C 280 824 V;
B 27.80 1.54 U G W U 3.75 T T 5 T G 100 0 V;
B 20.17 5.625 U G W U 1.71 T F 0 F 5 120 0 V;
B 32.08 4 U G H V 2.5 T F 0 T G 360 0 V;
B 33.17 1.04 U Q R H 6.5 T F 0 T C 164 31285 V;
A 22.92 11.585 U G CC U .04 T F 0 F C 80 1949 V;
B 54.42 .5 V P K H 3.96 T F 0 F C 180 314 V;
B 42.50 4.915 V P W U 3.165 T F 0 T C 52 1442 V;
B 29.08 .89 U C C U 2.46 C F 0 T C 420 0 V;
```

Applying unordered rule set to classify objects

- An unordered set of rules → three situations:
 - Matching to rules indicating the same class.
 - Multiple matching to rules from different classes.
 - No matching to any rule.
- An example:
- $e1 = \{(Age=m), (Job=p), (Period=6), (Income=3000), (Purpose=K)\}$
 - rule 3: if $(Period \in [3.5, 12.5])$ then $(Dec=d)$ [2]
 - Exact matching to rule 3. → Class $(Dec=d)$
- $e2 = \{(Age=m), (Job=p), (Period=2), (Income=2600), (Purpose=M)\}$
 - No matching!

Solving conflict situations

- LERS classification strategy (Grzymala 94)
 - Multiple matching
 - Two factors: $Strength(R)$ – number of learning examples correctly classified by R and final class $Support(Y_i)$:

$$\sum_{\text{matching rules } R \text{ for } Y_i} Strength(R)$$
 - Partial matching
 - Matching factor $MF(R)$ and

$$\sum_{\text{partially match. rules } R \text{ for } Y_i} MF(R) \cdot Strength(R)$$
- $e2 = \{(Age=m), (Job=p), (Period=2), (Income=2600), (Purpose=M)\}$
 - Partial matching to rules 2, 4 and 5 for all with $MF = 0.5$
 - $Support(r) = 0.5 \cdot 2 = 1$; $Support(d) = 0.5 \cdot 2 + 0.5 \cdot 2 = 2$
- Alternative approaches – e.g. nearest rules (Stefanowski 95)
- Instead of MF use a kind of normalized distance x to conditions of r

Some experiments

- Analysing strategies (total accuracy in [%]):

data set	all	multiple	exact
large soybean	87.9	85.7	79.2
election	89.4	79.5	71.8
hsv2	77.1	70.5	59.8
concretes	88.9	82.8	81.0
breast cancer	67.1	59.3	51.2
imidazolium	53.3	44.8	34.4
lymphography	85.2	73.6	67.6
oncology	83.8	82.4	74.1
buses	98.0	93.5	90.8
bearings	96.4	90.9	87.3

- Comparing to other classification approaches
 - Depends on the data
 - Generally → similar to decision trees

Variations of inducing minimal sets of rules

- Sequential vs. simultaneous covering of data.
- General-to-specific vs. specific-to-general; begin search from single most general vs. many most specific starting hypotheses.
- Generate-and-test vs. example driven (as in AQ).
- Pre-pruning vs. post-pruning of rules
- What evaluation functions to use?
- ...

Different perspectives of rule application

- In a descriptive perspective
 - To present, analyse the relationships between values of attributes, to explain and understand classification patterns
- In a prediction/classification perspective,
 - To predict value of decision class for new (unseen) object)

Perspectives are different;
Moreover rules are evaluated in a different ways!

Evaluating single rules

- rule r (if P then Q) derived from DT , examples U .

	Q	$\neg Q$	
P	n_{PQ}	$n_{P\neg Q}$	n_P
$\neg P$	$n_{\neg PQ}$	$n_{\neg P\neg Q}$	$n_{\neg P}$
	n_Q	$n_{\neg Q}$	n

- Reviews of measures, e.g.
 - Yao Y.Y, Zhong N., An analysis of quantitative measures associated with rules, In: Proc. the 3rd Pacific-Asia Conf. on Knowledge Discovery and Data Mining, LNAI 1574, Springer, 1999, pp. 479-488.
 - Hilderman R.J., Hamilton H.J, Knowledge Discovery and Measures of Interest. Kluwer, 2002.
- Support of rule r $G(P \wedge Q) = \frac{n_{PQ}}{n}$ Coverage $AS(P|Q) = \frac{n_{PQ}}{n_Q}$
- Confidence of rule r $AS(Q|P) = \frac{n_{PQ}}{n_P}$ and others ...

Other descriptive measures

**Change of support – confirmation of supporting Q by a premise P
(Piatetsky-Shapiro)**

$$CS(Q|P) = AS(Q|P) - G(Q)$$

where $G(Q) = \frac{n_Q}{n}$

Interpretation: Zakres wartości od -1 do +1 ; Różnica między prawdopodobieństwami a priori i a posterior; dodatnie wartości wystąpienie przesłanki P powoduje konkluzję Q; ujemna wartość wskazuje że nie ma wpływu.

Degree of independence:

$$IND(Q,P) = \frac{G(P \wedge Q)}{G(P) \cdot G(Q)}$$

Aggregated measures

Based on previous measures:

Significance of a rule (propozycja Yao i Liu)

$$S(Q|P) = AS(Q|P) \cdot IND(Q,P)$$

Klosgen's measure of interest

$$K(Q|P) = G(P)^\alpha \cdot (AS(Q|P) - G(Q))$$

Michalski's weighted sum

$$WSC(Q|P) = w_1 \cdot AS(Q|P) + w_2 \cdot AS(P|Q)$$

The relative risk (Ali, Srikant):

$$r(Q|P) = \frac{AS(Q|P)}{AS(Q|\neg P)}$$

Descriptive requirements to single rules

- In descriptive perspective users may prefer to discover rules which should be:
 - strong / general – high enough rule coverage $AS(P/Q)$ or support.
 - accurate – sufficient accuracy $AS(Q/P)$.
 - simple (e.g. which are in a limited number and have short condition parts).
 - Number of rules should not be too high.
- Covering algorithms biased towards minimum set of rules - containing only a limited part of potentially 'interesting' rules.
 - We need another kind of rule induction algorithms!

Explore algorithm (Stefanowski, Vanderpooten)

- Another aim of rule induction
 - to extract from data set inducing **all rules** that *satisfy* some *user's requirements* connected with *his interest* (regarding, e.g. the strength of the rule, level of confidence, length, sometimes also emphasis on the syntax of rules).
- Special technique of exploration the space of possible rules:
 - Progressively generation rules of increasing size using in the most efficient way some 'good' pruning and stopping condition that reject unnecessary candidates for rules.
- Similar to adaptations of Apriori principle for looking frequent itemsets [AIS94]; Brute [Etzioni]

Explore – some algorithmic details

```

procedure Explore (LS: list of conditions;
  SC: stopping conditions; var R:
  set_of_rules);
begin
  R ← ∅;
  Good_Candidates(LS,R); {LS - ordered
  list of  $c_1, c_2, \dots, c_n$ }
  Q ← LS; {create a queue Q}
  while Q ≠ ∅ do
    begin
      select the first conjunction C from Q ;
      Q ← Q \ {C};
      Extend(C,LC); {LC - list of extended
      conjunctions}
      Good_Candidates(LC,R);
      Q ← Q ∪ C; {place all conjunctions from
      LC at the end of Q}
    end
end.
  
```

```

procedure Extend(C : conjunction, var L : list of
  conjunctions);
  {This procedure puts in list L extensions of
  conjunction C that are possible candidates
  for rules}
begin
  Let k be the size of C and h be the highest index
  of elementary conditions involved in C;
  L ← { $C \wedge c_{h+i}$  where  $ch+i \in LS$  and such that all the
  k subconjunctions of  $C \wedge c_{h+i}$  of size k and
  involving  $c_{h+i}$  belong to Q,  $i=1, \dots, n-h$ }
end
procedure Good_Candidates(LC : list of
  conjunctions, var R - set of rules );
  {This procedure prunes list LC discarding:
  - conjunctions whose extension cannot give rise
  to rules due to SC,
  - conjunctions corresponding to rules which are
  already stored in R
  
```

Various sets of rules (Stefanowski and Vanderpooten 1994)

Table 1: The illustrative set of learning exam

- A minimal set of rules (LEM2):

rule 1.	if $(q_1 = 2) \wedge (q_3 = 1)$ then $(d = 1)$	{1, 2, 3, 4, 5}	5/8
rule 2.	if $(q_1 = 1)$ then $(d = 1)$	{6, 7}	2/8
rule 3.	if $(q_3 = 2) \wedge (q_6 = 2)$ then $(d = 1)$	{6, 8}	2/8
rule 4.	if $(q_1 = 3)$ then $(d = 2)$	{9, 10, 11, 13, 14}	5/7
rule 5.	if $(q_3 = 3)$ then $(d = 2)$	{15}	1/7
rule 6.	if $(q_3 = 2) \wedge (q_4 = 1) \wedge (q_6 = 1)$ then $(d = 2)$	{12}	1/7

- A „satisfactory” set of rules (Explore):

Let us assume that the user's level of interest to the possible strength of a rule by assigning a value $l = 50\%$ in *SC*.

Explore gives the following decision rules:

rule 1.	if $(q_2 = 3)$ then $(d = 1)$	{1, 2, 3, 6, 7}	5/8
rule 2.	if $(q_1 = 2) \wedge (q_3 = 1)$ then $(d = 1)$	{1, 2, 3, 4, 5}	5/8
rule 3.	if $(q_1 = 3)$ then $(d = 2)$	{9, 10, 11, 13, 14}	5/7
rule 4.	if $(q_4 = 2)$ then $(d = 2)$	{10, 13, 14, 15}	4/7

No.	q_1	q_2	q_3	q_4	q_5	q_6	d
1	2	3	1	3	1	2	1
2	2	3	1	1	1	1	1
3	2	3	1	3	2	1	1
4	2	1	1	1	1	1	1
5	2	2	1	1	2	2	1
6	1	3	2	3	1	2	1
7	1	3	2	3	2	1	1
8	2	1	2	1	2	2	1
9	3	1	1	3	1	2	2
10	3	1	2	2	2	1	2
11	3	1	1	3	2	2	2
12	2	1	2	1	2	1	2
13	3	2	4	2	1	1	2
14	3	2	4	2	2	1	2
15	2	2	3	2	1	2	2
16	2	2	2	1	1	1	1
17	2	2	2	1	1	1	2

A diagnostic case study

- A fleet of homogeneous 76 buses (AutoSan H9-21) operating in an inter-city and local transportation system.
- The following symptoms characterize these buses :
 - s_1 – maximum speed [km/h],
 - s_2 – compression pressure [Mpa],
 - s_3 – blacking components in exhaust gas [%],
 - s_4 – torque [Nm],
 - s_5 – summer fuel consumption [l/100km],
 - s_6 – winter fuel consumption [l/100km],
 - s_7 – oil consumption [l/1000km],
 - s_8 – maximum horsepower of the engine [km].

Experts' classification of buses:

1. Buses with engines in a good technical state – further use (46 buses),
2. Buses with engines in a bad technical state – requiring repair (30 buses).

LEM2 algorithm – (sequential covering)

- A **minimal** set of *discriminating* decision rules
 1. if ($s_2 \geq 2.4$ MPa) & ($s_7 < 2.1$ l/1000km) then (technical state=good) [46]
 2. if ($s_2 < 2.4$ MPa) then (technical state=bad) [29]
 3. if ($s_7 \geq 2.1$ l/1000km) then (technical state=bad) [24]
- The prediction accuracy ('leaving-one-out' reclassification test) is equal to 98.7%.

Another set of rules (EXPLORE)

All decision rules with min. SC1 threshold (rule coverage > 50%):

1. if (s1>85 km/h) then (technical state=good) [34]
2. if (s8>134 km) then (technical state=good) [26]
3. if (s2≥2.4 MPa) & (s3<61 %) then (technical state=good) [44]
4. if (s2≥2.4 MPa) & (s4>444 Nm) then (technical state=good) [44]
5. if (s2≥2.4 MPa) & (s7<2.1 //1000km) then (technical state=good) [46]
6. if (s3<61 %) & (s4>444 Nm) then (technical state=good) [42]
7. if (s1≤77 km/h) then (technical state=bad) [25]
8. if (s2<2.4 MPa) then (technical state=bad) [29]
9. if (s7≥2.1 //1000km) then (technical state=bad) [24]
10. if (s3≥61 %) & (s4≤444 Nm) then (technical state=bad) [28]
11. if (s3≥61 %) & (s8<120 km) then (technical state=bad) [27]

The prediction accuracy - 98.7%

Descriptive vs. classification properties (Explore)

Data set	Stopping conditions		Number of rules	Average rule length (# cond.)	Average rule strength (# exam.)	Classification accuracy (%)
	SC1	SC2				
Iris	All rules		80	2.1	6.03	92.67
	5%	---	35	1.89	12.23	92.67
	10%	---	22	1.86	17.27	92
	15%	---	20	1.85	18.4	90
	20%	---	15	1.8	21.6	83.33
	25%	---	14	1.79	22.36	78.67
	30%	---	6	1.83	33.83	60.67
	Minimum	rule set	23	1.91	11	95.33
Tic-tac-toe	All rules		3858	4.63	4.27	91.35
	5%	5	16	3	60.25	97.19
	10%	5	16	3	60.25	96.14
	15%	5	2	3	50	---
	20%	5	0	---	---	---
	30%	5	0	---	---	---
		Minimum	rule set	24	3.67	40.83
Voting	All rules		1502	4.723	10.61	95.87
	5%	4	231	3.6	45.86	94.51
	10%	4	138	3.3	66.96	94.5
	15%	4	104	3.1	79.61	93.8
	20%	4	82	3.1	89.87	94
	25%	4	67	3.1	96.99	93.32
	30%	4	50	3.1	104.7	93.31
	Minimum	rule set	26	3.69	43.77	95.87
Election	All rules		>200000	---	---	---
	10%	---	828	3.48	26.91	89.39
	15%	---	87	3.05	33.82	87.37
	20%	---	8	2.38	53.75	73.88
	25%	---	2	1.5	79	32.96
	30%	---	1	1	105	23.64
		Minimum	rule set	48	3.27	21.176

- Tuning a proper value of stopping condition SC (rule coverage) leads to sets of rules which are „satisfactory” with respect to a number of rules, average rule length and average rule strength without decreasing too much the classification accuracy.

Preference ordered data

- MCDA vs. traditional classification (ML & Stat):
 - Attributes with preference ordered domains → criteria.
 - Ordinal classes rather than nominal labels.
 - „Semantic correlation” between values of criteria, and classes.
 - For objects x, y if $a(x) \preceq a(y)$ then their labels $\lambda(x) \preceq \lambda(y)$
- Possible inconsistency

Client	Month salary	Account status	Credit risk
A	9000	high	low
B	4000	medium	medium
C	5500	medium	high

- Dominance based rough set approach to handle it
 - Greco S., Matarazzo B., Slowinski R.

Dominance based decision rules

- Induced from rough approximations of unions of classes (upward and downward):
 - *certain* $D \geq$ -decision rules, supported by objects $\in Cl_i^{\geq}$ without ambiguity:

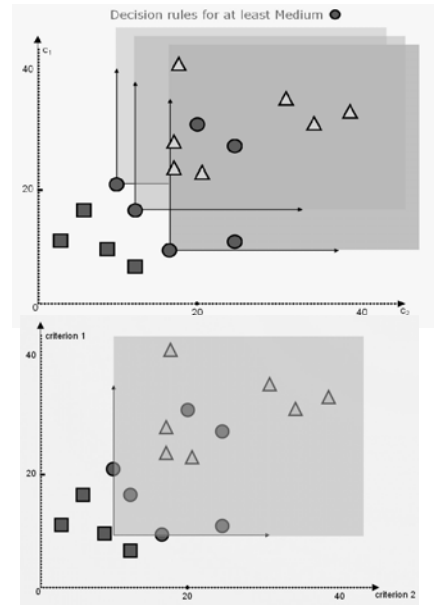
if $q_1(x) \succeq_{q_1} r_{q_1}$ and $q_2(x) \succeq_{q_2} r_{q_2}$ and ... $q_p(x) \succeq_{q_p} r_{q_p}$ then $x \in Cl_i^{\geq}$
 - *possible* $D \geq$ -decision rules, supported by objects $\in Cl_i^{\geq}$ and ambiguous ones from its upper approximation:

if $q_1(x) \succeq_{q_1} r_{q_1}$ and $q_2(x) \succeq_{q_2} r_{q_2}$ and ... $q_p(x) \succeq_{q_p} r_{q_p}$, then x possibly $\in Cl_i^{\geq}$
 - *certain* $D \leq$ -decision rules, supported by objects $\in Cl_i^{\leq}$ without ambiguity:

if $q_1(x) \preceq_{q_1} r_{q_1}$ and $q_2(x) \preceq_{q_2} r_{q_2}$ and ... $q_p(x) \preceq_{q_p} r_{q_p}$, then $x \in Cl_i^{\leq}$

Algorithms for inducing dominance based rules

- Greco, Slowinski, Stefanowski, Blaszczynski, Dembczyński and others – a number of proposals
- Minimal sets of rules:
 - DOMLEM → adaptation of ideas behind MODLEM.
- DOMApriori → richer set of rules
- Robust rules → syntax based on an object from data table.
 - All rules → modifications of boolean reasoning
 - Glance → incremental learning.



Software from PUT

File Show Calculate Apply Rules Report Window Help

View All Attributes

Attributes

System of Classes

Quality of Approximation

Units of Classes

File Options: Rule Type All Length: Filter

Statistics: Average Length: 2.32 Average Strength: 5.11 Max. Strength: 25 Generated Rules: 30 Displayed Rules: 30 Calculation time: 00:00:00:07

Number	Condition	Decision	Stren.	Relat.
1	$! \text{leaf} \rightarrow 46$	decision at most 1	4	23.51 %
2	$(\text{leaf} \rightarrow 0.8) \& (\text{type} \rightarrow 1.15)$	decision at most 1	2	11.76 %
3	$(\text{leaf} \rightarrow 1.2) \& (\text{type} \rightarrow 0.88)$	decision at most 1	2	11.76 %
4	$(\text{leaf} \rightarrow 0.75) \& (\text{type} \rightarrow 0.3)$	decision at most 1	2	11.76 %
5	$(\text{leaf} \rightarrow 0.75) \& (\text{type} \rightarrow 0.9) \& (\text{leaf} \rightarrow 1)$	decision at most 1	2	11.76 %
6	$(\text{leaf} \rightarrow 0.34) \& (\text{type} \rightarrow 1.2) \& (\text{leaf} \rightarrow 2)$	decision at most 1	2	11.76 %
7	$! \text{leaf} \rightarrow 1.25 \& (\text{type} \rightarrow 0.82)$	decision at most 1	1	5.88 %
8	$(\text{leaf} \rightarrow 0.34) \& (\text{type} \rightarrow 0.88) \& (\text{leaf} \rightarrow 1.38)$	decision at most 1	1	5.88 %
9	$(\text{leaf} \rightarrow 0.82) \& (\text{type} \rightarrow 1.8) \& (\text{leaf} \rightarrow 0.98) \& (\text{leaf} \rightarrow 1)$	decision at most 1	3	17.65 %
10	$(\text{leaf} \rightarrow 0.82) \& (\text{type} \rightarrow 1.3) \& (\text{leaf} \rightarrow 1)$	decision at most 1	1	5.88 %
11	$! \text{leaf} \rightarrow 1.95$	decision at most 2	1	2.44 %
12	$(\text{leaf} \rightarrow 0.82) \& (\text{type} \rightarrow 1)$	decision at most 2	1	5.88 %
13	$! \text{leaf} \rightarrow 2 \& (\text{leaf} \rightarrow 1.15)$	decision at most 2	13	31.71 %
14	$(\text{leaf} \rightarrow 0.95) \& (\text{leaf} \rightarrow 1.2)$	decision at most 2	4	9.76 %
15	$(\text{leaf} \rightarrow 1.15) \& (\text{leaf} \rightarrow 1)$	decision at most 2	25	60.98 %
16	$(\text{leaf} \rightarrow 1.2) \& (\text{leaf} \rightarrow 1) \& (\text{leaf} \rightarrow 1)$	decision at most 2	13	29.71 %
17	$(\text{leaf} \rightarrow 0.37) \& (\text{leaf} \rightarrow 0.37)$	decision at most 2	1	2.44 %
18	$(\text{leaf} \rightarrow 1.08) \& (\text{leaf} \rightarrow 2)$	decision at most 3	10	50.00 %
19	$(\text{leaf} \rightarrow 0.9) \& (\text{leaf} \rightarrow 0.34)$	decision at most 3	3	15.00 %

Supporting Examples:

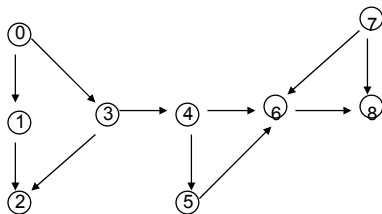
leaf	type	leaf	type	leaf	type	leaf	type	leaf	type	leaf	type	leaf	type	leaf	type	leaf	type	
0	1.15	0.94	1.3	1	1.06	1.56	1.3	1	2	0.71	0.91	1	1.21	1.14	1	1.3	1.3	
1	1.15	0.94	1.3	1	1.3	1.21	1.15	1	3	0.86	1	0.86	1.1	1.07	1	0.91	1	
2	1.15	0.94	1.3	0.83	1.11	1.56	1	1.07	3	0.86	0.82	0.86	0.9	1	1	1	1	
3	1.15	0.94	1.3	0.83	1.11	1.56	1	1.07	3	0.86	0.82	0.86	0.9	1	1	1	1	
4	1.15	0.94	1.3	1	1.11	1.06	1	1	3	0.86	0.82	0.86	1	0.95	1	0.91	1.08	
5	1.15	0.94	1.3	0.96	1.11	1.06	1.15	1	2	0.71	1	0.7	1.3	1	2	0.82	1	
6	1.15	0.94	1.05	1	1.3	1.95	1.15	1	4	0.86	1	0.7	1.1	1.07	1	1.1	1.24	1.23
7	1.15	0.94	1.3	1	1.3	1.06	1.15	0.87	3	0.86	1.13	0.86	1.21	1.14	1	0.91	1	1.23
8	1.15	0.94	1.15	1	1.25	1.21	1	0.87	3	1	1	1	1	1	2	0.82	1.3	1.08
9	1.15	0.94	1.3	1	1.06	1.21	1.15	1	2	1	0.91	1	1.3	1	1	1.24	1.24	1
10	0.94	1.3	0.83	1	1	1	0.87	3	0.86	0.82	1.17	1	1	1	1	1.3	1	1
11	0.94	1.3	1	1	1	1	0.87	0.87	1	1	0.82	0.7	0.9	0.95	2	0.91	0.91	1
12	1.15	0.94	1.3	1	1	1	0.87	0.87	1	0.71	0.82	0.7	1.3	1.07	1	1.3	1	1.04

Learning First Order Rules

- Is object/attribute table sufficient data representation?
- Some limitations:
 - Representation expressiveness – unable to express relations between objects or object elements. ,
 - *background knowledge* sometimes is quite complicated.
- Can learn sets of rules such as
 - $Parent(x,y) \rightarrow Ancestor(x,y)$
 - $Parent(x,z) \text{ and } Ancestor(z,y) \rightarrow Ancestor(x,y)$
- Research field of **Inductive Logic Programming**.

Why ILP? (slide of S.Matwin)

- **expressiveness of logic as representation** (Quinlan)



- can't represent this graph as a fixed length vector of attributes
- can't represent a "transition" rule:

A can-reach B if A link C, and C can-reach B
without variables

FINITE ELEMENT MESH DESIGN

Given a geometric **structure** and **loadings/boundary conditions**
Find an **appropriate resolution** for a finite element mesh

Examples: ten structures with appropriate meshes (cca. 650 edges)

Background knowledge

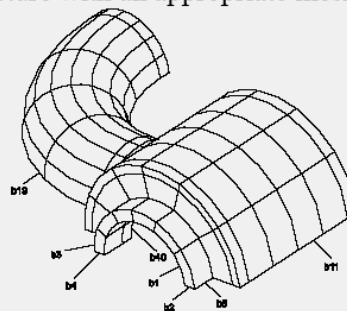
- Properties of edges (short, loaded, two-side-fixed, ...)
- Relations between edges (neighbor, opposite, equal)

ILP systems applied: GOLEM, CLAUDIEN

Many interesting rules discovered (according to expert evaluation)

Finite element mesh design (ctd.)

Example structure with an appropriate mesh



Example rules

```
mesh(Edge, 7) ← usual_length(Edge),  
neighbour_xy(Edge, EdgeY), two_side_fixed(EdgeY),  
neighbour_zx(EdgeZ, Edge), not_loaded(EdgeZ)  
mesh(Edge, N) ← equal(Edge, Edge2), mesh(Edge2, N)
```

Application areas

- Medicine
- Economy, Finance
- Environmental cases
- Engineering
 - Control engineering and robotics
 - Technical diagnostics
 - Signal processing and image analysis
- Information sciences
- Social Sciences
- Molecular Biology
- Chemistry and Pharmacy
- ...

Where to find more?

- T. Mitchell *Machine Learning* New York: McGraw-Hill, 1997.
- I. H. Witten & Eibe Frank *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations* San Francisco: Morgan Kaufmann, 1999.
- Michalski R.S., Bratko I., Kubat M. *Machine learning and data mining*; J. Wiley. 1998.
- Clark, P., & Niblett, T. (1989). The CN2 induction algorithm. *Machine Learning*, 3, 261–283.
- Cohen W. Fast effective rule induction. *Proc. of the 12th Int. Conf. on Machine Learning 1995*. 115–123
- R.S. Michalski, I. Mozetic, J. Hong and N. Lavrac, The multi-purpose incremental learning system AQ15 and its testing application to three medical domains, *Proceedings of i4AAI 1986, 1041-1045, (1986)*.
- J.W. Grzymala-Busse, LERS-A system for learning from example-s based on rough sets, In *Intelligent Decision Support: Handbook of Applications and Advances of Rough Sets Theory*, (Edited by R.Slowinski), pp. 3-18
- Michalski R.S.: A theory and methodology of inductive learning. W Michalski R.S, Carbonell J.G., Mitchell T.M. (red.) *Machine learning: An Artificial Intelligence Approach*, Morgan Kaufmann Publishers, Los Altos (1983),.
- J.Stefanowski: On rough set based approaches to induction of decision rules, w: A. Skowron, L. Polkowski (red.), *Rough Sets in Knowledge Discovery Vol 1*, Physica Verlag, Heidelberg, 1998, 500-529.
- J.Stefanowski, The rough set based rule induction technique for classification problems, w: *Proceedings of 6th European Conference on Intelligent Techniques and Soft Computing, Aachen, EUFIT 98, 1998, 109-113*.
- J. Furnkranz . Separate-and-conquer rule learning. *Artificial Intelligence Review*, 13(1):3–54, 1999.

Where to find more - 2

- P. Clark and R. Boswell. Rule induction with CN2: Some recent improvements. In *Proceedings of the 5th European Working Session on Learning (EWSL-91)*, pp. 151–163, 1991.
- Grzymala-Busse J.W.: Managing uncertainty in machine learning from examples. *Proceedings of 3rd Int. Symp. on Intelligent Systems*, Wigry 1994 .
- Cendrowska J.: PRISM, an algorithm for inducing modular rules. *Int. J. Man-Machine Studies*, 27 (1987), 349-370.
- Frank, E., & Witten, I. H. (1998). Generating accurate rule sets without global optimization. *Proc. of the 15th Int. Conf. on Machine Learning (ICML-98)* (pp. 144–151).
- J. Furnkranz and P. Flach. An analysis of rule evaluation metrics. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pp. 202–209,
- S. M. Weiss and N. Indurkha. Lightweight rule induction. In *Proc. of the 17th Int. Conference on Machine Learning (ICML-2000)*, pp. 1135–1142,
- J.Stefanowski, D.Vanderpooten: Induction of decision rules in classification and discovery-oriented perspectives, *International Journal of Intelligent Systems*, vol. 16 no. 1, 2001, 13-28.
- J.W.Grzymala-Busse, J.Stefanowski: Three approaches to numerical attribute discretization for rule induction, *International Journal of Intelligent Systems*, vol. 16 no. 1, 2001, 29-38.
- P. Domingos. Unifying instance-based and rule-based induction. *Machine Learning*, 24:141–168, 1996.
- R. Holte. Very simple classification rules perform well on most commonly used datasets. *Machine Learning*, 11:63–91, 1993.

Any questions, remarks?

