

Output: knowledge representation

Most of these slides (used with permission) are based on the book:

Data Mining: Practical Machine Learning Tools and Techniques
by I. H. Witten, E. Frank, M. A. Hall, and C. J. Pal

1

Output: Knowledge representation

- Tables
- Linear models
- Trees
- Rules
- Classification rules
- Association rules
- Rules with exceptions
- More expressive rules
- Instance-based representation
- Clusters

2

2

Output: representing structural patterns

- Many different ways of representing patterns
 - Decision trees, rules, instance-based, ...
- Also called “knowledge” representation
- Representation determines inference method
- Understanding the output is the key to understanding the underlying learning methods
- Different types of output for different learning problems (e.g., classification, regression, ...)

3

3

Decision tables

- Simplest way of representing output:
 - Use the format that is used for representing the input!
- Decision table for the weather problem:

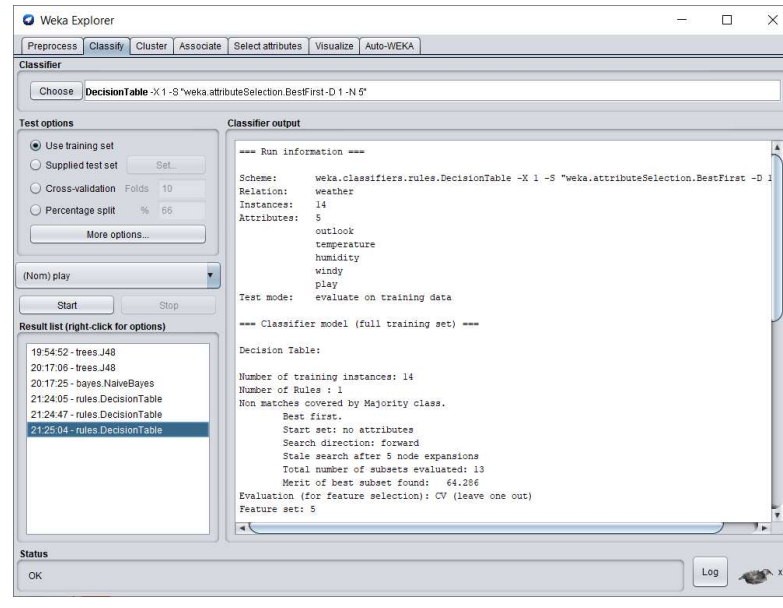
Outlook	Humidity	Play
Sunny	High	No
Sunny	Normal	Yes
Overcast	High	Yes
Overcast	Normal	Yes
Rainy	High	No
Rainy	Normal	No

- Main problem: selecting the right attributes

4

4

Decision table classifier – uses majority class



5

Linear models

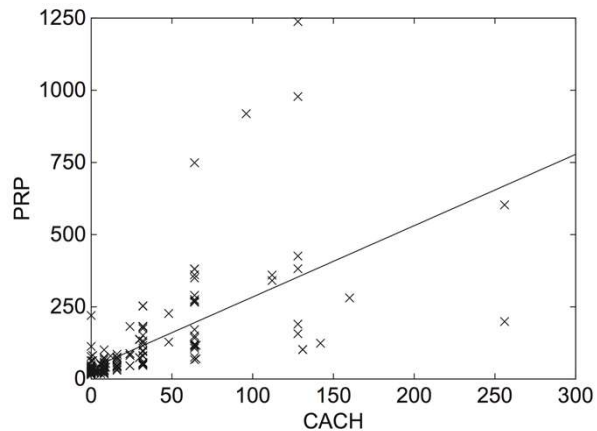
- Another simple representation
- Traditionally primarily used for regression:
 - Inputs (attribute values) and output are all numeric
- Output is the sum of the weighted input attribute values
- The trick is to find good values for the weights
- There are different ways of doing this, which we will consider later; the most famous one is to minimize the squared error

6

6

A linear regression function for the CPU performance data

PRP: Performance CACH: Cache



$$\text{PRP} = 37.06 + 2.47\text{CACH}$$

7

7

Linear models for classification

- Binary classification (Logistic regression)
- Line *separates* the two classes
 - Decision boundary - defines where the decision changes from one class value to the other
- Prediction is made by plugging in observed values of the attributes into the expression
 - Predict one class if output ≥ 0 , and the other class if output < 0
- Boundary becomes a high-dimensional plane (*hyperplane*) when there are multiple attributes

8

8

Logistic regression for classification

Classifier
Choose: SimpleLogistic -I 0 -M 500 -H 50 -W 0.0

Test options
 Use training set
 Supplied test set
 Cross-validation Folds: 10
 Percentage split %: 66
 More options...
 (Nom) play
 Start Stop

Classifier output

```

=== Run information ===
Scheme: weka.classifiers.functions.SimpleLogistic -I 0 -M 500 -H 50 -W 0.0
Relation: weather.symbolic
Instances: 14
Attributes: 5
  outlook
  temperature
  humidity
  windy
  play
Test mode: evaluate on training data

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

SimpleLogistic:
Class yes :
-1.9 +
[outlook=overcast] * 1.85 +
[temperature=mild] * 0.75 +
[humidity=normal] * 1.65 +
[windy=FALSE] * 1.26

Class no :
1.9 +
[outlook=overcast] * -1.85 +
[temperature=mild] * -0.75 +
[humidity=normal] * -1.65 +
[windy=FALSE] * -1.26

Time taken to build model: 0.01 seconds
  
```

Result list (right-click for options)

- 22:43:49 - rules DecisionTable
- 22:51:31 - functions LinearRegression
- 22:51:50 - functions LinearRegression
- 23:00:11 - functions SimpleLogistic
- 23:01:09 - functions SimpleLogistic
- 23:01:52 - functions SimpleLogistic

Status
OK Log

9

Separating setosas from versicolors - Preprocessing

Filter
Choose: RemoveWithValues -S 0.0 -C last-L 3

Current relation
Relation: iris-weka.filters.unsupervised.instance.RemoveWithValues...
Instances: 100
Attributes: 5
Sum of weights: 100

Attributes
 All None Invert Pattern
 1 sepalength
 2 sepalwidth
 3 petalength
 4 petalwidth
 5 class

Selected attribute

No.	Label	Count	Weight
1	Iris-setosa	50	50.0
2	Iris-versicolor	50	50.0
3	Iris-virginica	0	0.0

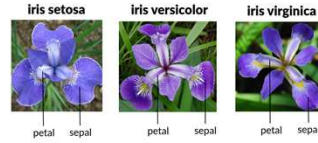
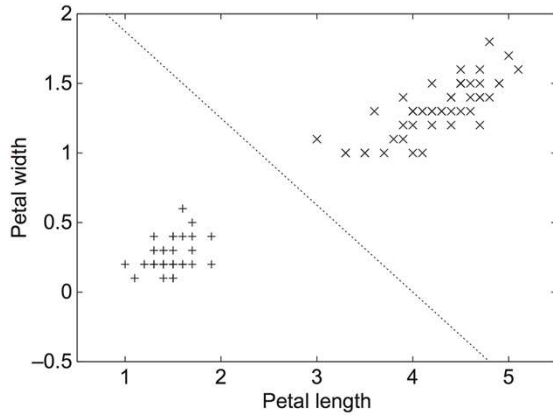
Class: class (Nom) Visualize All

Status
OK Log

10

10

Separating setosas from versicolors



$$2.0 - 0.5\text{PETAL-LENGTH} - 0.8\text{PETAL-WIDTH} = 0$$

11

11

Separating setosas from versicolors – iris2class.arff

The screenshot shows the Weka Explorer interface. The classifier selected is 'SimpleLogistic-10-M 500-H 50-W 0.0'. The 'Classifier output' window displays the following summary:

```

Time taken to test model on training data: 0 seconds

--- Summary ---
Correctly Classified Instances 100      100 %
Incorrectly Classified Instances 0      0 %
Kappa statistic 1
Mean absolute error 0.1415
Root mean squared error 0.1555
Relative absolute error 28.3078 %
Root relative squared error 31.1007 %
Total Number of Instances 100

--- Detailed Accuracy By Class ---

```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	FRC Area	Cls
Weighted Avg.	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	Iris

The 'Confusion Matrix' section shows:

```

--- Confusion Matrix ---
a b <-- Classified as
50 0 | a = Iris-setosa
0 50 | b = Iris-versicolor

```

The 'Result list' on the left shows various classifiers, with '21:39:42 - functions.SimpleLogistic' selected.

12

12

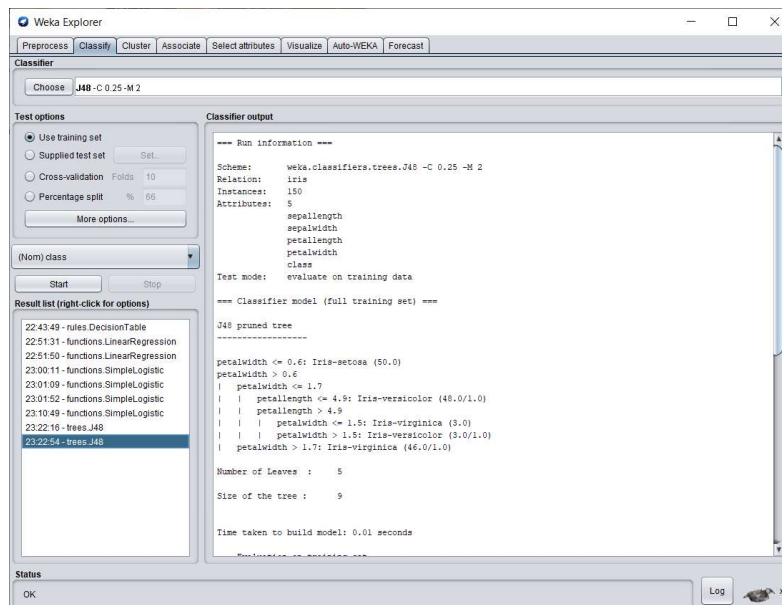
Decision trees

- “Divide-and-conquer” approach produces tree
- Nodes involve testing a particular attribute
- Usually, attribute value is compared to constant
- Leaves assign classification, set of classifications, or probability distribution to instances
- Unknown instance is routed down the tree

13

13

Decision trees – J48 algorithm



The screenshot shows the Weka Explorer interface with the J48 classifier selected. The classifier output window displays the following information:

```
--- Run Information ---
Scheme: weka.classifiers.trees.J48 -C 0.25 -M 2
Relation: iris
Instances: 150
Attributes: 5
  sepalwidth
  sepalwidth
  petalwidth
  petalwidth
  class
Test mode: evaluate on training data

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

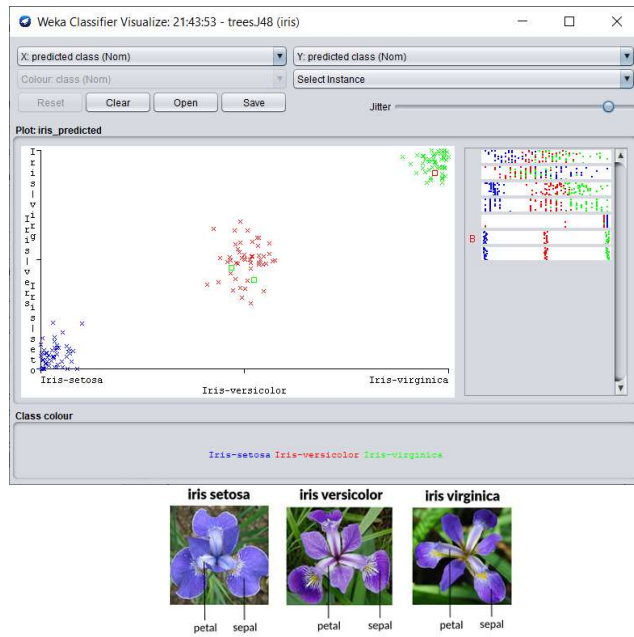
J48 pruned tree
-----
petalwidth <= 0.6: Iris-setosa (50.0)
petalwidth > 0.6
 | petalwidth <= 1.7
 | | petalwidth <= 4.9: Iris-versicolor (48.0/1.0)
 | | | petalwidth > 4.9
 | | | | petalwidth <= 1.5: Iris-virginica (3.0)
 | | | | petalwidth > 1.5: Iris-versicolor (3.0/1.0)
 | | | petalwidth > 1.7: Iris-virginica (46.0/1.0)

Number of Leaves : 5
Size of the tree : 9
Time taken to build model: 0.01 seconds
```

14

14

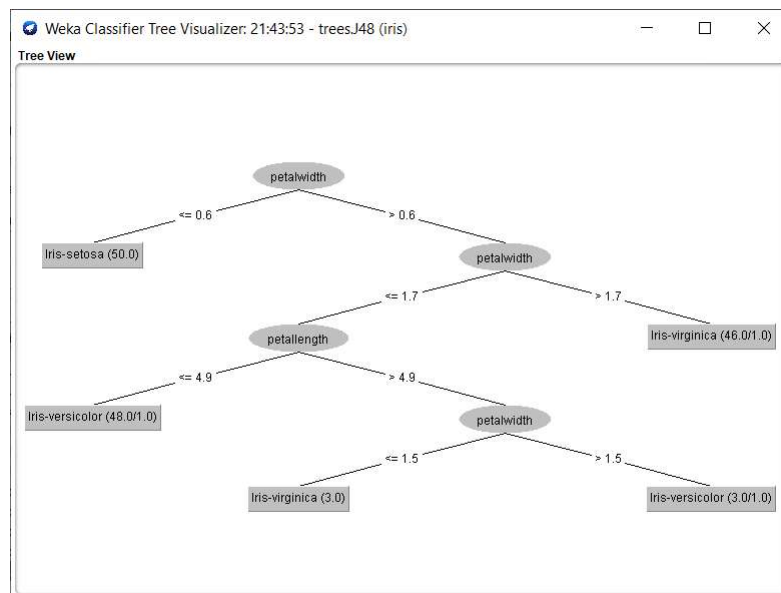
Interactive tree construction I



15

15

Interactive tree construction II



16

16

Nominal and numeric attributes in trees

- **Nominal:**
number of children usually equal to number values
⇒ attribute won't get tested more than once
- Other possibility: division into two subsets
- **Numeric:**
test whether value is greater or less than constant
⇒ attribute may get tested several times
 - Other possibility: three-way split (or multi-way split)
 - Integer: *less than, equal to, greater than*
 - Real: *below, within, above*

17

17

Missing values

- Does absence of value have some significance?
- Yes ⇒ "missing" is a separate value
- No ⇒ "missing" must be treated in a special way
 - Solution: assign instance to most popular branch
 - Solution: Use Weka filter for filling in missing values (See Weka presentation)

18

18

Predicting CPU performance

- Example: 209 different computer configurations

	Cycle time (ns)	Main memory (Kb)		Cache (Kb)	Channels		Performance
	MYCT	MMIN	MMAx	CACH	CHMIN	CHMAX	PRP
1	125	256	6000	256	16	128	198
2	29	8000	32000	32	8	32	269
...							
208	480	512	8000	32	0	0	67
209	480	1000	4000	0	0	0	45

- Linear regression function

$$PRP = -56.1 + 0.049 MYCT + 0.015 MMIN + 0.006 MMAx + 0.630 CACH - 0.270 CHMIN + 1.46 CHMAX$$

19

19

Predicting CPU performance

The screenshot shows the Weka Explorer interface with the Linear Regression classifier selected. The classifier output window displays the following information:

```

Scheme: weka.classifiers.functions.LinearRegression -S 0 -R 1.0E-8 -num-decimal-places 4
Relation: cpu
Instances: 209
Attributes: 7
  MYCT
  MMIN
  MMAx
  CACH
  CHMIN
  CHMAX
  class
Test mode: evaluate on training data

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

Linear Regression Model

class =
  0.0491 * MYCT +
  0.0152 * MMIN +
  0.0056 * MMAx +
  0.6298 * CACH +
  1.4599 * CHMAX +
  -56.075

Time taken to build model: 0.07 seconds

=== Evaluation on training set ===
    
```

The result list on the left shows the selected configuration: 21:48:29 - functions.LinearRegression.

20

20

Trees for numeric prediction

- *Regression*: the process of computing an expression that predicts a numeric quantity
- *Regression tree*: “decision tree” where each leaf predicts a numeric quantity
 - Predicted value is average value of training instances that reach the leaf
- *Model tree*: “regression tree” with linear regression models at the leaf nodes

21

21

Regression tree for the CPU data – M5P Algo.

The screenshot shows the Weka Explorer interface with the M5P classifier selected. The Classifier output pane displays the following information:

```
=== Run information ===
Scheme:   weka.classifiers.trees.M5P -M 4.0
Relation: cpu
Instances: 209
Attributes: 7
  MICT
  MMIN
  MMAX
  CACH
  CEMIN
  CEMAX
  class

Test mode: evaluate on training data

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

M5 pruned model tree:
(using smoothed linear models)

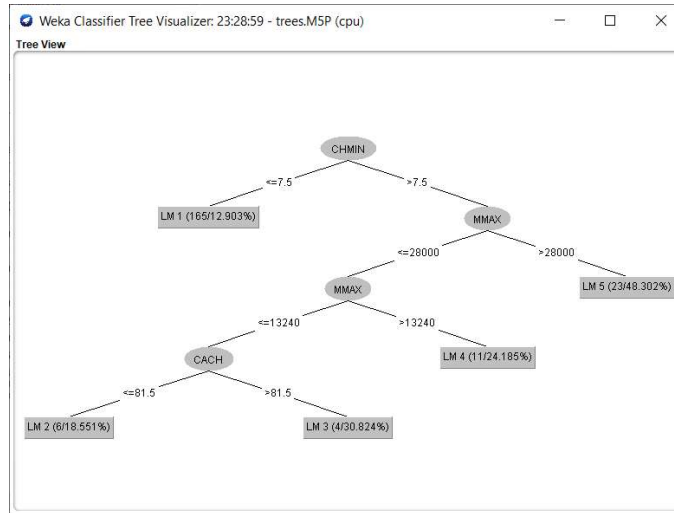
CEMIN <= 7.5 : IM1 (165/12.9038)
CEMIN > 7.5 :
|  | MMAX <= 28000 :
|  | | MMAX <= 13240 :
|  | | | CACH <= 81.5 : IM2 (6/18.5518)
|  | | | CACH > 81.5 : IM3 (4/30.8248)
|  | | MMAX > 13240 : IM4 (12/24.1894)
|  | MMAX > 28000 : IM5 (23/48.3024)

IM num: 1
class =
-0.0055 * MICT
+ 0.0013 * MMIN
```

22

22

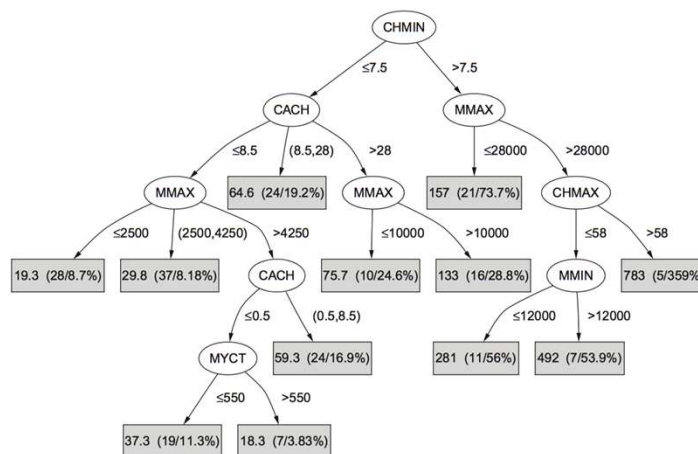
Regression tree for the CPU data – M5P Algo.



23

23

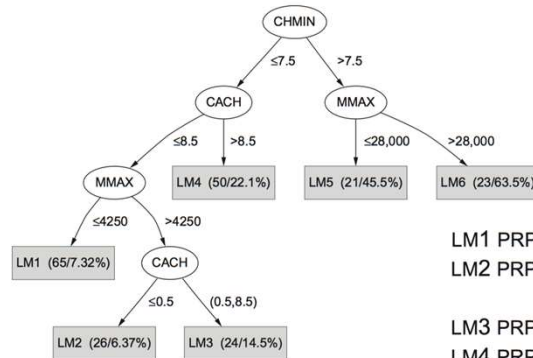
Regression tree for the CPU data



24

24

Model tree for the CPU data



$$\begin{aligned} \text{LM1 PRP} &= 8.29 + 0.004 \text{ MMAX} + 2.77 \text{ CHMIN} \\ \text{LM2 PRP} &= 20.3 + 0.004 \text{ MMIN} - 3.99 \text{ CHMIN} \\ &\quad + 0.946 \text{ CHMAX} \\ \text{LM3 PRP} &= 38.1 + 0.012 \text{ MMIN} \\ \text{LM4 PRP} &= 19.5 + 0.002 \text{ MMAX} + 0.698 \text{ CACH} \\ &\quad + 0.969 \text{ CHMAX} \\ \text{LM5 PRP} &= 285 + 1.46 \text{ MYCT} + 1.02 \text{ CACH} \\ &\quad - 9.39 \text{ CHMIN} \\ \text{LM6 PRP} &= -65.8 + 0.03 \text{ MMIN} - 2.94 \text{ CHMIN} \\ &\quad + 4.98 \text{ CHMAX} \end{aligned}$$

25

25

Classification rules

- Popular alternative to decision trees
- *Antecedent* (pre-condition): a series of tests (just like the tests at the nodes of a decision tree)
- Tests are usually logically ANDed together (but may also be general logical expressions)
- *Consequent* (conclusion): classes, set of classes, or probability distribution assigned by rule
- Individual rules are often logically ORed together
 - Conflicts arise if different conclusions apply

26

26

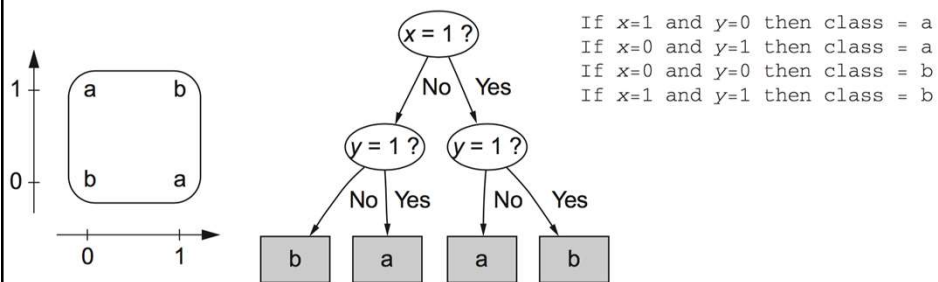
From trees to rules

- Easy: converting a tree into a set of rules
 - One rule for each leaf:
 - Antecedent contains a condition for every node on the path from the root to the leaf
 - Consequent is class assigned by the leaf
 - Produces rules that are unambiguous
 - Doesn't matter in which order they are executed
 - But: resulting rules are unnecessarily complex
 - Pruning to remove redundant tests/rules

27

27

From rules to trees – not straightforward The exclusive-or problem



- Even if the rule involves two attributes, split on one attribute first to get a sub-tree

28

28

Interpreting rules

- What if two or more rules conflict?
 - Give no conclusion at all?
 - Go with rule that is most popular on training data?
 - ...
- What if no rule applies to a test instance?
 - Give no conclusion at all?
 - Go with class that is most frequent in training data?
 - ...

29

29

Association rules

- Association rules...
 - ... can predict any attribute and combinations of attributes
 - ... are not intended to be used together as a set
- Problem: immense number of possible associations
 - Output needs to be restricted to show only the most predictive associations
 - ⇒ only those with high *support* and high *confidence*

30

30

Support and confidence of a rule

- Support: number of instances predicted correctly
- Confidence: number of correct predictions, as proportion of all instances that rule applies to
- Example: 4 cool days with normal humidity

```
If temperature = cool then
humidity = normal
```

⇒ Support = 4, confidence = 100%

- Normally: minimum support and confidence pre-specified (e.g. 58 rules with support ≥ 2 and confidence $\geq 95\%$ for weather data)

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

31

31

Interpreting association rules

- Interpretation is not obvious:

```
If windy = false and play = no then outlook = sunny
and humidity = high
```

is *not* the same as

```
If windy = false and play = no then outlook = sunny
If windy = false and play = no then humidity = high
```

- It means that the following also holds:

```
If humidity = high and windy = false and play = no
then outlook = sunny
```

32

32

Association rules – weather.arff

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize Auto-WEKA

Choose Apriori-N10-T0-C0.9-D0.05-U1.0-M0.1-S1.0-c-1

Start Stop

Associator output

Result list (right-click...)

22:15:27 - Apriori
22:16:52 - Apriori

Minimum support: 0.15 (2 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 17

Generated sets of large itemsets:

Size of set of large itemsets L(1): 12
Size of set of large itemsets L(2): 47
Size of set of large itemsets L(3): 39
Size of set of large itemsets L(4): 6

Best rules found:

1. outlook=overcast 4 ==> play=yes 4 <conf:(1)> lift:(1.56) lev:(0.1) [1] conv:(1.43)
2. temperature=cool 4 ==> humidity=normal 4 <conf:(1)> lift:(2) lev:(0.14) [2] conv:(2)
3. humidity=normal windy=FALSE 4 ==> play=yes 4 <conf:(1)> lift:(1.56) lev:(0.1) [1] conv:(1.43)
4. outlook=sunny play=no 3 ==> humidity=high 3 <conf:(1)> lift:(2) lev:(0.11) [1] conv:(1.5)
5. outlook=sunny humidity=high 3 ==> play=no 3 <conf:(1)> lift:(2.8) lev:(0.14) [1] conv:(1.93)
6. outlook=rainy play=yes 3 ==> windy=FALSE 3 <conf:(1)> lift:(1.78) lev:(0.09) [1] conv:(1.29)
7. outlook=rainy windy=FALSE 3 ==> play=yes 3 <conf:(1)> lift:(1.56) lev:(0.08) [1] conv:(1.07)
8. temperature=cool play=yes 3 ==> humidity=normal 3 <conf:(1)> lift:(2) lev:(0.11) [1] conv:(1.5)
9. outlook=sunny temperature=hot 2 ==> humidity=high 2 <conf:(1)> lift:(2) lev:(0.07) [1] conv:(1)
10. temperature=hot play=no 2 ==> outlook=sunny 2 <conf:(1)> lift:(2.8) lev:(0.09) [1] conv:(1.29)

Status
OK Log X0

33

Rules with exceptions

- Idea: allow rules to have *exceptions*
- Example: rule for iris data

If petal-length ≥ 2.45 and petal-length < 4.45 then Iris-versicolor

- New instance:

Sepal Length	Sepal Width	Petal Length	Petal Width	Type
5.1	3.5	2.6	0.2	?

- Modified rule:

If petal-length ≥ 2.45 and petal-length < 4.45 then Iris-versicolor
EXCEPT if petal-width < 1.0 then Iris-setosa

34

34

More on exceptions

- `Default...except if...then...`
is logically equivalent to
`if...then...else`
(where the “else” specifies what the “default” does)
- But: exceptions offer a psychological advantage
 - Assumption: defaults and tests early on apply more widely than exceptions further down
 - **Exceptions reflect special cases**
- Advantages of rules
 - Rules can be updated incrementally
 - Easy to incorporate new data
 - Easy to incorporate domain knowledge
 - People often think in terms of exceptions
 - Each conclusion can be considered just in the context of rules and exceptions that lead to it

35

35

Rules involving relations

- So far: all rules involved comparing an attribute-value to a constant (e.g. `temperature < 45`)
- These rules are called “propositional” because they have the same expressive power as propositional logic
- What if problem involves relationships between examples (e.g. family tree problem from above)?
 - Can’t be expressed with propositional rules
 - More expressive representation required

Propositional logic example

P: a grizzly is a bear.
Q: a bear is a mammal.
R: a grizzly is a mammal.

$P \ \& \ Q \ \rightarrow \ R$

Each proposition symbol represents an atomic statement.

Atomic statements can be combined with logical connectives to form compound statements.

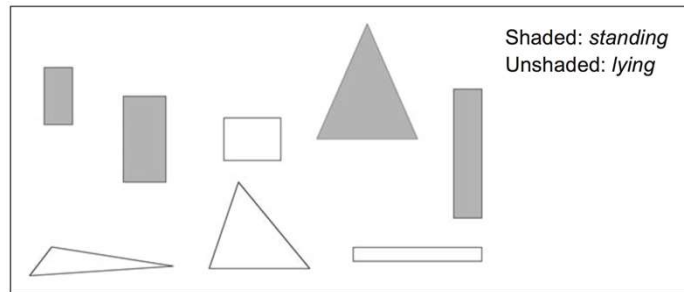
$\sim P, P \ \& \ Q, P \ \vee \ Q, P = Q, P \ \rightarrow \ Q$

36

36

The shapes problem

- Target concept: *standing up*
- Shaded: *standing*
Unshaded: *lying*



37

37

A propositional solution

Width	Height	Sides	Class
2	4	4	Standing
3	6	4	Standing
4	3	4	Lying
7	8	3	Standing
7	6	3	Lying
2	9	4	Standing
9	1	4	Lying
10	2	3	Lying

```
If width  $\geq$  3.5 and height  $<$  7.0  
then lying  
If height  $\geq$  3.5 then standing
```

38

38

Using relations between attributes

- Comparing attributes with each other enables rules like this:

```
If width > height then lying  
If height > width then standing
```

- This description generalizes better to new data
- Standard relations: =, <, >
- But: searching for relations between attributes can be costly
- Simple solution: add extra attributes
(e.g., a binary attribute “*is width < height?*”)

39

39

Instance-based representation

- Simplest form of learning: *rote learning*
 - Training instances are searched for instance that most closely resembles new instance
 - The instances themselves represent the knowledge
 - Also called *instance-based learning*
- Similarity function defines what’s “learned”
- Instance-based learning is *lazy learning*
- Methods: *nearest-neighbor*, *k-nearest-neighbor*, ...

40

40

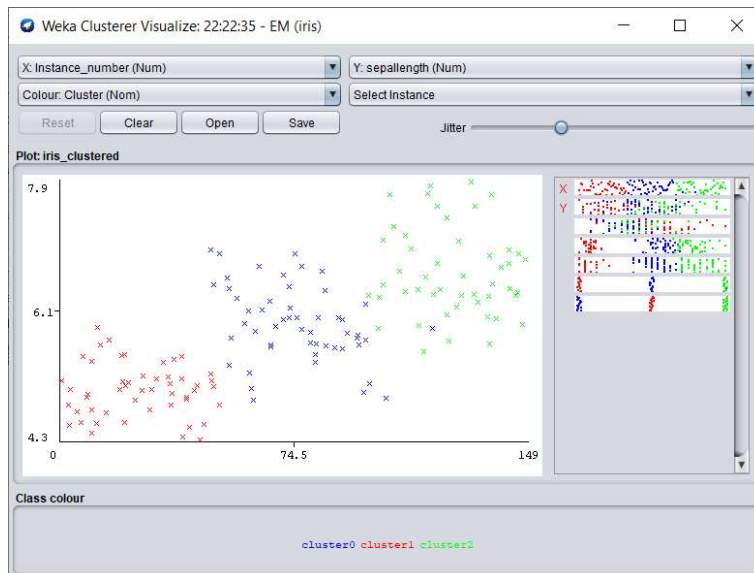
The distance function

- Simplest case: one numeric attribute
 - Distance is the difference between the two attribute values involved (or a function thereof)
- Several numeric attributes: normally, Euclidean distance is used and attributes are normalized
- Nominal attributes: distance is set to 1 if values are different, 0 if they are equal
- Are all attributes equally important?
 - Weighting the attributes might be necessary

41

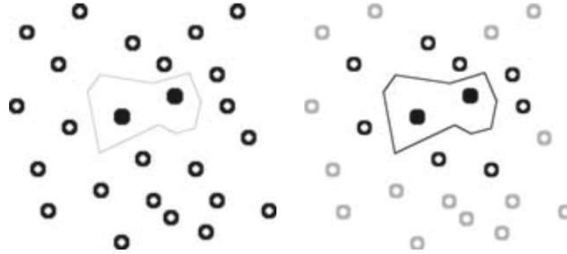
41

IRIS data clusters



42

Learning prototypes

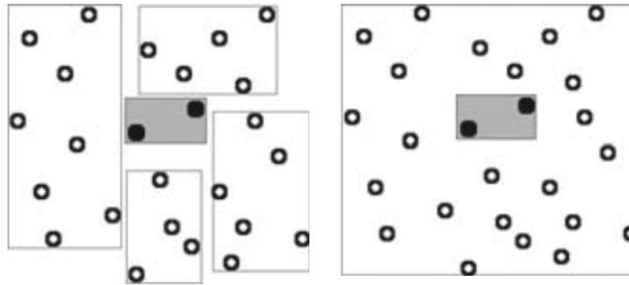


- Only those instances involved in a decision need to be stored
- Noisy instances should be filtered out

43

43

Rectangular generalizations



- Nearest-neighbor rule is used outside rectangles
- Rectangles are rules! (But they can be more conservative than “normal” rules.)
- Nested rectangles are rules with exceptions

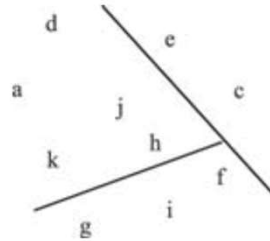
44

44

Representing clusters I

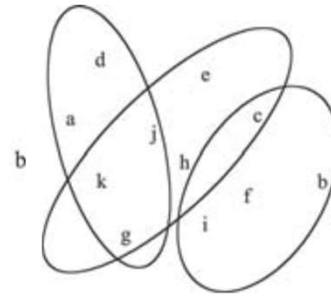
Simple 2-D representation

One cluster per example



Venn diagram

Multiple clusters per example



45

45

Representing clusters II

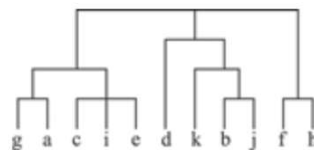
Probabilistic assignment

Probability of belonging to each cluster

	1	2	3
a	0.4	0.1	0.5
b	0.1	0.8	0.1
c	0.3	0.3	0.4
d	0.1	0.1	0.8
e	0.4	0.2	0.4
f	0.1	0.4	0.5
g	0.7	0.2	0.1
h	0.5	0.4	0.1
...			

Dendrogram

Hierarchical clusters



46

46