

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:

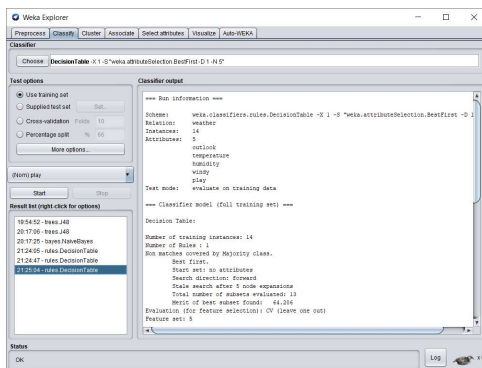
Outlook	Humidity	Play
Sunny	High	No
Sunny	Normal	Yes
Overcast	High	Yes
Overcast	Normal	Yes
Rainy	High	No
Rainy	Normal	No

- Main problem: selecting the right attributes

4

4

Decision table classifier – uses majority class



5

Linear models

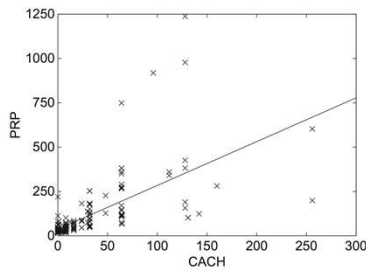
- Another simple representation
- Traditionally primarily used for regression:
 - Inputs (attribute values) and output are all numeric
- Output is the sum of the weighted input attribute values
- The trick is to find good values for the weights
- There are different ways of doing this, which we will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



$$PRP = 37.06 + 2.47CACH$$

7

7

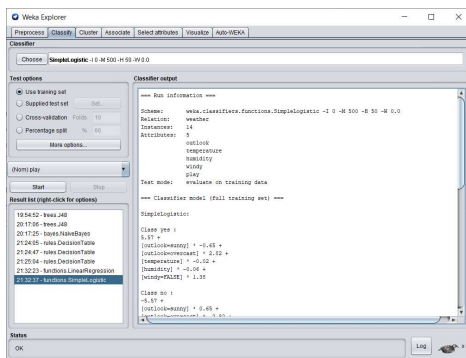
Linear models for classification

- 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

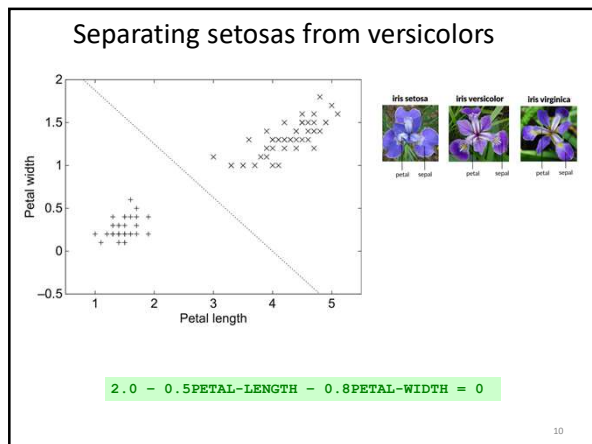
8

8

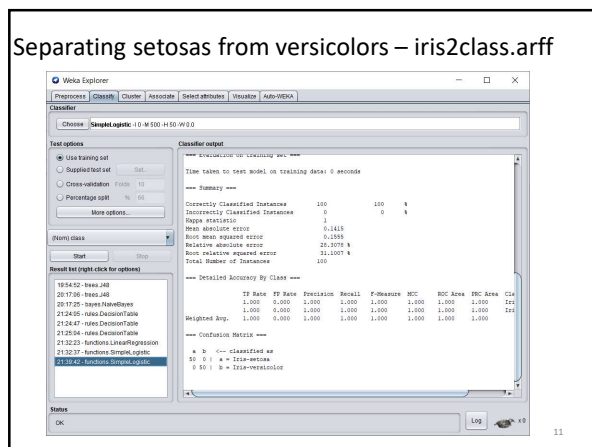
Logistic regression for classification



9



10

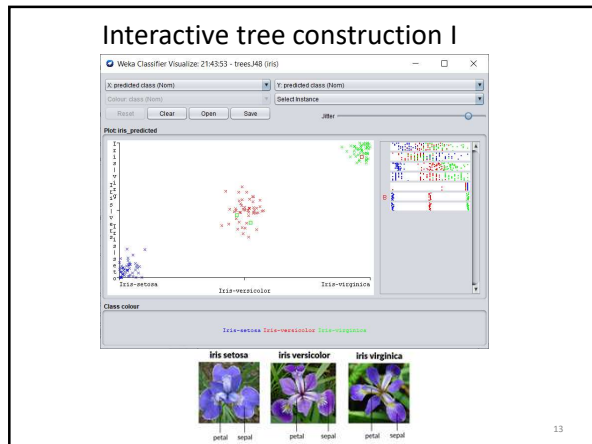


11

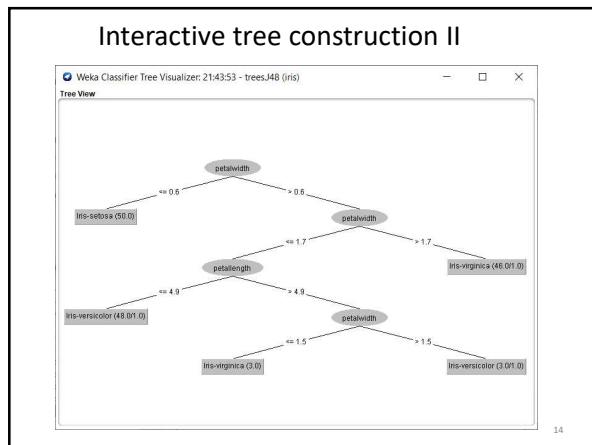
Decision trees

- “Divide-and-conquer” approach produces tree
- Nodes involve testing a particular attribute
- Usually, attribute value is compared to constant
- Other possibilities:
 - Comparing values of two attributes
 - Using a function of one or more attributes
- Leaves assign classification, set of classifications, or probability distribution to instances
- Unknown instance is routed down the tree

12



13



14

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*

15

Missing values

- Does absence of value have some significance?
- Yes \Rightarrow "missing" is a separate value
- No \Rightarrow "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)

16

16

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

17

17

Predicting CPU performance

- Example: 209 different computer configurations

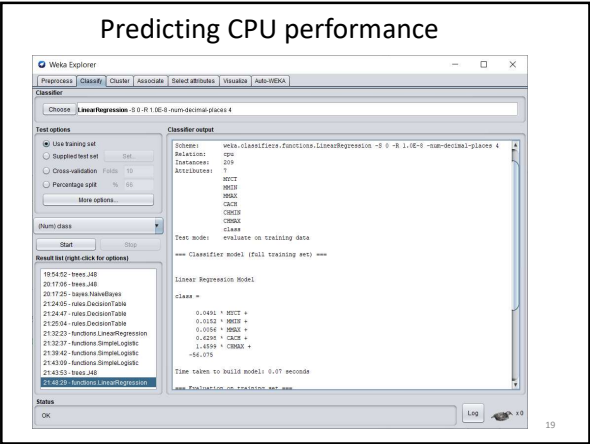
	Cycle time (ns)		Main memory (Kb)		Cache (Kb)	Channels		Performance
	MYCT	MMIN	MMAX	CACH	CHMIN	CHMAX	PRP	
1	125	256	6000	256	16	128	198	
2	29	8000	32000	32	8	32	269	
...								
208	480	512	8000	32	0	0	67	
209	480	1000	4000	0	0	0	45	

- Linear regression function

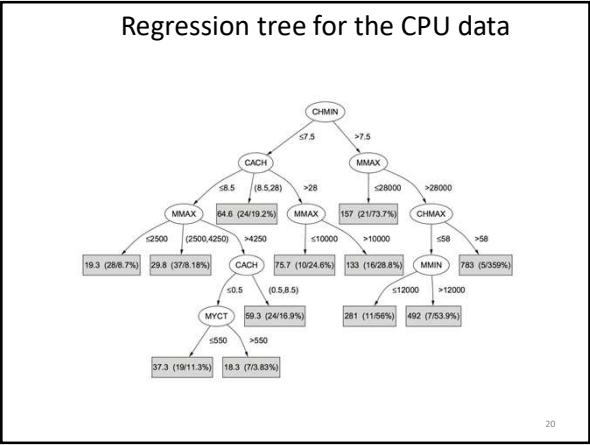
$$\text{PRP} = -56.1 + 0.049 \text{ MYCT} + 0.015 \text{ MMIN} + 0.006 \text{ MMAX} + 0.630 \text{ CACH} - 0.270 \text{ CHMIN} + 1.46 \text{ CHMAX}$$

18

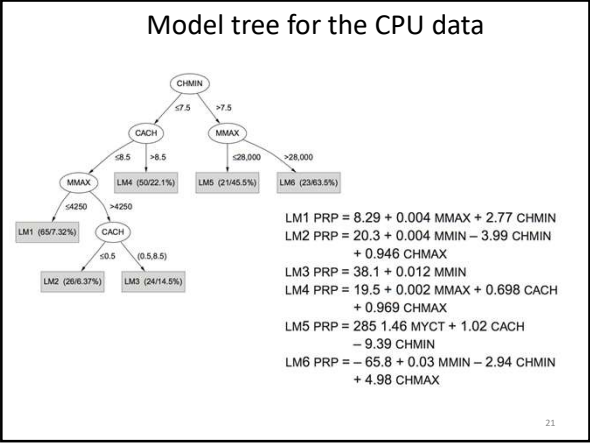
18



19



20



21

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

22

22

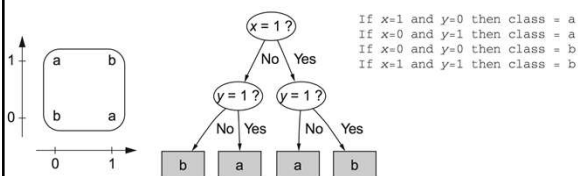
From trees to rules

- 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

23

23

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

24

24

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?
 - ...

25

25

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*

26

26

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

27

27

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

28

28

Association rules – weather.arff

29

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

30

30

Advantages of using exceptions

- 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
 - Locality property is important for understanding large rule sets
 - “Normal” rule sets do not offer this advantage

31

31

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

32

32

Rules involving relations

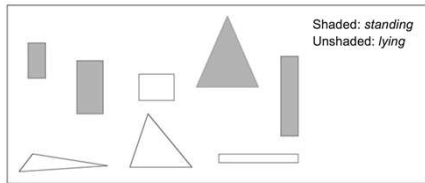
- So far: all rules involved comparing an attribute-value to a constant (e.g. temperature < 45)
- These rules are called “propositional” because they have the same expressive power as propositional logic
- What if problem involves relationships between examples (e.g. family tree problem from above)?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

33

33

The shapes problem

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



34

34

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 ≥ 3.5 and height < 7.0
then lying
If height ≥ 3.5 then standing
```

35

35

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?*")

36

36

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*, ...

37

37

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

38

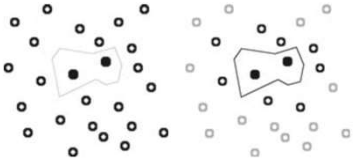
38

IRIS data clusters



39

Learning prototypes

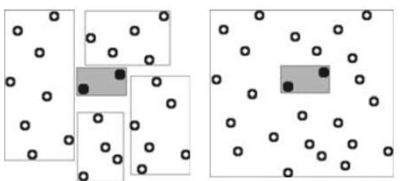


- Only those instances involved in a decision need to be stored
- Noisy instances should be filtered out

40

40

Rectangular generalizations

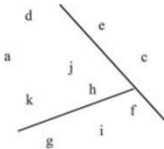
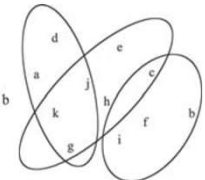


- Nearest-neighbor rule is used outside rectangles
- Rectangles are rules! (But they can be more conservative than "normal" rules.)
- Nested rectangles are rules with exceptions

41

41

Representing clusters I

<i>Simple 2-D representation</i>	<i>Venn diagram</i>
One cluster per example	Multiple clusters per example
	

42

42

Representing clusters II

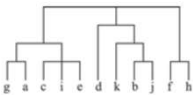
Probabilistic assignment

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



43
