

# Decision Trees

From ID3 to C4.5

Lecture 3 – Part 3

Marina Santini

Acknowledgements

Slides borrowed and adapted from:  
*Data Mining* by I. H. Witten, E. Frank and M. A. Hall

2016

Decision Trees - Part 3

1

## Implementation:

Real machine learning schemes

- Decision trees
  - ♦ From ID3 to C4.5 (pruning, numeric attributes, ...)

2016

Decision Trees - Part 3

2

## Decision trees

- Extending ID3:
  - to permit numeric attributes: *straightforward*
  - to deal sensibly with missing values: *trickier*
  - stability for noisy data: *requires pruning mechanism*
- End result: C4.5 (Quinlan)
  - Best-known and (probably) most widely-used learning algorithm
  - Commercial successor: C5.0

2016

Decision Trees - Part 3

3

## Numeric attributes

- Standard method: binary splits
  - ♦ E.g. temp < 45
- Unlike nominal attributes, every attribute has many possible split points
- Solution is straightforward extension:
  - ♦ Evaluate info gain (or other measure) for every possible split point of attribute
  - ♦ Choose “best” split point
  - ♦ Info gain for best split point is info gain for attribute
- Computationally more demanding

2016

Decision Trees - Part 3

4

## Missing values

- Split instances with missing values into pieces
  - ♦ A piece going down a branch receives a weight proportional to the popularity of the branch
  - ♦ weights sum to 1
- Info gain works with fractional instances
  - ♦ use sums of weights instead of counts
- During classification, split the instance into pieces in the same way
  - ♦ Merge probability distribution using weights

2016

Decision Trees - Part 3

5

## Pruning

- Prevent overfitting to noise in the data
- “Prune” the decision tree
- Two strategies:
  - *Postpruning*  
take a fully-grown decision tree and discard unreliable parts
  - *Prepruning*  
stop growing a branch when information becomes unreliable
- Postpruning preferred in practice—prepruning can “stop early”

2016

Decision Trees - Part 3

6

### Prepruning

- Based on statistical significance test
  - Stop growing the tree when there is no *statistically significant* association between any attribute and the class at a particular node
- Most popular test: *chi-squared test*
- ID3 used chi-squared test in addition to information gain
  - Only statistically significant attributes were allowed to be selected by information gain procedure

2016 Decision Trees - Part 3 7

### Early stopping

	a	b	class
1	0	0	0
2	0	1	1
3	1	0	1
4	1	1	0

- Pre-pruning may stop the growth process prematurely: *early stopping*
- Classic example: XOR/Parity-problem
  - No *individual* attribute exhibits any significant association to the class
  - Structure is only visible in fully expanded tree
  - Prepruning won't expand the root node
- But: XOR-type problems rare in practice
- And: prepruning faster than postpruning

2016 Decision Trees - Part 3 8

### Postpruning

- First, build full tree
- Then, prune it
- Fully-grown tree shows all attribute interactions
- Problem: some subtrees might be due to chance effects
- Two pruning operations:
  - Subtree replacement*
  - Subtree raising*
- Possible strategies:
  - error estimation
  - significance testing
  - MDL principle

2016 Decision Trees - Part 3 9

### Subtree replacement

- Bottom-up*
- Consider replacing a tree only after considering all its subtrees

2016 Decision Trees - Part 3 10

### Subtree raising

- Delete node
- Redistribute instances
- Slower than subtree replacement (*Worthwhile?*)

2016 Decision Trees - Part 3 11

### Estimating error rates

- Prune only if it does not increase the estimated error
- Error on the training data is NOT a useful estimator (*would result in almost no pruning*)
- Use hold-out set for pruning ("reduced-error pruning")
- C4.5's method
  - Derive confidence interval from training data
  - Use a heuristic limit, derived from this, for pruning
  - Standard Bernoulli-process-based method
  - Shaky statistical assumptions (based on training data)

2016 Decision Trees - Part 3 12

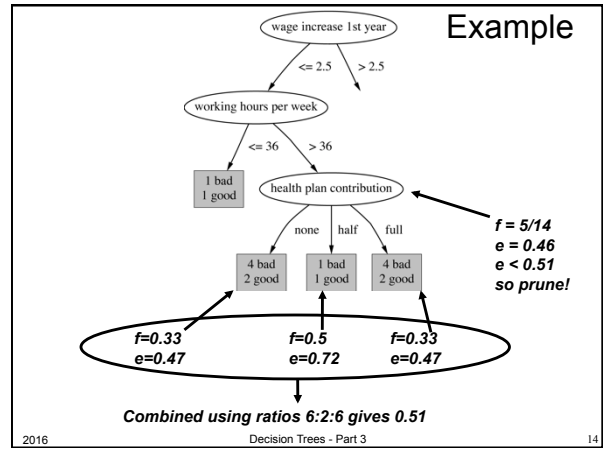
### C4.5's method

- Error estimate for subtree is weighted sum of error estimates for all its leaves
- Error estimate for a node:

$$e = (f + \frac{z^2}{2N} + z \sqrt{\frac{f - \frac{f^2}{N} + \frac{z^2}{4N^2}}{N}}) / (1 + \frac{z^2}{N})$$

- If  $c = 25\%$  then  $z = 0.69$  (from normal distribution)
- $f$  is the error on the training data
- $N$  is the number of instances covered by the leaf

### Example



The end