**Data Mining and Knowledge Discovery: Practice Notes**

Discussion about decision trees

dr. Petra Kralj Novak
Petra.Kralj.Novak@ijs.si
15.11.2018

1

---

## Discussion

- How much is the information gain for the "attribute" Person? How would it perform on the test set?
- How do we compute entropy for a target variable that has three values? Lenses = {hard=4, soft=5, none=13}
- What would be the classification accuracy of our decision tree if we pruned it at the node *Astigmatic*?
- What are the stopping criteria for building a decision tree?
- How would you compute the information gain for a numeric attribute?

2

---

## Discussion about decision trees

→ • How much is the information gain for the "attribute" Person? How would it perform on the test set?
- How do we compute entropy for a target variable that has three values? Lenses = {hard=4, soft=5, none=13}
- What would be the classification accuracy of our decision tree if we pruned it at the node *Astigmatic*?
- What are the stopping criteria for building a decision tree?
- How would you compute the information gain for a numeric attribute?

3

---

## Information gain of the "attribute" Person



On training set
- As many values as there are examples
- Each leaf has exactly one example
- E(1/1, 0/1) = 0 (entropy of each leaf is zero)
- The weighted sum of entropies is zero
- The information gain is maximum (as much as the entropy of the entire training set)

On testing set
- The values from the testing set do not appear in the tree

4

---

## Discussion about decision trees

- How much is the information gain for the "attribute" Person? How would it perform on the test set?
→ • How do we compute entropy for a target variable that has three values? Lenses = {hard=4, soft=5, none=13}
- What would be the classification accuracy of our decision tree if we pruned it at the node *Astigmatic*?
- What are the stopping criteria for building a decision tree?
- How would you compute the information gain for a numeric attribute?

5

---

## Entropy{hard=4, soft=5, none=13}=

$= E(4/22, 5/22, 13/22)$

$= -\sum p_i * \log_2 p_i$

$= -4/22 * \log_2 4/22 - 5/22 * \log_2 5/22 - 13/22 * \log_2 13/22$
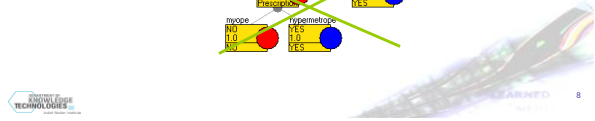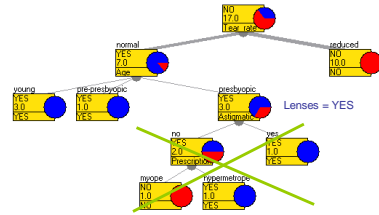
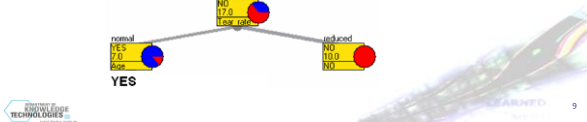$= 1.38$

6

## Discussion about decision trees

- How much is the information gain for the "attribute" Person? How would it perform on the test set?
- How do we compute entropy for a target variable that has three values? Lenses = {hard=4, soft=5, none=13}
→ • What would be the classification accuracy of our decision tree if we pruned it at the node Astigmatic?
- What are the stopping criteria for building a decision tree?
- How would you compute the information gain for a numeric attribute?

## Decision tree pruning



Lenses = YES

## These two trees are equivalent

## Classification accuracy of the pruned tree

| Person | Age | Prescription | Astigmatic | Tear_rate | Lenses |
|--------|-----|--------------|------------|-----------|--------|
| P3 | young | hypermetrope | no | normal | YES |
| P9 | pre-presbyopic | myope | no | normal | YES |
| P12 | pre-presbyopic | hypermetrope | no | reduced | NO |
| P13 | pre-presbyopic | myope | yes | normal | YES |
| P15 | pre-presbyopic | hypermetrope | yes | normal | NO |
| P16 | pre-presbyopic | hypermetrope | yes | reduced | NO |
| P23 | presbyopic | hypermetrope | yes | normal | NO |

$Ca = (3+2)/ (3+2+2+0) = 71\%$



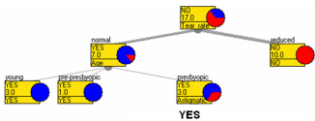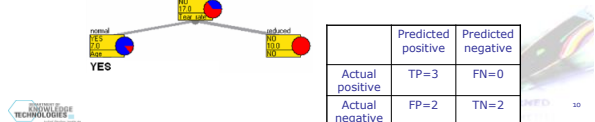| | Predicted positive | Predicted negative |
|--------|--------|--------|
| Actual positive | TP=3 | FN=0 |
| Actual negative | FP=2 | TN=2 |

## Discussion about decision trees

- How much is the information gain for the "attribute" Person? How would it perform on the test set?
- How do we compute entropy for a target variable that has three values? Lenses = {hard=4, soft=5, none=13}
- What would be the classification accuracy of our decision tree if we pruned it at the node Astigmatic?
→ • What are the stopping criteria for building a decision tree?
- How would you compute the information gain for a numeric attribute?

## Stopping criteria for building a decision tree

- ID3
  - "Pure" nodes (entropy =0)
  - Out of attributes
- J48 (C4.5)
  - Minimum number of instances in a leaf constraint

## Discussion about decision trees

- How much is the information gain for the "attribute" Person? How would it perform on the test set?
- How do we compute entropy for a target variable that has three values? Lenses = {hard=4, soft=5, none=13}
- What would be the classification accuracy of our decision tree if we pruned it at the node Astigmatic?
- What are the stopping criteria for building a decision tree?
- → How would you compute the information gain for a numeric attribute?

13

## Information gain of a numeric attribute

| Age | Lenses |
|-----|--------|
| 67 | YES |
| 52 | YES |
| 63 | NO |
| 26 | YES |
| 65 | NO |
| 23 | YES |
| 65 | NO |
| 25 | YES |
| 26 | YES |
| 57 | NO |
| 49 | NO |
| 23 | YES |
| 39 | NO |
| 55 | NO |
| 53 | NO |
| 38 | NO |
| 67 | YES |
| 54 | NO |
| 29 | YES |
| 46 | NO |
| 44 | YES |
| 32 | NO |
| 39 | NO |
| 45 | YES |

14

## Information gain of a numeric attribute

| Age | Lenses |
|-----|--------|
| 67 | YES |
| 52 | YES |
| 63 | NO |
| 26 | YES |
| 65 | NO |
| 23 | YES |
| 65 | NO |
| 25 | YES |
| 26 | YES |
| 57 | NO |
| 49 | NO |
| 23 | YES |
| 39 | NO |
| 55 | NO |
| 53 | NO |
| 38 | NO |
| 67 | YES |
| 54 | NO |
| 29 | YES |
| 46 | NO |
| 44 | YES |
| 32 | NO |
| 39 | NO |
| 45 | YES |

**Sort by Age** →

| Age | Lenses |
|-----|--------|
| 23 | YES |
| 23 | YES |
| 25 | YES |
| 26 | YES |
| 26 | YES |
| 29 | YES |
| 32 | NO |
| 38 | NO |
| 39 | NO |
| 39 | NO |
| 44 | YES |
| 45 | YES |
| 46 | NO |
| 49 | NO |
| 52 | YES |
| 53 | NO |
| 54 | NO |
| 55 | NO |
| 57 | NO |
| 63 | NO |
| 65 | NO |
| 65 | NO |
| 67 | YES |
| 67 | YES |

15

## Information gain of a numeric attribute

| Age | Lenses |
|-----|--------|
| 67 | YES |
| 52 | YES |
| 63 | NO |
| 26 | YES |
| 65 | NO |
| 23 | YES |
| 65 | NO |
| 25 | YES |
| 26 | YES |
| 57 | NO |
| 49 | NO |
| 23 | YES |
| 39 | NO |
| 55 | NO |
| 53 | NO |
| 38 | NO |
| 67 | YES |
| 54 | NO |
| 29 | YES |
| 46 | NO |
| 44 | YES |
| 32 | NO |
| 39 | NO |
| 45 | YES |

**Sort by Age** →

| Age | Lenses |
|-----|--------|
| 23 | YES |
| 23 | YES |
| 25 | YES |
| 26 | YES |
| 26 | YES |
| 29 | YES |
| 32 | NO |
| 38 | NO |
| 39 | NO |
| 39 | NO |
| 44 | YES |
| 45 | YES |
| 46 | NO |
| 49 | NO |
| 52 | YES |
| 53 | NO |
| 54 | NO |
| 55 | NO |
| 57 | NO |
| 63 | NO |
| 65 | NO |
| 65 | NO |
| 67 | YES |
| 67 | YES |

**Define possible splitting points** →

| Age | Lenses |
|-----|--------|
| 23 | YES |
| 23 | YES |
| 25 | YES |
| 26 | YES |
| 26 | YES |
| 29 | YES |
| 32 | NO |
| 38 | NO |
| 39 | NO |
| 39 | NO |
| 44 | YES |
| 45 | YES |
| 46 | NO |
| 49 | NO |
| 52 | YES |
| 53 | NO |
| 54 | NO |
| 55 | NO |
| 57 | NO |
| 63 | NO |
| 65 | NO |
| 65 | NO |
| 67 | YES |
| 67 | YES |

16

## Information gain of a numeric attribute

| Age | Lenses |     |
|-----|--------|-----|
| 23 | YES | |
| 23 | YES | |
| 25 | YES | |
| 26 | YES | |
| 26 | YES | |
| 29 | YES | → 30.5 |
| 32 | NO | |
| 38 | NO | |
| 39 | NO | |
| 39 | NO | → 41.5 |
| 44 | YES | |
| 45 | YES | → 45.5 |
| 46 | NO | |
| 49 | NO | → 50.5 |
| 52 | YES | → 52.5 |
| 53 | NO | |
| 54 | NO | |
| 55 | NO | |
| 57 | NO | |
| 63 | NO | |
| 65 | NO | |
| 65 | NO | → 66 |
| 67 | YES | |
| 67 | YES | |

17

## Information gain of a numeric attribute

| Age | Lenses |     |
|-----|--------|-----|
| 23 | YES | |
| 23 | YES | |
| 25 | YES | |
| 26 | YES | |
| 26 | YES | |
| 29 | YES | → 30.5 |
| 32 | NO | |
| 38 | NO | |
| 39 | NO | |
| 39 | NO | → 41.5 |
| 44 | YES | |
| 45 | YES | → 45.5 |
| 46 | NO | |
| 49 | NO | → 50.5 |
| 52 | YES | → 52.5 |
| 53 | NO | |
| 54 | NO | |
| 55 | NO | |
| 57 | NO | |
| 63 | NO | |
| 65 | NO | |
| 65 | NO | → 66 |
| 67 | YES | |
| 67 | YES | |

**Age**
<30.5 ⟋   ⟍ >=30.5
6/24        18/24

$E(6/6 , 0/6) = 0$     $E(5/18 , 13/18) = 0.85$

18

3

## Information gain of a numeric attribute

| Age | Lenses |
|-----|--------|
| 23 | YES |
| 23 | YES |
| 25 | YES |
| 26 | YES |
| 26 | YES |
| 29 | YES |
| 32 | NO |
| 38 | NO |
| 39 | NO |
| 39 | NO |
| 44 | YES |
| 45 | YES |
| 46 | NO |
| 49 | NO |
| 52 | YES |
| 53 | NO |
| 54 | NO |
| 55 | NO |
| 57 | NO |
| 63 | NO |
| 65 | NO |
| 65 | NO |
| 67 | YES |
| 67 | YES |

→ 30.5
→ 41.5
→ 45.5
→ 50.5
→ 52.5
→ 66

$E(S) = E(11/24, 13/24) = 0.99$

Age
<30.5      >=30.5
6/24          18/24

$E(6/6, 0/6) = 0$      $E(5/18, 13/18) = 0.85$

InfoGain $(S, Age_{30.5})$=

$= E(S) - \sum p_v E(p_v)$

$= 0.99 - (6/24*0 + 18/24*0.85)$

$= 0.35$

19

## Information gain of a numeric attribute

| Age | Lenses |
|-----|--------|
| 23 | YES |
| 23 | YES |
| 25 | YES |
| 26 | YES |
| 26 | YES |
| 29 | YES |
| 32 | NO |
| 38 | NO |
| 39 | NO |
| 39 | NO |
| 44 | YES |
| 45 | YES |
| 46 | NO |
| 49 | NO |
| 52 | YES |
| 53 | NO |
| 54 | NO |
| 55 | NO |
| 57 | NO |
| 63 | NO |
| 65 | NO |
| 65 | NO |
| 67 | YES |
| 67 | YES |

→ 30.5
→ 41.5
→ 45.5
→ 50.5
→ 52.5
→ 66

Age
<30.5      >=30.5

InfoGain $(S, Age_{30.5}) = 0.35$

Age                    Age
<41.5  >=41.5      <45.5  >=45.5

Age                    Age
<50.5  >=50.5      <52.5  >=52.5

Age
<66  >=66

20

## Decision trees

• Many possible decision trees

$$\sum_{i=0}^{k} 2^i (k-i) = -k + 2^{k+1} - 2$$

– k is the number of binary attributes
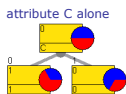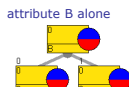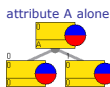
• Heuristic search with information gain
• Information gain is short-sighted

21

## Trees are shortsighted (1)

| A | B | C | A xor B |
|---|---|---|---------|
| 1 | 1 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 |
| 1 | 1 | 1 | 0 |
| 1 | 0 | 1 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 1 | 0 |

• Three attributes:
  A, B and C
• Target variable is a logical combination attributes A and B class = A xor B
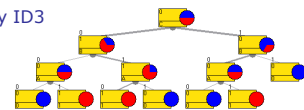• Attribute C is random w.r.t. the target variable
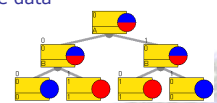
22

## Trees are shortsighted (2)

attribute A alone

attribute B alone

attribute C alone

Attribute C has the highest information gain!

23

## Trees are shortsighted (3)

• Decision tree by ID3

• The real model behind the data

24

4

## Overcoming shortsightedness of decision trees

- Random forests
  (Breinmann & Cutler, 2001)
  - A random forest is a set of decision trees
  - Each tree is induced from a bootstrap sample of examples
  - For each node of the tree, select among a subset of attributes
  - All the trees vote for the classification
  - See also ansemble learning
- ReliefF for attribute estimation
  (Kononenko el al., 1997)

- More on 14.12.2018 by Martin Žnidaršič

25