

# Data Mining and Knowledge Discovery: Practice Notes

dr. Petra Kralj Novak

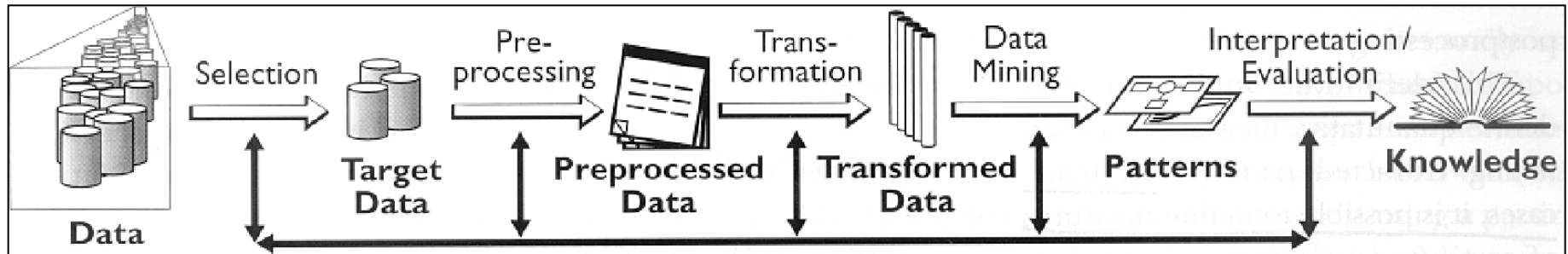
[Petra.Kralj.Novak@ijs.si](mailto:Petra.Kralj.Novak@ijs.si)

2016/11/15

- Prof. Nada Lavrač:
  - Data mining overview
  - Advanced topics
- Dr. Petra Kralj Novak
  - Data mining basis



# Keywords



- Data
  - Attribute, example, attribute-value data, target variable, class, discretization
- Algorithms
  - Decision tree induction, entropy, information gain, overfitting, Occam's razor, model pruning, naïve Bayes classifier, KNN, association rules, support, confidence, numeric prediction, regression tree, model tree, heuristics vs. exhaustive search, predictive vs. descriptive DM
- Evaluation
  - Train set, test set, accuracy, confusion matrix, cross validation, true positives, false positives, ROC space, error, precision, recall

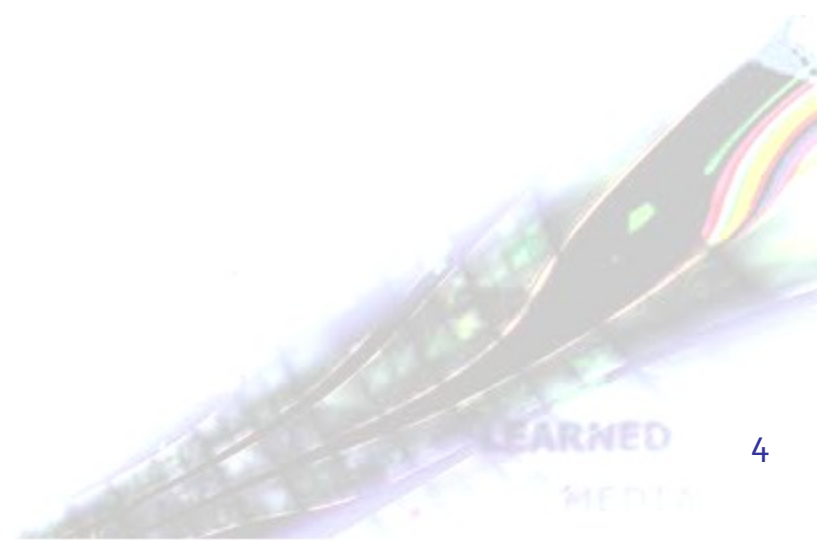
# Decision tree induction

Given

- Attribute-value data with nominal target variable

Induce

- A decision tree and estimate its performance on new data



# Attribute-value data

(nominal)  
target  
variable

attributes

Person	Age	Prescription	Astigmatic	Tear_Rate	Lenses
P1	young	myope	no	normal	<b>YES</b>
P2	young	myope	no	reduced	<b>NO</b>
P3	young	hypermetrope	no	normal	<b>YES</b>
P4	young	hypermetrope	no	reduced	<b>NO</b>
P5	young	myope	yes	normal	<b>YES</b>
P6	young	myope	yes	reduced	<b>NO</b>
P7	young	hypermetrope	yes	normal	<b>YES</b>
P8	young	hypermetrope	yes	reduced	<b>NO</b>
P9	pre-presbyopic	myope	no	normal	<b>YES</b>
P10	pre-presbyopic	myope	no	reduced	<b>NO</b>
P11	pre-presbyopic	hypermetrope	no	normal	<b>YES</b>
P12	pre-presbyopic	hypermetrope	no	reduced	<b>NO</b>
P13	pre-presbyopic	myope	yes	normal	<b>YES</b>
P14	pre-presbyopic	myope	yes	reduced	<b>NO</b>
P15	pre-presbyopic	hypermetrope	yes	normal	<b>NO</b>
P16	pre-presbyopic	hypermetrope	yes	reduced	<b>NO</b>
P17	presbyopic	myope	no	normal	<b>NO</b>
P18	presbyopic	myope	no	reduced	<b>NO</b>
P19	presbyopic	hypermetrope	no	normal	<b>YES</b>
P20	presbyopic	hypermetrope	no	reduced	<b>NO</b>
P21	presbyopic	myope	yes	normal	<b>YES</b>
P22	presbyopic	myope	yes	reduced	<b>NO</b>
P23	presbyopic	hypermetrope	yes	normal	<b>NO</b>
P24	presbyopic	hypermetrope	yes	reduced	<b>NO</b>

examples

classes  
=  
values of  
the  
(nominal)  
target  
variable

# Decision tree induction (ID3)

Given:

Attribute-value data with nominal target variable

Divide the data into training set (S) and test set (T)

---

Induce a decision tree on training set S:

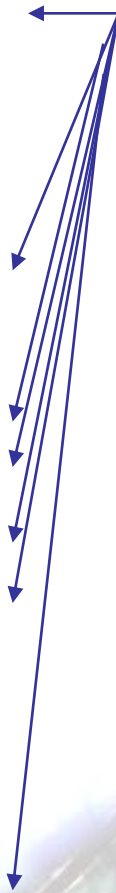
1. Compute the entropy  $E(S)$  of the set S
2. **IF**  $E(S) = 0$
3. The current set is "clean" and therefore a leaf in our tree
4. **IF**  $E(S) > 0$
5. Compute the information gain of each attribute  $\text{Gain}(S, A)$
6. The attribute A with the highest information gain becomes the root
7. Divide the set S into subsets  $S_i$  according to the values of A
8. Repeat steps 1-7 on each  $S_i$

Test the model on the test set T

# Training and test set

Person	Age	Prescription	Astigmatic	Tear_Rate	Lenses
P1	young	myope	no	normal	<b>YES</b>
P2	young	myope	no	reduced	<b>NO</b>
P3	young	hypermetrope	no	normal	<b>YES</b>
P4	young	hypermetrope	no	reduced	<b>NO</b>
P5	young	myope	yes	normal	<b>YES</b>
P6	young	myope	yes	reduced	<b>NO</b>
P7	young	hypermetrope	yes	normal	<b>YES</b>
P8	young	hypermetrope	yes	reduced	<b>NO</b>
P9	pre-presbyopic	myope	no	normal	<b>YES</b>
P10	pre-presbyopic	myope	no	reduced	<b>NO</b>
P11	pre-presbyopic	hypermetrope	no	normal	<b>YES</b>
P12	pre-presbyopic	hypermetrope	no	reduced	<b>NO</b>
P13	pre-presbyopic	myope	yes	normal	<b>YES</b>
P14	pre-presbyopic	myope	yes	reduced	<b>NO</b>
P15	pre-presbyopic	hypermetrope	yes	normal	<b>NO</b>
P16	pre-presbyopic	hypermetrope	yes	reduced	<b>NO</b>
P17	presbyopic	myope	no	normal	<b>NO</b>
P18	presbyopic	myope	no	reduced	<b>NO</b>
P19	presbyopic	hypermetrope	no	normal	<b>YES</b>
P20	presbyopic	hypermetrope	no	reduced	<b>NO</b>
P21	presbyopic	myope	yes	normal	<b>YES</b>
P22	presbyopic	myope	yes	reduced	<b>NO</b>
P23	presbyopic	hypermetrope	yes	normal	<b>NO</b>
P24	presbyopic	hypermetrope	yes	reduced	<b>NO</b>

Put 30% of examples in a separate test set



# Test set

Person	Age	Prescription	Astigmatic	Tear_Rate	Lenses
P3	young	hypermetrope	no	normal	<b>YES</b>
P9	pre-presbyopic	myope	no	normal	<b>YES</b>
P12	pre-presbyopic	hypermetrope	no	reduced	<b>NO</b>
P13	pre-presbyopic	myope	yes	normal	<b>YES</b>
P15	pre-presbyopic	hypermetrope	yes	normal	<b>NO</b>
P16	pre-presbyopic	hypermetrope	yes	reduced	<b>NO</b>
P23	presbyopic	hypermetrope	yes	normal	<b>NO</b>

Put these data away and do not look at them in the training phase!





# Training set

Person	Age	Prescription	Astigmatic	Tear_Rate	Lenses
P1	young	myope	no	normal	<b>YES</b>
P2	young	myope	no	reduced	<b>NO</b>
P4	young	hypermetrope	no	reduced	<b>NO</b>
P5	young	myope	yes	normal	<b>YES</b>
P6	young	myope	yes	reduced	<b>NO</b>
P7	young	hypermetrope	yes	normal	<b>YES</b>
P8	young	hypermetrope	yes	reduced	<b>NO</b>
P10	pre-presbyopic	myope	no	reduced	<b>NO</b>
P11	pre-presbyopic	hypermetrope	no	normal	<b>YES</b>
P14	pre-presbyopic	myope	yes	reduced	<b>NO</b>
P17	presbyopic	myope	no	normal	<b>NO</b>
P18	presbyopic	myope	no	reduced	<b>NO</b>
P19	presbyopic	hypermetrope	no	normal	<b>YES</b>
P20	presbyopic	hypermetrope	no	reduced	<b>NO</b>
P21	presbyopic	myope	yes	normal	<b>YES</b>
P22	presbyopic	myope	yes	reduced	<b>NO</b>
P24	presbyopic	hypermetrope	yes	reduced	<b>NO</b>

# Decision tree induction (ID3)

Given:

Attribute-value data with nominal target variable

Divide the data into training set (S) and test set (T)

---

Induce a decision tree on training set S:

1. Compute the entropy  $E(S)$  of the set S
2. **IF**  $E(S) = 0$
3. The current set is “clean” and therefore a leaf in our tree
4. **IF**  $E(S) > 0$
5. Compute the information gain of each attribute  $\text{Gain}(S, A)$
6. The attribute A with the highest information gain becomes the root
7. Divide the set S into subsets  $S_i$  according to the values of A
8. Repeat steps 1-7 on each  $S_i$

Test the model on the test set T

# Information gain

number of examples in the subset  $S_v$   
(probability of the branch)

$$\sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \cdot E(S_v)$$

number of examples in set  $S$

set  $S$       attribute  $A$

$$\text{Gain}(S, A) = E(S) -$$

# Entropy

$$E(S) = - \sum_{c=1}^N p_c \cdot \log_2 p_c$$

- Calculate the following entropies:

$$E(0,1) =$$

$$E(1/2, 1/2) =$$

$$E(1/4, 3/4) =$$

$$E(1/7, 6/7) =$$

$$E(6/7, 1/7) =$$

$$E(0.1, 0.9) =$$

$$E(0.001, 0.999) =$$

# Entropy

$$E(S) = - \sum_{c=1}^N p_c \cdot \log_2 p_c$$

- Calculate the following entropies:

$$E(0,1) = 0$$

$$E(1/2, 1/2) = 1$$

$$E(1/4, 3/4) = 0.81$$

$$E(1/7, 6/7) = 0.59$$

$$E(6/7, 1/7) = 0.59$$

$$E(0.1, 0.9) = 0.47$$

$$E(0.001, 0.999) = 0.01$$

# Entropy

$$E(S) = - \sum_{c=1}^N p_c \cdot \log_2 p_c$$

- Calculate the following entropies:

$$E(0,1) = 0$$

$$E(1/2, 1/2) = 1$$

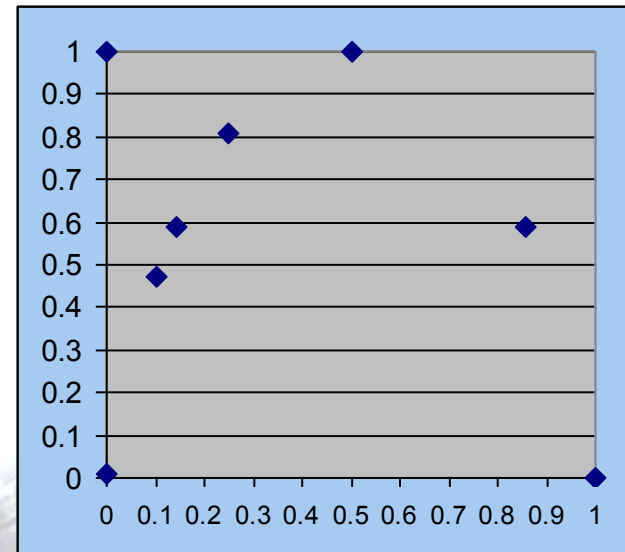
$$E(1/4, 3/4) = 0.81$$

$$E(1/7, 6/7) = 0.59$$

$$E(6/7, 1/7) = 0.59$$

$$E(0.1, 0.9) = 0.47$$

$$E(0.001, 0.999) = 0.01$$



# Entropy

$$E(S) = - \sum_{c=1}^N p_c \cdot \log_2 p_c$$

- Calculate the following entropies:

$$E(0,1) = 0$$

$$E(1/2, 1/2) = 1$$

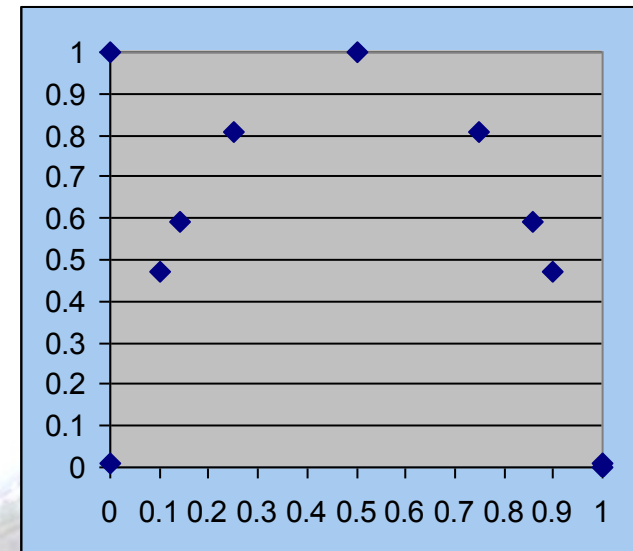
$$E(1/4, 3/4) = 0.81$$

$$E(1/7, 6/7) = 0.59$$

$$E(6/7, 1/7) = 0.59$$

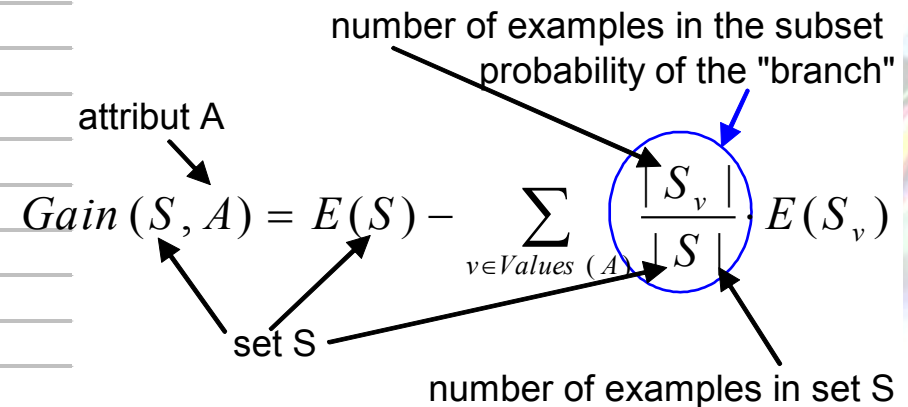
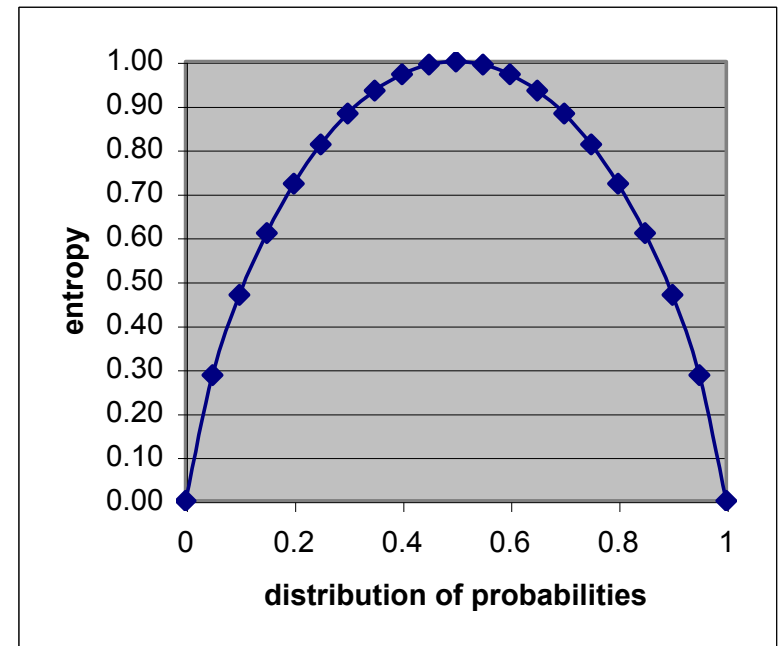
$$E(0.1, 0.9) = 0.47$$

$$E(0.001, 0.999) = 0.01$$



# Entropy and information gain

probability of class 1	probability of class 2	entropy $E(p_1, p_2) = -p_1 \cdot \log_2(p_1) - p_2 \cdot \log_2(p_2)$
$p_1$	$p_2 = 1 - p_1$	
0	1	0.00
0.05	0.95	0.29
0.10	0.90	0.47
0.15	0.85	0.61
0.20	0.80	0.72
0.25	0.75	0.81
0.30	0.70	0.88
0.35	0.65	0.93
0.40	0.60	0.97
0.45	0.55	0.99
0.50	0.50	1.00
0.55	0.45	0.99
0.60	0.40	0.97
0.65	0.35	0.93
0.70	0.30	0.88
0.75	0.25	0.81
0.80	0.20	0.72
0.85	0.15	0.61
0.90	0.10	0.47
0.95	0.05	0.29
1	0	0.00





# Decision tree induction (ID3)

Given:

Attribute-value data with nominal target variable

Divide the data into training set (S) and test set (T)

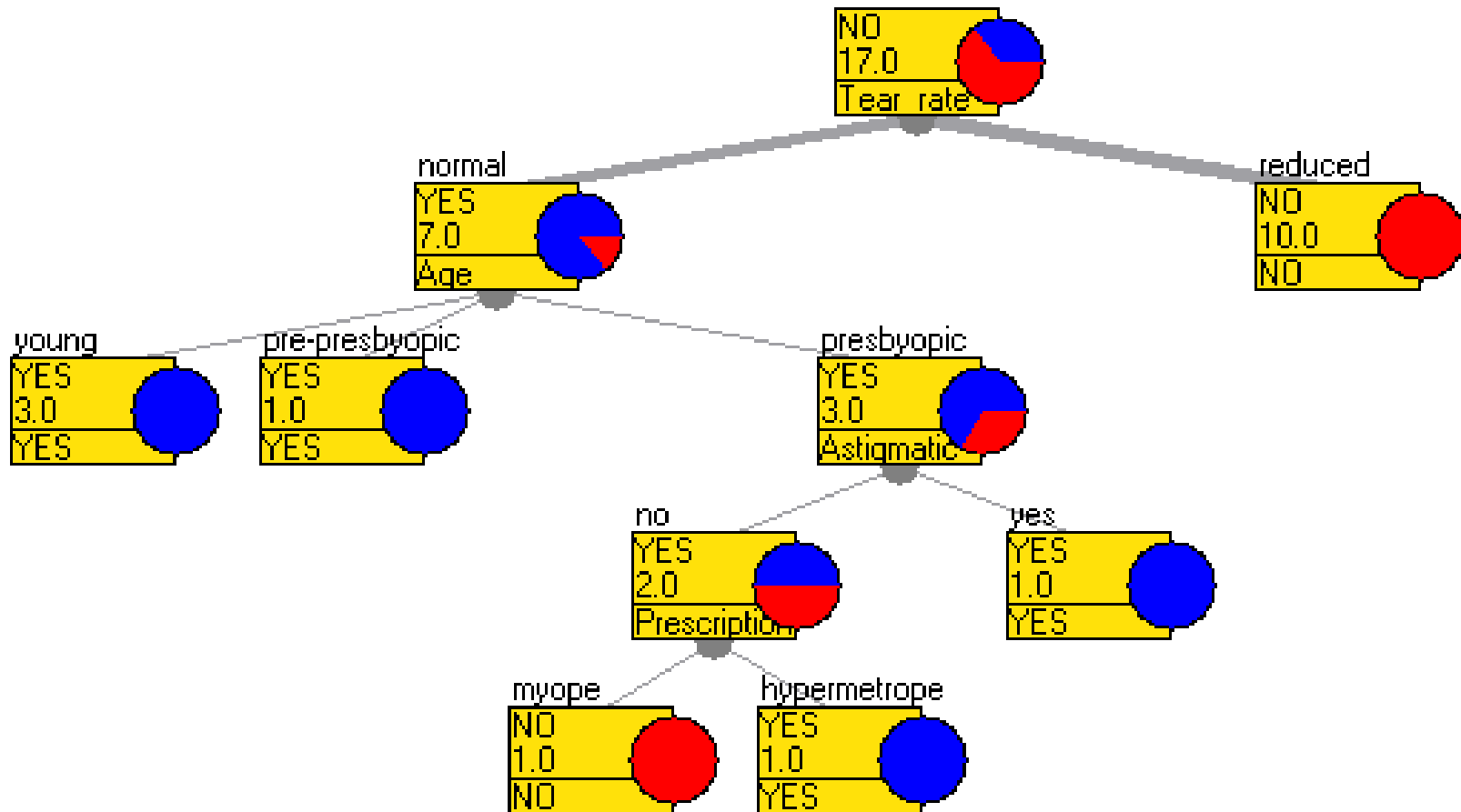
---

Induce a decision tree on training set S:

1. Compute the entropy  $E(S)$  of the set S
2. **IF**  $E(S) = 0$
3. The current set is "clean" and therefore a leaf in our tree
4. **IF**  $E(S) > 0$
5. Compute the information gain of each attribute  $\text{Gain}(S, A)$
6. The attribute A with the highest information gain becomes the root
7. Divide the set S into subsets  $S_i$  according to the values of A
8. Repeat steps 1-7 on each  $S_i$

Test the model on the test set T

# Decision tree



# Confusion matrix

		predicted	
		Predicted positive	Predicted negative
actual	Actual positive	TP	FN
	Actual negative	FP	TN

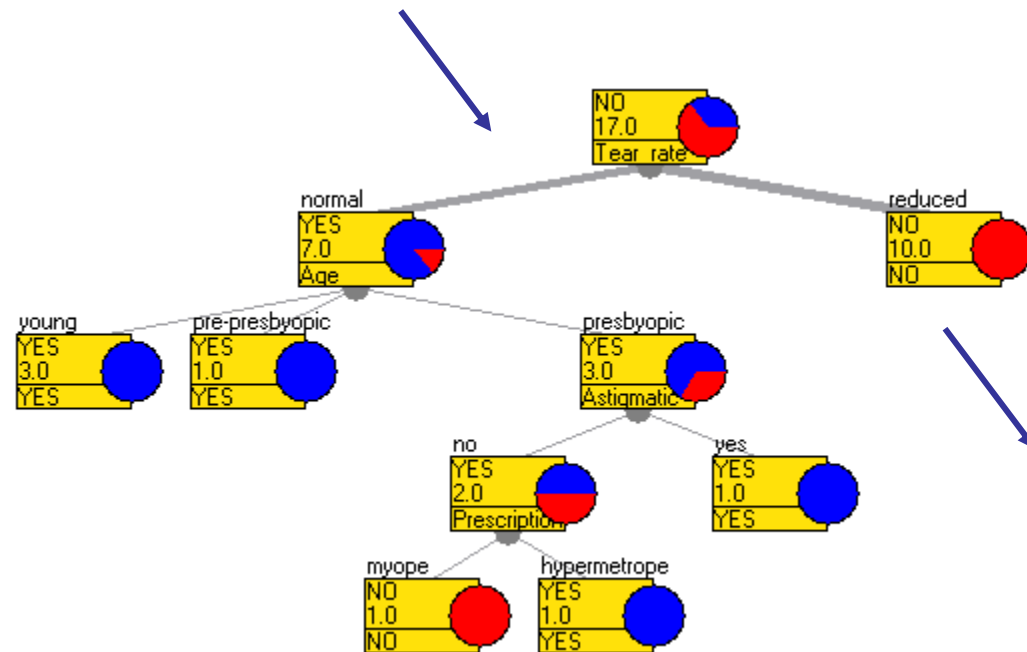
- Confusion matrix is a matrix showing actual and predicted classifications
- Classification measures can be calculated from it, like classification accuracy
  - =  $\#(\text{correctly classified examples}) / \#(\text{all examples})$
  - =  $(TP+TN) / (TP+TN+FP+FN)$



# Evaluating decision tree accuracy

Person	Age	Prescription	Astigmatic	Tear_Rate	Lenses
P3	young	hypermetrope	no	normal	<b>YES</b>
P9	pre-presbyopic	myope	no	normal	<b>YES</b>
P12	pre-presbyopic	hypermetrope	no	reduced	<b>NO</b>
P13	pre-presbyopic	myope	yes	normal	<b>YES</b>
P15	pre-presbyopic	hypermetrope	yes	normal	<b>NO</b>
P16	pre-presbyopic	hypermetrope	yes	reduced	<b>NO</b>
P23	presbyopic	hypermetrope	yes	normal	<b>NO</b>

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



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

# Discussion

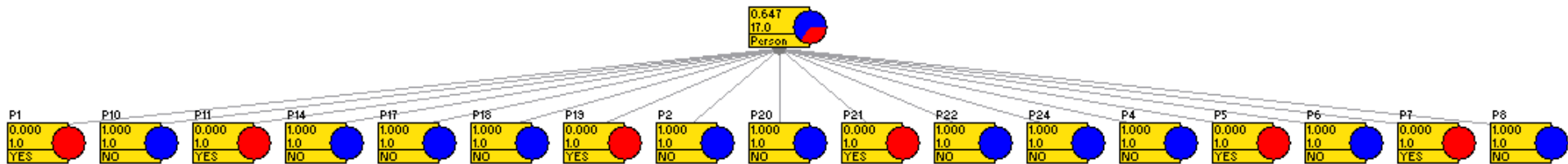
- How much is the information gain for the “attribute” Person? How would it perform on the test set?
- How do we compute entropy for a target variable that has three values? Lenses = {hard=4, soft=5, none=13}
- What would be the classification accuracy of our decision tree if we pruned it at the node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



# 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?

# Information gain of the "attribute" Person



## On training set

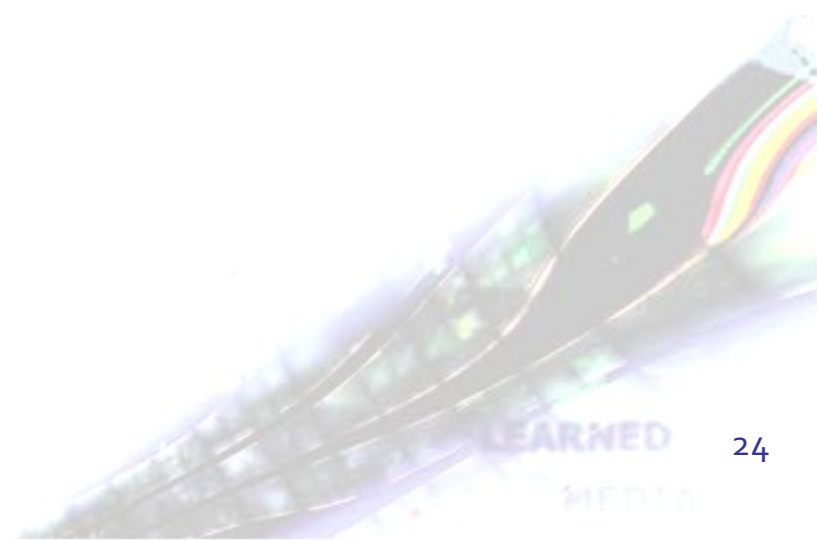
- As many values as there are examples
- Each leaf has exactly one example
- $E(1/1, 0/1) = 0$  (entropy of each leaf is zero)
- The weighted sum of entropies is zero
- The information gain is maximum (as much as the entropy of the entire training set)

## On testing set

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

# 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?





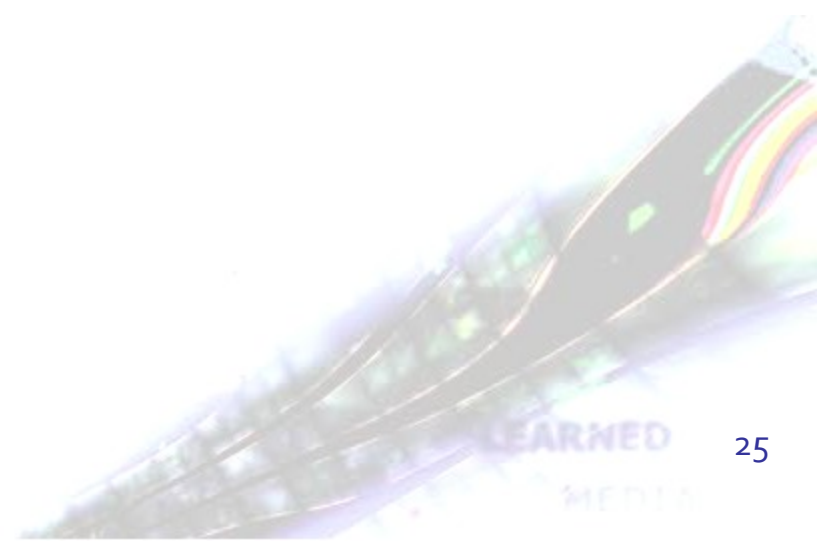
Entropy {hard=4, soft=5, none=13} =

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

$$= -\sum p_i \cdot \log_2 p_i$$

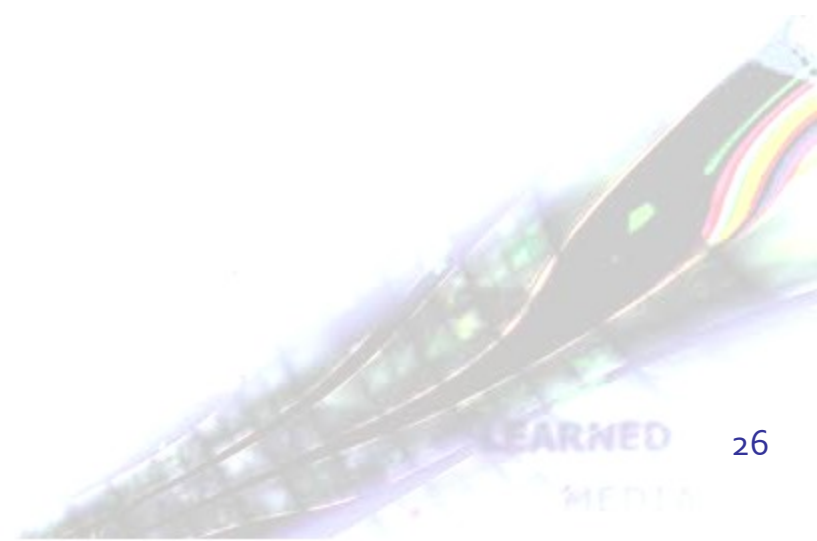
$$= -4/22 * \log_2 4/22 - 5/22 * \log_2 5/22 - 13/22 * \log_2 13/22$$

$$= 1.38$$

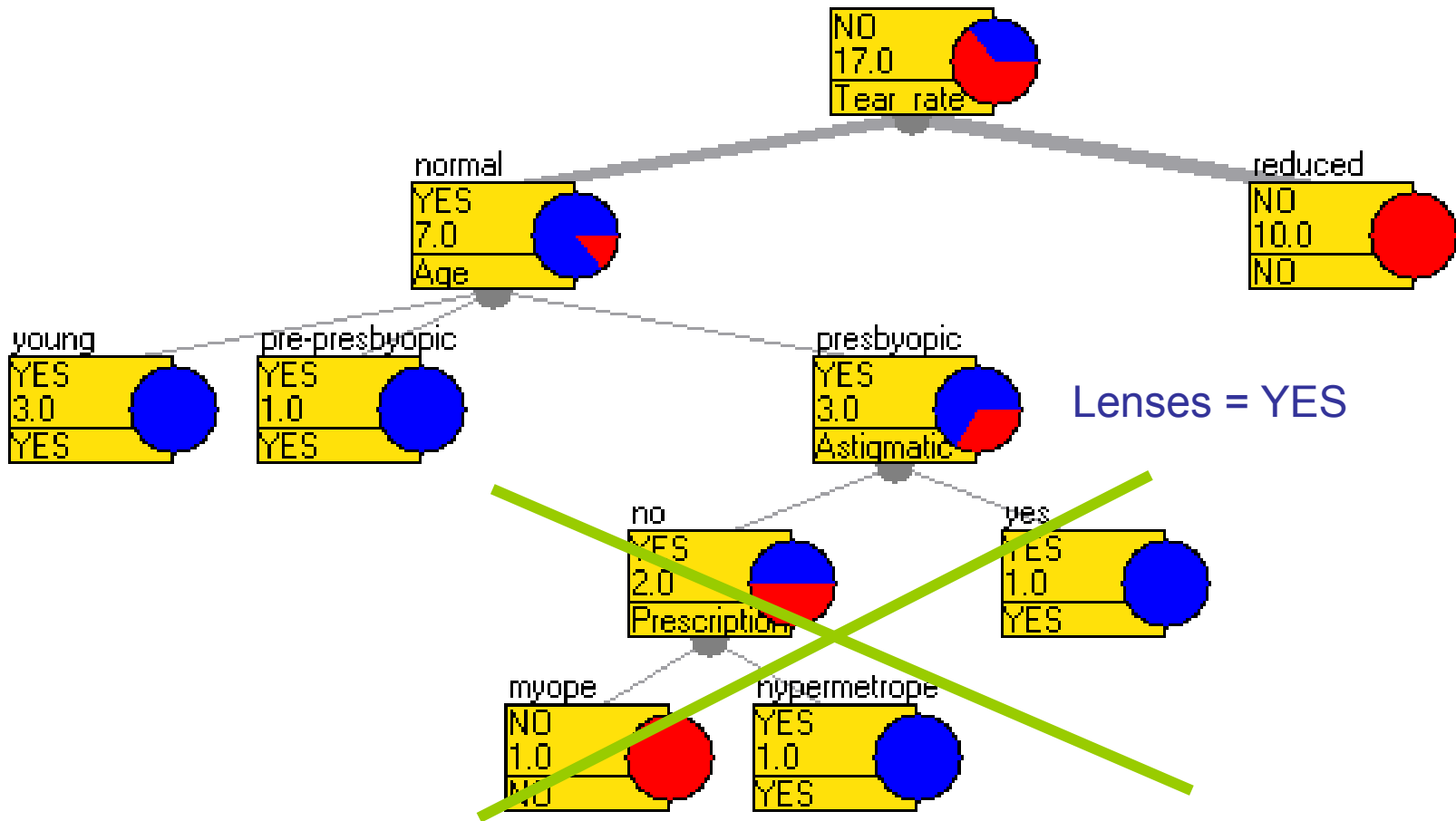


# Discussion about decision trees

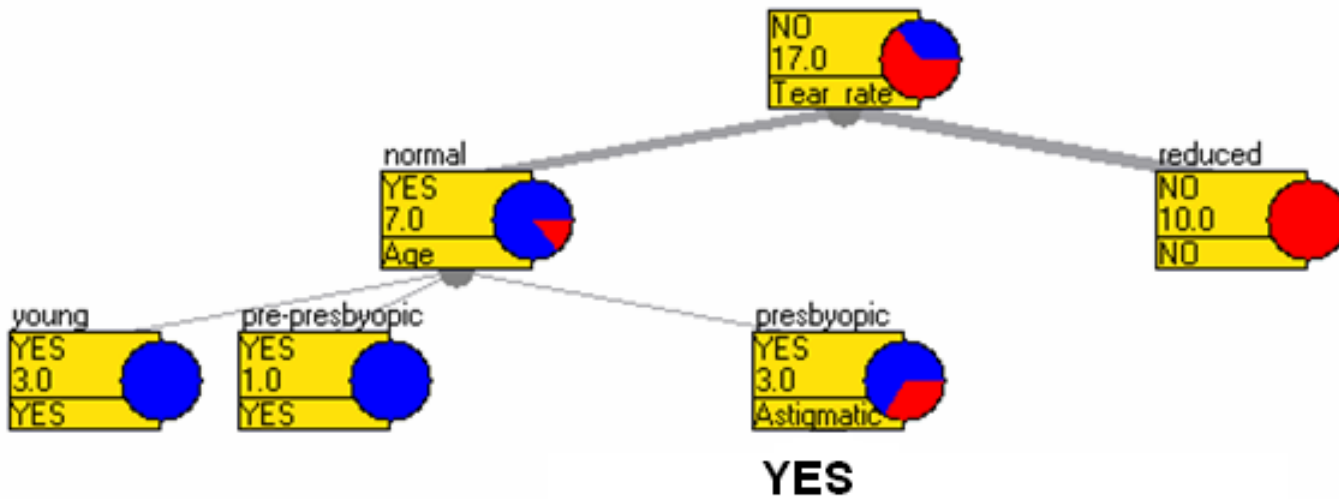
- How much is the information gain for the “attribute” Person? How would it perform on the test set?
- How do we compute entropy for a target variable that has three values? Lenses = {hard=4, soft=5, none=13}
- • What would be the classification accuracy of our decision tree if we pruned it at the node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



# 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

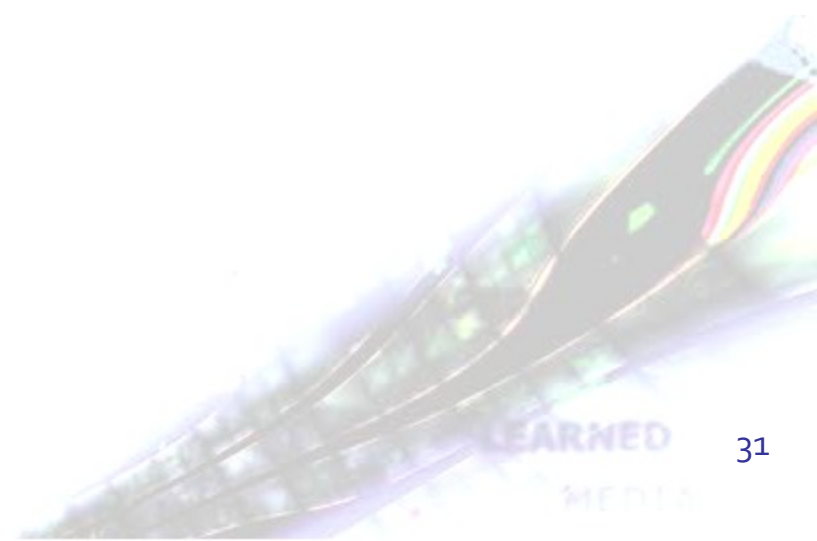
# Discussion about decision trees

- 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?



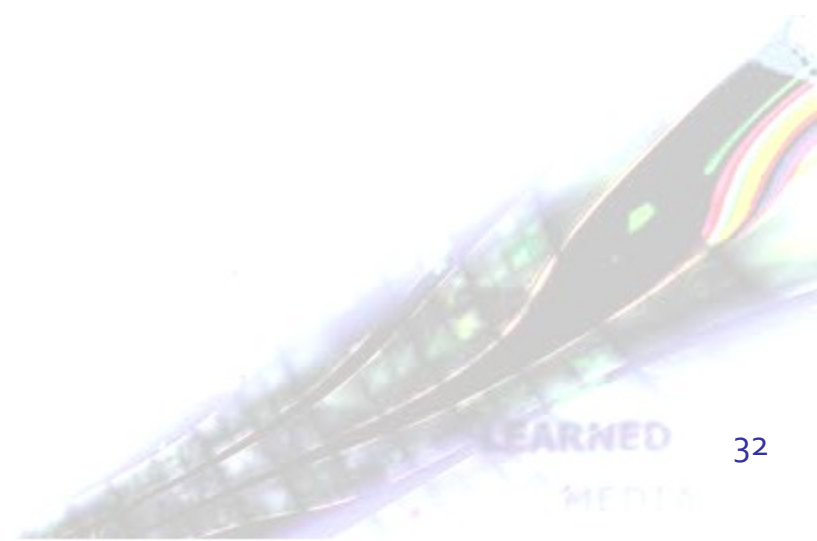
# Stopping criteria for building a decision tree

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



# Discussion about decision trees

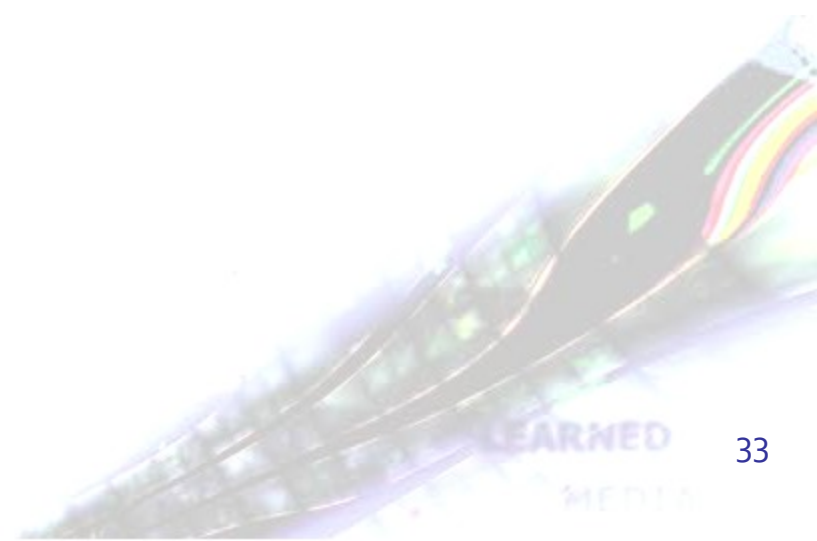
- How much is the information gain for the “attribute” Person? How would it perform on the test set?
- How do we compute entropy for a target variable that has three values? Lenses = {hard=4, soft=5, none=13}
- What would be the classification accuracy of our decision tree if we pruned it at the node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





# 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



# 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

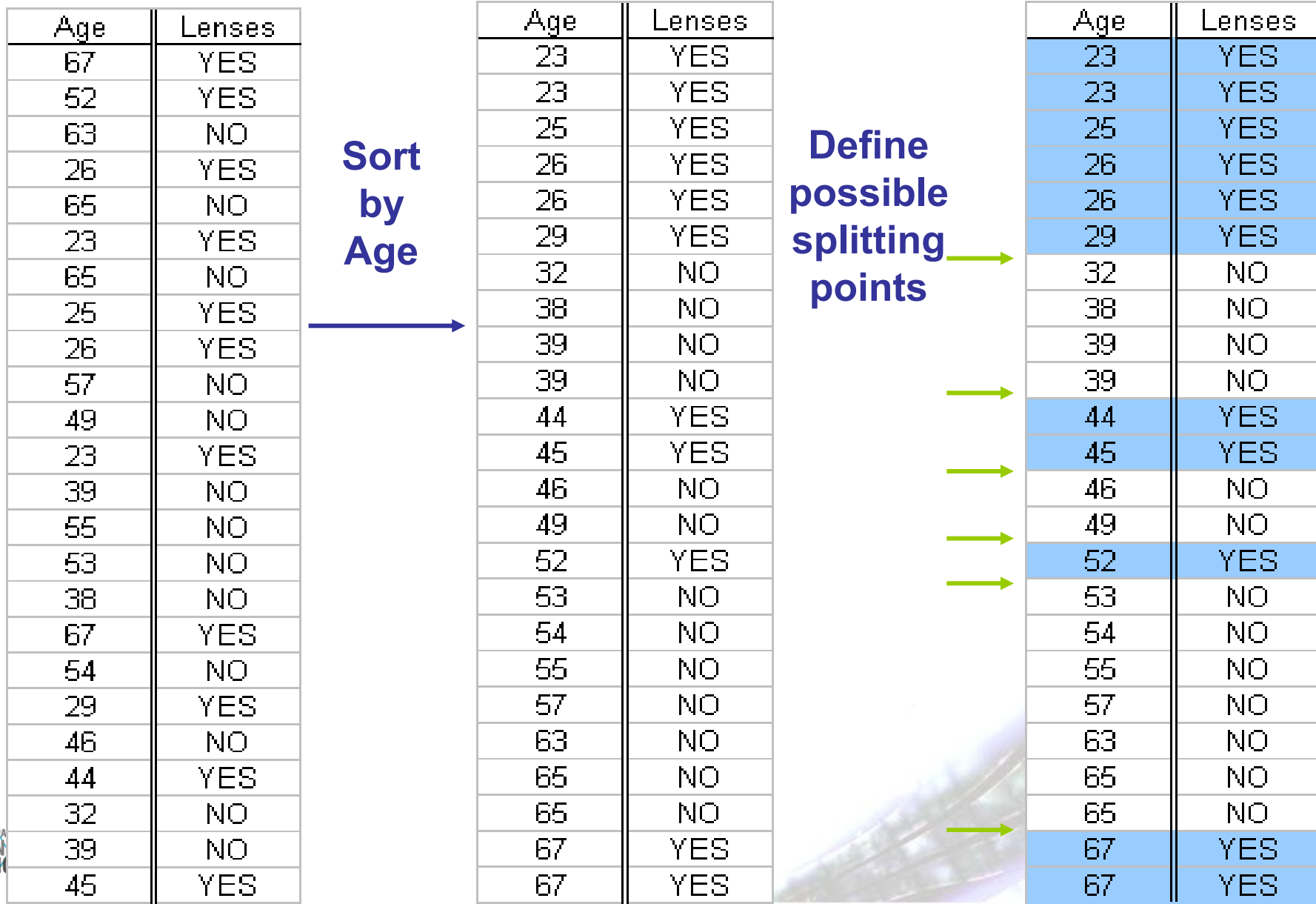
**Sort  
by  
Age**



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

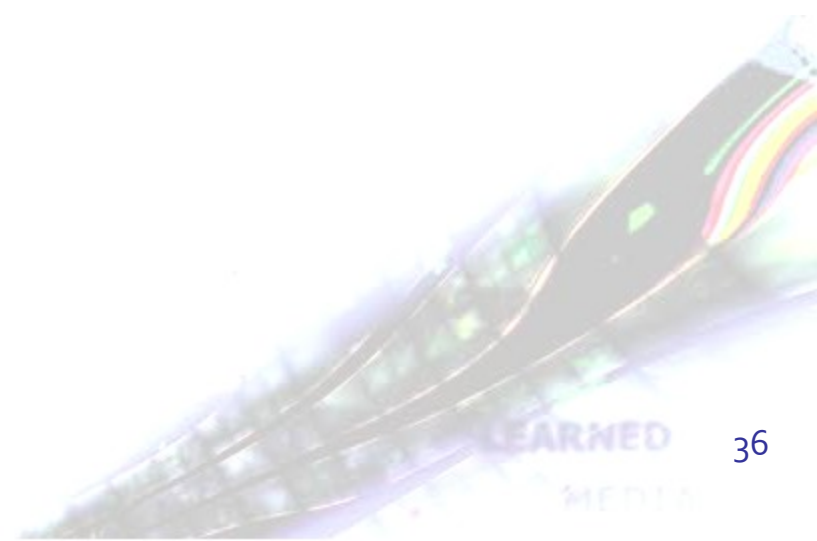


# Information gain of a numeric attribute



# Information gain of a numeric attribute

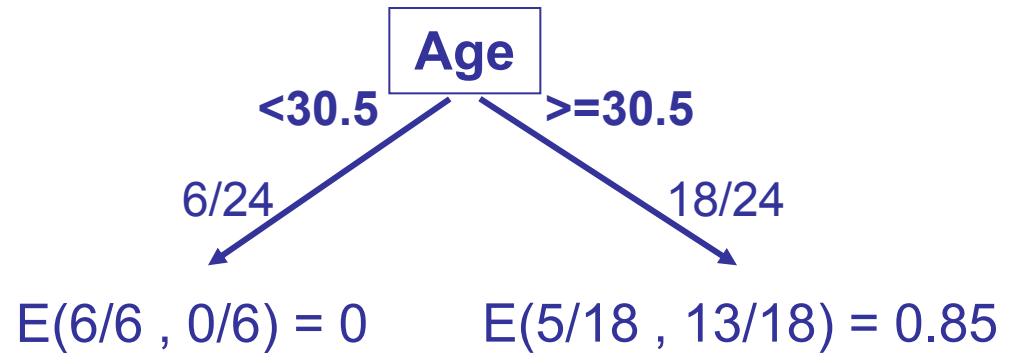
Age	Lenses	
23	YES	
23	YES	
25	YES	
26	YES	
26	YES	
29	YES	→ 30.5
32	NO	
38	NO	
39	NO	
39	NO	→ 41.5
44	YES	
45	YES	→ 45.5
46	NO	
49	NO	→ 50.5
52	YES	→ 52.5
53	NO	
54	NO	
55	NO	
57	NO	
63	NO	
65	NO	
65	NO	
67	YES	→ 66
67	YES	



# 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

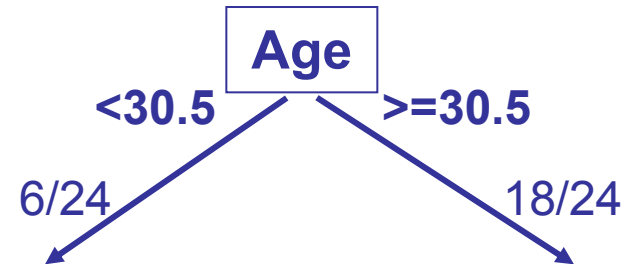


# Information gain of a numeric attribute

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

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



$$E(6/6, 0/6) = 0$$

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

$$\text{InfoGain}(S, \text{Age}_{30.5}) =$$

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

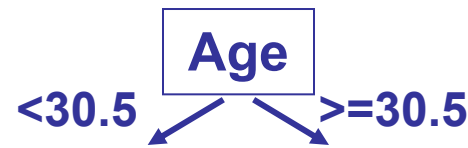
$$= 0.99 - (6/24 * 0 + 18/24 * 0.85)$$

$$= 0.35$$

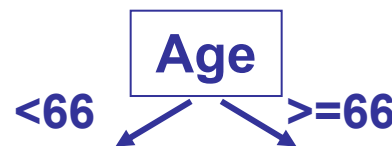
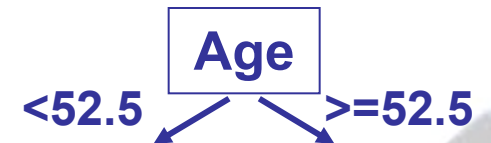
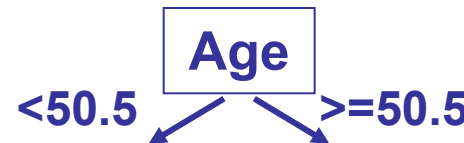
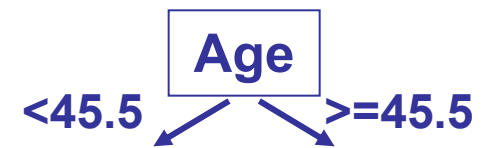
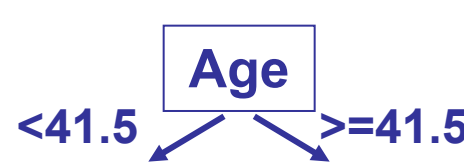
# Information gain of a numeric attribute

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



$\text{InfoGain}(S, \text{Age}_{30.5}) = 0.35$



# 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



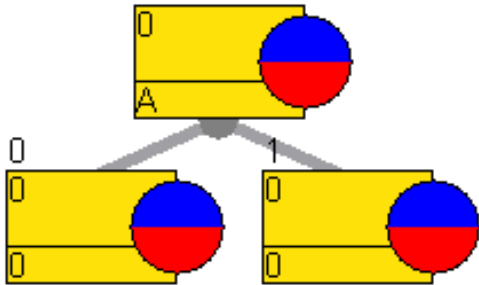
# Trees are shortsighted (1)

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

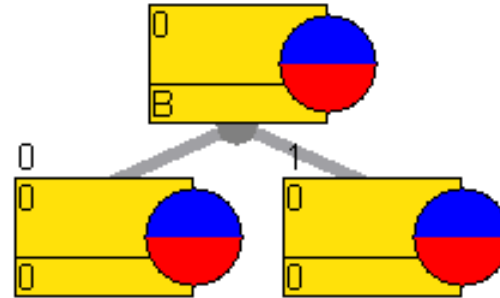
- Three attributes:  
A, B and C
- Target variable is a logical combination attributes A and B  
class = A xor B
- Attribute C is random w.r.t. the target variable

# Trees are shortsighted (2)

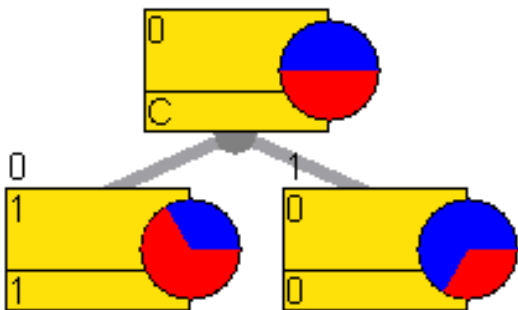
attribute A alone



attribute B alone



attribute C alone



Attribute C has the highest information gain!



# 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 ensemble learning

- ReliefF for attribute estimation

(Kononenko et al., 1997)

