

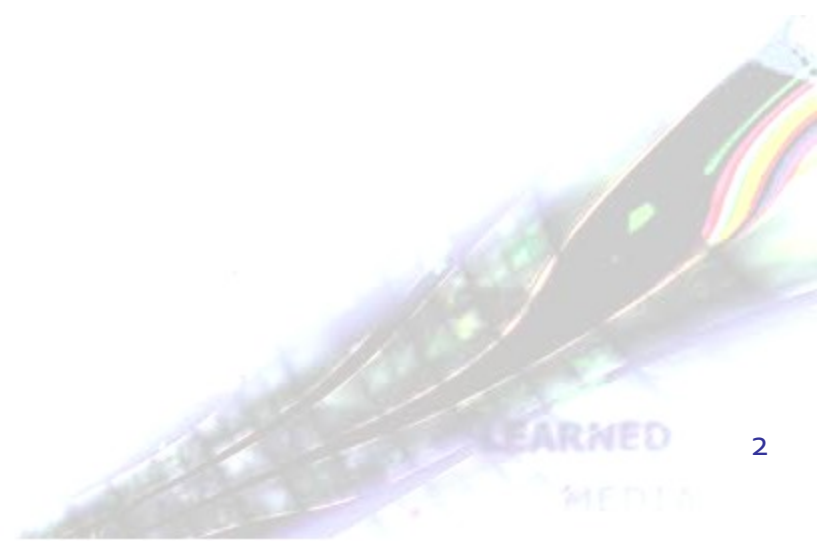
# Data Mining and Knowledge Discovery: Practice Notes

Petra Kralj Novak

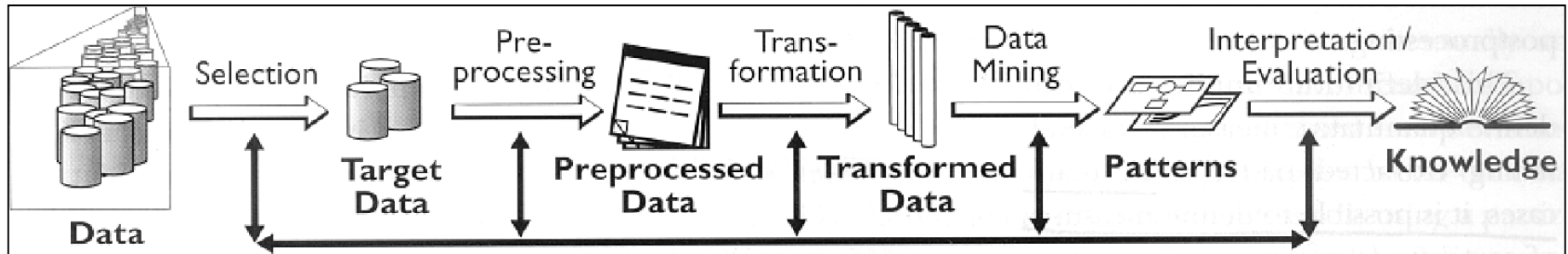
[Petra.Kralj.Novak@ijs.si](mailto:Petra.Kralj.Novak@ijs.si)

2011/11/20

- 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

- **Data mining**

- Heuristics vs. exhaustive search, decision tree induction, entropy, information gain, overfitting, Occam's razor, model pruning, naïve Bayes classifier, KNN, association rules, support, confidence, predictive vs. descriptive DM, numeric prediction, regression tree, model tree

- **Evaluation**

- Train set, test set, accuracy, confusion matrix, cross validation, true positives, false positives, ROC space, error

# Practice plan

---

- 2012/11/20: Predictive data mining 1
  - Decision trees
  - Evaluating classifiers 1: separate test set, confusion matrix, classification accuracy
  - Hands on Weka 1: Just a taste of Weka
- 2012/12/4: Predictive data mining 2
  - Discussion on decision trees
  - Naïve Bayes classifier
  - Evaluating classifiers 2: Cross validation
  - Numeric prediction
  - Hands on Weka 2: Classification and numeric prediction
- 2012/12/4: Descriptive data mining
  - Discussion on classification
    - Association rules
    - Hands on Weka 3: Descriptive data mining
    - Discussion about seminars and exam
- 2013/1/15: Written exam, seminar proposal discussion
- 2013/2/12: Data mining seminar presentations

# 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:

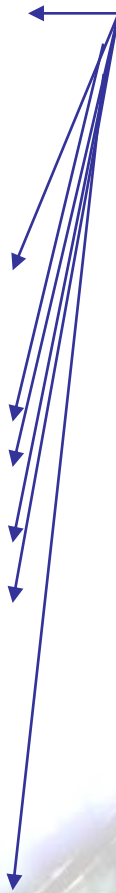
1. Compute the entropy  $E(S)$  of the set S
2. **IF**  $E(S) = 0$
3. The current set is "clean" and therefore a leaf in our tree
4. **IF**  $E(S) > 0$
5. Compute the information gain of each attribute  $\text{Gain}(S, A)$
6. The attribute A with the highest information gain becomes the root
7. Divide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

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

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

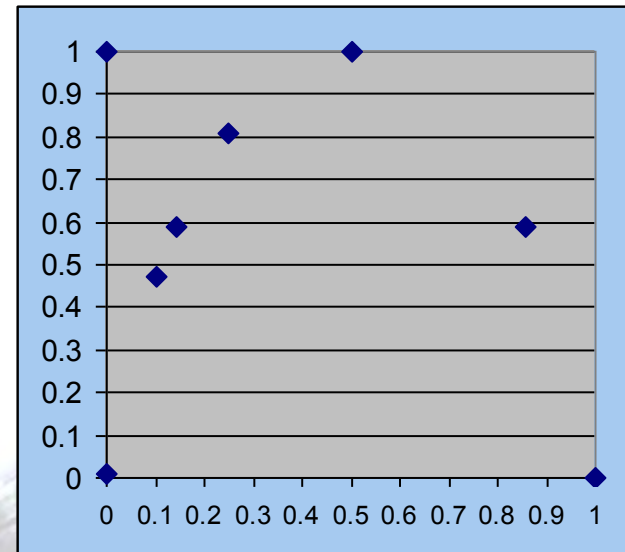
$$E(1/4, 3/4) = 0.81$$

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

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

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

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



# Entropy

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

- Calculate the following entropies:

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

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

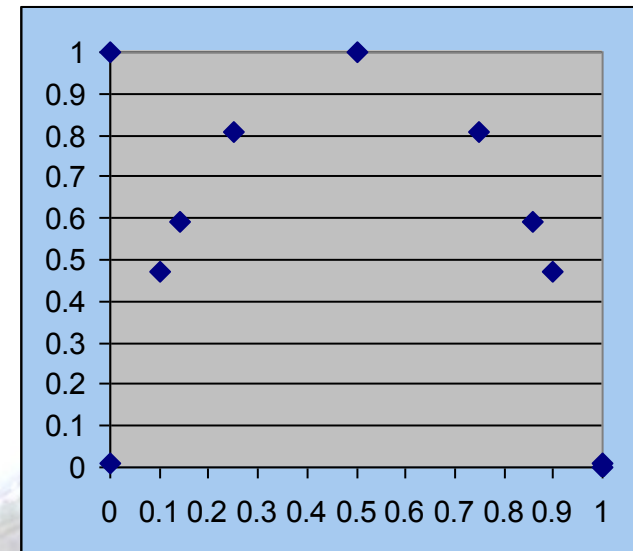
$$E(1/4, 3/4) = 0.81$$

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

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

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

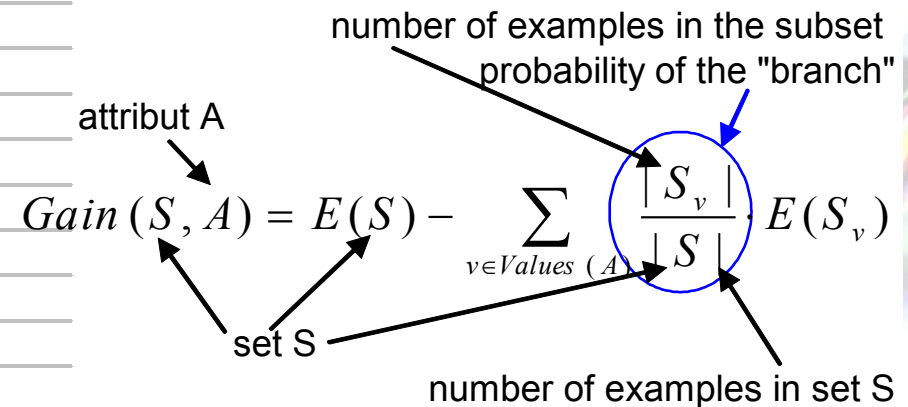
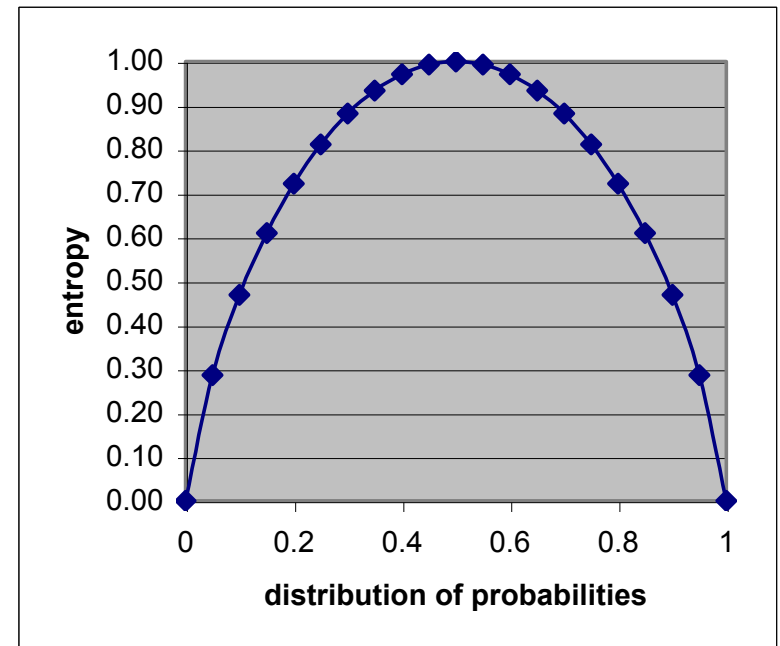
$$E(0.001, 0.999) = 0.01$$





# 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)

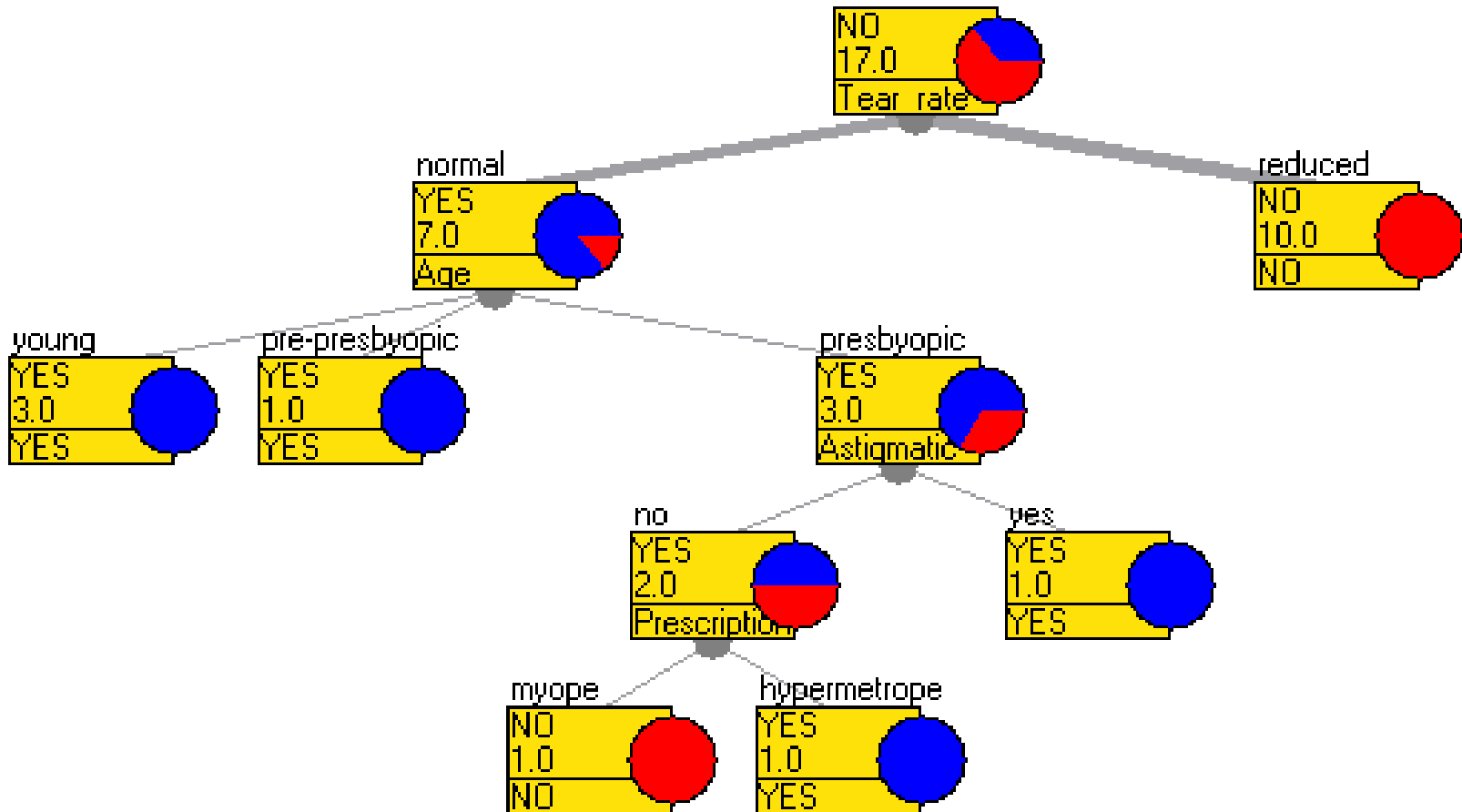
---

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

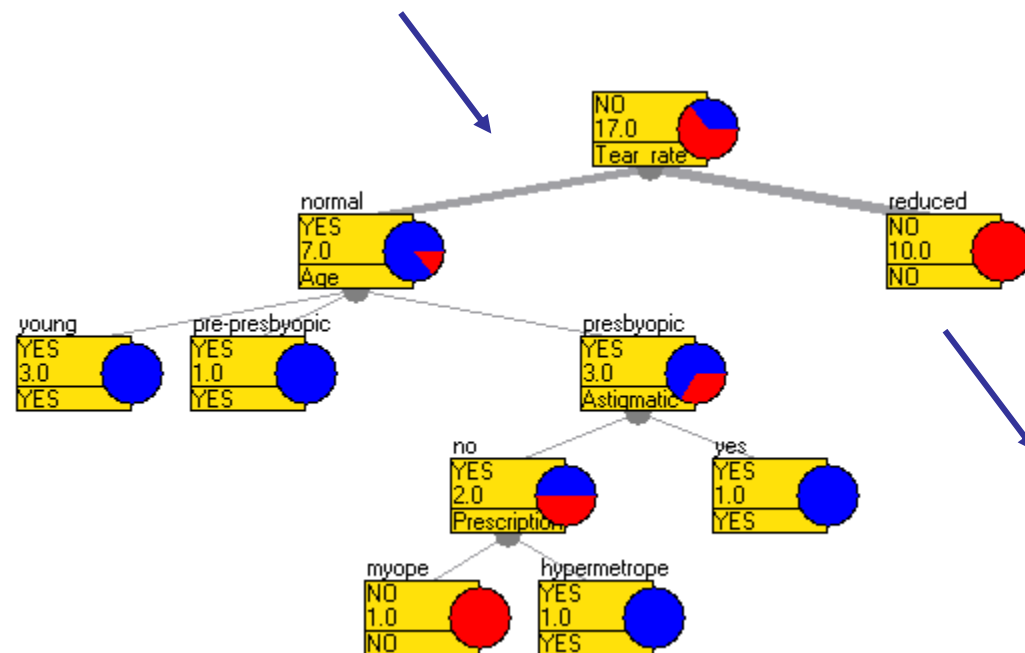
		predicted	
		Predicted positive	Predicted negative
actual	Actual positive	TP	FN
	Actual negative	FP	TN

- 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?

