

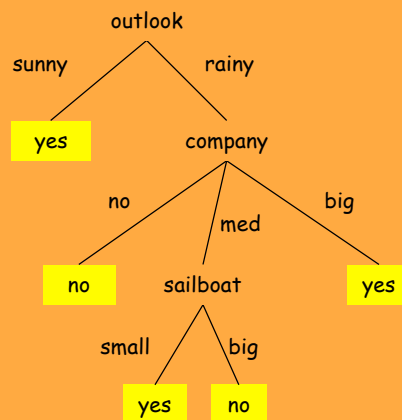
Induction of Decision Trees

Blaž Zupan, Ivan Bratko

magix.fri.uni-lj.si/predavanja/uisp

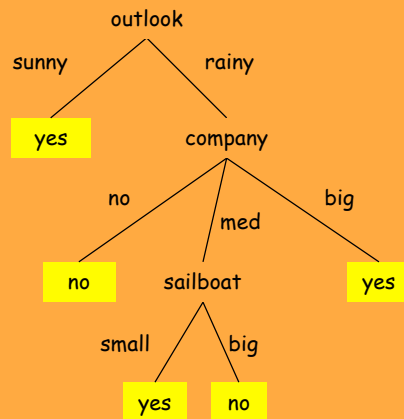
An Example Data Set and Decision Tree

#	Attribute			Class
	Outlook	Company	Sailboat	Sail?
1	sunny	big	small	yes
2	sunny	med	small	yes
3	sunny	med	big	yes
4	sunny	no	small	yes
5	sunny	big	big	yes
6	rainy	no	small	no
7	rainy	med	small	yes
8	rainy	big	big	yes
9	rainy	no	big	no



Classification

#	Attribute			Class
	Outlook	Company	Sailboat	Sail?
1	sunny	no	big	?
2	rainy	big	small	?



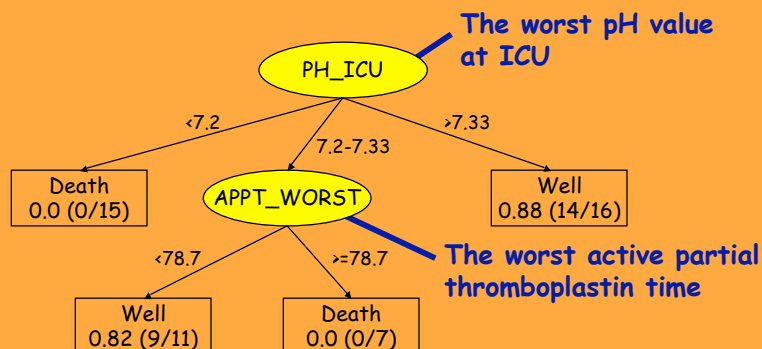
Induction of Decision Trees

- Data Set (Learning Set)
 - Each example = Attributes + Class
- Induced description = Decision tree
- TDIDT
 - Top Down Induction of Decision Trees
- Recursive Partitioning

Some TDIDT Systems

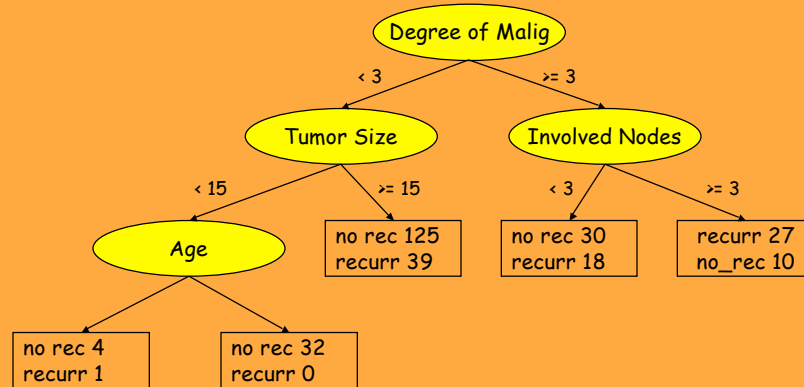
- ID3 (Quinlan 79)
- CART (Breiman et al. 84)
- Assistant (Cestnik et al. 87)
- C4.5 (Quinlan 93)
- See5 (Quinlan 97)
- ...
- Orange (Demšar, Zupan 98-03)

Analysis of Severe Trauma Patients Data



PH_ICU and APPT_WORST are exactly the two factors (theoretically) advocated to be the most important ones in the study by Rotondo et al., 1997.

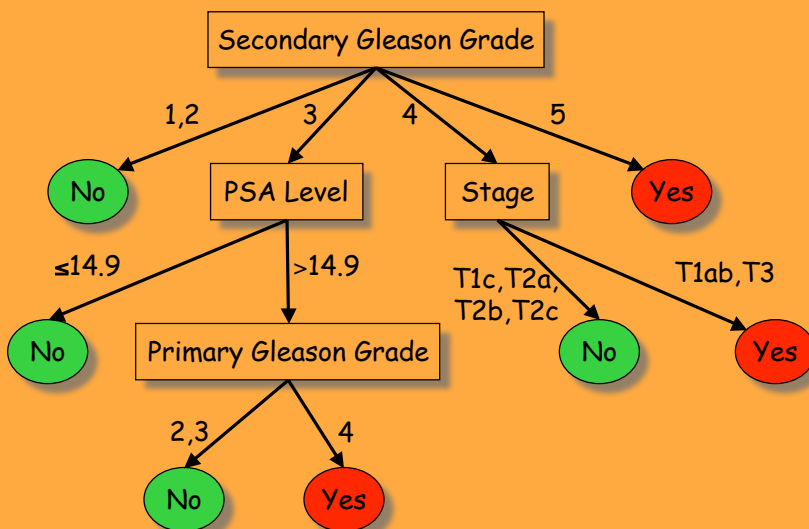
Breast Cancer Recurrence



Tree induced by Assistant Professional

Interesting: Accuracy of this tree compared to medical specialists

Prostate cancer recurrence

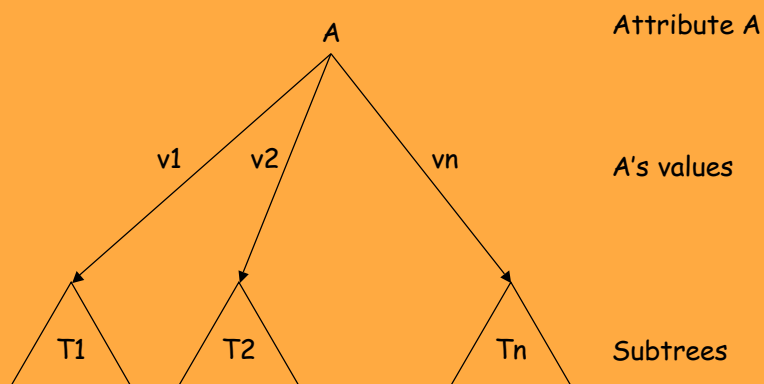


TDIDT Algorithm

- Also known as ID3 (Quinlan)
- To construct decision tree T from learning set S :
 - **If** all examples in S belong to some class C **Then** make leaf labeled C
 - **Otherwise**
 - select the "most informative" attribute A
 - partition S according to A 's values
 - recursively construct subtrees T_1, T_2, \dots , for the subsets of S

TDIDT Algorithm

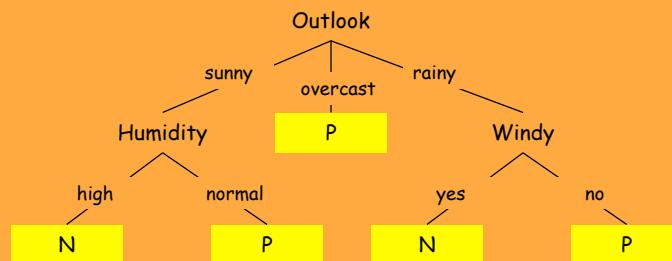
- Resulting tree T is:



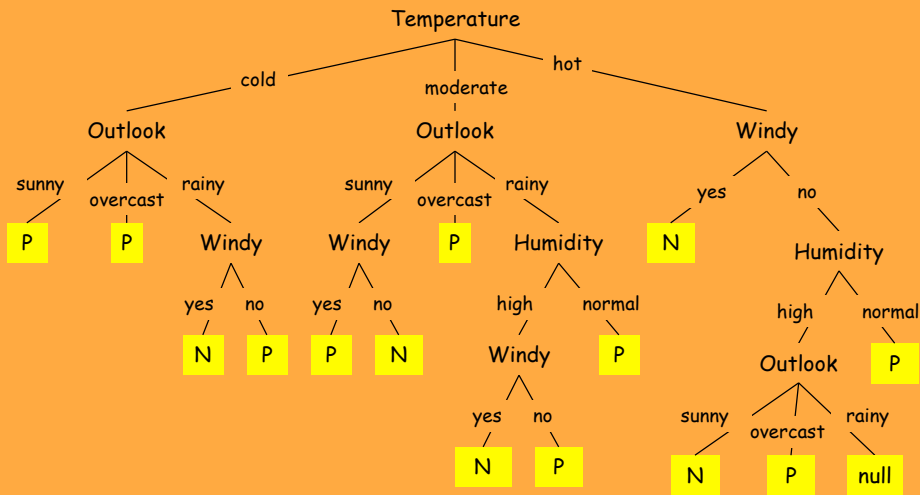
Another Example

#	Attribute				Class
	Outlook	Temperature	Humidity	Windy	
1	sunny	hot	high	no	N
2	sunny	hot	high	yes	N
3	overcast	hot	high	no	P
4	rainy	moderate	high	no	P
5	rainy	cold	normal	no	P
6	rainy	cold	normal	yes	N
7	overcast	cold	normal	yes	P
8	sunny	moderate	high	no	N
9	sunny	cold	normal	no	P
10	rainy	moderate	normal	no	P
11	sunny	moderate	normal	yes	P
12	overcast	moderate	high	yes	P
13	overcast	hot	normal	no	P
14	rainy	moderate	high	yes	N

Simple Tree



Complicated Tree



Attribute Selection Criteria

- **Main principle**
 - Select attribute which partitions the learning set into subsets as "pure" as possible
- **Various measures of purity**
 - Information-theoretic
 - Gini index
 - χ^2
 - ReliefF
 - ...
- **Various improvements**
 - probability estimates
 - normalization
 - binarization, subsetting

Information-Theoretic Approach

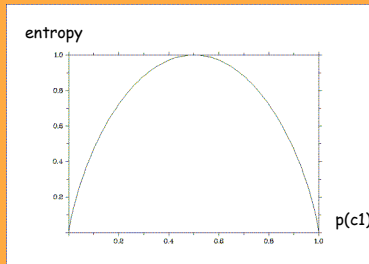
- To classify an object, a certain information is needed
 - I , information
- After we have learned the value of A , we only need some remaining amount of information to classify the object
 - I_{res} , residual information
- Gain
 - $Gain(A) = I - I_{res}(A)$
- The most 'informative' attribute is the one that minimizes I_{res} , *i.e.*, maximizes $Gain$

Entropy

- The average amount of information I needed to classify an object is given by the entropy measure

$$I = - \sum_c p(c) \log_2 p(c)$$

- For a two-class problem:



Residual Information

- After applying attribute A , S is partitioned into subsets according to values v of A
- I_{res} is equal to weighted sum of the amounts of information for the subsets

$$I_{res} = - \sum_v p(v) \sum_c p(c|v) \log_2 p(c|v)$$

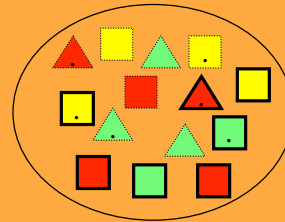
Triangles and Squares

#	Attribute			Shape
	Color	Outline	Dot	
1	green	dashed	no	triange
2	green	dashed	yes	triange
3	yellow	dashed	no	square
4	red	dashed	no	square
5	red	solid	no	square
6	red	solid	yes	triange
7	green	solid	no	square
8	green	dashed	no	triange
9	yellow	solid	yes	square
10	red	solid	no	square
11	green	solid	yes	square
12	yellow	dashed	yes	square
13	yellow	solid	no	square
14	red	dashed	yes	triange

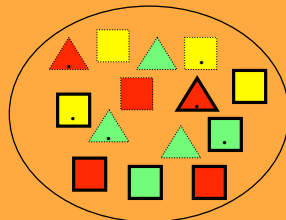
Triangles and Squares

#	Attribute			Shape
	Color	Outline	Dot	
1	green	dashed	no	triangle
2	green	dashed	yes	triangle
3	yellow	dashed	no	square
4	red	dashed	no	square
5	red	solid	no	square
6	red	solid	yes	triangle
7	green	solid	no	square
8	green	dashed	no	triangle
9	yellow	solid	yes	square
10	red	solid	no	square
11	green	solid	yes	square
12	yellow	dashed	yes	square
13	yellow	solid	no	square
14	red	dashed	yes	triangle

Data Set:
A set of classified objects



Entropy



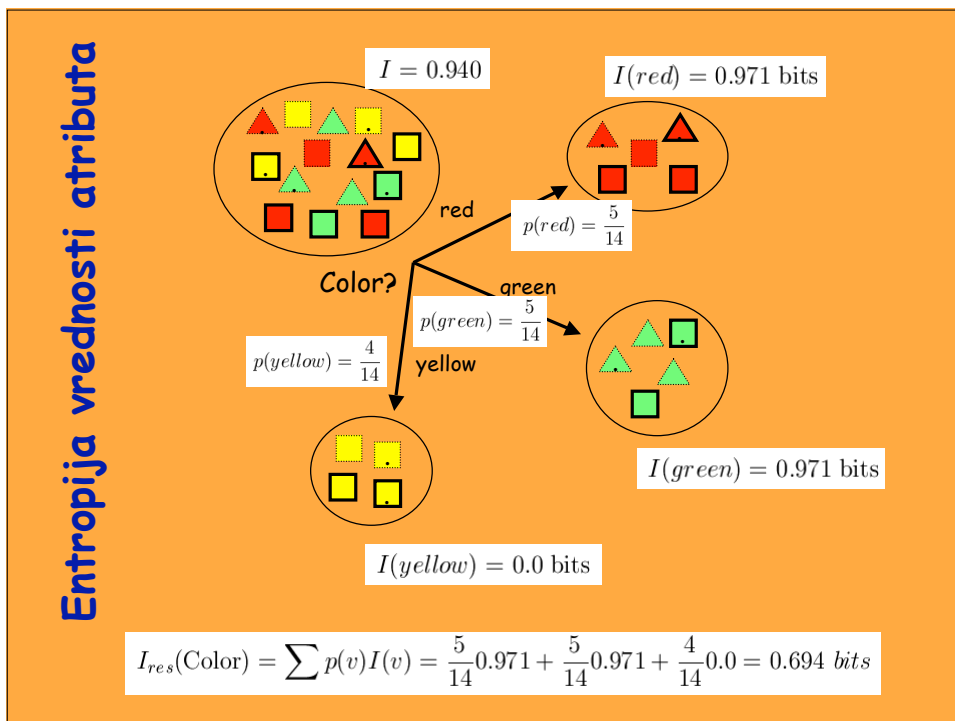
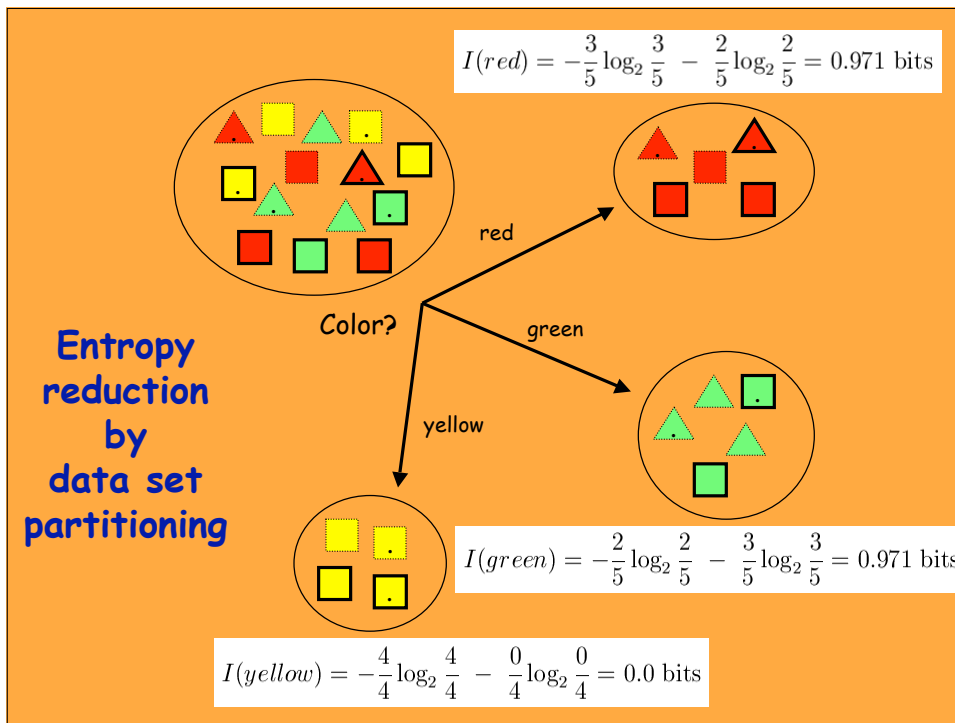
- 5 triangles
- 9 squares
- class probabilities

$$p(\square) = \frac{9}{14}$$

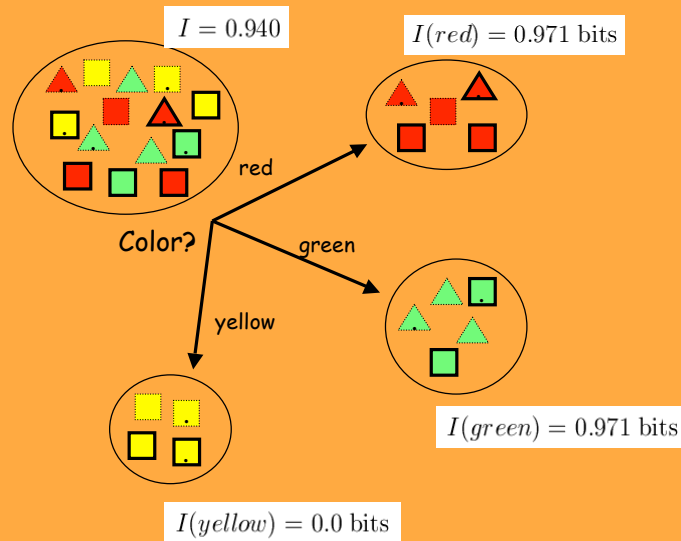
$$p(\triangle) = \frac{5}{14}$$

- entropy

$$I = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.940 \text{ bits}$$



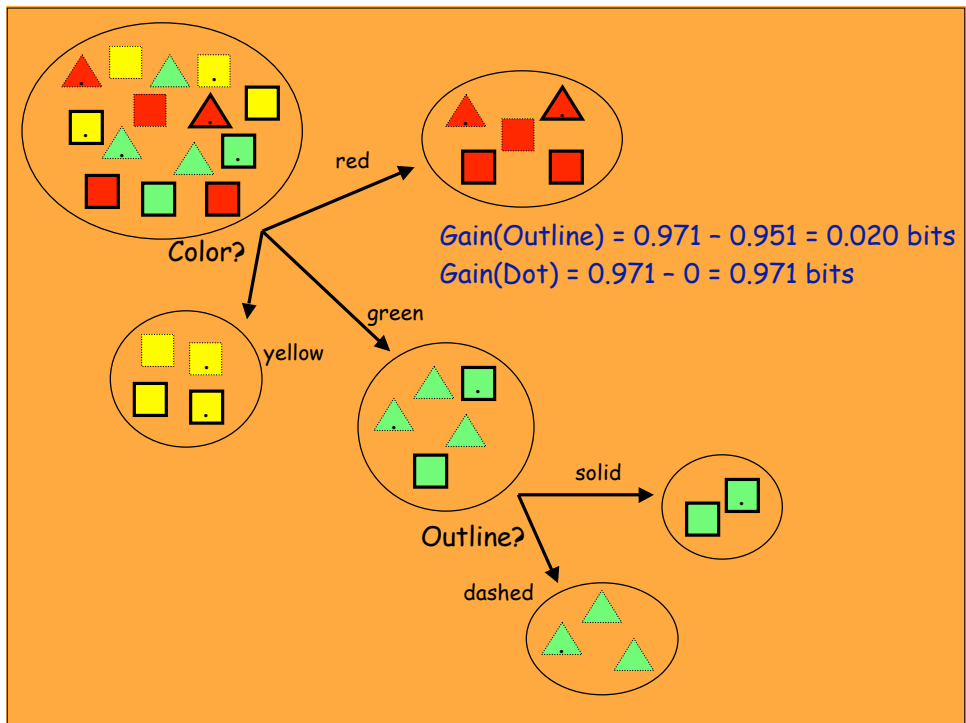
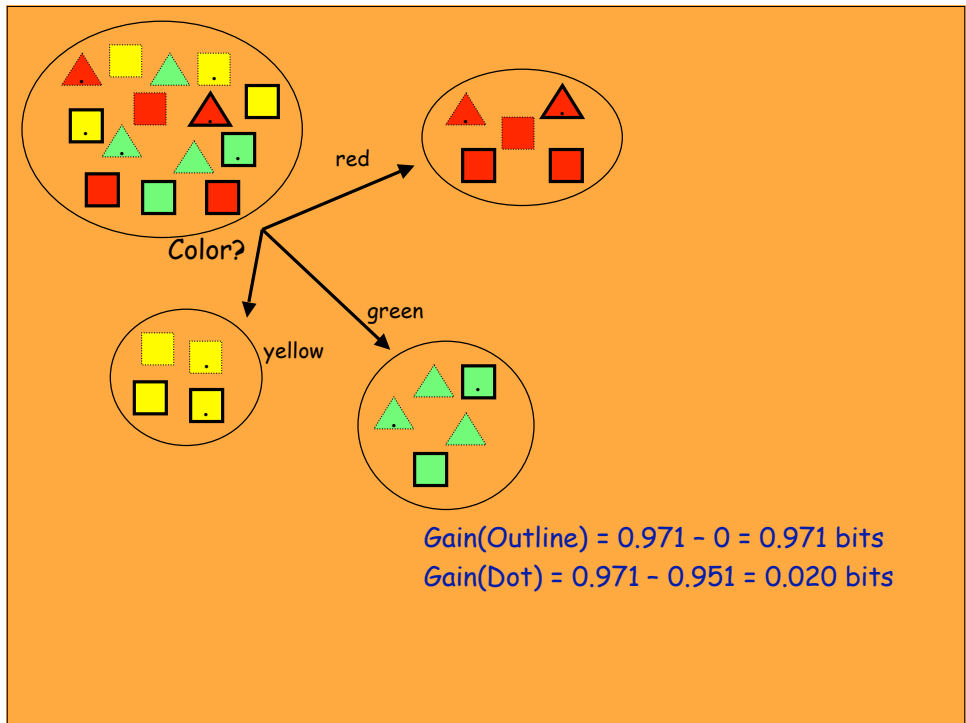
Information Gain

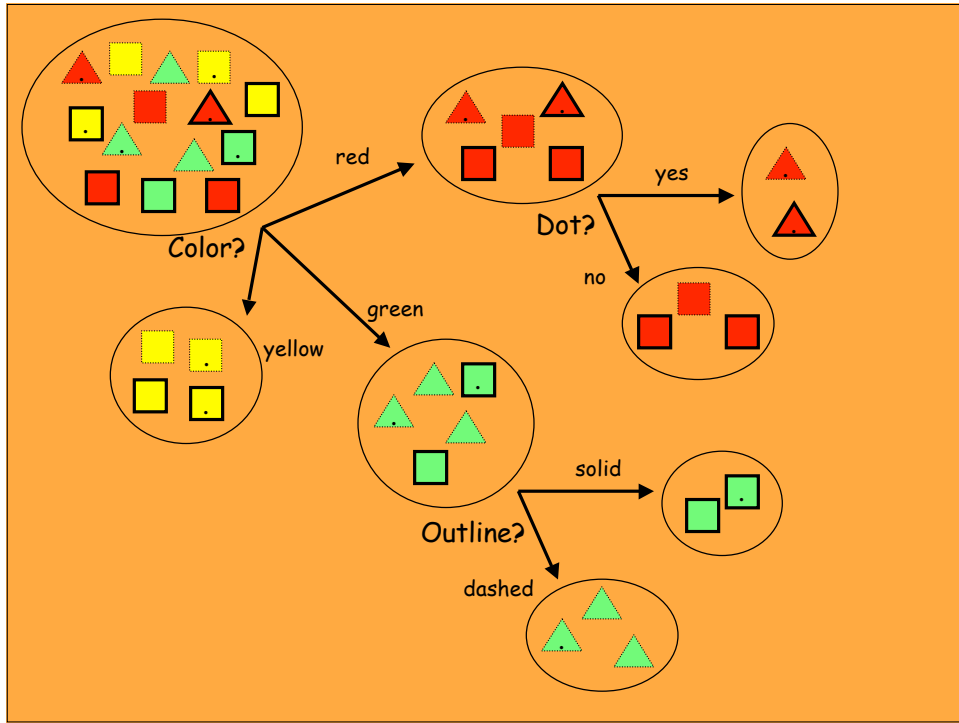


$$\text{Gain}(\text{Color}) = I - I_{\text{res}}(\text{Color}) = 0.940 - 0.694 = 0.246 \text{ bits}$$

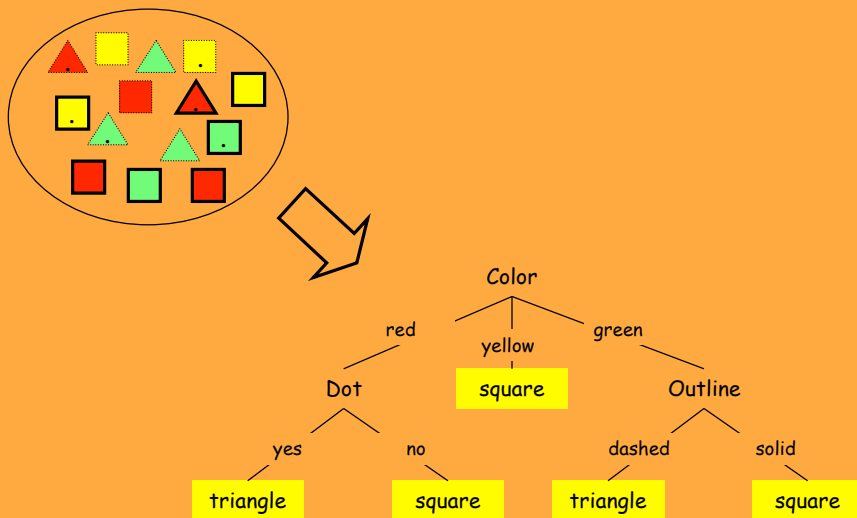
Information Gain of The Attribute

- Attributes
 - $\text{Gain}(\text{Color}) = 0.246$
 - $\text{Gain}(\text{Outline}) = 0.151$
 - $\text{Gain}(\text{Dot}) = 0.048$
- Heuristics: attribute with the highest gain is chosen
- This heuristics is local (local minimization of impurity)





Decision Tree



A Defect of *Ires*

- *Ires* favors attributes with many values
- Such attribute splits S to many subsets, and if these are small, they will tend to be pure anyway
- One way to rectify this is through a corrected measure of **information gain ratio**.

Information Gain Ratio

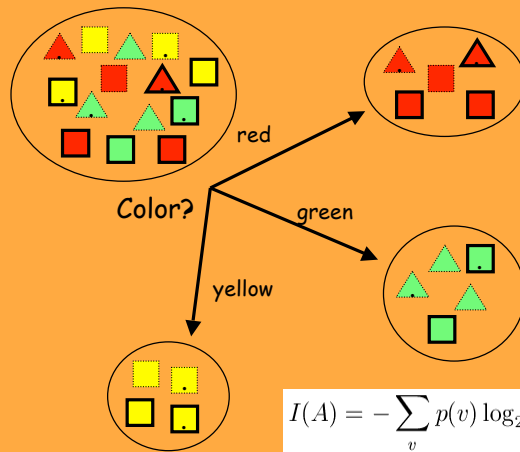
- $I(A)$ is amount of information needed to determine the value of an attribute A

$$I(A) = - \sum_v p(v) \log_2(p(v))$$

- Information gain ratio

$$\text{GainRatio}(A) = \frac{\text{Gain}(A)}{I(A)} = \frac{I - I_{res}(A)}{I(A)}$$

Information Gain Ratio



$$I(A) = - \sum_v p(v) \log_2(p(v))$$

$$I(\text{Color}) = -\frac{5}{14} \log_2 \frac{5}{14} - \frac{5}{14} \log_2 \frac{5}{14} - \frac{4}{14} \log_2 \frac{4}{14} = 1.58 \text{ bits}$$

$$\text{GainRatio}(\text{Color}) = \frac{\text{Gain}(\text{Color})}{I(\text{Color})} = \frac{0.940 - 0.694}{1.58} = 0.156$$

Information Gain and Information Gain Ratio

A	v(A)	Gain(A)	GainRatio(A)
Color	3	0.247	0.156
Outline	2	0.152	0.152
Dot	2	0.048	0.049

Gini Index

- Another sensible measure of impurity (i and j are classes)

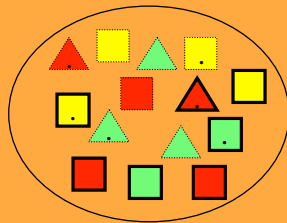
$$Gini = \sum_{i \neq j} p(i)p(j)$$

- After applying attribute A, the resulting Gini index is

$$Gini(A) = \sum_v p(v) \sum_{i \neq j} p(i|v)p(j|v)$$

- Gini can be interpreted as expected error rate

Gini Index



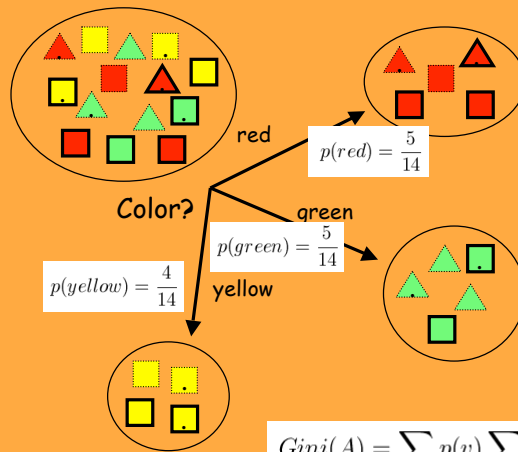
$$p(\square) = \frac{9}{14}$$

$$p(\triangle) = \frac{5}{14}$$

$$Gini = \sum_{i \neq j} p(i)p(j)$$

$$Gini = \frac{9}{14} \times \frac{5}{14} = 0.230$$

Gini Index for Color



$$Gini(A) = \sum_v p(v) \sum_{i \neq j} p(i|v)p(j|v)$$

$$Gini(\text{Color}) = \frac{5}{14} \times \left(\frac{3}{5} \times \frac{2}{5}\right) + \frac{5}{14} \times \left(\frac{2}{5} \times \frac{3}{5}\right) + \frac{4}{14} \times \left(\frac{4}{4} \times \frac{0}{4}\right) = 0.171$$

Gain of Gini Index

$$Gini = \frac{9}{14} \times \frac{5}{14} = 0.230$$

$$Gini(\text{Color}) = \frac{5}{14} \times \left(\frac{3}{5} \times \frac{2}{5}\right) + \frac{5}{14} \times \left(\frac{2}{5} \times \frac{3}{5}\right) + \frac{4}{14} \times \left(\frac{4}{4} \times \frac{0}{4}\right) = 0.171$$

$$GiniGain(\text{Color}) = 0.230 - 0.171 = 0.058$$

Three Impurity Measures

A	Gain(A)	GainRatio(A)	GiniGain(A)
Color	0.247	0.156	0.058
Outline	0.152	0.152	0.046
Dot	0.048	0.049	0.015

- These impurity measures assess the effect of a single attribute
- Criterion "most informative" that they define is local (and "myopic")
- It does not reliably predict the effect of several attributes applied jointly

Orange: Shapes Data Set

shape.tab

Color	Outline	Dot	Shape
d	d	d	d
			class
green	dashed	no	triange
green	dashed	yes	triange
yellow	dashed	no	square
red	dashed	no	square
red	solid	no	square
red	solid	yes	triange
green	solid	no	square
green	dashed	no	triange
yellow	solid	yes	square
red	solid	no	square
green	solid	yes	square
yellow	dashed	yes	square
yellow	solid	no	square

Orange: Impurity Measures

```
import orange
data = orange.ExampleTable('shape')

gain = orange.MeasureAttribute_info
gainRatio = orange.MeasureAttribute_gainRatio
gini = orange.MeasureAttribute_gini

print
print "%15s %-8s %-8s %-8s" % ("name", "gain", "g ratio", "gini")
for attr in data.domain.attributes:
    print "%15s %4.3f %4.3f %4.3f" % \
        (attr.name, gain(attr, data), gainRatio(attr, data), gini(attr, data))
```

name	gain	g ratio	gini
Color	0.247	0.156	0.058
Outline	0.152	0.152	0.046
Dot	0.048	0.049	0.015

Orange: orngTree

```
import orange, orngTree
data = orange.ExampleTable('shape')

tree = orngTree.TreeLearner(data)
orngTree.printTxt(tree)

print '\nWith contingency vector:'
orngTree.printTxt(tree,
    internalNodeFields=['contingency'],
    leafFields=['contingency'])
```

```
Color green:
| Outline dashed: triange (100.0%)
| Outline solid: square (100.0%)
Color yellow: square (100.0%)
Color red:
| Dot no: square (100.0%)
| Dot yes: triange (100.0%)

With contingency vector:

Color (<5, 9>) green:
| Outline (<3, 2>) dashed: triange (<3, 0>)
| Outline (<3, 2>) solid: square (<0, 2>)
Color (<5, 9>) yellow: square (<0, 4>)
Color (<5, 9>) red:
| Dot (<2, 3>) no: square (<0, 3>)
| Dot (<2, 3>) yes: triange (<2, 0>)
```

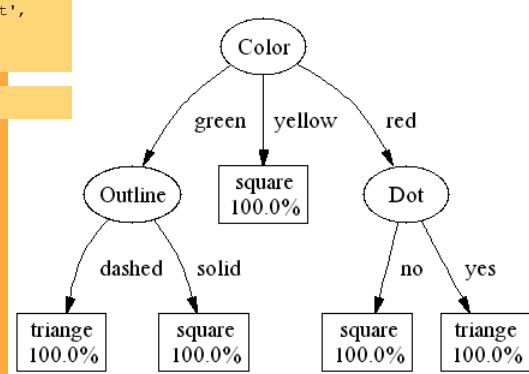
Orange: Saving to DOT

```
import orange, orngTree
data = orange.ExampleTable('shape')

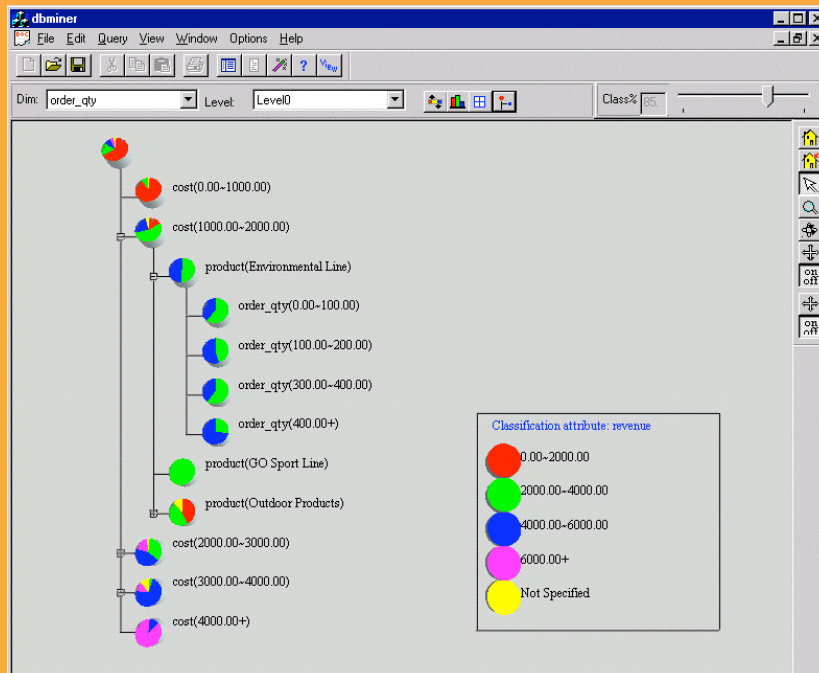
tree = orngTree.TreeLearner(data)

orngTree.printDot(tree, 'shape.dot',
                  leafShape='box',
                  internalNodeShape='ellipse')

> dot -Tgif shape.dot > shape.gif
```



DB Miner: visualization



SGI MineSet: visualization

