

# Decision Tree Models



# Decision Tree Model

Decision trees are especially attractive for a data mining environment for three reasons.

- Due to their **intuitive** representation, they are easy to assimilate by humans.
- They can be constructed **relatively fast** compared to other methods.
- The **accuracy** of decision tree classifiers is comparable or superior to other models.

# Decision Tree Learning

# Building Decision Trees

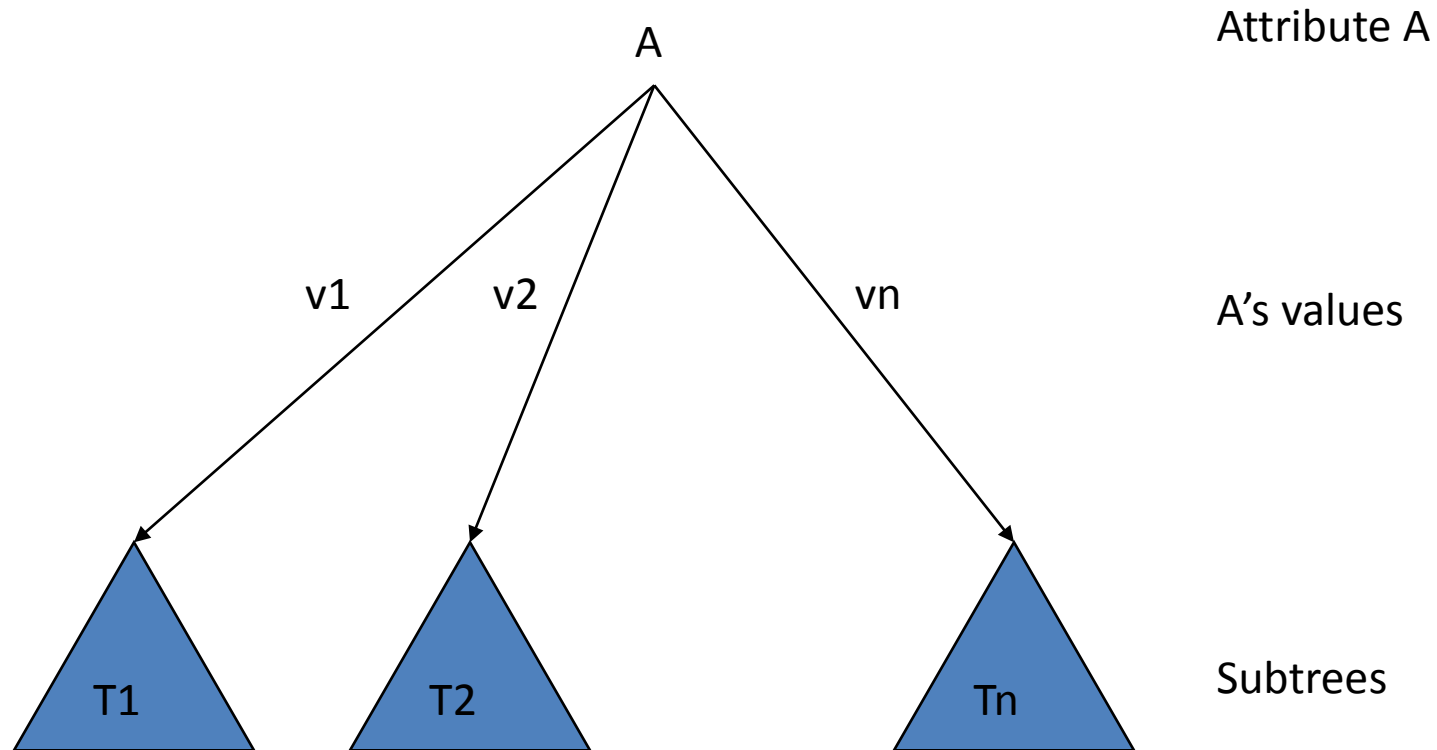
Two steps for the decision tree learning

- (1) Growing phase → maximizing the purity of the leaves
- (2) (Post) pruning phase → minimizing the "true" error rate

\* Pruning is an additional step to avoid the over-dependence to the growing sample(i.e.,overfitting)

# Recursive partitioning

Resulting tree T is:



# Issues in Decision Tree Learning

## Choosing the splitting criterion

- Impurity based criteria
- Information gain
- Statistical measures of association...

## Binary or multiway splits

- Multiway split: 1 value of the splitting attribute = 1 leaf
- Binary split: finding the best binary grouping
- Grouping only the leaves which are similar regarding the classes distribution

## Finding the right sized tree

- Pre-pruning
- Post-pruning

# Splitting Criteria

# I. Splitting Criterion

- **Central Idea** : Select attribute which partitions the learning set into subsets as “pure” as possible
- A partition is PURE if all of the observations in it belong to the same class.

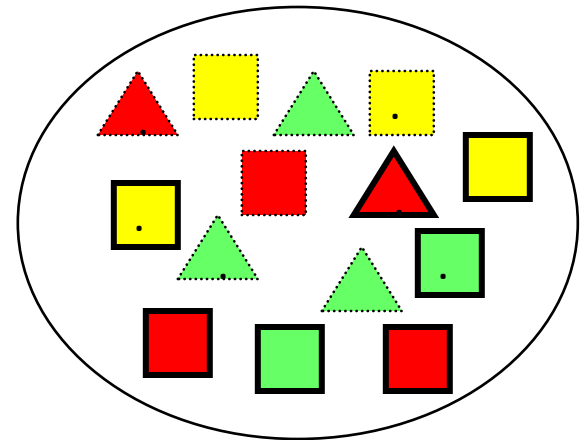


# Example: Triangles and Squares

#	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

Data Set:

A set of classified objects



# I. Entropy & Information Gain – C4.5

## Shannon entropy

Measure of uncertainty

$$E(Y) = - \sum_{k=1}^K \frac{n_{k.}}{n} \times \log_2 \left( \frac{n_{k.}}{n} \right)$$

## Condition entropy

Expected entropy of Y knowing the values of X

$$E(Y / X) = - \sum_{l=1}^L \frac{n_{.l}}{n} \sum_{k=1}^K \frac{n_{kl}}{n_{.l}} \times \log_2 \left( \frac{n_{kl}}{n_{.l}} \right)$$

## Information gain

Reduction of uncertainty

$$G(Y / X) = E(Y) - E(Y / X)$$

## (Information) Gain ratio

Favors the splits with low number of leaves

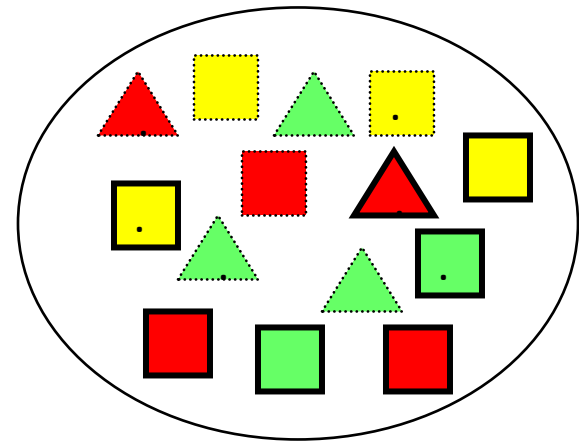
$$GR(Y / X) = \frac{E(Y) - E(Y / X)}{E(X)}$$

# Example: Entropy of given dataset

- 5 triangles
- 9 squares
- class probabilities

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

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

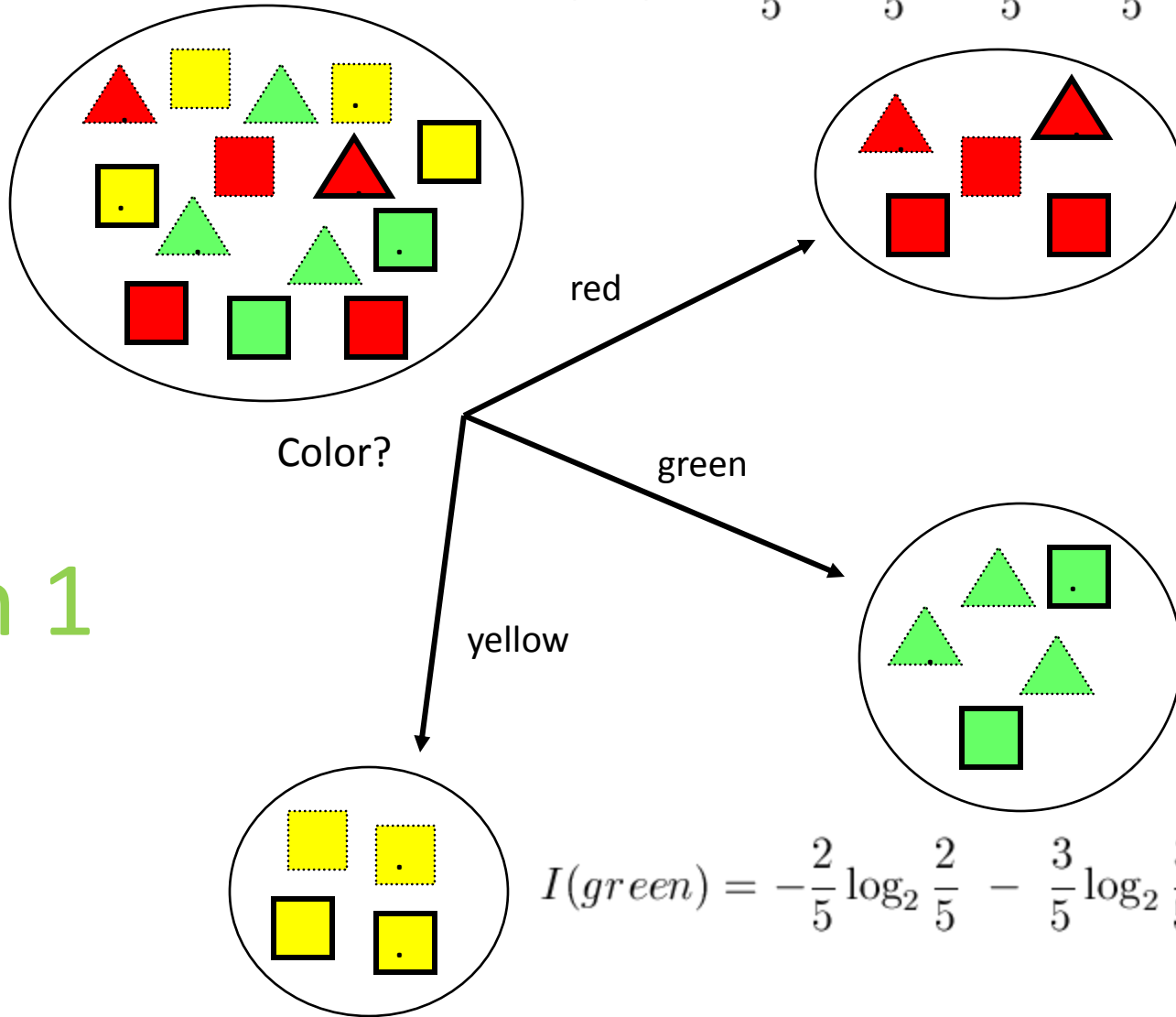


- entropy

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

$$I(\text{red}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.971 \text{ bits}$$

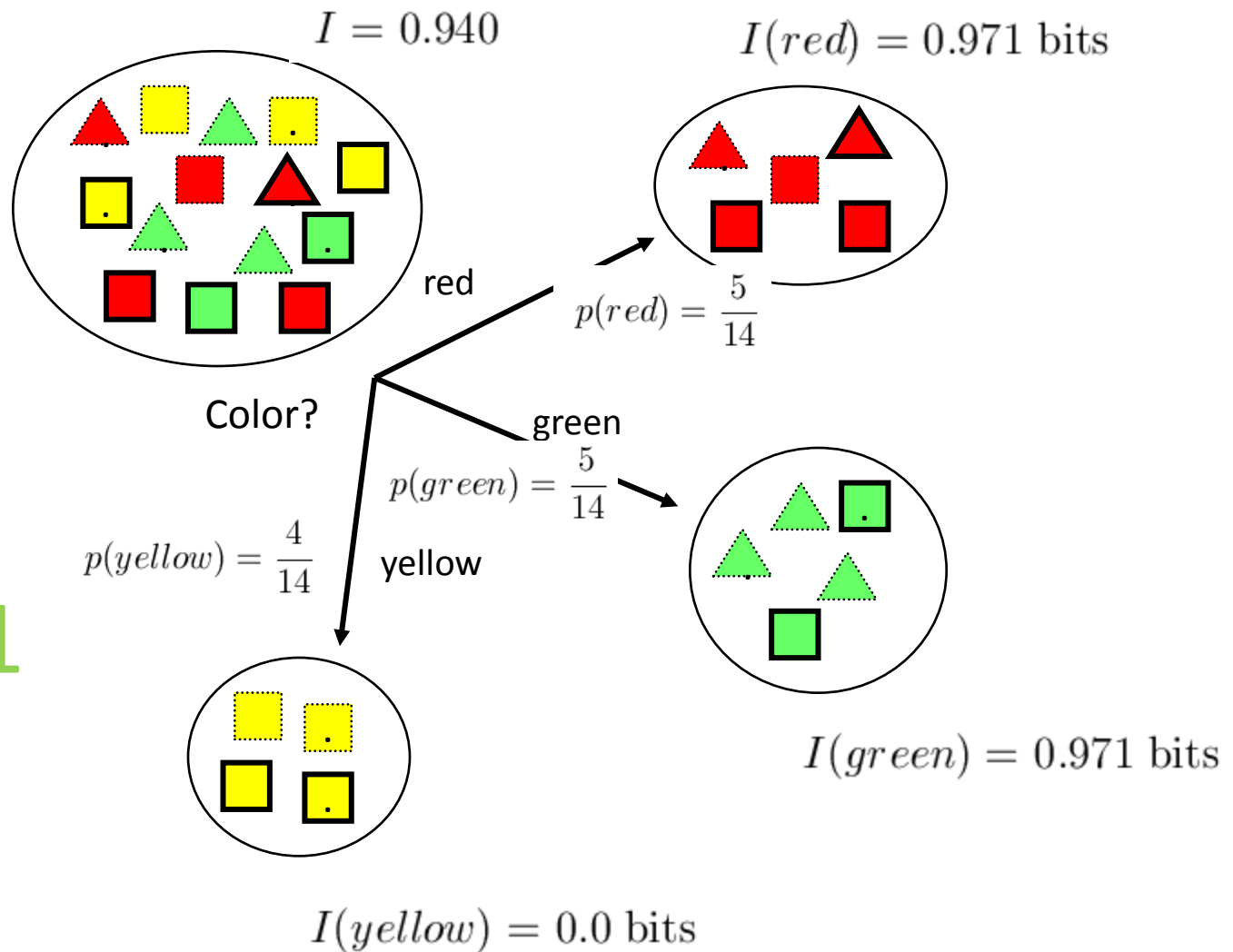
Depth 1



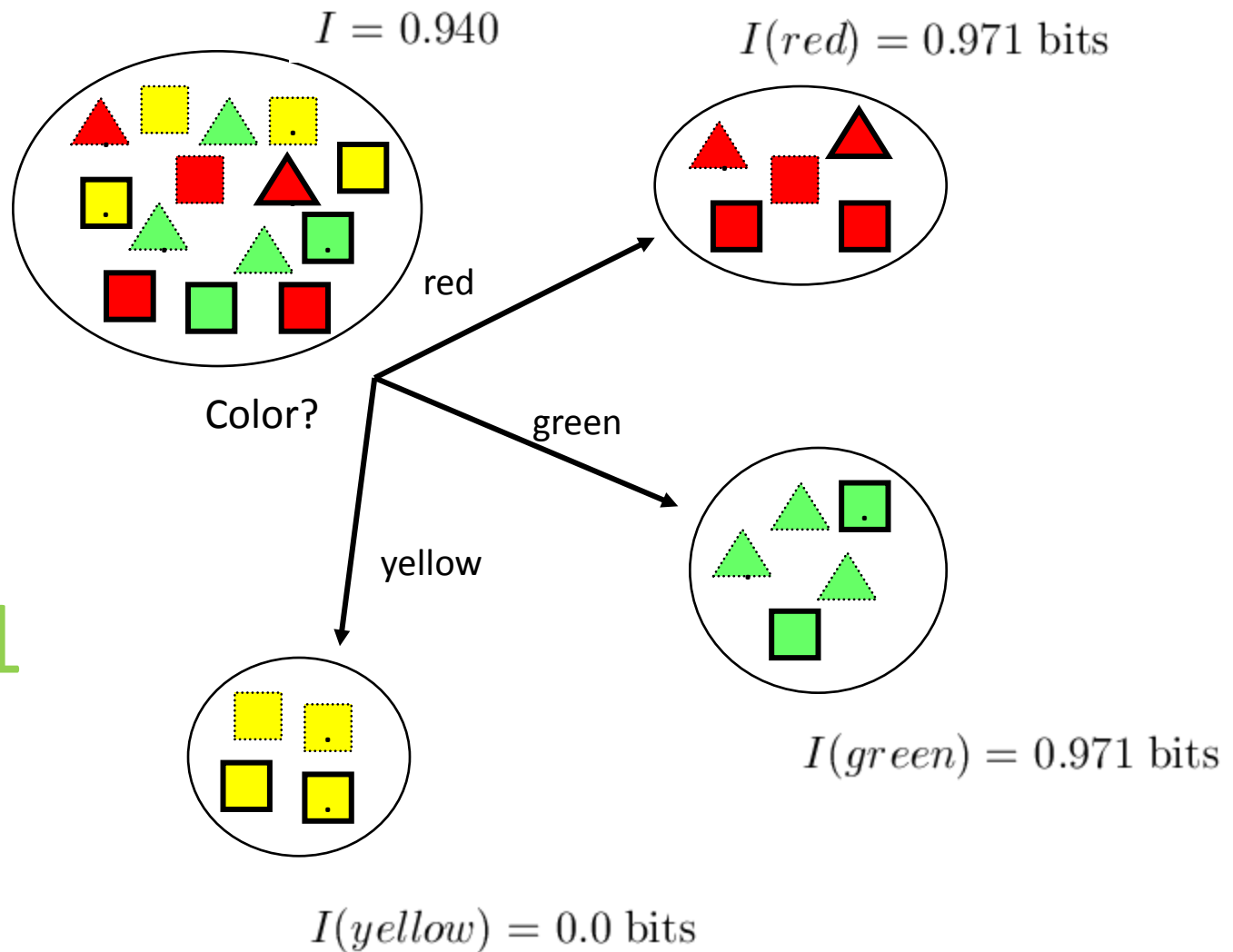
$$I(\text{green}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971 \text{ bits}$$

$$I(\text{yellow}) = -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} = 0.0 \text{ bits}$$

Depth 1



$$I_{res}(\text{Color}) = \sum p(v)I(v) = \frac{5}{14}0.971 + \frac{5}{14}0.971 + \frac{4}{14}0.0 = 0.694 \text{ bits}$$

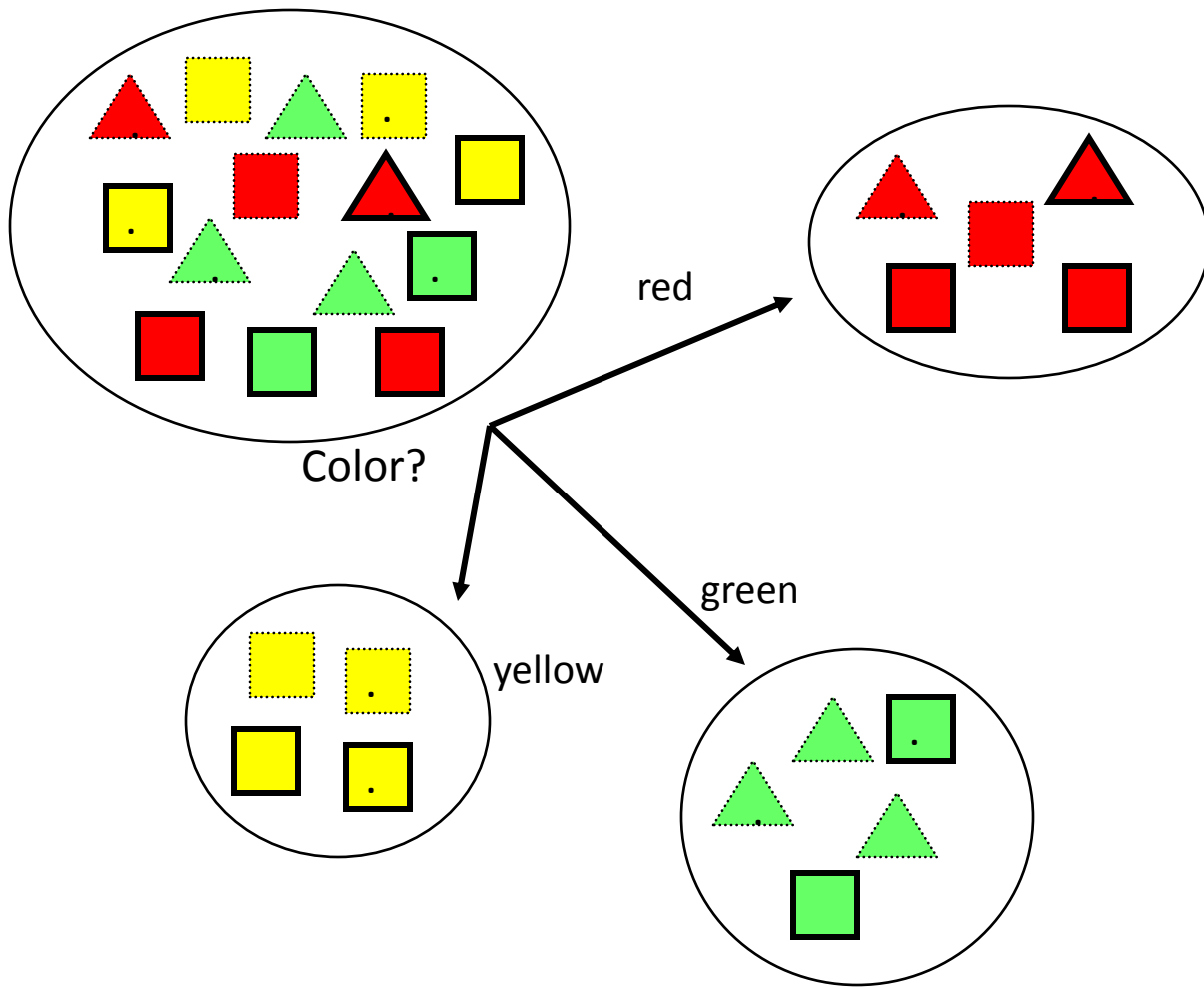


Depth 1

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

# Depth1 :Information Gains

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

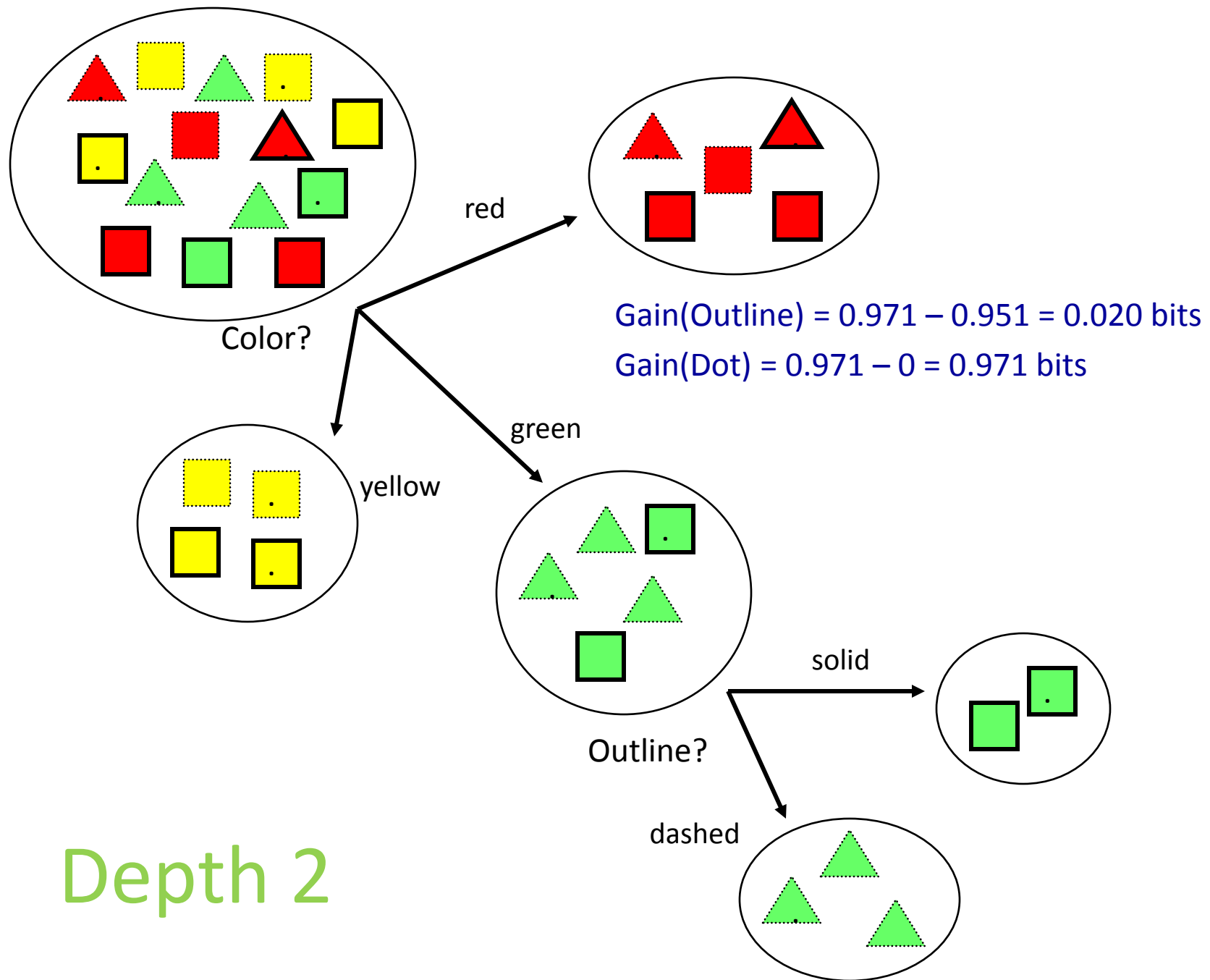


$$\text{Gain(Outline)} = 0.971 - 0 = 0.971 \text{ bits}$$

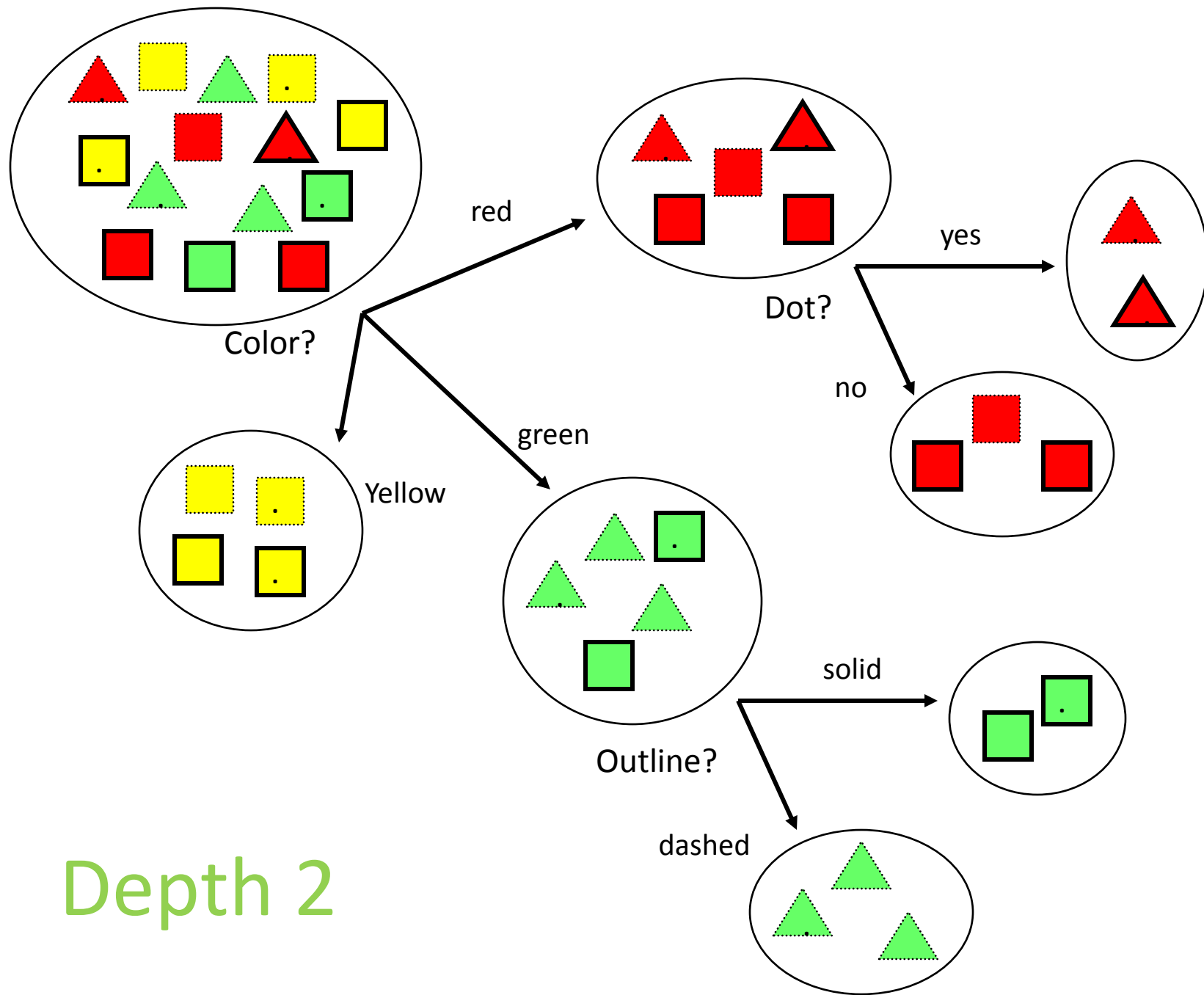
$$\text{Gain(Dot)} = 0.971 - 0.951 = 0.020 \text{ bits}$$

Depth 2



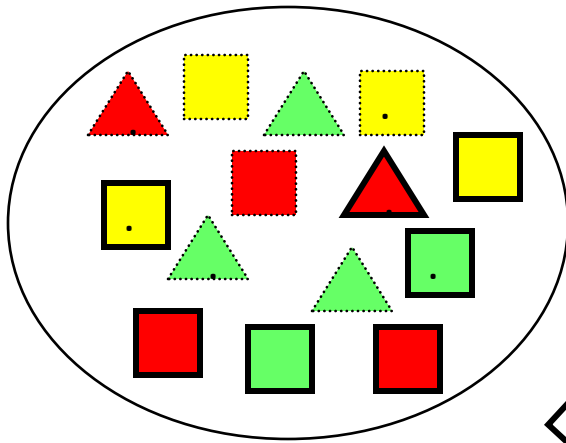


Depth 2

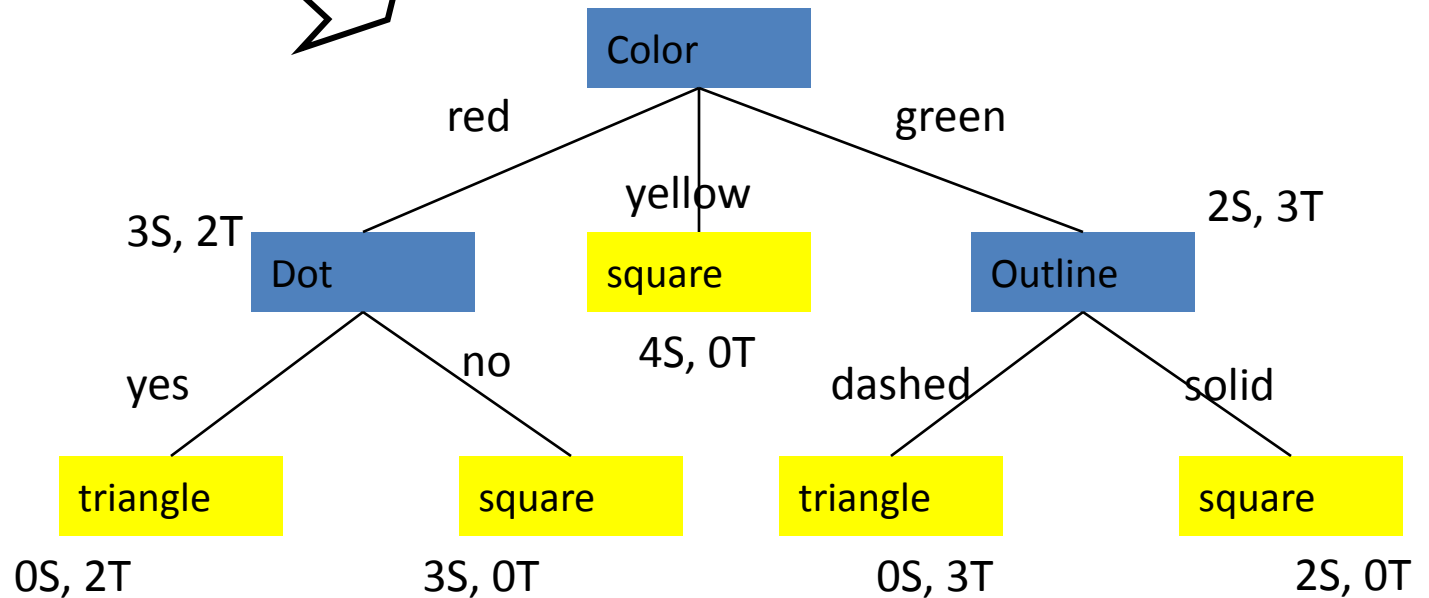


Depth 2

# Final Decision Tree



9S, 5T



# I. Gini Gain – CART

**Gini index**

Measure of impurity

$$I(Y) = - \sum_{k=1}^K \frac{n_{k.}}{n} \times \left( 1 - \frac{n_{k.}}{n} \right)$$

**Conditional impurity**

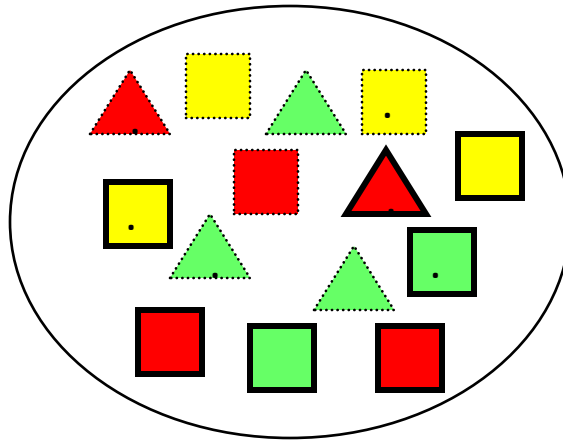
Average impurity of Y conditionally to X

$$I(Y / X) = - \sum_{l=1}^L \frac{n_{.l}}{n} \sum_{k=1}^K \frac{n_{kl}}{n_{.l}} \times \left( 1 - \frac{n_{kl}}{n_{.l}} \right)$$

**Gain**

$$D(Y / X) = I(Y) - I(Y / X)$$

# Gini Index



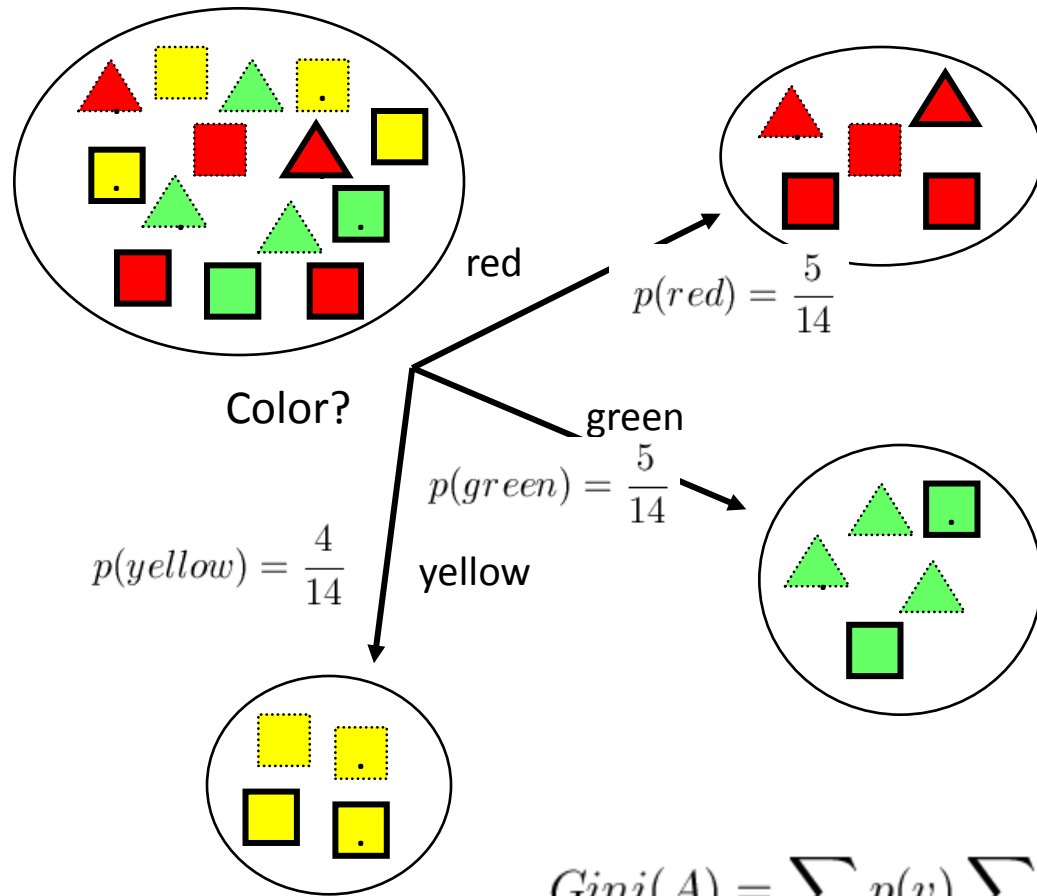
$$p(\square) = \frac{9}{14}$$

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

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

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

# 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

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

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

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

# I. Un-biased measures

- Allows to alleviate the data fragmentation problem
- Gain Ratio corrects the bias of the information gain
- The Gini reduction in impurity is biased in favor of variables with more levels  
(but the CART algorithm constructs necessarily a binary decision tree)

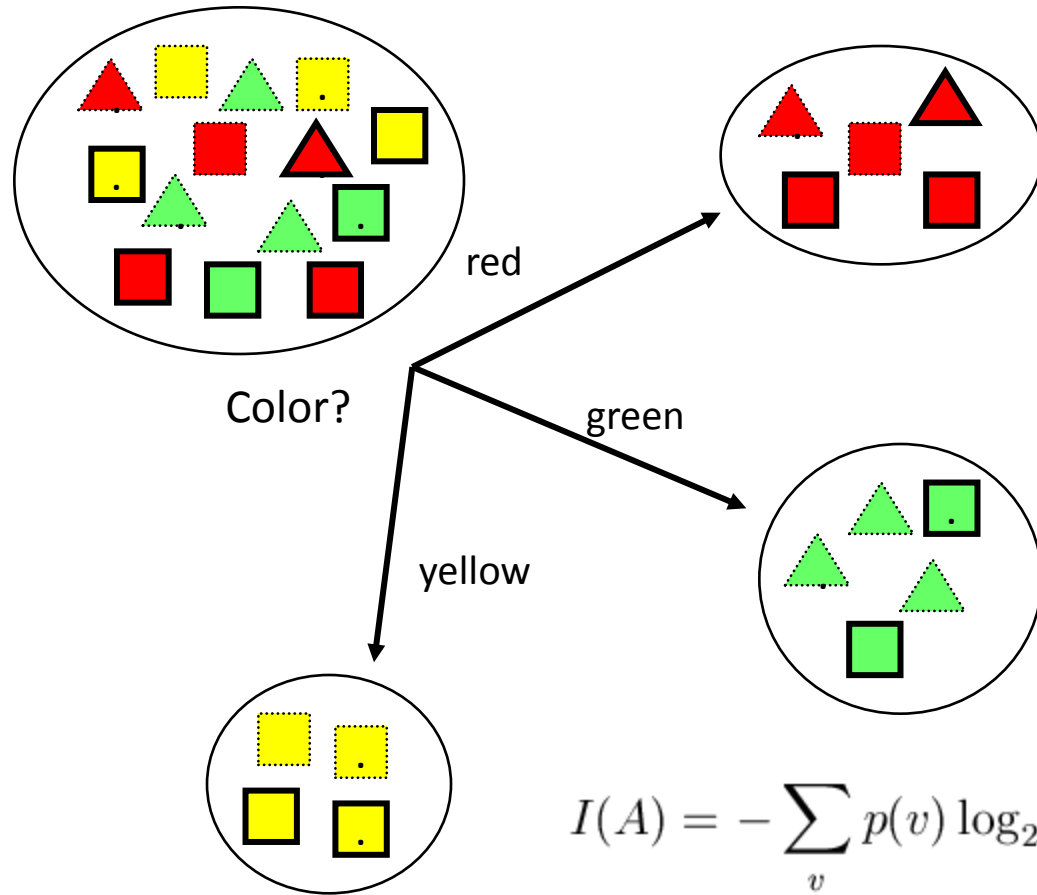


# Problems with Information Gain

Attributes which have a large number of possible values -> leads to **many child nodes**.

- Information gain is biased towards choosing attributes with a large number of values
- This may result in *overfitting* (selection of an attribute that is non-optimal for prediction)

# 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

<b>A</b>	<b><math> v(A) </math></b>	<b><i>Gain(A)</i></b>	<b><i>GainRatio(A)</i></b>
Color	3	0.247	0.156
Outline	2	0.152	0.152
Dot	2	0.048	0.049

# Three Impurity Measures

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

# Merging/Splitting Process

## II. Multi-way splitting in C4.5

- The prediction rules are easy to read
- Data fragmentation problem, especially for small dataset
- “Large” decision tree with a high number of leaves
- “Low depth” decision tree

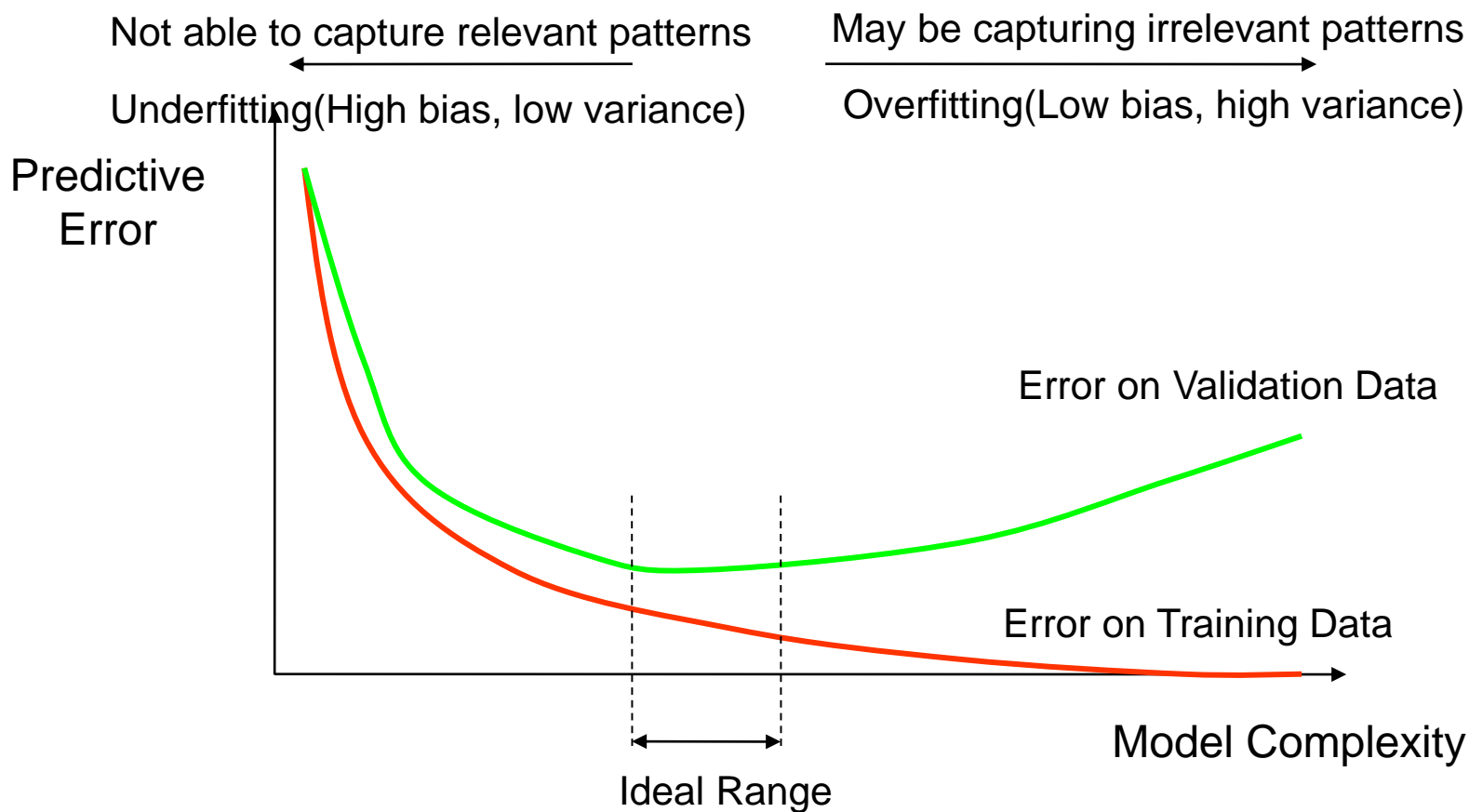
## II. Binary splitting in CART

- This grouping allows to overcome the bias of the splitting measure used
- The data fragmentation problem is alleviated
- “High depth” decision tree (CART uses a post pruning process for remedy to this)
- Merging into two groups is not always relevant !

Finding the Right sized tree



# III. Finding the right sized tree



**Bias** ~ how powerful is the model

**Variance** ~ how sensitive the model is to the training set

### III. Avoid overfitting

Pruning is an additional step in order to avoid the over-dependence to the growing sample (i.e., overfitting)

# III. Pre-pruning: Stopping the growing process

## Confidence and support criteria

- Group purity: confidence threshold
- Support criterion: min. size node to split, min. instances in leaves

# III. Postpruning in CART

- Split data into training & validation sets
- Build full decision tree from training set
- For every non-leaf node N
  - Prune subtree rooted by N, replace with majority class.
  - Test accuracy of pruned tree on validation set, that is, check if the pruned tree performs no worse than the original over the validation set.
  - Greedily remove the subtree that results in greatest improvement in accuracy on validation set

# Cost-Complexity Pruning

- Definition: Cost-Complexity Criterion

$$R_\alpha = MC + \alpha/T$$

- $MC$  = misclassification rate
  - Relative to # misclassifications in root node.
- $T$  = # leaves (terminal nodes)
- You get a credit for lower  $MC$ .
- But you *also* get a penalty for more leaves.
- Let  $T_0$  be the biggest tree.
- Find sub-tree of  $T_\alpha$  of  $T_0$  that minimizes  $R_\alpha$ .
  - Optimal trade-off of accuracy and complexity.

| If  $\alpha$  is small, the penalty for having a large number of terminal nodes is small and  $T(\alpha)$  tends to be large. | For  $\alpha$  sufficiently large, the minimizing subtree  $T(\alpha)$  will consist of the root node only

# What is the Optimal Size?

- Note that  $\alpha$  is a free parameter in:

$$R_{\alpha} = MC + \alpha L$$

- What value of  $\alpha$  should we choose?
  - $\alpha=0 \rightarrow$  maximum tree  $T_0$  is best.
  - $\alpha=\text{big} \rightarrow$  You never get past the root node.
  - Truth lies in the middle.
- Use cross-validation to select optimal  $\alpha$  (size)

### III. Postpruning in C4.5

- No Validation dataset for pruning
- Predicted (pessimistic) error rate  
= upper bound of the confidence  
limit of the resubstitution error  
rate

# III. Postpruning in C4.5

- Make estimates of the error based on the whole training data.
- Build full decision tree from training set
- For every non-leaf node N
  - The majority class is chosen to represent the node
  - Count the number of errors,  $\#e/N$  = error rate
  - Establish a confidence interval and use the upper limit, pessimistic estimate of the error rate
  - Compare such estimate with the combined estimate of the error estimates for the leaves
  - If the first is smaller replace the node by a leaf



# CART vs C4.5

Method	CART	C4.5
Splitting criterion	Gini index	Information Gain (Gain Ratio)
Merging process	Binary grouping	No merging 1 value = 1 leaf
Determining the right sized tree (overall)	Min. size to node split Min. instances in leaves Confidence threshold Tree depth	
Determining the right sized tree (specific)	Post-pruning Cost complexity pruning	Post-pruning Error based pruning
Recommended when...	Classification performance - Reliability No complicated settings	Small dataset Not sensitive to the settings

# Pros & Cons

# Advantages over other approaches

- No formal distributional assumptions
- Can automatically fit highly non-linear interactions
- Automatic variable selection
- Handle missing values
- Very easy to interpret if the tree is small
- The terminal nodes suggest a natural clustering
- Same tool for regression and classification
- Handle categorical predictors naturally
- Quick to fit, even for large problems

# Disadvantages

- *Accuracy* – newer methods can have lower error rates than CART & C4.5
- *Instability* – if we change the data a little, the tree picture can change a lot