# Classification:

# Decision Trees

# Outline

- Top-Down Decision Tree Construction

- Choosing the Splitting Attribute
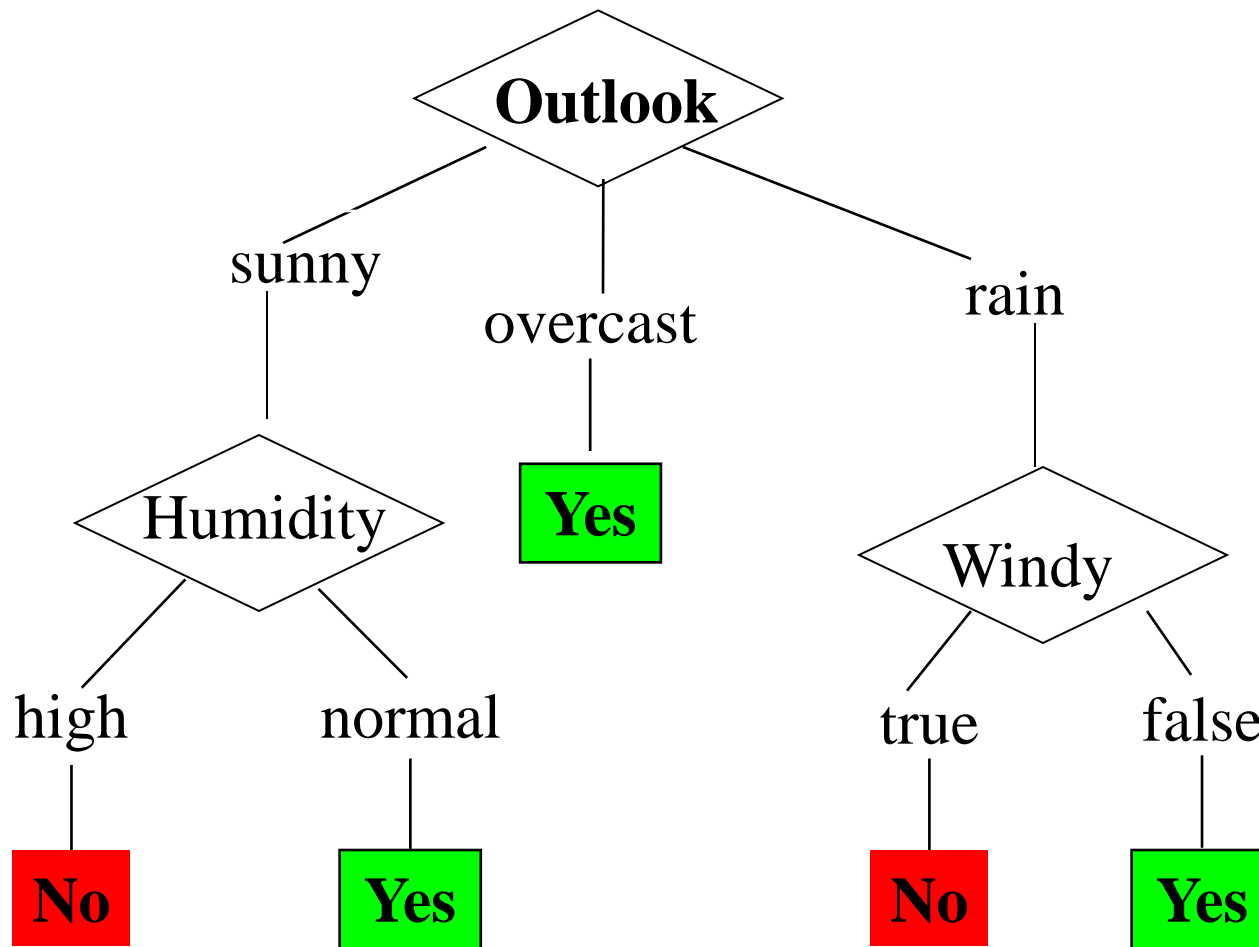
- Information Gain and Gain Ratio

- Gini index

# DECISION TREE

- An internal node is a test on an attribute.

- A branch represents an outcome of the test, e.g., Color=red.

- A leaf node represents a class label or class label distribution.

- At each node, one attribute is chosen to split training examples into distinct classes as much as possible

- A new case is classified by following a matching path to a leaf node.

# Weather Data: Play or not Play?

| Outlook | Temperature | Humidity | Windy | Play? |
|---------|-------------|----------|-------|-------|
| sunny | hot | high | false | No |
| sunny | hot | high | true | No |
| overcast | hot | high | false | Yes |
| rain | mild | high | false | Yes |
| rain | cool | normal | false | Yes |
| rain | cool | normal | true | No |
| overcast | cool | normal | true | Yes |
| sunny | mild | high | false | No |
| sunny | cool | normal | false | Yes |
| rain | mild | normal | false | Yes |
| sunny | mild | normal | true | Yes |
| overcast | mild | high | true | Yes |
| overcast | hot | normal | false | Yes |
| rain | mild | high | true | No |

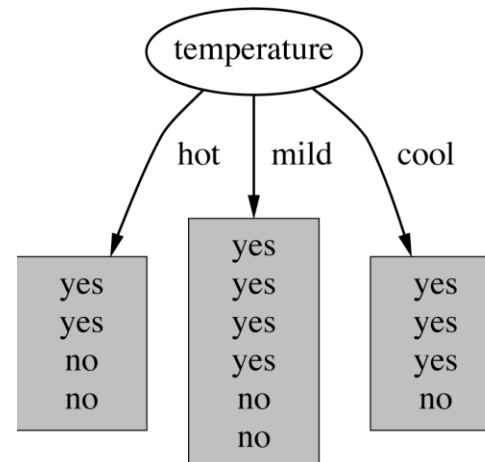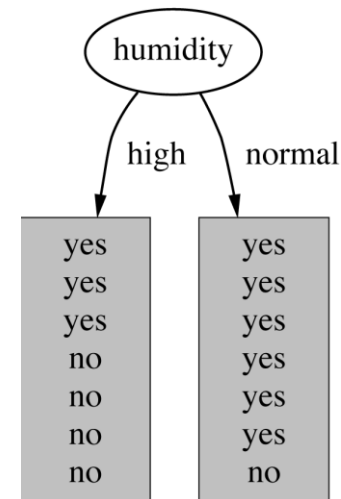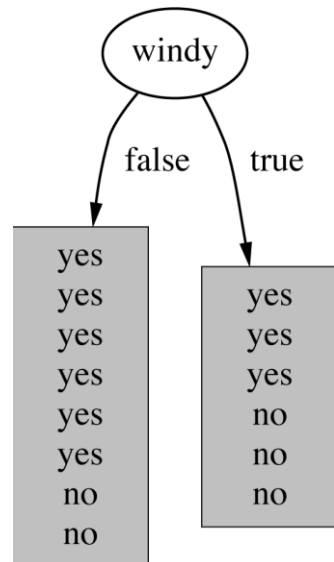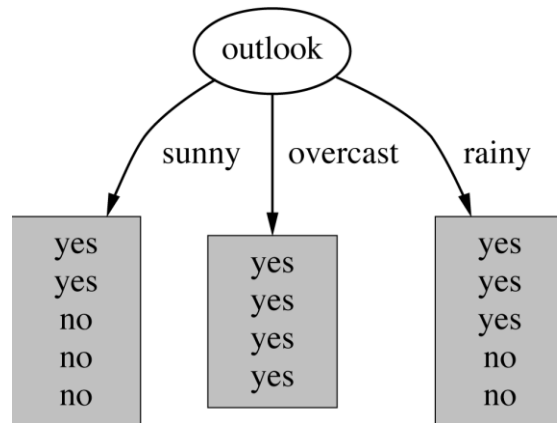# Example Tree for "Play?"

# Building Decision Tree [Q93]

- Top-down tree construction

  - At start, all training examples are at the root.

  - Partition the examples recursively by choosing one attribute each time.

- Bottom-up tree pruning

  - Remove subtrees or branches, in a bottom-up manner, to improve the estimated accuracy on new cases.

# Choosing the Splitting Attribute

- At each node, available attributes are evaluated on the basis of separating the classes of the training examples. A Goodness function is used for this purpose.

- Typical goodness functions:

  - information gain (ID3/C4.5)

  - information gain ratio

  - Gini index

# Which attribute to select?

# A criterion for attribute selection

- Which is the best attribute?

  - The one which will result in the smallest tree

  - Heuristic: choose the attribute that produces the "purest" nodes

- Popular *impurity criterion*: *information gain*

  - Information gain increases with the average purity of the subsets that an attribute produces

- Strategy: choose attribute that results in greatest information gain

# Computing information

- Information is measured in *bits*

  - Given a probability distribution, the info required to predict an event is the distribution's *entropy*

  - Entropy gives the information required in bits (this can involve fractions of bits!)

- Formula for computing the entropy:

$$\text{entropy}(p_1, p_2, ..., p_n) = -p_1 \log p_1 - p_2 \log p_2 ... - p_n \log p_n$$

# *Claude Shannon

*"Father of information theory"*

**Born: 30 April 1916**

**Died: 23 February 2001**

Claude Shannon, who has died aged 84, perhaps more than anyone laid the groundwork for today's digital revolution. His exposition of information theory, stating that all information could be represented mathematically as a succession of noughts and ones, facilitated the digital manipulation of data without which today's information society would be unthinkable.

Shannon's master's thesis, obtained in 1940 at MIT, demonstrated that problem solving could be achieved by manipulating the symbols 0 and 1 in a process that could be carried out automatically with electrical circuitry. That dissertation has been hailed as one of the most significant master's theses of the 20th century. Eight years later, Shannon published another landmark paper, *A Mathematical Theory of Communication*, generally taken as his most important scientific contribution.



RIP Claude Shannon © 2001 Stan Rowin

Claude Shannon © Stanley Rowin

Shannon applied the same radical approach to cryptography research, in which he later became a consultant to the US government.

Many of Shannon's pioneering insights were developed before they could be applied in practical form. He was truly a remarkable man, yet unknown to most of the world.

# Example: attribute "Outlook", 1

| Outlook | Temperature | Humidity | Windy | Play? |
|---------|-------------|----------|-------|-------|
| **sunny** | hot | high | false | No |
| **sunny** | hot | high | true | No |
| **overcast** | hot | high | false | Yes |
| **rain** | mild | high | false | Yes |
| **rain** | cool | normal | false | Yes |
| **rain** | cool | normal | true | No |
| **overcast** | cool | normal | true | Yes |
| **sunny** | mild | high | false | No |
| **sunny** | cool | normal | false | Yes |
| **rain** | mild | normal | false | Yes |
| **sunny** | mild | normal | true | Yes |
| **overcast** | mild | high | true | Yes |
| **overcast** | hot | normal | false | Yes |
| **rain** | mild | high | true | No |

# Example: attribute "Outlook", 2

- **"Outlook" = "Sunny":**

$$\text{info}([2,3]) = \text{entropy}(2/5, 3/5) = -2/5\log(2/5) - 3/5\log(3/5) = 0.971 \text{ bits}$$

- **"Outlook" = "Overcast":**

$$\text{info}([4,0]) = \text{entropy}(1,0) = -1\log(1) - 0\log(0) = 0 \text{ bits}$$

*Note: log(0) is not defined, but we evaluate 0\*log(0) as zero*

- **"Outlook" = "Rainy":**

$$\text{info}([3,2]) = \text{entropy}(3/5, 2/5) = -3/5\log(3/5) - 2/5\log(2/5) = 0.971 \text{ bits}$$

- **Expected information for attribute:**

$$\text{info}([3,2],[4,0],[3,2]) = (5/14) \times 0.971 + (4/14) \times 0 + (5/14) \times 0.971$$

$$= 0.693 \text{ bits}$$

# Computing the information gain

- Information gain:

(information before split) – (information after split)

$$\text{gain("Outlook")} = \text{info}([9,5]) - \text{info}([2,3],[4,0],[3,2]) = 0.940 - 0.693$$
$$= 0.247 \text{ bits}$$

- Compute for attribute "Humidity"

# Example: attribute "Humidity"

- **"Humidity" = "High":**

$$\text{info}([3,4]) = \text{entropy}(3/7, 4/7) = -3/7 \log(3/7) - 4/7 \log(4/7) = 0.985 \text{ bits}$$

- **"Humidity" = "Normal":**

$$\text{info}([6,1]) = \text{entropy}(6/7, 1/7) = -6/7 \log(6/7) - 1/7 \log(1/7) = 0.592 \text{ bits}$$

- **Expected information for attribute:**

$$\text{info}([3,4],[6,1]) = (7/14) \times 0.985 + (7/14) \times 0.592 \quad = 0.79 \text{ bits}$$

- **Information Gain:**

$$\text{info}([9,5]) - \text{info}([3,4],[6,1]) = 0.940 - 0.788 = 0.152$$

# Computing the information gain

- Information gain:

(information before split) − (information after split)

$$\text{gain}(\text{"Outlook"}) = \text{info}([9,5]) - \text{info}([2,3],[4,0],[3,2]) = 0.940 - 0.693$$

$$= 0.247 \text{ bits}$$
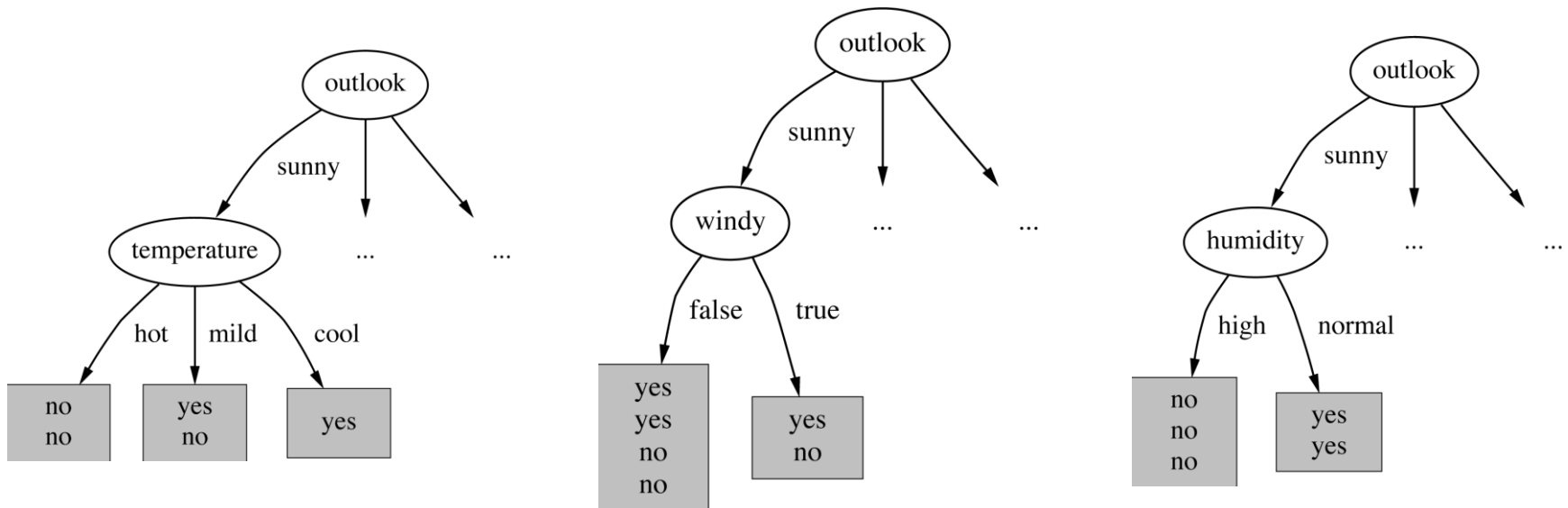
- Information gain for attributes from weather data:

$$\text{gain}(\text{"Outlook"}) = 0.247 \text{ bits}$$

$$\text{gain}(\text{"Temperature"}) = 0.029 \text{ bits}$$

$$\text{gain}(\text{"Humidity"}) = 0.152 \text{ bits}$$

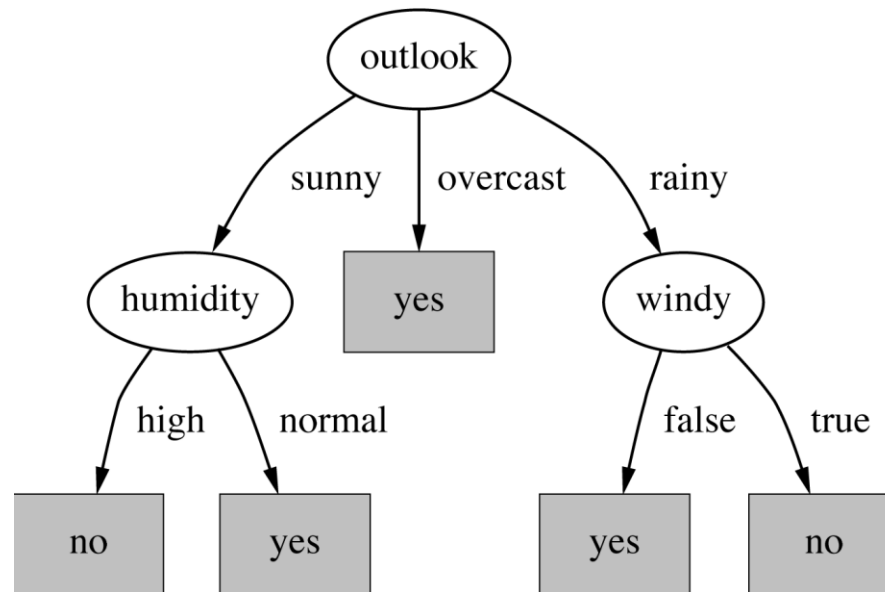$$\text{gain}(\text{"Windy"}) = 0.048 \text{ bits}$$

# Continuing to split



$$gain("Humidity") = 0.971 bits$$

$$gain("Temperature") = 0.571 bits$$

$$gain("Windy") = 0.020 bits$$

# The final decision tree



- Note: not all leaves need to be pure; sometimes identical instances have different classes

  ⇒ Splitting stops when data can't be split any further

# *Wish list for a purity measure

- Properties we require from a purity measure:

  - When node is pure, measure should be zero

  - When impurity is maximal (i.e. all classes equally likely), measure should be maximal

  - Measure should obey *multistage property* (i.e. decisions can be made in several stages):

  $$\text{measure}([2,3,4]) = \text{measure}([2,7]) + (7/9) \times \text{measure}([3,4])$$

- Entropy is a function that satisfies all three properties!

# *Properties of the entropy

- The multistage property:

$$\text{entropy}(p,q,r) = \text{entropy}(p,q+r) + \frac{(q+r)}{(p+q+r)} \times \text{entropy}(\frac{q}{q+r}, \frac{r}{q+r})$$

- Simplification of computation:

$$\text{info}([2,3,4]) = -2/9 \times \log(2/9) - 3/9 \times \log(3/9) - 4/9 \times \log(4/9)$$

$$= [-2\log 2 - 3\log 3 - 4\log 4 + 9\log 9]/9$$

- Note: instead of maximizing info gain we could just minimize information
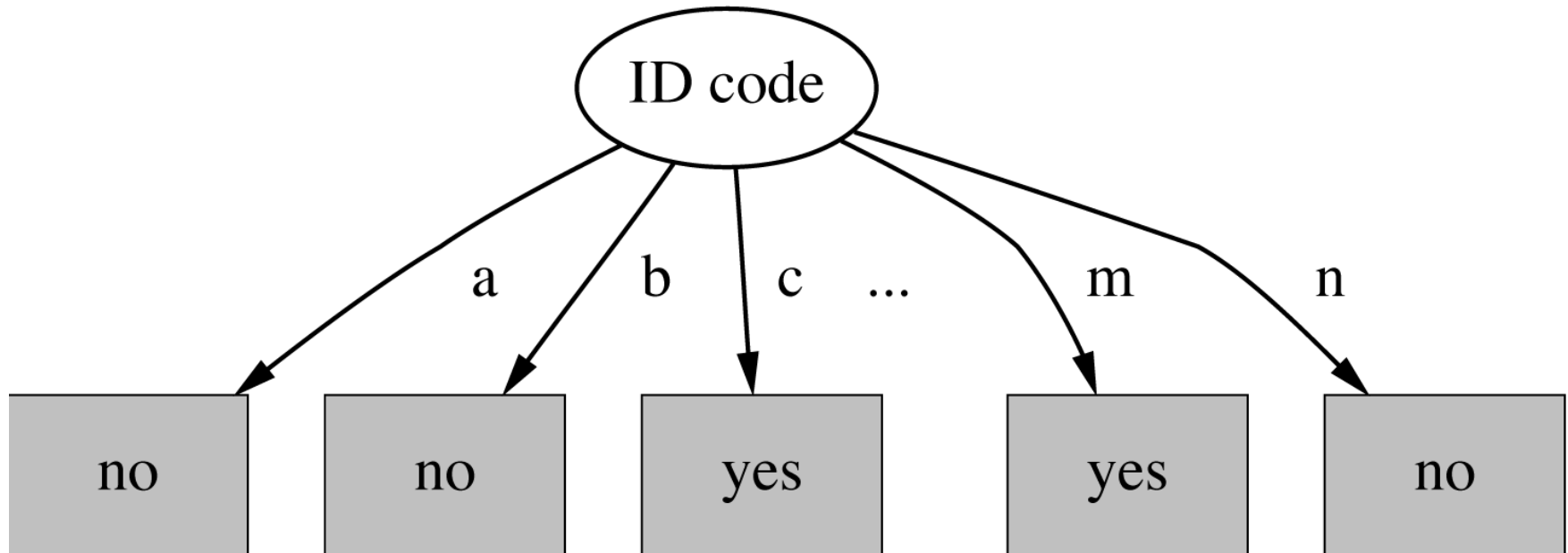
# Highly-branching attributes

- Problematic: attributes with a large number of values (extreme case: ID code)

- Subsets are more likely to be pure if there is a large number of values

  ⇒ 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)

# Weather Data with ID code

| ID | Outlook | Temperature | Humidity | Windy | Play? |
|---|---|---|---|---|---|
| **A** | sunny | hot | high | false | No |
| **B** | sunny | hot | high | true | No |
| **C** | overcast | hot | high | false | Yes |
| **D** | rain | mild | high | false | Yes |
| **E** | rain | cool | normal | false | Yes |
| **F** | rain | cool | normal | true | No |
| **G** | overcast | cool | normal | true | Yes |
| **H** | sunny | mild | high | false | No |
| **I** | sunny | cool | normal | false | Yes |
| **J** | rain | mild | normal | false | Yes |
| **K** | sunny | mild | normal | true | Yes |
| **L** | overcast | mild | high | true | Yes |
| **M** | overcast | hot | normal | false | Yes |
| **N** | rain | mild | high | true | No |

# Split for ID Code Attribute



Entropy of split = 0
(since each leaf node is "pure", having only one case.

Information gain is maximal for ID code

# Gain ratio

- *Gain ratio*: a modification of the information gain that reduces its bias on high-branch attributes

- Gain ratio should be
  - Large when data is evenly spread
  - Small when all data belong to one branch

- Gain ratio takes number and size of branches into account when choosing an attribute

  - It corrects the information gain by taking the *intrinsic information* of a split into account (i.e. how much info do we need to tell which branch an instance belongs to)

# Gain Ratio and Intrinsic Info.

- Intrinsic information: entropy of distribution of instances into branches

$$IntrinsicInfo(S,A) \equiv -\sum_i \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}.$$

- *Gain ratio* (Quinlan'86) normalizes info gain by:

$$GainRatio(S,A) = \frac{Gain(S,A)}{IntrinsicInfo(S,A)}.$$

# Computing the gain ratio

- Example: intrinsic information for ID code

$$\text{info}([1,1,\ldots,1) = 14 \times (-1/14 \times \log 1/14) = 3.807 \text{ bits}$$

- **Importance of attribute decreases as intrinsic information gets larger**

- Example of gain ratio:

$$\text{gain\_ratio}("\text{Attribute}") = \frac{\text{gain}("\text{Attribute}")}{\text{intrinsic\_info}("\text{Attribute}")}$$

- Example:
$$\text{gain\_ratio}("\text{ID\_code}") = \frac{0.940 \text{ bits}}{3.807 \text{ bits}} = 0.246$$

# Gain ratios for weather data

| Outlook | | | Temperature | | |
|---|---|---|---|---|---|
| Info: | 0.693 | | Info: | 0.911 | |
| Gain: 0.940-0.693 | 0.247 | | Gain: 0.940-0.911 | 0.029 | |
| Split info: info([5,4,5]) | 1.577 | | Split info: info([4,6,4]) | 1.362 | |
| Gain ratio: 0.247/1.577 | 0.156 | | Gain ratio: 0.029/1.362 | 0.021 | |

| Humidity | | | Windy | | |
|---|---|---|---|---|---|
| Info: | 0.788 | | Info: | 0.892 | |
| Gain: 0.940-0.788 | 0.152 | | Gain: 0.940-0.892 | 0.048 | |
| Split info: info([7,7]) | 1.000 | | Split info: info([8,6]) | 0.985 | |
| Gain ratio: 0.152/1 | 0.152 | | Gain ratio: 0.048/0.985 | 0.049 | |

# More on the gain ratio

- "Outlook" still comes out top

- However: "ID code" has greater gain ratio

  - Standard fix: *ad hoc* test to prevent splitting on that type of attribute

- Problem with gain ratio: it may overcompensate

  - May choose an attribute just because its intrinsic information is very low

  - Standard fix:

    - First, only consider attributes with greater than average information gain

    - Then, compare them on gain ratio

# *CART Splitting Criteria: Gini Index

- If a data set T contains examples from n classes, gini index, gini(T) is defined as

$$gini(T) = 1 - \sum_{j=1}^{n} p_j^2$$

  where $p_j$ is the relative frequency of class j in T.

  gini(T) is minimized if the classes in T are skewed.

# *Gini Index

- After splitting T into two subsets T1 and T2 with sizes N1 and N2, the gini index of the split data is defined as

$$gini_{split}(T) = \frac{N_1}{N} gini(T_1) + \frac{N_2}{N} gini(T_2)$$

- The attribute providing smallest $gini_{split}(T)$ is chosen to split the node.

# Discussion

- Algorithm for top-down induction of decision trees ("ID3") was developed by Ross Quinlan

    - Gain ratio just one modification of this basic algorithm

    - Led to development of C4.5, which can deal with numeric attributes, missing values, and noisy data

- Similar approach: CART

- There are many other attribute selection criteria! (But almost no difference in accuracy of result.)

# Summary

- Top-Down Decision Tree Construction

- Choosing the Splitting Attribute

- Information Gain biased towards attributes with a large number of values

- Gain Ratio takes number and size of branches into account when choosing an attribute