# CS412 Office Hours

April 5, 2019

# Decision Tree Induction: Algorithm

❑ Basic algorithm

  ❑ Tree is constructed in a **top-down, recursive, divide-and-conquer manner**

  ❑ At start, all the training examples are at the root

  ❑ Examples are partitioned recursively based on selected attributes

  ❑ On each node, attributes are selected based on the training examples on that node, and a heuristic or statistical measure (e.g., **information gain**)

❑ Conditions for stopping partitioning

  ❑ All samples for a given node belong to the same class

  ❑ There are no remaining attributes for further partitioning

  ❑ There are no samples left

❑ Prediction

  ❑ **Majority voting** is employed for classifying the leaf

Key issue: How to select the best attribute to partition on?

# Three Attribute Selection Measures

❑ **Information gain**: $Gain(A) = Info(D) - Info_A(D)$

  ❑ difference between entropy of the response variable H(Y) and the conditional entropy H(Y|A)

  ❑ biased towards multivalued attributes

❑ **Gain ratio**: $SplitInfo_A(D) = -\sum_{j=1}^{v} \frac{|D_j|}{|D|} \times \log_2(\frac{|D_j|}{|D|})$

  ❑ normalized information gain

  ❑ tends to prefer unbalanced splits in which one partition is much smaller than the others

❑ **Gini index**: $\Delta gini(A) = gini(D) - gini_A(D)$

  ❑ is biased to multivalued attributes; has difficulty when # of classes is large

  ❑ tends to favor tests that result in equal-sized partitions and purity in both partitions

# Information Gain: An Attribute Selection Measure

❑ Select the attribute with the highest information gain (used in typical decision tree induction algorithm: ID3)

❑ Let $p_i$ be the probability that an arbitrary tuple in D belongs to class $C_i$, estimated by $|C_{i, D}|/|D|$

❑ Expected information (entropy) needed to classify a tuple in D:

$$Info(D) = -\sum_{i=1}^{m} p_i \log_2(p_i)$$

❑ Information needed (after using A to split D into v partitions) to classify D:

$$Info_A(D) = \sum_{j=1}^{v} \frac{|D_j|}{|D|} \times Info(D_j)$$

❑ Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

# Example: Attribute Selection with Information Gain

❑ Class P: buys_computer = "yes"
❑ Class N: buys_computer = "no"

$$Info(D) = I(9,5) = -\frac{9}{14}\log_2(\frac{9}{14}) - \frac{5}{14}\log_2(\frac{5}{14}) = 0.940$$

| age | $p_i$ | $n_i$ | I($p_i$, $n_i$) |
|-----|-------|-------|-----------------|
| <=30 | 2 | 3 | 0.971 |
| 31…40 | 4 | 0 | 0 |
| >40 | 3 | 2 | 0.971 |

| age | income | student | credit_rating | buys_computer |
|-----|--------|---------|---------------|---------------|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | yes |
| 31…40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

$$Info_{age}(D) = \frac{5}{14}I(2,3) + \frac{4}{14}I(4,0)$$

$$+ \frac{5}{14}I(3,2) = 0.694$$

$\frac{5}{14}I(2,3)$ means "age <=30" has 5 out of 14 samples, with 2 yes'es and 3 no's. Hence

$$Gain(age) = Info(D) - Info_{age}(D) = 0.246$$

Similarly, we can get

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit\_rating) = 0.048$$

# Gain Ratio: A Refined Measure for Attribute Selection

❑ Information gain measure is biased towards attributes with a large number of values

❑ Gain ratio: Overcomes the problem (as a normalization to information gain)

$$SplitInfo_A(D) = -\sum_{j=1}^{v} \frac{|D_j|}{|D|} \times \log_2(\frac{|D_j|}{|D|})$$

   ❑ GainRatio(A) = Gain(A)/SplitInfo(A)

   ❑ The attribute with the maximum gain ratio is selected as the splitting attribute

   ❑ Gain ratio is used in a popular algorithm C4.5 (a successor of ID3) by R. Quinlan

❑ Example

   ❑ $SplitInfo_{income}(D) = -\frac{4}{14}\log_2\frac{4}{14} - \frac{6}{14}\log_2\frac{6}{14} - \frac{4}{14}\log_2\frac{4}{14} = 1.557$

   ❑ GainRatio(income) = 0.029/1.557 = 0.019

# Gini Index

❑ Gini index: Used in CART, and also in IBM IntelligentMiner

❑ If a data set $D$ contains examples from $n$ classes, gini index, $gini(D)$ is defined as

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

      ❑ $p_j$ is the relative frequency of class $j$ in $D$

❑ If a data set $D$ is split on $A$ into two subsets $D_1$ and $D_2$, the $gini$ index $gini(D)$ is defined as

   ❑ $gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$

❑ Reduction in Impurity:

   ❑ $\Delta gini(A) = gini(D) - gini_A(D)$

❑ The attribute which provides the smallest $gini_{split}(D)$ (or the largest reduction in impurity) is chosen to split the node (*need to enumerate all the possible splitting points for each attribute*)

# Computation of Gini Index

❑ Example:  D has 9 tuples in buys_computer = "yes" and 5 in "no"

$$gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

❑ Suppose the attribute income partitions D into 10 in $D_1$: {low, medium} and 4 in $D_2$

❑ $gini_{income \in \{low,\, medium\}}(D) = \frac{10}{14} gini(D_1) + \frac{4}{14} gini(D_2)$

$$= \frac{10}{14}\left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) + \frac{4}{14}\left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right) = 0.443$$

$$= Gini_{income \in \{high\}}(D)$$

❑ Gini$_{\{low,\, high\}}$ is 0.458; Gini$_{\{medium,\, high\}}$ is 0.450

❑ Thus, split on $income \in$ {low, medium} (i.e., also $\{high\}$) has the lowest Gini index

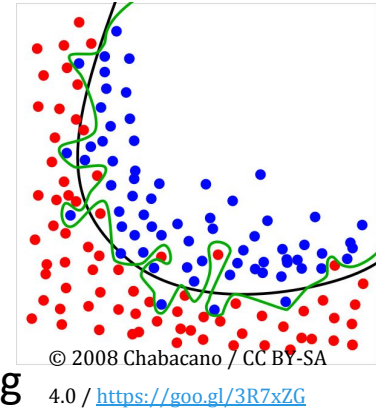| age | income | student | credit_rating | buys_computer |
|---|---|---|---|---|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | yes |
| 31…40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

# How to Handle Continuous-Valued Attributes?

❑ Method 1: Discretize continuous values and treat them as categorical values

  ❑ E.g., age: < 20, 20..30, 30..40, 40..50, > 50

❑ Method 2: Determine the **best split point** for continuous-valued attribute A

  ❑ Sort the value A in increasing order:, e.g., 15, 18, 21, 22, 24, 25, 29, 31, …

  ❑ *Possible split point:* the midpoint between *each pair of adjacent values*

    ❑ $(a_i + a_{i+1})/2$ is the midpoint between the values of $a_i$ and $a_{i+1}$

    ❑ e.g., (15+18)/2 = 16.5, 19.5, 21.5, 23, 24.5, 27, 30, …

  ❑ The point with the *maximum information gain* for A is selected as the **split-point** for A

❑ Split:  Based on split point P

  ❑ The set of tuples in D satisfying A ≤ P vs. those with A > P

# Tree Pruning

❑ Two approaches to avoid overfitting

   ❑ <u>Prepruning</u>: *Halt tree construction early*—do not split a node if this would result in the goodness measure falling below a threshold

      ❑ Difficult to choose an appropriate threshold

      ❑ You can experiment with different thresholds for the programming assignment

   ❑ <u>Postpruning</u>: *Remove branches* from a "fully grown" tree—get a sequence of progressively pruned trees

      ❑ Use a set of data different from the training data to decide which is the "best pruned tree"

# Using Decision Trees for Large Datasets: RainForest

- Our goal is to reduce the number of times we read/write from the disk where the database resides.
- RainForest decouples the "counting" process and the split criteria computation process.
- Database access only happens during the construction of the AVC-groups
- After we get the AVC-groups, we do not need to refer to the original data points anymore

# RainForest: A Scalable Classification Framework

❑ The criteria that determine the quality of the tree can be computed separately

  ❑ Builds an AVC-list: **AVC (Attribute, Value, Class_label)**

❑ **AVC-set** (of an attribute $X$ )

  ❑ Projection of training dataset onto the attribute $X$ and class label where counts of individual class label are aggregated

❑ **AVC-group** (of a node $n$ )

  ❑ Set of AVC-sets of all predictor attributes at node $n$

❑ Then read the data again & do it similarly to generate the next level of the tree

| age | income | student | credit_rating | com |
|------|--------|---------|---------------|------|
| <=30 | high | no | fair | no |
| <=30 | high | no | excellent | no |
| 31…40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31…40 | low | yes | excellent | yes |
| <=30 | medium | no | fair | no |
| <=30 | low | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| <=30 | medium | yes | excellent | yes |
| 31…40 | medium | no | excellent | yes |
| 31…40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

**The Training Data**

**AVC-set on *Age***

| Age | Buy_Computer | |
|------|-----|-----|
| | yes | no |
| <=30 | 2 | 3 |
| 31..40 | 4 | 0 |
| >40 | 3 | 2 |

**AVC-set on *Income***

| income | Buy_Computer | |
|--------|-----|-----|
| | yes | no |
| high | 2 | 2 |
| medium | 4 | 2 |
| low | 3 | 1 |

**AVC-set on *Student***

| student | Buy_Computer | |
|---------|-----|-----|
| | yes | no |
| yes | 6 | 1 |
| no | 3 | 4 |

**AVC-set on *Credit_Rating***

| Credit rating | Buy_Computer | |
|------|-----|-----|
| | yes | no |
| fair | 6 | 2 |
| excellent | 3 | 3 |

**Its AVC Sets**