# Data Mining

## DECISION TREE

**Prof. Dr. Hikmat Ullah Khan**
**Department of Information Technology**

UNIVERSITY OF SARGODHA, SARGODHA

# Lesson from Holy Quran

فَاصْبِرْ لِحُكْمِ رَبِّكَ وَلَا تُطِعْ مِنْهُمْ آثِمًا أَوْ كَفُورًا ٢٤

So be patient for the decision of your Lord and do not obey from among them a sinner or ungrateful (disbeliever)
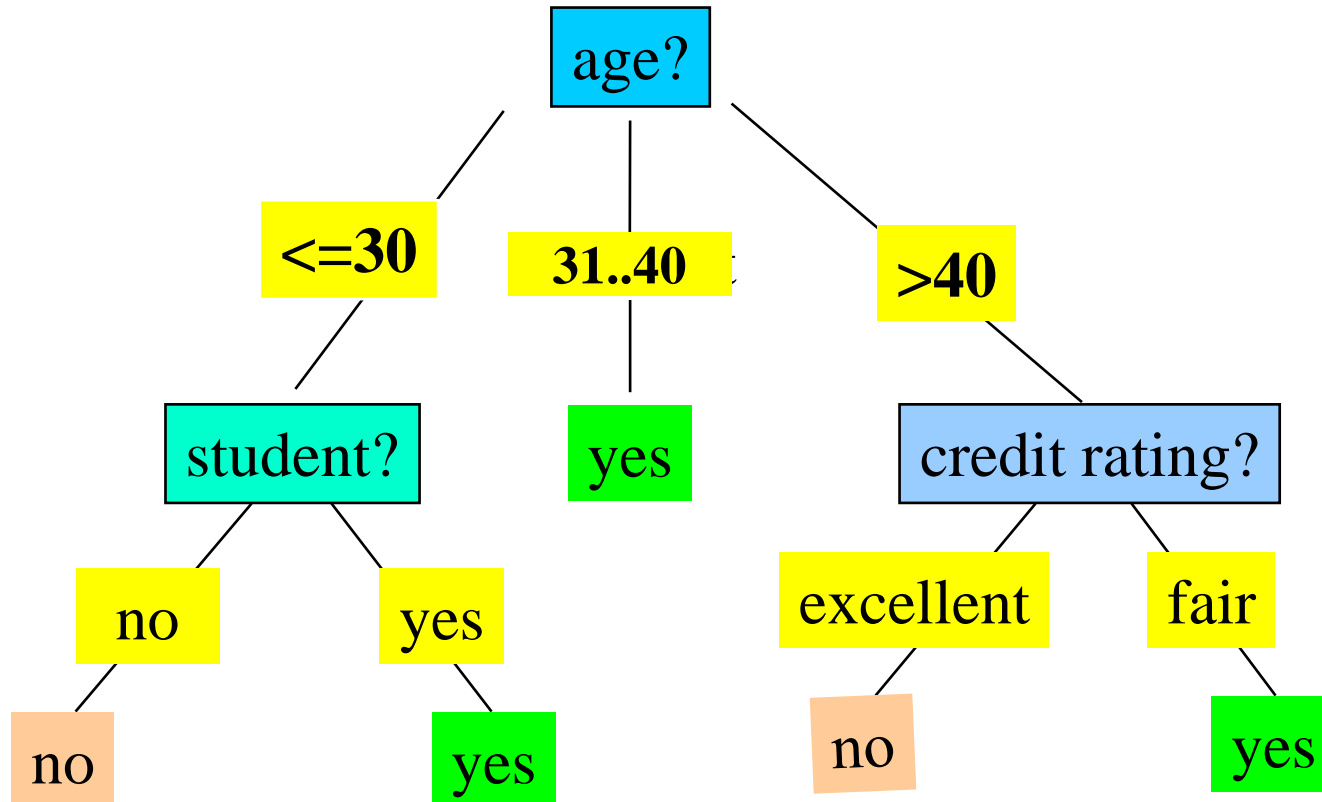
mwoch.tumblr.com

# Agenda

- Decision Tree
  - Introduction
  - Applying model
  - Properties
  - Attribute Selection
    - Information Gain
    - Entropy
    - Gini
    - Classification Error
  - Advantages/Dis-advantage

# Data to Decision Tree

- Training data set: Buys_computer

- Resulting tree:

```
              age?
         /     |      \
      <=30   31..40    >40
       /       |         \
   student?   yes    credit rating?
    /    \            /          \
   no    yes      excellent      fair
   /       \         |             \
  no       yes       no            yes
```

| 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 |

# Example of a Decision Tree

categorical   categorical   continuous   class

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

Training Data

*Splitting Attributes*

Refund
Yes → NO
No → MarSt

MarSt
Single, Divorced → TaxInc
Married → NO

TaxInc
< 80K → NO
> 80K → YES

Model:  Decision Tree

# Decision Tree Classification Task

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 1 | Yes | Large | 125K | No |
| 2 | No | Medium | 100K | No |
| 3 | No | Small | 70K | No |
| 4 | Yes | Medium | 120K | No |
| 5 | No | Large | 95K | Yes |
| 6 | No | Medium | 60K | No |
| 7 | Yes | Large | 220K | No |
| 8 | No | Small | 85K | Yes |
| 9 | No | Medium | 75K | No |
| 10 | No | Small | 90K | Yes |

Training Set

Tree Induction algorithm

Induction

Learn Model

Model

Decision Tree

Apply Model

Deduction

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

Test Set

# DEDUCTION

**generalization**

All birds fly.

**specific**

The cardinal is a bird.

The cardinal is a bird that flies.

**conclusion.**

# INDUCTION

**specific**

The cardinal is a bird that flies.

**specific**

The robin is a bird that flies.

All birds fly.

**generalization**

# Apply Model to Test Data

☐ Start from the root of tree.



Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

# Apply Model to Test Data

Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

# Apply Model to Test Data

Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|---------------|----------------|-------|
| No | Married | 80K | ? |

# Apply Model to Test Data

Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

# Apply Model to Test Data

Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

# Apply Model to Test Data

Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

Refund

Yes → NO

No → MarSt

Single, Divorced → TaxInc

Married → NO

TaxInc
< 80K → NO
> 80K → YES

Assign Cheat to "No"

# Properties of Decision Tree

- Basic algorithm (a greedy algorithm)
  - Tree is constructed in a
    - top-down
    - recursive
    - divide-and-conquer manner
  - At start, all the training examples are at the root
  - Attributes are categorical
    - (continuous are discretized in advance)
  - Test attributes are selected on the basis on statistical measure (e.g., information gain)
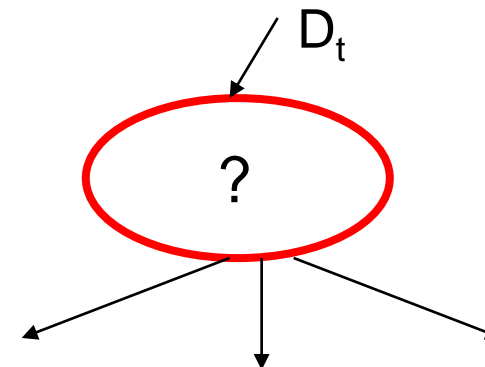
# Decision Tree Induction

- Many Algorithms:
  - Hunt's Algorithm (one of the earliest)
  - CART
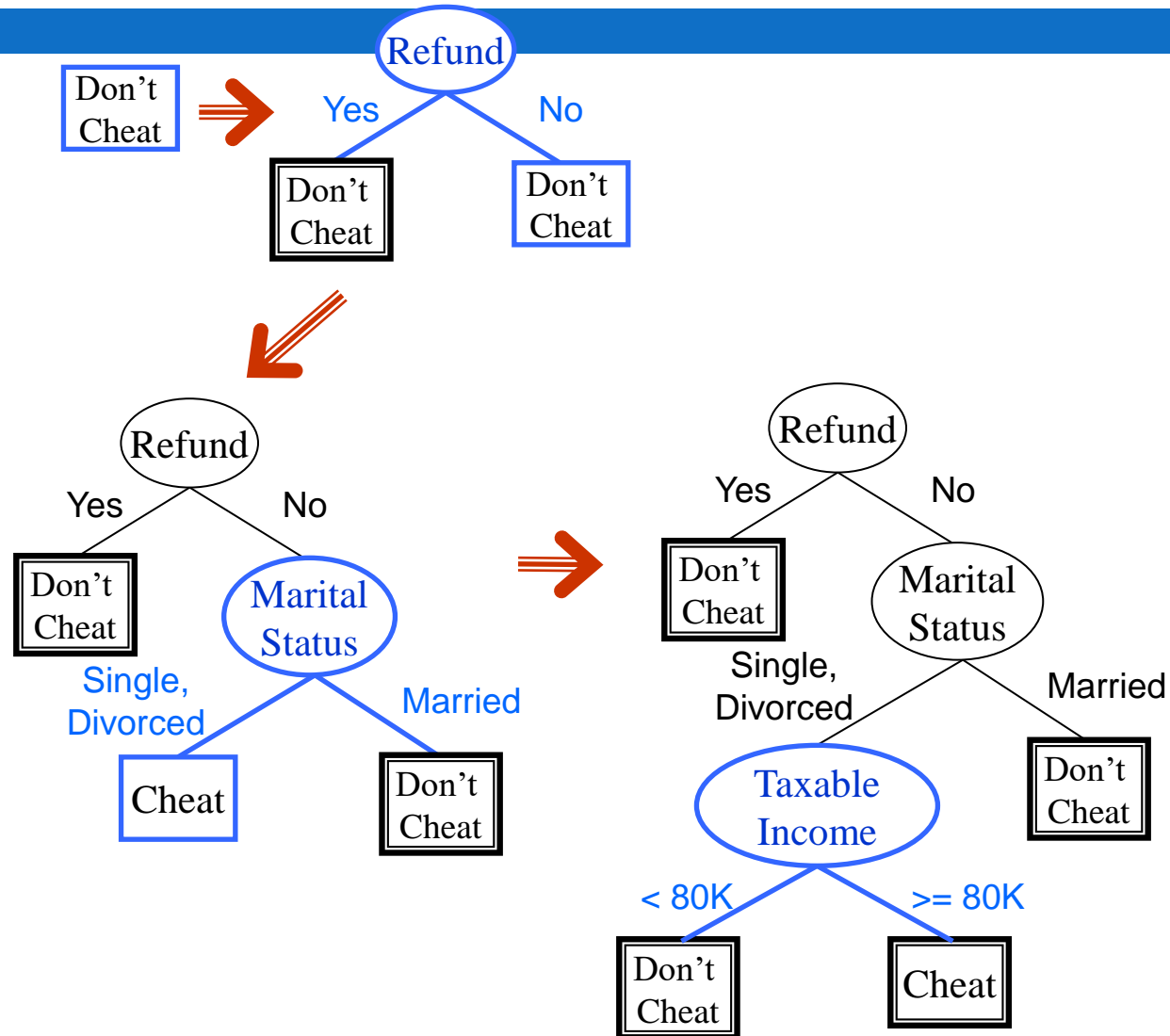  - ID3, C4.5, C5.0
  - SLIQ,SPRINT

# General Structure

- Let $D_t$ be the set of training records that reach a node t

- General Procedure:

  - If $D_t$ contains records that belong the <span style="color:red">same class</span> $y_t$, then t is a leaf node labeled as $y_t$

  - If $D_t$ contains records that belong to <span style="color:red">more than one class</span>, **use an attribute test to split the data into smaller subsets**. Recursively apply the procedure to each subset.

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

$D_t$

?

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

# Tree Induction

□ Greedy strategy.

  ◘ Split the records based on an attribute test that optimizes certain criterion.

□ Issues

  ◘ Determine how to split the records

    ■ How to specify the attribute test condition?

    ■ How to determine the best split?

  ◘ Determine when to stop splitting

# Tree Induction

- ☐ Greedy strategy.
  - ◘ Split the records based on an attribute test that optimizes certain criterion.

- ☐ Issues
  - ◘ Determine how to split the records
    - ■ How to specify the attribute test condition?
    - ■ How to determine the best split?
  - ◘ Determine when to stop splitting
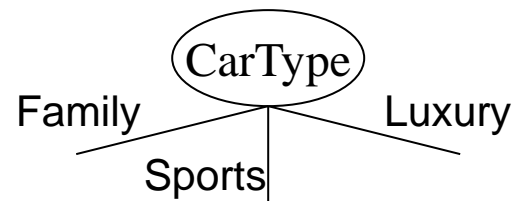
# How to Specify Test Condition?

- Depends on attribute types
  - Nominal
  - Ordinal
  - Continuous

- Depends on number of ways to split
  - 2-way split
  - Multi-way split

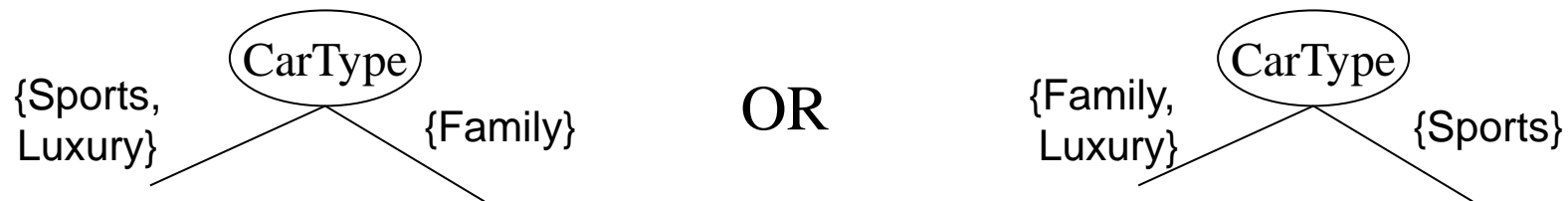# Splitting Based on Nominal Attributes

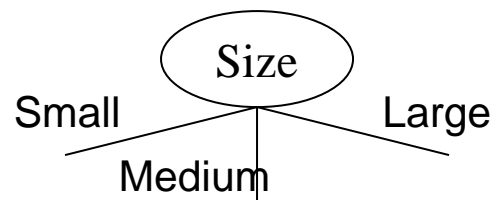☐ Multi-way split: Use as many partitions as distinct values.

```
            (CarType)
Family    /    |    \    Luxury
              Sports
```

☐ Binary split:  Divides values into two subsets.
                  Need to find optimal partitioning.

```
{Sports,  (CarType)              {Family,  (CarType)
Luxury}  /        \  {Family}    Luxury}  /        \  {Sports}
                         OR
```

# Splitting Based on Nominal Attributes

□ Multi-way split: Use as many partitions as distinct values.

Size
Small          Large
Medium

□ Binary split:  Divides values into two subsets.
                 Need to find optimal partitioning.

Size
{Small, Medium}          {Large}

OR

Size
{Medium, Large}          {Small}

□ What about this split?

Size
{Small, Large}          {Medium}

# Splitting Based on Continuous Attributes

Taxable Income > 80K?

Yes          No

Taxable Income?

< 10K          > 80K

[10K,25K)     [25K,50K)     [50K,80K)

(i) Binary split                    (ii) Multi-way split

# Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.

- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting

# Brief Review of Entropy

□ Entropy is the measure of uncertainty associated with a random measure

▫ High entropy ->  high uncertainty

▫ Low entropy -> low uncertainty

□ It is also known as measure of dispersion

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

# Attribute Selection Measure: Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- 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

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

# Examples for computing Entropy

$$Entropy(t) = -\sum_j p(j \mid t) \log_2 p(j \mid t)$$

| C1 | 0 |
|----|---|
| C2 | 6 |

P(C1) = 0/6 = 0     P(C2) = 6/6 = 1

Entropy = − 0 log 0 − 1 log 1 = − 0 − 0 = 0

| C1 | 1 |
|----|---|
| C2 | 5 |

| C1 | 2 |
|----|---|
| C2 | 4 |

# Examples for computing Entropy

$$Entropy(t) = -\sum_j p(j \mid t) \log_2 p(j \mid t)$$

| C1 | 0 |
|----|---|
| C2 | 6 |

P(C1) = 0/6 = 0     P(C2) = 6/6 = 1

Gini = 1 – P(C1)$^2$ – P(C2)$^2$ = 1 – 0 – 1 = 0

| C1 | 1 |
|----|---|
| C2 | 5 |

P(C1) = 1/6        P(C2) = 5/6

Entropy = – (1/6) log$_2$ (1/6) – (5/6) log$_2$ (5/6) = 0.65

| C1 | 2 |
|----|---|
| C2 | 4 |

P(C1) = 2/6        P(C2) = 4/6

Entropy = – (2/6) log$_2$ (2/6) – (4/6) log$_2$ (4/6) = 0.92

# Attribute Selection: 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 | 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 |

# Attribute Selection: 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 |

# Attribute Selection: 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 |

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

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

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

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

| 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 |

# Gini Index (CART, IBM IntelligentMiner)

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

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

where $p_j$ is the relative frequency of class *j* in *D*

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

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

# Examples for computing GINI

$$GINI(t) = 1 - \sum_j [p(j \mid t)]^2$$

| C1 | 0 |
|----|---|
| C2 | 6 |

P(C1) = 0/6 = 0     P(C2) = 6/6 = 1

Gini = 1 – P(C1)$^2$ – P(C2)$^2$ = 1 – 0 – 1 = 0

| C1 | 1 |
|----|---|
| C2 | 5 |

| C1 | 2 |
|----|---|
| C2 | 4 |

# Examples for computing GINI

$$GINI(t) = 1 - \sum_j [p(j \mid t)]^2$$

| | |
|---|---|
| C1 | **0** |
| C2 | **6** |

P(C1) = 0/6 = 0     P(C2) = 6/6 = 1

Gini = 1 – P(C1)$^2$ – P(C2)$^2$ = 1 – 0 – 1 = 0

| | |
|---|---|
| C1 | **1** |
| C2 | **5** |

P(C1) = 1/6        P(C2) = 5/6

Gini = 1 – (1/6)$^2$ – (5/6)$^2$ = 0.278

| | |
|---|---|
| C1 | **2** |
| C2 | **4** |

P(C1) = 2/6        P(C2) = 4/6

Gini = 1 – (2/6)$^2$ – (4/6)$^2$ = 0.444

# Third measure of Classification Error

$$\text{Entropy}(t) = -\sum_{i=0}^{c-1} p(i|t) \log_2 p(i|t),$$

$$\text{Gini}(t) = 1 - \sum_{i=0}^{c-1} [p(i|t)]^2,$$

$$\text{Classification error}(t) = 1 - \max_i [p(i|t)],$$

# Three in One

| Node $N_1$ | Count |
|---|---|
| Class=0 | 0 |
| Class=1 | 6 |

$$\text{Gini} = 1 - (0/6)^2 - (6/6)^2 = 0$$
$$\text{Entropy} = -(0/6)\log_2(0/6) - (6/6)\log_2(6/6) = 0$$
$$\text{Error} = 1 - \max[0/6, 6/6] = 0$$

| Node $N_2$ | Count |
|---|---|
| Class=0 | 1 |
| Class=1 | 5 |

$$\text{Gini} = 1 - (1/6)^2 - (5/6)^2 = 0.278$$
$$\text{Entropy} = -(1/6)\log_2(1/6) - (5/6)\log_2(5/6) = 0.650$$
$$\text{Error} = 1 - \max[1/6, 5/6] = 0.167$$

| Node $N_3$ | Count |
|---|---|
| Class=0 | 3 |
| Class=1 | 3 |

$$\text{Gini} = 1 - (3/6)^2 - (3/6)^2 = 0.5$$
$$\text{Entropy} = -(3/6)\log_2(3/6) - (3/6)\log_2(3/6) = 1$$
$$\text{Error} = 1 - \max[3/6, 3/6] = 0.5$$

# Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.

- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting

# Stopping Criteria for Tree Induction

- Stop expanding a node when all the records belong to the same class

# Decision Tree Based Classification

- Advantages:
  - Extremely fast at classifying unknown records
  - Easy to interpret for small-sized trees
  - Accuracy is comparable to other classification techniques for many simple data sets

- Disadvantages
  - Not scalable (add one attribute, all tree needed to be computed again)
  - Not good accuracy for large dataset
  - Not robust (less handling of large attributes)

# WEKA complete Book

- WEKA provides Wiki for all the concepts of Machine Learning and data mining

- https://www.cs.waikato.ac.nz/ml/weka/book.html

- WEKA examples for Decision Tree has been uploaded as reading material

# Examples

- **Search** "Data Mining Lecture -- Decision Tree | Solved Example (Eng-Hindi)"
- URL:
  - https://www.youtube.com/watch?v=cKl7WV_EKDU

  You All should solve the complete example of Weather data

  30 min video

  2 hour solution
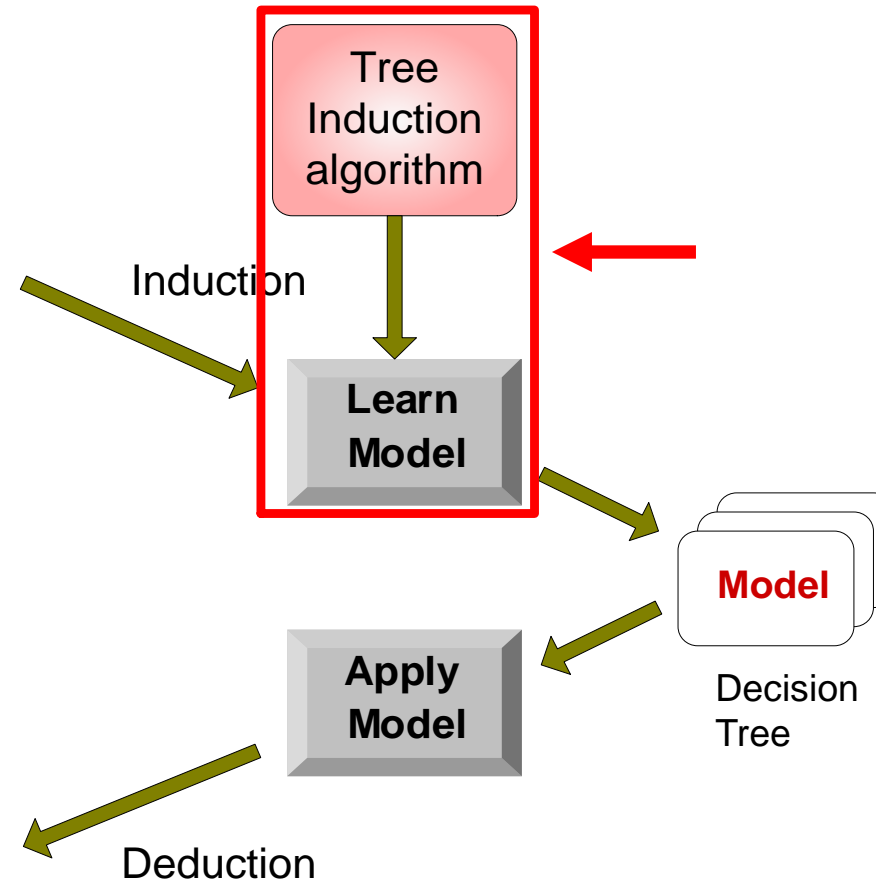
# DT Classification Task (optional)

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 1 | Yes | Large | 125K | No |
| 2 | No | Medium | 100K | No |
| 3 | No | Small | 70K | No |
| 4 | Yes | Medium | 120K | No |
| 5 | No | Large | 95K | Yes |
| 6 | No | Medium | 60K | No |
| 7 | Yes | Large | 220K | No |
| 8 | No | Small | 85K | Yes |
| 9 | No | Medium | 75K | No |
| 10 | No | Small | 90K | Yes |

Training Set

Tree Induction algorithm

Induction

Learn Model

Model

Decision Tree

Apply Model

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

Test Set

Deduction

# Every one can has DT in his mind for every task

The leader of a company needs to have a decision tree in his head - if this happens, we go this way, but if it winds up like that, then we go this other way.

— Sean Parker —

AZ QUOTES