# CS4104 Applied Machine Learning

Decision Tree Classifier

# Classification Techniques

- Decision Tree based Methods

- Rule-based Methods

- Memory based reasoning

- Neural Networks

- Naïve Bayes and Bayesian Belief Networks

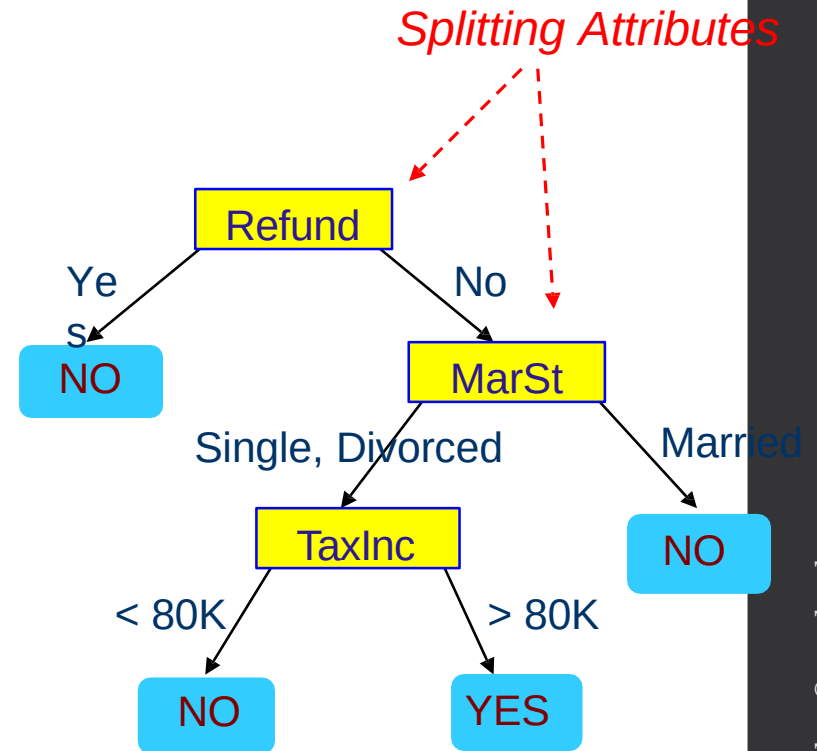- Support Vector Machines

# 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

Single, Divorced          Married

TaxInc          NO

< 80K          > 80K

NO          YES

Model:  Decision Tree

3

# Another Example of Decision Tree

| 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** |

*categorical*   *categorical*   *continuous*   *class*

MarSt

Married → NO

Single, Divorced → Refund

Refund: Yes → NO

Refund: No → TaxInc

TaxInc: < 80K → NO

TaxInc: > 80K → YES

There could be more than one tree that fits the same data!

zeshan.khan@nu.edu.pk

# 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

| 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

Apply
Model

Deduction

5

# Apply Model to Test Data

Start from the root of tree.

Test Data

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

Refund

Yes → NO

No → MarSt

MarSt
Single, Divorced → TaxInc
Married → NO

TaxInc
< 80K → NO
> 80K → YES

6

# Apply Model to Test Data

Test Data

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

Refund

Yes → NO

No → MarSt

MarSt:
- Single, Divorced → TaxInc
- Married → NO

TaxInc:
- < 80K → NO
- > 80K → YES

7

# Apply Model to Test Data

Test Data

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

Refund

Yes → NO

No → MarSt

MarSt

Single, Divorced → TaxInc

Married → NO

TaxInc

< 80K → NO

> 80K → YES

8

# 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

< 80K → NO

> 80K → YES

# Apply Model to Test Data

Test Data

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

```
        Refund
       /      \
    Yes        No
    /            \
  NO            MarSt
              /        \
    Single, Divorced   Married
          /               \
       TaxInc             NO
       /    \
   < 80K    > 80K
    /          \
   NO         YES
```

10

# Apply Model to Test Data

## Test Data

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

Refund
- Yes → NO
- No → MarSt

MarSt
- Single, Divorced → TaxInc
- Married → NO

TaxInc
- < 80K → NO
- > 80K → YES

Assign Cheat to "No"

# Exercise

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Single | 80K | **?** |
| Yes | Single | 100K | **?** |
| Yes | Married | 1K | **?** |
| Yes | Divorced | 50K | **?** |

# Decision Tree Induction

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

# Decision Tree Based Classification: Advantages

- Inexpensive to construct

- 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

14

# 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 |

*Note:*
*Outlook is the*
 *Forecast,*
*no relation to*
*Microsoft*
*email program*

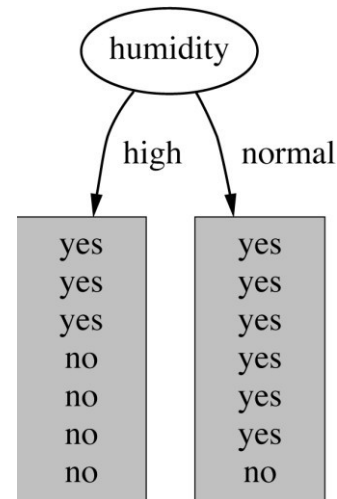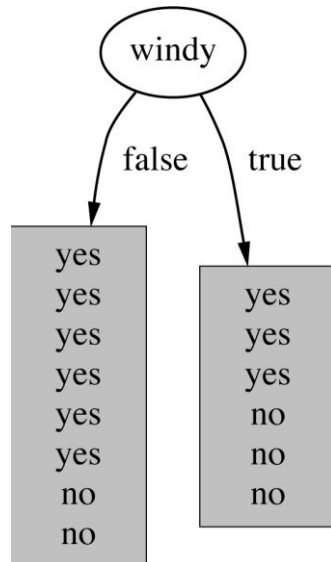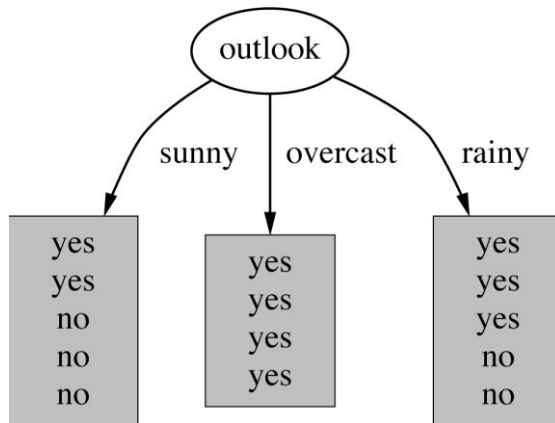# Example Tree for "Play?"

zeshan.khan@nu.edu.pk

# Building Decision Tree

- 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)
  - accuracy
  - gini index
  - others (information gain ratio)

# Which attribute to select?

zeshan.khan@nu.edu.pk

# 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:

# Example: attribute "Outlook"

- "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:

$$\text{info}([2,3],[4,0],[3,2]) = (5/14) \times 0.971 + (4/14) \times 0 + (5/14) \times 0.971 = 0.693 \text{ bits}$$

zeshan.khan@nu.edu.pk

# Computing the information gain

- Information gain:

(information before split) − (information after split)

gain("Outlook") ▢ info([9,5]) - info([2,3],[4,0],[3,2]) ▢ 0.940 - 0.693 ▢ 0.247 bits

- Information gain for attributes from weather data:

gain("Outlook") ▢ 0.247 bits

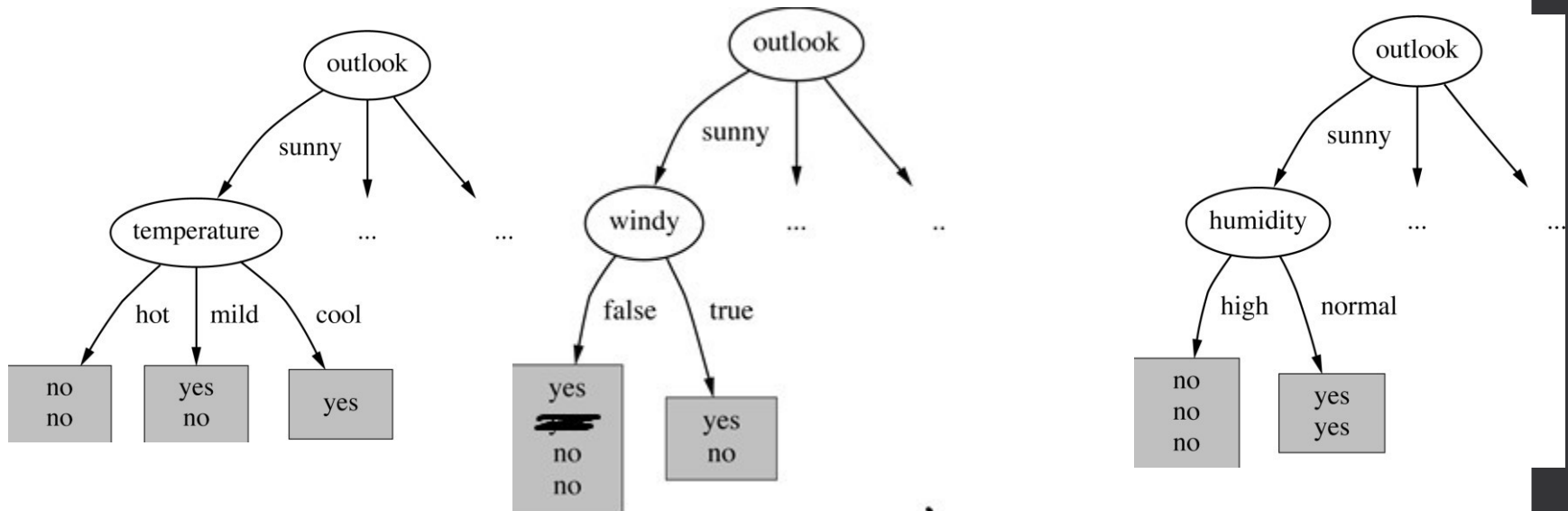gain("Temperature") ▢ 0.029

bits gain("Humidity") ▢

0.152 bits gain("Windy") ▢

0.048 bits

# Continuing to split



gain("Temperature") ▫ 0.571 bits

gain("Humidity") ▫ 0.971 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

# CART Splitting Criteria: Gini Index

- index, gini(T) is defined as

- where is the relative frequency of class j in T. gini(T) is minimized if the classes in T are skewed.

27

# Measure of Impurity: GINI

- Gini Index for a given node t :


- (NOTE: $p( j \mid t)$ is the relative frequency of class j at node t).

  – Maximum (1 - 1/$n_c$) when records are equally distributed among all classes, implying least interesting information

  – Minimum (0.0) when all records belong to one class,

  - implying most interesting information

| C1 | 0 |
|----|---|
| C2 | 6 |
| **Gini=0.000** | |

| C1 | 1 |
|----|---|
| C2 | 5 |
| **Gini=0.278** | |

| C1 | 2 |
|----|---|
| C2 | 4 |
| **Gini=0.444** | |

| C1 | 3 |
|----|---|
| C2 | 3 |
| **Gini=0.500** | |

# 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

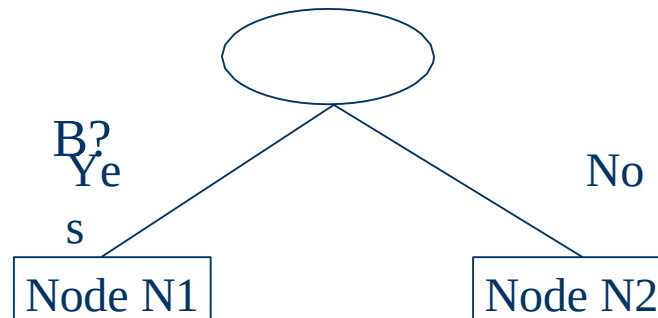| C1 | 2 |
|----|---|
| C2 | 4 |

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

P(C1) = 2/6

# Splitting Based on GINI

- Used in CART, SLIQ, SPRINT.
- When a node p is split into k partitions (children), the quality of split is computed as,

- where,

- $n_i$ = number of records at child i

- n = number of records at node p

# Binary Attributes: Computing GINI Index

- Splits into two partitions
- Effect of Weighing partitions:
  - Larger and Purer Partitions are sought for.

| | Parent |
|---|---|
| C1 | 6 |
| C2 | 6 |
| **Gini = 0.500** | |

B?

Yes                                                        No

Node N1                                                   Node N2

Gini(N1)
$= 1 - (5/7)^2 - (2/7)^2$
$= 0.408$

Gini(N2)
$= 1 - (1/5)^2 - (4/5)^2$
$= 0.32$

| | N1 | N2 |
|---|---|---|
| C1 | 5 | 1 |
| C2 | 2 | 4 |

Gini(Children)
$= 7/12 * 0.408 + 5/12 * 0.32$
$= 0.371$

# Categorical Attributes: Computing Gini Index

- For each distinct value, gather counts for each class in the dataset
- Use the count matrix to make decisions

Multi-way split

| CarType | | | |
|---|---|---|---|
| | Family | Sports | Luxury |
| C1 | 1 | 2 | 1 |
| C2 | 4 | 1 | 1 |
| Gini | 0.393 | | |

Two-way split
(find best partition of values)

| CarType | |
|---|---|
| {Sports, Luxury} | {Family} |
| C1 | 3 | 1 |
| C2 | 2 | 4 |
| Gini | 0.400 | |

| CarType | |
|---|---|
| {Sports} | {Family, Luxury} |
| C1 | 2 | 2 |
| C2 | 1 | 5 |
| Gini | 0.419 | |

# Continuous Attributes: Computing Gini Index

- Use Binary Decisions based on one value
- Several Choices for the splitting value
  - Number of possible splitting values = Number of distinct values
- Each splitting value has a count matrix associated with it
  - Class counts in each of the partitions, A < v and A $\geq$ v
- Simple method to choose best v
  - For each v, scan the database to gather count matrix and compute its Gini index
  - Computationally Inefficient! Repetition of work.

| 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 |

Taxable Income > 80K?

Yes    No

33

# Continuous Attributes: Computing Gini Index...

- For efficient computation: for each attribute,
  - Sort the attribute on values
  - Linearly scan these values, each time updating the count matrix and computing gini index
  - Choose the split position that has the least gini index

| Cheat | No | | No | | No | | Yes | | Yes | | Yes | | No | | No | | No | | No | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Taxable Income** | | | | | | | | | | | | | | | | | | | | |
| | 60 | | 70 | | 75 | | 85 | | 90 | | 95 | | 100 | | 120 | | 125 | | 220 | |
| | 55 | | 65 | | 72 | | 80 | | 87 | | 92 | | 97 | | 110 | | 122 | | 172 | | 230 |
| | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > |
| Yes | 0 | 3 | 0 | 3 | 0 | 3 | 0 | 3 | 1 | 2 | 2 | 1 | 3 | 0 | 3 | 0 | 3 | 0 | 3 | 0 |
| No | 0 | 7 | 1 | 6 | 2 | 5 | 3 | 4 | 3 | 4 | 3 | 4 | 3 | 4 | 4 | 3 | 5 | 2 | 6 | 1 | 7 | 0 |
| Gini | 0.420 | | 0.400 | | 0.375 | | 0.343 | | 0.417 | | 0.400 | | *0.300* | | 0.343 | | 0.375 | | 0.400 | | 0.420 | |

Sorted Values

Split

Positions

zeshan.khan@nu.edu.pk

# Splitting Criteria based on Classification Error

- Classification error at a node t :

$$Error(t) = 1 - \max_i P(i \mid t)$$

- Measures misclassification error made by a node.
  - Maximum $(1 - 1/n_c)$ when records are equally distributed among all classes, implying least interesting information
  - Minimum $(0.0)$ when all records belong to one class, implying
    most interesting information

# Examples for Computing Error

$$Error(t) = 1 - \max_{i} P(i \mid t)$$

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

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

Error = 1 – max (0, 1) = 1 – 1 = 0

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

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

Error = 1 – max (1/6, 5/6) = 1 – 5/6 = 1/6

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

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

Error = 1 – max (2/6, 4/6) = 1 – 4/6 = 1/3

# Misclassification Error vs Gini

A?

Yes → Node N1

No → Node N2

|  | Parent |
|----|--------|
| C1 | 7 |
| C2 | 3 |
| Gini = 0.42 | |

Gini(N1)
$= 1 - (3/3)^2 - (0/3)^2$
$= 0$

Gini(N2)
$= 1 - (4/7)^2 - (3/7)^2$
$= 0.489$

|  | N1 | N2 |
|----|----|----|
| C1 | 3 | 4 |
| C2 | 0 | 3 |
| Gini=0.361 | | |

Gini(Children)
$= 3/10 * 0$
$+ 7/10 * 0.489$
$= 0.342$

Gini improves !!

zeshan.khan@nu.edu.pk

# 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

- Stop expanding a node when all the records have similar attribute values

- Early termination