# Data Mining
# Classification: Basic Concepts, Decision Trees, and Model Evaluation

第 四 章

分类的基本概念、决策树和模型评估

# Classification: Definition

- Given a collection of records (*training set* )
  - Each record contains a set of *attributes*, one of the attributes is the *class*.
- Find a *model* for class attribute as a function of the values of other attributes.
- Goal: <u>previously unseen</u> records should be assigned a class as accurately as possible.
  - A *test set* is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.
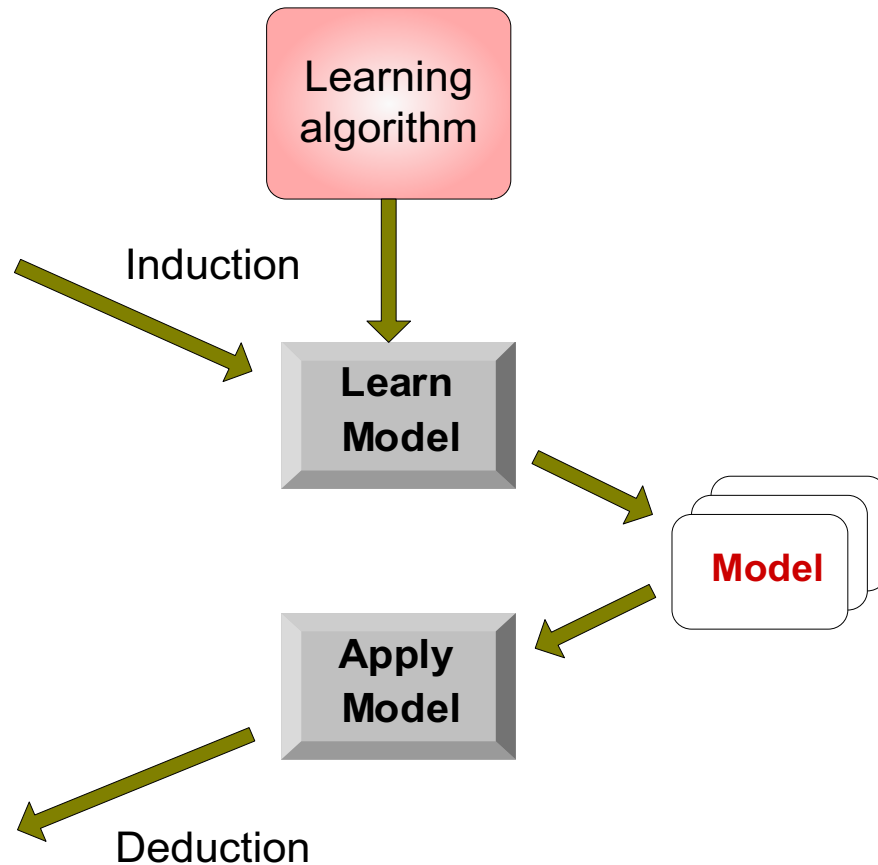
# Illustrating Classification Task

| 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**

Learning algorithm

Induction

**Learn Model**

**Model**

**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

# Examples of Classification Task

- Predicting tumor cells as benign良性的 or malignant恶性的

- Classifying credit card transactions as legitimate合法的 or fraudulent欺骗的

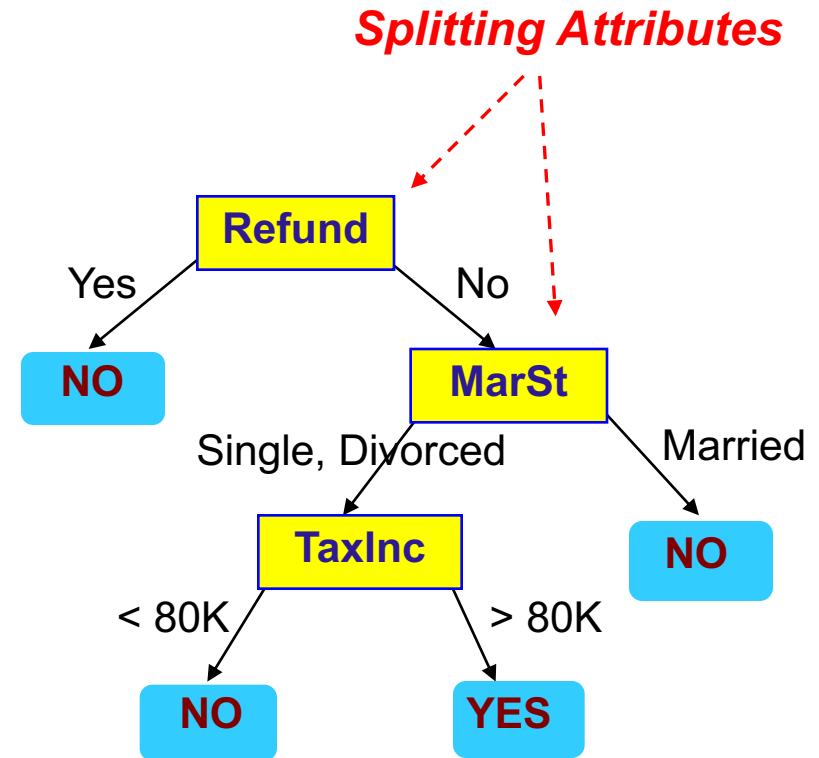- Categorizing news stories as finance财经, weather, entertainment娱乐, sports, etc

# Classification Techniques

- Decision Tree based Methods 决策树
- Rule-based Methods 基于规则的方法
- Memory based reasoning 基于记忆的推理
- Neural Networks 神经网络
- Naïve Bayes 朴素贝叶斯

   Bayesian Belief Networks 贝叶斯信念网络
- Support Vector Machines 支持向量机
- K Nearest Neighbors K近邻算法

# Classification Techniques

- Decision Tree based Methods 决策树
- Rule-based Methods 基于规则的方法
- Memory based reasoning 基于记忆的推理
- Neural Networks 神经网络
- Naïve Bayes 朴素贝叶斯

  Bayesian Belief Networks 贝叶斯信念网络
- Support Vector Machines 支持向量机
- K Nearest Neighbors K近邻算法

# 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**

# 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**
- No → **TaxInc**
  - < 80K → **NO**
  - > 80K → **YES**

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

# 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

# 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**

Single, Divorced → **TaxInc**

Married → **NO**

< 80K → **NO**

> 80K → **YES**

# Apply Model to 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 ——→ **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

MarSt
- Single, Divorced → TaxInc
- Married → NO

TaxInc
- < 80K → NO
- > 80K → YES

# Apply Model to Test Data

**Test Data**

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

# Apply Model to 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**

Assign Cheat to "No"

# 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

| 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

# Decision Tree Induction

● Many Algorithms:

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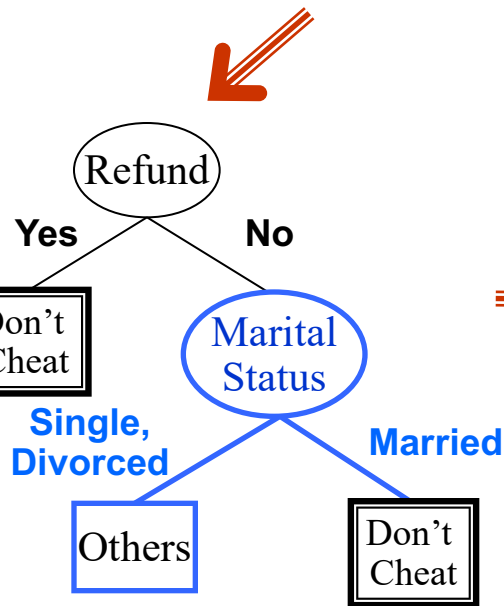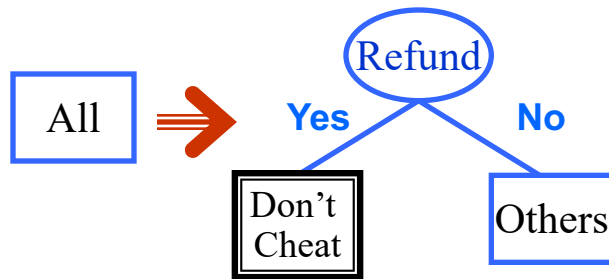- – SLIQ, SPRINT

# General Structure of Hunt's Algorithm

- Let $D_t$ be the set of training records that reach a node t
- General Procedure:
  - If $D_t$ contains records that belong the same class $y_t$, then t is a leaf node labeled as $y_t$
  - If $D_t$ is an empty set, then t is a leaf node labeled by the default class, $y_d$
  - If $D_t$ contains records that belong to more than one class, 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 |

$D_t$

**?**

# Hunt's Algorithm

| 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** |

All ⟹

Refund
- Yes → Don't Cheat
- No → Others

Refund
- Yes → Don't Cheat
- No → Marital Status
  - Single, Divorced → Others
  - Married → Don't Cheat

⟹

Refund
- Yes → Don't Cheat
- No → Marital Status
  - Single, Divorced → Taxable Income
    - < 80K → Don't Cheat
    - >= 80K → Cheat
  - Married → Don't Cheat

# 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?
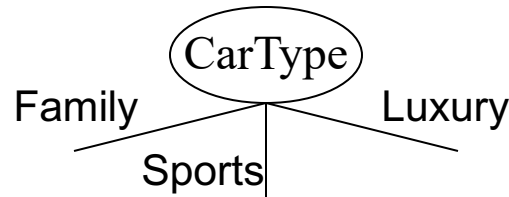    - How to determine the best split?
  - Determine when to stop splitting

# Tree Induction

● Greedy strategy.

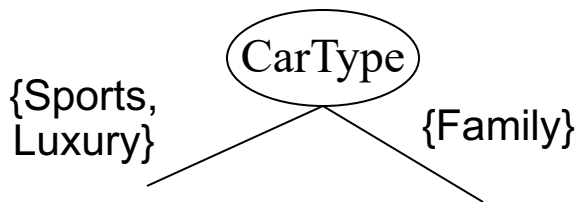– Split the records based on an attribute test that optimizes certain criterion.

● 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 Attribute Test Condition?

- Depends on attribute types
  - Nominal 标称
  - Ordinal 有序
  - Continuous 连续

- 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.



- Binary split:  Divides values into two subsets. Need to find optimal partitioning.
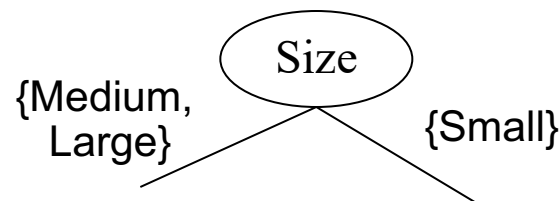
# Splitting Based on Ordinal Attributes

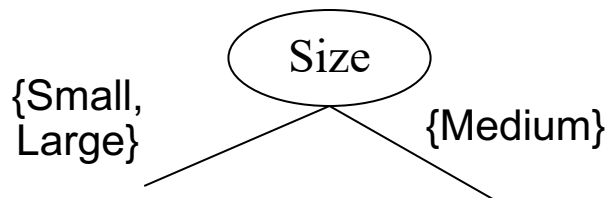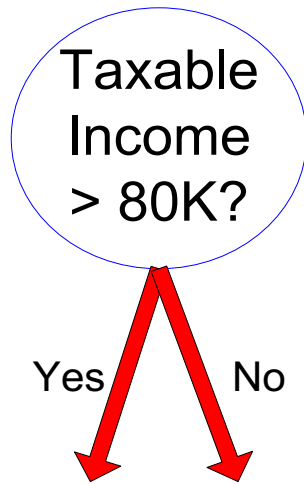- Multi-way split: Use as many partitions as distinct values.

```
        Size
Small    |    Large
       Medium
```

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

```
          Size                    OR                    Size
{Small,        {Large}                     {Medium,          {Small}
Medium}                                     Large}
```

- What about this split?

```
                Size
{Small,              {Medium}
Large}
```
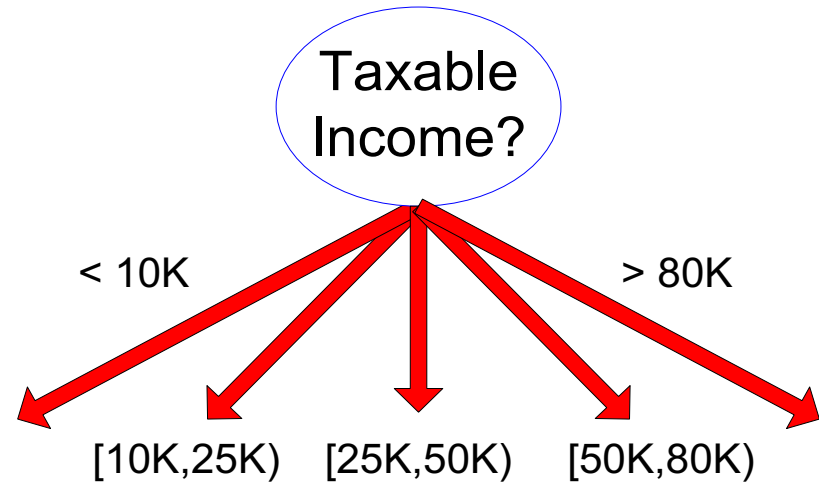
# Splitting Based on Continuous Attributes

- Different ways of handling
  - Discretization to form an ordinal categorical attribute
    - Static – discretize once at the beginning
    - Dynamic – ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.

  - Binary Decision: $(A < v)$ or $(A \geq v)$
    - consider all possible splits and finds the best cut
    - can be more computational intensive

# Splitting Based on Continuous Attributes

Taxable
Income
> 80K?

Yes    No

(i) Binary split

Taxable
Income?

< 10K    > 80K

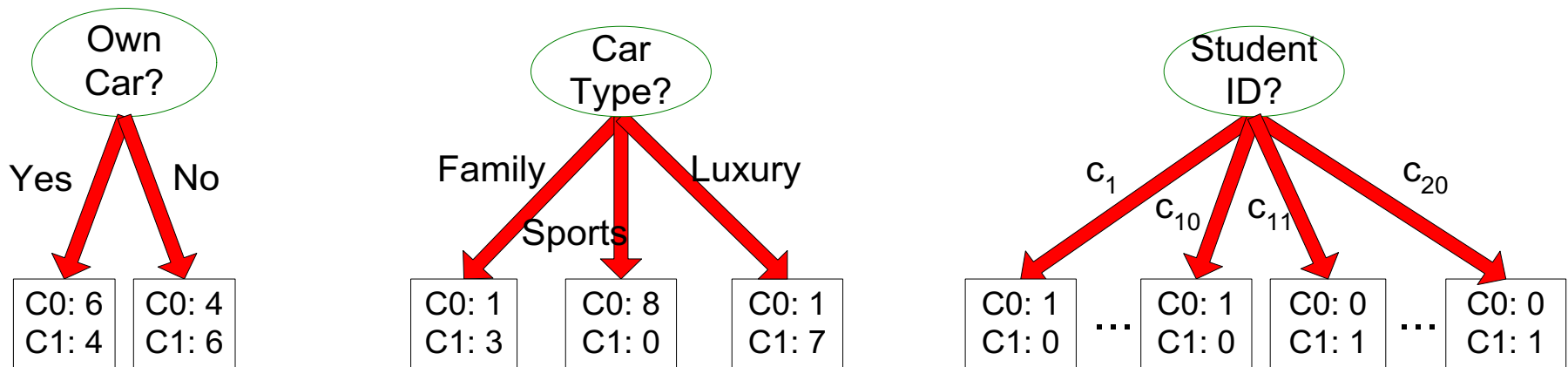[10K,25K)    [25K,50K)    [50K,80K)

(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

# How to determine the Best Split

**Before Splitting: 10 records of class 0,**
**10 records of class 1**

Own Car?
Yes / No

| C0: 6 | C0: 4 |
| C1: 4 | C1: 6 |

Car Type?
Family / Sports / Luxury

| C0: 1 | C0: 8 | C0: 1 |
| C1: 3 | C1: 0 | C1: 7 |

Student ID?
$c_1$ / $c_{10}$ / $c_{11}$ / $c_{20}$

| C0: 1 | | C0: 1 | C0: 0 | | C0: 0 |
| C1: 0 | ... | C1: 0 | C1: 1 | ... | C1: 1 |

**Which test condition is the best?**

# How to determine the Best Split

- Greedy approach:
  - Nodes with homogeneous (同质的) class distribution are preferred
- Need a measure of node impurity (不纯性):

C0: 5
C1: 5

**Non-homogeneous,**

**High degree of impurity**

C0: 9
C1: 1

**Homogeneous,**

**Low degree of impurity**

# Measures of Node Impurity

- Gini Index (基尼指数)

- Entropy (信息熵)

- Misclassification error (分类误差)

# Measure of Impurity: GINI

● Gini Index for a given node t :

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

(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) 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

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

# Splitting Based on GINI

- Used in CART(Classification and Regression Tree), SLIQ, SPRINT.

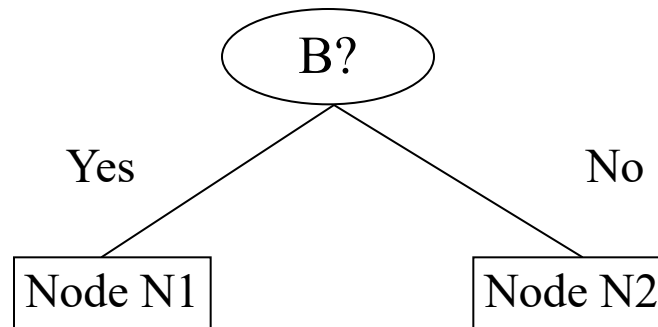- When a node p is split into k partitions (children), the quality of split is computed as,

$$GINI_{split} = \sum_{i=1}^{k} \frac{n_i}{n} GINI(i)$$

**因为子女结点所包含的样本（记录）数不同，所以给子女结点赋予权重，使得样本越多的子女结点的影响越大**

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=0.371** | | |

**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

Two-way split
(find best partition of attribute values)

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

| 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 ≥ v
- Simple method to choose best v
  - For each v, scan the database to gather count matrix (O(N)) and compute its Gini index
  - Computationally Inefficient ($O(N^2)$)! 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

# Continuous Attributes: Computing Gini Index...

- For efficient computation: for each attribute,
  - Sort the attribute on values (O(NlogN))
  - 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** | | | | | | | | | | | | | | | | | | | | |
| Sorted Values | 60 | | 70 | | 75 | | 85 | | 90 | | 95 | | 100 | | 120 | | 125 | | 220 | |
| Split Positions | 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 | |

# Alternative Splitting Criteria based on INFO

● Information Entropy at a given node t:

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

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

– Measures homogeneity (同质性) of a node.
  ◆ Maximum ($\log n_c$) when records are equally distributed among all classes implying least information
  ◆ Minimum (0) when all records belong to one class, implying most information

– Entropy based computations are similar to the GINI index computations

# 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 |

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

# Splitting Based on INFO...

● Information Gain (信息增益):

$$GAIN_{split} = Entropy(p) - \left( \sum_{i=1}^{k} \frac{n_i}{n} Entropy(i) \right)$$

Parent Node, p (with n records) is split into k partitions;

$n_i$ is number of records in partition i

– Measures Reduction in Entropy achieved because of the split. Choose the split that achieves most reduction (maximizes GAIN)

– Used in ID3 (Iterative Dichotomiser, 迭代二分器，罗斯-昆兰)

– Disadvantage: Tends to prefer splits that result in large numbers of partitions, each being small but pure.

# Splitting Based on INFO...

● Gain Ratio (增益率):

$$GainRATIO_{split} = \frac{GAIN_{Split}}{SplitINFO} \qquad SplitINFO = -\sum_{i=1}^{k} \frac{n_i}{n} \log \frac{n_i}{n}$$

    Parent Node, p (with n records) is split into k partitions

    $n_i$ is the number of records in partition i

– Adjusts Information Gain by the entropy of the partitioning (SplitINFO). Higher entropy partitioning (large numbers of small partitions) is penalized!

– Used in C4.5 (Classifier 4.5，罗斯-昆兰)

– Designed to overcome the disadvantage of Information Gain

# 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) 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

# Comparison among Splitting Criteria

**For a 2-class problem:**

# Misclassification Error vs Gini



A?

Yes — Node N1

No — Node N2

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

**Error(Parent)=0.3**

**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=0.342** | | |

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

**Gini improves !!**
**But Error NO!!**

**Error(Children)=(3/10)*0+(7/10)*(3/7)=0.3**

# 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 <span style="color:red">same (or similar) attribute values</span>

- Early termination (to be discussed later)

# Decision Tree Based Classification

● Advantages:

– Inexpensive to construct

– Extremely fast at classifying unknown records

– <span style="color:red">Easy to interpret</span> for small-sized trees

– Accuracy is comparable to other classification techniques for many simple data sets

# Example: C4.5

- Simple depth-first construction.
- <span style="color:red">Uses Gain Rate</span>
- Sorts Continuous Attributes at each node.
- Needs entire data to fit in memory.
- Unsuitable for Large Datasets.
  - Needs out-of-core sorting.

# 分类算法与模型的分析

- 分类模型设计中的实际问题
- 分类模型的评估

# 一、分类模型设计中的实际问题
# Practical Issues of Classification

- 1、Underfitting拟合不足 and Overfitting过分拟合

- 2、Missing Values

- 3、Costs of Classification

# 1、 Underfitting and Overfitting (Example)



500 circular and 500 triangular data points.

Circular points:

$0.5 \leq \text{sqrt}(x_1^2+x_2^2) \leq 1$

Triangular points:

$\text{sqrt}(x_1^2+x_2^2) < 0.5$ or

$\text{sqrt}(x_1^2+x_2^2) > 1$

# 1、 Underfitting and Overfitting



**Underfitting**: when model is too simple, both training and test errors are large

# Overfitting due to Noise



**Decision boundary (决策边界) is distorted by noise point**

# Overfitting due to Insufficient Examples

考虑表4-5中的五个训练记录，表中所有的记录都是正确标记的，对应的决策树在图4-26中。
尽管它的训练误差为0，但是它的检验误差却高达30%。

表 4-5　哺乳动物分类的训练集样本

| 名称 | 体温 | 胎生 | 4 条腿 | 冬眠 | 类标号 |
|------|------|------|--------|------|--------|
| 蝾螈 | 冷血 | 否 | 是 | 是 | 否 |
| 虹鳉 | 冷血 | 是 | 否 | 否 | 否 |
| 鹰 | 恒温 | 否 | 否 | 否 | 否 |
| 弱夜鹰 | 恒温 | 否 | 否 | 是 | 否 |
| 鸭嘴兽 | 恒温 | 否 | 是 | 是 | 是 |

图 4-26　根据表 4-5 中的数据集建立的决策树

人、大象和海豚都被误分类，因为决策树把恒温但不冬眠的脊柱动物划分为非哺乳动物。决
策树做出这样的分类决策是因为只有一个训练记录（鹰）具有这些特性。这个例子清楚地表明，
当决策树的叶结点没有足够的代表性样本时，很可能做出错误的预测。

# Notes on Overfitting

- 事实：Overfitting results in decision trees that are more complex than necessary

- 结果：Training error no longer provides a good estimate of how well the tree will perform on previously unseen records (test set)

- 结论：Need new ways for estimating errors

# Estimating Generalization Errors

- Training errors (训练误差): error on training set ($\Sigma$ e(t) )
- Generalization errors (泛化误差): error on test set ($\Sigma$ e'(t))

- Methods for estimating generalization errors:
  - Optimistic approach:  e'(t) = e(t)
  - Pessimistic approach:
    - For each leaf node: e'(t) = (e(t)+0.5)
    - Total errors: e'(T) = e(T) + N $\times$ 0.5 (N: number of leaf nodes)
    - For a tree with 30 leaf nodes and 10 errors on training (out of 1000 instances):
      Training error = 10/1000 = 1%
      Generalization error = (10 + 30$\times$0.5)/1000 = 2.5%

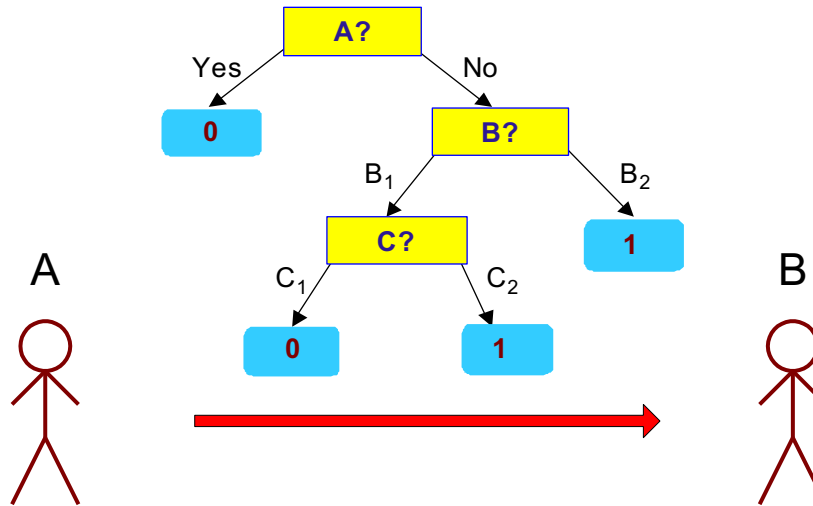# Occam's Razor（奥卡姆剃刀）-如无必要，勿增实体

# Occam's Razor（奥卡姆剃刀）

- **原则**：Given two models of similar generalization errors, **one should prefer the simpler model over the more complex model**

- **原因**：For complex models, **there is a greater chance that it was fitted accidentally by errors in data**

- **结论**：Therefore, **one should include model complexity when evaluating a model**

# Minimum Description Length (MDL)
## 最小描述长度



- Cost(Model,Data) = Cost(Data|Model) + Cost(Model)
  - **Cost is the number of bits needed for encoding**.
  - **Search for the least costly model.**
- Cost(Data|Model) encodes the misclassification errors.
- Cost(Model) uses node encoding (number of children) plus splitting condition encoding.

# How to Address Overfitting

- **Pre-Pruning (Early Stopping Rule)先剪枝（提前终止规则）**
  - Stop the algorithm before it becomes a fully-grown tree
  - Typical stopping conditions for a node:
    - ◆ Stop if all instances belong to the same class
    - ◆ Stop if all the attribute values are the same
  - More restrictive conditions:
    - ◆ Stop if **number of instances** is less than some user-specified threshold
    - ◆ Stop if **class distribution of instances** are independent of the available features (e.g., using $\chi^2$ test)
    - ◆ Stop if expanding the current node does not improve **impurity measures** (e.g., Gini, information gain or gain rate)

# How to Address Overfitting…

● Post-pruning（后剪枝）

– Grow decision tree to its entirety (完全增长)

– **Trim the nodes of the decision tree in a bottom-up fashion**

– If generalization error improves (减小) after trimming, replace sub-tree by a leaf node

– Class label of leaf node is determined from majority class of instances in the sub-tree

– Can use MDL(最小描述长度) for post-pruning

# Example of Post-Pruning

**Training Error (Before splitting) = 10/30**

**Pessimistic error = (10 + 0.5)/30 = 10.5/30**

| Class = Yes | 20 |
|---|---|
| Class = No | 10 |
| Error = 10/30 | |

**Training Error (After splitting) = 9/30**

**Pessimistic error (After splitting)**

$$= (9 + 4 \times 0.5)/30 = 11/30$$
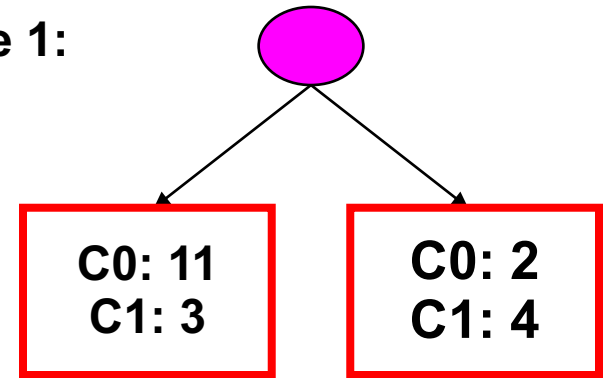
A?

**PRUNE based on pessimistic error!**

A1     A2     A3     A4

| Class = Yes | 8 |
|---|---|
| Class = No | 4 |

| Class = Yes | 3 |
|---|---|
| Class = No | 4 |

| Class = Yes | 4 |
|---|---|
| Class = No | 1 |

| Class = Yes | 5 |
|---|---|
| Class = No | 1 |

# Examples of Post-pruning

– ## Optimistic error?

**Don't prune for both cases**

– ## Pessimistic error?

**Don't prune case 1, prune case 2**

**Case 1:**



| C0: 11 | C0: 2 |
|---|---|
| C1: 3 | C1: 4 |

**Case 2:**



| C0: 14 | C0: 2 |
|---|---|
| C1: 3 | C1: 2 |

# 前四章作业

- 第二章习题：2，6，14，18，19，25
- 第三章习题：5，7，9，12
- 第四章习题：2，3，7，8

# Handling Missing Attribute Values

- Missing values affect decision tree construction in three different ways:
    - Affects how impurity measures are computed

      树节点中包括缺失特征值的记录，影响其不纯性计算
    - Affects how to distribute instance with missing value to child nodes

      父节点中包括缺失特征值的记录，影响该记录如何分配到孩子节点
    - Affects how a test instance with missing value is classified

      测试记录中包括缺失特征值的记录，影响测试记录的分类

# Computing Impurity Measure

| Tid | Refund | Marital Status | Taxable Income | Class |
|-----|--------|----------------|----------------|-------|
| 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 | **?** | Single | 90K | **Yes** |

**Missing value**

**Before Splitting:**
Entropy(Parent)
= -0.3 log(0.3)-(0.7)log(0.7) = 0.8813

|  | Class = Yes | Class = No |
|--|-------------|------------|
| Refund=Yes | **0** | **3** |
| Refund=No | **2** | **4** |
| Refund=? | **1** | **0** |

**Split on Refund:**

**Entropy(Refund=Yes) = 0**

**Entropy(Refund=No)**
**= -(2/6)log(2/6) – (4/6)log(4/6) = 0.9183**

**Entropy(Children)**
**= 0.3 (0) + 0.6 (0.9183) = 0.551**

**Gain = 0.9 × (0.8813 – 0.551) = 0.3303**

# Distribute Instances

| Tid | Refund | Marital Status | Taxable Income | Class |
|-----|--------|----------------|----------------|-------|
| 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 |

| Tid | Refund | Marital Status | Taxable Income | Class |
|-----|--------|----------------|----------------|-------|
| 10 | ? | Single | 90K | Yes |

**Refund**

Yes ← → No

| Class=Yes | 0 |
|-----------|---|
| Class=No | 3 |

| Cheat=Yes | 2 |
|-----------|---|
| Cheat=No | 4 |

**Refund**

Yes ← → No

| Class=Yes | 0 + 3/9 |
|-----------|---------|
| Class=No | 3 |

| Class=Yes | 2 + 6/9 |
|-----------|---------|
| Class=No | 4 |

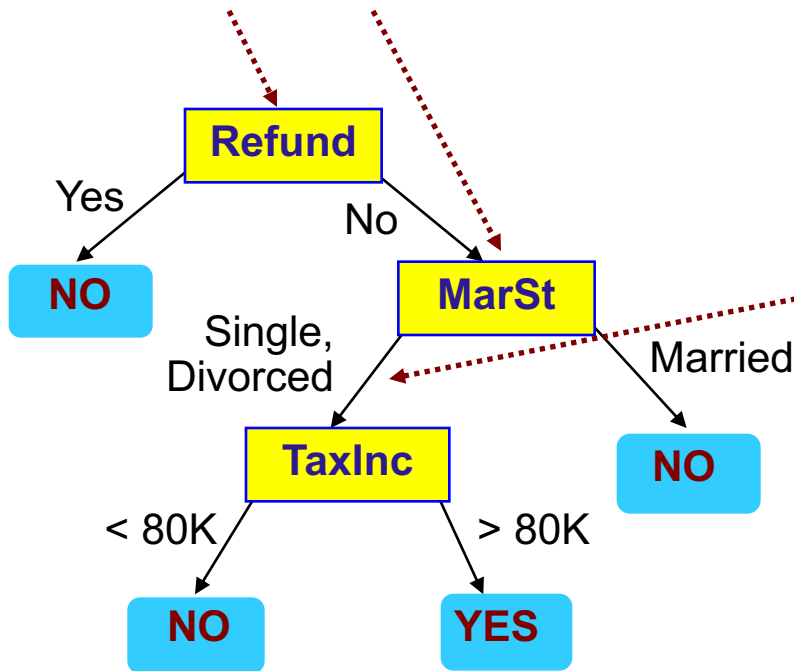**Probability that Refund=Yes is 3/9**

**Probability that Refund=No is 6/9**

**Assign record to the left child with weight = 3/9 and to the right child with weight = 6/9**

# Classify Instances

**New record:**

| Tid | Refund | Marital Status | Taxable Income | Class |
|-----|--------|----------------|----------------|-------|
| 11 | No | ? | 85K | ? |

|  | Married | Single | Divorced | Total |
|--|---------|--------|----------|-------|
| Class=No | 3 | 1 | 0 | 4 |
| Class=Yes | 6/9 | 1 | 1 | 2.67 |
| Total | 3.67 | 2 | 1 | 6.67 |

**Refund**

Yes → **NO**

No → **MarSt**

Single, Divorced → **TaxInc**

Married → **NO**

< 80K → **NO**

> 80K → **YES**

**Probability that Marital Status = Married is 3.67/6.67**

**Probability that Marital Status ={Single,Divorced} is 3/6.67**

# Classify Instances

**New record:**

| Tid | Refund | Marital Status | Taxable Income | Class |
|-----|--------|----------------|----------------|-------|
| 11 | No | ? | 85K | ? |

|  | Married | Single | Divorced | Total |
|---|---------|--------|----------|-------|
| Class=No | 3 | 1 | 0 | 4 |
| Class=Yes | 0 | 1+6/9 | 1 | 2.67 |
| Total | 3 | 2.67 | 1 | 6.67 |

**Refund**

Yes → **NO**

No → **MarSt**

Single, Divorced → **TaxInc**

Married → **NO**

< 80K → **NO**

> 80K → **YES**

**Probability that Marital Status = Married is 3/6.67**

**Probability that Marital Status ={Single,Divorced} is 3.67/6.67**

# Other Issues

- Data Fragmentation
- Search Strategy
- Expressiveness
- Tree Replication

# Data Fragmentation 数据碎片

- Number of instances gets smaller as you traverse down the tree

- Number of instances at the leaf nodes could be too small to make any statistically significant decision
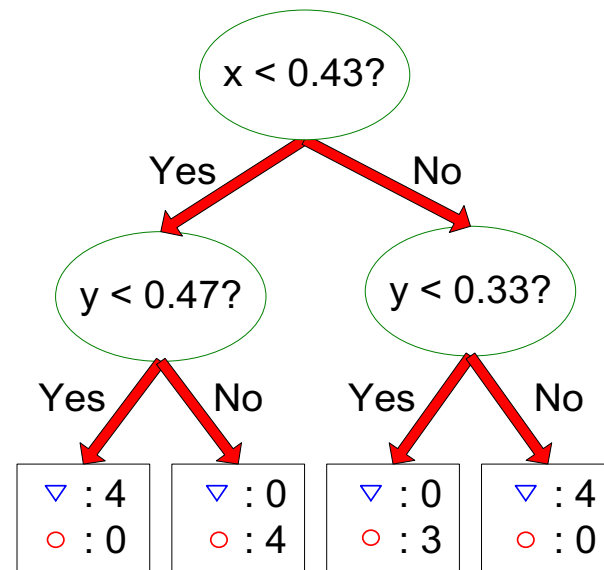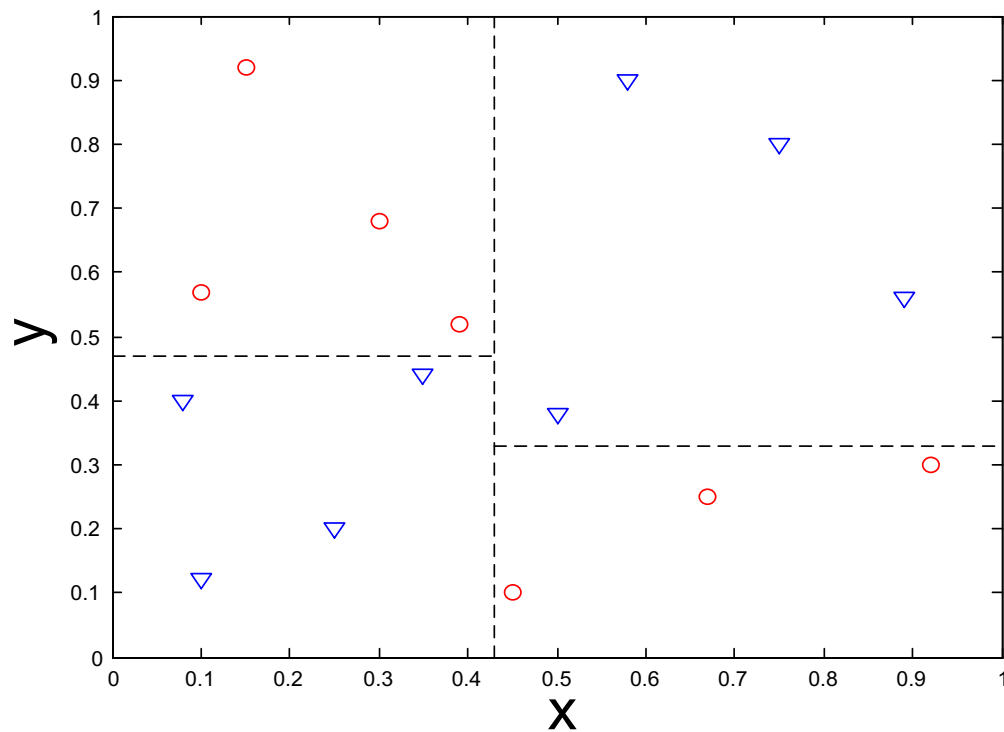
# Search Strategy

- Finding an optimal decision tree is NP-hard

- The algorithm presented so far uses a greedy, top-down, recursive partitioning strategy to induce a reasonable solution

- Other strategies?
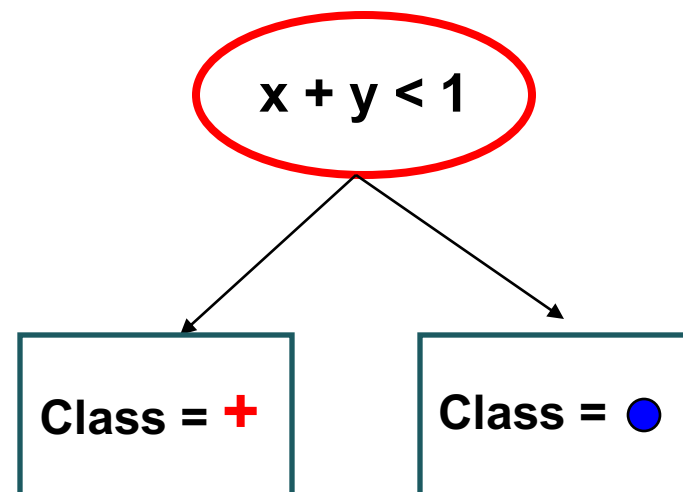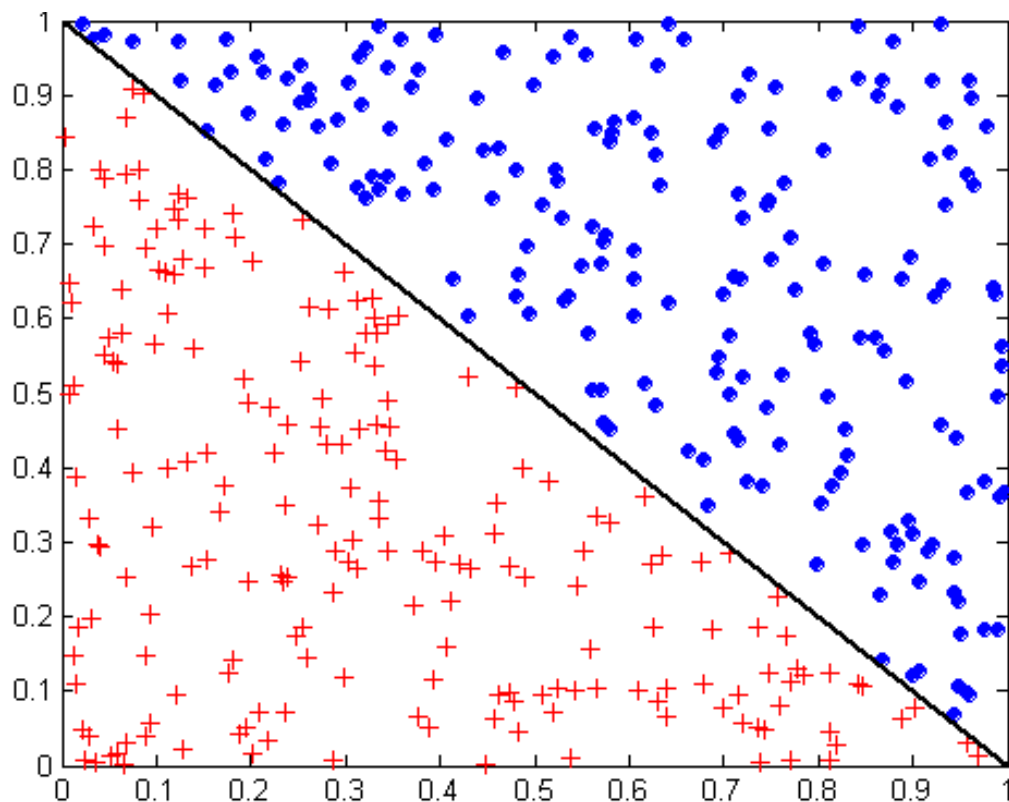    - Bottom-up
    - Bi-directional

# Expressiveness

- Decision tree provides expressive representation for learning discrete-valued function
    - But they do not generalize well to certain types of Boolean functions
        - ◆ Example: parity function:
            - Class = 1 if there is an even number of Boolean attributes with truth value = True
            - Class = 0 if there is an odd number of Boolean attributes with truth value = True
        - ◆ For accurate modeling, must have a complete tree

- Not expressive enough for modeling continuous variables
    - Particularly when test condition involves only a single attribute at-a-time
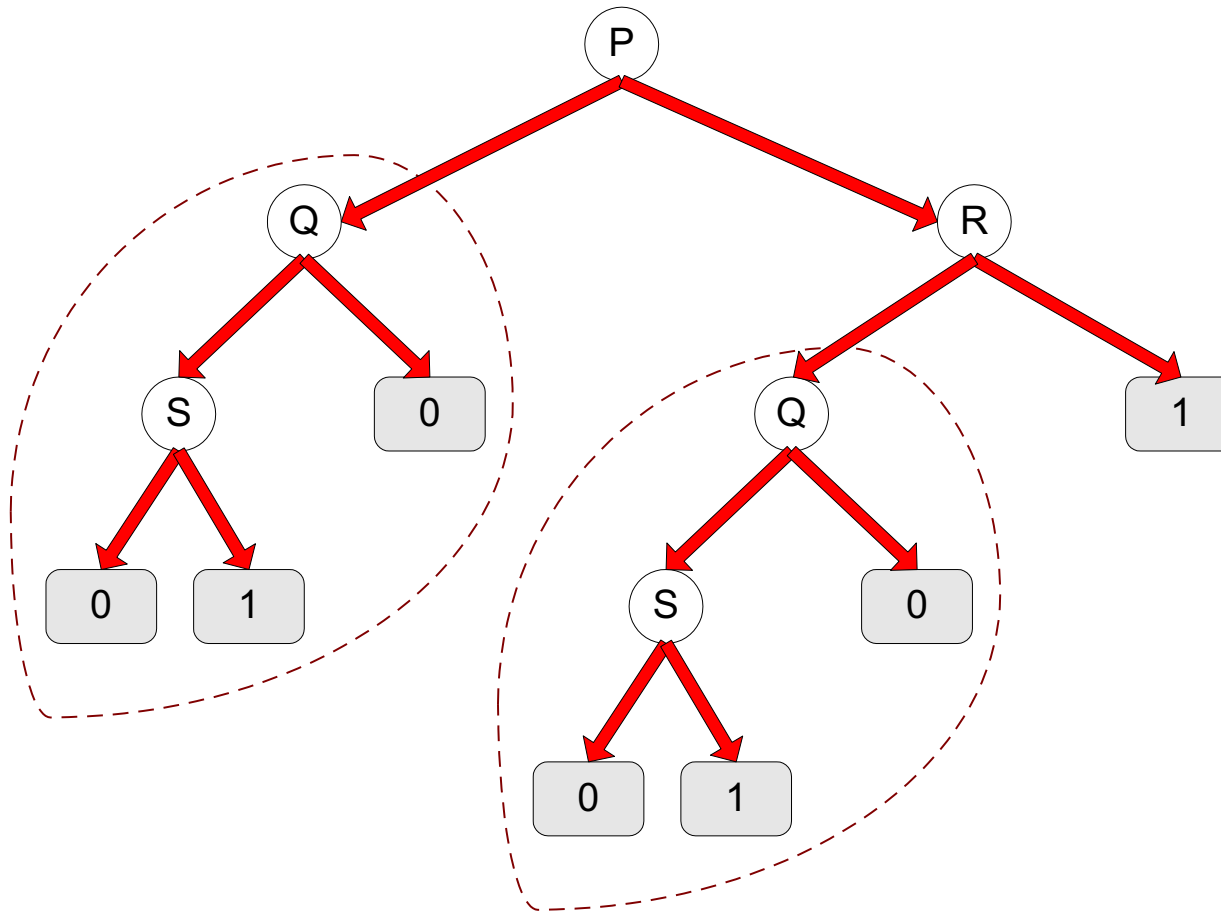
# Decision Boundary 决策边界



- **Border line between two neighboring regions of different classes is known as decision boundary**

- **Decision boundary is parallel to axes because test condition involves a single attribute at-a-time**

# **Oblique Decision Trees 斜决策树**



- **Test condition may involve multiple attributes**

- **More expressive representation 更强的表达能力，产生更紧凑的决策树**

- **Finding optimal test condition is computationally expensive**

# Tree Replication 树的复制



- **Same subtree appears in multiple branches**
- **Subtree raising** 子树提升

# 二、分类模型的评估 **Model Evaluation**

- **Metrics for Performance Evaluation**
  - How to evaluate the performance of a model?

- **Methods for Performance Evaluation**
  - How to obtain reliable estimates?

- **Methods for Model Comparison**
  - How to compare the relative performance among competing models?

# Model Evaluation

- ## **Metrics for Performance Evaluation**
  - How to evaluate the performance of a model?

- ## **Methods for Performance Evaluation**
  - How to obtain reliable estimates?

- ## **Methods for Model Comparison**
  - How to compare the relative performance among competing models?

# Metrics for Performance Evaluation

- Focus on the predictive capability of a model
  - Rather than how fast it takes to classify or build models, scalability, etc.
- Confusion Matrix混淆矩阵:

| | PREDICTED CLASS | | |
|---|---|---|---|
| ACTUAL CLASS | | Class=Yes | Class=No |
| | Class=Yes | a | b |
| | Class=No | c | d |

**a: TP (true positive)**

**b: FN (false negative)**

**c: FP (false positive)**

**d: TN (true negative)**

# Metrics for Performance Evaluation...

| | PREDICTED CLASS | | |
|---|---|---|---|
| **ACTUAL CLASS** | | Class=Yes | Class=No |
| | Class=Yes | a (TP) | b (FN) |
| | Class=No | c (FP) | d (TN) |

- Most widely-used metric:  **Accuracy**

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

# Limitation of Accuracy

- Consider a 2-class problem
  - Number of Class 0 examples = 9990
  - Number of Class 1 examples = 10

- If model predicts everything to be class 0, accuracy is 9990/10000 = 99.9 %
  - Accuracy is misleading because model does not detect any class 1 example

# Cost Matrix

| | PREDICTED CLASS | | |
|---|---|---|---|
| ACTUAL CLASS | C(i\|j) | **Class=Yes** | **Class=No** |
| | **Class=Yes** | C(Yes\|Yes) | C(No\|Yes) |
| | **Class=No** | C(Yes\|No) | C(No\|No) |

C(i|j): Cost of misclassifying class j example as class i

# Computing Cost of Classification

| Cost Matrix | PREDICTED CLASS | | |
|---|---|---|---|
| | C(i\|j) | + | - |
| ACTUAL CLASS | + | -1 | 100 |
| | - | 1 | 0 |

| Model M$_1$ | PREDICTED CLASS | | |
|---|---|---|---|
| | | + | - |
| ACTUAL CLASS | + | 150 | 40 |
| | - | 60 | 250 |

| Model M$_2$ | PREDICTED CLASS | | |
|---|---|---|---|
| | | + | - |
| ACTUAL CLASS | + | 250 | 45 |
| | - | 5 | 200 |

Accuracy = 80%
Cost = 3910

Accuracy = 90%
Cost = 4255

# Cost vs Accuracy

| Count | PREDICTED CLASS | |
|---|---|---|
| | Class=Yes | Class=No |
| **ACTUAL CLASS** Class=Yes | a | b |
| Class=No | c | d |

Accuracy is proportional to cost if
1. C(Yes|No)=C(No|Yes) = q
2. C(Yes|Yes)=C(No|No) = p

$N = a + b + c + d$

$Accuracy = (a + d)/N$

| Cost | PREDICTED CLASS | |
|---|---|---|
| | Class=Yes | Class=No |
| **ACTUAL CLASS** Class=Yes | p | q |
| Class=No | q | p |

$Cost = p (a + d) + q (b + c)$

$\quad = p (a + d) + q (N - a - d)$

$\quad = q N - (q - p)(a + d)$

$\quad = N [q - (q-p) \times Accuracy]$

# Cost-Sensitive Measures

| | PREDICTED CLASS | | |
|---|---|---|---|
| ACTUAL CLASS | | Class=Yes | Class=No |
| | Class=Yes | a (TP) | b (FN) |
| | Class=No | c (FP) | d (TN) |

$$\text{Precision (p)} = \frac{a}{a+c}$$

$$\text{Recall (r)} = \frac{a}{a+b}$$

$$\text{F-measure (F)} = \frac{2rp}{r+p} = \frac{2a}{2a+b+c}$$

r和p的调和平均值2/(1/r+1/p)

● Precision is biased towards C(Yes|Yes) & C(Yes|No)
● Recall is biased towards C(Yes|Yes) & C(No|Yes)
● F-measure is biased towards all except C(No|No)

$$\text{Weighted Accuracy} = \frac{w_1 a + w_4 d}{w_1 a + w_2 b + w_3 c + w_4 d}$$

# Model Evaluation

- Metrics for Performance Evaluation
  - How to evaluate the performance of a model?

- Methods for Performance Evaluation
  - How to obtain reliable estimates?

- Methods for Model Comparison
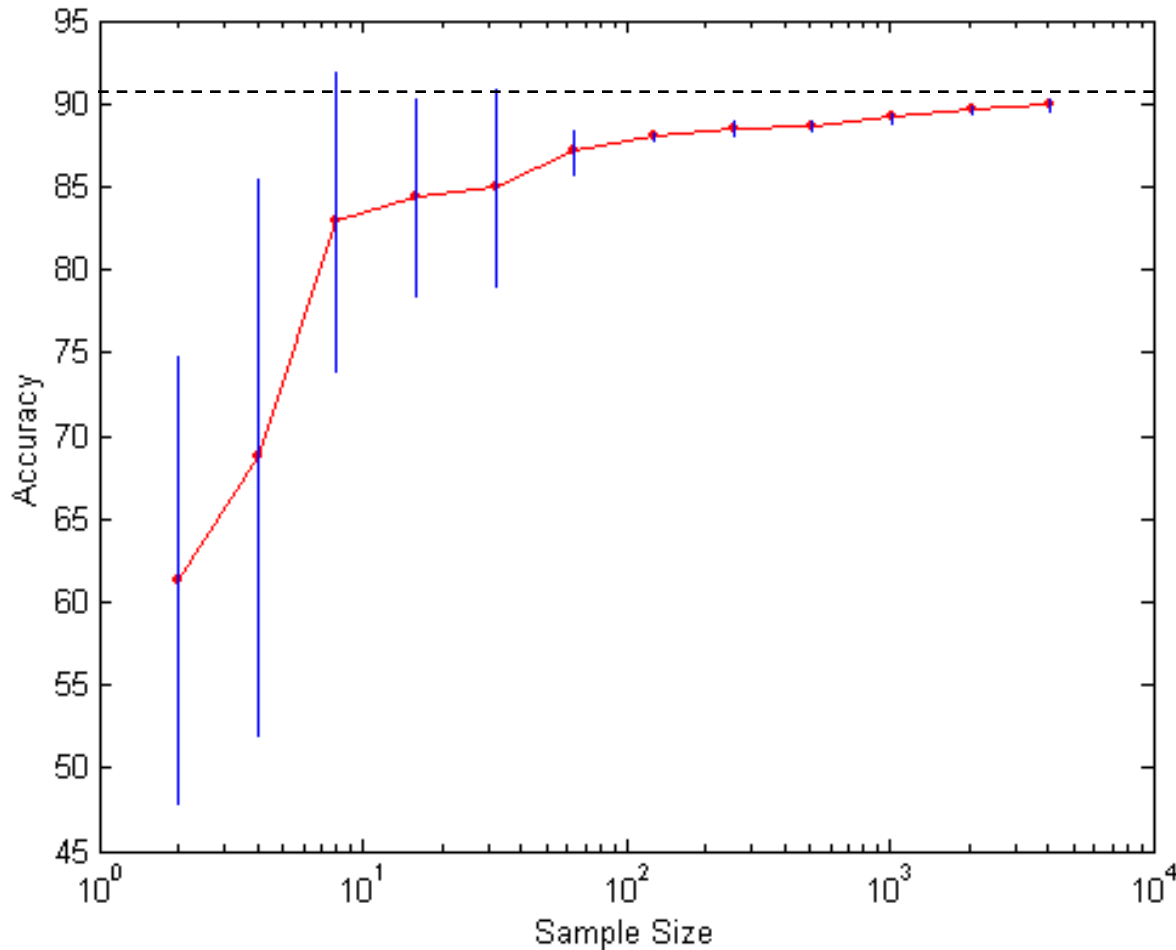  - How to compare the relative performance among competing models?

# Methods for Performance Evaluation

- How to obtain a reliable estimate of performance?

- Performance of a model may depend on other factors besides the learning algorithm:
  - Class distribution
  - Cost of misclassification
  - Size of training and test sets

# Learning Curve



- Learning curve shows how accuracy changes with varying sample size

- Requires a sampling schedule for creating learning curve:
  - Arithmetic sampling (Langley, et al)
  - Geometric sampling (Provost et al)

Effect of small sample size:
- Bias in the estimate
- Variance of estimate

# Methods of Estimation

- Holdout 保持方法
  - Reserve 2/3 for training and 1/3 for testing
- Random subsampling 随机二次抽样
  - Repeated holdout
- **Cross validation** 交叉验证
  - Partition data into k disjoint subsets
  - k-fold: train on k-1 partitions, test on the remaining one
  - Leave-one-out: k=n
- Stratified sampling 分层抽样
  - Oversampling(过采样) vs undersampling（欠采样）
- Bootstrap 自助法
  - Sampling with replacement
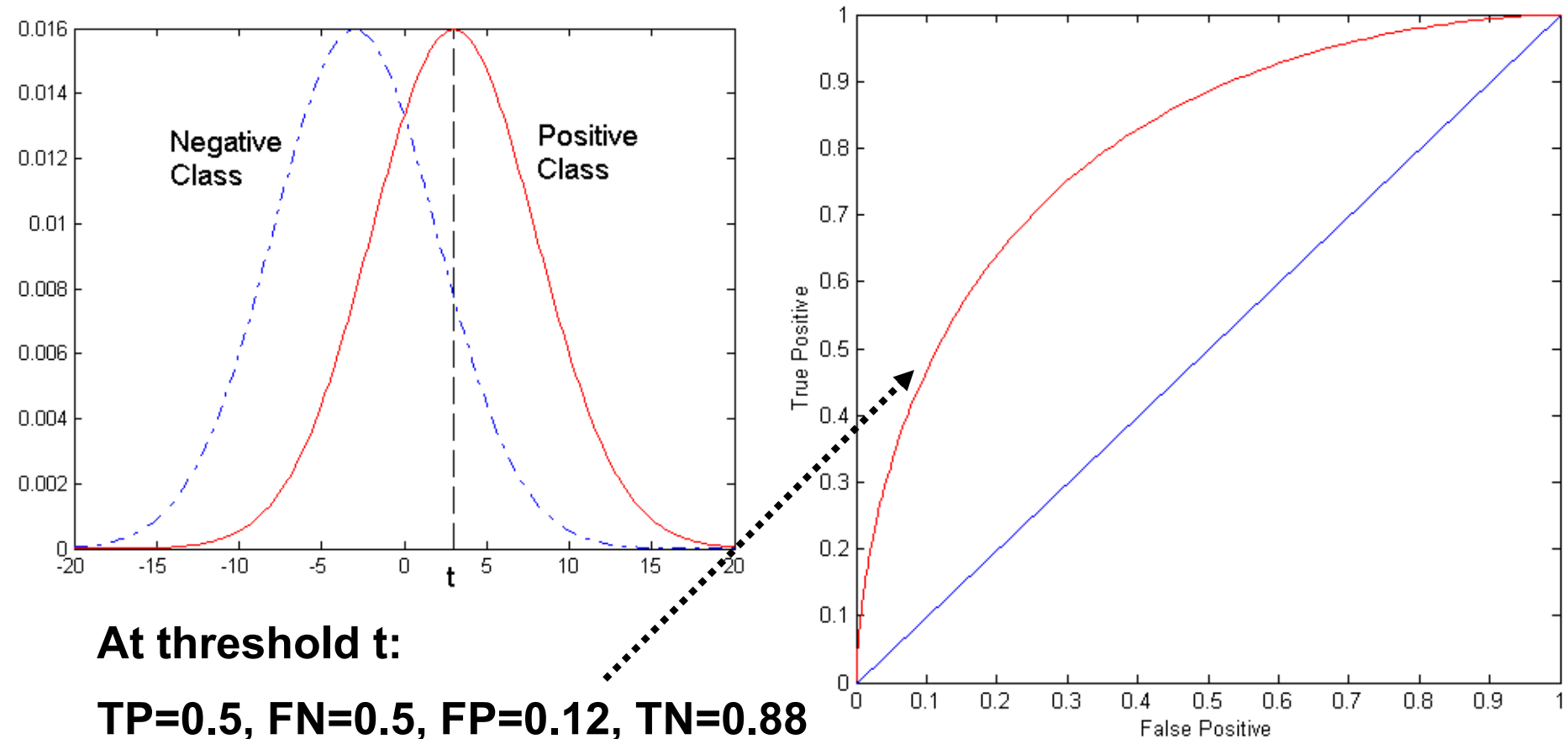
# ROC (Receiver Operating Characteristic)

- Developed in 1950s for signal detection theory to analyze noisy signals
    - Characterize the trade-off between positive hits and false alarms
- ROC curve plots TP (on the y-axis) against FP (on the x-axis)
- Performance of each classifier represented as a point on the ROC curve
    - changing **the threshold of algorithm**, **sample distribution** or **cost matrix** changes the location of the point
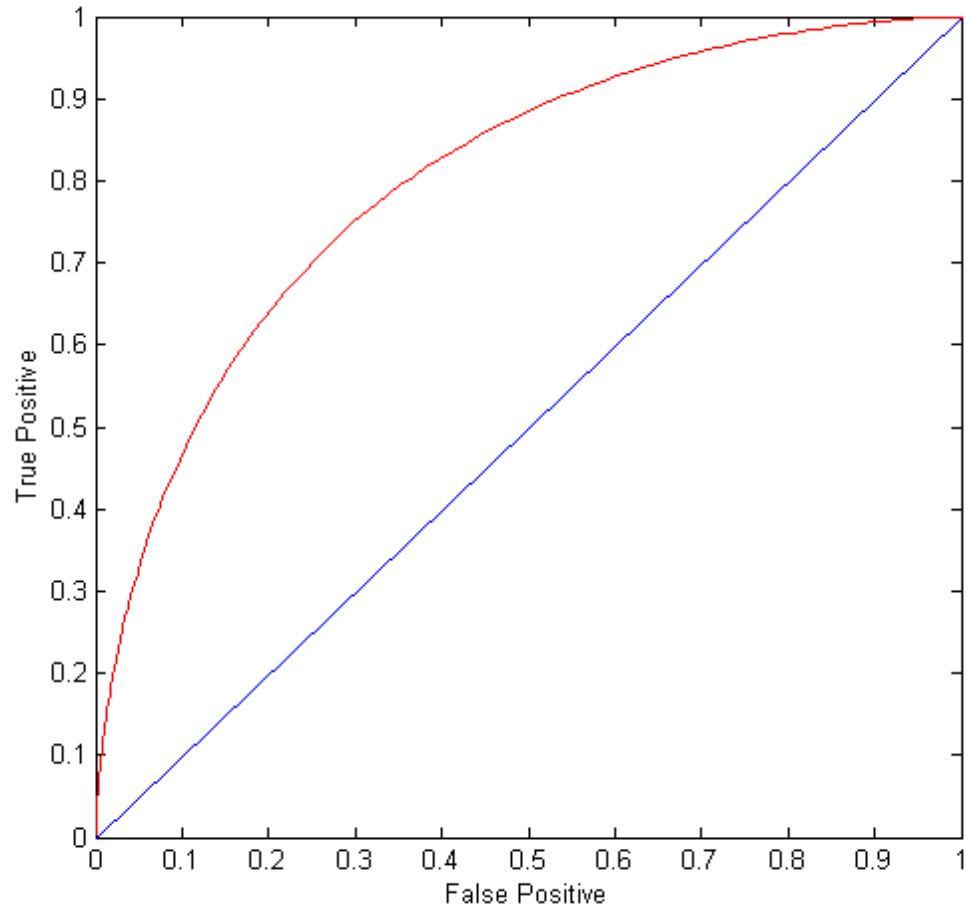
# ROC Curve

- 1-dimensional data set containing 2 classes (positive and negative)

- any points located at x > t is classified as positive



At threshold t:

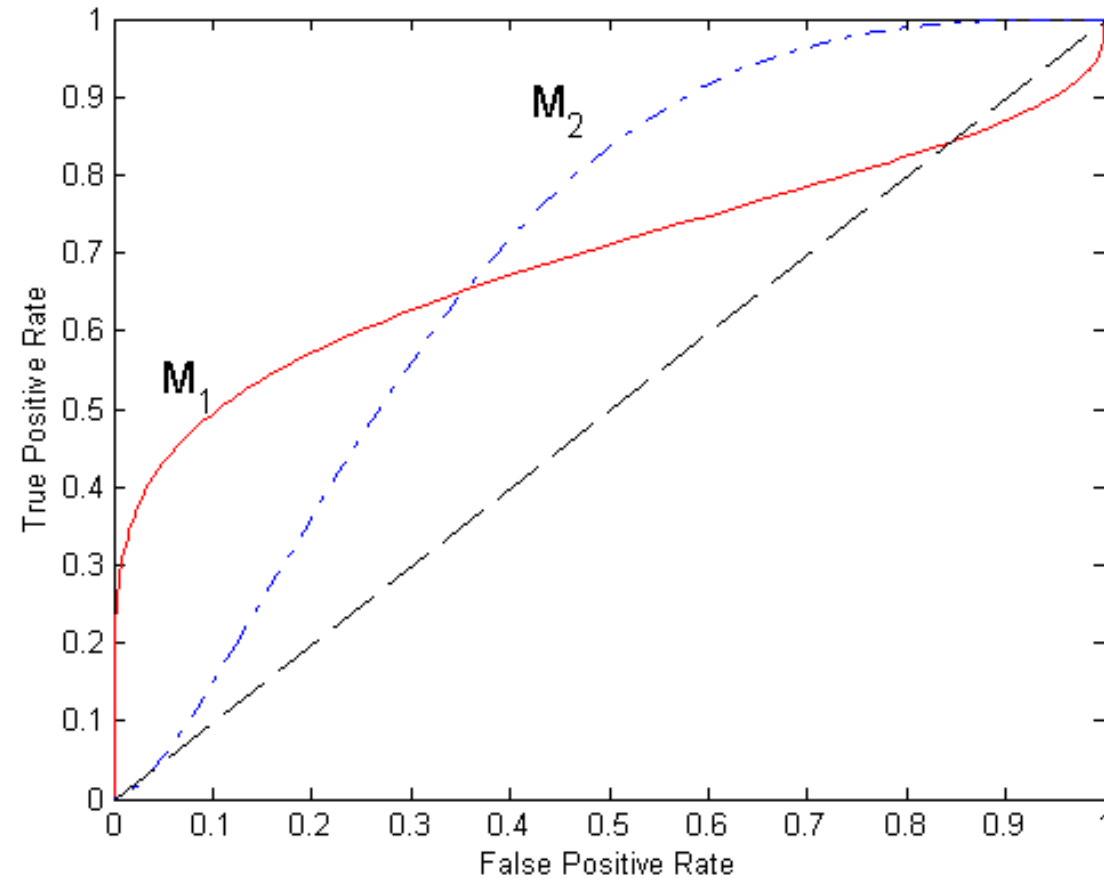TP=0.5, FN=0.5, FP=0.12, TN=0.88

# ROC Curve

(TP,FP):

- (0,0): declare everything
    to be negative class
- (1,1): declare everything
    to be positive class
- (1,0): ideal

- Diagonal line:
    - Random guessing
    - Below diagonal line:
        ◆ prediction is opposite of
        the true class

# Using ROC for Model Comparison



- No model consistently outperform the other
  - $M_1$ is better for small FPR
  - $M_2$ is better for large FPR

- Area Under the ROC curve
  - Ideal:
    - Area = 1
  - Random guess:
    - Area = 0.5

# How to Construct an ROC curve

| Instance | P(+|A) | True Class |
|----------|--------|------------|
| 1 | 0.95 | + |
| 2 | 0.93 | + |
| 3 | 0.87 | - |
| 4 | 0.85 | - |
| 5 | 0.85 | - |
| 6 | 0.85 | + |
| 7 | 0.76 | - |
| 8 | 0.53 | + |
| 9 | 0.43 | - |
| 10 | 0.25 | + |

正样本中被预测为阳性的比例

负样本中被预测为阳性的比例

- Use classifier that produces posterior probability for each test instance P(+|A)

- Sort the instances according to P(+|A) in decreasing order

- Apply threshold at each unique value of P(+|A)

- Count the number of TP, FP, TN, FN at each threshold

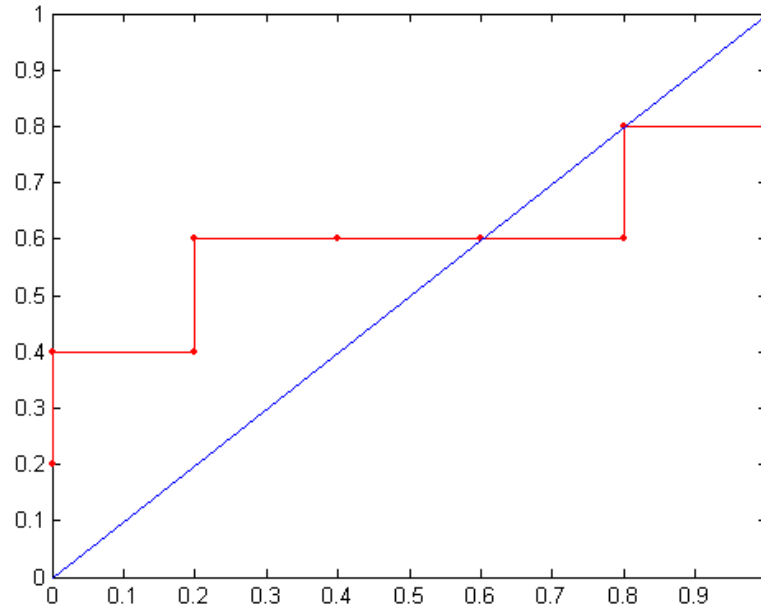- TP rate, TPR = TP/(TP+FN)

- FP rate, FPR = FP/(FP + TN)

# How to construct an ROC curve

| Class | + | - | + | - | - | - | + | - | + | + | |
|-------|------|------|------|------|------|------|------|------|------|------|------|
| **Threshold >=** | 0.25 | 0.43 | 0.53 | 0.76 | 0.85 | 0.85 | 0.85 | 0.87 | 0.93 | 0.95 | 1.00 |
| TP | 5 | 4 | 4 | 3 | 3 | 3 | 3 | 2 | 2 | 1 | 0 |
| FP | 5 | 5 | 4 | 4 | 3 | 2 | 1 | 1 | 0 | 0 | 0 |
| TN | 0 | 0 | 1 | 1 | 2 | 3 | 4 | 4 | 5 | 5 | 5 |
| FN | 0 | 1 | 1 | 2 | 2 | 2 | 2 | 3 | 3 | 4 | 5 |
| TPR | 1 | 0.8 | 0.8 | 0.6 | 0.6 | 0.6 | 0.6 | 0.4 | 0.4 | 0.2 | 0 |
| FPR | 1 | 1 | 0.8 | 0.8 | 0.6 | 0.4 | 0.2 | 0.2 | 0 | 0 | 0 |

• TP rate, TPR = TP/(TP+FN)

• FP rate, FPR = FP/(FP + TN)

## ROC Curve:

# Sensitivity & Specificity

敏感性 =TP/(TP+FN)    敏感性也称真阳性率TPR

**敏感性越高，则漏诊率越低**

特异性 = TN/(TN+FP)   特异性也称真阴性率TNR=1-FPR

**特异性越高，则误诊率越低**

ROC曲线

横坐标: 1-specificity    FPR
纵坐标:  sensitivity      TPR

# Model Evaluation

- Metrics for Performance Evaluation
  - How to evaluate the performance of a model?

- Methods for Performance Evaluation
  - How to obtain reliable estimates?

- Methods for Model Comparison
  - How to compare the relative performance among competing models?

# Test of Significance

- Given two models:
  - Model M1: accuracy = 85%, tested on 30 instances
  - Model M2: accuracy = 75%, tested on 5000 instances

- Can we say M1 is better than M2?
  - How much confidence can we place on accuracy of M1 and M2?
  - Can the difference in performance measure be explained as a result of random fluctuations in the test set?
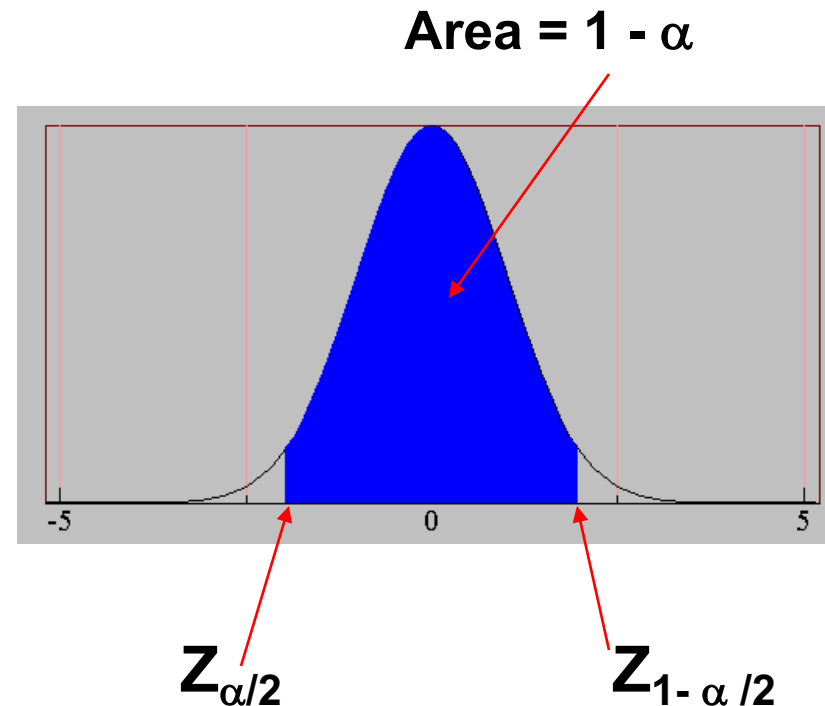
# Confidence Interval for Accuracy

- Prediction can be regarded as a Bernoulli trial
  - A Bernoulli trial has 2 possible outcomes
  - Possible outcomes for prediction: correct or wrong
  - Collection of Bernoulli trials has a Binomial distribution:
    - $x \sim \text{Bin}(N, p)$     x: number of correct predictions
    - e.g:   Toss a fair coin 50 times, how many heads would turn up?
      Expected number of heads = $N \times p = 50 \times 0.5 = 25$

- Given x (# of correct predictions) or equivalently, acc=x/N, and N (# of test instances),

  Can we predict p (true accuracy of model)?

# Confidence Interval for Accuracy

- For large test sets (N > 30),

  – acc has a normal distribution with mean p and variance p(1-p)/N

$$P(Z_{\alpha/2} < \frac{acc - p}{\sqrt{p(1-p)/N}} < Z_{1-\alpha/2})$$

$$= 1 - \alpha$$

**Area = 1 - α**

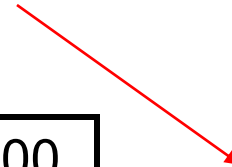**Z$_{\alpha/2}$**   **Z$_{1-\alpha/2}$**

- Confidence Interval for p:

$$p = \frac{2 \times N \times acc + Z_{\alpha/2}^2 \pm \sqrt{Z_{\alpha/2}^2 + 4 \times N \times acc - 4 \times N \times acc^2}}{2(N + Z_{\alpha/2}^2)}$$

# Confidence Interval for Accuracy

- Consider a model that produces an accuracy of 80% when evaluated on 100 test instances:

  - N=100, acc = 0.8

  - Let 1-$\alpha$ = 0.95 (95% confidence)

  - From probability table, $Z_{\alpha/2}$=1.96

| 1-$\alpha$ | Z |
|------|------|
| 0.99 | 2.58 |
| 0.98 | 2.33 |
| 0.95 | 1.96 |
| 0.90 | 1.65 |

| N | 50 | 100 | 500 | 1000 | 5000 |
|------|------|------|------|------|------|
| p(lower) | 0.670 | 0.711 | 0.763 | 0.774 | 0.789 |
| p(upper) | 0.888 | 0.866 | 0.833 | 0.824 | 0.811 |

# Comparing Performance of 2 Models

- Given two models, say M1 and M2, which is better?
  - M1 is tested on D1 (size=n1), found error rate = $e_1$
  - M2 is tested on D2 (size=n2), found error rate = $e_2$
  - Assume D1 and D2 are independent
  - If n1 and n2 are sufficiently large, then

$$e_1 \sim N(\mu_1, \sigma_1)$$
$$e_2 \sim N(\mu_2, \sigma_2)$$

  - Approximate: $\hat{\sigma}_i = \sqrt{\dfrac{e_i(1-e_i)}{n_i}}$

# Comparing Performance of 2 Models

- To test if performance difference is statistically significant:  d = e1 – e2
  - $d \sim N(d_t, \sigma_t)$   where $d_t$ is the true difference
  - Since D1 and D2 are independent, their variance adds up:

$$\sigma_t^2 = \sigma_1^2 + \sigma_2^2 \cong \hat{\sigma}_1^2 + \hat{\sigma}_2^2$$

$$= \frac{e1(1-e1)}{n1} + \frac{e2(1-e2)}{n2}$$

  - At (1-$\alpha$) confidence level,   $d_t = d \pm Z_{\alpha/2} \hat{\sigma}_t$

# An Illustrative Example

- Given: M1: n1 = 30, e1 = 0.15
  - M2: n2 = 5000, e2 = 0.25
- d = |e2 – e1| = 0.1   (2-sided test)

$$\hat{\sigma}_d = \frac{0.15(1-0.15)}{30} + \frac{0.25(1-0.25)}{5000} = 0.0043$$

- At 95% confidence level, $Z_{\alpha/2}$=1.96

$$d_t = 0.100 \pm 1.96 \times \sqrt{0.0043} = 0.100 \pm 0.128$$

=> Interval contains 0 => difference may not be statistically significant

# Comparing Performance of 2 Algorithms

- Each learning algorithm may produce k models:
  - L1 may produce M11 , M12, …, M1k
  - L2 may produce M21 , M22, …, M2k

- If models are generated on the same test sets D1,D2, …, Dk (e.g., via cross-validation)
  - For each set: compute $d_j = e_{1j} - e_{2j}$
  - $d_j$ has mean $d_t$ and variance $\sigma_t$
  - Estimate:

$$\hat{\sigma}_t^2 = \frac{\sum\limits_{j=1}^{k}(d_j - \overline{d})^2}{k(k-1)}$$

$$d_t = d \pm t_{1-\alpha, k-1}\, \hat{\sigma}_t$$