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

# **Classification: Definition**

- Given a collection of examples (*dataset*)
  - Each example (record) contains a set of attributes (features), one of the attributes is the class.
- Find a *model* (or procedure) to predict the class attribute as a function of the values of other attributes.
- Goal: <u>previously unseen</u> examples should be classified as accurately as possible (generalization).

## **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** 

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

# **Examples of Classification Task**

- Predicting tumor cells as benign or malignant
- Classifying credit card transactions as legitimate or fraudulent



- Classifying secondary structures of protein as alpha-helix, beta-sheet, or random coil
- Categorizing news stories as finance, weather, entertainment, sports, etc



# **Classification Techniques**

- Decision Tree based Methods
- Rule-based Methods
- Memory based reasoning
- Linear Classifiers and Neural Networks
- Statistical Methods: eg.Naïve Bayes
- Kernel-Based Methods (e.g. Support Vector Machines)

#### **Measures of Classification Performance**

- Confusion Matrix
- Two-class examples

	PREDICTED CLASS		
		Class=Yes	Class=No
ACTUAL	Class=Yes	а	b
CLASS	Class=No	С	d

a: TP (true positive)b: FN (false negative)c: FP (false positive)d: TN (true negative)

#### **Measures of Classification Performance**

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

• Most widely-used metric:

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

#### **Other Measures**

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

Precision (p)  $= \frac{a}{a+c}$ Recall (r)  $= \frac{a}{a+b}$ F-measure (F1)  $= \frac{2rp}{r+p} = \frac{2a}{2a+b+c}$ 

• F1 measure is closer to the minimum between p and r

# **Estimating Generalization**

- General idea: use another dataset not used for training
- Holdout
  - Split the dataset: a% for training and (100-a)% for testing (depends on randomness)
  - Class analogy should be the same in training and test
- **k-fold Cross validation** (k-fold CV) (usually k=10)
  - Partition the dataset into k disjoint subsets (keep class analogy)
  - k-fold: train on k-1 partitions, test on the remaining one
  - Compute average performance on k-folds (+ std)
  - Expensive: requires training of k models
  - Depends on randomness (although to a smaller degree)
  - Leave-one-out: k=n (deterministic split, heavy computation)
- Stratified cross-validation (random partitioning into folds)

# **Comparing Performance of two methods**

- Or of the same method with different parameters
- If models are generated on the same test sets D1,D2, ..., Dk (e.g., via cross-validation)
- Each method produces k models:
  - 1<sup>st</sup> method produces M11 , M12, ..., M1k
  - 2nd method produces M21 , M22, ..., M2k
  - For each set: compute  $d_j = acc_{1j} acc_{2j}$
  - $d_j$  has mean  $d_t$  and variance  $\sigma_t$  $\hat{\sigma}_i^2 = \frac{\sum_{j=1}^k (d_j - d)^2}{k(k-1)}$

$$d_{t} = d \pm t_{t_{1-lpha,k-1}} \hat{\sigma}$$

Interval for *d* contains 0 => difference
 may not be statistically significant

# The generalization problem

- Occam's Razor: Prefer the simplest model that fits well to the data
- Very simple models underfit (bad generalization)
- Complex models overfit (bad generalization)
- For a given dataset and type of model, there is an optimal model complexity providing optimal generalization

#### **Example of a Decision Tree**



Model: Decision Tree



**Class Labels at Leaves** 

Dataset

#### **Another Example of Decision Tree**





# 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** 

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



Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?











### **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** 

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** 



# **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<sub>t</sub> be the set of training records that reach a node t
- General Procedure:
  - If D<sub>t</sub> contains records that belong the same class y<sub>t</sub>, then t is a leaf node labeled as y<sub>t</sub>
  - If D<sub>t</sub> is an empty set, then t is a leaf node labeled by the default class, y<sub>d</sub>
  - If D<sub>t</sub> 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





## **Tree Induction**

- Successive Partitioning of the dataset based on attribute tests
- Greedy strategy.
  - Split the records based on the 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**

- Successive Partitioning of the dataset based on attribute tests
- 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.



 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.



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



## **Splitting Based on Continuous Attributes**

- Different ways of handling
  - Discretization to form a 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  $(A \ge v)$ 
    - consider all possible splits and finds the best cut
    - can be more compute intensive

#### **Splitting Based on Continuous Attributes**



### **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



Which test condition is the best?

#### How to determine the Best Split

- Nodes with homogeneous class distribution are preferred
- Need a measure of node impurity:

Non-homogeneous, High degree of impurity C0: 9 C1: 1

Homogeneous,

Low degree of impurity

#### **Measures of Node Impurity**

• Gini Index

#### Entropy

Misclassification error

#### How to Find the Best Split



# **Measure of Impurity: GINI**

• Gini Index for a given node n :

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

(NOTE: p(j | n) is the relative frequency of class j at node n).

- Maximum (1 1/n<sub>c</sub>) 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 GINI**

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

C1	0
C2	6

P(C1) = 0/6 = 0 P(C2) = 6/6 = 1Gini = 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)<sup>2</sup> - (5/6)<sup>2</sup> = 0.278

C1	2
C2	4

P(C1) = 2/6 P(C2) = 4/6Gini = 1 - (2/6)<sup>2</sup> - (4/6)<sup>2</sup> = 0.444

# **Splitting Based on GINI**

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

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

where,  $n_i =$  number of records at child i, n = number of records at node n.

# **Binary Attributes:Computing GINI Index**

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



#### **Continuous Attributes: Computing Gini Index**

- Use Binary Decisions based on one value
- Simple method to choose best v
  - For each v, scan the database to gather count matrix and compute its Gini index
  - Computationally Inefficient (may be used if the computational load is affordable)

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



#### **Alternative Splitting Criteria**

• Entropy at a given node n:

$$Entropy(n) = -\sum_{j} p(j \mid n) \log p(j \mid n)$$

(NOTE: p(j | n) is the relative frequency of class j at node n).

- Measures homogeneity of a node.
  - Maximum (log n<sub>c</sub>) when records are equally distributed among all classes implying least information
  - Minimum (0.0) when all records belong to one class, implying most information
- Entropy based computations are similar to the GINI index computations

## **Examples for computing Entropy**

$$Entropy(n) = -\sum_{j} p(j | n) \log_2 p(j | n)$$

C1	0
C2	6

P(C1) = 0/6 = 0 P(C2) = 6/6 = 1Entropy = -0 log 0 - 1 log 1 = -0 - 0 = 0

C1	1
C2	5

P(C1) = 1/6 P(C2) = 5/6 Entropy =  $-(1/6) \log_2 (1/6) - (5/6) \log_2 (1/6) = 0.65$ 

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 Information Gain**

• Information Gain:

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

Parent Node, p is split into k partitions;

n<sub>i</sub> 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)
- Disadvantage: Tends to prefer splits that result in large number of partitions, each being small but pure.

#### **Splitting Based on Information Gain**

• Gain Ratio:

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

Parent Node, p 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 number of small partitions) is penalized!
- Designed to overcome the disadvantage of Information Gain

#### **Splitting Criteria based on Classification Error**

• Classification error at a node t :

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

Measures misclassification error made by a node.

- Maximum (1 1/n<sub>c</sub>) 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 = 1Error = 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) = 4/6Error = 1 - max (2/6, 4/6) = 1 - 4/6 = 1/3

## **Comparison among Splitting Criteria**

For a 2-class problem:



#### **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 (needs an estimate of generalization)

#### **Decision Tree Based Classification**

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

#### Example: C4.5

- Simple depth-first construction.
- Uses Information Gain
- Sorts Continuous Attributes at each node.

 You can download the software from: <u>http://www.cse.unsw.edu.au/~quinlan/c4.5r8.tar.gz</u>

# The generalization problem (again)

- Occam's Razor: Prefer the simplest model that fits well to the data
- Very simple models underfit (bad generalization)
- Complex models overfit (bad generalization)
- For a given dataset and type of model, there is an optimal model complexity providing optimal generalization

## Underfitting and Overfitting (Example)



500 circular and 500 triangular data points.

Circular points:

 $0.5 \le sqrt(x_1^2 + x_2^2) \le 1$ 

Triangular points:  $sqrt(x_1^2+x_2^2) > 0.5$  or  $sqrt(x_1^2+x_2^2) < 1$ 

## **Underfitting and Overfitting**



#### **Overfitting due to Noise**



#### Decision boundary is distorted by noise point

#### **Insufficient Examples**



Lack of data points in the lower half of the diagram affects generalization: makes it difficult to predict correctly the class labels of that region

 Insufficient number of examples in the region causes the decision tree to predict the test examples using other examples that are irrelevant to the classification task

# Minimum Description Length (MDL)

- An alternative formulation of Occam's razor
- Cost(Model,Data) = Cost(Data|Model) + Cost(Model)
   Cost is the number of bits needed for encoding.
- Cost(Data|Model) encodes the misclassification errors in the dataset.
- Cost(Model) uses node encoding (number of children) plus splitting condition encoding (penalty term).
- Search for the least costly model.

## **How to Address Overfitting**

#### • Pre-Pruning (Early Stopping Rule)

- 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 expanding the current node does not improve impurity measures (e.g., Gini or information gain).

## How to Address Overfitting...

#### • Post-pruning

- Grow decision tree to its entirety
- Trim the nodes of the decision tree in a bottom-up fashion
- If generalization error improves 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
- Estimate generalization error using a validation set (hold out some training examples for validation)
- Validation set is different from test set.

# **Decision Boundary (continuous attributes)**



 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**



- Test condition may involve multiple attributes
- More expressive representation
- Finding optimal test condition is computationally expensive