

Classification Techniques (2)

Overview

Previous Lecture

- Classification Problem
- Classification based on Regression
- Distance-based Classification (KNN)

This Lecture

- Classification using Decision Trees
- Classification using Rules
- Quality of Classifiers

Data Mining Lecture 4: Classification 2

2

Classification Using Decision Trees

- A **partitioning based** technique
 - Divides the search space into rectangular regions
- Each tuple is placed into a class based on the region within which it falls
- Internal nodes associated with attribute and arcs with values for that attribute
- DT approaches differ in how the tree is built
- Algorithms: **Hunt's, ID3, C4.5, CART**

Data Mining Lecture 4: Classification 2

3

Decision Tree

Given:

- $D = \{t_1, \dots, t_n\}$ where $t_i = \langle t_{i1}, \dots, t_{in} \rangle$
- Database schema contains $\{A_1, A_2, \dots, A_n\}$
- Classes $C = \{C_1, \dots, C_m\}$

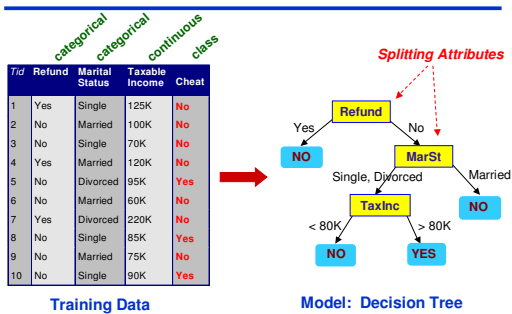
Decision or **Classification Tree** is a tree associated with D such that

- Each internal node is labeled with attribute, A_i
- Each arc is labeled with predicate which can be applied to attribute at parent
- Each leaf node is labeled with a class, C_j

Data Mining Lecture 4: Classification 2

4

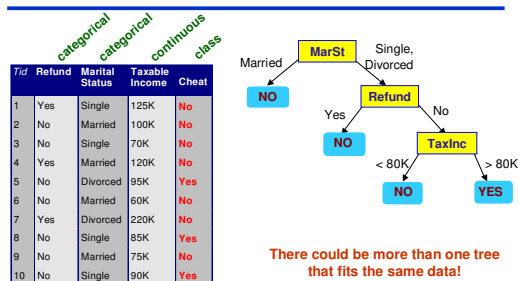
Example of a Decision Tree



Data Mining Lecture 4: Classification 2

5

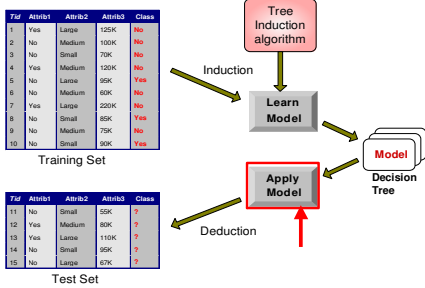
Another Example of Decision Tree



Data Mining Lecture 4: Classification 2

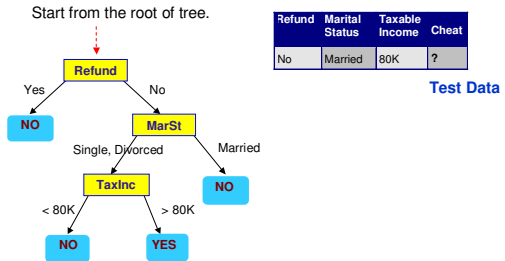
6

Decision Tree Classification Task



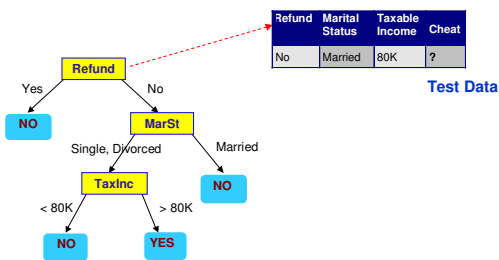
Data Mining Lecture 4: Classification 2 7

Apply Model to Test Data



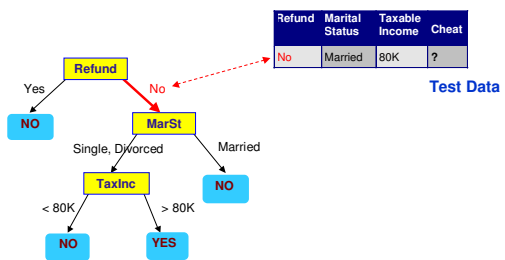
Data Mining Lecture 4: Classification 2 8

Apply Model to Test Data



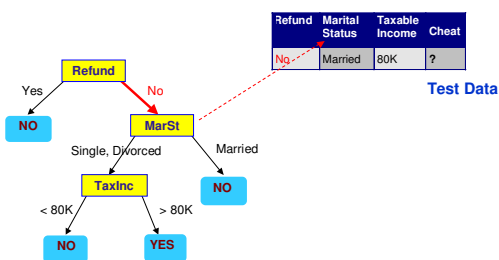
Data Mining Lecture 4: Classification 2 9

Apply Model to Test Data



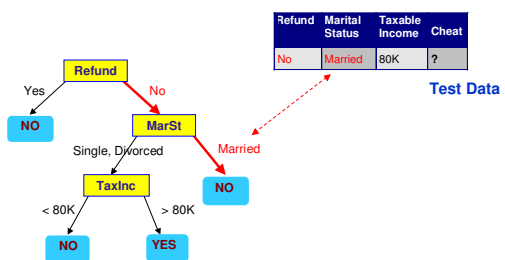
Data Mining Lecture 4: Classification 2 10

Apply Model to Test Data



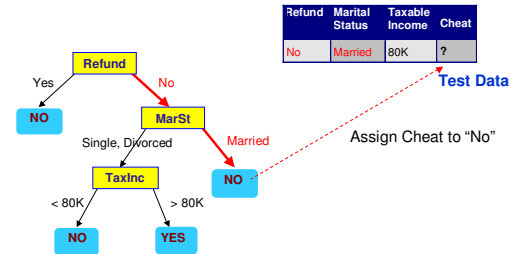
Data Mining Lecture 4: Classification 2 11

Apply Model to Test Data



Data Mining Lecture 4: Classification 2 12

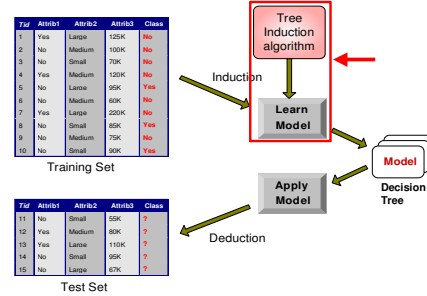
Apply Model to Test Data



Data Mining Lecture 4: Classification 2

13

Decision Tree Classification Task

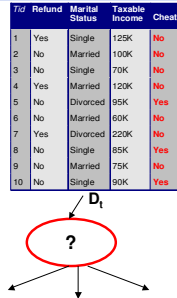


Data Mining Lecture 4: Classification 2

14

General Structure of Hunt's Algorithm

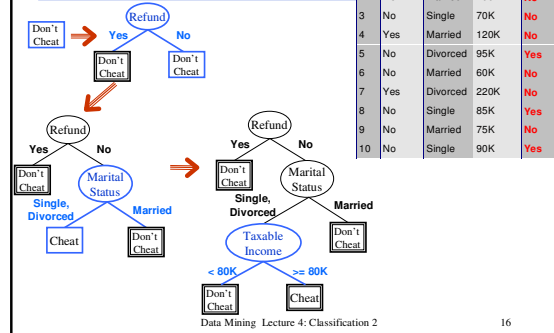
- Let D_t be the set of training records that reach a node t
- General Procedure:
 - If D_t contains records that belong to the same class y_t , then t is a leaf node labeled as y_t .
 - If D_t is an empty set, then t is a leaf node labeled by the default class, y_d .
 - If D_t contains records that belong to more than one class, then use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.



Data Mining Lecture 4: Classification 2

15

Hunt's Algorithm



Data Mining Lecture 4: Classification 2

16

Decision Tree Induction

```

Input:
D //Training data
Output:
T //Decision Tree
DTBuild Algorithm:
//Simplistic algorithm to illustrate naive approach to building DT
T = {};
Determine best splitting criterion;
T = Create root node and label with splitting attribute;
T = Add arc to root node for each split predicate and label;
for each arc do
    D = Database created by applying splitting predicate to D;
    if stopping point reached for this path then
        T' = Create leaf node and label with appropriate class;
    else
        T' = DTBuild(D);
    T = Add T' to arc;
    
```

Data Mining Lecture 4: Classification 2

17

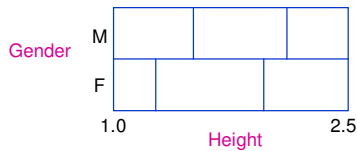
Decision Tree Induction

- 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

Data Mining Lecture 4: Classification 2

18

DT Split Areas



Data Mining Lecture 4: Classification 2

19

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

Data Mining Lecture 4: Classification 2

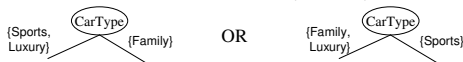
20

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.

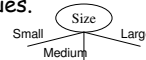


Data Mining Lecture 4: Classification 2

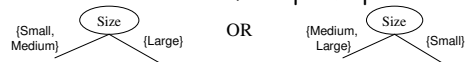
21

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.



- What about this split?



Data Mining Lecture 4: Classification 2

22

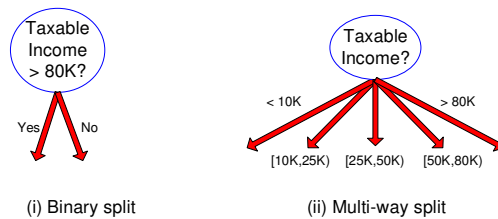
Splitting Based on Continuous Attributes

- 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 $(A \geq v)$
 - considers all possible splits and finds the best cut
 - can be more compute intensive

Data Mining Lecture 4: Classification 2

23

Splitting Based on Continuous Attributes



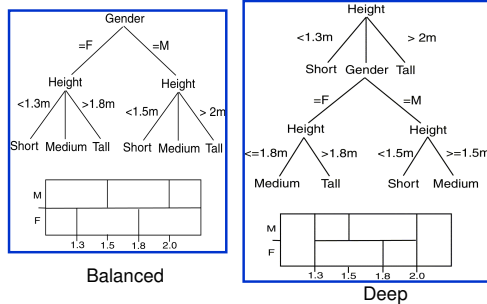
(i) Binary split

(ii) Multi-way split

Data Mining Lecture 4: Classification 2

24

Comparing Decision Trees



Data Mining Lecture 4: Classification 2

25

DT Induction Issues that affect Performance

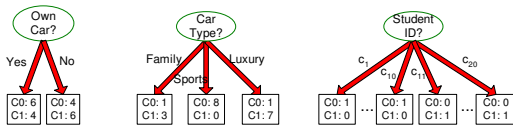
- Choosing Splitting Attributes
- Ordering of Splitting Attributes
- Split Points
- Tree Structure
- Stopping Criteria
- Training Data (size of)
- Pruning

Data Mining Lecture 4: Classification 2

26

How to determine the Best Split

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



Which test condition is the best?

Data Mining Lecture 4: Classification 2

27

How to determine the Best Split

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

C0: 5
C1: 5

Non-homogeneous,
High degree of impurity

C0: 9
C1: 1

Homogeneous,
Low degree of impurity

Data Mining Lecture 4: Classification 2

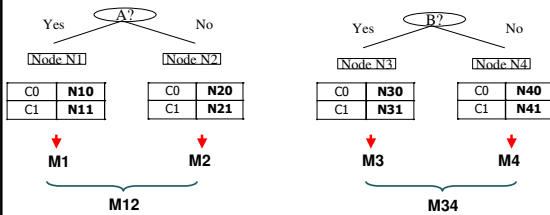
28

How to Find the Best Split

Before Splitting:

C0	N00
C1	N01

 → M0



Gain = M0 - M12 vs M0 - M34

Data Mining Lecture 4: Classification 2

29

Measure of Impurity: GINI

- Gini Index for a given node t :

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

(NOTE: $p(j|t)$ 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	

Data Mining Lecture 4: Classification 2

30

Examples for computing GINI

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

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

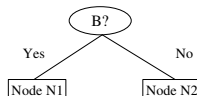
- Used in **CART**
- When a node p is split into k partitions (children), the quality of split is computed as:

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

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	

Gini(N1)
 = 1 - (5/6)² - (2/6)²
 = 0.194

Gini(N2)
 = 1 - (1/6)² - (4/6)²
 = 0.528

N1		N2	
C1	5	1	
C2	2	4	
Gini=0.333			

Gini(Children)
 = 7/12 * 0.194 +
 5/12 * 0.528
 = 0.333

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		CarType	
	[Sports, Luxury]	(Family)	(Sports, Luxury]	(Family, Luxury]
C1	3	1	2	2
C2	2	4	1	5
Gini	0.400		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.

T/A	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



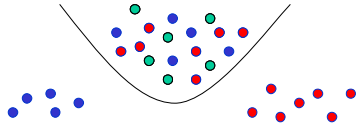
Continuous Attributes: Computing GINI Index

- 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	Taxable Income																			
	No	No	No	Yes	Yes	Yes	No	No	No	No										
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	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	4	3	5	2	6	1	7	0
Gini	0.420	0.400	0.375	0.345	0.417	0.400	0.390	0.345	0.375	0.400	0.420									

Information

Decision Tree Induction is often based on Information Theory



Data Mining Lecture 4: Classification 2

37

DT Induction

- When all the marbles in the bowl are mixed up, little information is given.
- When the marbles in the bowl are all from one class and those in the other two classes are on either side, more information is given.

Use this approach with DT Induction !

Data Mining Lecture 4: Classification 2

38

Information/Entropy

Given probabilities p_1, p_2, \dots, p_s whose sum is 1, Entropy is defined as:

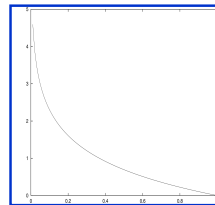
$$H(p_1, p_2, \dots, p_s) = \sum_{i=1}^s (p_i \log(1/p_i))$$

- Entropy measures the amount of randomness or surprise or uncertainty.
- Goal in classification
 - no surprise
 - entropy = 0

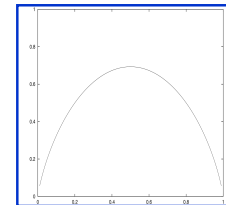
Data Mining Lecture 4: Classification 2

39

Entropy



log (1/p)



H(p, 1-p)

Data Mining Lecture 4: Classification 2

40

ID3

- Creates a decision tree using information theory concepts and tries to reduce the expected number of comparisons.
- ID3 chooses to split on an attribute that gives the highest information gain:

$$Gain(D, S) = H(D) - \sum_{i=1}^s P(D_i) H(D_i)$$

Data Mining Lecture 4: Classification 2

41

Height Example Data

Name	Gender	Height	Output1	Output2
Kristina	F	1.60	Short	Medium
Jim	M	2.02	Tall	Medium
Maggie	F	1.90	Medium	Tall
Martha	F	1.88	Medium	Tall
Stephanie	F	1.71	Short	Medium
Bob	M	1.85	Medium	Medium
Kathy	F	1.60	Short	Medium
Dave	M	1.72	Short	Medium
Worth	M	2.12	Tall	Tall
Steven	M	2.10	Tall	Tall
Debbie	F	1.78	Medium	Medium
Todd	M	1.95	Medium	Medium
Kim	F	1.89	Medium	Tall
Amy	F	1.81	Medium	Medium
Wynette	F	1.75	Medium	Medium

Data Mining Lecture 4: Classification 2

42

ID3 Example (Output1)

- Starting state entropy:
 $4/15 \log(15/4) + 8/15 \log(15/8) + 3/15 \log(15/3) = 0.4384$
- Gain using gender:
 - Female: $3/9 \log(9/3) + 6/9 \log(9/6) = 0.2764$
 - Male: $1/6 (\log 6/1) + 2/6 \log(6/2) + 3/6 \log(6/3) = 0.4392$
 - Weighted sum: $(9/15)(0.2764) + (6/15)(0.4392) = 0.34152$
 - Gain: $0.4384 - 0.34152 = 0.09688$
- Gain using height:
 $0.4384 - (2/15)(0.301) = 0.3983$
- Choose height as first splitting attribute

C4.5 Algorithm

- ID3 favors attributes with large number of divisions (is vulnerable to overfitting)
- Improved version of ID3:
 - Missing Data
 - Continuous Data
 - Pruning
 - Rules
 - GainRatio:

$$GainRatio(D, S) = \frac{Gain(D, S)}{H\left(\frac{|D_1|}{|D|}, \dots, \frac{|D_n|}{|D|}\right)}$$

- Takes into account the cardinality of each split area

CART: Classification and Regression Trees

- Creates a Binary Tree
- Uses entropy to choose the best splitting attribute and point
- Formula to choose split point, s , for node t :

$$\Phi(s/t) = 2P_L P_R \sum_{j=1}^m |P(C_j | t_L) - P(C_j | t_R)|$$

- P_L, P_R probability that a tuple in the training set will be on the left or right side of the tree.

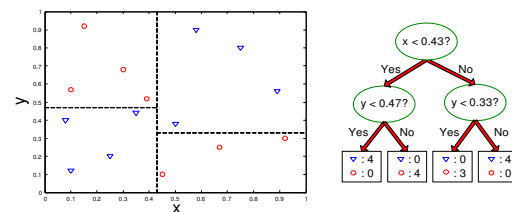
CART Example

- At the start, there are six choices for split point (right branch on equality):
 - $\Phi(\text{Gender}) = 2(6/15)(9/15)(2/15 + 4/15 + 3/15) = 0.224$
 - $\Phi(1.6) = 0$
 - $\Phi(1.7) = 2(2/15)(13/15)(0 + 8/15 + 3/15) = 0.169$
 - $\Phi(1.8) = 2(5/15)(10/15)(4/15 + 6/15 + 3/15) = 0.385$
 - $\Phi(1.9) = 2(9/15)(6/15)(4/15 + 2/15 + 3/15) = 0.256$
 - $\Phi(2.0) = 2(12/15)(3/15)(4/15 + 8/15 + 3/15) = 0.32$
- Split at 1.8

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

Decision Boundary



- 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

$x + y < 1$

Class = + Class = -

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

Data Mining Lecture 4: Classification 2 49

Tree Replication

- Same subtree appears in multiple branches

Data Mining Lecture 4: Classification 2 50

Classification Using Rules

- Perform classification using If-Then rules
- **Classification Rule:** $r = \langle a, c \rangle$
 - Antecedent, Consequent
- May generate rules from other techniques (DT, NN) or generate directly.
- Algorithms: **Gen, RX, 1R, PRISM**

Data Mining Lecture 4: Classification 2 51

Generating Rules from Decision Trees

```

Input:
  T //Decision Tree
Output:
  R //Rules
Gen Algorithm:
  //Illustrate simple approach to generating classification rules from a DT
  R = {}
  for each path from root to a leaf in T do
    a = True
    for each internal node do
      a = a ∧ (label of parent node in path combined with label of incident arc)
      c = label of leaf node
    R = R ∪ r = < a, c >
  
```

Data Mining Lecture 4: Classification 2 52

Generating Rules Example

```

{ < (Height ≤ 1.6m), Short >,
  < ((Height > 1.6m) ∧ (Height ≤ 1.7m)), Short >,
  < ((Height > 1.7m) ∧ (Height ≤ 1.8m)), Medium >,
  < ((Height > 1.8m) ∧ (Height ≤ 1.9m)), Medium >,
  < ((Height > 1.9m) ∧ (Height ≤ 2m) ∧ (Height ≤ 1.95m)), Medium >,
  < ((Height > 1.9m) ∧ (Height ≤ 2m) ∧ (Height > 1.95m)), Tall >,
  < (Height > 2m), Tall > }
  
```

Data Mining Lecture 4: Classification 2 53

1R Algorithm

```

Input:
  D //Training data
  R //Attributes to consider for rules
  C //Classes
Output:
  R //Rules
1R Algorithm:
  //1R algorithm generates rules based on one attribute
  R = {};
  for each A ∈ R do
    R_A = {};
    for each possible value, v_i, of A do
      //v_i may be a range rather than a specific value
      for each C_j ∈ C, find count(C_j);
      // Here count is the number of occurrences of this class for this attribute
      let C_n be the class with the largest count;
      R_A = R_A ∪ ((A = v_i) → (class = C_n));
    ERR_A = number of tuples incorrectly classified by R_A;
  R = R_n where ERR_A is minimum;
  
```

Data Mining Lecture 4: Classification 2 54

1R Example

Option	Attribute	Rules	Errors	Total Errors
1	Gender	F → Medium	3/9	6/15
		M → Tall	3/6	
2	Height	(0,1.6] → Short	0/2	1/15
		(1.6,1.7] → Short	0/2	
		(1.7,1.8] → Medium	0/3	
		(1.8,1.9] → Medium	0/4	
		(1.9,2.0] → Medium	1/2	
		(2.0,∞) → Tall	0/2	

Data Mining Lecture 4: Classification 2

55

PRISM Algorithm

```

Algorithm 0.1
Input:
D // Training data
C // Classes
Output:
R // Rules
PRISM Algorithm:
// PRISM algorithm generates rules based on best attribute-value pairs
R = ∅;
for each Cj ∈ C do
  repeat
    T = D; // All instances of class Cj will be systematically removed from T
    p = true; // Create new rule with empty left hand side
    r = ( ! p then Cj );
    repeat
      for each attribute A value v pair found in T do
        calculate ( |{tuples ∈ T with A=values(Cj)}| );
        find A = v that maximizes this value;
        p = p ∧ (A = v);
        T = {tuples in T that satisfy A = v};
      until all tuples in T belong to Cj;
    D = D - T;
    R = R ∪ r;
  until there are no tuples in D which belong to Cj;
    
```

Data Mining Lecture 4: Classification 2

56

PRISM Example

Gender = F 0/9
 Gender = M 3/6

If Gender = M then Class = Tall.
 If Gender = M and Height in range ? then Class = Tall.

Height <= 1.6 0/0
 1.6 < Height <= 1.7 0/1
 1.7 < Height <= 1.8 0/0
 1.8 < Height <= 1.9 0/1
 1.9 < Height <= 2.0 1/2
 2.0 < Height 2/2

If Gender = M and Height > 2 then Class = Tall.

Data Mining Lecture 4: Classification 2

57

Decision Tree vs. Rules

- Tree has an implied order in which splitting is performed.
- Tree is created based on looking at all classes.
- Rules have no ordering of predicates.
- Only need to look at one class to generate its rules.

Data Mining Lecture 4: Classification 2

58

Metrics for Performance Evaluation...

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

- Most widely-used metric:

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

Data Mining Lecture 4: Classification 2

59

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

Data Mining Lecture 4: Classification 2

60

Estimating Classifier Accuracy

IDEA: Randomly select sampled partitions of the training data to estimate accuracy

- **Holdout method:**
 - Partition known data into two independent sets
 - Training set (usually 2/3 of data)
 - Test set (remaining 1/3)
 - Estimate of the accuracy of classifier is pessimistic
- **Random sub-sampling:**
 - Repeat the holdout method k times;
 - Overall accuracy estimate is taken as the average estimates obtained by the process.

Data Mining Lecture 4: Classification 2

61

Estimating Classifier Accuracy

- **K-fold cross-validation:**
 - Partition known data S , into k mutually exclusive subsets (or "folds") S_1, S_2, \dots, S_k of approximately equal size;
 - Use each S_i as a test set
 - Accuracy estimate is the overall number of correct classifications divided by the total number of samples in the initial data
- **Leave-one-out:**
 - K-fold cross-validation with k set to $|S|$.

Data Mining Lecture 4: Classification 2

62

Increasing Classifier Accuracy

Bagging:

- each classifier "votes";
- winner class wins classification.

Boosting:

- each classifier "votes";
- votes are combined based on weights obtained by the estimates of each classifier's accuracy;
- winner class wins classification.

Data Mining Lecture 4: Classification 2

63

Is Accuracy enough to judge a Classifier?

In practice, there are also other considerations

- Speed
- Robustness (influence of noisy data)
- Scalability (number of I/O operations)
- Interpretability of classification output

Data Mining Lecture 4: Classification 2

64