CSE4334/5334 Data Mining Classification: Decision Tree

Chengkai Li

Department of Computer Science and Engineering University of Texas at Arlington Fall 2020 (Slides courtesy of Pang-Ning Tan, Michael Steinbach and Vipin Kumar, and Jiawei Han, Micheline Kamber and Jian Pei)

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

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 vs. Prediction

Classification

- o Predicts categorical class labels
- Most suited for nominal attributes
- o Less effective for ordinal attributes

Prediction

- models continuous-valued functions or ordinal attributes,
 i.e., predicts unknown or missing values
- o E.g., Linear regression



Supervised vs. Unsupervised Learning

Supervised learning (e.g., classification)

- Supervision: The training data (observations, measurements, etc.) are accompanied by labels indicating the class of the observations
- New data is classified based on the training set

Unsupervised learning (e.g., clustering)

- The class labels of training data is unknown
- Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

Classification Techniques

- Decision Tree based Methods
- Rule-based Methods
- Memory based reasoning
- Neural Networks
- Naïve Bayes and Bayesian Belief Networks
- Support Vector Machines



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



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

Decision Tree Classification Task



Test Set

Yes

No

No

Yes

No

No

Yes

No

No

No

Refund

1

2

3

4

5

6

7

8

9

10

No

No

No

No

Yes

Yes

Apply Model to Test Data



Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

Apply Model to Test Data





Test Data





Apply Model to Test Data Test Data Refund Marital Taxable Status Cheat Income 80K ? No Married Refund Yes No NO MarSt Married Single, Divorced TaxInc NO < 80K >= 80K YES NO

Apply Model to Test Data Test Data Refund Marital Taxable Status Cheat Income 80K No Married ? Refund Yes No NO MarSt Assign Cheat to "No" Married Single, Divorced TaxInc NO < 80K >= 80K YES NO

Tid

1

2

3

4

5

6

7

8

9

10

No

No

No

No

Yes

Divorced

Sinale

Married

Test Set

90K

40K

80K

?

? ?

Yes

Decision Tree Classification Task







Decision Tree Induction



Large search space

- Exponential size, with respect to the set of attributes
- Finding the optimal decision tree is computationally infeasible

Efficient algorithm for accurate suboptimal decision tree

• Greedy strategy

• Grow the tree by making locally optimally decisions in selecting the attributes

Decision Tree Induction



Many Algorithms:

0 Hunt's Algorithm (one of the earliest)

• CART

o ID3, C4.5

○ SLIQ,SPRINT

General Structure of Hunt's Algorithm

- $\circ~$ 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 majority class among the records of Dt's parent node.
 - If D_t contains records that have identical values on all attributes but the class attribute, then t is a leaf node labeled by the majority class among D_t 's records.
 - If none of the above conditions is satisfied, 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



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?
- o Determine when to stop splitting



How to Specify Test Condition?

- Depends on attribute types
- o Categorical vs. Numeric
 - Categorical attributes: Nominal, Ordinal
 - Numeric attributes: Interval, Ratio
- o Discrete vs. Continuous
- Depends on number of ways to split
- o 2-way split
- 0 Multi-way split

Splitting Based on Nominal Attributes



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

Sports





Splitting Based on Continuous Attribute

- 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.
- o Binary Decision: $(A \le v)$ or $(A \ge v)$
 - consider all possible splits and finds the best cut
 - can be more compute intensive

Splitting Based on Continuous Attribute



(i) Binary split

(ii) Multi-way split

Tree Induction



Greedy strategy.

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

Issues

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



Which test condition is the best?





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

C0: 5	
C1: 5	

C0: 9 C1: 1

Non-homogeneous,

High degree of impurity

Homogeneous,

Low degree of impurity







Misclassification error





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).

- \circ 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		C1	1	I	C1	2	C1	3
C2	6		C2	5		C2	4	C2	3
Gini=0.000		Gini=	0.278		Gini=	0.444	Gini=	0.500	



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

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)	$)^{2} - (5/6)^{2} = 0.278$
P(C1) = 2/6	P(C2) = 4/6

C1	2
C2	4

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

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,

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

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

Binary Attributes: Computing GINI Index

B?

No

Node N2

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

Yes

Node N1

	Farent
C1	6
C2	6
Gini	= 0.500
-	

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



amily uxury 2 5

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

Multi-way split

Tw	o-way sp	olit
find best	partition	of values

	CarType									
	Family Sports Luxu									
C1	1	2	1							
C2	4 1 1									
Gini	0.393									

	CarT	уре		CarT	У
	{Sports, Luxury}	{Family}		{Sports}	{ L
C1	3	1	C1	2	
C2	2	4	C2	1	
Gini	0.4	00	Gini	0.4	1

Continuous Attributes: Computing Gini Index



- o 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, $A \le v$ and A > v
- o Simple method to choose best v
 - For each v, scan the database to gather count matrix and compute its Gini index
 - o 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



Continuous Attributes: Computing Gini Index...



For efficient computation: for each attribute,

- o 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		Nc)	N	0	Ye	s	Ye	s	Ye	es	N	•	N	ο	N	0		No	
			Taxable Income																				
			60	ĺ	70)	7	5	85	5	90)	9	5	10	00	12	20	1:	25		220	
Sorted Values	\rightarrow	5	5	6	5	7	2	8	0	8	7	9	2	9	7	11	0	12	22	17	72	23	0
Split Positions		<=	>	<=	>	 	>	 	>	 	>	<=	>	<	>	<=	>	<=	>	<=	>	<=	>
	Yes	0	3	0	3	0	3	0	3	1	2	2	1	3	0	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 <u>0.300</u> 0.343 0.375									375	0.4	00	0.4	20									



Entropy at a given node t:

41

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

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

• Measures homogeneity of a node.

- Maximum (log n_c) 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



$$Entropy(t) = -\sum_{j} p(j \mid t) \log_2 p(j \mid t)$$

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) = 5/6Entropy = - (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 = $-(2/6) \log_2 (2/6) - (4/6) \log_2 (4/6) = 0.92$



Why is that $0 \log 0 = 0$?

$$\lim_{x \to 0} x \log_2(x) = \lim_{x \to 0} \frac{\frac{\ln(x)}{\ln(2)}}{x^{-1}} = \lim_{x \to 0} \frac{\frac{x^{-1}}{\ln(2)}}{-x^{-2}} = \lim_{x \to 0} \frac{-x}{\ln(2)} = 0$$

L'Hospital's Rule (Wikipedia)

$$\begin{split} &\lim_{x \to c} f(x) = \lim_{x \to c} g(x) = 0 \text{ or } \pm \infty \text{, and} \\ &\lim_{x \to c} \frac{f'(x)}{g'(x)} \text{ exists, and} \\ &g'(x) \neq 0 \text{ for all } x \text{ in } I \text{ with } x \neq c \text{,} \end{split}$$

then

If

$$\lim_{x \to c} \frac{f(x)}{g(x)} = \lim_{x \to c} \frac{f'(x)}{g'(x)}.$$

Splitting Based on INFO ...



Information Gain:

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

Parent Node, p 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)
- o Used in ID3 and C4.5
- Disadvantage: Tends to prefer splits that result in large number of partitions, each being small but pure.

Splitting Based on INFO ...



Gain Ratio:

 $GainRATIO_{split} = \frac{GAIN_{Split}}{SplitINFC}$

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

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!
- Used in C4.5
- 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 (0.0) 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:



Misclassification Error vs Gini



	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=0.342
 3
 3

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

Gini improves !!



Tree Induction



Greedy strategy.

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

Issues

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

Stopping Criteria for Tree Induction

- Stop expanding a node (i.e., making the node a leaf node) when one of the following conditions is met:
- 1) The (information) gain of best split is below a threshold (e.g., all or most of the records in the node belong to the same class).
- 2) All or most of the records have similar attribute values.
- 3) The number of records in the node is below a threshold. (Extreme case: the node is empty.)
- Note that these criteria are generalization of the special cases in the Hunt's algorithm.
- What to give to the node as the class attribute value prediction? Majority voting based on the classes of the records in the node or in its parent node if it doesn't have enough records.



Advantages:

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