

Data Mining

Classification: Basic Concepts, Decision Trees, and Model Evaluation

Lecture Notes for Chapter 4 Part II

Introduction to Data Mining by

Tan, Steinbach, Kumar

Adapted by Qiang Yang (2010)

Continuous Attribute: Binary Split for Temperature?

Outlook	Tempreature	Humidity	Windy	Class
Sunny	40	high	false	N
sunny	37	high	true	N
overcast	34	high	false	P
rain	26	high	false	P
rain	15	normal	false	P
rain	13	normal	true	N
overcast	17	normal	true	P
sunny	28	high	false	N
sunny	25	normal	false	P
rain	23	normal	false	P
sunny	27	normal	true	P
overcast	22	high	true	P
overcast	40	normal	false	P
rain	31	high	true	N

Finding the best split

- | Sort the Temperature attribute
- | For each possible binary split, calculate the information gain
 - That is, calculate the entropy: $-p(P) \cdot \log p(P) - p(N) \cdot \log p(N)$
 - Select the smallest one
- | Let the value be L. Two branches:
Temperature < L, and Temperature >= L.

Measures of Node Impurity

- | Gini Index
- | Entropy (already covered)
- | Misclassification error

How to Find the Best Split: let M be the measure

Before Splitting:

C0	N00
C1	N01

→ M0

A?

Yes

No

Node N1

Node N2

C0

N10

C1

N11

C0

N20

C1

N21

M1

M2

M12

B?

Yes

No

Node N3

Node N4

C0

N30

C1

N31

C0

N40

C1

N41

M3

M4

M34

Gain = M0 – M12 vs M0 – M34

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	

Examples for computing GINI

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

C1	0
C2	6

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Gini = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$$

C1	1
C2	5

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

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

C1	2
C2	4

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

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

Splitting Based on GINI

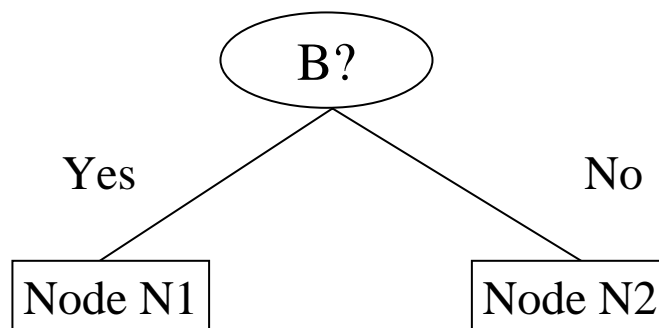
- | Used in CART, SLIQ, SPRINT.
- | When a node p is split into k partitions (children), the quality of split is computed as,

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

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

Binary Attributes: Computing GINI Index

- | 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)

$$= 1 - (5/7)^2 - (2/7)^2$$
$$= 0.194$$

Gini(N2)

$$= 1 - (1/5)^2 - (4/5)^2$$
$$= 0.528$$

	N1	N2
C1	5	1
C2	2	4
Gini=0.333		

Gini(Children)

$$= 7/12 * 0.194 +$$
$$5/12 * 0.528$$
$$= 0.333$$

Multi-way Splits: Computing Gini Index

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

Multi-way split

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

Two-way split
(find best partition of values)

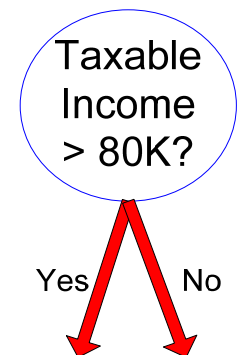
	CarType	
	{Sports, Luxury}	{Family}
C1	3	1
C2	2	4
Gini	0.400	

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

Continuous Attributes: Computing Gini Index

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

<i>Tid</i>	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	No		No		No		Yes		Yes		Yes		No		No		No		No			
			Taxable Income																					
Sorted Values Split Positions	→		60		70		75		85		90		95		100		120		125		220			
			55		65		72		80		87		92		97		110		122		172		230	
			<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>
		Yes	0	3	0	3	0	3	0	3	1	2	2	1	3	0	3	0	3	0	3	0	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		0.400		0.420	

Training Set: Build a Decision Tree 1

Outlook	Tempreature	Humidity	Windy	Class
Sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

Classification Error

- | Classification error at a node t :

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

- | 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

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

C1	0
C2	6

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Error = 1 - \max(0, 1) = 1 - 1 = 0$$

C1	1
C2	5

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$Error = 1 - \max(1/6, 5/6) = 1 - 5/6 = 1/6$$

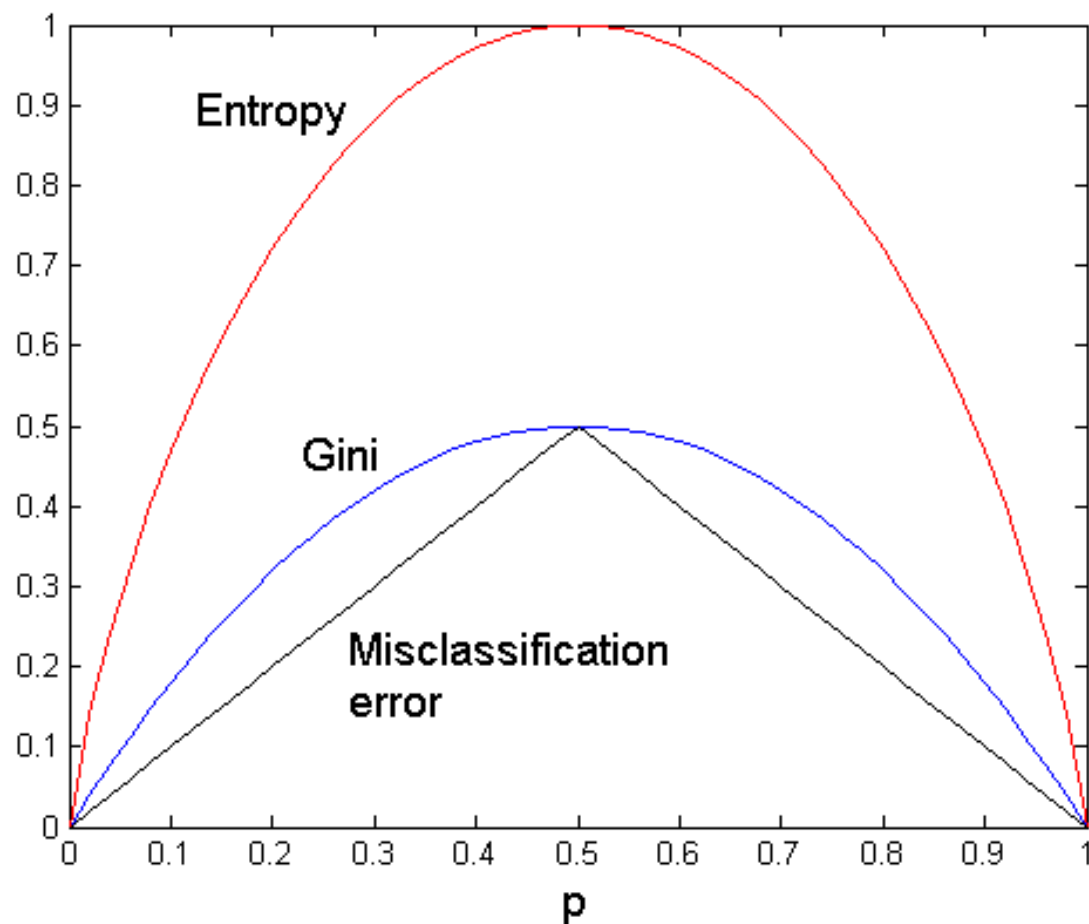
C1	2
C2	4

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$Error = 1 - \max(2/6, 4/6) = 1 - 4/6 = 1/3$$

Comparison among Splitting Criteria

For a 2-class problem:



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 (to be discussed later)

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

Example: C4.5

- | Simple depth-first construction.
- | Uses Information Gain
- | Sorts Continuous Attributes at each node.
- | Needs entire data to fit in memory.
- | Unsuitable for Large Datasets.
 - Needs out-of-core sorting.
- | You can download the software from:
<http://www.cse.unsw.edu.au/~quinlan/c4.5r8.tar.gz>