

Brian Holliday
Professor

Intro to Data Mining
30 March 2020

Project 5 Part 1
GINI Index

①

B	+	-	$P(t_i)$	$GINI(t_i)$
$t_1 = T$	3	1	4/10	$1 - (3/4)^2 - (1/4)^2 = 3/4$
$t_2 = F$	1	5	6/10	$1 - (1/6)^2 - (5/6)^2 = 5/8$

② $GINI_{split} = \sum_{i=1}^n P(t_i) \cdot GINI(t_i) = (4/10)(3/4) + (6/10)(5/8)$
 $= 0.71667$

A	+	-	$P(t_i)$	$GINI(t_i)$
$t_1 = T$	4	3	7/10	$1 - (4/7)^2 - (3/7)^2 = 25/49$
$t_2 = F$	0	3	3/10	$1 - (0/3)^2 - (3/3)^2 = 0$

$GINI_{split} = \sum P(t_i) \cdot GINI(t_i) = (7/10)(25/49) + (3/10)(0)$
 $= \frac{7}{10} \cdot \frac{25}{49} = \frac{175}{490} = 0.3571$

We would choose B because the GINI Index is smaller therefore closer to having pure nodes
 Misclassification Error

②

B	+	-	$P(t_i)$	Error(t_i)
$t_1 = T$	3	1	4/10	$1 - 1/3 = 2/3$
$t_2 = F$	1	5	6/10	$1 - 5/6 = 1/6$

Error Split = $(4/10)(2/3) + (6/10)(1/6) = 0.36667$

A	+	-	$P(t_i)$	Error(t_i)
$t_1 = T$	4	3	7/10	$1 - (3/7) = 4/7$
$t_2 = F$	0	3	3/10	$1 - 0/3 = 1$

Error = $(7/10)(4/7) + (3/10)(1) = 0.40$

Based on the classification error we would choose B

3

~~3~~

A

A: A: 10/10

True 7/10

False 3/10

$\frac{4}{7}^+$

$\frac{3}{7}^-$

$\frac{0}{3}^+$

$\frac{3}{3}^-$

B

10/10

True 4/10

False 6/10

$\frac{3}{4}^+$

$\frac{1}{4}^-$

$\frac{1}{6}^+$

$\frac{5}{6}^-$