

*Journal of Korean
Data & Information Science Society
2002, Vol. 13, No.1 pp. 97-104*

A Study of Combined Splitting Rules in Regression Trees¹⁾

Yung-Seop Lee²⁾

Abstract

Regression trees, a technique in data mining, are constructed by splitting function-a independent variable and its threshold. Lee(2002) considered one-sided purity(OSP) and one-sided extreme(OSE) splitting criteria for finding a interesting node as early as possible. But these methods cannot be crossed each other in the same tree. They are just concentrated on OSP or OSE separately in advance.

In this paper, a new splitting method, which is the combination and extension of OSP and OSE, is proposed. By these combined criteria, we can select the nodes by considering both pure and extreme in the same tree. These criteria are not the generalized one of the previous criteria but another option depending on the circumstance.

: , , , CART .

1.

가 , ,
가 . 1990 CRM 가 ,
CRM 가 .
가 . , 가 ,
가 가 .

1. This work is supported by the Dongguk University research fund.
2. Lecturer, Department of Statistics, Dongguk University, Seoul, 110-715, Korea.
E-mail: yung@dongguk.edu

CART (Breiman et al., 1984)가
 CHAID (Hatigan, 1975), CART (Breiman et al., 1984), C4.5 (Quinlan, 1993)가
 CHAID (Chi-Square Automatic Interaction Detection) AID (Morgan and Sonquist, 1963)
 (multiway splits)가 가 (pre-pruning).

C4.5 (machine learning) χ^2 (binary splits),
 가 (post-pruning) (entropy)

가 CART, C4.5 가 post-pruning (Gini index)
 (classification trees),

(regression trees)
 $i(O)$, $i(L)$, $i(R)$
 n, n_L, n_R , Δi (1)

$$\Delta i = i(O) - \left(\frac{n_L}{n} i(L) + \frac{n_R}{n} i(R) \right) \quad (1)$$

가 가 (, 가 Δi)
 CART (splitting method) (1)

가 가 (threshold)
 (residual sum of squares)

Lee(2001)

Lee(2002) , Lee(2002)가

2.

Lee(2002)

가

가 가

, CART (Breiman et al., 1984), Bagging (Breiman, 1996), AdaBoost (Freund and Schapire, 1996), Bumping (Tibshirani and Knight, 1996)

Lee(2002)

가

1

(Y ,)

(X_1 , (X_2),

(X_3 ,

(X_4)

가

(CART) Lee(2002)

가

1

1

CART

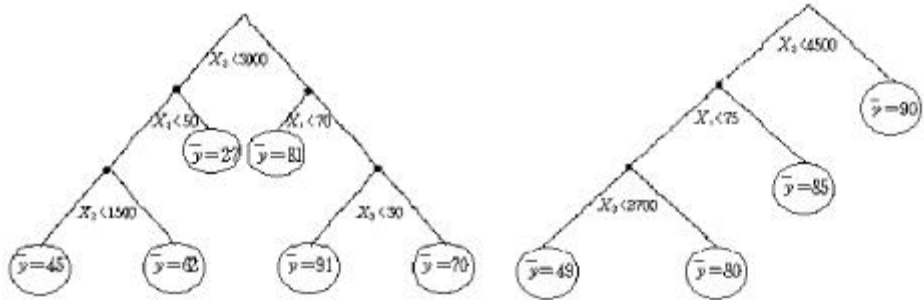
(one-sided extremes, OSE)

(R^2)

가

CART가 OSE

가



CART

OSE

1:

1 CART \bar{y} 가 가

X_2 가 3000 , X_1

70, $X_3 = 30$ 가 $(\bar{y} = 91)$ 가
 OSE X_2 가 4500 가
 $(\bar{y} = 90)$ 가
 Lee(2002) (one-sided purity, OSP)
 OSE 가
 OSP 가 OSE

3.

$$1) \min \left(\frac{\hat{\mu}_L}{\hat{\sigma}_L^2}, \frac{\hat{\mu}_R}{\hat{\sigma}_R^2} \right)$$

$$2) \min \left(\frac{\hat{\sigma}_L^2}{\hat{\mu}_L}, \frac{\hat{\sigma}_R^2}{\hat{\mu}_R} \right)$$

1) 2) (X) (MLE)
 1) 가 2) 가
 가 가
 가 가
 가 가
 (coefficient of variation, CV) 2)

4.

CART (Breiman et al., 1984)

Lee(2002) OSP OSE
 506 (tract) 13
 (median)
 1
 < 4.1> Lee(2002) < 3.2>

OSP , 가
 . < 4.2> Lee(2002) OSE $\max(\hat{\mu}_L, \hat{\mu}_R)$
 가 가 . Lee(2002)

1.

CRIM	
ZN	25,000
INDUS	
CHAS	(: 1, : 0)
NOX	(pphm)
RM	
AGE	1940
DIS	5
RAD	
TAX	(\$10,000)
PTRATIO	
B	
LATAT	(%)
MEDV	가 ()

4.3> < 4.4> m (506) 1) 2) sz (%) . <
 .
 < 4.3> $R^2 = 0.22$
 . < 4.4> $R^2 = 0.68$ CART, OSE, OSP
 , $m = 17$ < 4.1> OSP
 가 . $m = 45.10$ <
 4.2> OSE 가
 () . OSP OSE가
 가 . Lee(2002) OSP OSE

< 4.4>

가

가

5.

가

가

Lee(2002)

OSP OSE

OSP OSE

가

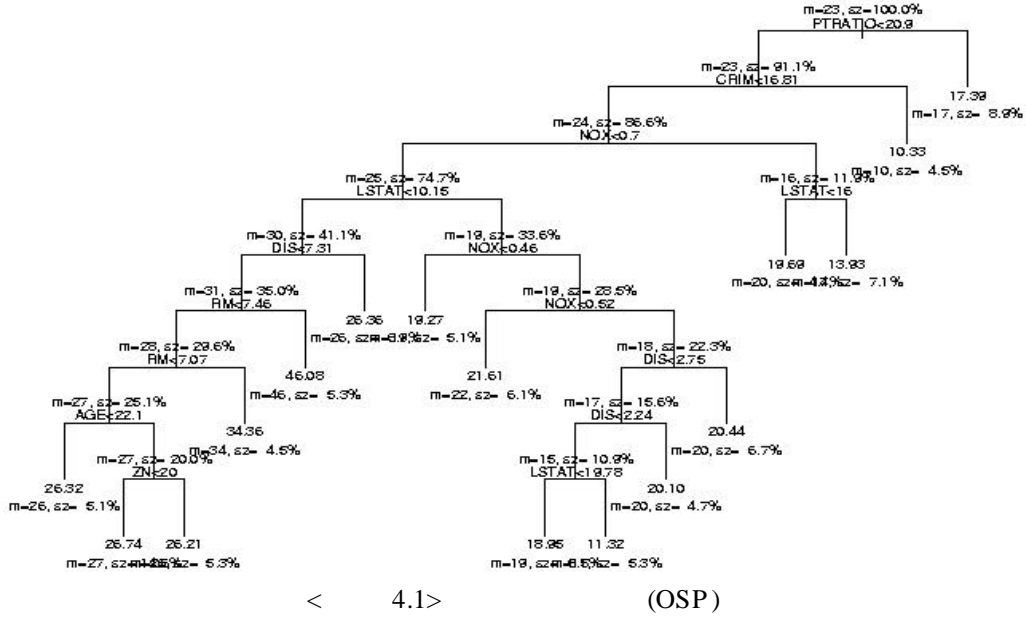
가

가

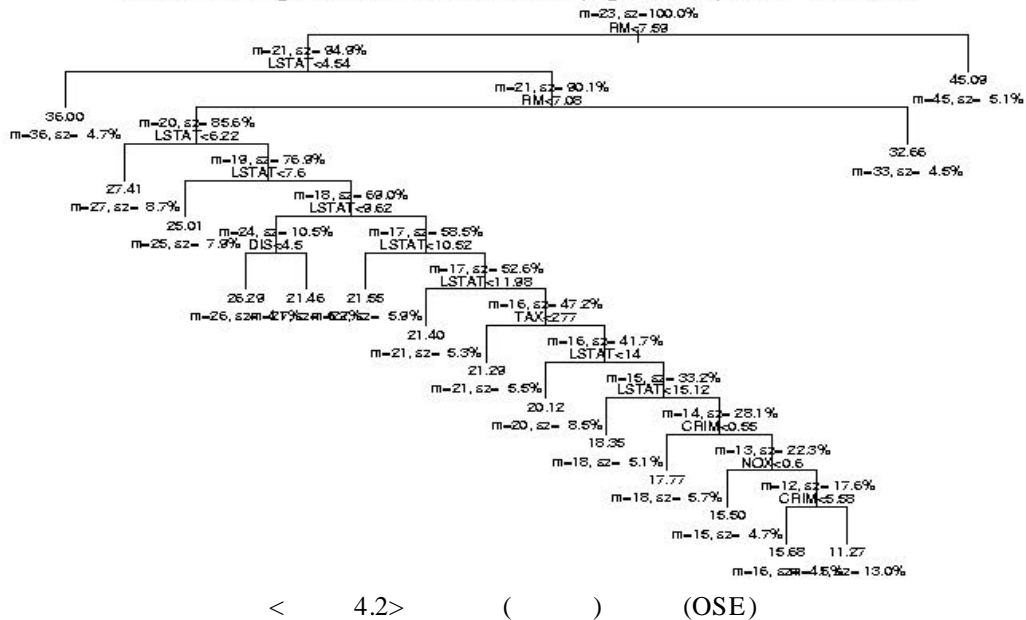
1. Breiman, L. (1996), Bagging predictors, *Machine Learning*, 24, 123- 140.
2. Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984), *Classification and Regression Trees*, Pacific Grove, CA: Wadsworth.
3. Freund, Y. and Schapire, R. E. (1996), A Decision-Theoretic Generalization of On-Line Learning and Application to Boosting, *Journal of Computer and System Science*, 55, 119- 139.
4. Hartigan, J. A. (1975), *Clustering Algorithms*, New York, NY: John Wiley & Sons, Inc..
5. Lee, Y-S. (2001), New Splitting Criteria for Classification Trees, *The Korean Communications in Statistics*, 8, 885-894.
6. Lee, Y-S. (2002), Interesting Node Finding Criteria for Regression Trees, *The Korean Journal of Applied Statistics* (in process).
7. Morgan, J. N., and Sonquist, J. A. (1963), Problems in the Analysis of Survey Data, and a Proposal, *Journal of the American Statistical Association*, 58, 415-434.
8. Quinlan, J. R. (1993), *C4.5: Programs for Machine Learning*, San Mateo, CA: Morgan Kaufmann.
9. Tibshirani, R. and Knight, K. (1996), Model Search and Inference via

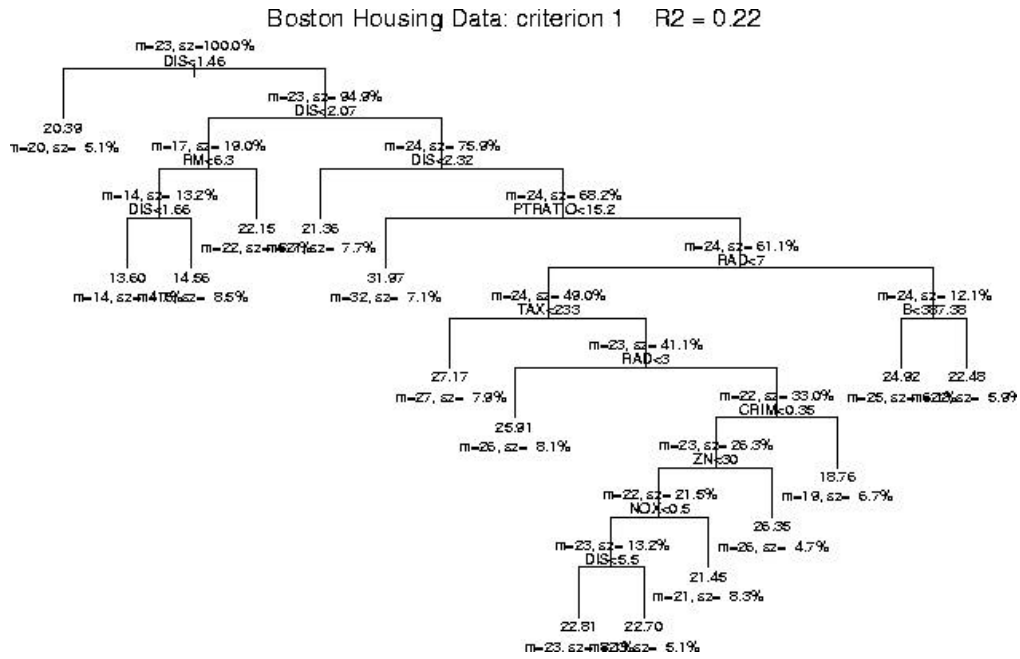
Bootstrap Bumping, Technical Report, Dept. of Statistics, U. of Toronto.

Boston Housing Data: one-sides purity,OSP R2 = 0.75

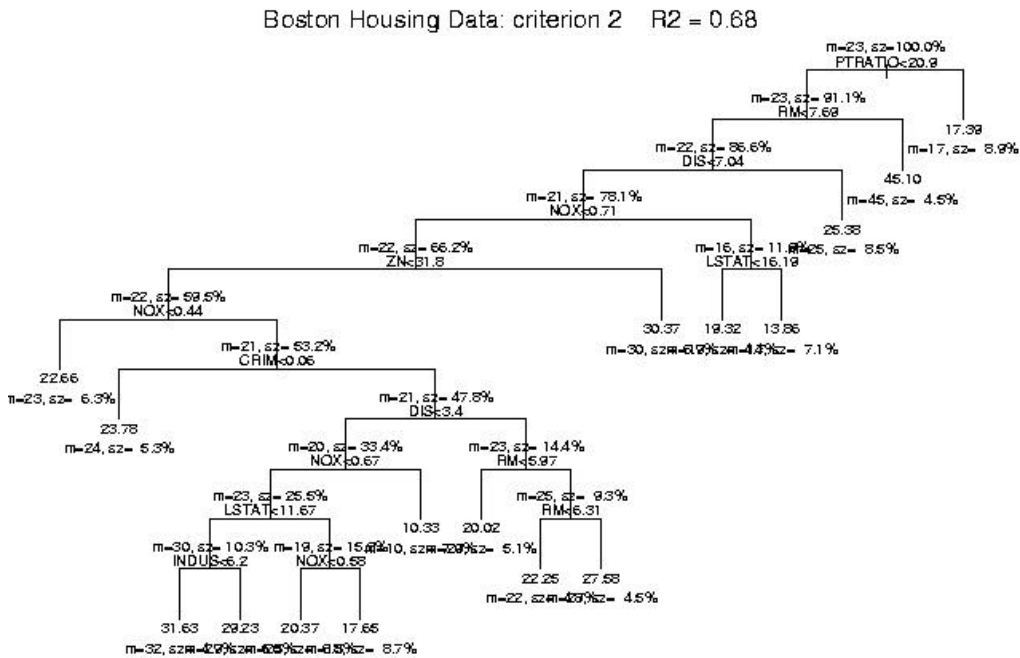


Boston Housing Data: one-sides extremes(high means),OSE R2 = 0.78





< 4.3> Criterion 1



< 4.4> Criterion 2