

\*                      \*\*  
,

# Missing Values in Classification Trees

Young-Hoon Kwak\*, Ji-Hyun Kim\*\*

## Abstract

There are many statistical methods for classification or discriminant analysis. The classification tree is one of them. There are several algorithms for classification tree. In this paper a comparison study is done of the performance of those algorithms for classification tree when we have missing values in predictors at hand. With the misclassification error rate as a performance measure, different algorithms are compared for different patterns of missing values.

---

: 2001 12

\* ( )

\*\*

•  
 (Classification Tree)  
 , 가 . 가  
 , 가  
 . (Missing Completely  
 At Random, MCAR) (Missing At Random, MAR) ,  
 (Non-ignorable) 가 .  
 가  
 가 S-PLUS, QUEST, SAS, SPSS  
 (simulation) (algorithm) ,  
 •  
 가 (training  
 stage)  
 가 ,  
 Breiman et al. (1984) CART (Classification And Regression Tree)  
 가 (surrogate splits) .  
**S-PLUS**  
 S-PLUS 'na.tree.replace()' 가 . 가  
 가 .  
 가 a, b, c , 가  
 a, b, c, NA .  
 가 1, 1, 2, 3, 3, 4, 5  
 가 , 1, 2, 3, 4, 5  
 NA , 1 2  
 1 가 ,  
 . 가 S-PLUS

**QUEST**

QUEST

QUEST

가

Loh et al. (1997)

(node mean)

(node mode)

$a, b, c$

가

가

가

**SAS**

SAS

(surrogate rule)

가

CART

**SPSS**

SPSS

al.(1997)

CART

C5.0

C5.0

Michael et

CART

(binary split)

C5.0

가 가

C5.0

가

CART

가

가

가

가

C5.0

가

가

가

S-PLUS, QUEST, SAS

SPSS

3가

(waveform)

(2000)

(Breiman et al. 1984)  $h_1(t), h_2(t), h_3(t)$

$h_1(t), h_2(t), h_3(t)$

21

1

$u$  , 21  $\epsilon_1, \dots, \epsilon_{21}$  .  
 $x$  .  
 $x_m = uh_1(m) + (1-u)h_2(m) + \epsilon_m$   
 2 .  
 $x_m = uh_1(m) + (1-u)h_3(m) + \epsilon_m$   
 3 .  
 $x_m = uh_2(m) + (1-u)h_3(m) + \epsilon_m$   
 200 600 . 10  
 , 600 3000 .  
 3000 .  
 10 .

Random, MCAR) 가 가 (Missing Completely At  
 (Non-ignorable) 가 (Missing At Random, MAR) 가 가

1. MCAR :  $x$   $h_1(t), h_2(t), h_3(t)$  가 가  $x_7, x_{11}, x_{15}$   
 (surrogate)  
 10, 20, 30, 40, 50%

2. MAR :  $h(t)$  가  
 (0,1) 가  
 $x_{15}$  , 1 ,  $h_1$   $x_7$  2  
 $x_7$  , 2  $h_3$   $x_{11}$  2  
 $x_{15}$  2 ,  $x_{11}$  .  
 20%, 40%, 60%, 80%, 100%

3. (Non-Ignorable case) : 1  $x_{15}$  가 2  $x_{15}$   
 , 2  $x_7$  2  $x_7$   
 , 3  $x_{11}$  2  $x_{11}$  .  
 20%, 40%, 60%, 80%, 100%

1. S-PLUS

< 1>      < 3>      10  
 . MCAR      가  
 가      가      가      가      .      2  
 0 50%      가      가      가      . MAR      가  
 S-PLUS      가      가      가      .

< 1> S-PLUS (MCAR )

		10%	20%	30%	40%	50%
	0.2887	0.4465	0.3696	0.3410	0.3489	0.3534
	0.0247	0.0610	0.0281	0.0293	0.0255	0.0256

< 2> S-PLUS (MAR )

		20%	40%	60%	80%	100%
	0.2887	0.4783	0.4990	0.5124	0.5758	0.6232
	0.0247	0.0186	0.0392	0.0355	0.0559	0.0219

< 3> S-PLUS (NI )

		20%	40%	60%	80%	100%
	0.2887	0.4619	0.4933	0.5467	0.5866	0.6045
	0.0247	0.0583	0.0322	0.0523	0.0400	0.0472

2. QUEST

QUEST      MCAR      가      가 S-PLUS      가  
 . MAR      가      ,      가  
 가      가      .  
 가      .

< 4> QUEST (MCAR )

		10%	20%	30%	40%	50%
	0.2717	0.2781	0.2852	0.2960	0.2984	0.3063
	0.0174	0.0217	0.0223	0.0261	0.0156	0.0256

< 5> QUEST (MAR )

		20%	40%	60%	80%	100%
	0.2717	0.2724	0.2895	0.3080	0.3237	0.4427
	0.0174	0.0138	0.0151	0.0227	0.0241	0.0232

< 6> QUEST (NI )

		20%	40%	60%	80%	100%
	0.2717	0.28242	0.2921	0.3018	0.3112	0.4912
	0.0174	0.0124	0.0197	0.0147	0.0118	0.1018

3. SAS

SAS QUEST , MCAR 가  
 가 S-PLUS QUEST , MAR  
 , 가 가

< 7> SAS (MCAR )

		10%	20%	30%	40%	50%
	0.2784	0.2802	0.2823	0.2935	0.2922	0.3059
	0.0162	0.0210	0.0195	0.0227	0.0176	0.0272

< 8> SAS (MAR )

		20%	40%	60%	80%	100%
	0.2784	0.2673	0.2731	0.3095	0.3097	0.3486
	0.0162	0.0231	0.0109	0.0245	0.0215	0.0762

< 9> SAS (NI )

		20%	40%	60%	80%	100%
	0.2784	0.2821	0.2855	0.3122	0.3357	0.4199
	0.0162	0.0168	0.0164	0.0260	0.0201	0.0557



- [1] Breiman, L., Friedman, J., Olshen, R., and Stone, C. (1984). *Classification and Regression Trees*. Chapman and Hall, New York, NY.
- [2] Loh, W.-Y., and Shih, Y.-S. (2000). A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms, *Machine Learning Journal*.
- [3] Loh, W.-Y. and Shih, Y.-S. (1997). Split selection methods for classification trees. *Statistica Sinica*.
- [4] Venables, W. N. and Ripley, B. D. (1997). *Modern Applied Statistics with S-Plus*. Springer, New York, NY, 2nd edition.
- [5] Chambers, J. M. and Hastie, T. J. (1990). *Statistical Models in S*. Wadsworth.
- [6] Berry, M. J. A, Linoff, Gordon (1997). *Data Mining Techniques For Marketing, Sales, and Customer Support*, John Wiley & Sons, Inc.
- [7] (2000). CART