

Paper 351-2012

Kass Adjustments in Decision Trees on Binary/Interval Target

Manoj Immadi and Dr. Goutam Chakraborty, Oklahoma State University, Stillwater, OK, USA

ABSTRACT

Kass adjustment, a commonly used option, maximizes independence between the two branches of a decision tree after any split, but how well these adjustments work on interval and binary target variables is under-researched. This paper describes split search algorithm in decision trees for selecting useful inputs, followed by comparing decision tree results with and without Kass adjustments for a binary target variable and an interval target variable. In addition, the paper also provides insights to Kass adjustments and Bonferroni correction in decision trees.

INTRODUCTION

Decision trees are one of the most important classification algorithms in data mining and machine learning techniques. Decision trees are simple rules which are easy to explain and interpret. Thus, decision trees are one of the most commonly used predictive modeling algorithms.

Before attempting to delineate the differences between the Kass adjustment or Bonferroni correction in decision trees, it is advisable to understand how decision trees work, particularly with respect to how split search algorithm works. When predicting new cases, the decision tree algorithm scores them using simple prediction rules. Decision tree uses Split search algorithm to select the useful inputs. Split search algorithm is explained in detail later in this paper. Typically, after selection of appropriate input variables, a decision tree builds the model on the selected variables and prunes the tree to optimize the complexity.

A simple decision tree is shown below:

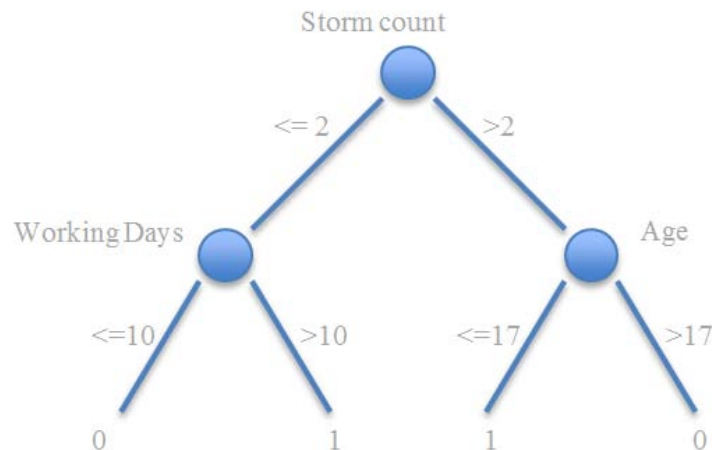


Figure 1. Example of a Decision Tree

The above decision tree can be interpreted by means of very simple rules. Target 0 indicates no admit and 1 indicates admit in the sample dataset. Then storm count, age and working days are the input variables to predict the binary admit (how many people are admitted in a hospital) target variable.

If the storm count is greater than 2 and the age of a person is less than or equal to 17, the decision tree model predicts that the person is likely to be admitted. Similarly, a storm count less than or equal to 2 and working days greater than 10 also indicate that the person is likely to be admitted. By looking at the decision tree diagram as shown above and/or English rules (as enumerated in the write-up), we can interpret the decision tree very easily.

KASS ADJUSTMENTS

Before getting into the details of KASS adjustments, it is important to understand how split search algorithm works in decision trees. To select useful inputs, decision tree use split - search algorithm. Split search starts by selecting an input for creating splitting rule in the training data. If the selected input is an interval variable, each and every value in the input variable serve as a potential split point for predicting target variable or if the input variable is categorical input variable, the average value of the target is taken for each value in categorical input variable and the average value will be served as a potential split point.

After selecting the split point, two groups of data (based on the two branches of the tree) will be generated. All the records with smaller than the value of the split point is typically in the left branch of the tree and the records with greater than the value of the split point are on the right branch of the tree. These two groups forms 2x2 contingency table assuming that the target is a binary variable. Pearson chi-squared statistic is used to quantify the independence, large value of chi squared statistic indicates that the proportion of zeros and ones in two branches is very different. A large difference also indicates it is a good split since the results in the p-values being very small. The log worth is calculated based on the p-value and used as the statistic for selecting split point.

$$\text{logworth} = -\log(\text{chi} - \text{squared } p - \text{value})$$

Logworth must exceed threshold value in order for a split to occur. By default, the threshold value of log-worth is 0.7 in SAS® Enterprise miner.

The algorithm is much more complex than explained above. There are several other factors makes algorithm more complicated including:

1. Tree settings, such as the minimum number of observations for a split and the minimum number of observations in leaf and so on.
2. If an input variable is an interval variable, number of possible split point will increase since the number of levels increase, the likelihood of obtaining significant value also increases. In this way, an input with a multitude of unique input values has a higher chance of having larger logworth than an input with less number of unique input values.

Kass adjustments and Bonferroni corrections are based on Boole's inequality, also known as the union bound. Boole's inequality says that in any countable set of events, probability that at least one of those events happens is no greater than the sum of the probabilities of individual events.

For n countable set of events $x_1, x_2, x_3...x_n$; we have:

$$P(U_n x_n) \leq \sum_n P(x_n)$$

Bonferroni correction is used to counteract the problem of multiple comparisons which is based on Boole's inequality explained above. Statisticians solved the problem of combining the results from multiple statistical tests with the aid of Bonferroni correction.

Since each split point corresponds to a statistical test, Bonferroni corrections are automatically applied to the logworth calculations. These corrections are called Kass adjustments, which penalize the inputs with many split points by reducing the logworth. It reduces the logworth to an amount equal to log of the number of distinct input values. These adjustments allow for a fair comparison of variables with high multitude and low multitude of unique values in input variables. There are special ways of handling the data with missing values. Two sets of Kass adjusted logworth values are typically calculated. These details are out of the scope of this paper.

After the first split, the significance of the secondary and the following splits depends on the significance of the earlier splits. The split search algorithm again goes through a multiple comparison problem. To compensate, this algorithm increases the threshold by an amount related to the number of splits.

DATA

The dataset is related to weather and the health care usage facilities. The dataset is created from 5 different datasets made available to student teams participating in the A2011 analytics shootout challenge. The final dataset contains the information about weather, hospital admits, storm and population in a given particular area code and for each age group. The target variable 'admits' contains the number of admits in a given particular area and in a particular age group and for a particular DRG (type of disease). This paper discusses about the Kass adjustments hence it will not go through the details of data preparation. A few of the important things are mentioned and relevant diagrams are shown below.

'Admits' is an interval target variable, from 'Admits' a new binary target variable '*isAdmit*' is created, which serves the purpose of comparing the results of interval and binary target variables with and without Kass adjustments. Many other new variables are created from the existing variables.

'AdmitsProp' a new interval target variable is created from admits (existing target variable) and population (total population in that particular area) is created. 'AdmitsProp' contain the proportion of admits to the population in the given area. This new variable has been created by looking at the normalized distribution charts of Admits in the data. The new 'AdmitsProp' target variable was used for Kass adjustments comparison purposes.

In the 'Admit' target variable 73.59% of observations do not have any admits in the dataset. 24.52% of observations have the number of admits as one and the rest 1.89 % of observations contains admits ranging from 2 to 14. The distribution of admits is highly right skewed, high number of admits are very rare and mostly occur in densely populated areas. In the Figure 2 we can have look at the number of admits and number of observations pertaining to Admit.

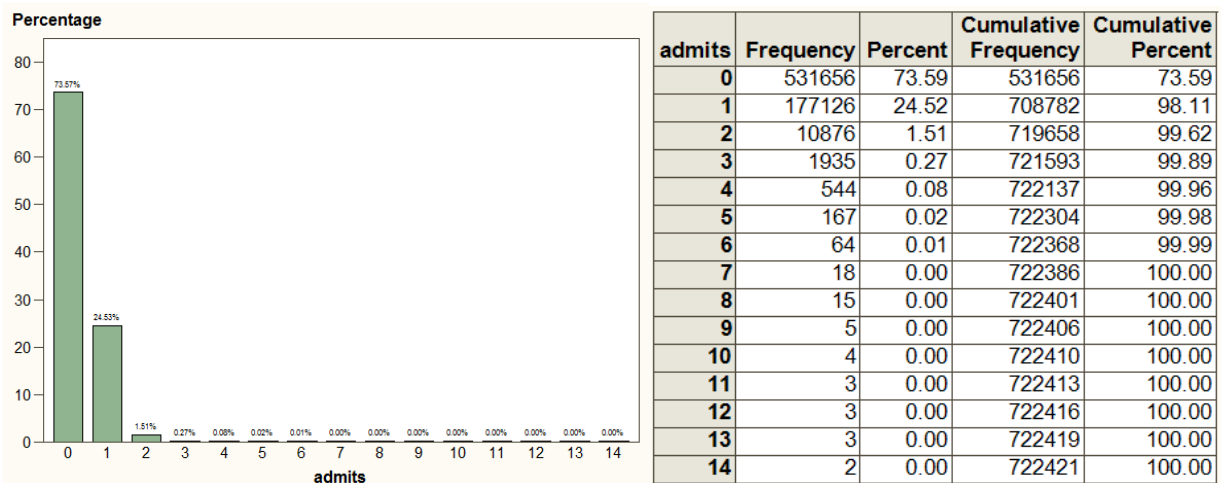


Figure 2. Admit (target variable) frequency chart

Data is normalized by calculating new target variable 'AdmitsProp'. The new variable is the proportion of admits to the population per 10,000 (ten thousand). Even though the distribution of admits proportion is also highly skewed with 99.49% of admits proportion concentrated to 0- 60 with a few exceptional cases like 2500, 5000, 10000 etc. The reason behind this is the high number of admits in low populated areas.

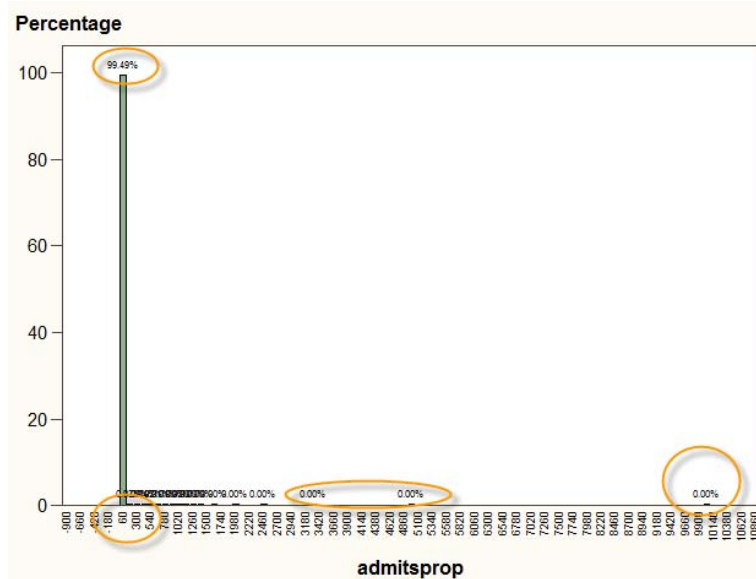


Figure 3. AdmitsProp variable distribution

For better results and for better comparison purposes the outliers were removed before proceeding further. The new interval target variable 'AdmitsProp' contains the values in the range of zero to two hundred only.

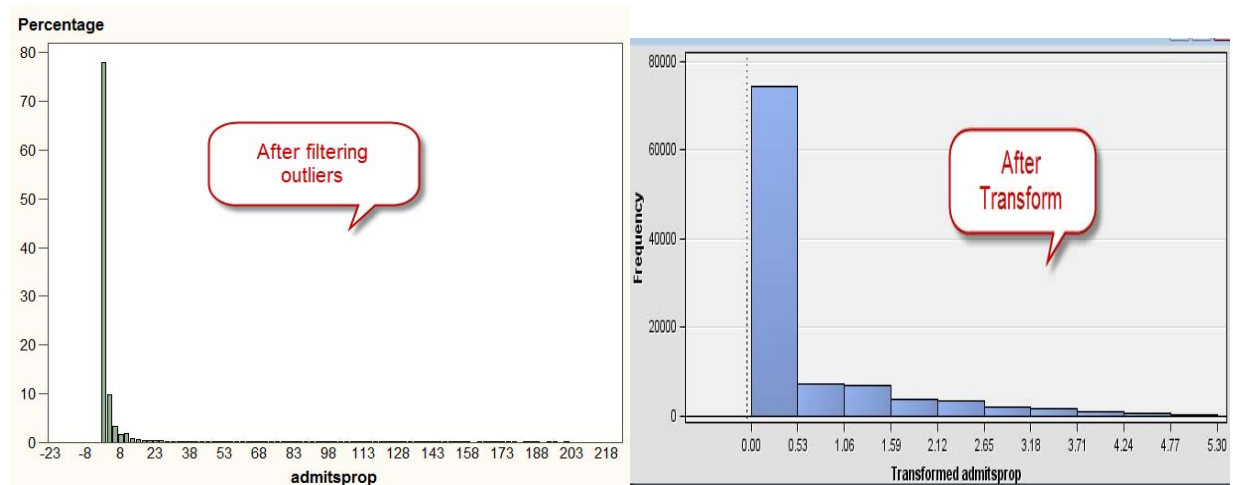


Figure 4. admitsprop after distribution after filtering and transforming

After the preparation of final dataset our first task in this dataset is to look for any inconsistencies, errors or extreme values in the data. Frequency distribution, descriptive statistics, cross tab and multiplot nodes were used in order to analyze and various techniques used for fixing the issues. A few of the important things are noted down below:

- Missing values in storm are replaced with '0' indicates 'no storms'
- Missing values in weather are replaced by the average value of temperatures in that particular area code and in that particular month.
- Minimum incubation is deleted since every value is 1.
- We can see a high number of admits in months 2,3,11 and 12 and we can observe less admits in months 6 and 7. Admits has shown seasonality in the sample data.

- We have less percentage of missing values in sn0-sn5; ah0-ah5; al0-al4; atd0-atd4; dl0-dl4; htd0-htd4; mh0-mh4; ml0-ml4; prcp0-prcp4 and are replaced by the mean.
- 'AdmitsProp' only has 0.224% of values which are greater than 200, they were removed, thus improving skewness.
- Different transformations were applied to different variables in order to improve the model performance.

The following is the screen shot of the data preparation:

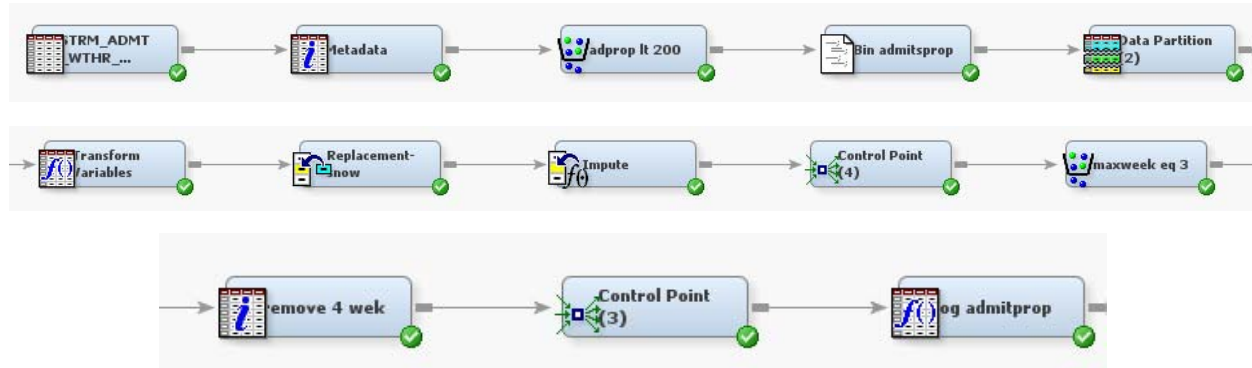


Figure 5. Data Preparation

Data understanding is performed by running the following nodes in the enterprise miner:

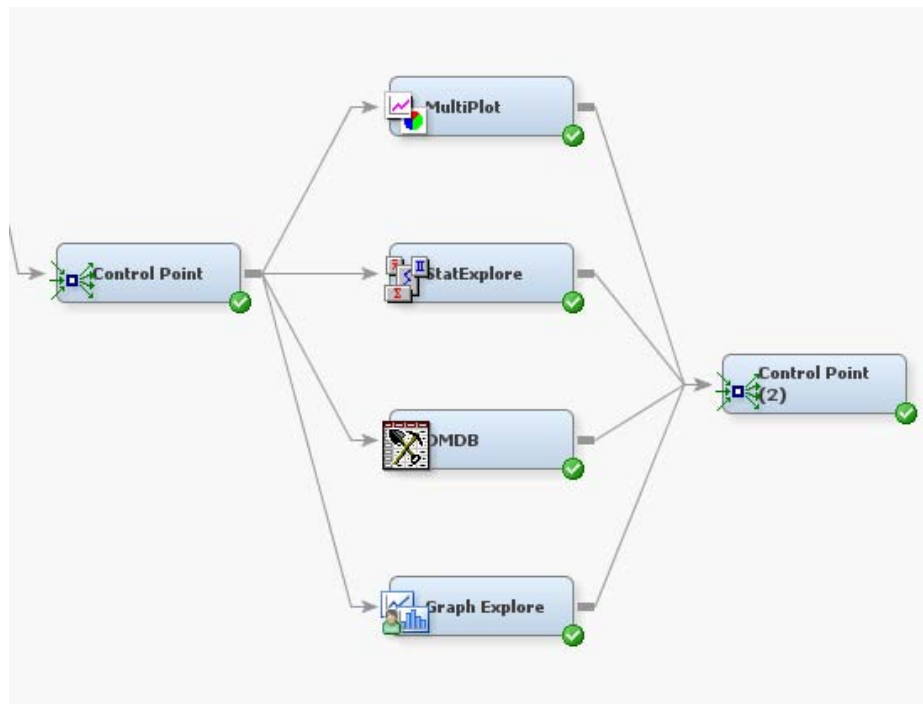


Figure 6. Data Understanding

METHODOLOGY

The comparison Kass adjustments/Bonferroni correction for binary/interval target variable is done based on the following methodology. Decision tree in enterprise comes with various options. Apart from changing the number of branches we have various options. For categorical input SAS® offers three different split worth criteria, they are ProbChiSq, Gini and Entropy. It is well-known that they give similar results if the number of levels in each categorical

input is similar. We also have option to turn on/off the Bonferroni Adjustment. By using this option we can compare how these adjustments work on binary/interval target variable.

My goal is to compare the results by turning on and off the Bonferroni Adjustment for binary and interval target variable, thus comparing the results.

The first task is to compare the results of binary target variable with and without Bonferroni Adjustment. The first three nodes Figure 7 are Entropy, Gini and ProbChiSq models with Bonferroni Adjustment and the other three nodes are without Bonferroni Adjustment.

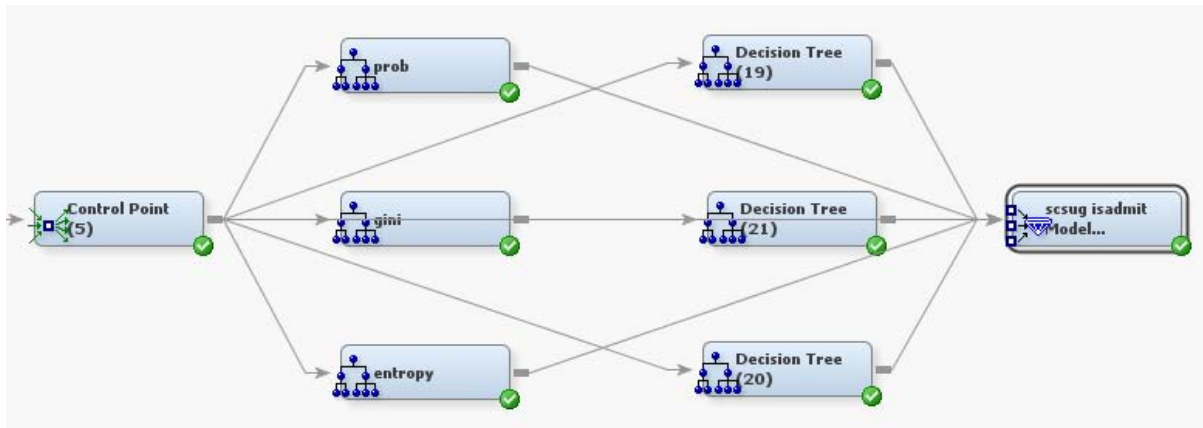


Figure 7. Binary target with/without Bonferroni Adjustment

The ProbChiSq with Bonferroni Adjustment is selected as the best model from the model comparison node. The properties of the selected model are:

Subtree	
Method	Assessment
Number of Leaves	1
Assessment Measure	Average Square Error
Assessment Fraction	0.25
Cross Validation	
Perform Cross Validation	No
Number of Subsets	10
Number of Repeats	1
Seed	12345
Observation Based Imp	
Observation Based Imp	No
Number Single Var Imp	5
P-Value Adjustment	
Bonferroni Adjustment	Yes
Time of Kass Adjustment	Before
Inputs	No
Number of Inputs	1
Split Adjustment	Yes
Output Variables	

Figure 8. Properties of ProbChiSq model with Bonferroni Adjustment

Bonferroni Adjustment in the properties panel is changed to no for the other three decision trees.

The next task is to compare the results interval target variable with and without Bonferroni Adjustment. The first three nodes in Figure 9 are Entropy, Gini and ProbChiSq models with Bonferroni Adjustment and the other three nodes are without Bonferroni Adjustment.

Please refer results section for the comparison how Bonferroni Adjustment/Kass adjustments improved the results of binary target variable and interval target variable.

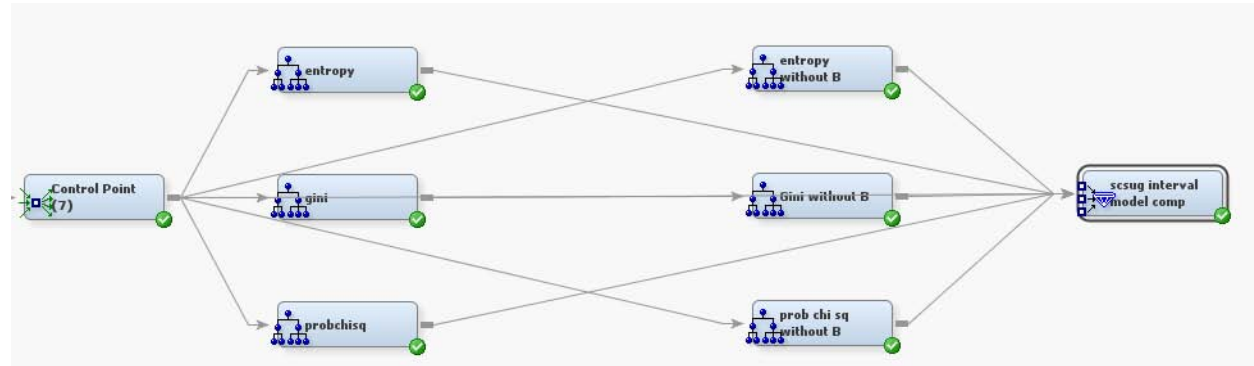


Figure 9. Interval target with/without Bonferroni Adjustment

Properties of ProbChiSq model without Bonferroni Adjustment are shown below:

Subtree	
Method	Assessment
Number of Leaves	1
Assessment Measure	Average Square Error
Assessment Fraction	0.25
Cross Validation	
Perform Cross Validation	No
Number of Subsets	10
Number of Repeats	1
Seed	12345
Observation Based Imp	
Observation Based Imp	No
Number Single Var Imp	5
P-Value Adjustment	
Bonferroni Adjustment	No
Time of Kass Adjustment	Before
Inputs	No
Number of Inputs	1
Split Adjustment	Yes
Output Variables	

Figure 10. ProbChiSq properties without Bonferroni Adjustment

RESULTS

For simplification purposes this paper compares the results selected model with/without Bonferroni Adjustment for each binary and interval target variable. Kass adjustments improved the results in all other models too. The models built for comparison purposes are very simple with properties, maximum branch of 2 and maximum depth of 6.

Results of binary target variable are shown below:

Selected Model	Predecessor Node	Model Node	Model Description	Target Variable	Valid: Average Squared Error	Train: Sum of Frequencies	Train: Sum of Case Weights Times Freq	Train: Misclassification Rate	Train: Maximum Absolute Error	Train: Sum of Squared Errors	Train: Average Squared Error
Y	Tree	Tree	prob	isAdmit	0.184616	282713	565426	0.274006	0.82774	108297.9	0.191533
	Tree2	Tree2	entropy	isAdmit	0.184638	282713	565426	0.273999	0.955224	108306	0.191548
	Tree23	Tree23	Entropy wit...	isAdmit	0.184638	282713	565426	0.273999	0.955224	108306	0.191548
	Tree22	Tree22	Probchisq ...	isAdmit	0.184678	282713	565426	0.27379	0.82774	108316.9	0.191567
	Tree24	Tree24	Gini Withou...	isAdmit	0.184678	282713	565426	0.27379	0.82774	108316.9	0.191567
	Tree3	Tree3	gini	isAdmit	0.184678	282713	565426	0.27379	0.82774	108316.9	0.191567

Figure 11. Results of binary target variable

Valid: Average Squared Error is the selection criteria for selecting the model, model comparison selected ProbChiSq with Bonferroni Adjustment as the best model. Valid average squared error of the selected model is 0.184616, but when you look at the train: average squared error used for the model building the average squared error of the selected model with and without Bonferroni are 0.191533 and 0.191567 respectively. Train average squared error has shown significant improvement by including Bonferroni Adjustment. Train average squared error improved 0.000037.

Results of interval target variable are shown below:

Selected Model	Predecessor Node	Model Node	Model Description	TARGET	Valid: Average Squared Error	Train: Sum of Frequencies	Train: Sum of Case Weights Times Freq	Train: Maximum Absolute Error	Train: Sum of Squared Errors	Train: Average Squared Error	Train: Root Average Squared Error
Y	Tree7	Tree7	probchisq	LOG_admitsprop	0.810779	282713	282713	4.94225	257210.6	0.909794	0.953831
	Tree8	Tree8	entropy	LOG_admitsprop	0.810779	282713	282713	4.94225	257210.6	0.909794	0.953831
	Tree9	Tree9	gini	LOG_admitsprop	0.810779	282713	282713	4.94225	257210.6	0.909794	0.953831
	Tree16	Tree16	prob chi sq without B	LOG_admitsprop	0.840064	282713	282713	4.94225	257218.5	0.909822	0.953846
	Tree17	Tree17	entropy without B	LOG_admitsprop	0.840064	282713	282713	4.94225	257218.5	0.909822	0.953846
	Tree18	Tree18	Gini without B	LOG_admitsprop	0.840064	282713	282713	4.94225	257218.5	0.909822	0.953846

Figure 12. Results of interval target variable

Valid: Average Squared Error is the selection criteria for selecting the model, model comparison selected ProbChiSq with Bonferroni Adjustment as the best model. Valid average squared error of the selected model is 0.810779, but when you look at the train: average squared error used for the model building the average squared error of the selected model with and without Bonferroni are 0.909794 and 0.909822 respectively. Train average squared error has shown significant improvement by including Bonferroni Adjustment. Train average square error is improved 0.000028.

CONCLUSION

Kass adjustments showed some significant improvement in both binary and interval target variables. But it has shown better improvement in the binary target variable than the interval target variable.

REFERENCES

- [1] "Using Decision Trees to Identify Medicare Part B Providers for Audit", poster SESUG -2002, by Noel McKetty, First Coast Service Options, Jacksonville, FL; Donna Mohr, University of North Florida, Jacksonville, FL; "<http://analytics.ncsu.edu/sesug/2002/PS08.pdf>"
- [2] "Decision Trees", by Andrew W. Moore Professor School of Computer Science Carnegie Mellon University "<http://www.autonlab.org/tutorials/dtree18.pdf>"
- [3] "Decision trees - what are they?" SAS® support community, "<http://support.sas.com/publishing/pubcat/chaps/57587.pdf>"

[4] "K Nearest neighbors Classifier & Decision Trees", "<http://www.ibms.sinica.edu.tw/~pan/classification/documents/KNN&DECISION%20TREE.ppt>"

[5] "Bonferroni Inequalities", "<http://mathworld.wolfram.com/BonferroniInequalities.html>"

[6] "Proof of Booles inequality", "<http://planetmath.org/encyclopedia/ProofOfBooleInequality.html>"

ACKNOWLEDGEMENTS

Thanks for IICH and Analytics2011 team for allowing me to work on this dataset.

CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the author at:

Manoj Immadi, Oklahoma State University, Stillwater, OK, Email: manoj.immadi@okstate.edu

Manoj Immadi is a master's student in Management Information Systems at Oklahoma State University. He has two years of professional experience as programmer analyst. He is SAS® Certified Advanced Programmer for SAS® 9 and Certified Predictive Modeler Using SAS® Enterprise Miner 6.1™

Dr. Goutam Chakraborty, Oklahoma State University, Stillwater OK, Email: goutam.chakraborty@okstate.edu

Goutam Chakraborty is a professor of marketing and founder of SAS® and OSU data mining certificate program at Oklahoma State University. He has published in many journals such as *Journal of Interactive Marketing*, *Journal of Advertising Research*, *Journal of Advertising*, *Journal of Business Research*, etc. He chaired the national conference for direct marketing educators in 2004 and 2005 and co-chaired the M2007 data mining conference. He is also a Business Knowledge Series instructor for SAS®.

TRADEMARKS

SAS® and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration. Other brand and product names are trademarks of their respective companies.