

Classification and Regression Trees as Alternatives to Regression

Mandy C. Phelps & Edgar C. Merkle

Department of Psychology, Fairmount College of Liberal Arts and Sciences

Abstract. Traditional statistical analyses, such as ANOVA and Regression, require that many assumptions be fulfilled in order to obtain accurate results. For example, there are assumptions of normally-distributed data, linear relationships between the dependent variable(s) and independent variables, and homogeneity of variance. In this presentation, we will describe the use of Classification and Regression Trees (CART) to sidestep the assumptions required by traditional analyses. CART has the added benefit of not requiring large sample sizes in order to obtain accurate results, although larger sample sizes are preferred. There is also a difference in the goal of CART compared to traditional analyses. CART is geared toward prediction, whereas traditional analyses are geared toward developing a specific model for your data. The poster will contain specific information about the procedures underlying CART, as well as an example involving data from a legibility of fonts study.

1. Introduction

Classification and Regression Trees are nonparametric therefore, they “can handle numerical data that are highly skewed or multi-modal, as well as categorical predictors with either ordinal or non-ordinal structure” [1]. This will save researchers time because they will not have to check the normality of their distribution, as well as other assumptions, which will allow them to focus more on the interpretation of their results. Researchers also have more leniencies in choosing which form of independent variables to use. Classification trees are capable of handling independent variables that are continuous, categorical, or a mixture of both, and unlike traditional statistics, actually perform better when there are numerous independent variables. The final tree will only include the independent variables that were predictive of the dependent variable; the other non-predictive independent variables will have no affect on the final results, as they do with more traditional statistics, such as regression.

Missing data also poses less of a problem with classification trees because classification trees have a built-in algorithm that readily handles missing data by replacing it with a surrogate variable that is most similar to the primary splitter [1]. The graphical presentation of the results also allows the researcher to see the complex interactions among the independent variables which is not as easily accomplished with traditional statistics. Lastly, the results of classification trees generalize well to new data when the necessary steps are taken in the growing and pruning process, which will be explained shortly. Now that some of the background has been covered it is time to move on to how classification trees actually work.

The structure of a classification tree includes parent nodes, child nodes, and terminal nodes. The parent node contains all observations that are to be divided into two groups. Child nodes are the nodes resulting from the splitting of the parent node. Each child node becomes a parent node itself and then splits into two additional child nodes. Nodes that satisfy the split criterion are labeled with a one and those that do not are labeled with a zero. This process continues until the tree has finished growing. Terminal nodes are the final nodes of the tree which contain the most heterogeneous subsets of the data and give the predictions for the observations that were ran down the tree.

Some important steps to follow when growing a tree according to Breiman [2] include a set of binary questions, a goodness of split criterion, a stopping rule, and determining how to assign every terminal node to a class. The binary questions determine if the observations goes to the right child node or left child node. They are set up as “yes” or “no” responses that if the observation meets the criterion then the observation goes to the left and if the observation does not meet the criterion it goes to the right. This process continues until all the observations have been have been divided into child nodes that are as different as possible. There is no set number of observations that are required to be in each node so they do not have to be divided equally, it is based on the properties of the observations and any criteria set by the researcher, such as a splitting criterion.

The next step is determining the goodness of split criterion. The objective here is to determine the split that results in the minimum amount of node impurity, which is the amount of difference between the observations

