Chapter 1 Tree-Based Modeling Techniques

Dileep Kumar G. Adama Science and Technology University, Ethiopia

ABSTRACT

Tree-based learning techniques are considered to be one of the best and most used supervised learning methods. Tree-based methods empower predictive models with high accuracy, stability, and ease of interpretation. Unlike linear models, they map non-linear relationships pretty well. These methods are adaptable at solving any kind of problem at hand (classification or regression). Methods like decision trees, random forest, gradient boosting are being widely used in all kinds of machine learning and data science problems. Hence, for every data analyst, it is important to learn these algorithms and use them for modeling. This chapter guide the learner to learn tree-based modeling techniques from scratch.

INTRODUCTION

Tree based learning techniques are considered to be one of the best and mostly used supervised learning methods. Tree based methods empower predictive models with high accuracy, stability and ease of interpretation. Unlike linear models, they map non-linear relationships pretty well. These methods are adaptable at solving any kind of problem at hand (classification or regression). Methods like decision trees, random forest, gradient boosting are being widely used in all kinds of machine learning and data science problems. Hence, for every data analyst, it is important to learn these algorithms and use them for modeling. This chapter guide the learner to learn tree-based modeling techniques from scratch.

DOI: 10.4018/978-1-5225-3534-8.ch001

DECISION TREE

Decision tree is a type of supervised learning algorithm (having a predefined target variable) that is mostly used in classification problems (James, Witten, Hastie, & Tibshirani, 2013). It works for both categorical and continuous input and output variables. In this technique, we split the population or sample into two or more homogeneous sets (or sub-populations) based on most significant splitter / differentiator in input variables.

Example 1: We have a sample of 30 students in a class with three variables Gender (Boy/Girl), Class (IX/X) and Height (5 to 6 ft). 15 out of these 30 play football in leisure time. Now, create create a model to predict who will play cricket during leisure period? In this problem, we need to segregate students who play football in their leisure time based on highly significant input variable among all three.

This is where decision tree helps, it will segregate the students based on all values of three variable and identify the variable, which creates the best homogeneous sets of students (which are heterogeneous to each other). In Figure 2, 3, and 4, you can see that variable Gender is able to identify best homogeneous sets compared to the other two variables.

As mentioned above, decision tree identifies the most significant variable and its value that gives best homogeneous sets of population.



Figure 1. Decision tree

16 more pages are available in the full version of this document, which may be purchased using the "Add to Cart"

button on the publisher's webpage: www.igi-

global.com/chapter/tree-based-modeling-techniques/207377

Related Content

Evolution of Subway Network Systems, Subway Accessibility, and Change of Urban Landscape: A Longitudinal Approach to Seoul Metropolitan Area Yena Songand Hyun Kim (2018). Intelligent Transportation and Planning: Breakthroughs in Research and Practice (pp. 1087-1111). www.irma-international.org/chapter/evolution-of-subway-network-systems-subway-accessibilityand-change-of-urban-landscape/197177

Fitting a Three-Phase Discrete SIR Model to New Coronavirus Cases in New

York State Kris H. Green (2021). International Journal of Data Analytics (pp. 59-74). www.irma-international.org/article/fitting-a-three-phase-discrete-sir-model-to-new-coronavirus-

cases-in-new-york-state/285468

Data, Information, and Knowledge: Developing an Intangible Assets Strategy

G. Scott Ericksonand Helen N. Rothberg (2015). *Handbook of Research on Organizational Transformations through Big Data Analytics (pp. 85-96).* www.irma-international.org/chapter/data-information-and-knowledge/122750

A Markov-Chain-Based Model for Group Message Distribution in Connected Networks

Peter Bajorskiand Michael Kurdziel (2020). *International Journal of Data Analytics* (pp. 13-29).

www.irma-international.org/article/a-markov-chain-based-model-for-group-message-distributionin-connected-networks/258918

An Innovative Approach to Solve Healthcare Issues Using Big Data Image Analytics

Ramesh R., Udayakumar E., Srihari K.and Sunil Pathak P. (2021). International Journal of Big Data and Analytics in Healthcare (pp. 15-25).

www.irma-international.org/article/an-innovative-approach-to-solve-healthcare-issues-using-bigdata-image-analytics/268415