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

Immersion and Interaction via Avatars within Google Street View: Opening Possibilities beyond Traditional Cultural Learning

Ya-Chun Shihand Molly Leonard (2014). *Packaging Digital Information for Enhanced Learning and Analysis: Data Visualization, Spatialization, and Multidimensionality (pp. 266-281).*

www.irma-international.org/chapter/immersion-and-interaction-via-avatars-within-google-streetview/80222

SBASH Stack Based Allocation of Sheer Window Architecture for Real Time Stream Data Processing

Devesh Kumar Laland Ugrasen Suman (2020). International Journal of Data Analytics (pp. 1-21).

www.irma-international.org/article/sbash-stack-based-allocation-of-sheer-window-architecturefor-real-time-stream-data-processing/244166

Comprehensive Analysis of State-of-the-Art CAD Tools and Techniques for Chronic Kidney Disease (CKD)

Mynapati Lakshmi Prasudha, Rakesh Kasumollaand Deepak Sukheja (2021). International Journal of Big Data and Analytics in Healthcare (pp. 1-12). www.irma-international.org/article/comprehensive-analysis-of-state-of-the-art-cad-tools-and-techniques-for-chronic-kidney-disease-ckd/287605

Prediction Length of Stay with Neural Network Trained by Particle Swarm Optimization

Azadeh Oliyaeiand Zahra Aghababaee (2017). *International Journal of Big Data and Analytics in Healthcare (pp. 21-38).*

www.irma-international.org/article/prediction-length-of-stay-with-neural-network-trained-by-particle-swarm-optimization/204446

Effect of Crude Oil Price Variations on Stocks With Special Reference to NSE

Shimna Jayaraj, T. Shenbagavalli, V. Vipanchi, B.S. Sudha, Juhi Jainand Biswaranjan Senapati (2024). *Data-Driven Decision Making for Long-Term Business Success (pp. 248-263).*

www.irma-international.org/chapter/effect-of-crude-oil-price-variations-on-stocks-with-special-reference-to-nse/335576