

Learning Decision Trees with R: A Step-by-Step Guide

Authored by
Mohammed loot

October 28, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning Decision Trees with R: A Step-by-Step Guide*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=4581>

The Power and Interpretability of Decision Trees

In the vast landscape of statistical modeling and machine learning, the [decision tree](#) remains a supremely powerful and highly interpretable model. This methodology systematically partitions a dataset into increasingly homogeneous subsets based on the values of input features, culminating in a hierarchical, tree-like structure of sequential decisions. Structurally, every internal node represents a test or condition applied to a specific attribute; each branch illustrates the outcome of that test; and the final, conclusive leaf node (often referred to as a [terminal node](#)) provides either a class label for classification tasks or a numerical predicted value in the case of [regression tree](#) analysis.

The inherent versatility of decision trees allows them to seamlessly accommodate both categorical and numerical data types. They are particularly prized for their transparent nature, which gives analysts the ability to trace the exact logical path leading to any prediction. This exceptional interpretability provides a distinct advantage over more opaque "black-box" models, making decision trees the preferred choice for applications where understanding the underlying mechanism of the decision-making process is just as critical as achieving high predictive accuracy. Their graphical nature serves as an intuitive communication tool for complex relationships.

For data scientists and practitioners utilizing the [R programming environment](#), the process of constructing and visualizing these models is efficiently managed through dedicated packages. While the core functionality for building the model resides in the powerful [rpart](#) package, the [rpart.plot](#) package is essential for producing clear, aesthetically refined, and highly informative visualizations of the resulting tree structures. This powerful tandem empowers users to not only develop robust predictive models but also to communicate their analytical insights effectively and persuasively to diverse audiences.

Essential R Packages: `rpart` and `rpart.plot`

To master the implementation of decision trees within [R](#), a fundamental understanding of two key libraries is required: `rpart` and `rpart.plot`. The [rpart](#) (Recursive Partitioning and Regression Trees) package serves as the foundational engine for tree construction. It implements a highly regarded algorithm used globally for generating both classification and [regression trees](#), offering numerous control parameters that allow users to precisely fine-tune the complexity and structure of the model. Key functions within this package include the primary model-fitting function, `rpart()`, and the crucial simplification function, `prune()`.

The `rpart()` function is the core utility for initializing the tree construction process. It necessitates a standard formula structure that explicitly defines the [response variable](#) (the target to be predicted) and the set of [predictor variables](#) (the features used for splitting), along with the designated

dataset. A critical challenge in developing effective decision trees is managing model complexity to mitigate the risk of overfitting the training data. This control is primarily exerted through the `control` argument within `rpart()`, which enables specification of parameters such as the minimum count of observations required in a node or the pivotal [complexity parameter \(cp\)](#)--a factor central to the pruning methodology.

Once the decision model has been successfully built, the [rpart.plot](#) package becomes essential for high-quality visualization. Although base R does include basic plotting capabilities for `rpart` objects, the dedicated `prp()` function provided by [rpart.plot](#) offers dramatically superior features for generating diagrams that are both more readable and significantly more informative. These advancements include refined node labeling, improved visual aesthetics, and the capacity to display rich statistical details directly within the plot, making the analytical interpretation of the tree structure far more intuitive and efficient for the analyst.

To ensure a smooth execution of the subsequent example, please confirm that both **rpart** and **rpart.plot** packages are correctly installed on your system. If installation is necessary, use the standard R commands: `install.packages("rpart")` and `install.packages("rpart.plot")`. Following installation, it is mandatory to load these libraries into your active **R** session using the `library()` function, as demonstrated in the first lines of the code block below.

Practical Example: Building and Pruning a Regression Tree

To provide a comprehensive illustration of constructing and plotting a [regression tree](#) in R, we will leverage the popular **Hitters** dataset, which is conveniently packaged within the [ISLR](#) package. This dataset contains detailed statistical information for 263 professional baseball players, covering critical metrics such as their salaries, years of professional experience, and total home runs. Our primary goal is to develop a regression tree model that utilizes a player's experience level (`Years`) and accumulated home runs (`HmRun`) as [predictor variables](#) to accurately estimate their `Salary`.

The initial step involves building a maximally complex tree, which is often unpruned or only minimally restricted. While this initial structure captures all potential data splits, it is typically overly complex and highly susceptible to overfitting the nuances of the training data. To generate a model that generalizes well to new, unseen data, we must employ a systematic pruning strategy centered around the [complexity parameter \(cp\)](#). The `cp` value governs the inherent trade-off between the tree's final size and its overall predictive performance. A very small `cp` encourages a large, highly complex tree, whereas a larger `cp` results in a simpler, more generalized structure. Our objective is to identify the optimal `cp` value that yields the lowest cross-validated error.

The code block presented below outlines the essential steps: first, constructing the initial regression tree; second, determining the statistically optimal `cp` value by analyzing the results of

cross-validation performed by the `rpart` function; and third, using this optimal value to prune the tree into a more parsimonious and generalized model. Finally, the default settings of the `prp()` function are used to generate the initial visualization of this refined tree structure.

```
library(ISLR)
```

```
library(rpart)
```

```
library(rpart.plot)
```

```
# Build the initial decision tree with a very small cp to allow for a complex tree structure
```

```
tree <- rpart(Salary ~ Years + HmRun, data=Hitters, control=rpart.control(cp=.0001))
```

```
# Identify the best cp value to use based on minimizing the cross-validation error (xerror)
```

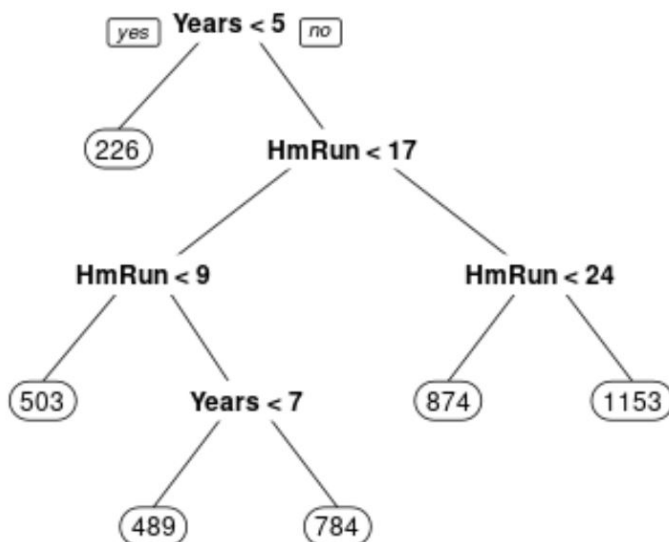
```
best <- tree$scptable[,"CP"]
```

```
# Produce a final, pruned tree based on the calculated optimal cp value
```

```
pruned_tree <- prune(tree, cp=best)
```

```
# Plot the pruned tree using the default prp() function settings
```

```
prp(pruned_tree)
```



The resulting default plot provides a foundational visualization of the pruned **decision tree** structure. While this initial diagram successfully outlines the primary decision paths, it often lacks the granular detail required for deep analytical scrutiny. To achieve enhanced clarity and integrate

more specific statistical metrics directly into the visualization, the ``prp()`` function offers a rich array of sophisticated customization options.

Advanced Visualization with ``prp()`` Customization

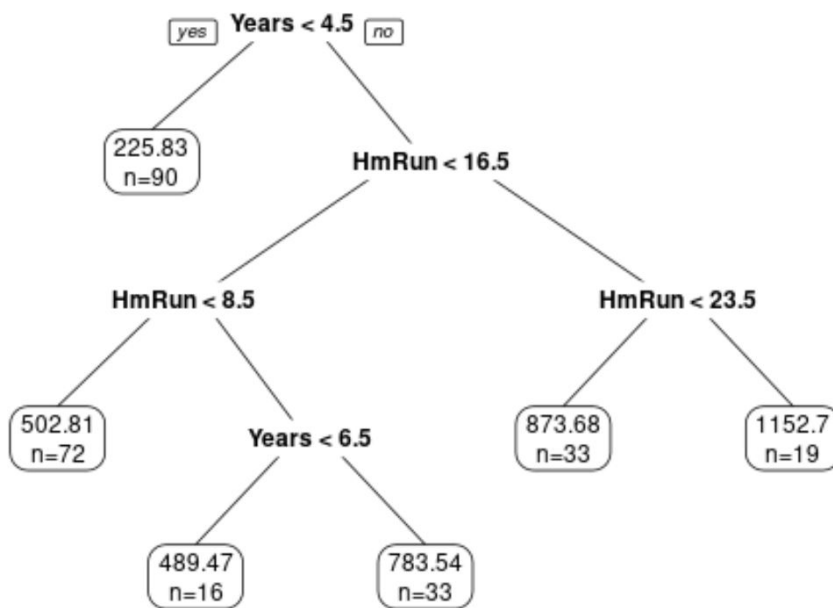
While the basic ``prp()`` output serves as an adequate initial representation, the true analytical value of the **rpart.plot** package is unlocked through its extensive customization capabilities. These powerful arguments empower users to meticulously control the visual appearance and the precise statistical information displayed within each node, significantly improving the tree's overall readability and analytical utility. Mastering these parameters is essential for transforming a basic plot into a publication-quality diagram that effectively communicates complex modeling insights.

Several key parameters are available for fine-tuning the plot. For instance, the ``faclen`` argument manages the displayed length of factor labels; ``extra`` is crucial as it determines what additional metrics are embedded within the nodes, such as the exact count of observations or percentages; ``roundint`` specifies whether numerical values should be forced into integers; and ``digits`` controls the necessary number of decimal places for displaying precise numerical outputs. By strategically adjusting this combination of arguments, analysts can generate a visualization that conveys highly precise details regarding every **terminal node** and the specific splits that lead to them, facilitating a deeper and more accurate interpretation.

The following code block demonstrates how these arguments are utilized to dramatically customize the visualization of the **decision tree**. This methodological approach ensures that all relevant data points and statistical measurements are presented clearly and accurately, thereby significantly aiding in a comprehensive understanding of the model's internal structure and its predictive mechanism.

Plot decision tree using custom arguments for enhanced clarity and detail

```
prp(pruned_tree,  
faclen=0, # Use full names for factor labels (0 means full length)  
extra=1, # Display the number of observations for each terminal node  
roundint=F, # Do not round the predicted values to integers  
digits=5) # Display 5 decimal places in the output for enhanced precision
```



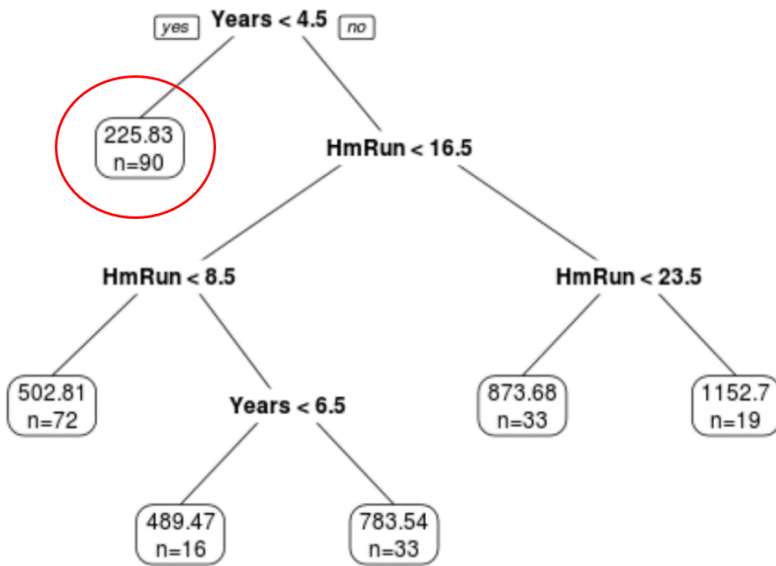
This customized visualization delivers a significantly richer and more analytically detailed view of the [regression tree](#). We can now clearly identify that the final model is composed of six distinct [terminal nodes](#). Each of these nodes represents a unique segment of the data and provides a specific predicted salary value. The inclusion of additional information--specifically, the exact number of players falling into each category and the precise predicted salary values--substantially elevates the interpretability and trustworthiness of the model's results.

Interpreting the Decision Tree and Deriving Predictions

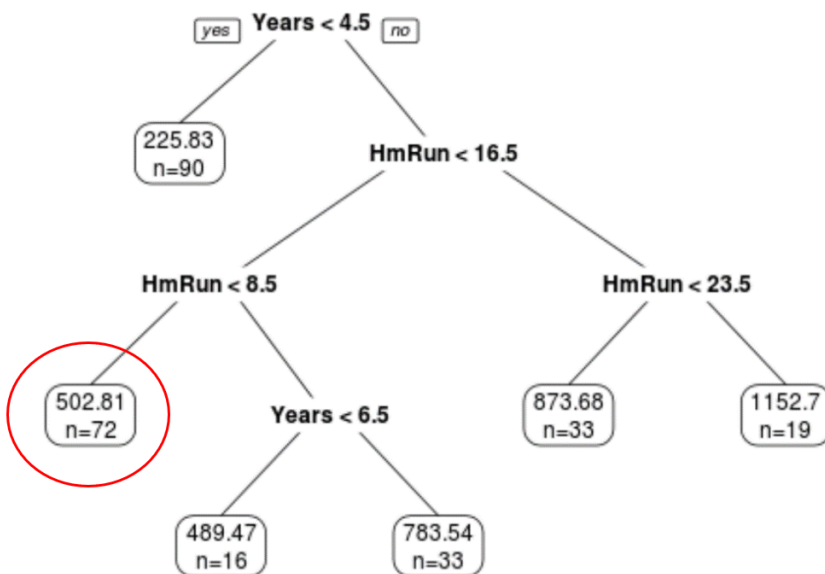
The process of interpreting a [decision tree](#) involves methodically tracing the sequential decision paths, starting from the root node at the top and proceeding downwards toward the final [terminal nodes](#). Every split encountered along a path is governed by a precise condition applied to one of the input [predictor variables](#). In the context of our baseball salary example, a highly influential initial split occurs based on the condition `Years < 4.5`. This critical bifurcation directs players with less than 4.5 years of experience down one distinct branch, while those with 4.5 years or more follow an alternative path, demonstrating the model's reliance on experience as a primary factor.

Each terminal node serves as the final output, summarizing two crucial pieces of information: the single predicted salary (the mean salary of all observations within that node) for players categorized into that segment, and the count of observations (the number of players from the dataset) belonging to that node. For a concrete example, by navigating the tree, we can isolate the terminal node encompassing players with less than 4.5 years of experience. This node explicitly reveals that 90 such players were present in the dataset, and their average predicted salary was approximately **\$225.83k**. This granular level of detail allows analysts to glean sharp insights into

precisely how specific player characteristics influence their final valuation.



One of the most compelling benefits of utilizing a [regression tree](#) is the straightforward ability to generate predictions for a new, unseen player based solely on their specific attributes (years of experience and home runs). The prediction mechanism is entirely intuitive: one simply follows the decision rules sequentially, starting from the root node until the path terminates at a leaf node. The value displayed at that final node constitutes the predicted salary. This transparent and easy-to-follow methodology makes decision trees exceptionally user-friendly for both individual prediction tasks and for articulating the rationale underpinning those predictions.



This inherent ease of both visualization and subsequent interpretation is the fundamental reason why [decision trees](#) are consistently favored across a wide spectrum of analytical contexts. They provide a clear, graphical mapping of complex decision rules, ensuring that the model's insights are accessible to both highly technical experts and non-technical stakeholders, thereby fostering superior communication and actionable intelligence.

Conclusion: Mastering Tree Visualization in R

This comprehensive tutorial successfully guided us through the entire workflow for developing, pruning, and effectively visualizing regression trees within the [R](#) environment. We skillfully utilized the foundational capabilities of the [rpart](#) package for precise model construction and then employed the sophisticated visualization features of the [rpart.plot](#) package, especially the versatile ``prp()`` function, to produce diagrams that are both clear and rich in information. Our illustrative example, based on the challenging **Hitters** dataset, provided a concrete demonstration of predicting player salaries utilizing critical performance indicators.

The capacity to interpret the decision-making logic of these models directly from their graphical representation is perhaps the most significant advantage they offer. By clearly understanding how distinct [predictor variables](#) contribute to the final prediction at every stage of the branching process, analysts gain profound insights into the underlying structure of their data. Furthermore, the robust customization options provided within the ``prp()`` function grant analysts the necessary flexibility to tailor their visualizations, ensuring that specific, high-impact aspects of the tree are prominently highlighted for maximum communicative effect.

Mastering the effective plotting of decision trees in R not only serves to expand your statistical analytical toolkit but also critically enhances your competence in communicating complex model findings in an accessible manner. We strongly encourage readers to continue experimenting with various datasets and exploring the full range of ``prp()`` arguments to fully appreciate the remarkable power and flexibility these essential packages bring to your future data analysis endeavors.

Further Learning Resources

For those dedicated to deepening their expertise in decision trees and related statistical modeling techniques, the following authoritative resources are highly recommended for continued study and practice:

Explore the official documentation for the [rpart](#) package on CRAN to understand its full range of functionalities and control parameters.

Consult the [rpart.plot](#) package documentation for an exhaustive list of ``prp()`` arguments and

advanced plotting techniques.

Review the [ISLR](#) package documentation, which accompanies the renowned textbook "An Introduction to Statistical Learning," for more detailed examples and theoretical background on decision trees and other statistical models.

Read relevant chapters on [decision tree](#) learning from foundational textbooks such as "An Introduction to Statistical Learning" or "Elements of Statistical Learning" for a deeper theoretical understanding of the methodology.