

# Introduction to Data Science: Treebased Methods

Héctor Corrada Bravo

University of Maryland, College Park, USA 2020-04-19



We saw in previous units the limitation of using linear methods for regression classification.

We saw in previous units the limitation of using linear methods for regression classification.

In this unit, we look at tree-based methods.

We saw in previous units the limitation of using linear methods for regression classification.

In this unit, we look at tree-based methods.

These are elegant and versatile methods that allow modeling of predictor space with regions that take complex, non-linear, shapes

We saw in previous units the limitation of using linear methods for regression classification.

In this unit, we look at tree-based methods.

These are elegant and versatile methods that allow modeling of predictor space with regions that take complex, non-linear, shapes

But still produce models that are interpretable.

We saw in previous units the limitation of using linear methods for regression classification.

In this unit, we look at tree-based methods.

These are elegant and versatile methods that allow modeling of predictor space with regions that take complex, non-linear, shapes

But still produce models that are interpretable.

We will concentrate on Regression and Decision Trees and their extension to Random Forests.

Consider a task where we are trying to predict a car's fuel consumption in miles per gallon based on the car's weight. A linear model in this case is not a good fit.



Let's take a look at what a regression tree estimates in this case.



The decision trees partitions the weight predictor into regions based on its value.



Outcome Y (mpg in this case) is predicted to be the mean *within each of the data partitions*.



Thus provides an empirical estimate of E[Y|X = x] where conditioning is given by this region partitioning.



## Tree models

Regression and decision trees operate by prediction an outcome variable Y by partitioning feature (predictor) space.

## Tree models

Regression and decision trees operate by prediction an outcome variable Y by partitioning feature (predictor) space.

In general, the regression tree model:

Partitions space into *J* non-overlapping regions, *R*<sub>1</sub>, *R*<sub>2</sub>, ..., *R*<sub>J</sub>.
For every observation that falls within region *R*<sub>j</sub>, predict response as mean of response for training observations in *R*<sub>j</sub>.

## Tree models

Regression and decision trees operate by prediction an outcome variable Y by partitioning feature (predictor) space.

In general, the regression tree model:

Partitions space into *J* non-overlapping regions, *R*<sub>1</sub>, *R*<sub>2</sub>, ..., *R*<sub>J</sub>.
For every observation that falls within region *R*<sub>j</sub>, predict response as mean of response for training observations in *R*<sub>j</sub>.

The important observation is that **Regression Trees create partitions** recursively

## **Tree Models**

For example, consider finding a good predictor j to partition space along its axis. A recursive algorithm would look like this:

• Find predictor j and value s that minimize RSS:

$$\sum_{i:\, x_i \in R_1(j,s))} (y_i - {\hat y}_{R_1})^2 + \sum_{i:\, x_i \in R_2(j,s))} (y_i - {\hat y}_{R_2})^2$$

Where  $R_1$  and  $R_2$  are regions resulting from splitting observations on predictor j and value s:

$$R_1(j,s) = X|X_j < s ext{ and } R_2(j,s)X|X_j \geq s$$

## Tree Models

For example, consider finding a good predictor j to partition space along its axis. A recursive algorithm would look like this:

• Find predictor j and value s that minimize RSS:

$$\sum_{i:\, x_i \in R_1(j,s))} (y_i - {\hat y}_{R_1})^2 + \sum_{i:\, x_i \in R_2(j,s))} (y_i - {\hat y}_{R_2})^2$$

• Apply recursively to regions  $R_1$  and  $R_2$ .

#### Tree Models

Within each region a prediction  $\hat{y}_{R_j}$  is made as the mean of the response Y of observations in  $R_j$ .



Consider building a model that used both horsepower and weight.

Here, value of the response Y is indicated by the size of the point.



This is what a decision tree would look like for these two predictors:



horsepower



*Quiz* What would this tree predict as mpg for an instance with variable values

- horsepower=85
- weight=2800

Classification, or decision trees, are used in classification problems, where the outcome is categorical.

Classification, or decision trees, are used in classification problems, where the outcome is categorical.

The same partitioning principle holds, but now, each region predicts the majority class for training observations within region.

Classification, or decision trees, are used in classification problems, where the outcome is categorical.

The same partitioning principle holds, but now, each region predicts the majority class for training observations within region.

The recursive partitioning method requires a score function to choose predictors (and values) to partition with.

Classification, or decision trees, are used in classification problems, where the outcome is categorical.

The same partitioning principle holds, but now, each region predicts the majority class for training observations within region.

The recursive partitioning method requires a score function to choose predictors (and values) to partition with.

A naive approach would looking for partitions that minimize training error.

Classification, or decision trees, are used in classification problems, where the outcome is categorical.

The same partitioning principle holds, but now, each region predicts the majority class for training observations within region.

The recursive partitioning method requires a score function to choose predictors (and values) to partition with.

A naive approach would looking for partitions that minimize training error.

Better performing approaches use more sophisticated metrics.

## **Decision Trees**

Let's look at how a classification tree performs on a credit card default dataset.





The predictor space

Suppose we have p predictor attributes  $X_1, \ldots, X_p$  and n observations.

The predictor space

Suppose we have p predictor attributes  $X_1, \ldots, X_p$  and n observations.

Each of the  $X_i$  can be

a) a numeric variable: there are n-1 possible splits b) an ordered factor (categorical variable): there are k-1 possible splits

c) an unordered factor:  $2^{k-1} - 1$  possible splits.

Learning Strategy

The general procedure for tree learning is the following:

**Grow**: an overly large tree using forward selection as follows: at each step, find the *best* split among all attributes. Grow until all terminal nodes either

(a) have < m (perhaps m = 1) data points (b) are "pure" (all points in a node have [almost] the same outcome).

Learning Strategy

The general procedure for tree learning is the following:

**Grow**: an overly large tree using forward selection

**Prune**: the tree back, creating a nested sequence of trees, decreasing in *complexity* 

Tree Growing

The recursive partitioning algorithm is as follows:

INITIALIZE All cases in the root node REPEAT Find optimal allowed split; Partition leaf according to split STOP Stop when pre-defined criterion is met

Tree Growing

An important issue in tree construction is how to use the training data to determine the binary splits of dataset  ${f X}$ 

**Tree Growing** 

An important issue in tree construction is how to use the training data to determine the binary splits of dataset  ${f X}$ 

The fundamental idea is to select each split of a subset so that the data in each of the descendent subsets are "purer" than the data in the parent subset.

#### Measures of impurity

Commonly used measures of impurity at a node i of a classification tree are

missclasification rate: 
$$\frac{1}{n_i} \sum_{j \in A_i} I(y_j \neq k_i) = 1 - \hat{p}_{ik_i}$$
  
entropy:  $\sum p_{ik} \log(p_{ik})$   
GINI index:  $\sum_{j \neq k} p_{ij} p_{ik} = 1 - \sum_k p_{ik}^2$ 

where  $k_i$  is the most frequent class in node i.

#### Measures of impurity

Commonly used measures of impurity at a node i of a classification tree are

missclasification rate: 
$$\frac{1}{n_i} \sum_{j \in A_i} I(y_j \neq k_i) = 1 - \hat{p}_{ik_i}$$
  
entropy:  $\sum p_{ik} \log(p_{ik})$   
GINI index:  $\sum_{j \neq k} p_{ij} p_{ik} = 1 - \sum_k p_{ik}^2$ 

where  $k_i$  is the most frequent class in node i.

In practice, the GINI index is preferred

For regression trees we use the residual sum of squares:

$$D = \sum_{ ext{cases } j} (y_j - \mu_{[j]})^2$$

where  $\mu_{[j]}$  is the mean values in the node that case j belongs to.

Good properties of Regression and Classification trees include:

- Decision trees are very "natural" constructs, in particular when the explanatory variables are catgorical (and even better when they are binary)
- Trees are easy to explain to non-data analysts
- The models are invariant under transformations in the predictor space

Good properties of Regression and Classification trees include:

- Multi-factor responses are easily dealt with
- The treatment of missing values is more satisfactory than for most other models
- The models go after interactions immediately, rather than as an afterthought
- Tree growth is much more efficient than described here

However, they do have important issues to address

- Tree space is huge, so we may need lots of data
- We might not be able to find the *best* model at all as it is a greedy algorithm
- It can be hard to assess uncertainty in inference about trees

However, they do have important issues to address

- Results can be quite variable (tree selection is not very stable)
- Simple trees usually don't have a lot of predictive power

Random Forests are a **very popular** approach that addresses these shortcomings via resampling of the training data.

Random Forests are a **very popular** approach that addresses these shortcomings via resampling of the training data.

Their goal is to improve prediction performance and reduce instability by *averaging* multiple decision trees (a forest constructed with randomness).

It uses two ideas to accomplish this. The first idea is *Bagging* (bootstrap aggregation)

General scheme:

- 1. Build many decision trees  $T_1, T_2, \ldots, T_B$  from training set
- 2. Given a new observation, let each  $T_j$  predict  $\hat{y}_j$
- 3. For regression: predict average  $\frac{1}{B}\sum_{j=1}^{B}\hat{y}_{j}$ , for classification: predict with majority vote (most frequent class)

How do we get many decision trees from a single training set?

Use the *bootstrap* resampling technique.



How do we get many decision trees from a single training set?

To create  $T_j, j = 1, \ldots, B$  from training set of size n:

a) create a bootstrap training set by sampling n observations from training set with replacement



How do we get many decision trees from a single training set?

To create  $T_j, j = 1, \ldots, B$  from training set of size n:

b) build a decision tree frombootstrap training set



The second idea used in Random Forests is to use a random selection of features to split when deciding partitions.

The second idea used in Random Forests is to use a random selection of features to split when deciding partitions.

Specifically, when building each tree  $T_j$ , at each recursive partition:

The second idea used in Random Forests is to use a random selection of features to split when deciding partitions.

Specifically, when building each tree  $T_j$ , at each recursive partition:

only consider a randomly selected subset of predictors to find best split.

The second idea used in Random Forests is to use a random selection of features to split when deciding partitions.

Specifically, when building each tree  $T_j$ , at each recursive partition:

only consider a randomly selected subset of predictors to find best split.

This reduces correlation between trees in forest, improving prediction accuracy.

Let's look at the same car dataset again and plot predicted vs. true miles per gallon given by a random forest and a regression tree.



50/61

#### Now let's look at the same plot on a *testing* dataset.



A disadvantage of random forests is that we lose interpretability.

A disadvantage of random forests is that we lose interpretability.

However, we can use the fact that a bootstrap sample was used to construct trees to measure *variable importance* from the random forest.

A disadvantage of random forests is that we lose interpretability.

However, we can use the fact that a bootstrap sample was used to construct trees to measure *variable importance* from the random forest.

Since we used bootstrap samples we can get out-of-bag (OOB) samples for each tree in the random forest.

When the bth tree is constructed, use the OOB samples as follows

1. Compute error rate for the OOB samples

2. For each predictor j:

a. permute its values in the OOB samples and recompute error rateb. calculate increase in error rate

Report increase in error rate over all bootstrap samples

Here is a table of *variable importance* for the random forest we just constructed.

	%IncMSE	IncNodePurity
cylinders	13.22	2238.51
displacement	17.79	2161.24
horsepower	17.16	1711.24
weight	20.54	3087.59
acceleration	13.64	442.72
year	42.08	1813.08

And a plot of variable importance



Variable Importance

%IncMSE

Tree-based methods are highly interpretable *prediction* models.

Tree-based methods are highly interpretable *prediction* models.

Some inferential tasks are possible (e.g., variable importance in random forests), but are much more limited than linear models.

Tree-based methods are highly interpretable *prediction* models.

Some inferential tasks are possible (e.g., variable importance in random forests), but are much more limited than linear models.

These methods are very commonly used across many application domains

Tree-based methods are highly interpretable *prediction* models.

Some inferential tasks are possible (e.g., variable importance in random forests), but are much more limited than linear models.

These methods are very commonly used across many application domains

Random Forests often perform at state-of-the-art for many tasks.