

Machine Learning for All: Tree based models

Anastasiya Yarygina

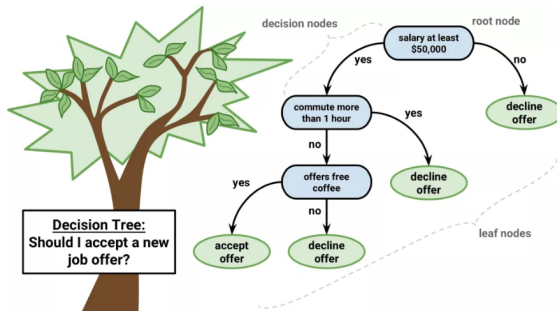
Friday, February 15, 2019

Decision Tree Algorithms in Machine Learning

In this lecture we will cover the following topics:

- ▶ Decision trees: Using tree-logic to make predictions
 - ▶ Regression single-tree models
 - ▶ Random Forest
 - ▶ If we have time: Boosting
- ▶ Examples:
 - ▶ Boston Housing, Kaggle
 - ▶ California Housing Prices, Kaggle
- ▶ Extra practice: Classification tree models using Iris dataset

What is a Decision Tree?



Tree-logic uses series of steps to come to a conclusion. Each decision is a **node**, and the final prediction is a **leaf node**.

Decision Trees in ML

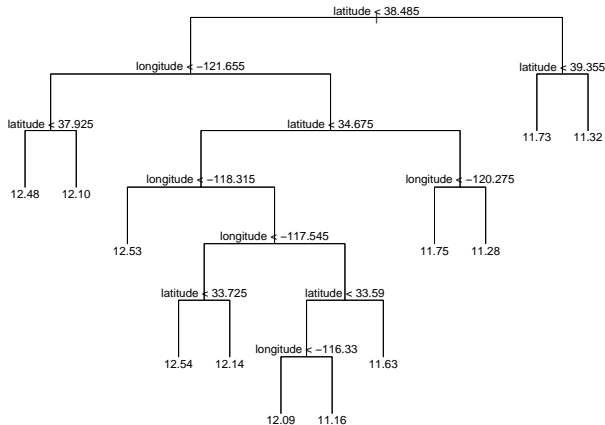
- ▶ **Decision Tree Algorithm** is a supervised learning algorithm that can be used for solving
 - ▶ **regression** and
 - ▶ continuous response variable
 - ▶ wage example from last week
 - ▶ **classification** problems
 - ▶ categorical or factor response variable
 - ▶ spam example from last week
- ▶ Classically, the name of this algorithm is **Decision Tree**
 - ▶ Some platforms like R use a modern term **CART** (Classification and Regression Trees)
- ▶ Objective: obtain **predictions**
 - ▶ of the **response variable** Y (dependent variable or output)
 - ▶ from the **input variables** X_1, X_2, \dots, X_n (features, predictors).

Predictions using Decision Trees

- ▶ **Key Idea:** Decision Tree **splits** the data into
 - ▶ two or more **homogeneous data segments**
 - ▶ based on the **best splitter**, which is a variable taken from the inputs X_1, X_2, \dots, X_n .
- ▶ Every time we split the sample we make a decision. Each decision is a **decision node**, and the final prediction is a **leaf node**.
- ▶ Example: fit a tree that predicts for each property **log price** using as inputs **longitude** and **latitude**.

Exmpale: California Housing dataset

We can fit a tree that predicts for each property **log price** using as inputs **longitude** and **latitude**.

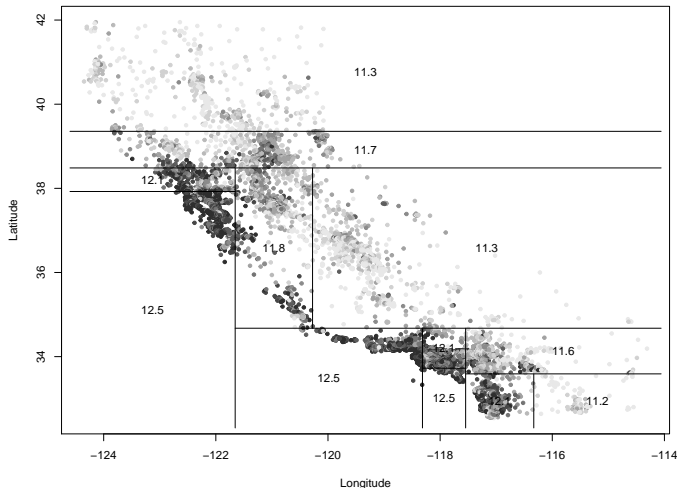


Exmpale: California Housing dataset

- ▶ The tree has 11 **decision nodes**, which are the nodes where the splitting of the data takes place.
 - ▶ Node 1: root node, split on latitude at 38.485
- ▶ And there are 12 **leaf nodes**
 - ▶ Leaf node 12: value 11.32
- ▶ The data space is partitioned in to 12 **homogeneous regions**.
- ▶ How do these **homogeneous regions** look like?

Exmpale: California Housing dataset

Overlay log price of properties on predicted partitions. Darker dots represent more expensive properties.



Exmpale: California Housing dataset

The tree model divided the **predictor space** (longitude and latitude in this case) into 12 **distinct** and **non-overlapping** rectangular **regions** $R_1, R_2, \dots, R_j, \dots, R_{12}$.

If there are more than two inputs, the data space is split in some kind of hyper-rectangles.

For every observation that falls into a given region R_j , the model assigns its **predicted value**, which is the **mean of the response** Y (log price in this case) for all observations in region R_j .

The regions with the log average value 12.5 are LA and the Bay Area.

Decision Tree Algorithm

To get homogeneous segments, the model makes **optimal splits**.

Each **optimal split** is made:

- ▶ at **certain value of some predictor** X_i ,
- ▶ so that the child set to the left of the split and the child set to the right of the split are **as homogeneous in response Y as possible**.

In regression trees **homogeneity** is measured by the **Sum of Squared Errors (SSE)**:

$$\text{sum}(y - \text{prediction}_{\text{left}})^2 + \text{sum}(y - \text{prediction}_{\text{right}})^2$$

Each **optimal split minimizes the SSE** to the left and to the right of the split.

Decision Tree Algorithm

Decision trees are fit in a **top-down, greedy** approach, which is also known as a **recursive binary splitting**.

- ▶ **Top-down**: it starts at the top of the tree
- ▶ **Greedy**: at each step the best split is made at *that* particular step, we do not look ahead and pick the split that will lead to a better tree in some future step.

Each split improves the fit of the tree (think of R^2 and adding new variables in a regression model).

The algorithm stops when:

- ▶ improvement in the fit is below some predefined threshold (default 0.01)
- ▶ number of observations in leafs is below some predefined threshold (default 5)

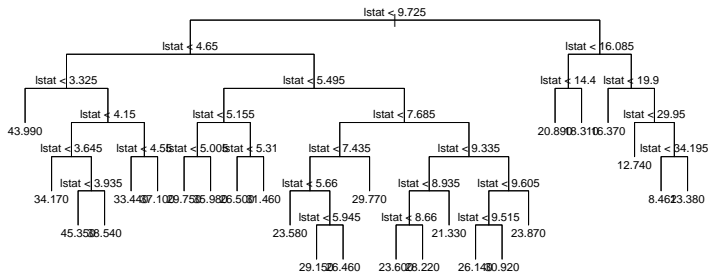
Practice: *tree* package and Boston Housing dataset

Fit a single tree model

- to predict **median value of properties** using
- **low income status** as predictor.

Practice: *tree* package and Boston Housing dataset

Fit a tree to predict median value of properties using low income status as predictor.



The big tree size is: 26

Pruning

Tree models are very flexible and **tend to overfit**.

To avoid overfitting trees are **pruned** by removing nodes and branches from the bottom up.

At each step we remove the split that contributes least to improvement in the fit.

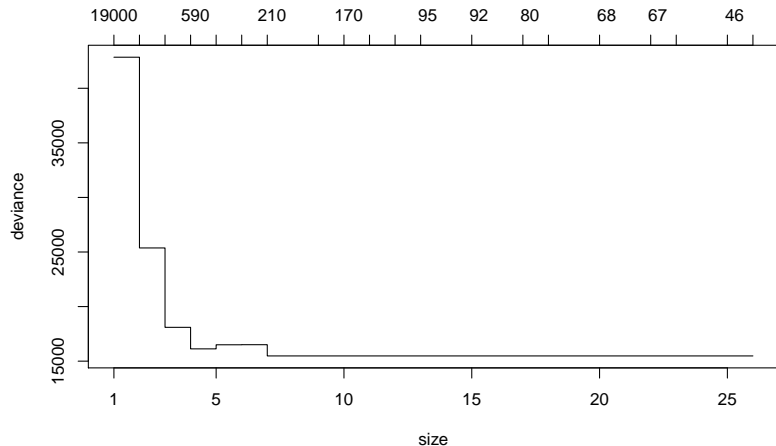
Pruning produces a set of **candidate trees** of different sizes.

What is the best tree size (what is the best number of nodes)?

We use **Cross Validation** to choose the best tree size.

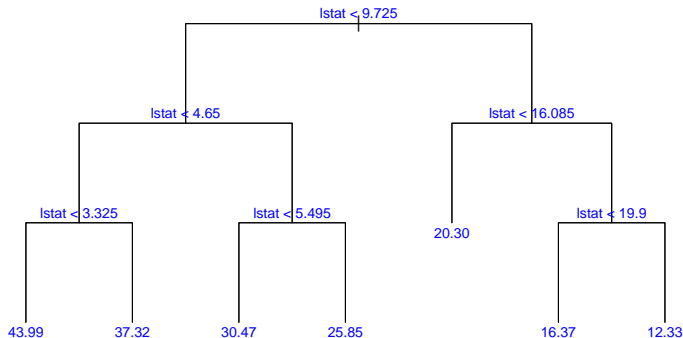
The tree with the best size has the smallest SSE.

Practice: *tree* package and Boston Housing dataset



CV and choose the size that minimizes SSE

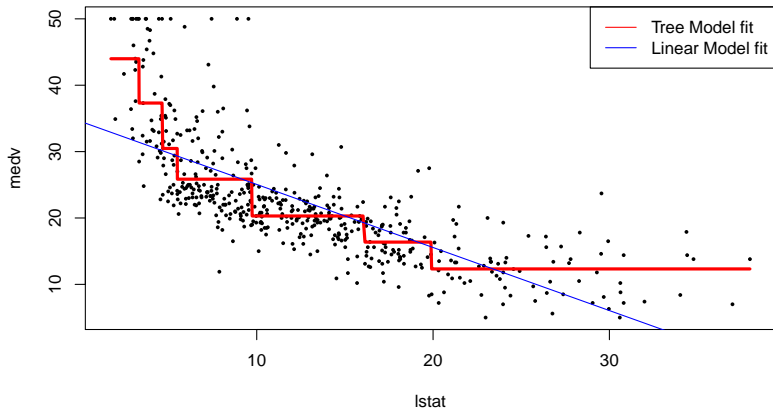
Practice: *tree* package and Boston Housing dataset



the size of the pruned tree is 7

Practice: *tree* package and Boston Housing dataset

Compare Tree Model fit and Linear Model fit

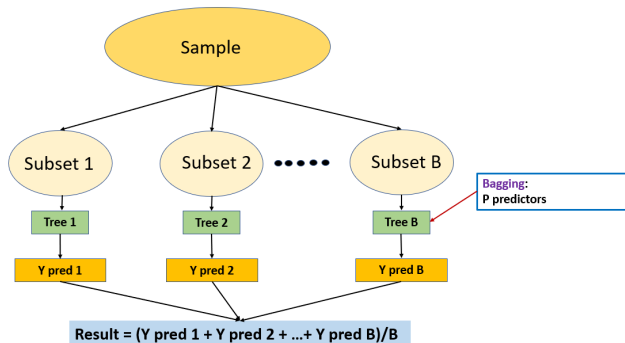


Aggregating Trees

- ▶ Single tree algorithms are effective in making reasonable predictions
- ▶ Single tree models are **weak learners**
- ▶ How can we get a **strong learner**? Aggregate predictions from many weak learners
- ▶ To improve predictions we can:
 - ▶ fit many tree models from the same data
 - ▶ and average predictions across these models.

Aggregating Trees: Bagging

- ▶ This is exactly what **Bagging** (or **Bootstrap aggregation**) procedures do. The steps are the following:
 - ▶ Sample (Bootstrap) B subsets of the data
 - ▶ Fit a tree to each subset to get B fitted trees
 - ▶ Average predictions across trees

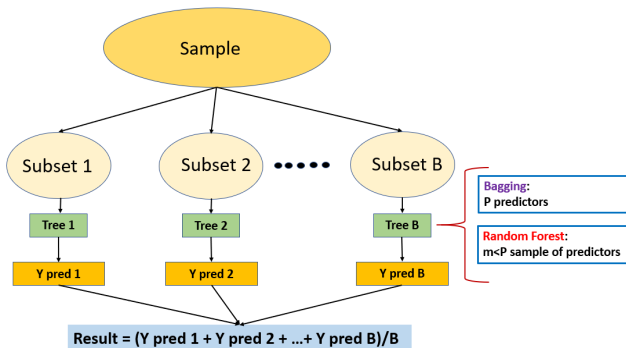


Aggregating Trees: Random Forest

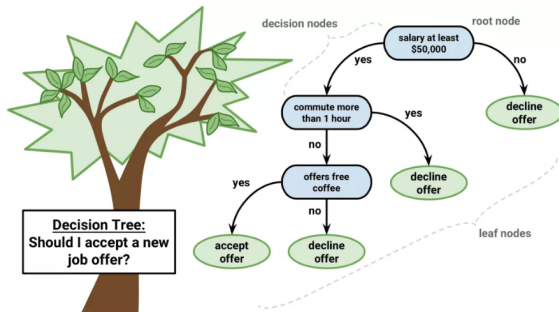
Random Forest is a special case of Bagging. It provides an improvement over bagged trees by way of a small tweak:

- ▶ Random Forest builds B trees on subsets of the data.
- ▶ But, for each split it randomly choose a **sample of m predictors** of all available p predictors (default $m = \sqrt{p}$).

Random Forest tuning parameters are B and m .



Aggregating Trees: Why is Random Forest better than Bagging?



Bagging: **all** p predictors \rightarrow **similar trees**

Random Forest: **selection of** $m < P$ predictors \rightarrow **different trees**

Practice: Regression Trees using California Housing data

Objectives:

- ▶ Fit Single tree model and Random Forest model
- ▶ Fit linear model
- ▶ Compare predictive capacity using Out of Sample (OOS) Mean Root Squared Error (MRSE), which is a $\sqrt{\text{SSE}}$.

We fit models on **training partition** and we evaluate their predictive capacity on **testing partition**.

Fit Single Tree model using *rpart* package

```
## size of big tree: 134
```

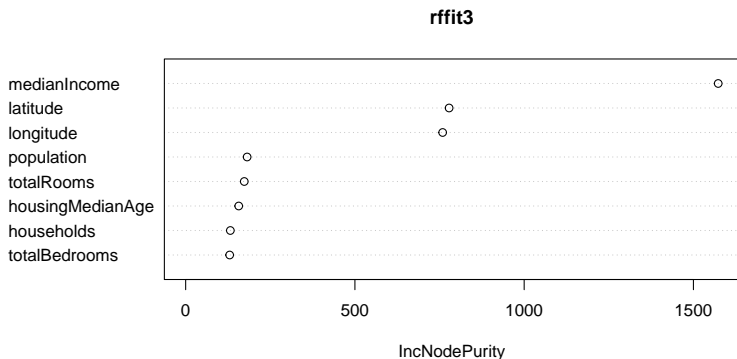
```
## RMSE Single Tree Model: 0.2937906
```

Fit Random Forest model using *randomForest* package

Ideally, fit many models with different tuning parameters. Choose the one with the best OOS RMSE.

```
## RMSE on test for RF m=3 ntree=50: 0.2472154
```

```
## Variable Importance RF
```



Fit Linear model and compare OOS RMSE

- ▶ Linear model OOS predictive capacity:

```
## RMSE on test for linear model: 0.3485098
```

- ▶ Now, let's compare OOS predictive capacity of all models

```
##      rmse_rpart  rmse_rf3   rmse_lm  
## [1,] 0.2937906 0.2472154 0.3485098
```

- ▶ Which model does the best job?
- ▶ Recall: we want models to have **low** RMSE

Takeaways

- ▶ Decision Trees are simple machine learning algorithms and they are easy to visualize.
- ▶ Trees are useful to solve regression and classification problems.
- ▶ However, single tree models tend to overfit.
- ▶ **Ensembling methods** such as Random Forest are good for improving predictive capacity of trees. They work growing many trees and combining predictions of the resulting ensemble of trees.
- ▶ Random Forest is among the state-of-the-art methods for supervised learning. However, it is computationally intensive and its results are difficult to interpret.

More on Ensembling methods: Boosting

Boosting builds many decision trees, but unlike Bagging or Random Forest, Boosting trees are grown **sequentially**. The steps are the following:

- ▶ Fit the model **tree#1** on the original data and save the residuals
- ▶ Fit the model **tree#2** on the residuals
- ▶ Update the initial model: **tree#3 = tree#1 + tree#2**
- ▶ Update the residuals
- ▶ Fit a new model **tree#4** on the residuals
- ▶ ...
- ▶ Repeat the process for a specified number of iterations

Updated trees are **shrunk** by the **shrinkage parameter** λ which controls the rate at which algorithm learns (default = 0.001 to 0.01).

Other tuning parameters: d the number of splits in each tree and B the number of trees to grow.

Practice: Fit Boosting model using *gbm*¹ package

Ideally, fit many models with different tuning parameters. Choose the one with the best OOS RMSE.

```
## RMSE on test for Boosting 4; 1000; 2: 0.2402014
```

```
## Variable Importance Boosting
```

| ## | var | rel.inf |
|---------------------|------------------|-----------|
| ## medianIncome | medianIncome | 45.594766 |
| ## longitude | longitude | 17.294545 |
| ## latitude | latitude | 17.137733 |
| ## population | population | 5.435569 |
| ## totalBedrooms | totalBedrooms | 4.283103 |
| ## totalRooms | totalRooms | 3.733635 |
| ## households | households | 3.392255 |
| ## housingMedianAge | housingMedianAge | 3.128394 |

¹GBM: Gradient Boost Machine

Compare OOS RMSE

- ▶ Let's compare predictive capacity of Single tree model, Random Forest, Boosting and Linear model:

```
##          rmse_rpart  rmse_rf3 rmse_gbm3  rmse_lm
## [1,]  0.2937906 0.2472154 0.2402014 0.3485098
```

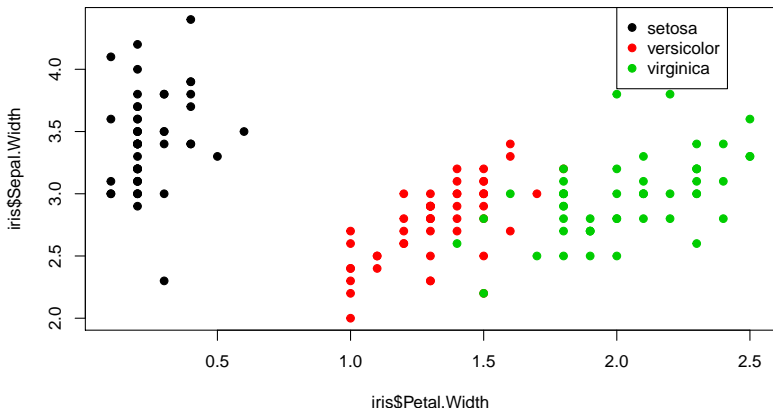
- ▶ Which model does the best job?
- ▶ Recall: we want models to have **low** RMSE

Extra: Classification problema using *iris* dataset

What happens if our problem is a **classification problem**?

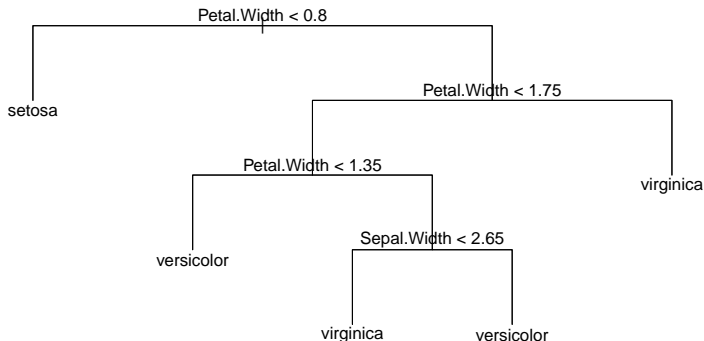
iris dataset: sepal and petal length and width, 150 plants, 3 species

- Setosa, Versicolor, Virginica.



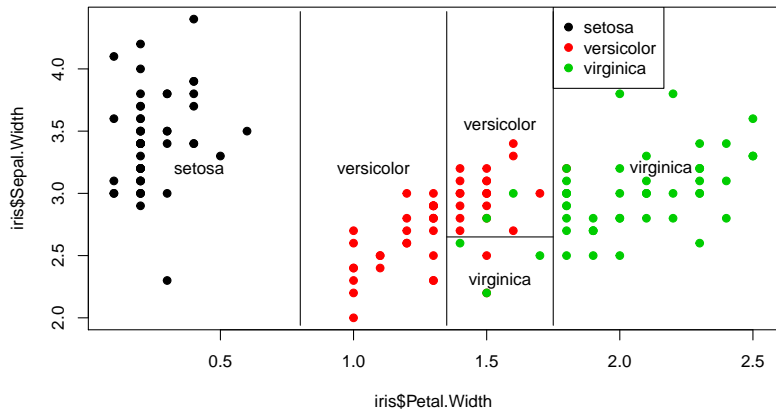
Exmpale: *iris* dataset

Fit a model that predicts species using as inputs sepal and petal width.



► Prediction is the **most common category**

Exmpale: *iris* dataset



Partitions are defined by the classification tree. The first node classifies plants with petal width < 0.8 as *setosa*. Next, all plants with petal width > 1.75 are *virginica*.

Extra practice: Classification Trees using *iris* dataset

Objectives:

- ▶ Fit Single Tree, Random Forest, Boosting models
- ▶ Fit multinomial model
- ▶ Compare predictive capacity using OOS **Accuracy**

Fit models on **training partition** and evaluate their predictive capacity on **testing partition**.

Fit Single Tree model

► Classification table Single Tree model

```
##          rpartfitpred
##          setosa versicolor virginica
## setosa          30          0          0
## versicolor       0          20          1
## virginica         0          1         23
```

Fit RF model

► Classification table RF

| ## | | rfritpred | | |
|----|------------|-----------|------------|-----------|
| ## | | setosa | versicolor | virginica |
| ## | setosa | 30 | 0 | 0 |
| ## | versicolor | 0 | 20 | 1 |
| ## | virginica | 0 | 1 | 23 |

Fit Boosting model

► Classification table Boosting model

| ## | | gbmfitpredcat | | |
|----|------------|---------------|----|----|
| ## | | 1 | 2 | 3 |
| ## | setosa | 30 | 0 | 0 |
| ## | versicolor | 0 | 19 | 2 |
| ## | virginica | 0 | 1 | 23 |

Fit Multinomial model

► Classification table Multinomial model

| ## | | mnfitpred | | |
|----|------------|-----------|------------|-----------|
| ## | | setosa | versicolor | virginica |
| ## | setosa | 30 | 0 | 0 |
| ## | versicolor | 0 | 20 | 1 |
| ## | virginica | 0 | 0 | 24 |

Compare OOS Accuracy

- ▶ Now, compare OOS predictive capacity of all models

```
##      rpart_acc      rf_acc gbm_acc      mn_acc
## [1,] 0.9733333 0.9733333      0.96 0.9866667
```

- ▶ Which model does the best job?
- ▶ Recall: we want models to have **high** accuracy