

# Lecture Outline

---

Motivation

Decision Trees

Splitting Criteria

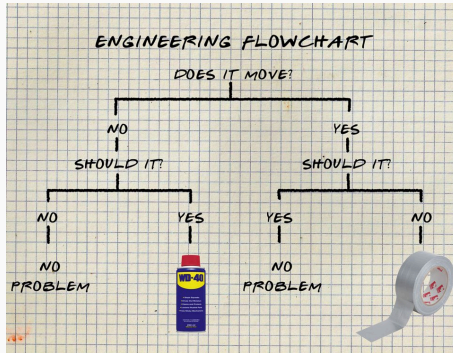
Stopping Conditions & Pruning

## Motivation

---

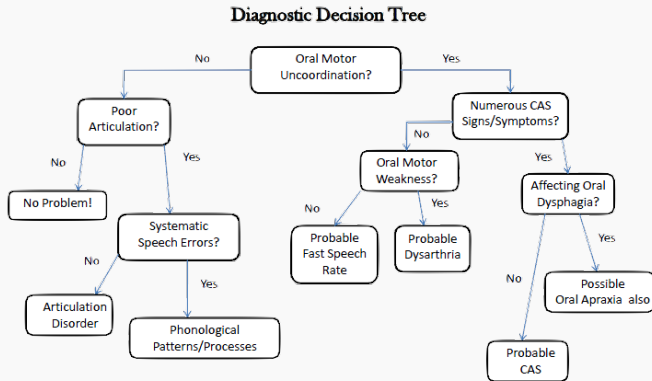
# Interpretable Models

But people in every walk of life have long been using interpretable models for differentiating between classes of objects and phenomena:



# Interpretable Models

But people in every walk of life have long been using interpretable models for differentiating between classes of objects and phenomena:



It turns out that the simple flow charts in our examples can be formulated as mathematical models for classification and these models have the properties we desire; they are:

1. interpretable by humans
2. have sufficiently complex decision boundaries
3. the decision boundaries are locally linear, each component of the decision boundary is simple to describe mathematically.

# The Geometry of Flow Charts

---

Flow charts whose graph is a tree (connected and no cycles) represents a model called a **decision tree**.

Formally, a **decision tree model** is one in which the final outcome of the model is based on a series of comparisons of the values of predictors against threshold values.

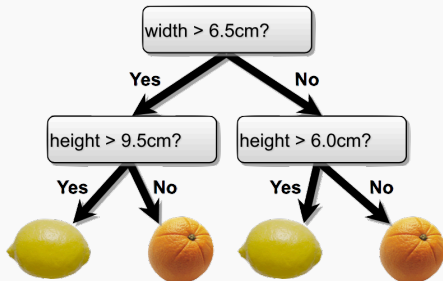
In a graphical representation (flow chart),

- ▶ the internal nodes of the tree represent attribute testing
- ▶ branching in the next level is determined by attribute value
- ▶ leaf nodes represent class assignments

# The Geometry of Flow Charts

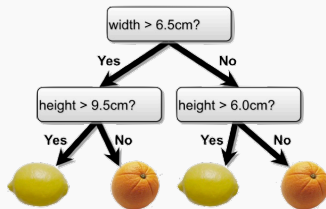
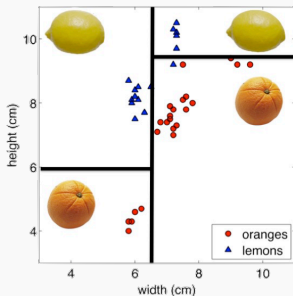
Flow charts whose graph is a tree (connected and no cycles) represents a model called a **decision tree**.

Formally, a **decision tree model** is one in which the final outcome of the model is based on a series of comparisons of the values of predictors against threshold values.



# The Geometry of Flow Charts

Every flow chart tree corresponds to a partition of the feature space by axis aligned lines or (hyper) planes. Conversely, every such partition can be written as a flow chart tree.



Each comparison and branching represents splitting a region in the feature space. Typically, at each iteration, we split once along one dimension (one predictor).



# Learning the Model

---

Given a training set, learning a decision tree model for binary classification means to produce an 'optimal' partition of the feature space with axis aligned linear boundaries, wherein each region is given a class label based on the largest class of the training points in that region.

# Learning the Model

---

Learning the smallest ‘optimal’ decision tree for any given set of data is NP complete for numerous simple definitions of ‘optimal’. Instead, we will seek a reasonably model using a greedy algorithm.

1. Start with an empty decision tree (undivided feature space)
2. Choose the ‘optimal’ predictor on which to split and choose the ‘optimal’ threshold value for splitting.
3. Recurse on on each new node until **stopping condition** is met

Now, we need only define our splitting criterion and stopping condition.

# Numerical vs Categorical Attributes

---

Note that the compare and branch method by which we defined regression tree works well for numerical features.

However, if a feature is categorical (with more than two possible values), comparisons like *feature < threshold* does not make sense.

A simple solution is to encode the values of a categorical feature using numbers and treat this feature like a numerical variable.

This is indeed what some computational libraries (e.g. `sklearn`) do, however, this method has drawbacks.

# Numerical vs Categorical Attributes

## Example

Suppose the feature we want to split on is **color**, and the values are: Red, Blue and Yellow. If we encode the categories numerically as:

$$\text{Red} = 0, \text{Blue} = 1, \text{Yellow} = 2$$

Then the possible non-trivial splits on **color** are

$$\{\{\text{Red}\}, \{\text{Blue}, \text{Yellow}\}\}, \quad \{\{\text{Red}, \text{Blue}\}, \{\text{Yellow}\}\}$$

But if we encode the categories numerically as:

$$\text{Red} = 2, \text{Blue} = 0, \text{Yellow} = 1$$

The possible splits are

$$\{\{\text{Blue}\}, \{\text{Yellow}, \text{Red}\}\}, \quad \{\{\text{Blue}, \text{Yellow}\}, \{\text{Red}\}\}$$

***Depending on the encoding, the splits we can optimize over can be different!***

## Splitting Criteria

---

# Optimality of Splitting

---

While there is no ‘correct’ way to define an optimal split, there are some common sensical guidelines for every splitting criterion:

- ▶ the regions in the feature space should grow progressively more pure with the number of splits. That is, we should see each region ‘specialize’ towards a single class.
- ▶ the fitness metric of a split should take a differentiable form (making optimization possible)
- ▶ we shouldn’t end up with empty regions - regions containing no training points.

# Classification Error

---

Suppose we have  $J$  number of predictors and  $K$  classes.

Suppose we select the  $j$ -th predictor and split a region containing  $N$  number of training points along the threshold  $t_j \in \mathbb{R}$ .

We can assess the quality of this split by measuring the classification error made by each newly created region,  $R_1, R_2$ :

$$\text{Error}(i|j, t_j) = 1 - \max_k p(k|R_i)$$

where  $p(k|R_i)$  is the proportion of training points in  $R_i$  that are labeled class  $k$ .

## Example

	Class 1	Class 2	Error( $i j, t_j$ )
$R_1$	0	6	$1 - \max\{6/6, 0/6\} = 0$
$R_2$	5	8	$1 - \max\{5/13, 8/13\} = 5/13$

We can now try to find the predictor  $j$  and the threshold  $t_j$  that minimizes the average classification error over the two regions, weighted by the population of the regions:

$$\min_{j, t_j} \left\{ \frac{N_1}{N} \text{Error}(1|j, t_j) + \frac{N_2}{N} \text{Error}(2|j, t_j) \right\}$$

where  $N_i$  is the number of training points inside region  $R_i$ .



# Gini Index

---

Suppose we have  $J$  number of predictors,  $N$  number of training points and  $K$  classes.

Suppose we select the  $j$ -th predictor and split a region containing  $N$  number of training points along the threshold  $t_j \in \mathbb{R}$ .

We can assess the quality of this split by measuring the purity of each newly created region,  $R_1, R_2$ . This metric is called the **Gini Index**:

$$\text{Gini}(i|j, t_j) = 1 - \sum_k p(k|R_i)^2$$

**Question:** What is the effect of squaring the proportions of each class? What is the effect of summing the squared proportions of classes within each region?

## Example

	Class 1	Class 2	Gini( $i j, t_j$ )
$R_1$	0	6	$1 - (6/6^2 + 0/6^2) = 0$
$R_2$	5	8	$1 - [(5/13)^2 + (8/13)^2] = 80/169$

We can now try to find the predictor  $j$  and the threshold  $t_j$  that minimizes the average Gini Index over the two regions, weighted by the population of the regions:

$$\min_{j, t_j} \left\{ \frac{N_1}{N} \text{Gini}(1|j, t_j) + \frac{N_2}{N} \text{Gini}(2|j, t_j) \right\}$$

where  $N_i$  is the number of training points inside region  $R_i$ .

# Information Theory

The last metric for evaluating the quality of a split is motivated by metrics of uncertainty in information theory.

Ideally, our decision tree should split the feature space into regions such that each region represents a single class. In practice, the training points in each region is distributed over multiple classes, e.g.:

	Class 1	Class 2
$R_1$	1	6
$R_2$	5	6

However, though both imperfect,  $R_1$  is clearly sending a stronger 'signal' for a single class (Class 2) than  $R_2$ .

# Information Theory

---

One way to quantify the strength of a signal in a particular region is to analyze the distribution of classes within the region. We compute the **entropy** of this distribution.

For a random variable with a discrete distribution, the entropy is computed by

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$

Higher entropy means the distribution is uniform-like (flat histogram) and thus values sampled it are ‘less predictable’ (all possible values are equally probable).

Lower entropy means the distribution has more defined peaks and valleys and thus values sampled from it are ‘more predictable’ (values around the peaks are more probable).

# Entropy

---

Suppose we have  $J$  number of predictors,  $N$  number of training points and  $K$  classes.

Suppose we select the  $j$ -th predictor and split a region containing  $N$  number of training points along the threshold  $t_j \in \mathbb{R}$ .

We can assess the quality of this split by measuring the entropy of the class distribution in each newly created region,  $R_1, R_2$ :

$$\text{Entropy}(i|j, t_j) = - \sum_k p(k|R_i) \log_2[p(k|R_i)]$$

**Note:** we are actually computing the conditional entropy of the distribution of training points amongst the  $K$  classes given that the point is in region  $i$ .

## Example

	Class 1	Class 2	Entropy( $i j, t_j$ )
$R_1$	0	6	$-(\frac{6}{6} \log_2 \frac{6}{6} + \frac{0}{6} \log_2 \frac{0}{6}) = 0$
$R_2$	5	8	$-(\frac{5}{13} \log_2 \frac{5}{13} + \frac{8}{13} \log_2 \frac{8}{13}) \approx 1.38$

We can now try to find the predictor  $j$  and the threshold  $t_j$  that minimizes the average entropy over the two regions, weighted by the population of the regions:

$$\min_{j, t_j} \frac{N_1}{N} \text{Entropy}(1|j, t_j) + \frac{N_2}{N} \text{Entropy}(2|j, t_j)$$

## Stopping Conditions & Pruning

# Stopping Conditions

---

Common simple stopping conditions:

- ▶ Don't split a region if all instances in the region belong to the same class
- ▶ Don't split a region if the number of instances in the sub-region will fall below pre-defined threshold
- ▶ Don't split a region if the total number of leaves in the tree will exceed pre-defined threshold

The appropriate thresholds can be determined by evaluating the model on a held-out data set or, better yet, via cross-validation.



# Stopping Conditions

---

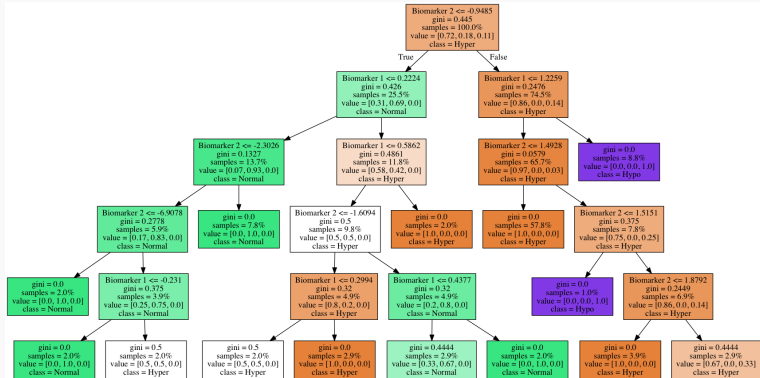
More restrictive stopping conditions:

- ▶ Don't split a region if the class distribution of the training points inside the region are independent of the predictors
- ▶ Compute the gain in purity, information or reduction in entropy of splitting a region  $R$

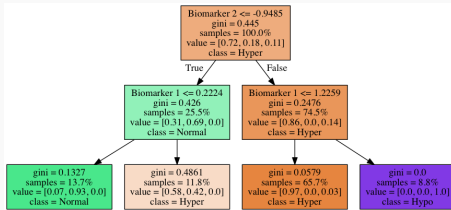
$$\text{Gain}(R) = \Delta(R) = m(R) - \frac{N_1}{N}m(R_1) - \frac{N_2}{N}m(R_2)$$

where  $m$  is a metric like the Gini Index or entropy.  
Don't split if the gain is less than some pre-defined threshold.

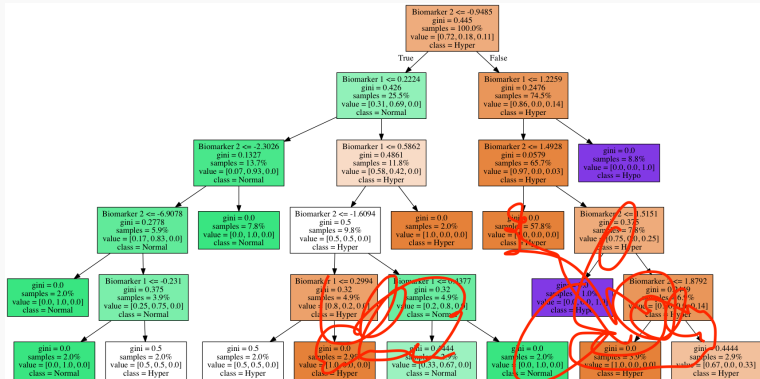
# Motivation for Pruning



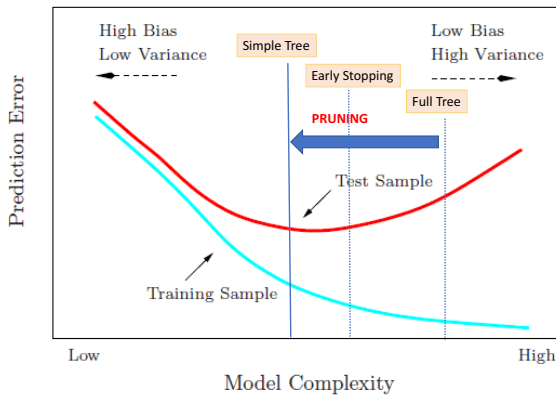
# Motivation for Pruning



# Motivation for Pruning



# Motivation for Pruning



# Pruning

The pruning algorithm:

1. Start with a full tree  $T_0$  (each leaf node contains exactly one training point)
2. Replace a subtree in  $T_0$  with a leaf node to obtain a pruned tree  $T_1$ . This subtree should be selected to minimize

$$\frac{\text{Error}(T_0) - \text{Error}(T_1)}{|T_0| - |T_1|}$$

3. Iterate this pruning process to obtain  $T_0, T_1, \dots, T_L$ , where  $T_L$  is the tree containing just the root of  $T_0$ .
4. Select the optimal tree  $T_i$  by cross validation.

**Note:** you might wonder where we are computing the cost-complexity  $C(T_i)$ . One can prove that this process is equivalent to explicitly optimizing  $C$ .

# Lecture Outline

---

Review

Decision Trees for Regression

Bagging

Random Forests

## Review

---

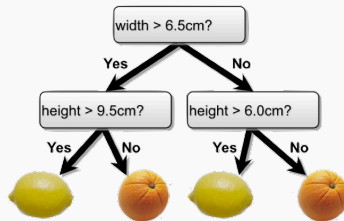
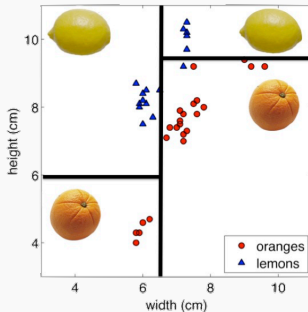


# Decision Trees

A **decision tree model** is an interpretable model in which the final output is based on a series of comparisons of the values of predictors against threshold values.

Graphically, decision trees can be represented by a flow chart.

Geometrically, the model partitions the feature space wherein each region is assigned a response variable value based on the training points contained in the region.



# Learning Algorithm

---

To learn a decision tree model, we take a greedy approach:

1. Start with an empty decision tree (undivided feature space)
2. Choose the 'optimal' predictor on which to split and choose the 'optimal' threshold value for splitting by applying a **splitting criterion**
3. Recurse on on each new node until **stopping condition** is met

For classification, we label each region in the model with the label of the class to which the plurality of the points within the region belong.

## Decision Trees for Regression

# Adaptations for Regression

---

With just two modifications, we can use a decision tree model for regression:

- ▶ The three splitting criteria we've examined each promoted splits that were pure - new regions increasingly specialized in a single class.

For classification, purity of the regions is a good indicator the performance of the model.

For regression, we want to select a splitting criterion that promotes splits that improves the predictive accuracy of the model as measured by, say, the MSE.

- ▶ For regression with output in  $\mathbb{R}$ , we want to label each region in the model with a real number - typically the average of the output values of the training points contained in the region.

# Learning Regression Trees

The learning algorithms for decision trees in regression tasks is:

1. Start with an empty decision tree (undivided feature space)
2. Choose a predictor  $j$  on which to split and choose a threshold value  $t_j$  for splitting such that the weighted average MSE of the new regions as smallest possible:

$$\operatorname{argmin}_{j,t_j} \left\{ \frac{N_1}{N} \operatorname{MSE}(R_1) + \frac{N_2}{N} \operatorname{MSE}(R_2) \right\}$$

or equivalently,

$$\operatorname{argmin}_{j,t_j} \left\{ \frac{N_1}{N} \operatorname{Var}[y|x \in R_1] + \frac{N_2}{N} \operatorname{Var}[y|x \in R_2] \right\}$$

where  $N_i$  is the number of training points in  $R_i$  and  $N$  is the number of points in  $R$ .

3. Recurse on on each new node until **stopping condition** is met

# Regression Trees Prediction

---

For any data point  $x_i$

1. Traverse the tree until we reach a leaf node.
2. Averaged value of the response variable  $y$ 's in the leaf (this is from the training set) is the  $\hat{y}_i$

## Stopping Conditions

---

Most of the stopping conditions, like maximum depth or minimum number of points in region, we saw last time can still be applied.

In the place of purity gain, we can instead compute accuracy gain for splitting a region  $R$

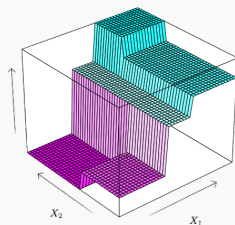
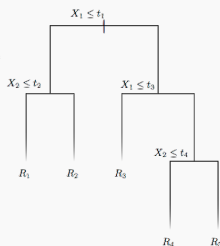
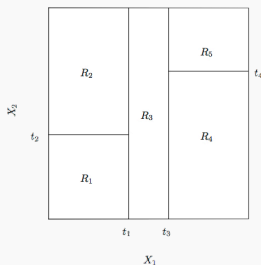
$$\text{Gain}(R) = \Delta(R) = MSE(R) - \frac{N_1}{N} MSE(R_1) - \frac{N_2}{N} MSE(R_2)$$

and stop the tree when the gain is less than some pre-defined threshold.

# Expressiveness of Decision Trees

We've seen that classification trees approximate boundaries in the feature space that separate classes.

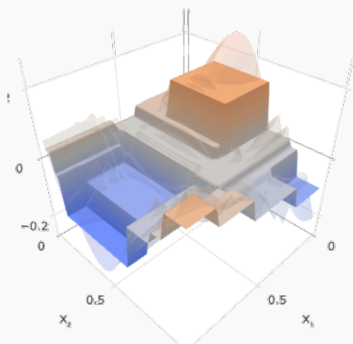
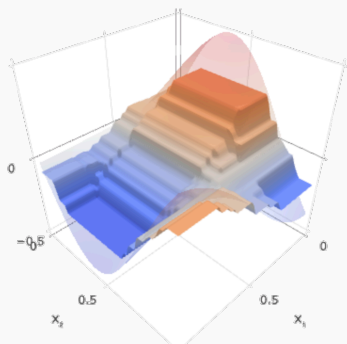
Regression trees, on the other hand, define **simple functions** or step functions, functions that are defined on partitions of the feature space and are constant over each part.





# Expressiveness of Decision Trees

For a fine enough partition of the feature space, these functions can approximate complex non-linear functions.



# Bagging

---

## Limitations of Decision Tree Models

---

Decision trees models are highly interpretable and fast to train, using our greedy learning algorithm.

However, in order to capture a complex decision boundary (or to approximate a complex function), we need to use a large tree (since each time we can only make axis aligned splits).

We've seen that large trees have high variance and are prone to overfitting.

For these reasons, in practice, decision tree models often underperforms when compared with other classification or regression methods.

# Bagging

---

One way to adjust for the high variance of the output of an experiment is to perform the experiment multiple times and then average the results.

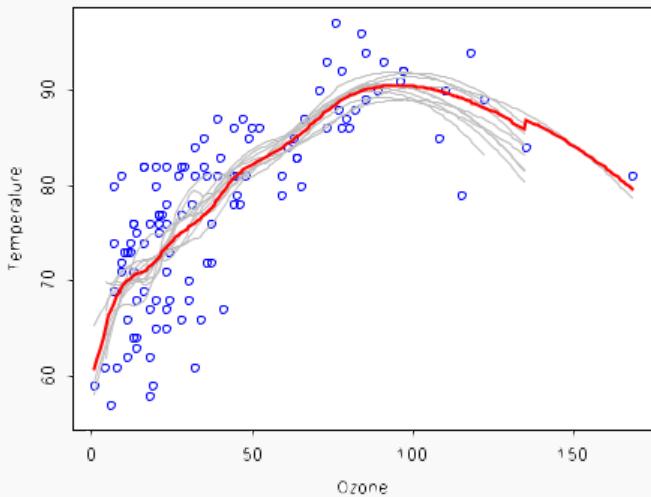
The same idea can be applied to high variance models:

1. **(Bootstrap)** we generate multiple samples of training data, via bootstrapping. We train a full decision tree on each sample of data.
2. **(Aggregate)** for a given input, we output the averaged outputs of all the models for that input.

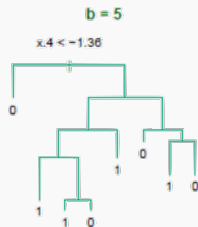
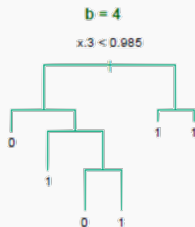
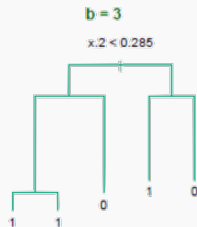
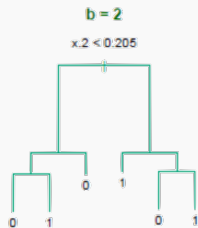
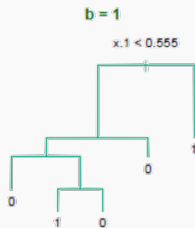
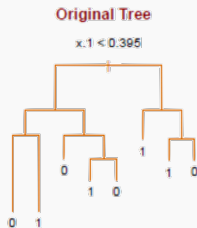
For classification, we return the class that is outputted by the plurality of the models.

This method is called **Bagging** (Breiman, 1996), short for, of course, Bootstrap Aggregating.

# Bagging



## Bagging



However, the major drawback of bagging (and other **ensemble methods** that we will study) is that the averaged model is no longer easily interpretable - i.e. one can no longer trace the 'logic' of an output through a series of decisions based on predictor values!

## Random Forests

---



**Random Forest** is a modified form of bagging that creates ensembles of independent decision trees.

To de-correlate the trees, we:

1. train each tree on a separate bootstrap sample of the full training set (same as in bagging)
2. for each tree, at each split, we **randomly** select a set of  $J'$  predictors from the full set of predictors.

From amongst the  $J'$  predictors, we select the optimal predictor and the optimal corresponding threshold for the split.

# Final Thoughts on Random Forests

---

- ▶ When the number of predictors is large, but the number of relevant predictors is small, random forests can perform poorly.

In each split, the chances of selected a relevant predictor will be low and hence most trees in the ensemble will be weak models.

# Final Thoughts on Random Forests

---

- Increasing the number of trees in the ensemble generally does not increase the risk of overfitting.

Again, by decomposing the generalization error in terms of bias and variance, we see that increasing the number of trees produces a model that is at least as robust as a single tree.

However, if the number of trees is too large, then the trees in the ensemble may become more correlated, increase the variance.