

CS 148 -

# Data Science Fundamentals

## Bagging and Random Forest

UCLA Computer Science

# Outline

---

- Review of Decision Trees
- Bagging
- Out of Bag Error (OOB)
- Variable Importance
- Random Forests

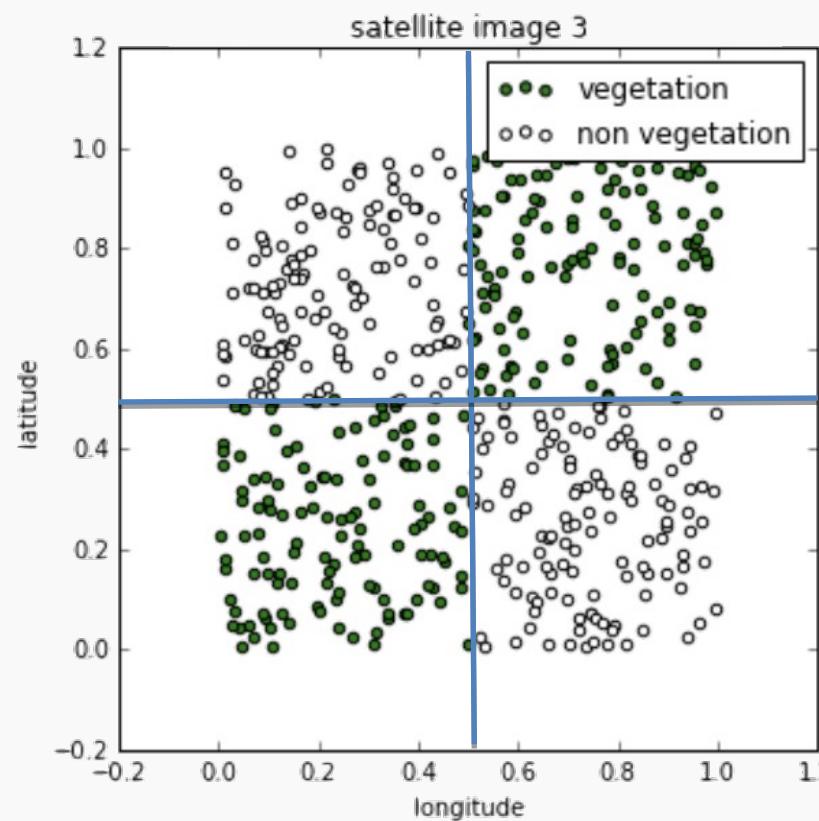
# Outline

---

- **Review of Decision Trees**
- Bagging
- Out of Bag Error (OOB)
- Variable Importance
- Random Forests

# Geometry of Data

Complicate decision boundaries can not be explained with LogRegression.



# Decision Trees

---

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 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.

For regression, we predict with the average of the output

# Splitting Criteria

---

The splitting criteria we've examined each minimize a loss function

- A. For classification**, purity of the regions is a good indicator the performance of the model. Entropy as splitting criteria minimizes the cross-entropy (greedy).
- B. For regression**, we want to select a splitting criterion that promotes splits that improves the predictive accuracy of the model as measured by the MSE.

# Stopping Conditions

---

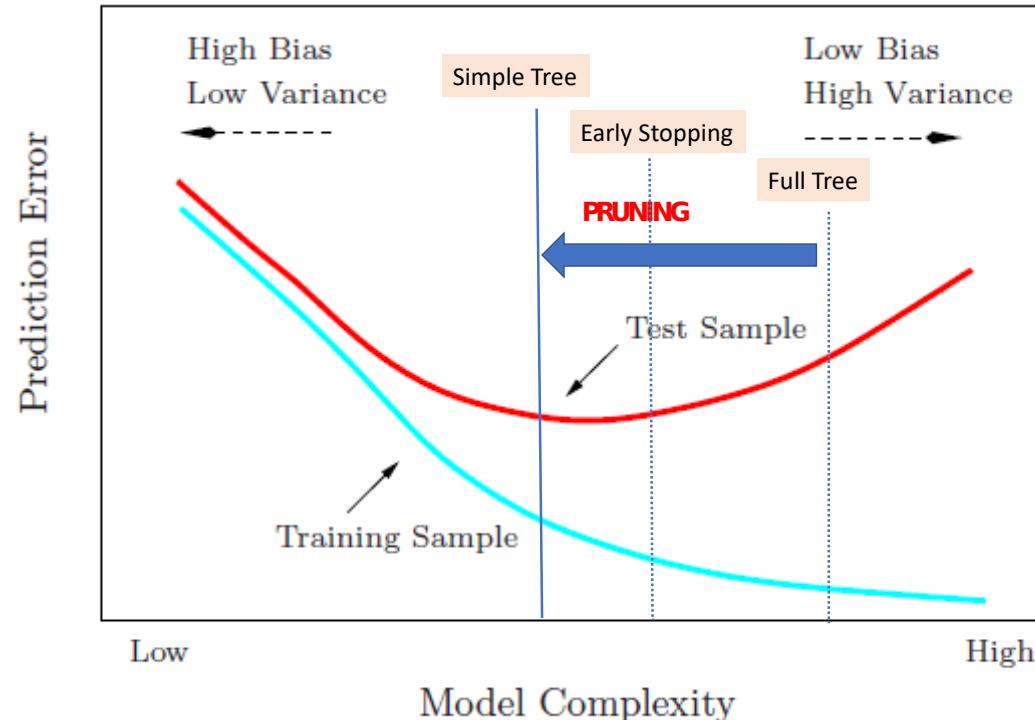
Common simple stopping conditions:

- 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 (**`min_samples_leaf`**).
- 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.

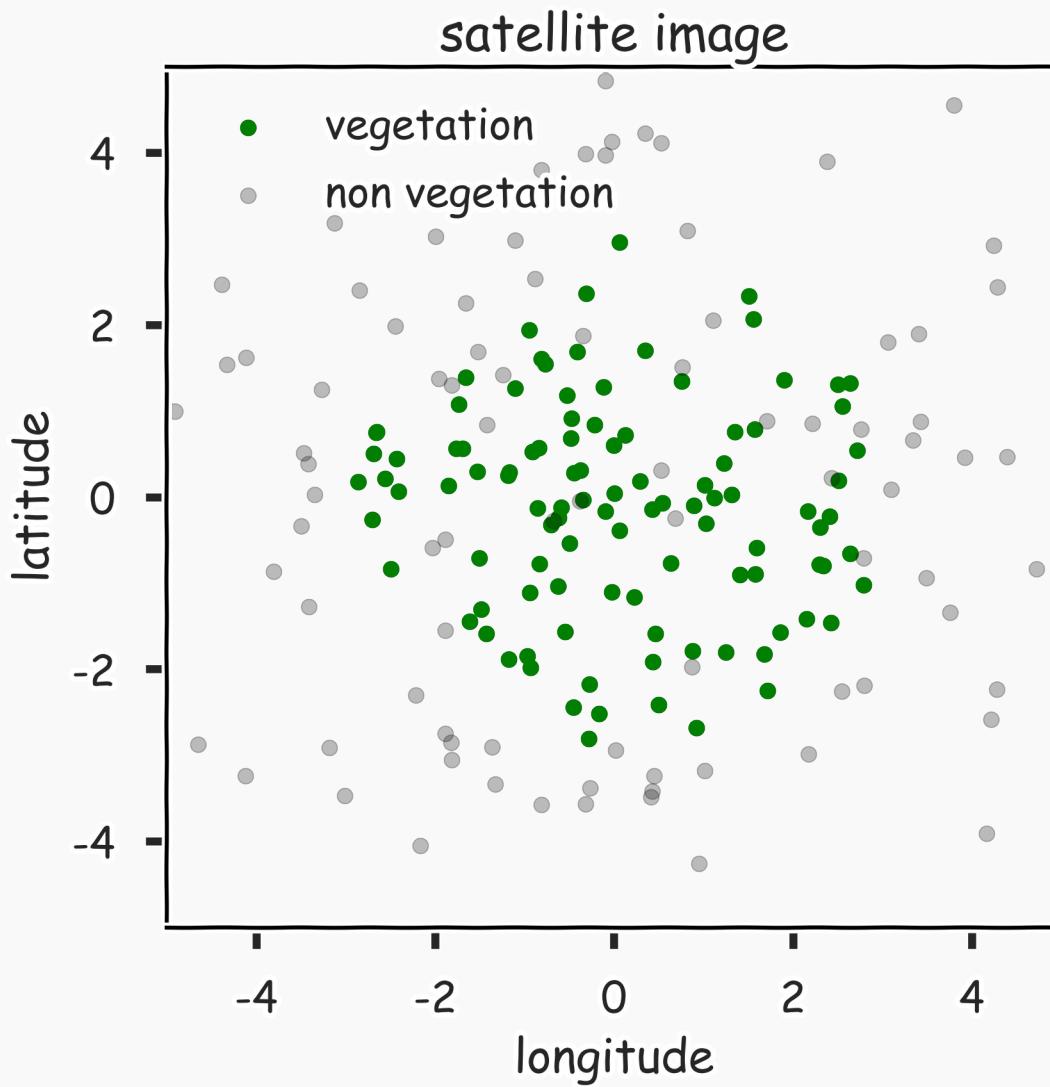
# Overfitting

Same issues as with classification trees. Avoid overfitting by pruning or limiting the depth of the tree and using CV.

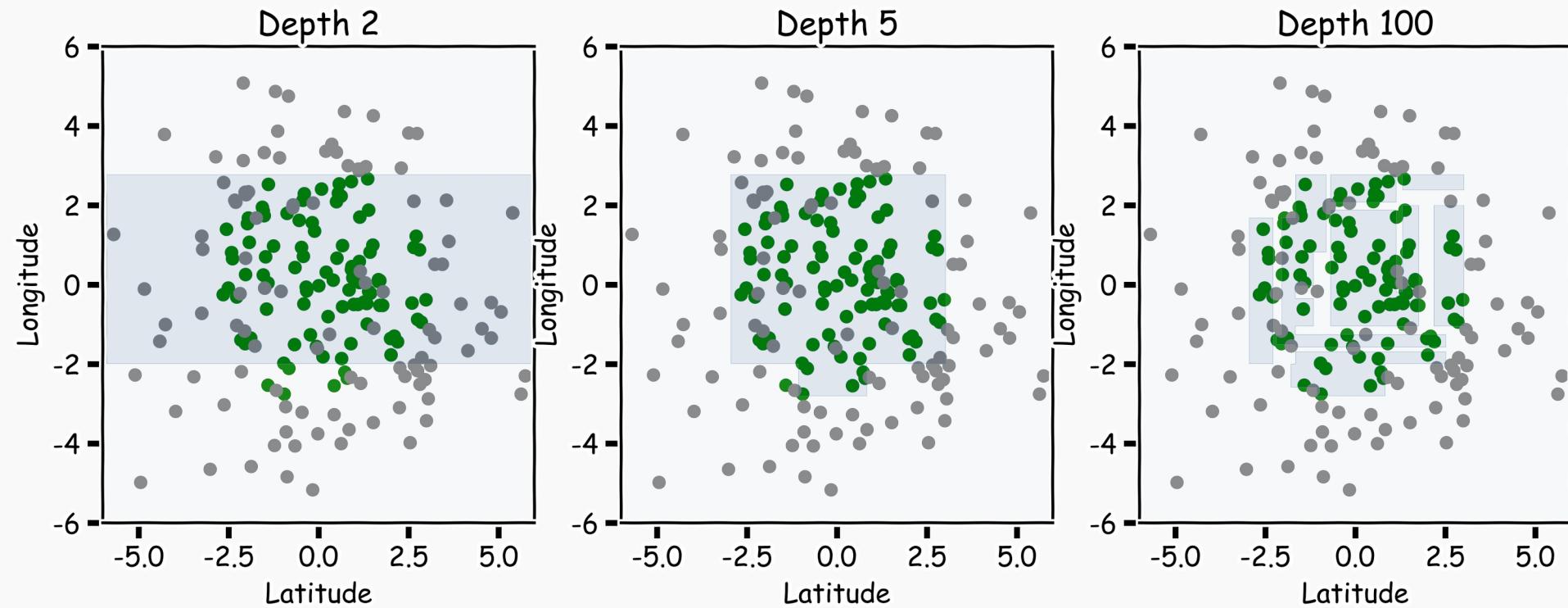


# Bagging

# Reduce the variance

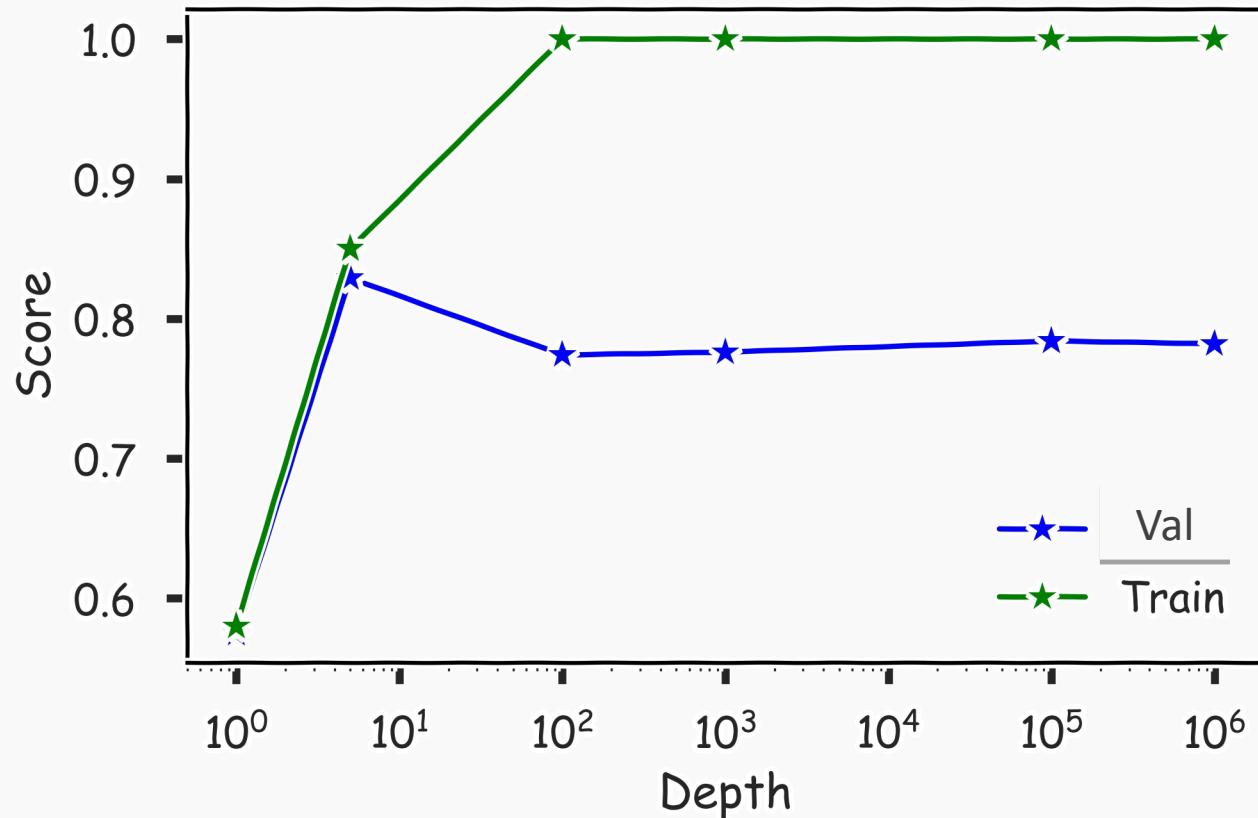


# Hyper-parameters: Depth

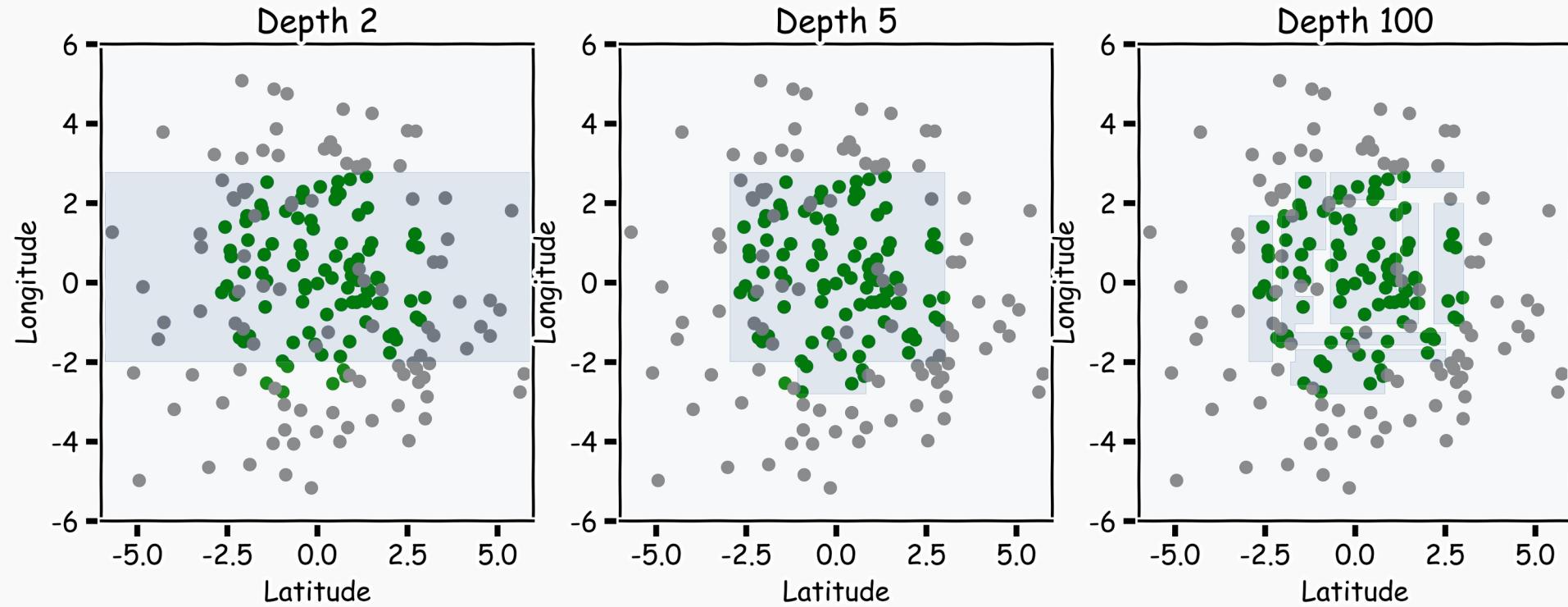


# Hyper-parameters: Depth

Use train/validation or cross validation to estimate the best depth.



# Magic realism: Bootstrap



# Limitations of Decision Tree Models

---

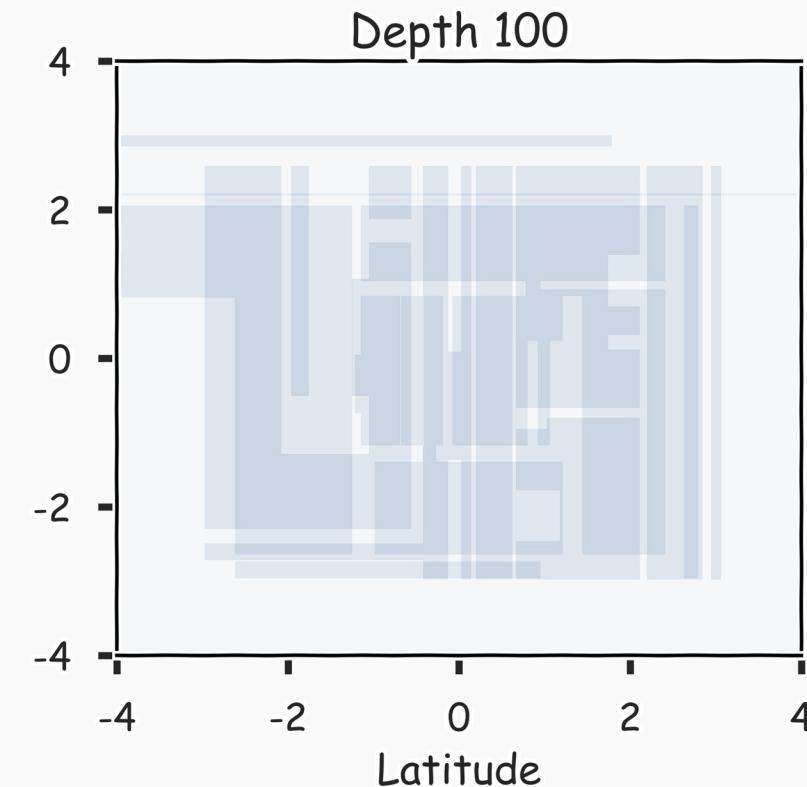
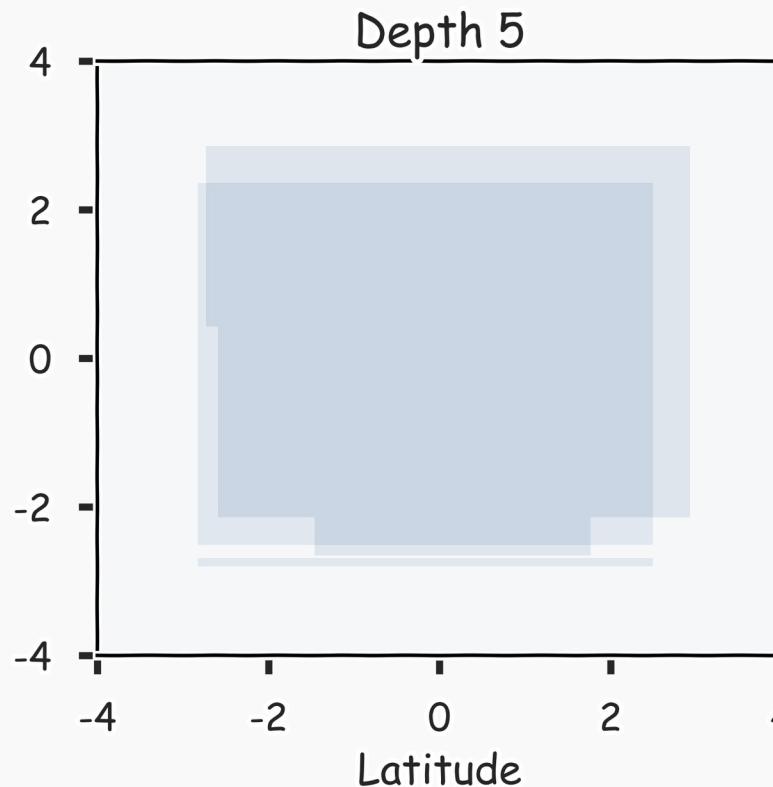
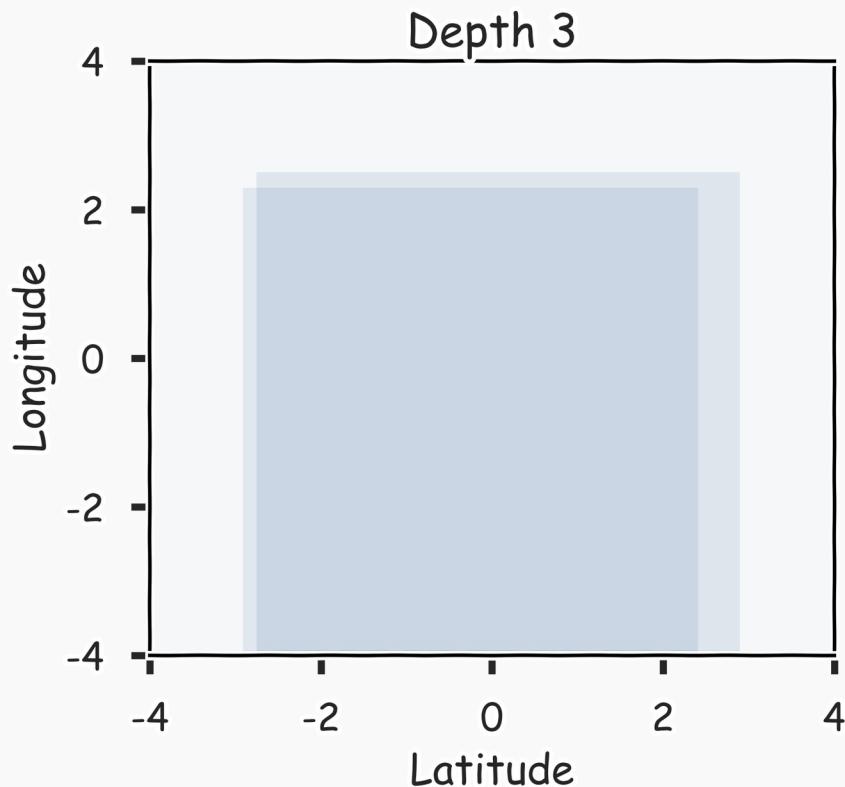
Decision trees models are highly interpretable and fast to train, using our greedy learning algorithm.

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.

# Combine them? 2 magic realisms



# 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. For regression we return the average of the outputs for each tree.

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

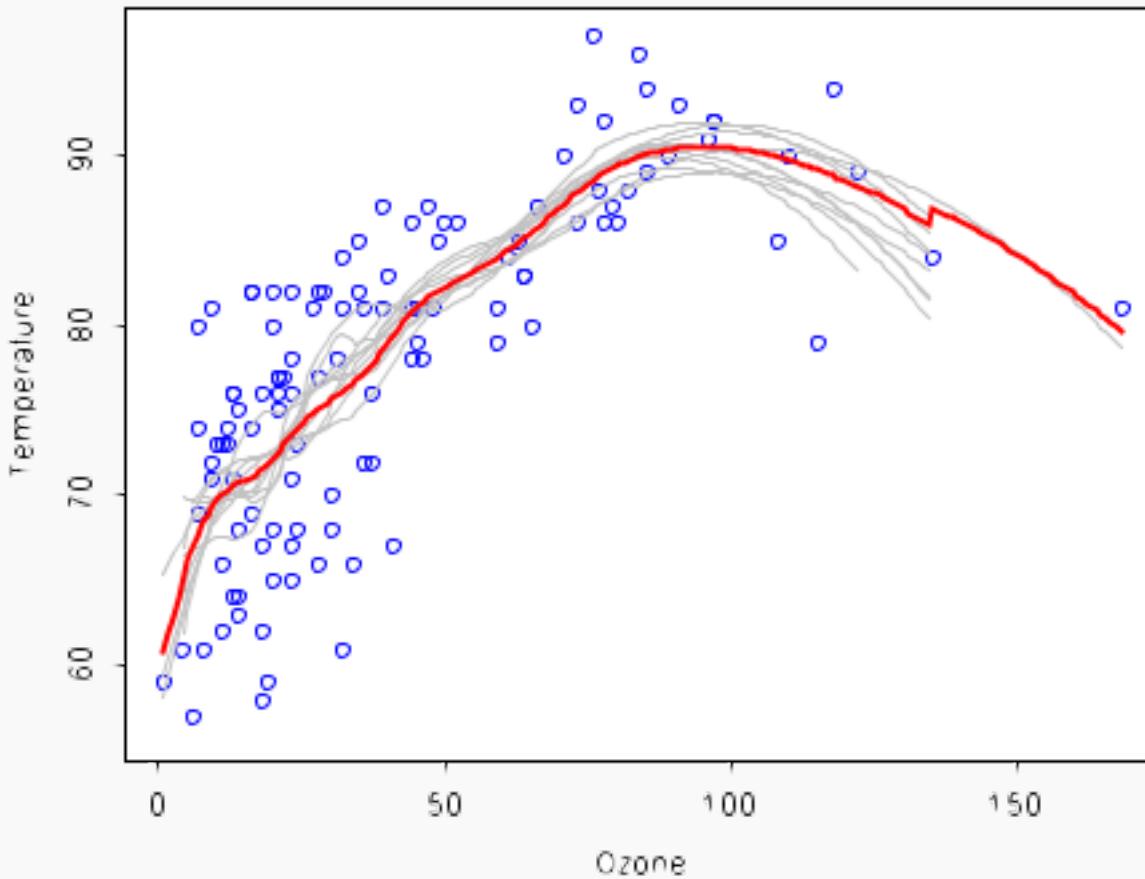
# Bagging

---

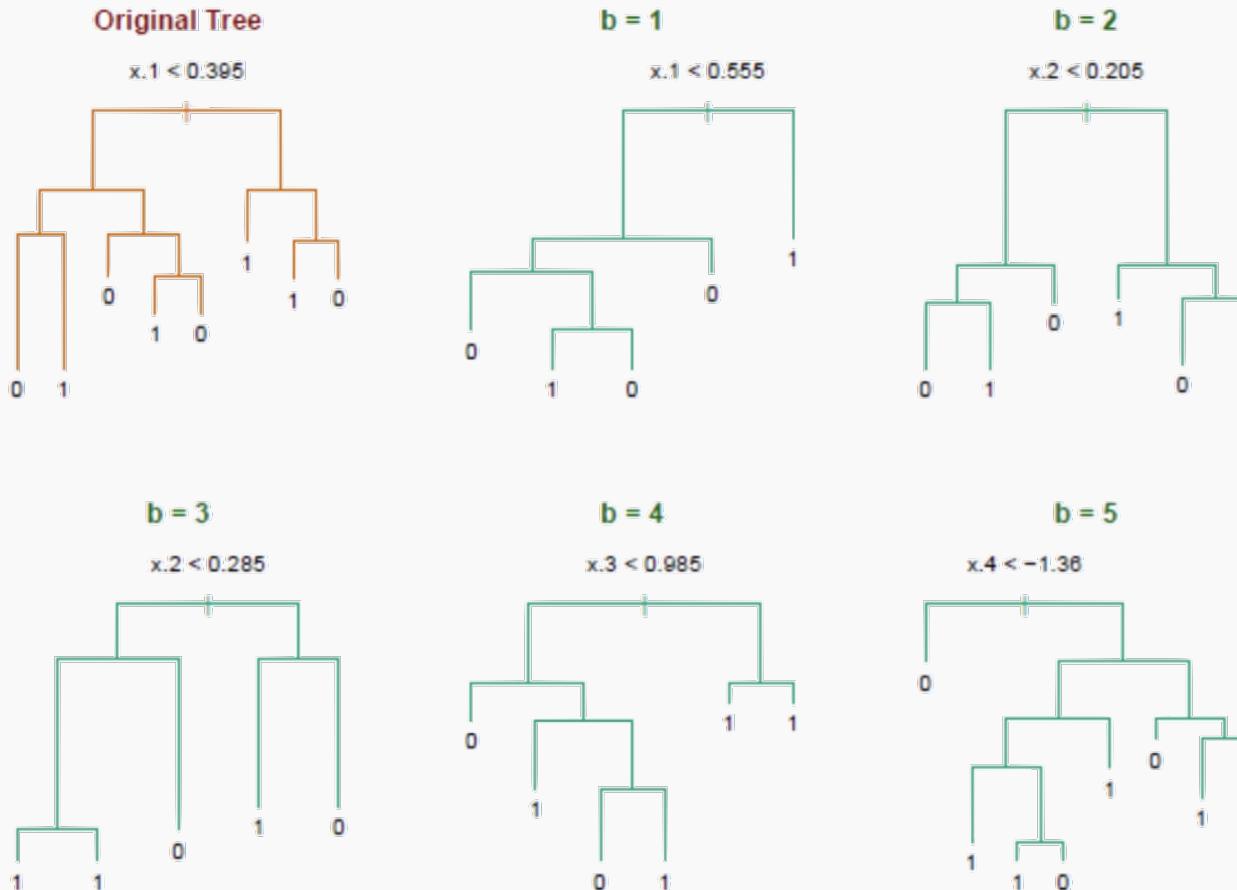
Note that bagging enjoys the benefits of:

1. High expressiveness - by using full trees each model is able to approximate complex functions and decision boundaries.
2. Low variance - averaging the prediction of all the models reduces the variance in the final prediction, assuming that we choose a sufficiently large number of trees.

# Bagging (regression)



# Bagging (classification)



# Bagging

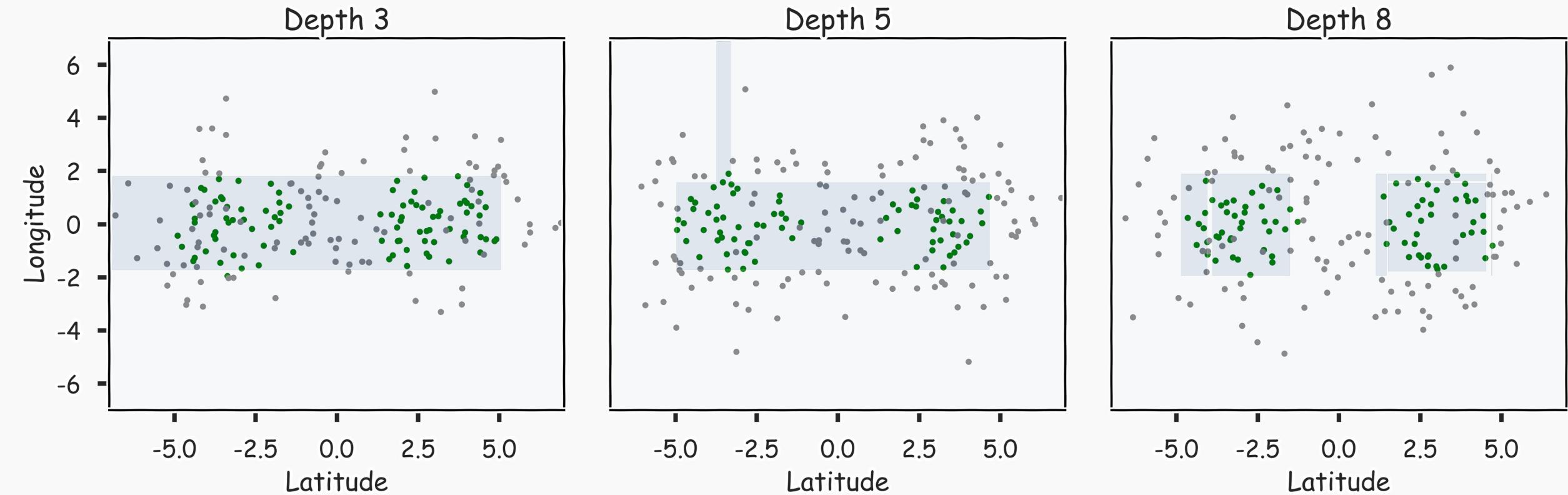
---

**Question:** Do you see any problems?

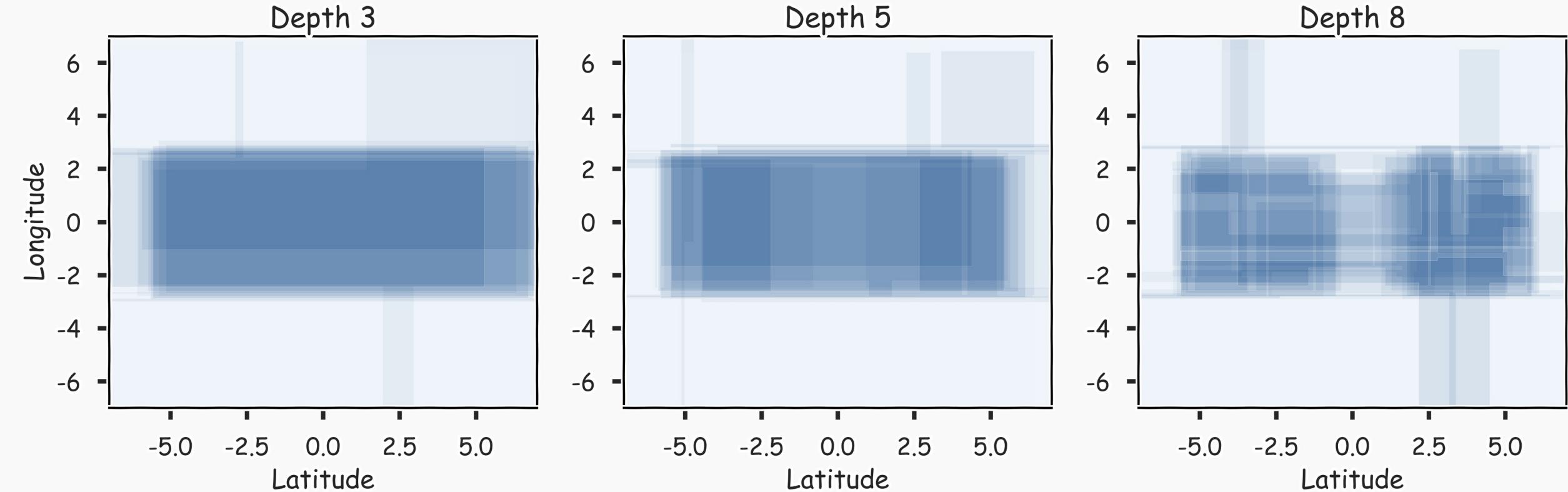
- Still some overfitting if the trees are too large.
- If trees are too shallow it can still underfits.
- Interpretability:

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!

# Case of underfitting



# Case of underfitting



# Bagging

---

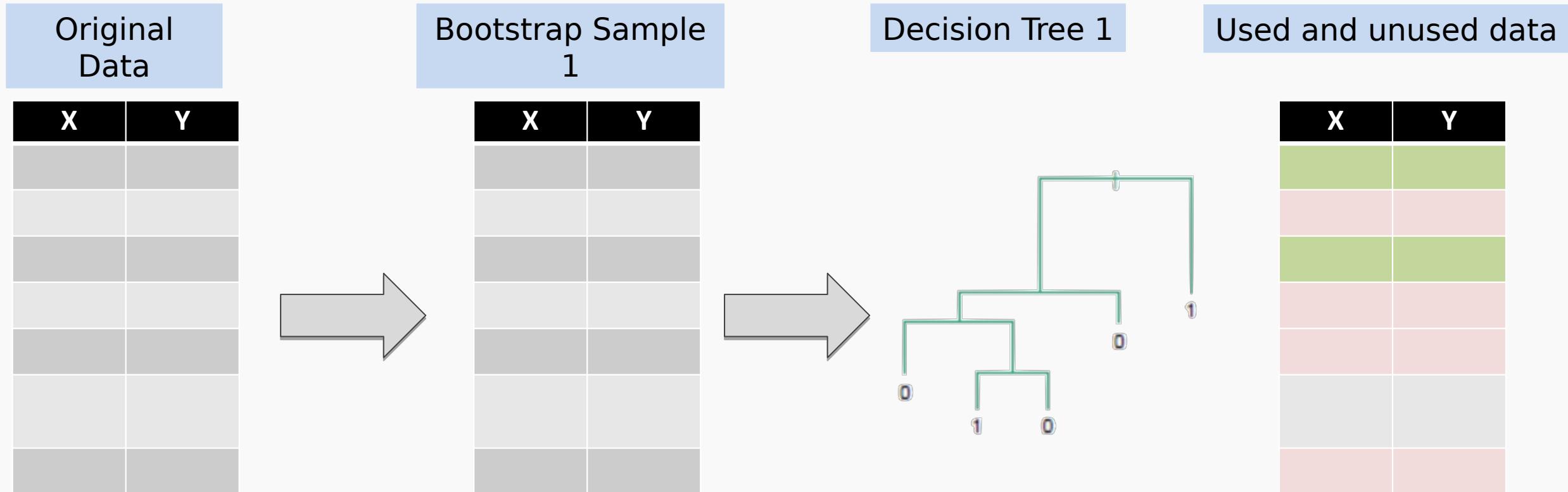
**Question:** Do you see any problems?

- Still some overfitting if the trees are too large
- If trees are too shallow it can still underfits.

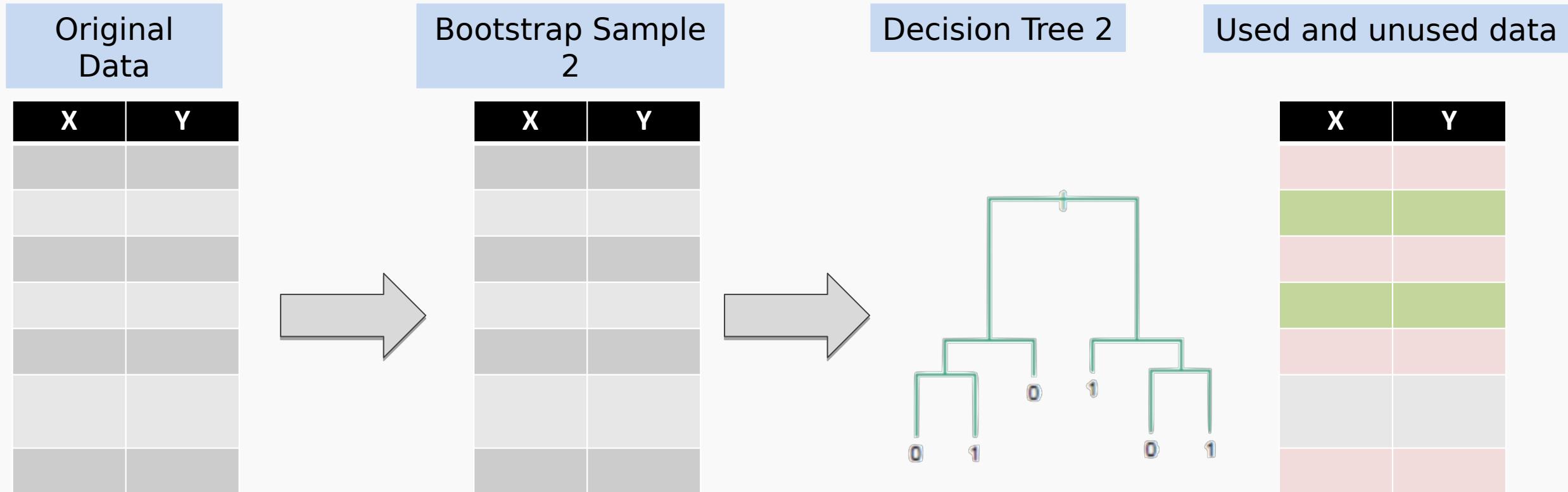
## Cross Validations

# Out-of-Bag Error

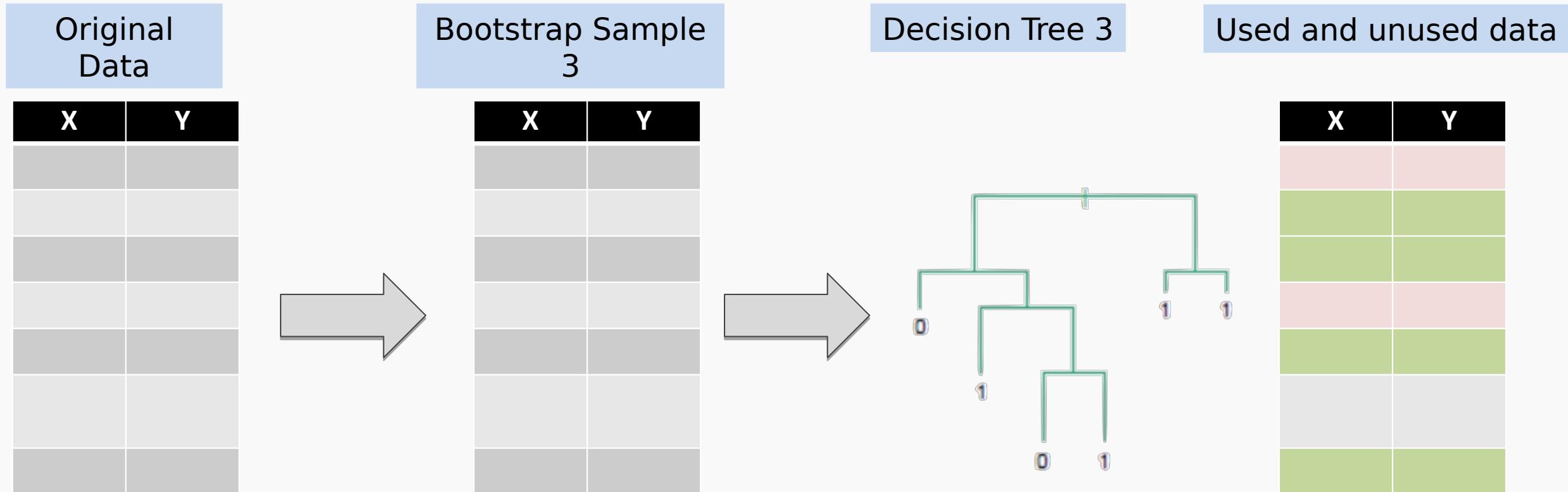
# Bagging



# Bagging



# Bagging



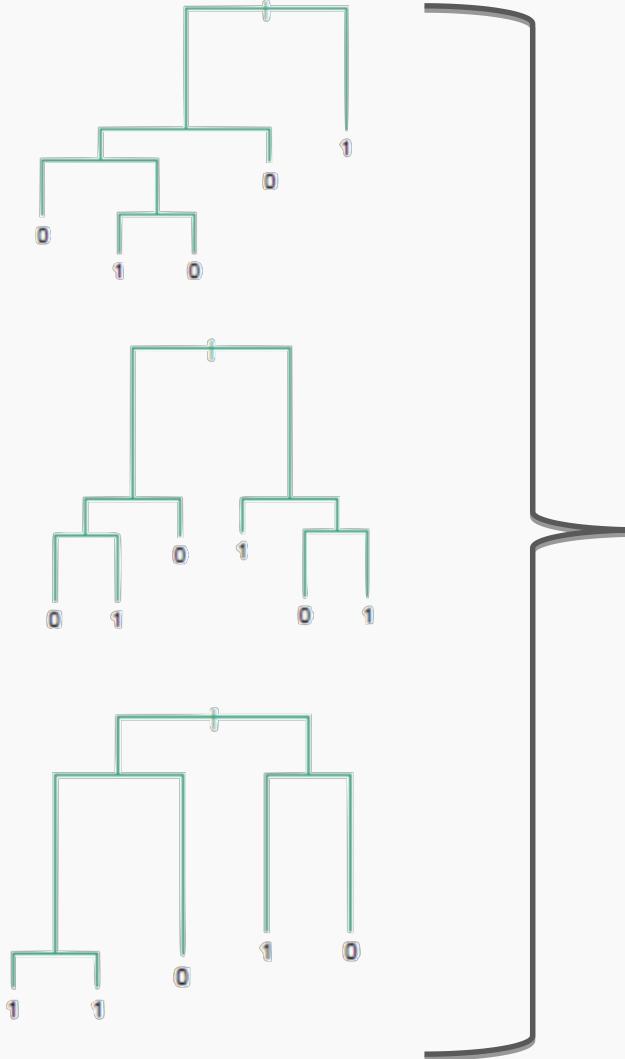
# Point-wise out-of-bag error

X	Y

# Point-wise out-of-bag error

X	Y

B Trees that did not see



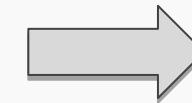
Classification

$$\hat{y}_{i,pw} = \text{majority } i$$

$$e_i = \mathbb{I}(\hat{y}_{i,pw} \neq y_i)$$

Regression

$$\hat{y}_{i,pw} = \sum_{j \in B} \hat{y}_{i,j}$$



$$e_i = (y_i - \hat{y}_{i,pw})^2$$

# OOB Error

We average the point-wise out-of-bag error over the full training set.

Classification

$$Error_{OOB} = \sum_i^n e_i = \sum_i^n \mathbb{I}(\hat{y}_{i,pw} \neq y_i)$$

Regression

$$Error_{OOB} = \sum_i^n e_i = \sum_i^n (y_i - \hat{y}_{i,pw})^2$$

# Out-of-Bag Error

---

Bagging is an example of an ***ensemble method***, a method of building a single model by training and aggregating multiple models.

With ensemble methods, we get a new metric for assessing the predictive performance of the model, the ***out-of-bag error***.

Given a training set and an ensemble of models, each trained on a bootstrap sample, we compute the ***out-of-bag error*** of the averaged model by

1. For each point in the training set, we average the predicted output for this point over the models whose bootstrap training set excludes this point. We compute the error or squared error of this averaged prediction. Call this the point-wise out-of-bag error.
2. We average the point-wise out-of-bag error over the full training set.

# Bagging

---

**Question:** Do you see any problems?

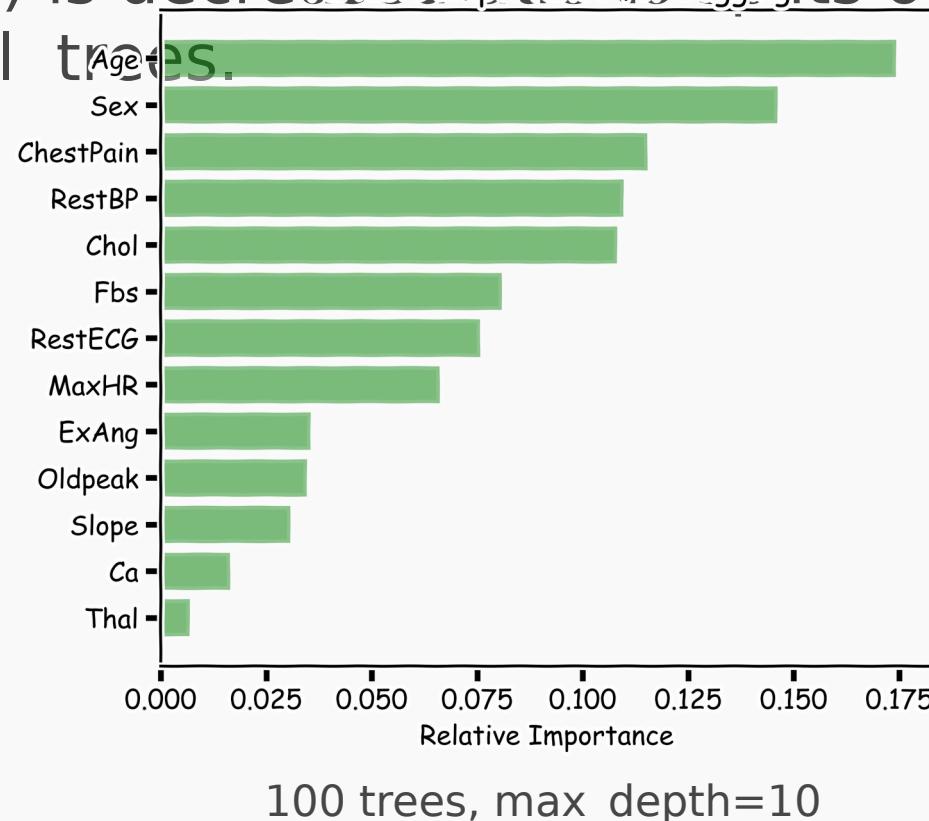
- Still some overfitting if the trees are too large.
- If trees are too shallow it can still underfits.
- **Interpretability:**

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!

# Variable Importance for Bagging

Bagging improves prediction accuracy at the expense of interpretability.

Calculate the total amount that the MSE (for regression) or Gini index (for classification) is decreased due to splits over a given predictor, averaged over all trees.



# Improving on Bagging

---

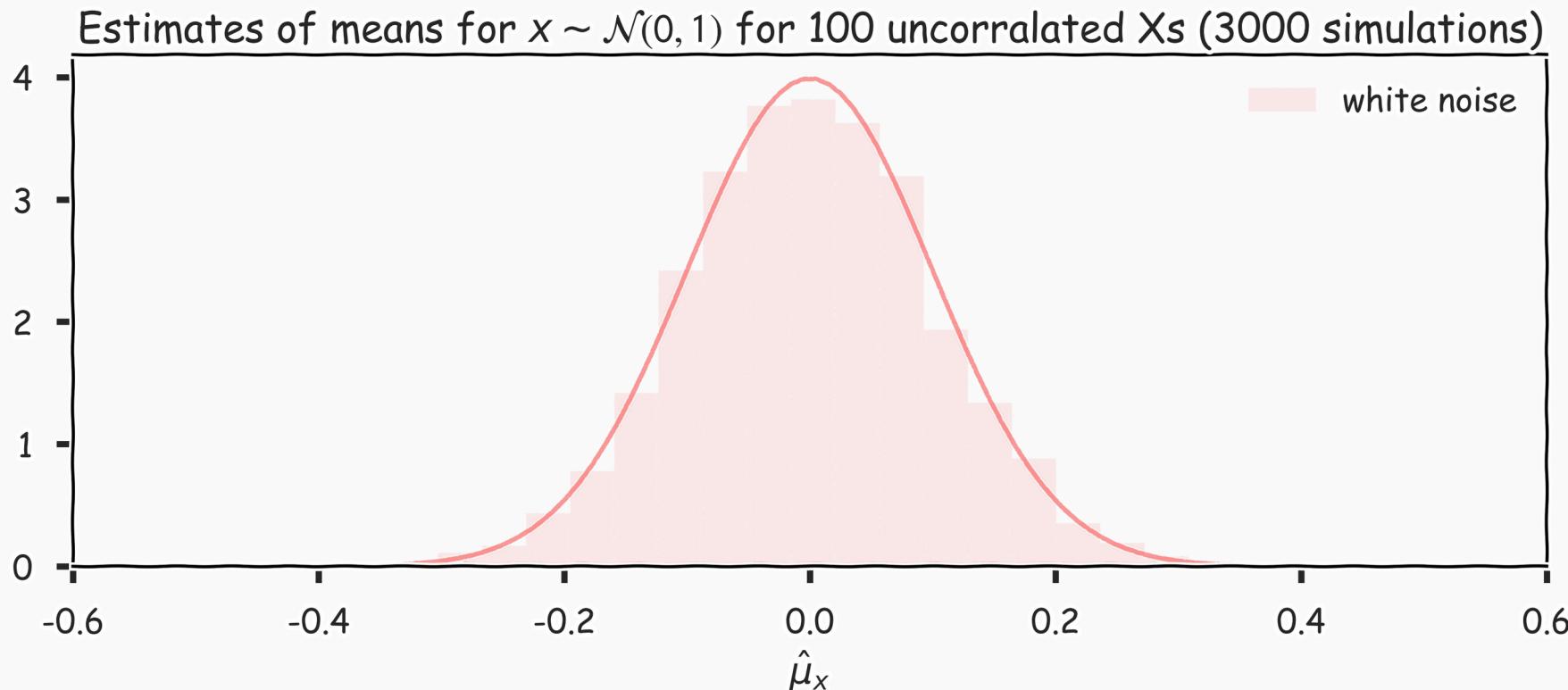
In practice, the ensembles of trees in Bagging tend to be highly correlated.

Suppose we have an extremely strong predictor, , in the training set amongst moderate predictors. Then the greedy learning algorithm ensures that most of the models in the ensemble will choose to split on in early iterations.

That is, each tree in the ensemble is identically distributed, with the expected output of the averaged model the same as the expected output of any one of the trees.

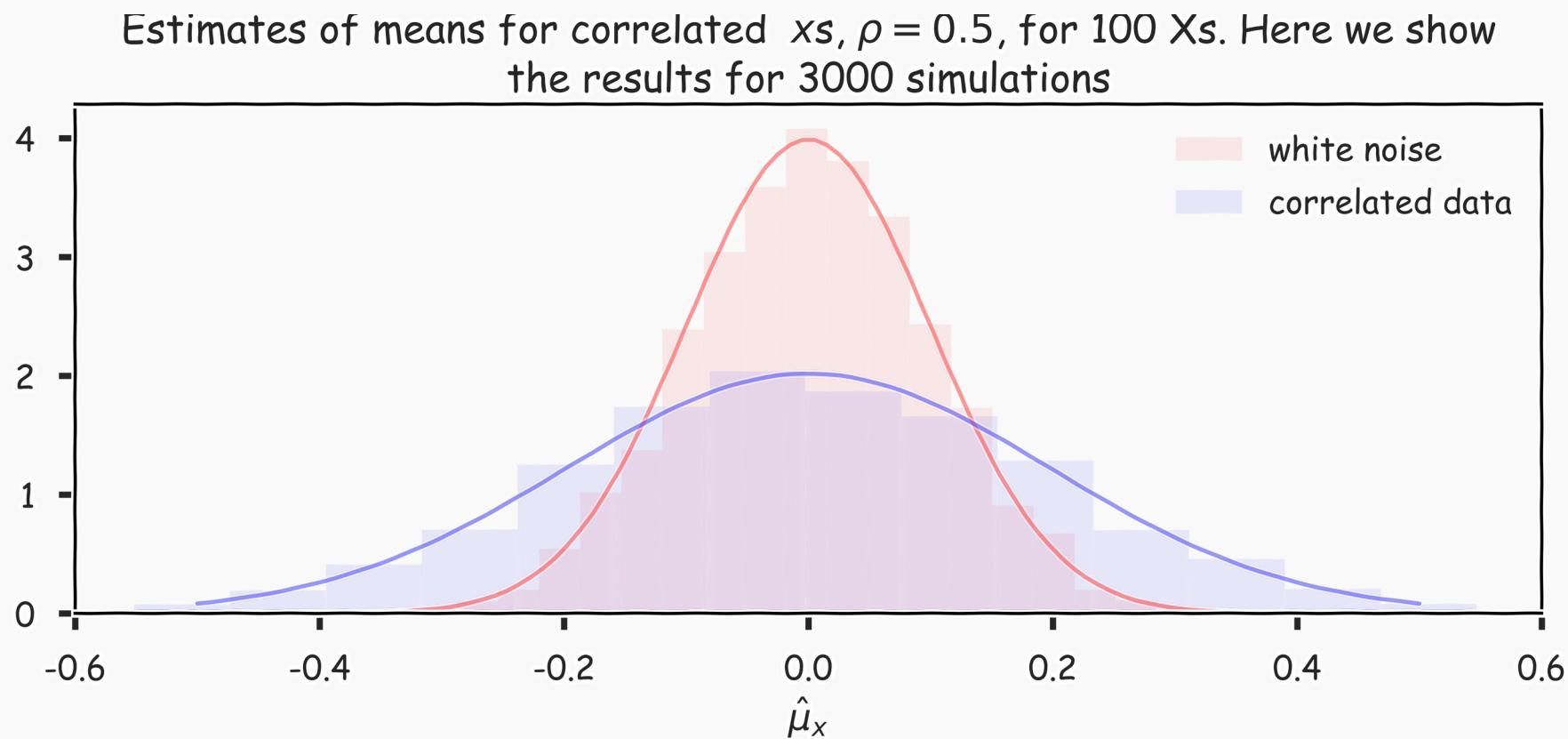
# Improving on Bagging

Recall, for  $n$  number of identically and independently distributed variable,  $X$ , with variance  $\sigma^2$ , the variance of the estimate of the mean is :



# Improving on Bagging

For number of identically but not independently distributed variables with pairwise correlation and variance , the variance of their mean is



# Bagging

---

**Question:** Do you see any problems?

- Still some overfitting if the trees are too large
- If trees are too shallow it can still underfits.
- interpretability
- 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 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 predictors from the full set of predictors.

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

# Tuning Random Forests

---

Random forest models have multiple hyper-parameters to tune:

1. the number of predictors to randomly select at each split
2. the total number of trees in the ensemble
3. the minimum leaf node size

In theory, each tree in the random forest is full, but in practice this can be computationally expensive (and added redundancies in the model), thus, imposing a minimum node size is not unusual.

# Tuning Random Forests

---

There are standard (default) values for each of random forest hyper-parameters recommended by long time practitioners, but generally these parameters should be tuned through **OOB** (making them data and problem dependent).

e.g. number of predictors to randomly select at each split:

- for classification
- for regression

Using out-of-bag errors, training and cross validation can be done in a single sequence - we cease training once the out-of-bag error stabilizes

# Variable Importance for RF

---

Same as with Bagging:

Calculate the total amount that the RSS (for regression) or Gini index (for classification) is decreased due to splits over a given predictor, averaged over all trees.

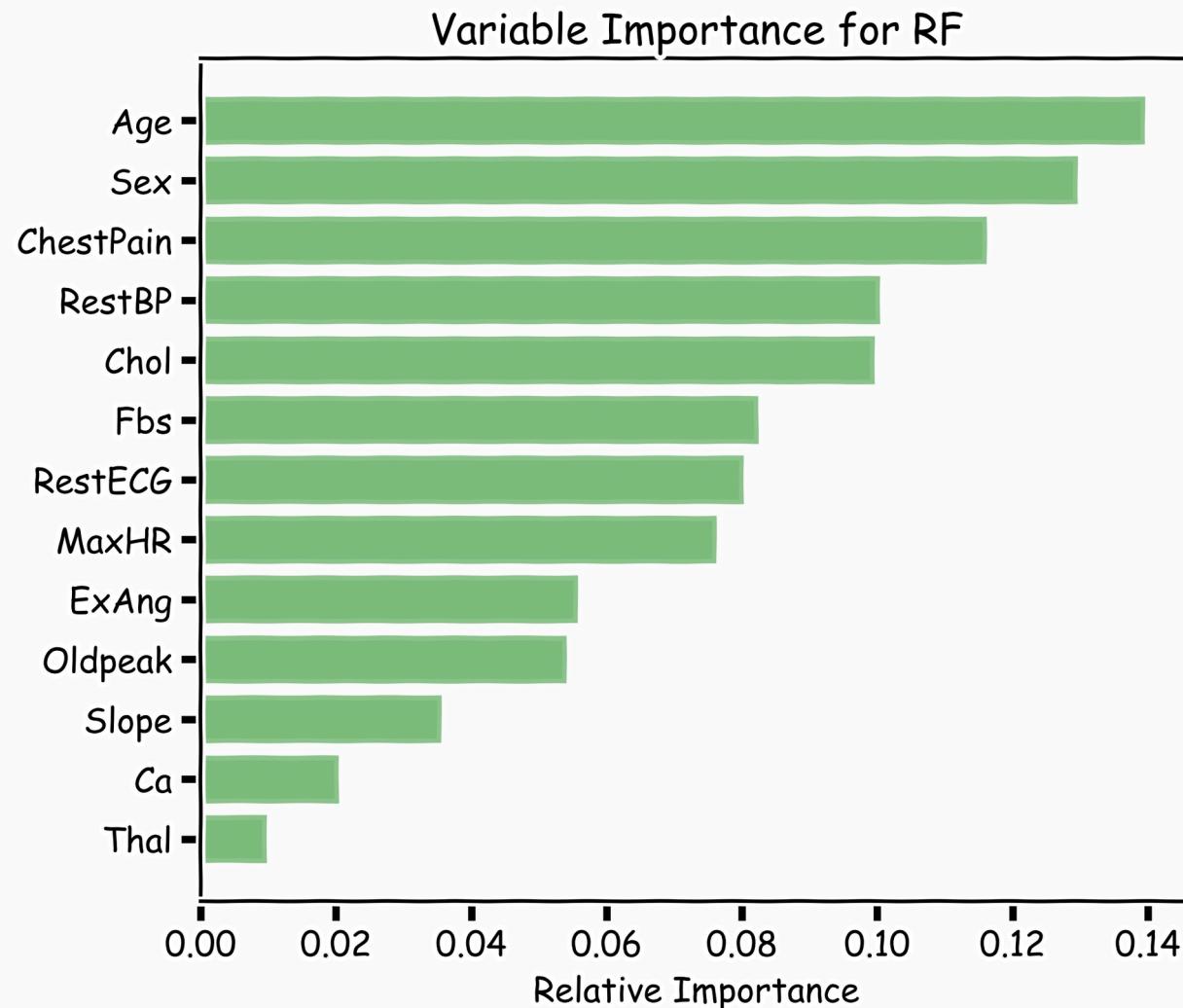
# Variable Importance for RF

---

## **Alternative:**

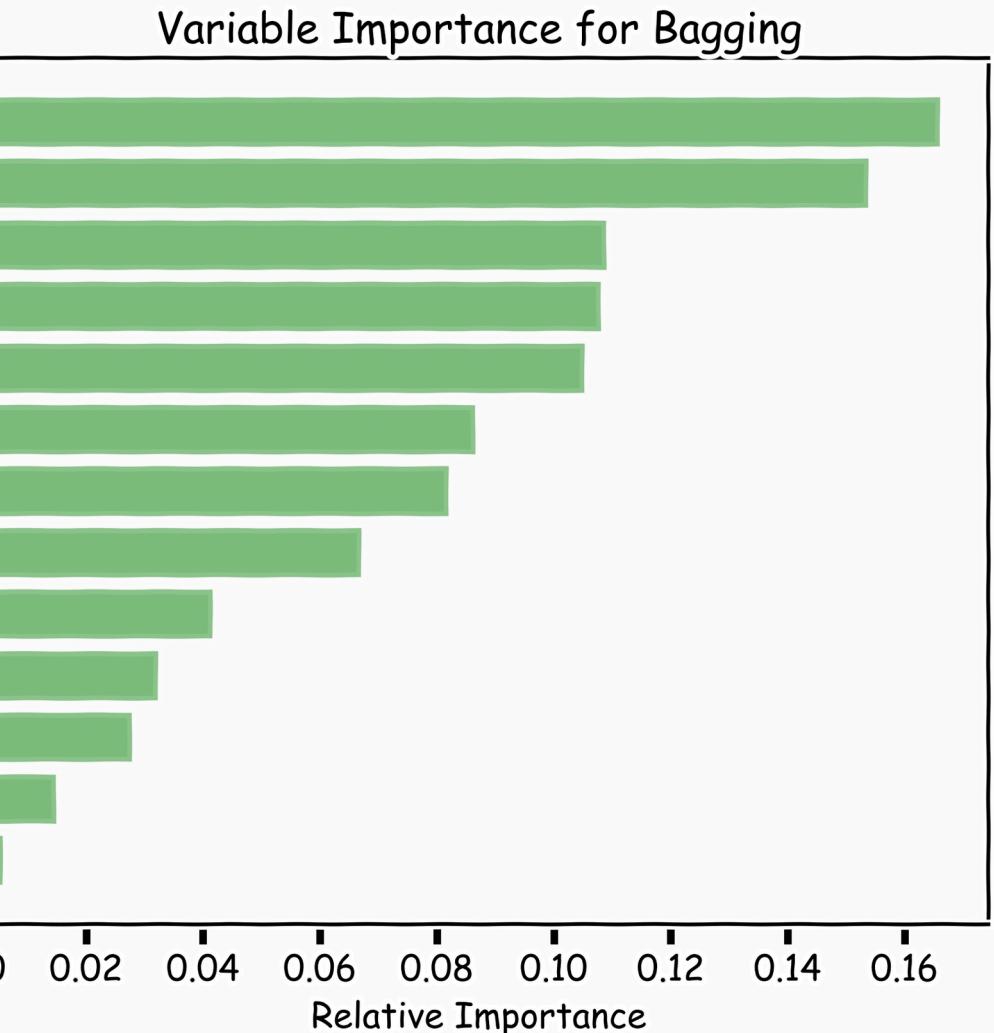
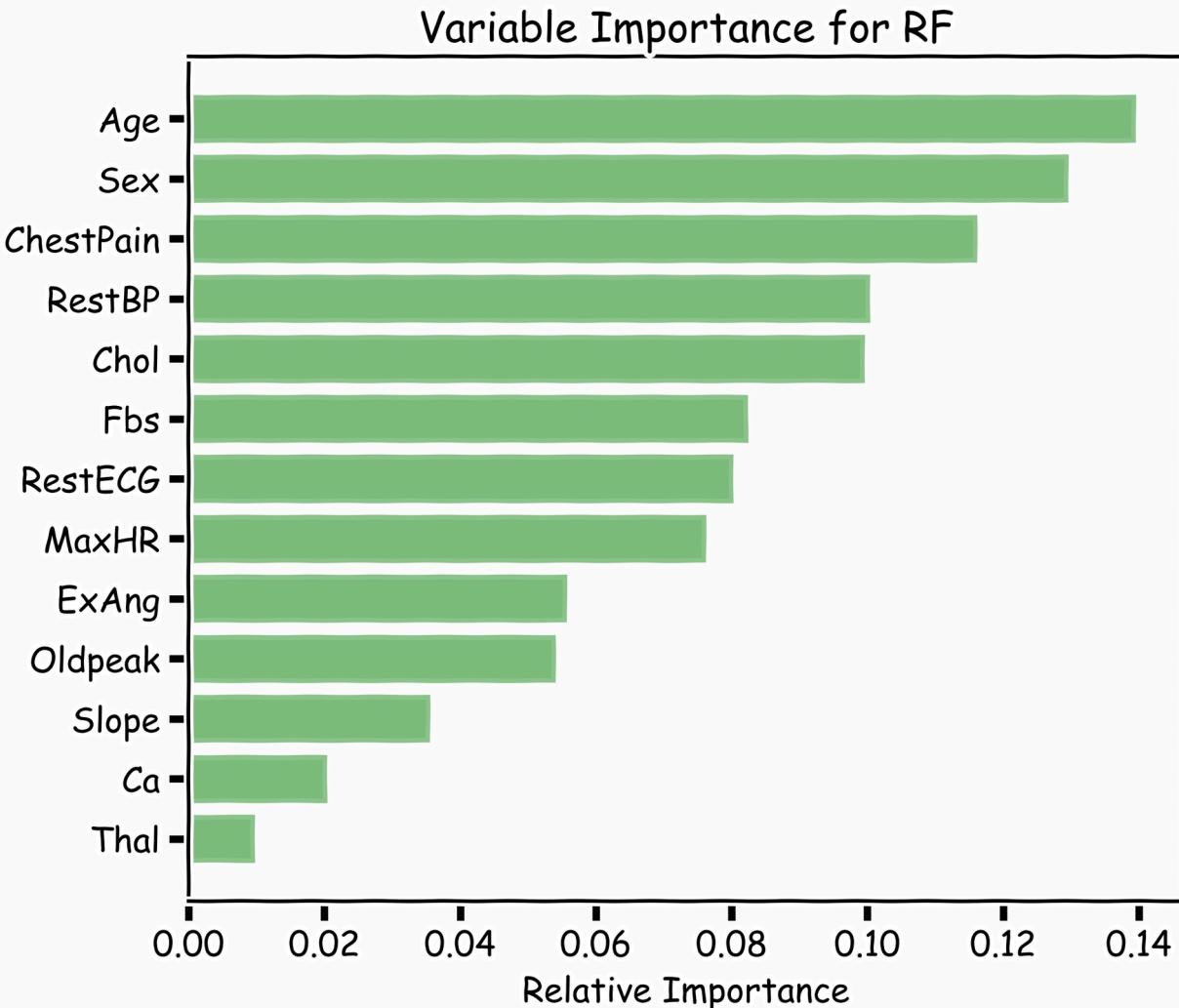
- Record the prediction accuracy on the *oob* samples for each tree.
- Randomly permute the data for column in the *oob* samples the record the accuracy again.
- The decrease in accuracy as a result of this permuting is averaged over all trees, and is used as a measure of the importance of variable in the random forest.

# Variable Importance for RF



100 trees, max\_depth=10

# Variable Importance for RF



100 trees, max\_depth=10

# 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.

**Question:** Why?

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

# Final Thoughts on Random Forests (cont.)

---

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.

# Final Thoughts on Random Forests (cont.)

---

## Probabilities:

- Random Forrest Classifier (and bagging) can return probabilities.
- **Question:** How?

# Next Lecture

---

- Unbalance dataset
- Weighted samples
- Categorical data
- Missing data
- Different implementations

**AND BOOSTING**