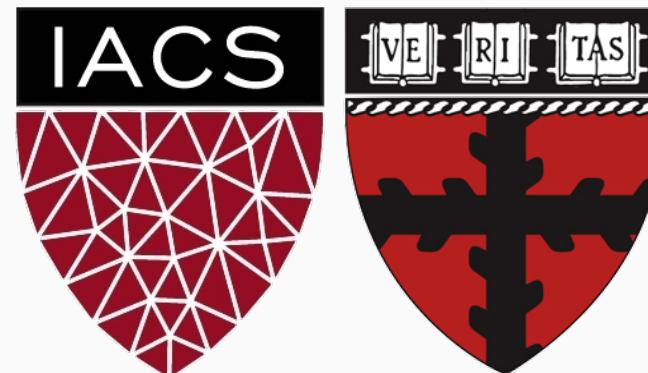


Lecture 16: Regression Trees, Bagging and Random Forest

CS109A Introduction to Data Science
Pavlos Protopapas and Kevin Rader



Outline

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



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 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:

1. The three splitting criteria we've examined each promoted splits that were pure - new regions increasingly specialized in a single class.
 - A. **For classification**, purity of the regions is a good indicator the performance of the model.
 - B. **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.
2. 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.

Adaptations for Regression

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

1. 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.
2. 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 features pace)
2. Choose a predictor j on which to split and choose a threshold value $t \downarrow j$ for splitting such that the weighted average MSE of the new regions as smallest possible:

$$\operatorname{argmin}_{\tau j, t \downarrow j} \square \{ N \downarrow 1 / N \operatorname{MSE}(R \downarrow 1) + N \downarrow 2 / N \operatorname{MSE}(R \downarrow 2) \}$$

or equivalently,

$$\operatorname{argmin}_{\tau j, t \downarrow j} \square \{ N \downarrow 1 / N \operatorname{Var}(y | x \in R \downarrow 1) + N \downarrow 2 / N \operatorname{Var}(y | x \in R \downarrow 2) \}$$

where $N \downarrow i$ is the number of training points in $R \downarrow i$ and N is the number of points in R .

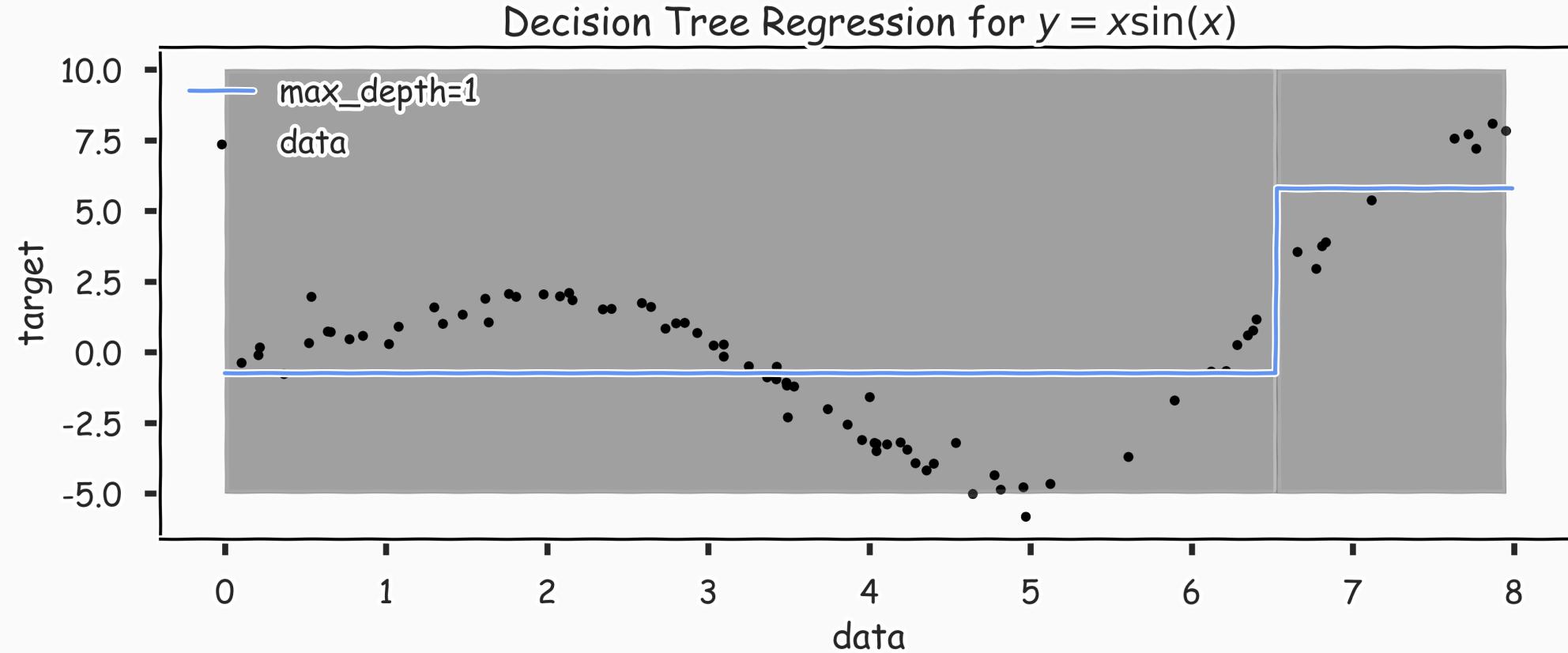
3. Recurse on each new node until ***stopping condition*** is met

Regression Trees Prediction

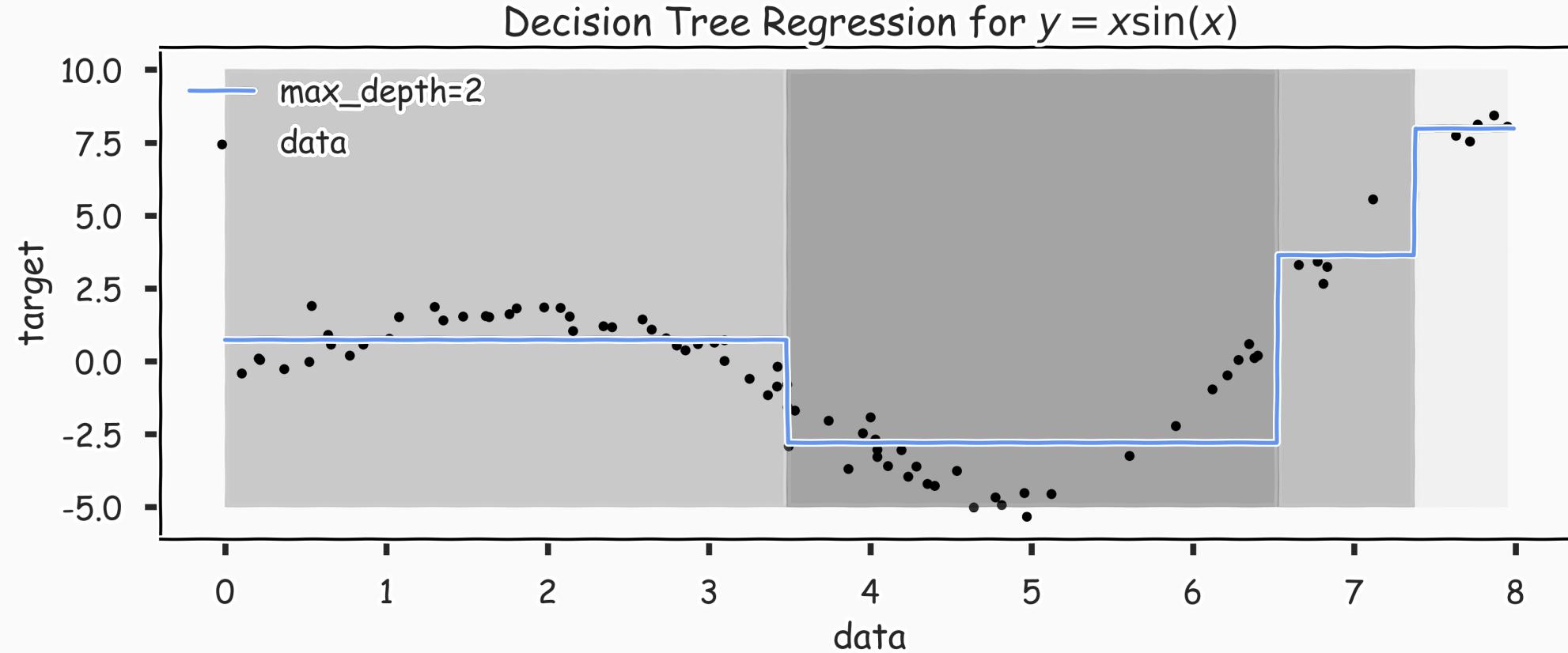
For any data point $x \downarrow i$

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 $y \downarrow i$.

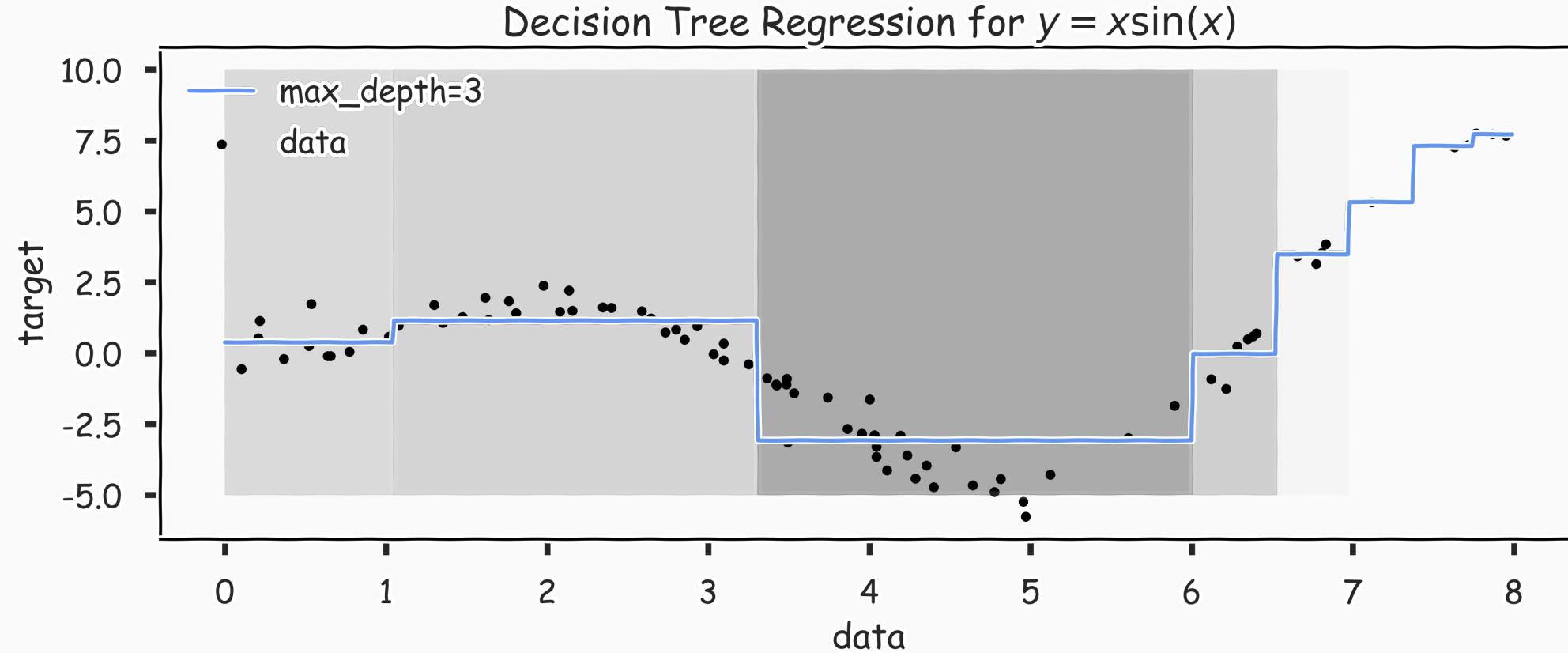
Regression Trees Prediction (grey scale represents MSE)



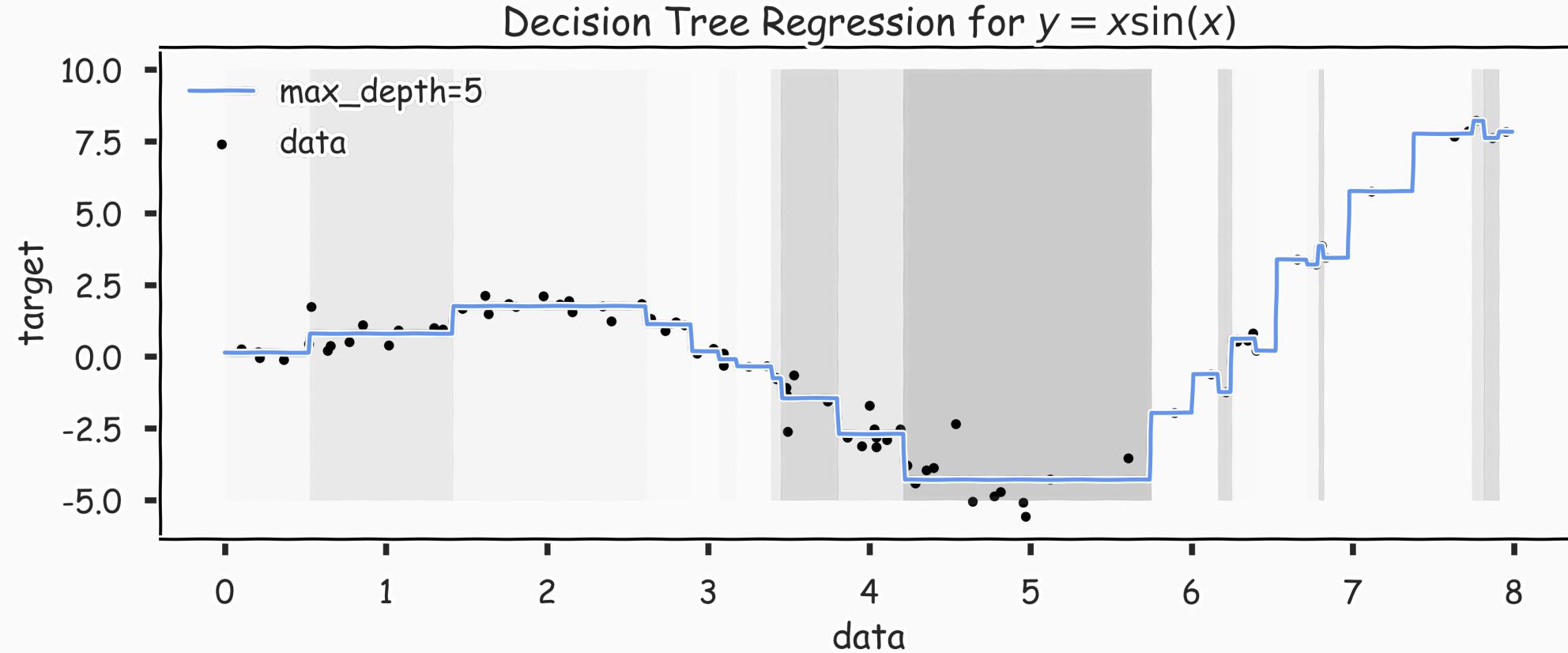
Regression Trees Prediction (grey scale represents MSE)



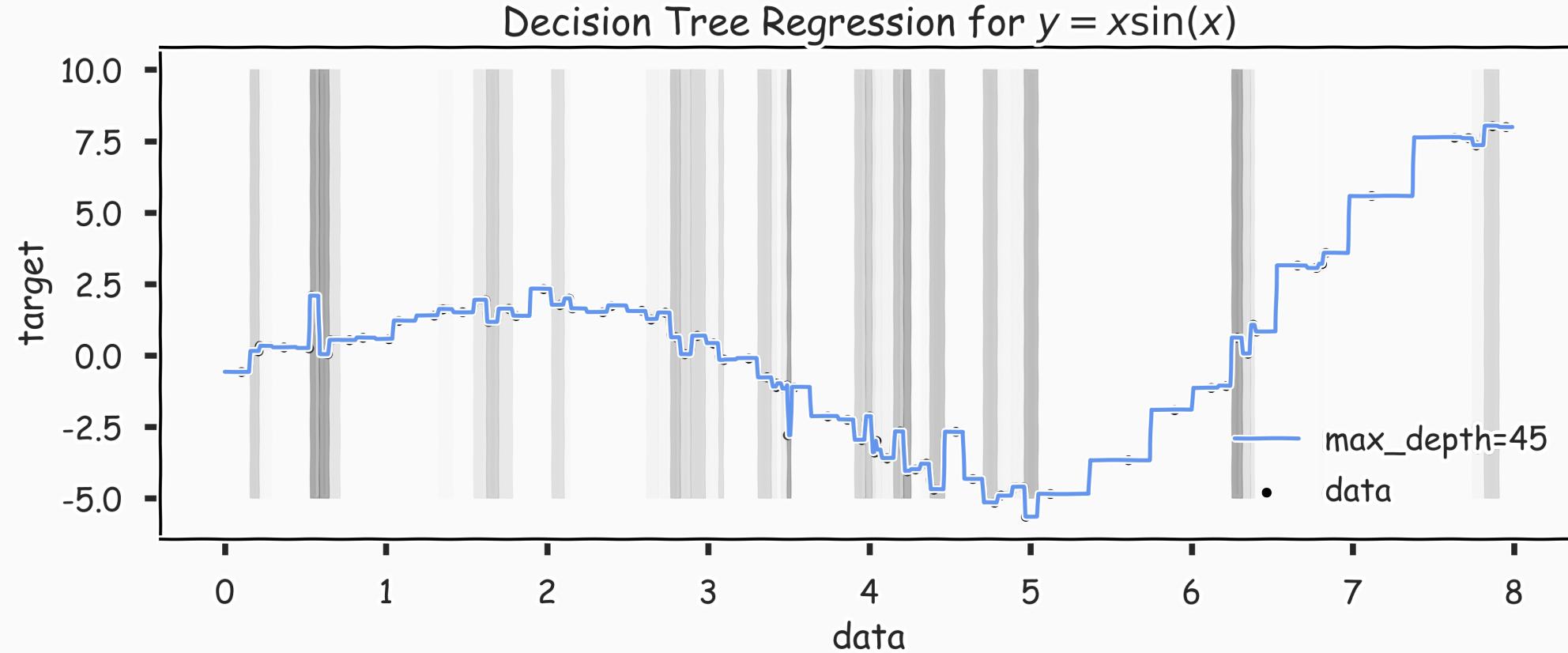
Regression Trees Prediction (grey scale represents MSE)



Regression Trees Prediction (grey scale represents MSE)



Regression Trees Prediction (grey scale represents MSE)



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.

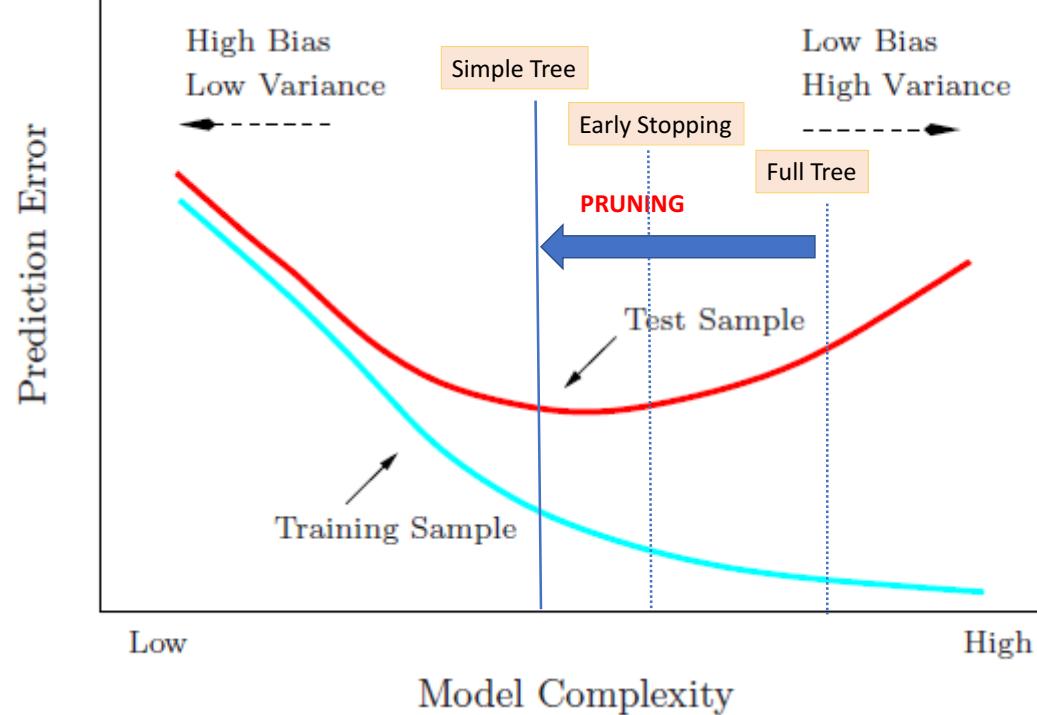
In the place of purity gain, we can instead compute accuracy gain for splitting a region R

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

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

Overfitting

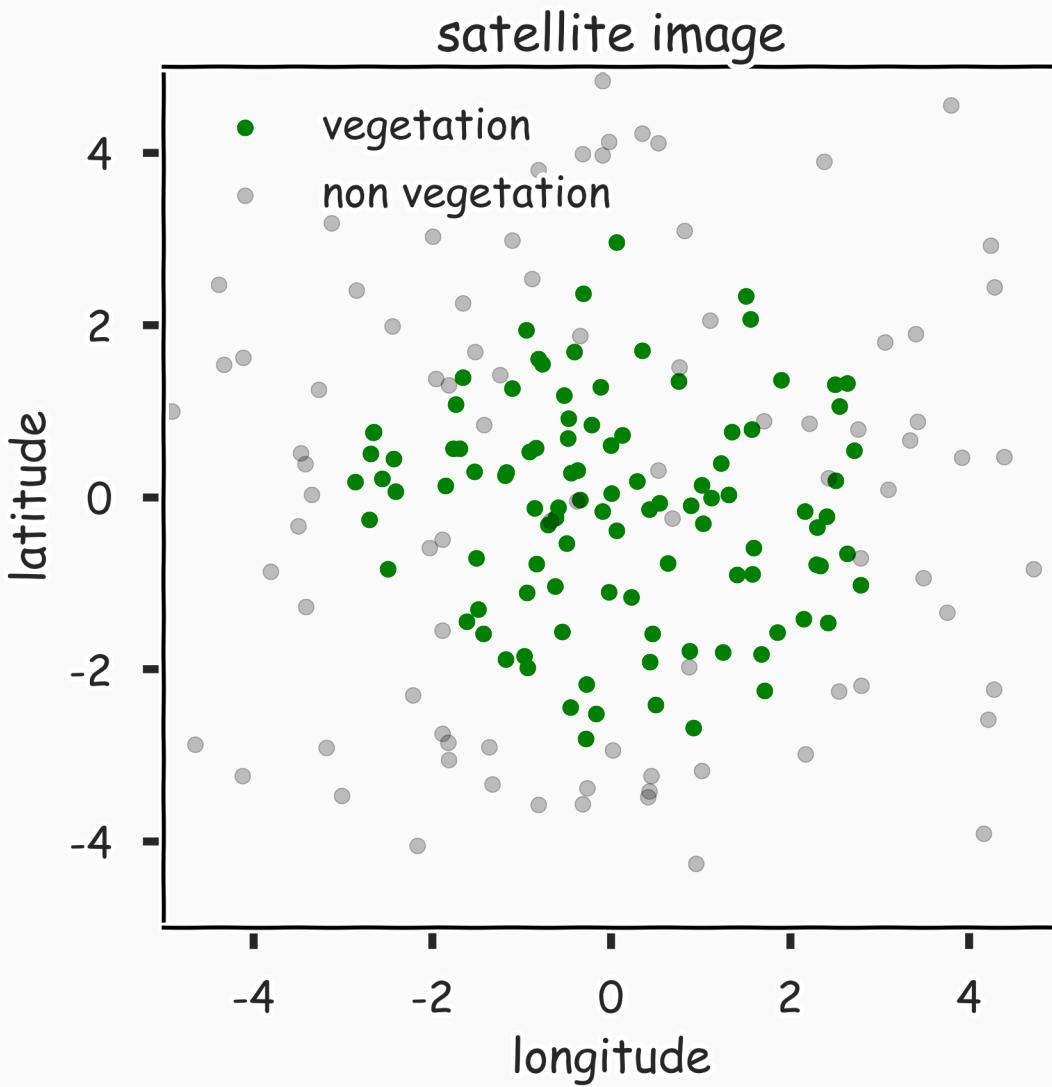
Same issues as with classification trees. Avoid overfitting by pruning or limiting the depth of the tree and using CV.



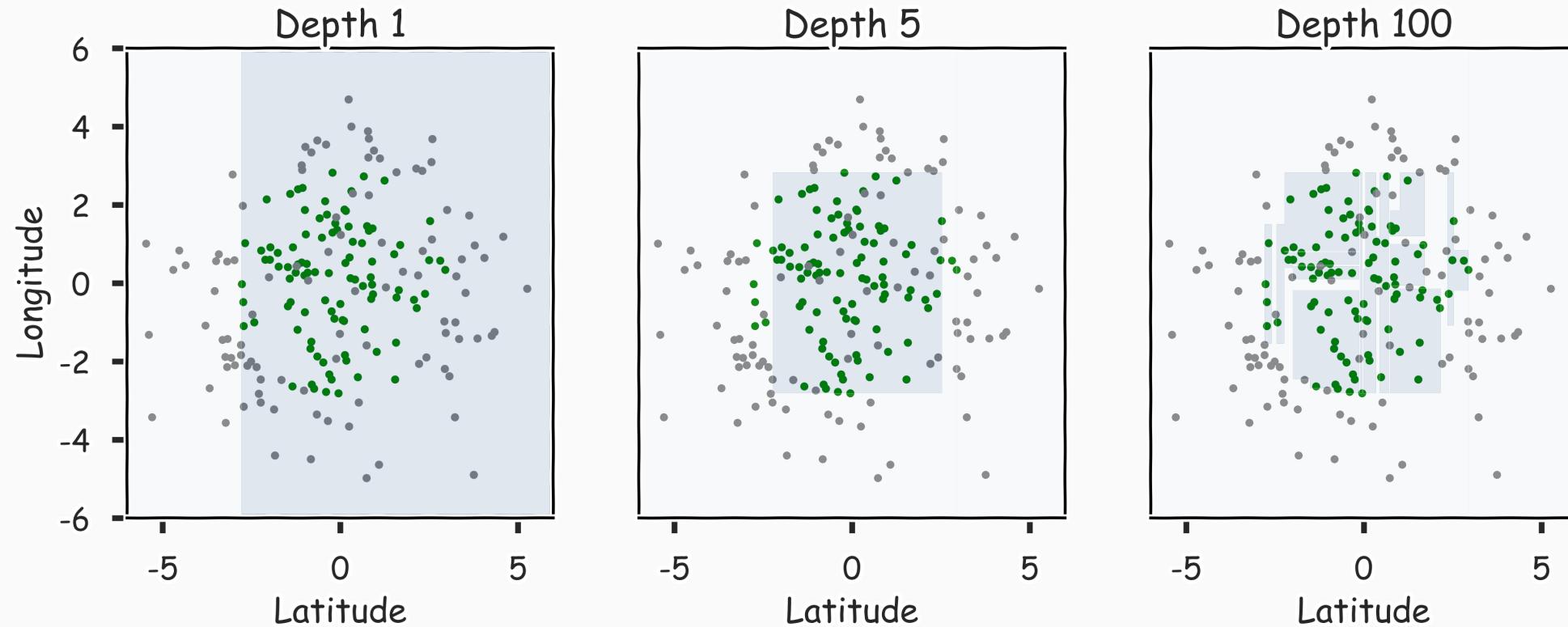
Bagging



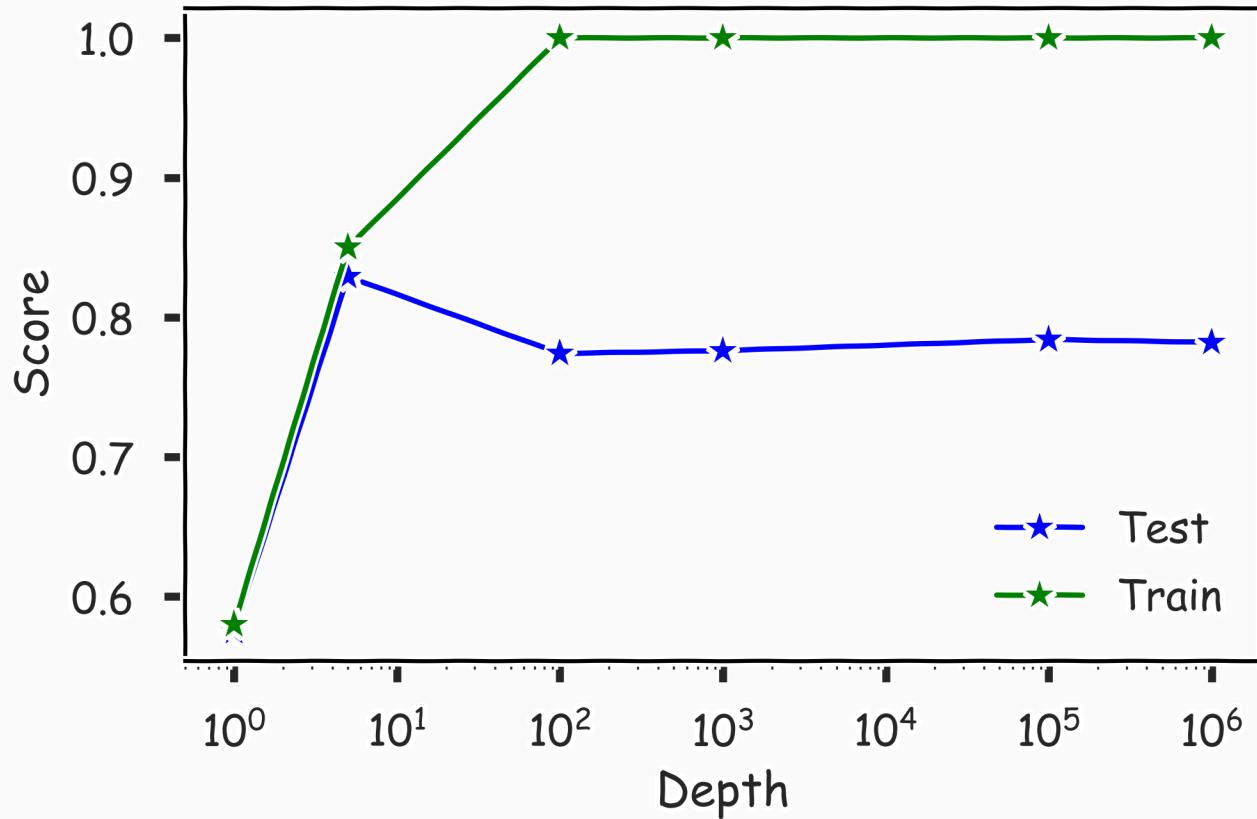
Reduce the variance



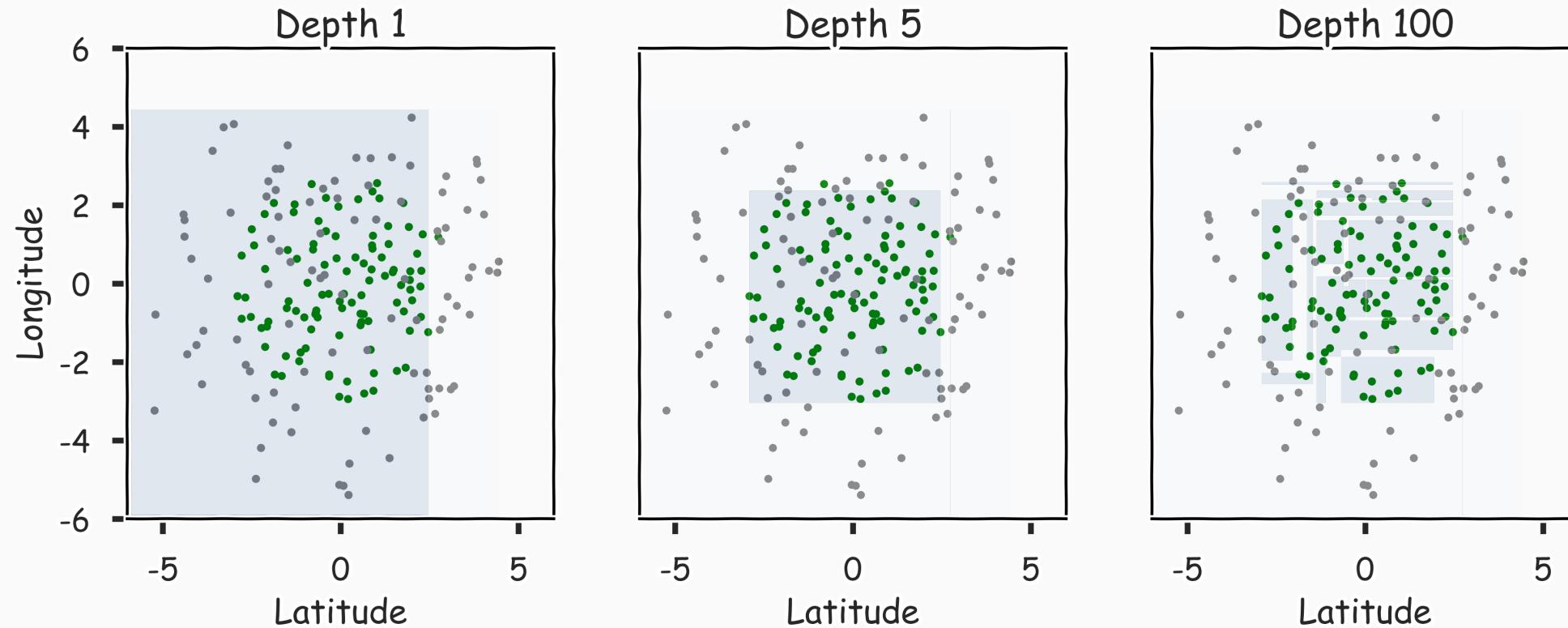
Hyper-parameters: Depth



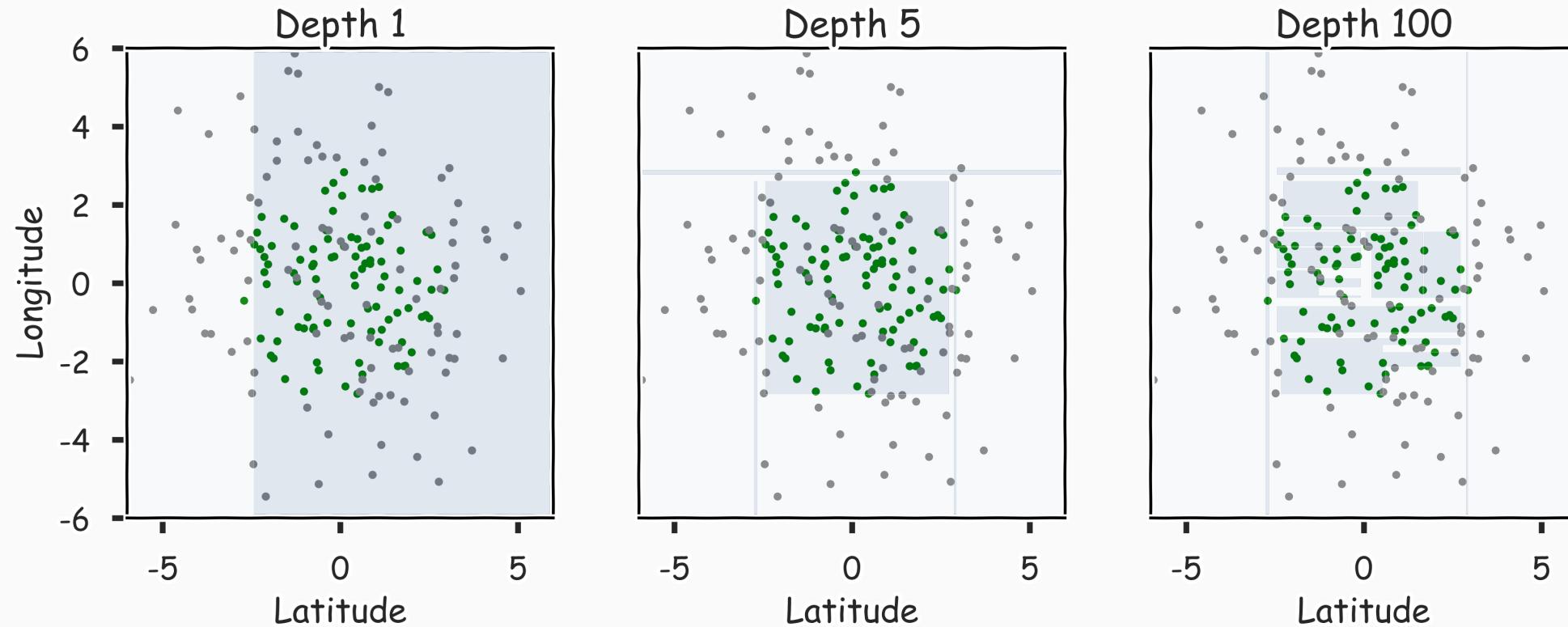
Hyper-parameters: Depth



Magic realism: Bootstrap



Magic realism: Bootstrap



Limitations of Decision Tree Models

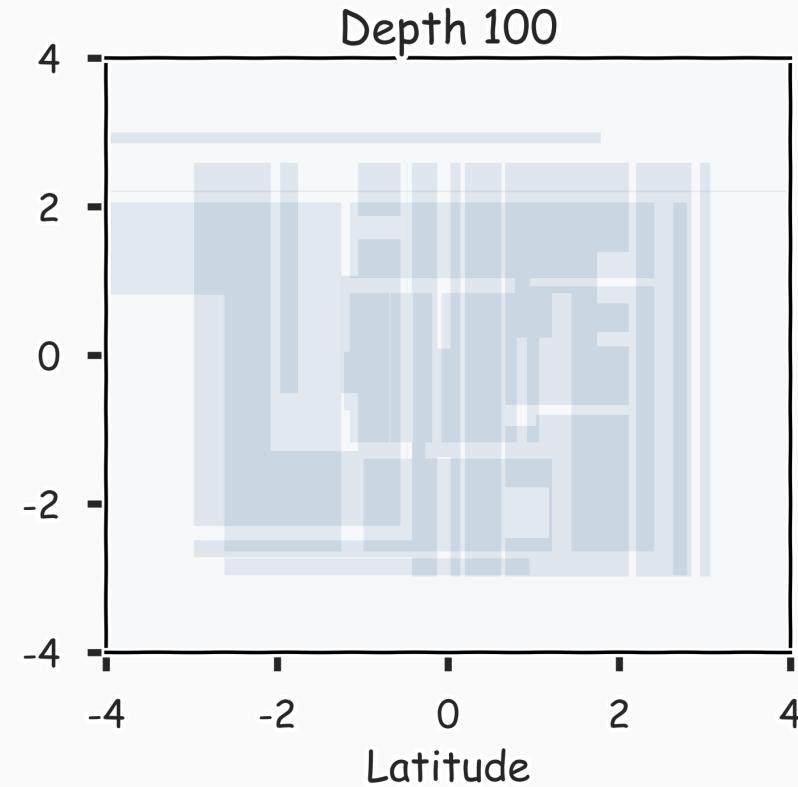
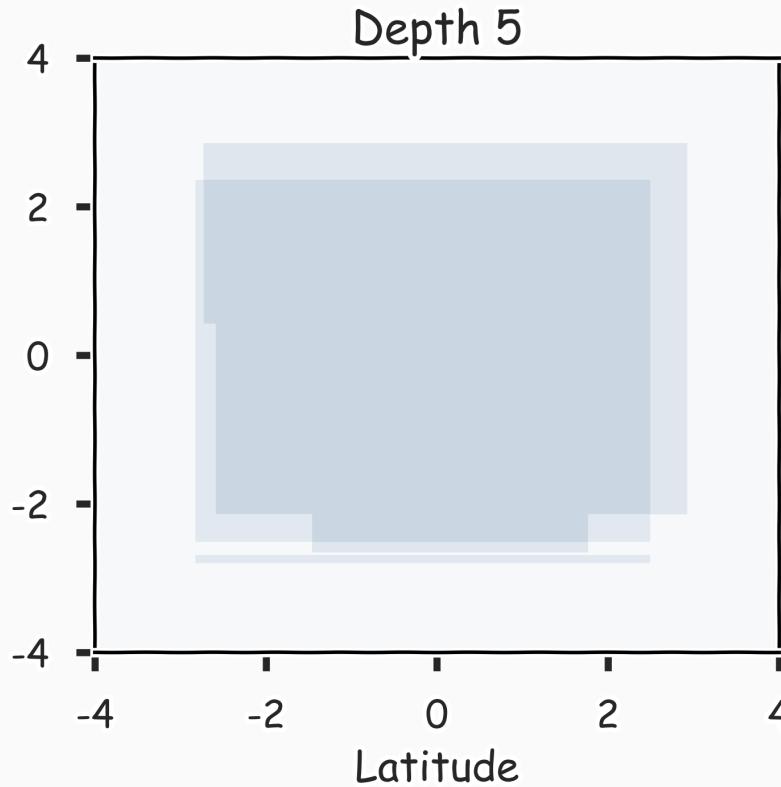
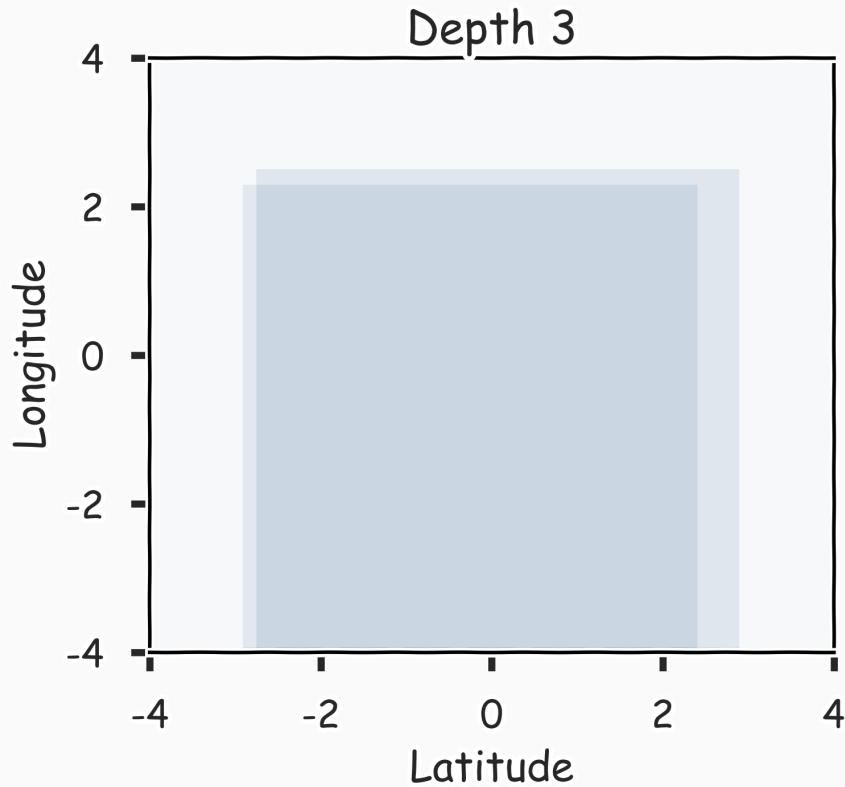
Decision trees models are highly interpretable and fast to train, using our greedy learning algorithm.

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).

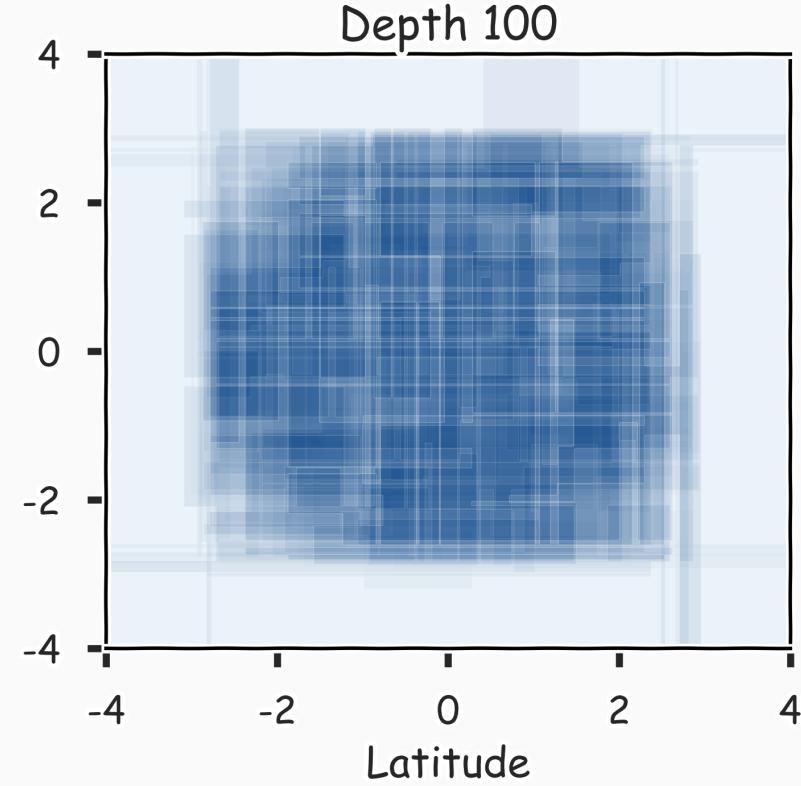
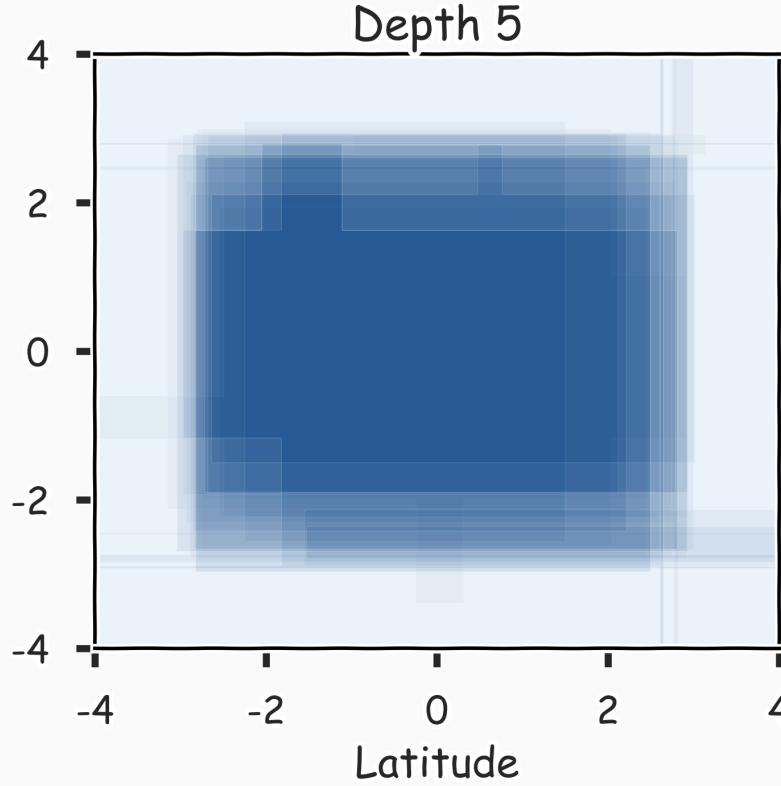
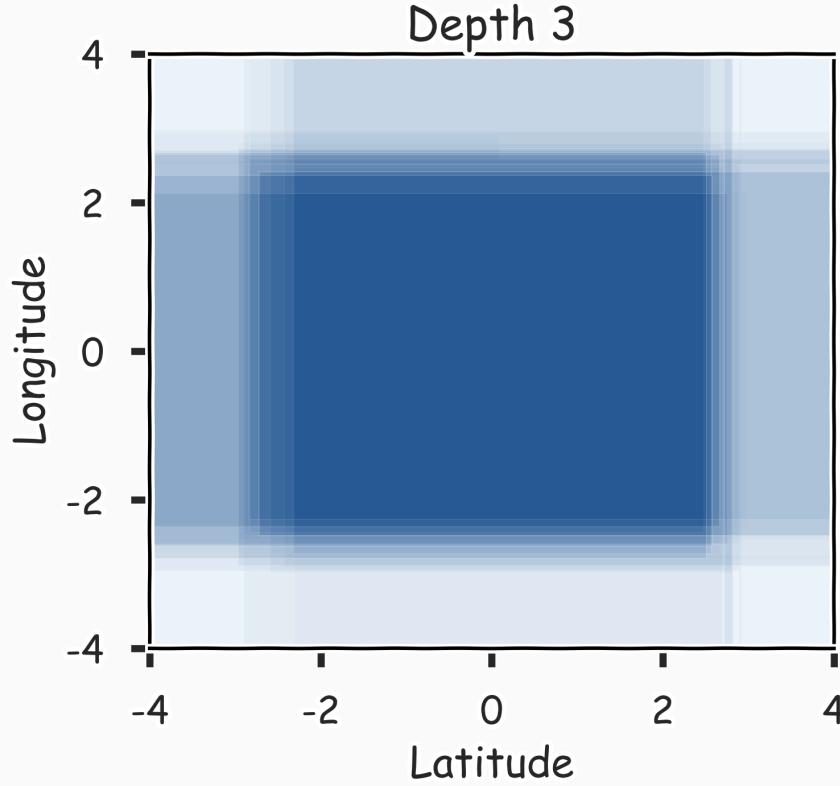
We've seen that large trees have high variance and are prone to overfitting.

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

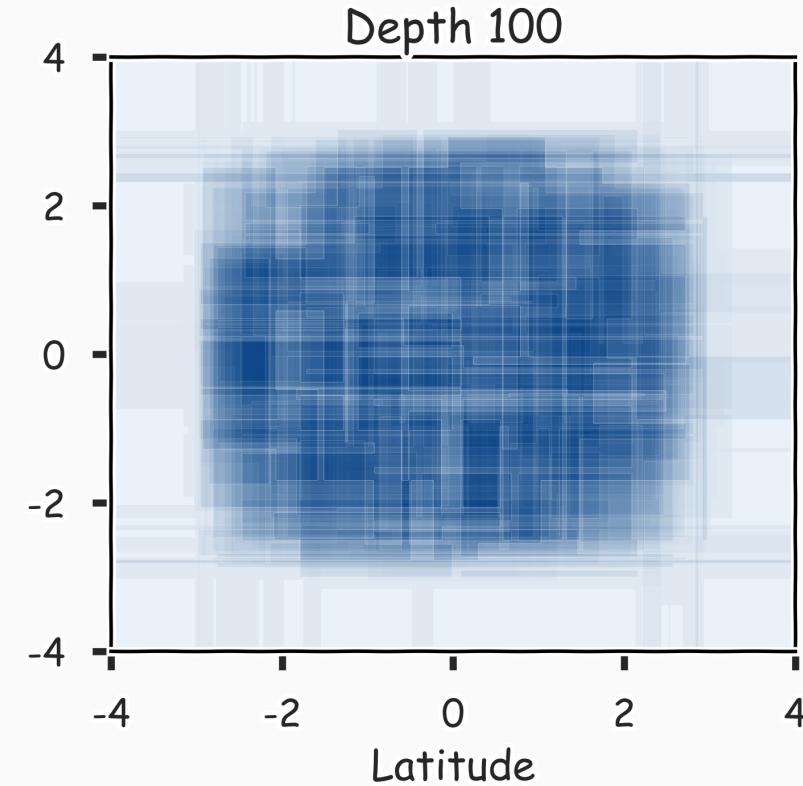
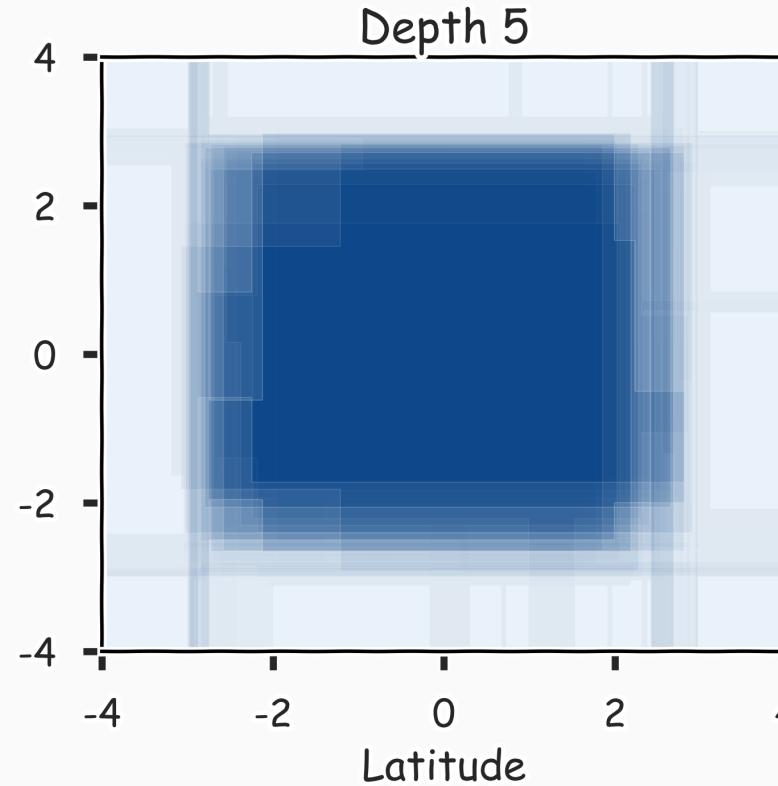
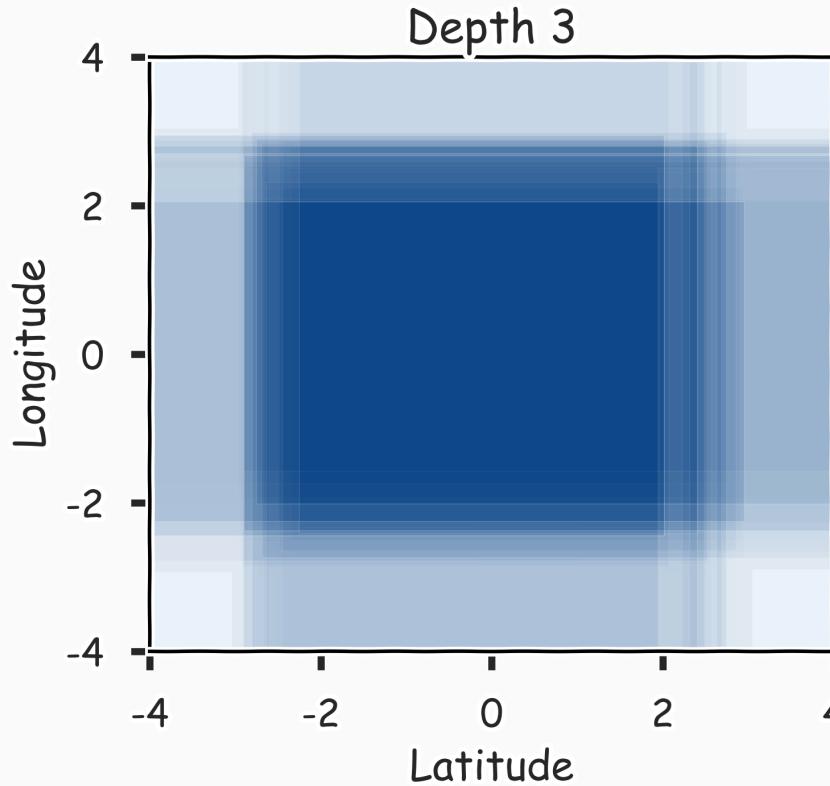
Combine them? 2 magic realisms



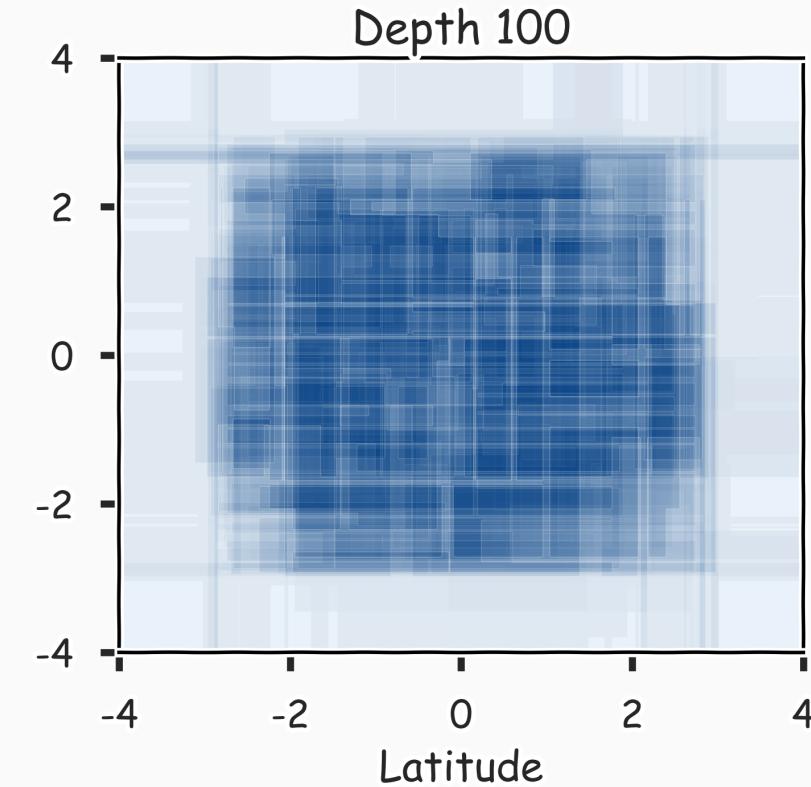
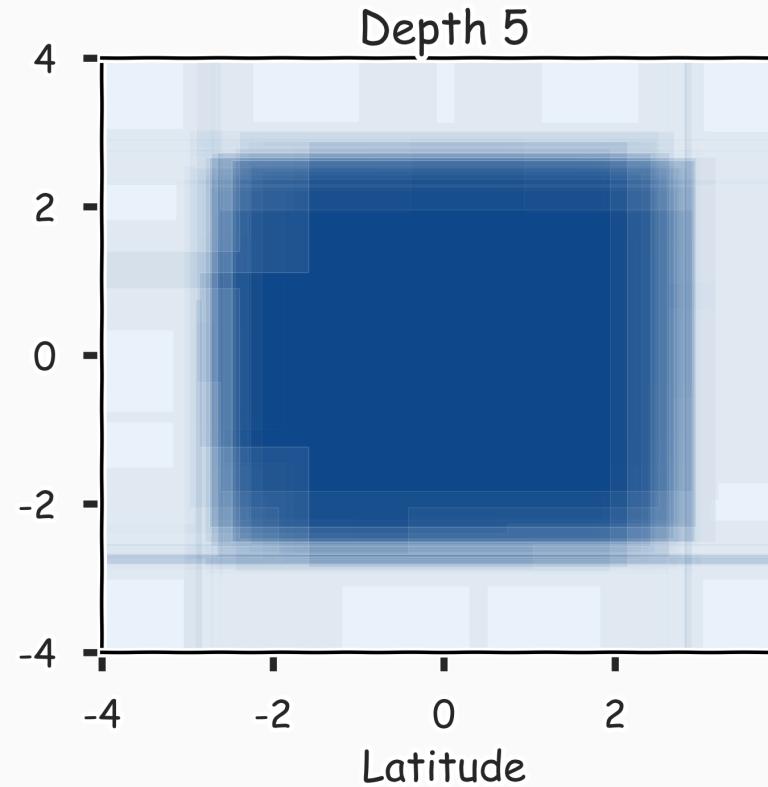
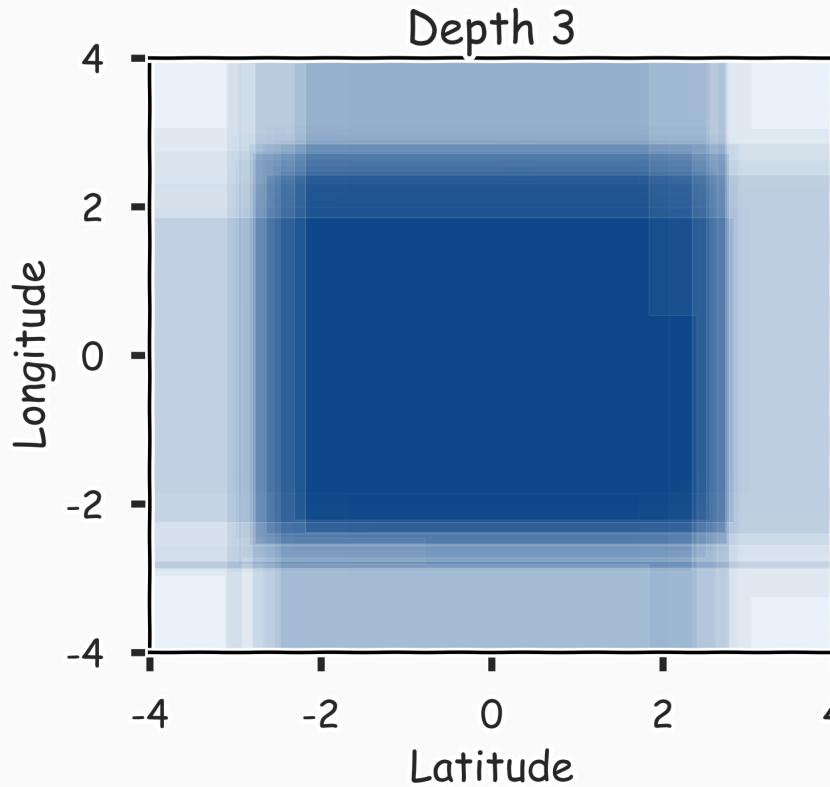
Combine them? 20 magic realisms



Combine them? 100 magic realisms



Combine them? 300 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.

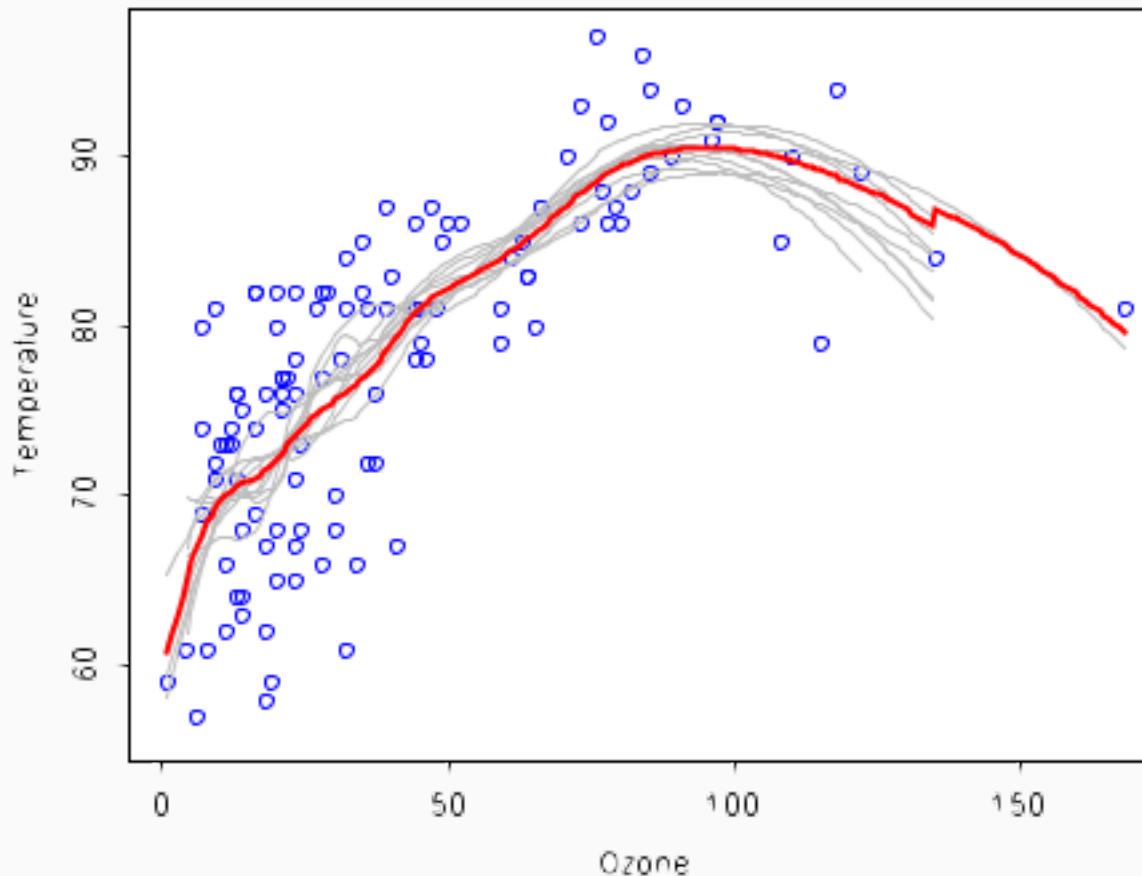
This method is called ***Bagging*** (Breiman, 1996), short for, of course, Bootstrap Aggregating.

Bagging

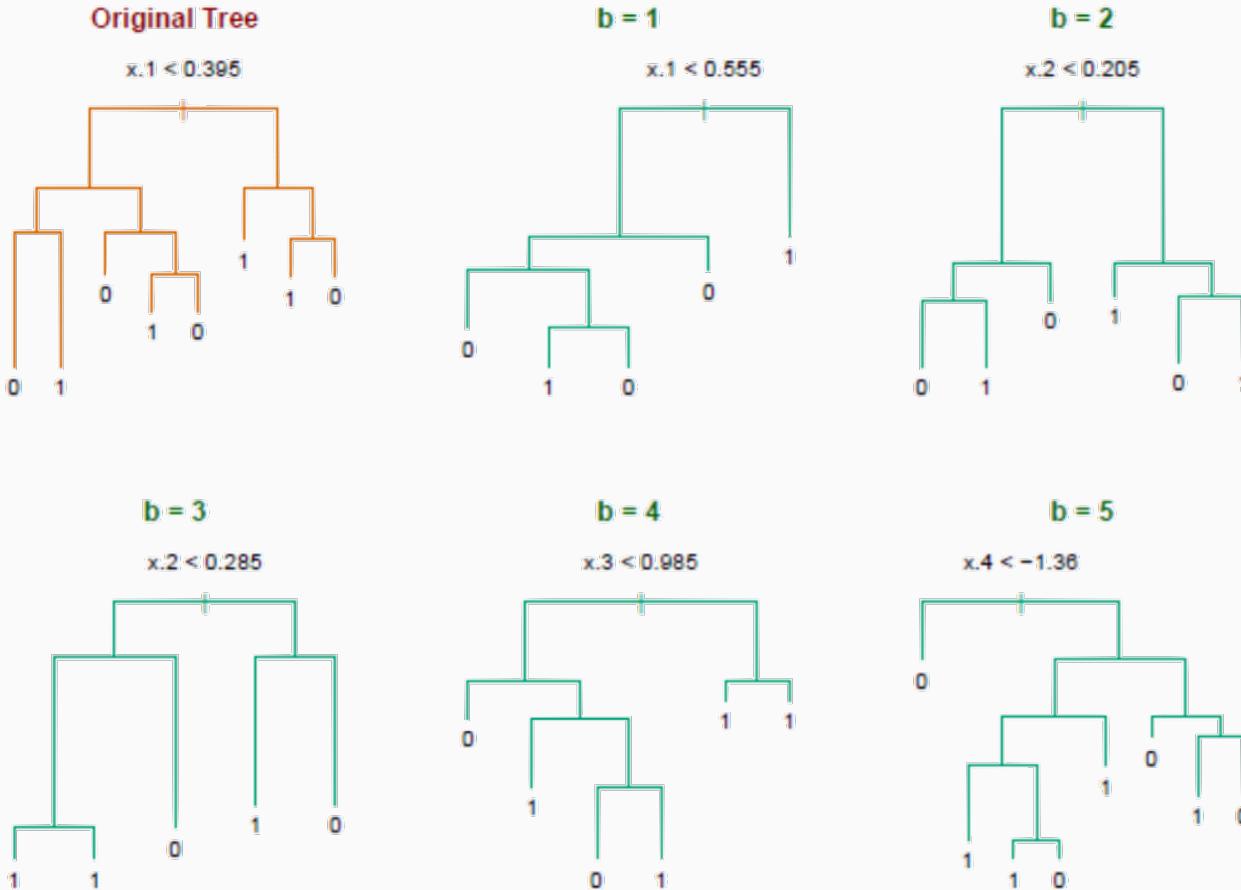
Note that bagging enjoys the benefits of:

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



Bagging

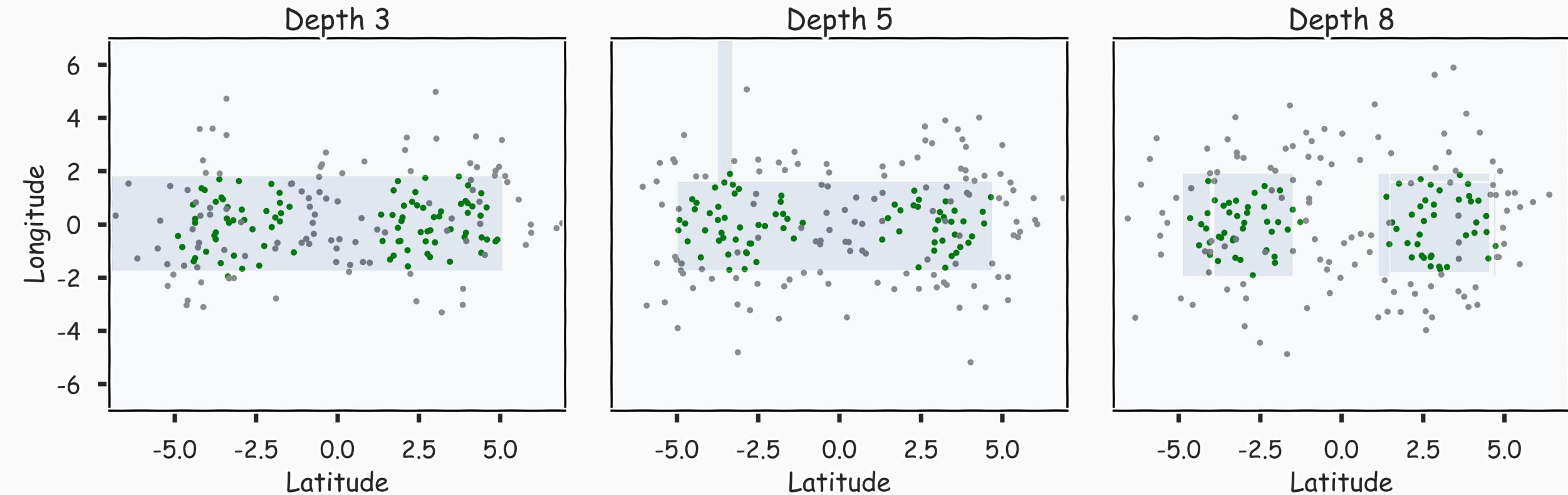


Bagging

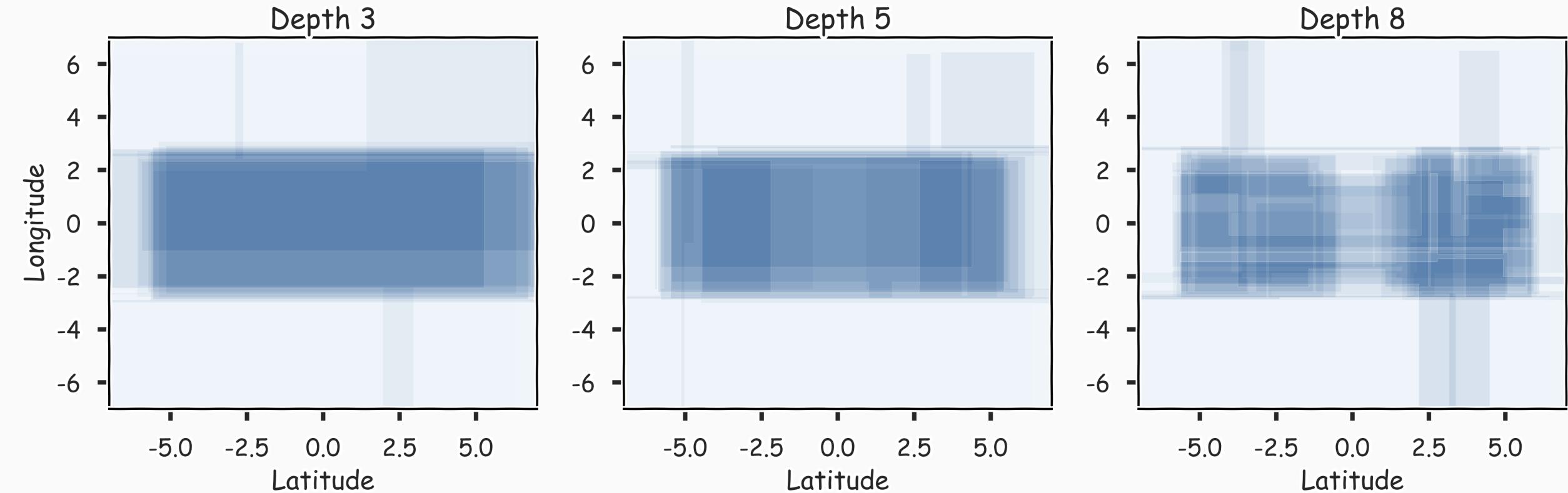
Question: Do you see any problems?

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

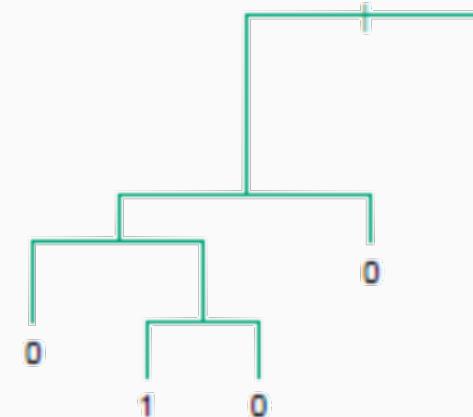
Original Data

X	Y
X↓1	y↓1
X↓2	y↓2
X↓3	y↓3
X↓4	y↓4
X↓5	y↓5
:	:
X↓n	y↓n

Bootstrap Sample 1

X	Y
X↓4	y↓4
X↓14	y↓14
X↓1	y↓1
X↓2	y↓2
X↓35	y↓35
:	:
X↓k	y↓k

Decision Tree 1



Used and unused data

X	Y
X↓1	y↓1
X↓2	y↓2
X↓3	y↓3
X↓4	y↓4
X↓5	y↓5
:	:
X↓n	y↓n

Bagging

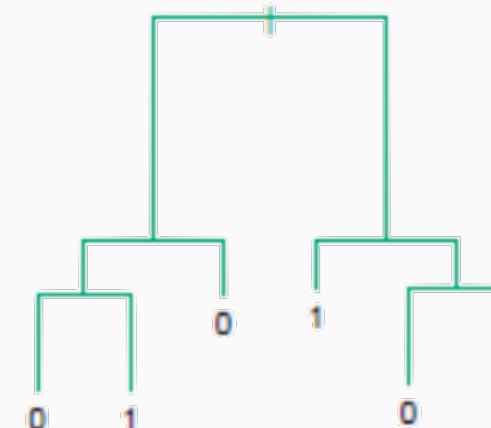
Original Data

X	Y
X↓1	y↓1
X↓2	y↓2
X↓3	y↓3
X↓4	y↓4
X↓5	y↓5
:	:
X↓n	y↓n

Bootstrap Sample 2

X	Y
X↓5	y↓5
X↓3	y↓3
X↓12	y↓12
X↓43	y↓43
X↓1	y↓1
:	:
X↓k	y↓k

Decision Tree 2



Used and unused data

X	Y
X↓1	y↓1
X↓2	y↓2
X↓3	y↓3
X↓4	y↓4
X↓5	y↓5
:	:
X↓n	y↓n

Bagging

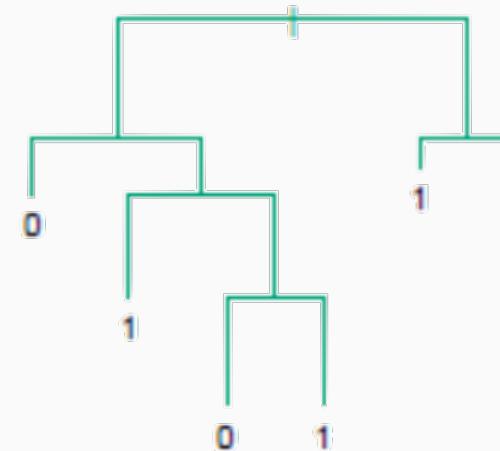
Original Data

X	Y
X↓1	y↓1
X↓2	y↓2
X↓3	y↓3
X↓4	y↓4
X↓5	y↓5
:	:
X↓n	y↓n

Bootstrap Sample 3

X	Y
X↓9	y↓9
X↓4	y↓4
X↓1	y↓1
X↓1	y↓1
X↓65	y↓65
:	:
X↓k	y↓k

Decision Tree 3



Used and unused data

X	Y
X↓1	y↓1
X↓2	y↓2
X↓3	y↓3
X↓4	y↓4
X↓5	y↓5
:	:
X↓n	y↓n

Point-wise out-of-bag error

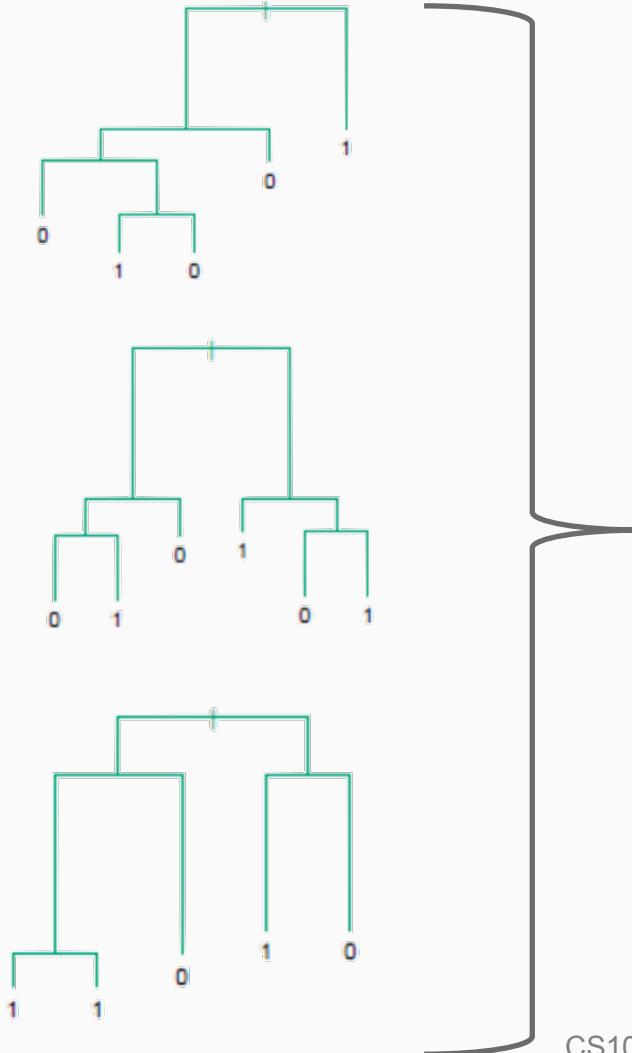
X	Y
$X \downarrow 1$	$y \downarrow 1$
$X \downarrow 2$	$y \downarrow 2$
$X \downarrow 3$	$y \downarrow 3$
\vdots	\vdots
$X \downarrow i$	$y \downarrow i$
\vdots	\vdots
$X \downarrow n$	$y \downarrow n$



Point-wise out-of-bag error

B Trees that did not see $\{X \downarrow i, y \downarrow i\}$

X	Y
$X \downarrow 1$	$y \downarrow 1$
$X \downarrow 2$	$y \downarrow 2$
$X \downarrow 3$	$y \downarrow 3$
:	:
$X \downarrow i$	$y \downarrow i$
:	:
$X \downarrow n$	$y \downarrow n$



Classification

$$y \downarrow i, pw = \text{majority}(\square y \downarrow i)$$

$$e \downarrow i = \mathbb{I}(y \downarrow i, pw \neq y \downarrow i)$$

Regression

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

$$e \downarrow i = (y \downarrow i - y \downarrow i, pw) / 2$$

OOB Error

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

Classification

$$\text{Error}_{\text{OOB}} = \sum_{i=1}^n e_{i,i} = \sum_{i=1}^n \mathbb{I}(y_{i,pw} \neq y_{i,i})$$

Regression

$$\text{Error}_{\text{OOB}} = \sum_{i=1}^n e_{i,i} = \sum_{i=1}^n (y_{i,pw} - y_{i,i})^2$$

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 modeled 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

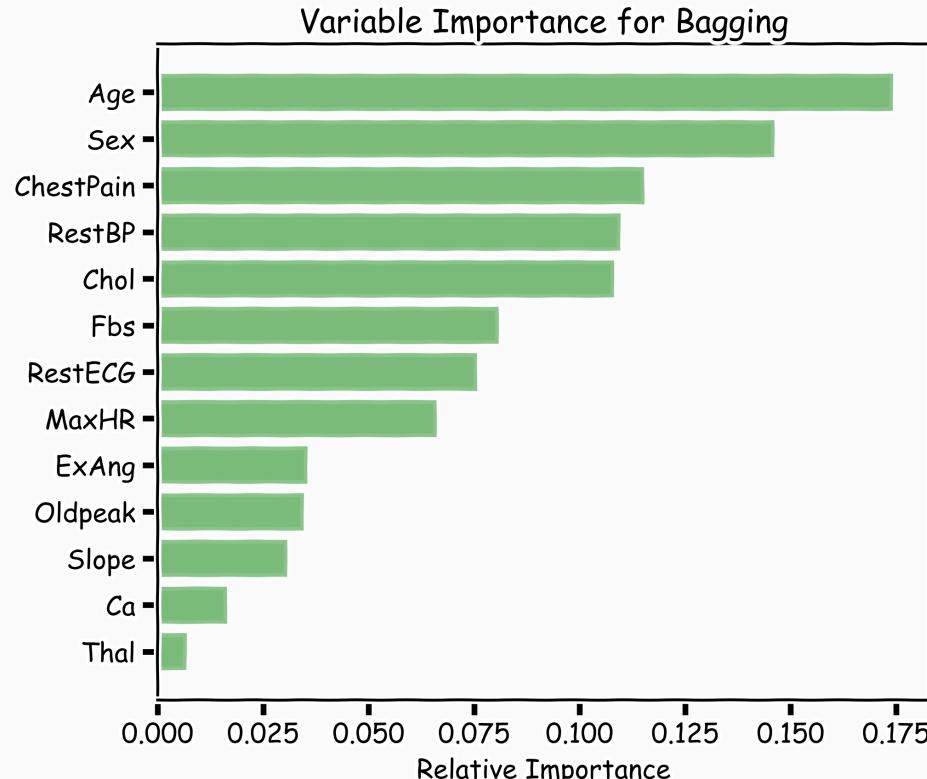
Question: Do you see any problems?

- 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 RSS (for regression) or Gini index (for classification) is decreased due to splits over a given predictor, averaged over all B trees.



100 trees, max_depth=10

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!

Improving on Bagging

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

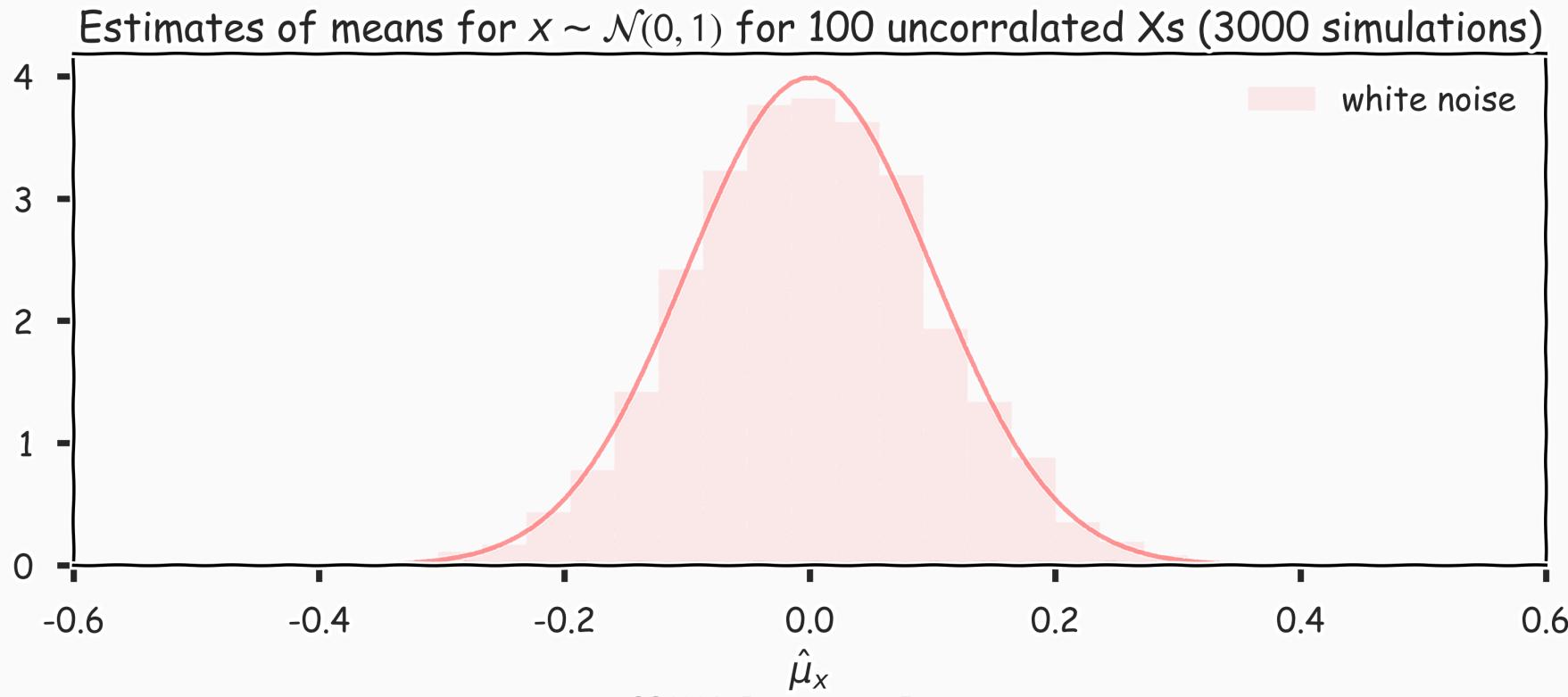
Suppose we have an extremely strong predictor, $x \downarrow j'$, in the training set amongst moderate predictors. Then the greedy learning algorithm ensures that most of the models in the ensemble will choose to split on $x \downarrow j'$ in early iterations.

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 B number of identically and independently distributed variable, X , with variance σ^2 , the variance of the estimate of the mean is :

$$\text{var}(\hat{\mu}_x) = \sigma^2 / B$$

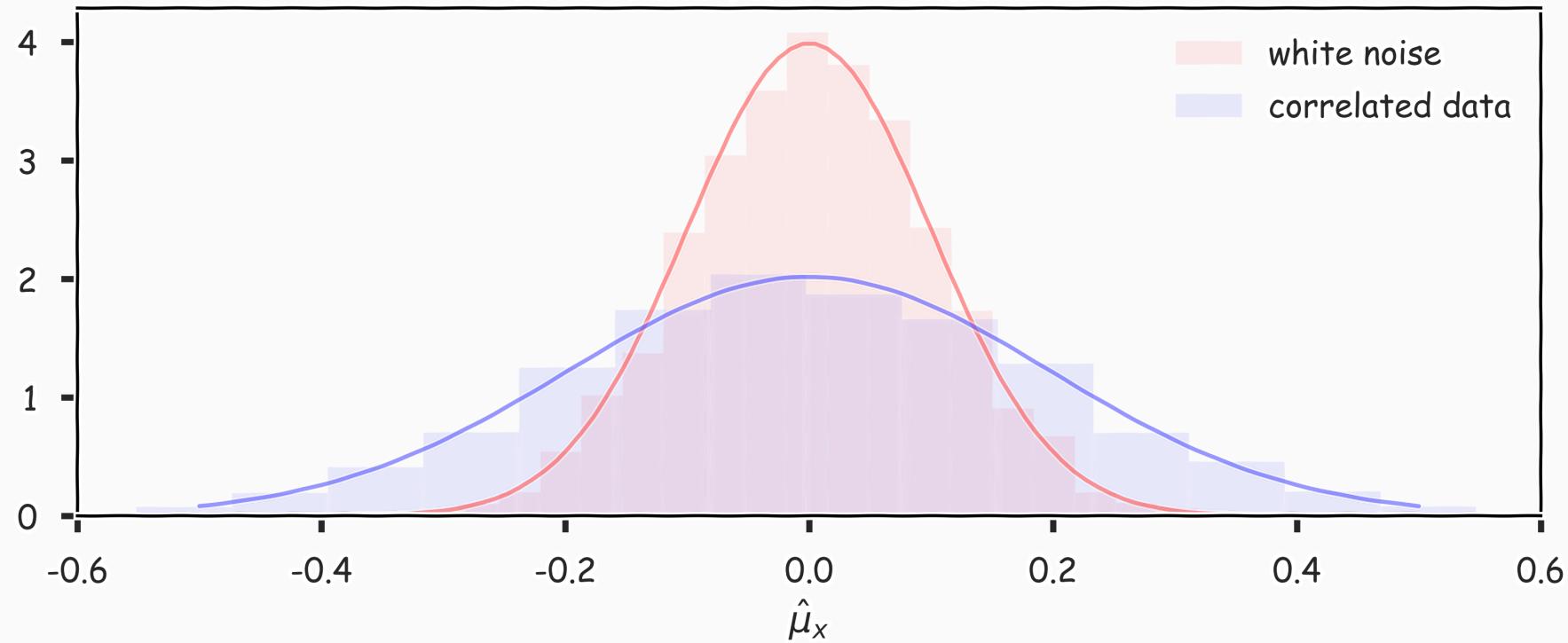


Improving on Bagging

For B number of identically but not independently distributed variables with pairwise correlation ρ and variance σ^2 , the variance of their mean is

$$\text{var}(\bar{\mu}_x) \propto \sigma^2 \rho / B$$

Estimates of means for correlated xs, $\rho = 0.5$, for 100 Xs. Here we show the results for 3000 simulations



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 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.

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:

- \sqrt{N}/j for classification
- $N/3$ 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:

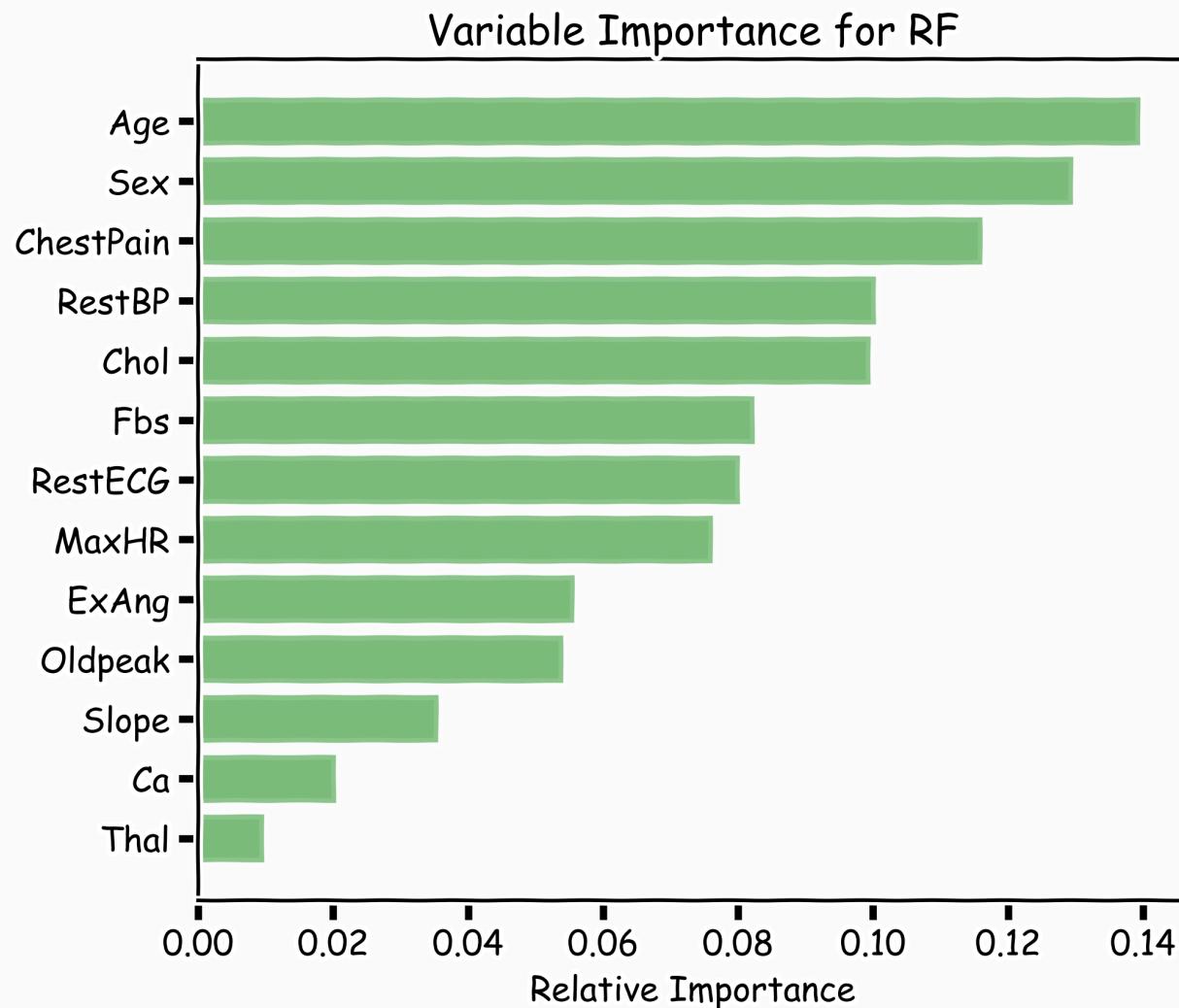
Calculate the total amount that the RSS (for regression) or Gini index (for classification) is decreased due to splits over a given predictor, averaged over all B trees.

Variable Importance for RF

Alternative:

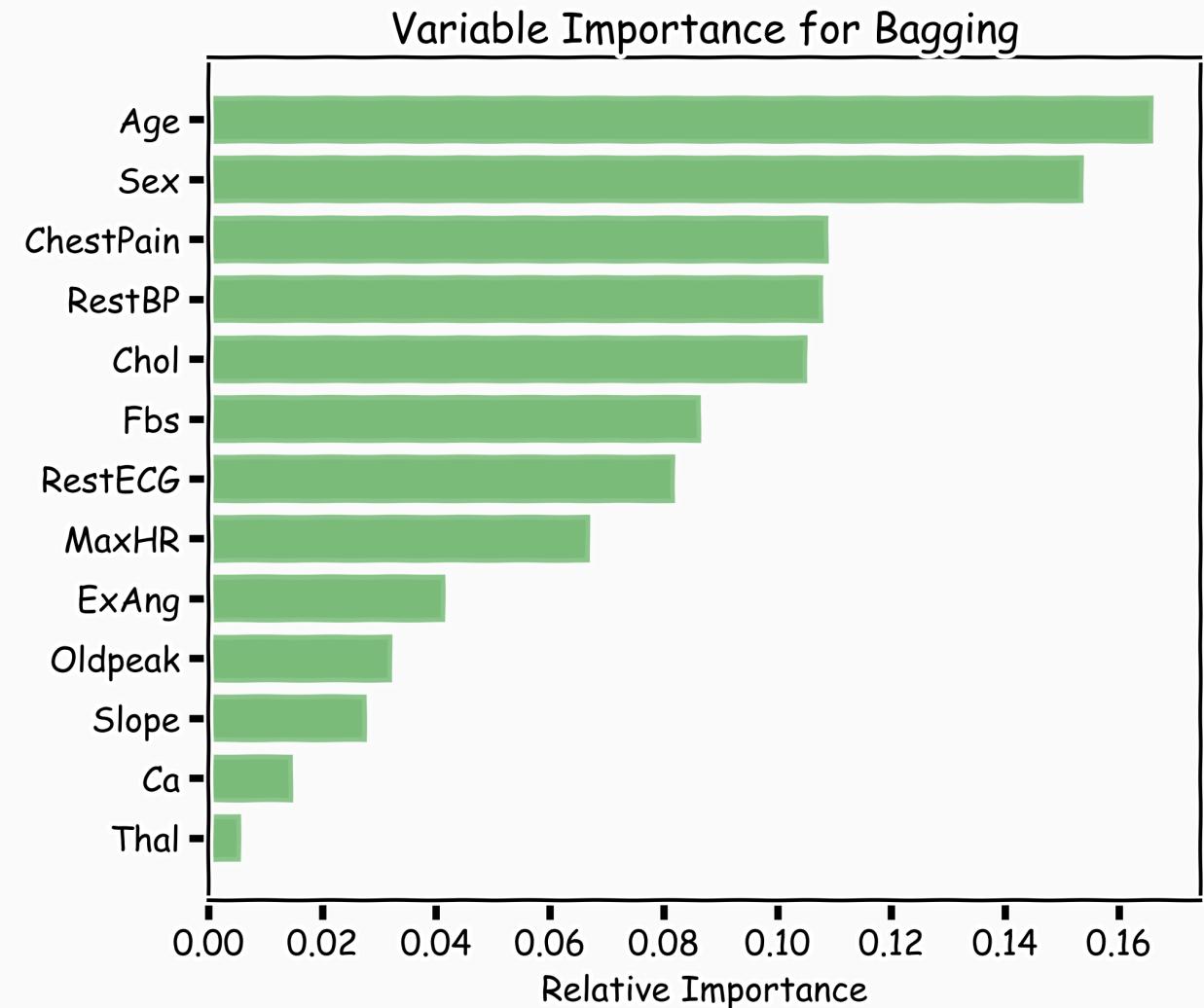
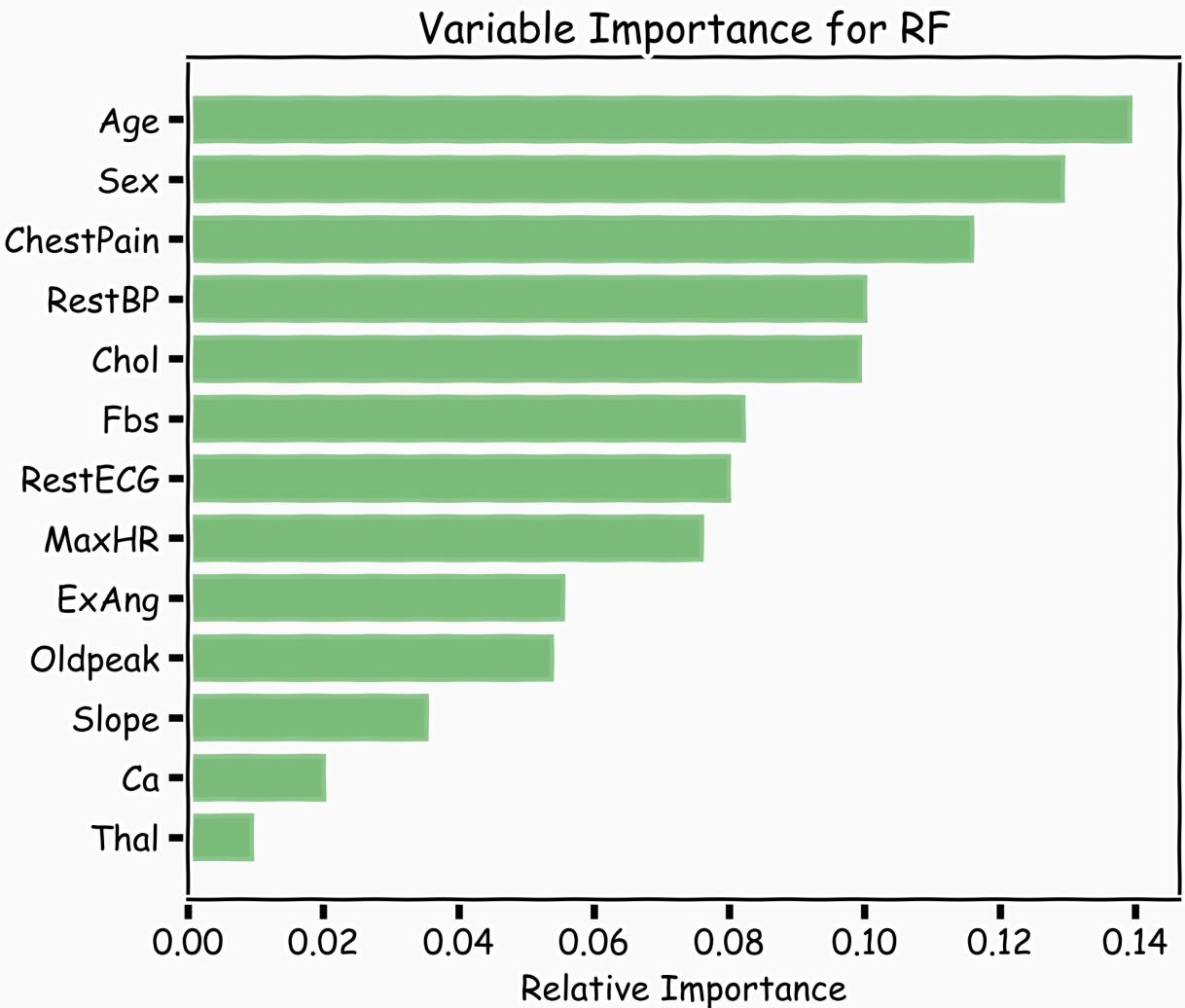
- Record the prediction accuracy on the *oob* samples for each tree.
- Randomly permute the data for column j in the *oob* samples the record the accuracy again.
- The decrease in accuracy as a result of this permuting is averaged over all trees, and is used as a measure of the importance of variable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

Unbalance dataset:

Weighted samples:

Categorical data:

Missing data: a-sec later today

Different implementations: a-sec later today