

The first step if you're new to machine learning.

Tutorial Data

Learn Tutorial

Intro to Machine Learning

Course step
1 of 7 ▼

Table of Contents

Introduction

Improving the Decision Tree

Continue

Machine learning works the same way. We'll start with a model called the Decision Tree. There are fancier models that give more accurate predictions. But decision trees are easy to understand, and they are the basic building block for some of the best models in data science.

For simplicity, we'll start with the simplest possible decision tree.

```
graph TD; A[Does house have more than 2 Bedrooms] -- No --> B[Predicted Price: $178000]; A -- Yes --> C[Predicted Price: $188000]
```

The diagram is a decision tree. The root node is a light blue rectangle with a black border containing the text "Does house have more than 2 Bedrooms". Two arrows originate from the bottom of this node. The left arrow is labeled "No" and points to a light blue rectangle with a black border containing the text "Predicted Price: \$178000". The right arrow is labeled "Yes" and points to a light blue rectangle with a black border containing the text "Predicted Price: \$188000".

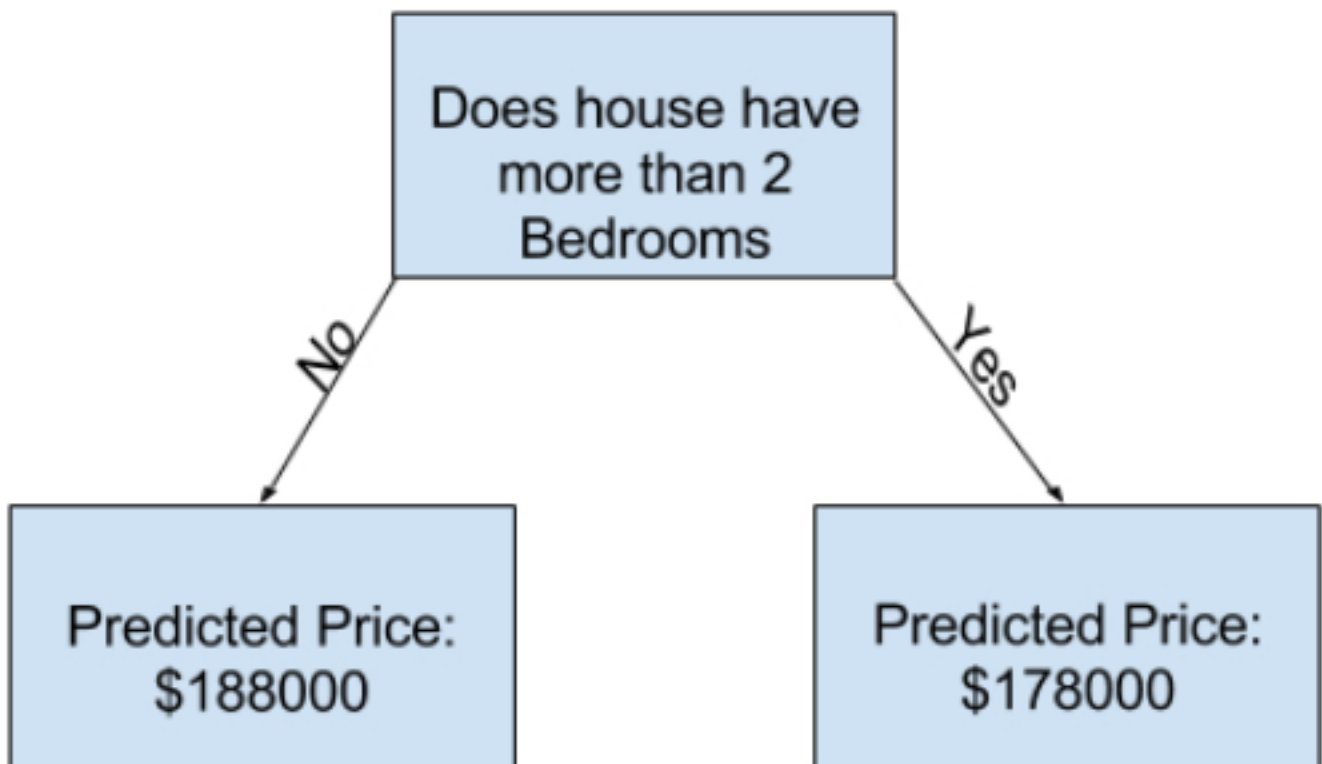
We use data to decide how to break the houses into two groups, and then again to determine the predicted price in each group. This step of capturing patterns from data is called **fitting** or **training** the model. The data used to **fit** the model is called the **training data**.

The details of how the model is fit (e.g. how to split up the data) is complex enough that we will save it for later. After the model has been fit, you can apply it to new data to **predict** prices of additional homes.

Which of the following two decision trees is more likely to result from fitting the real estate training data?

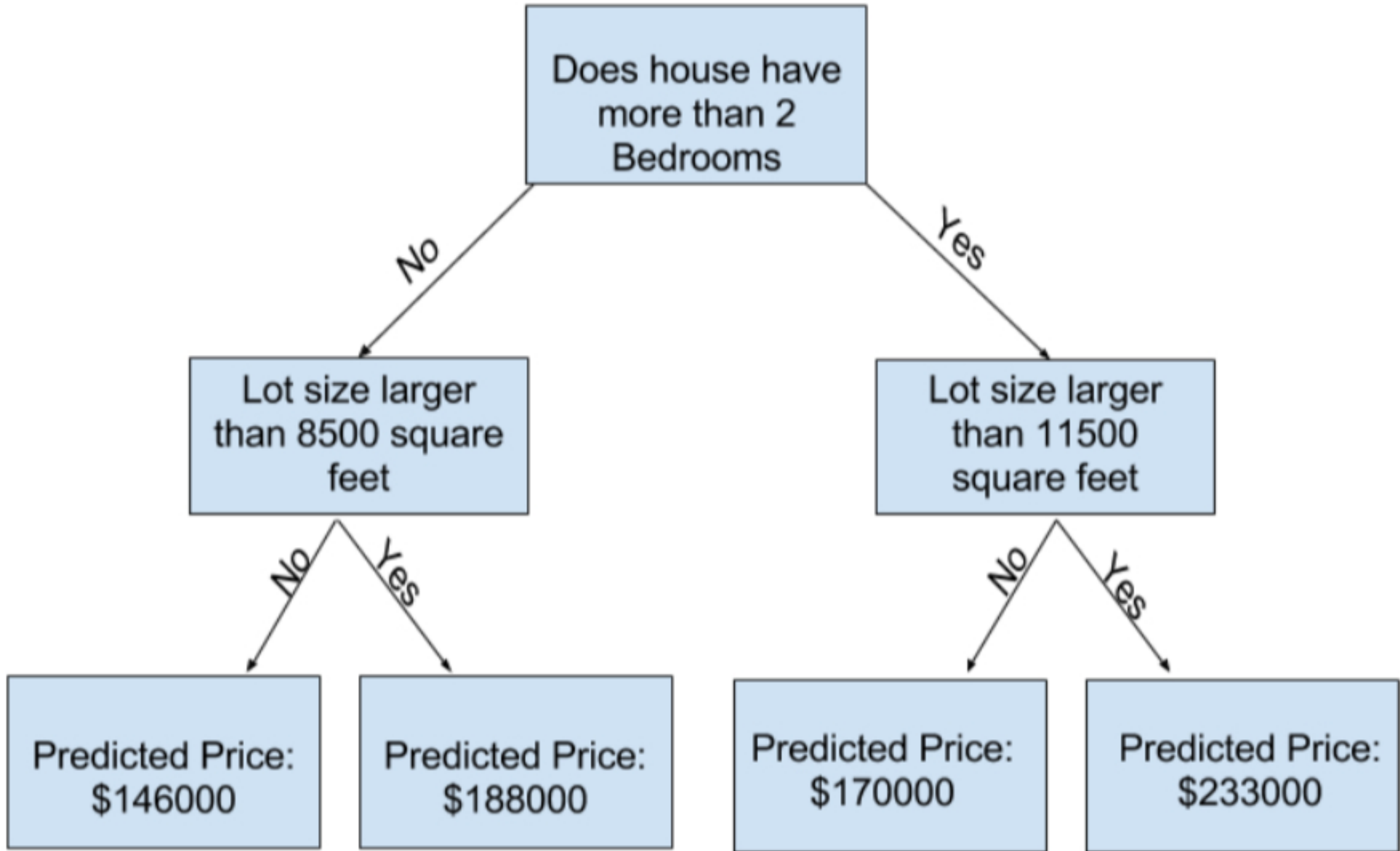
```
graph TD; A[Does house have more than 2 Bedrooms] -- No --> B[Predicted Price: $178000]; A -- Yes --> C[Predicted Price: $188000];
```

The diagram is a decision tree. The root node is a light blue rectangle containing the text "Does house have more than 2 Bedrooms". Two arrows originate from the bottom of this node. The left arrow is labeled "No" and points to a light blue rectangle containing "Predicted Price: \$178000". The right arrow is labeled "Yes" and points to a light blue rectangle containing "Predicted Price: \$188000".



The decision tree on the left (Decision Tree 1) probably makes more sense, because it captures the reality that houses with more bedrooms tend to sell at higher prices than houses with fewer bedrooms. The biggest shortcoming of this model is that it doesn't capture most factors affecting home price, like number of bathrooms, lot size, location, etc.

You can capture more factors using a tree that has more "splits." These are called "deeper" trees. A decision tree that also considers the total size of each house's lot might look like this:



You predict the price of any house by tracing through the decision tree, always picking the path corresponding to that house's characteristics. The predicted price for the house is at the bottom of the tree. The point at the bottom where we make a prediction is called a **leaf**.

The splits and values at the leaves will be determined by the data, so it's time for you to check out the data you will be working with.

Let's get more specific. It's time to **Examine Your Data.**

Have questions or comments? Visit the [course discussion forum](#) to chat with other learners.