

Plan:

1. Work through an ML example
2. Introduce overfitting

# Machine Learning: Example

Shannon E. Ellis, Ph.D  
UC San Diego



Department of Cognitive Science  
[sellis@ucsd.edu](mailto:sellis@ucsd.edu)

**Can we build a model that  
distinguishes a house in NY  
from a house in San Francisco?**

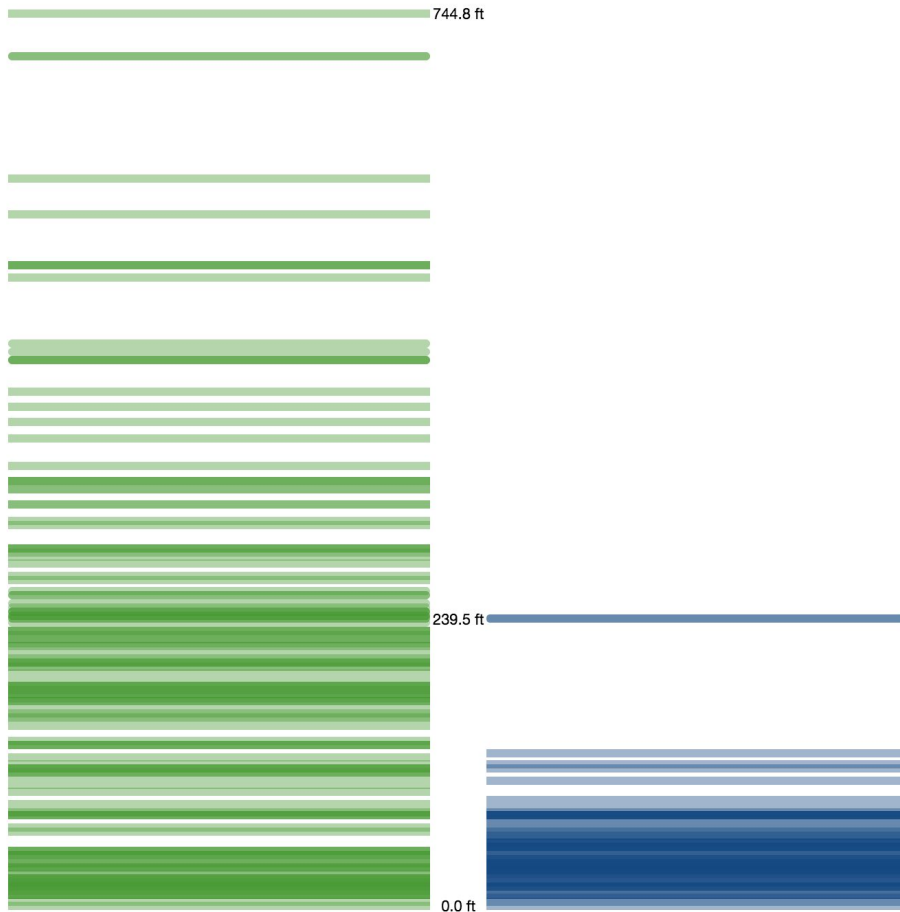
---

# First, some intuition

Let's say you had to determine whether a home is in **San Francisco** or in **New York**. In machine learning terms, categorizing data points is a **classification** task.

San Fran is hilly ...so elevation may be a helpful feature.

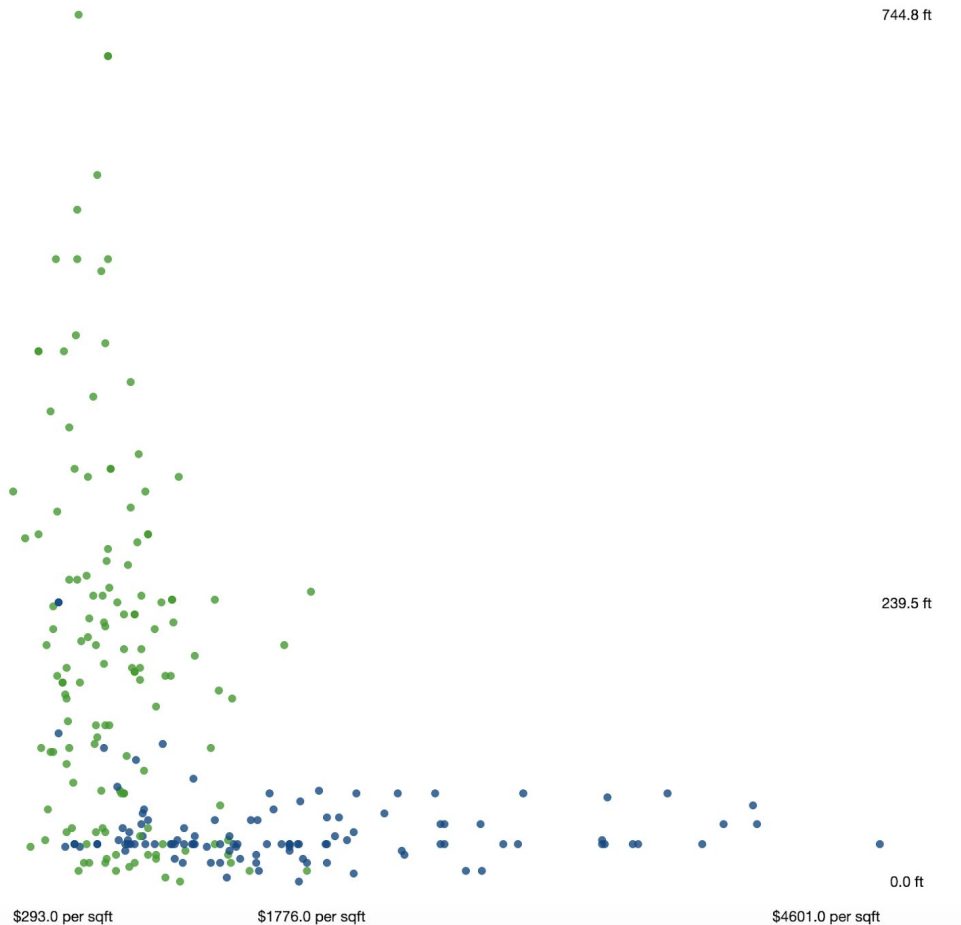
With the data here, homes > ~73m should be classified as San Fran homes



## Adding nuance

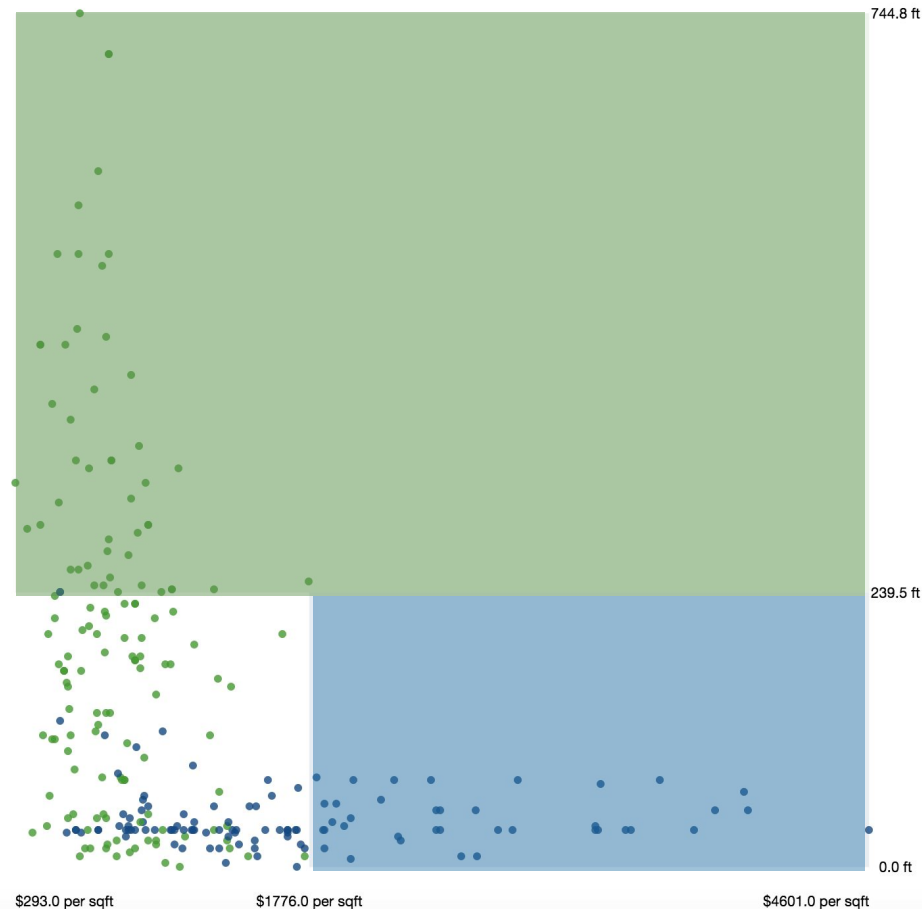
---

Elevation isn't a perfect feature for classification, so we can look at its relationship to other features, like *price per square foot*



## Drawing boundaries

Boundaries can be drawn so that if a house falls in the green box, it's classified as a San Fran home. Blue box, New York. Statistical learning figures out how to best draw these boxes.



Our training set will use 7 different **features**. At the right we see the **scatterplot matrix** of the relationship between these features.

Patterns are clear, but boundaries for delineation are not obvious.



**New York**  
**San Fran**

Our training set will use  
**features**. At the right  
**scatterplot matrix**  
relationship between

Patterns are clear, but  
delineation are not obvious

## And now, machine learning

---

Determining the best boundary is where  
**machine learning** comes in.

**Decision trees** are one example of  
machine learning method for classification  
tasks.

elevation





Finding better  
boundaries

---

We guessed ~73m  
before. Let's improve on  
that guess...





A **histogram** helps display frequency of homes by elevation more easily.

73m is the highest home in New York, but most of them have lower elevations

In machine learning, the splits are called **forks** and they split the data into **branches** based on some value.

The value that splits the branches is the **split point**. Homes to the left get categorized differently than those on the right.



## Your first fork

---

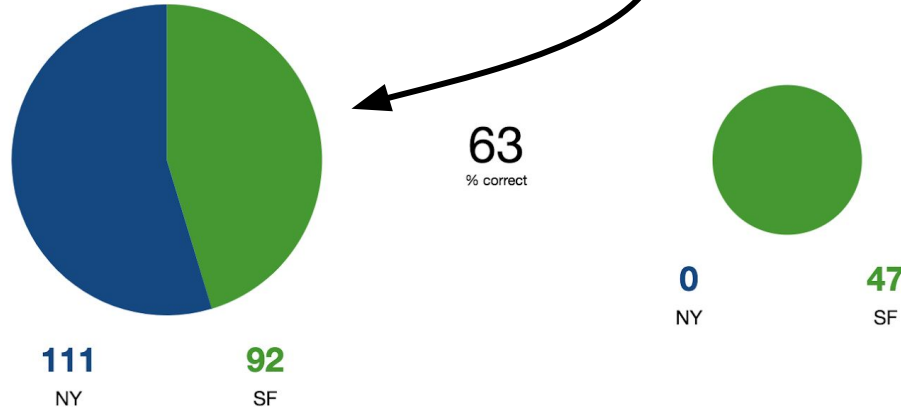
A decision tree uses if-then statements to define patterns in the data.

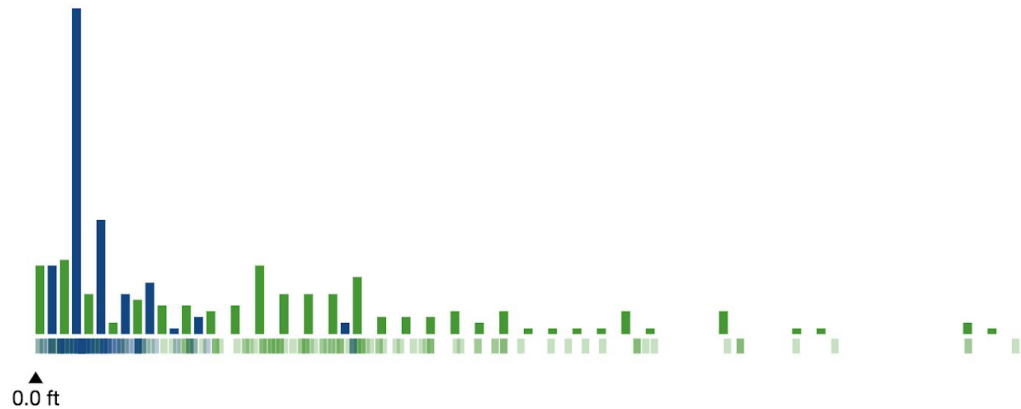
All the green over here are San Fran homes that were misclassified (**false negatives**)



## Tradeoffs

Splitting at ~73m incorrectly classifies some San Francisco homes as New York homes.

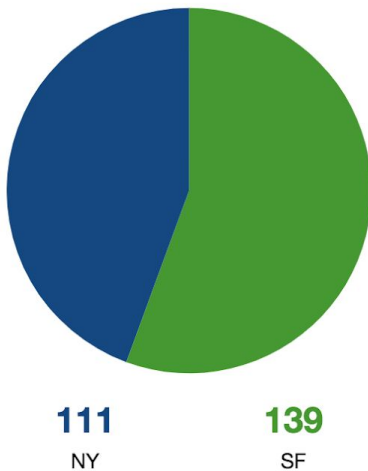




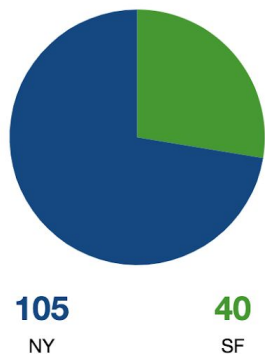
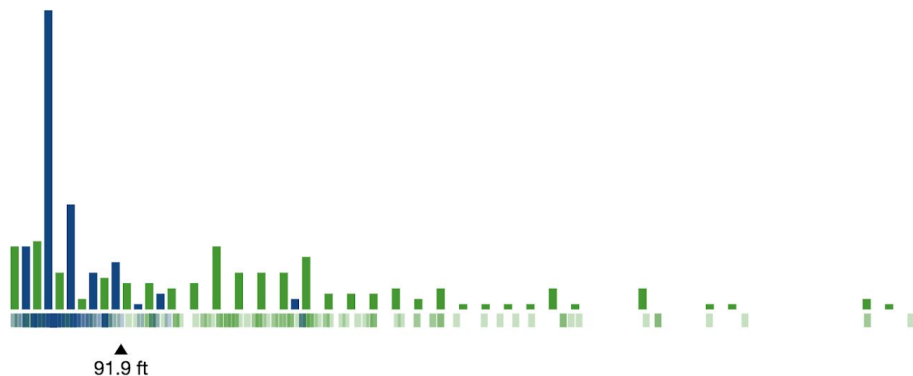
0  
NY

0  
SF

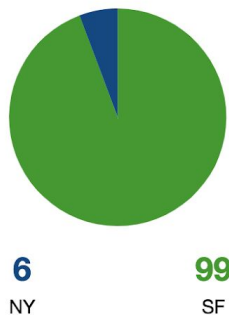
56  
% correct



If you split to capture *every* home in San Fran, you'll also get a bunch of New York homes (**false positives**)

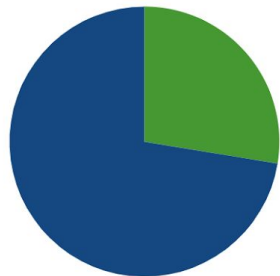


82  
% correct



## The best split

The best split point aims for branches that are as homogenous (pure) as possible



105

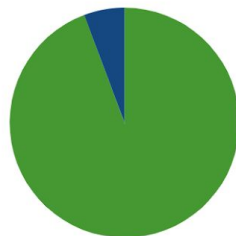
NY

40

SF

82

% correct

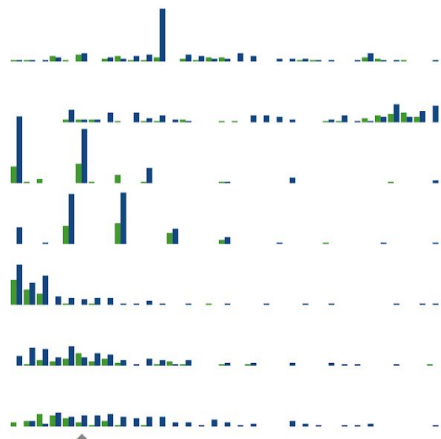


6

NY

99

SF



elevation

year built

bathrooms

bedrooms

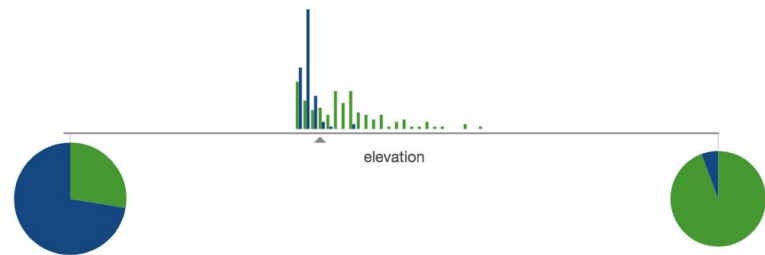
price

square feet

price per sqft

## Recursion

Additional split points are determined through repetition (**recursion**)



## Growing a tree

---

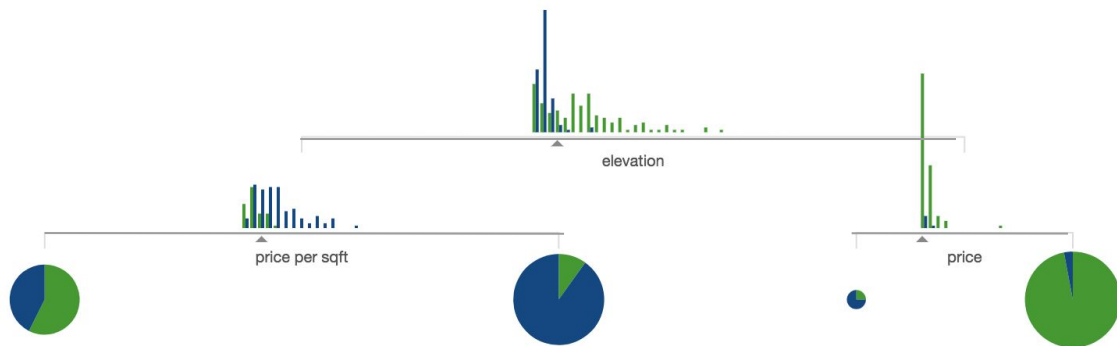
Additional forks add new information to improve **prediction accuracy**.

Accuracy: 82%

## Growing a tree

---

Additional forks add new information to improve **prediction accuracy**.

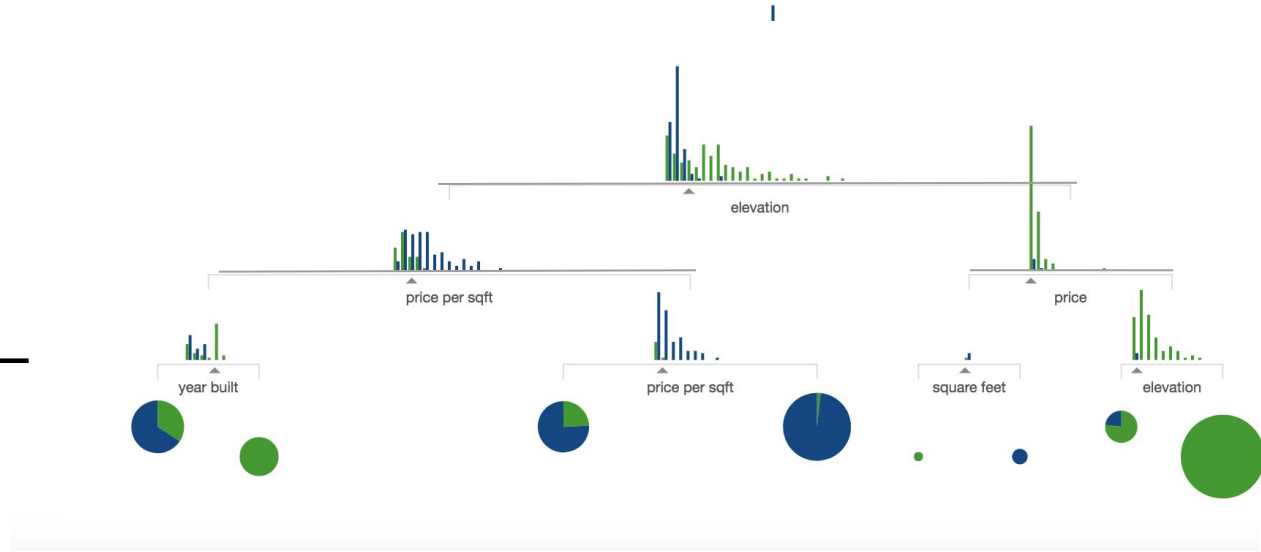


Accuracy: 86%

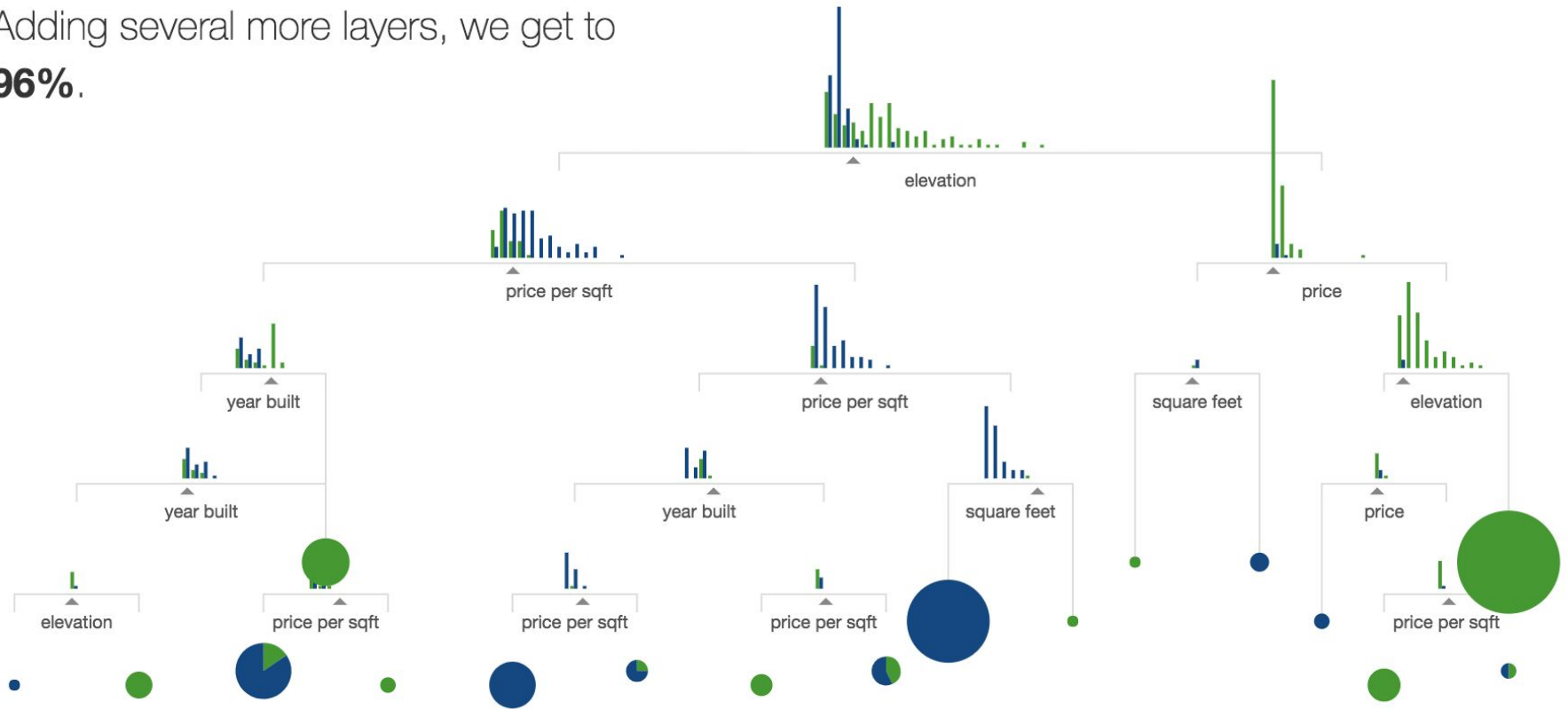


## Growing a tree

Additional forks add new information to improve **prediction accuracy**.

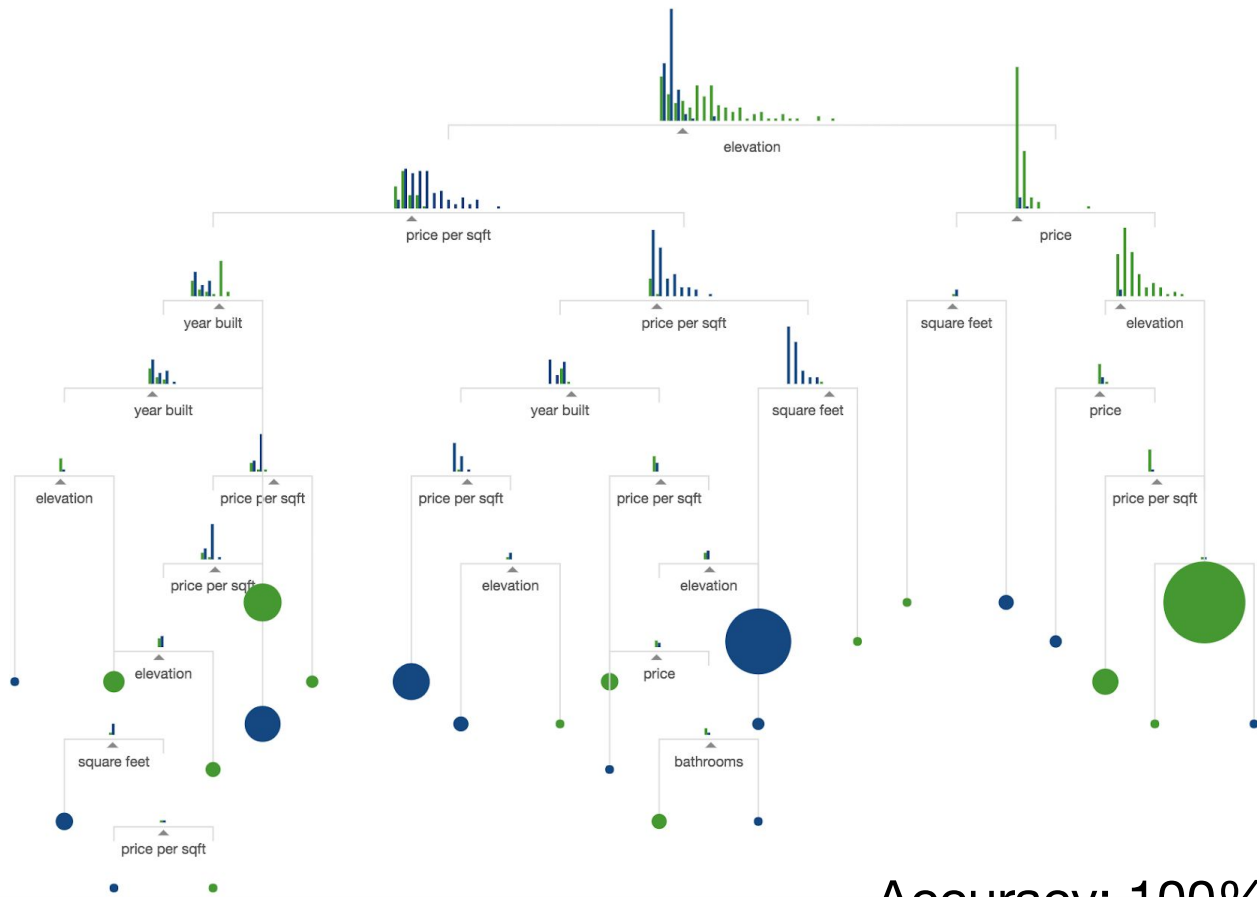


Adding several more layers, we get to  
**96%.**



Accuracy: 96%

It's possible to add  
branches until your  
model is **100%**  
**accurate.**

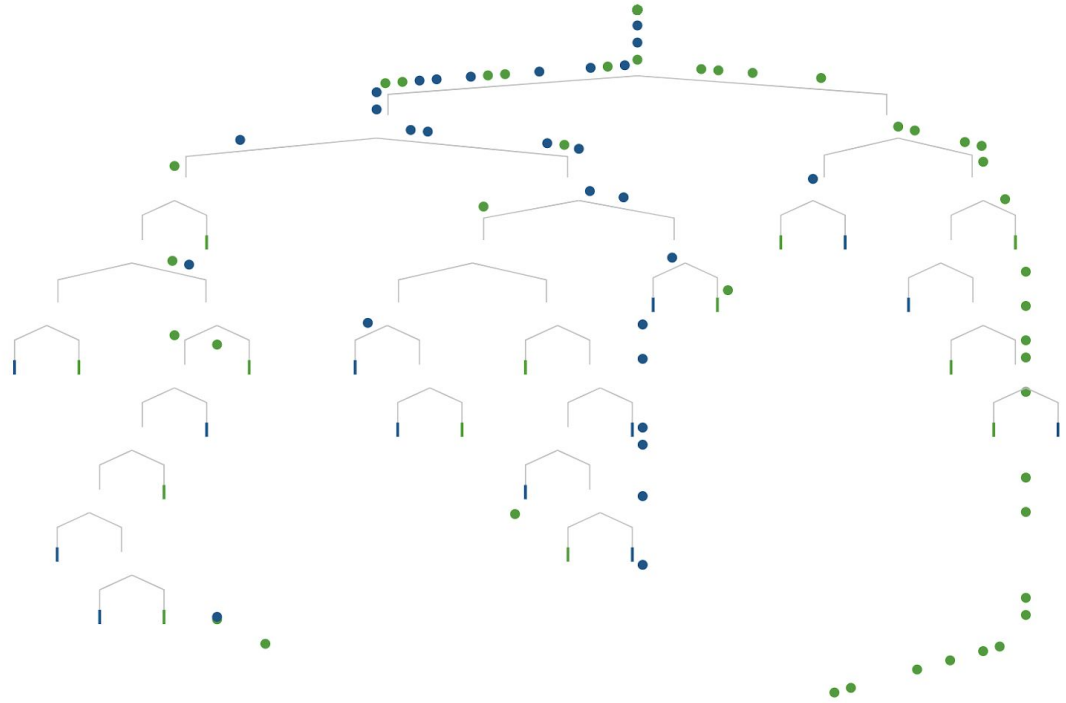


Accuracy: 100%

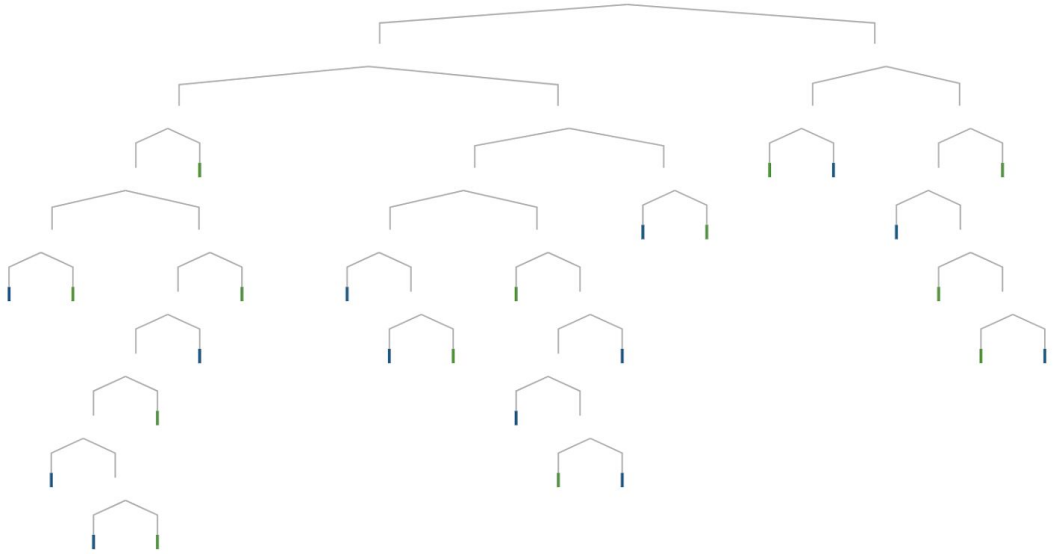
# Making predictions

The decision tree **model** can then predict which homes are in which city.

Here, we're using the **training data**.



Because our tree was trained on this data and we grew the tree to 100% accuracy, each house is perfectly sorted



111/111

Training Accuracy

100%

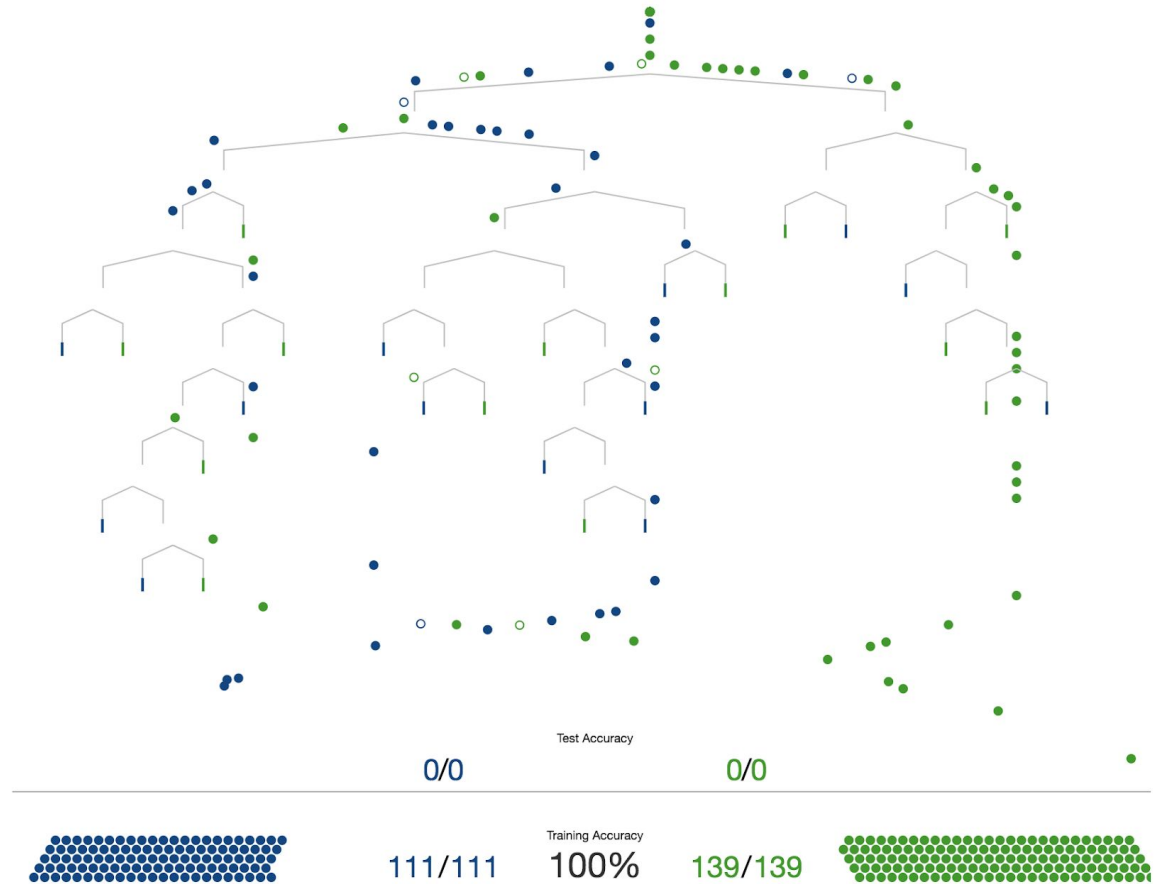
139/139

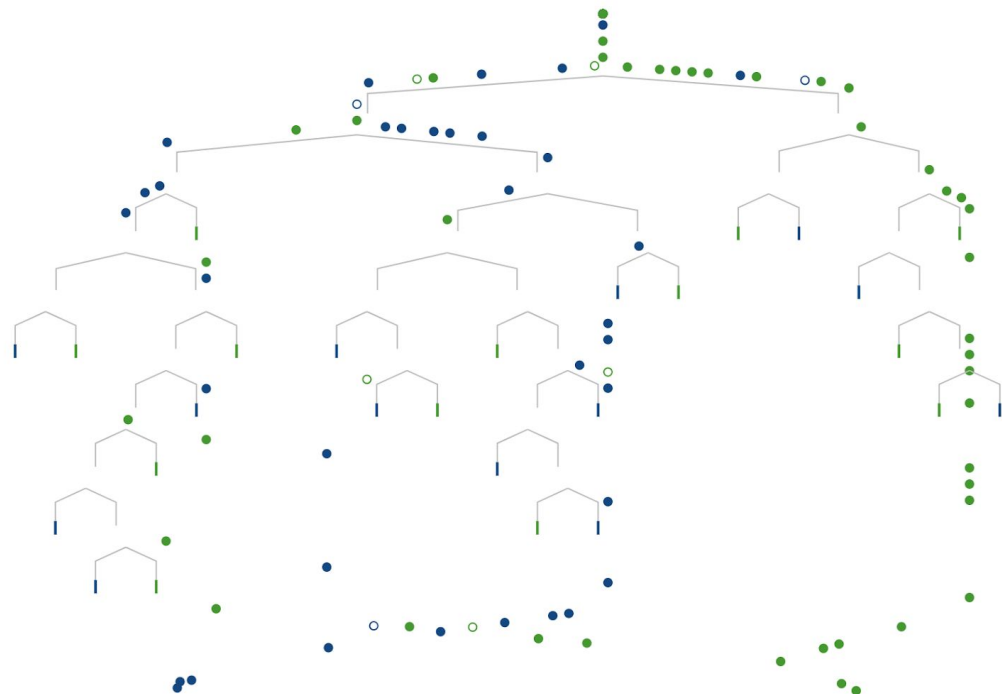


## Reality check

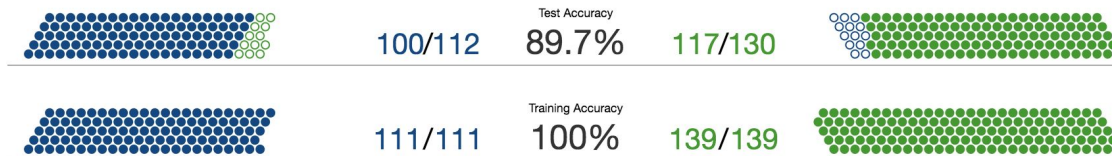
But...how does this tree  
on data that the model  
hasn't seen before?

The **test set** then  
makes it way through  
the decision tree.

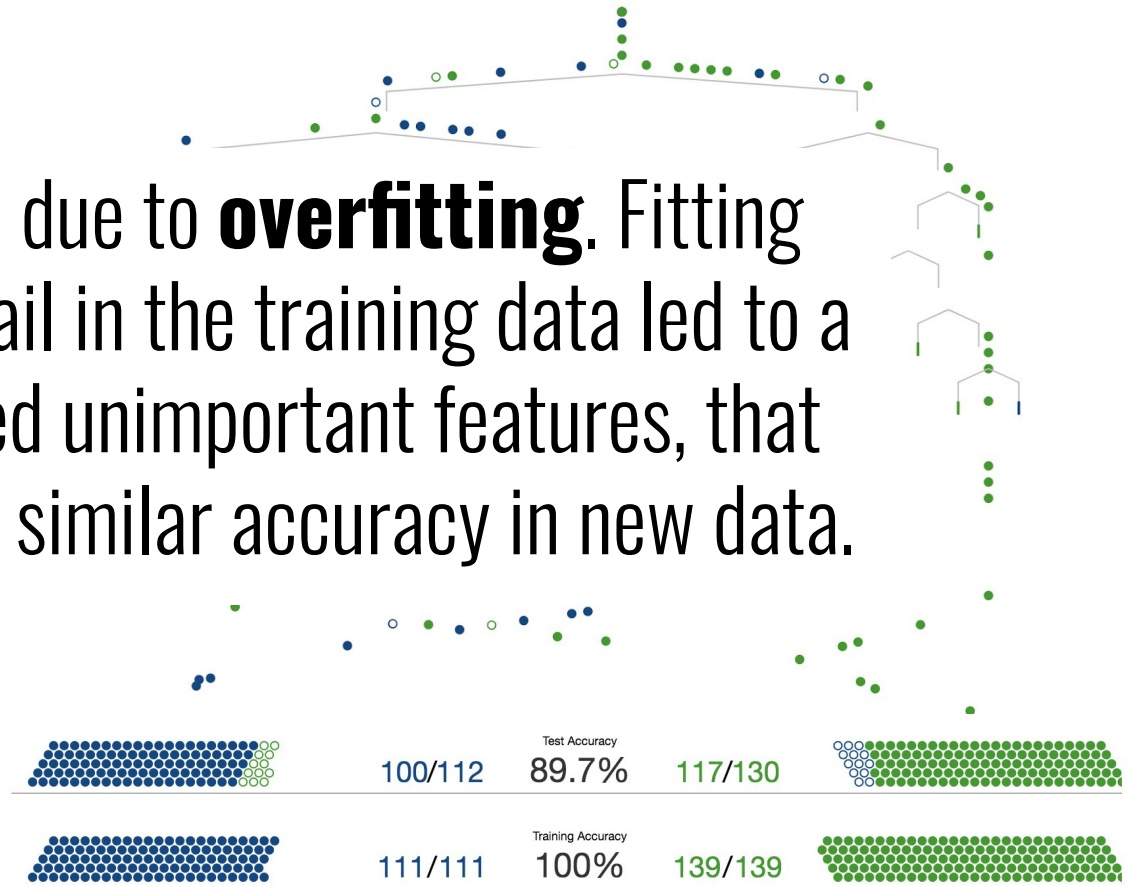




Ideally the tree should perform similarly on both known and unknown data



These errors are due to **overfitting**. Fitting every single detail in the training data led to a tree that modeled unimportant features, that did not allow for similar accuracy in new data.





# Recap

---

1. Machine learning identifies patterns using **statistical learning** and computers by unearthing **boundaries** in data sets. You can use it to make predictions.
2. One method for making predictions is called a decision trees, which uses a series of if-then statements to identify boundaries and define patterns in the data.
3. **Overfitting** happens when some boundaries are based on on *distinctions that don't make a difference*. You can see if a model overfits by having test data flow through the model.