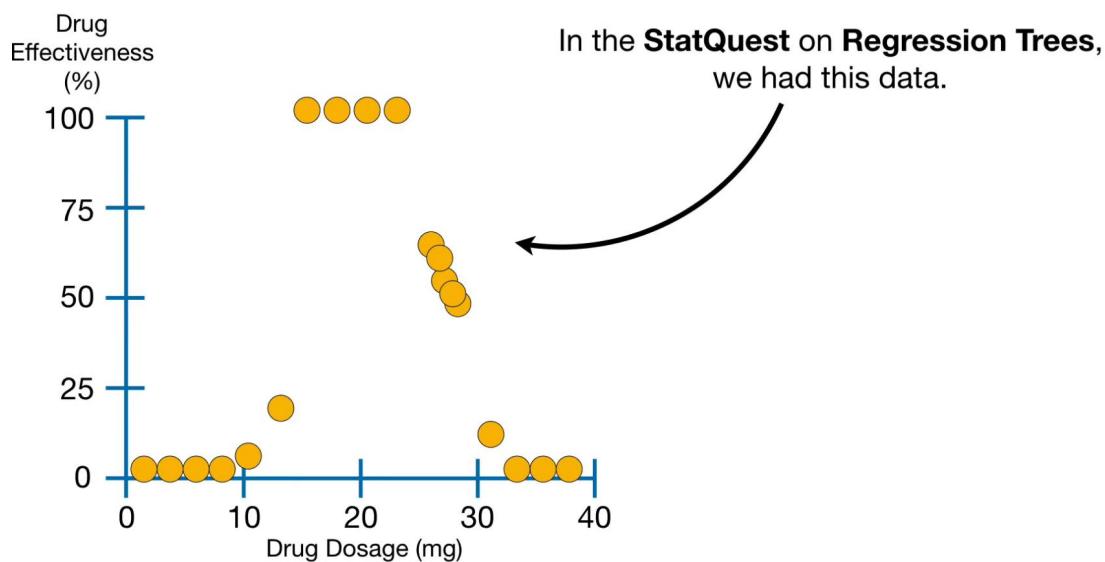


How to Prune Regression Trees

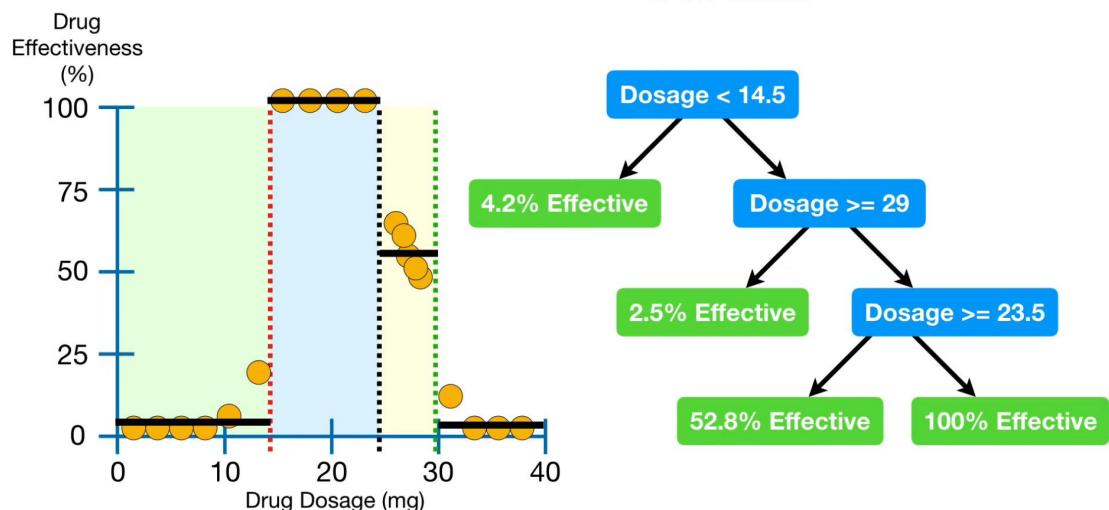
There are several methods for pruning **Regression Trees**.

The one we'll talk about in this '**Quest**' is called **Cost Complexity Pruning** aka **Weakest Link Pruning**.

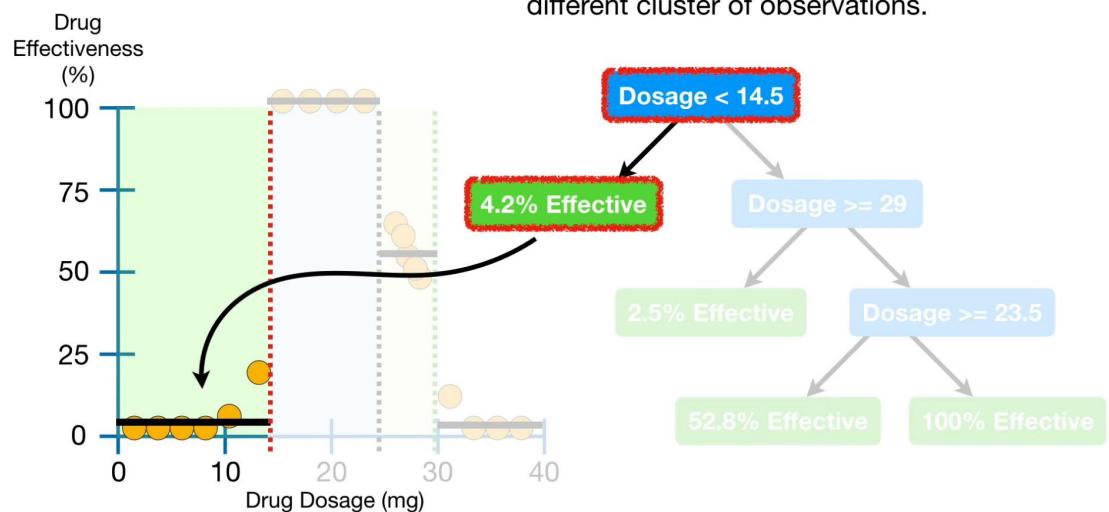
We'll start by giving a general overview of how **Cost Complexity Pruning** works, and then we'll describe how it is used to build **Regression Trees**.



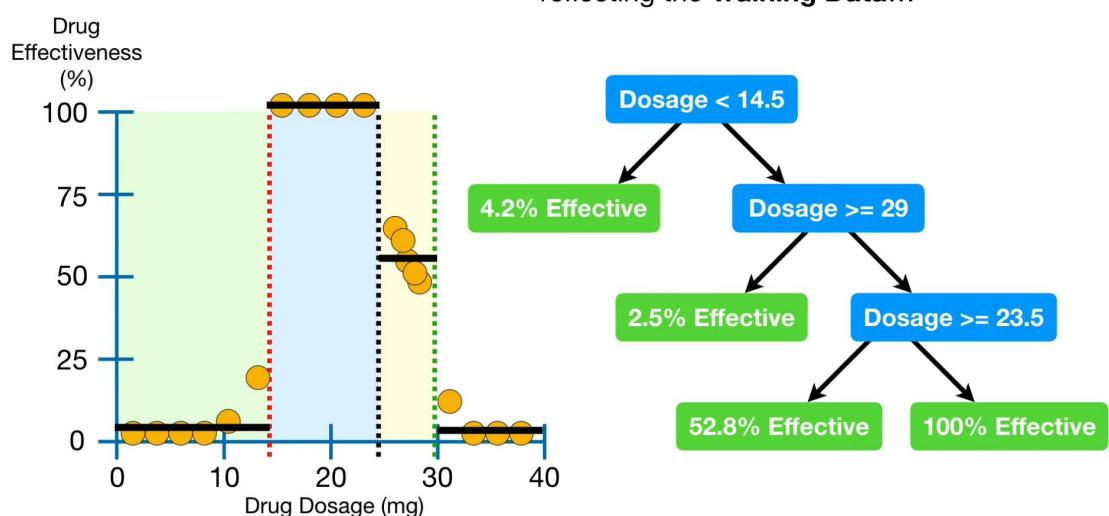
We then fit a **Regression Tree**
to the data...



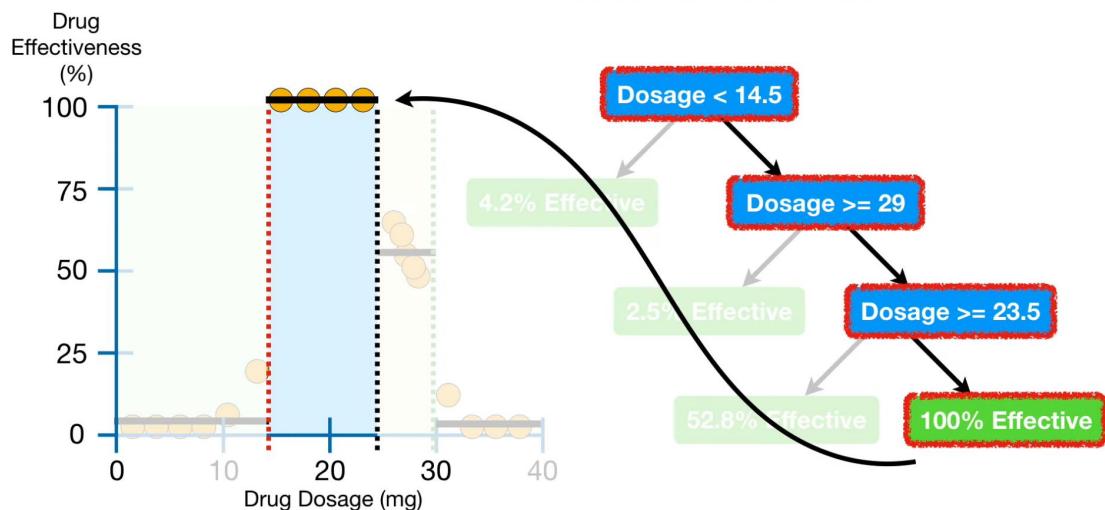
...and each leaf corresponded to the
average **Drug Effectiveness** from a
different cluster of observations.



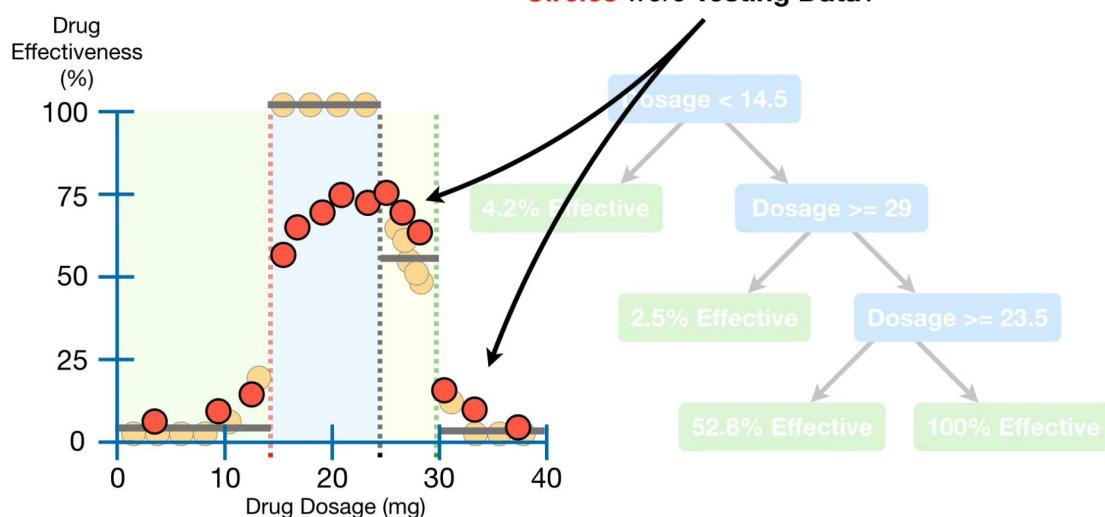
This tree does a pretty good job
reflecting the **Training Data**...



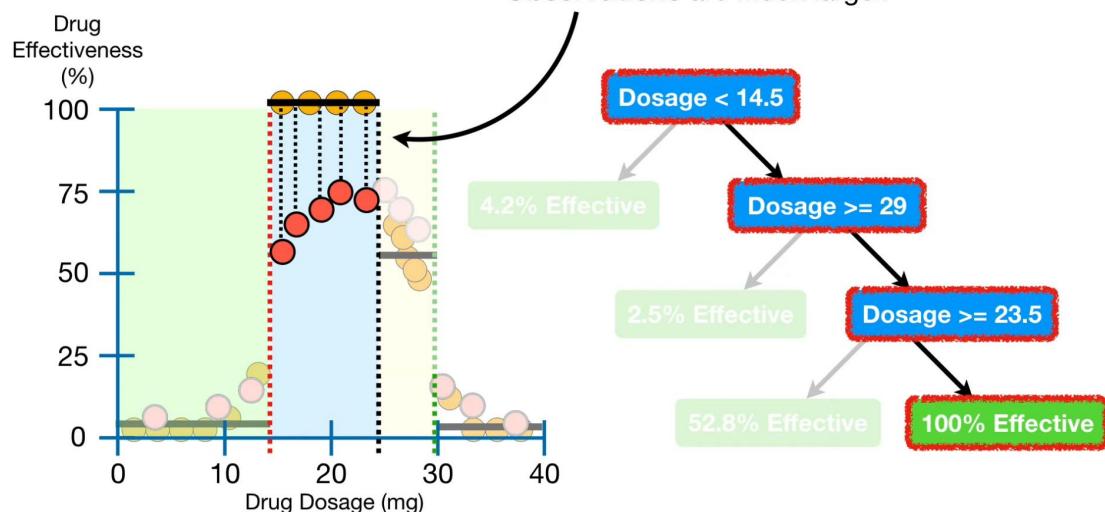
...because each leaf represents a value
that is close to the data.

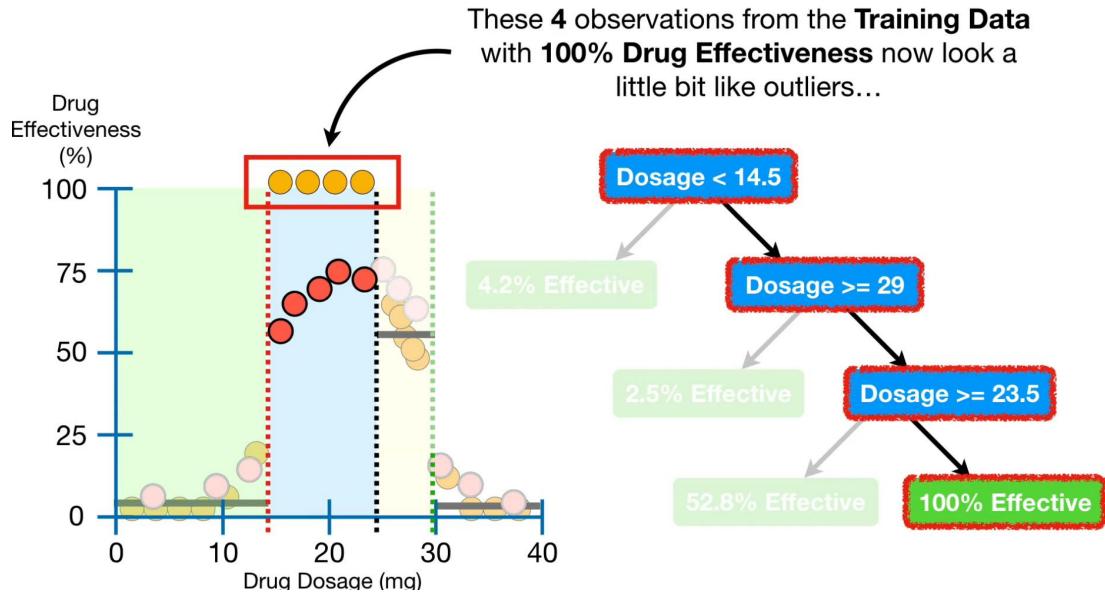


However, what if these Red Circles were Testing Data?

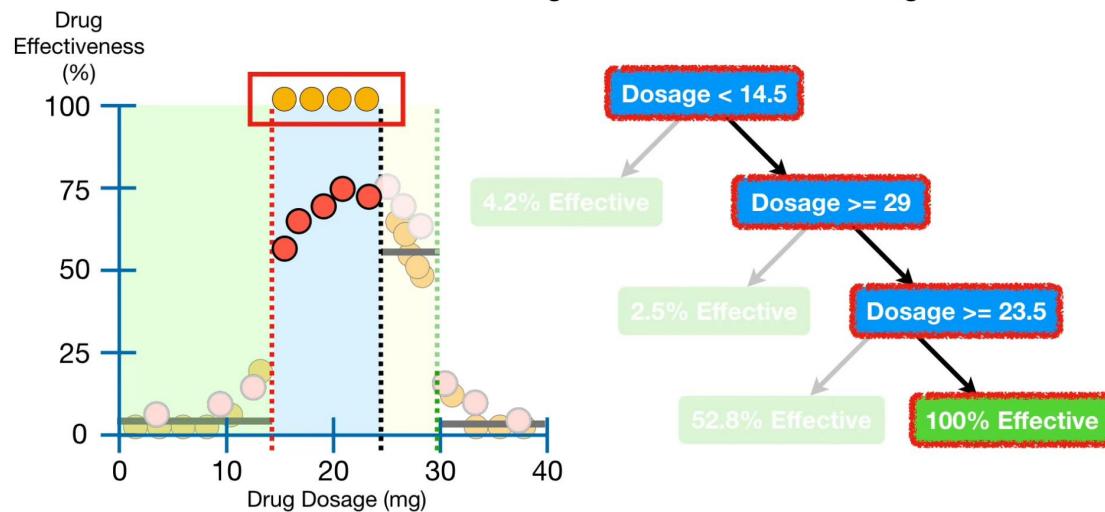


And the Residuals for these Observations are much larger.

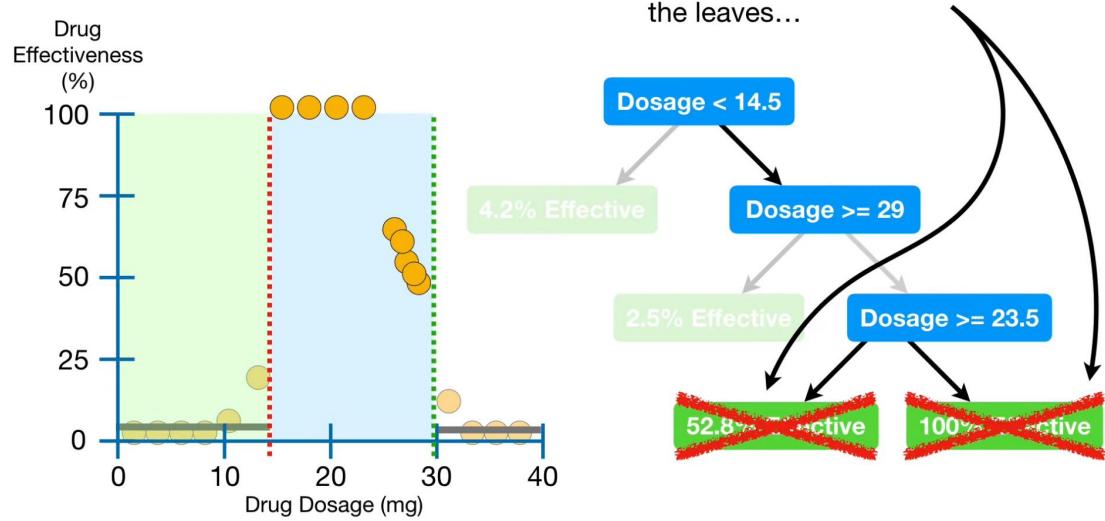




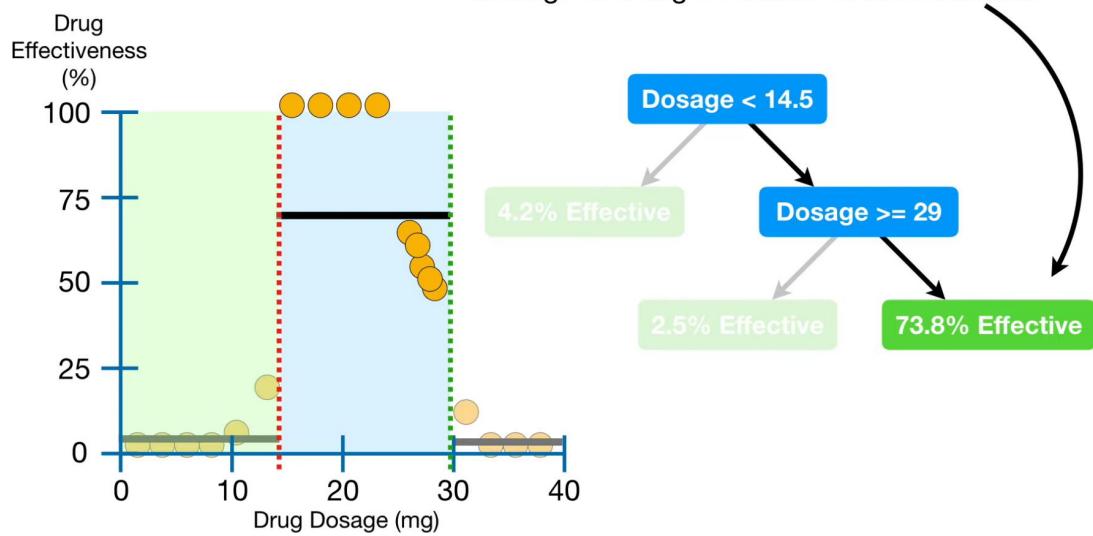
...and that means that maybe we **Overfit** the **Regression Tree** to the **Training Data**.

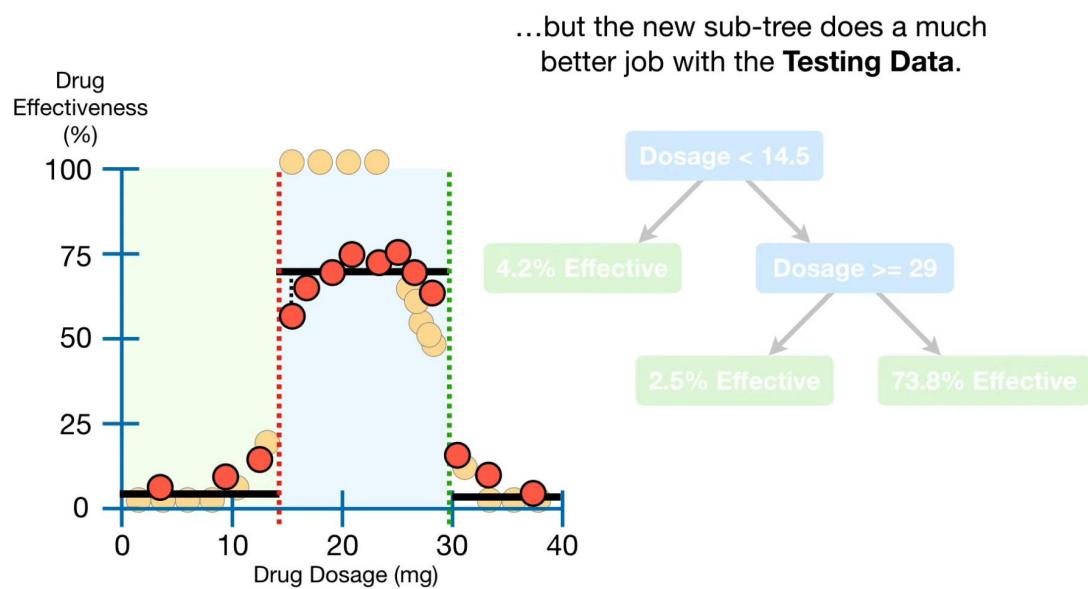
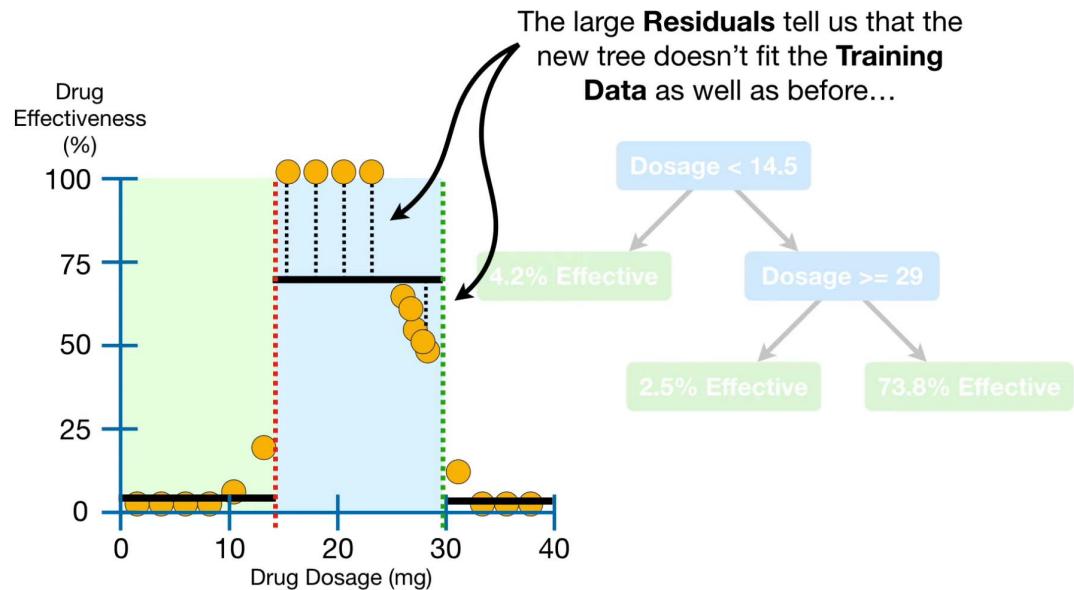


One way to prevent **Overfitting a Regression Tree** to the **Training Data** is to remove some of the leaves...

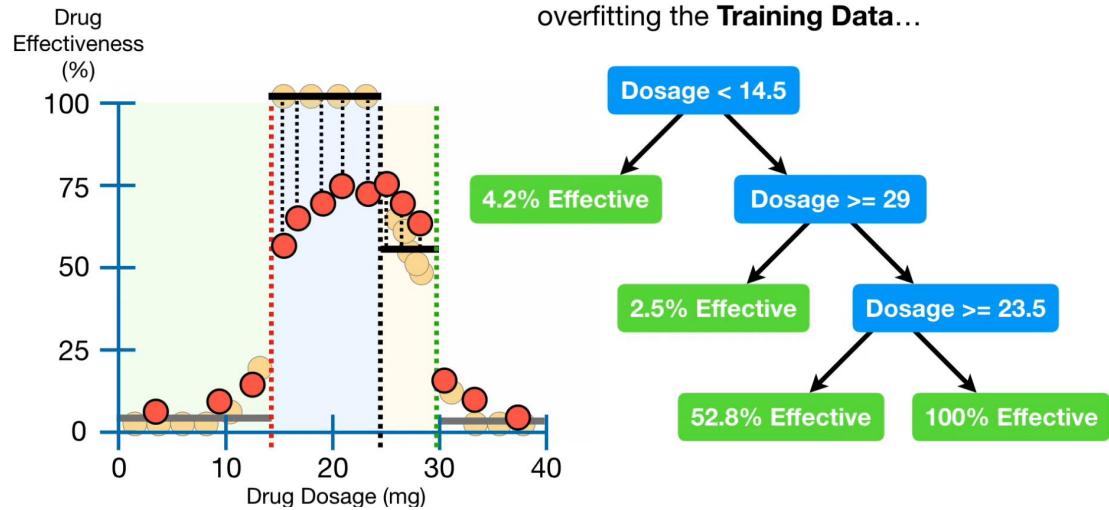


...and replace the split with a leaf that is the average of a larger number of observations.

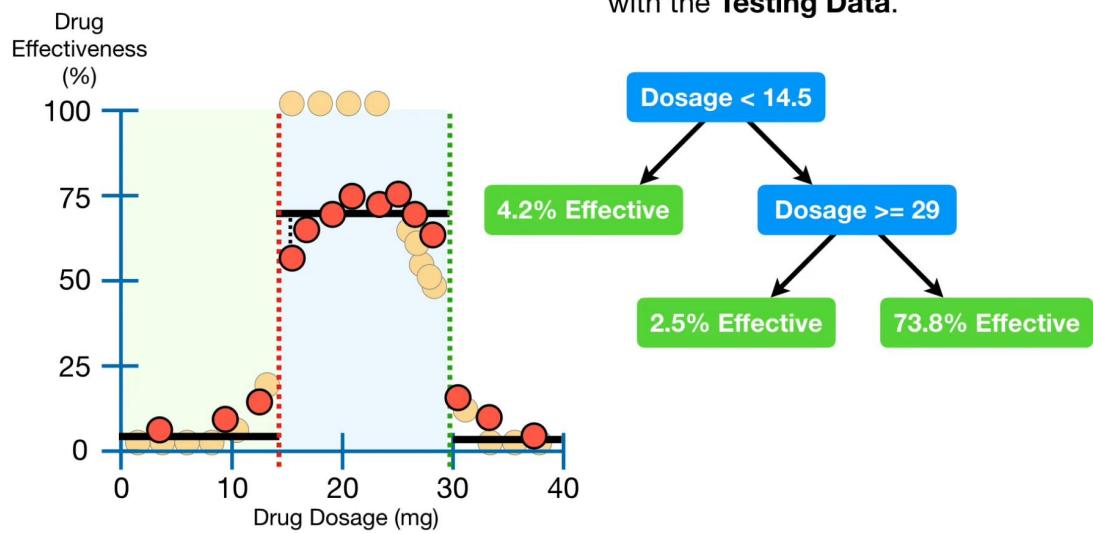


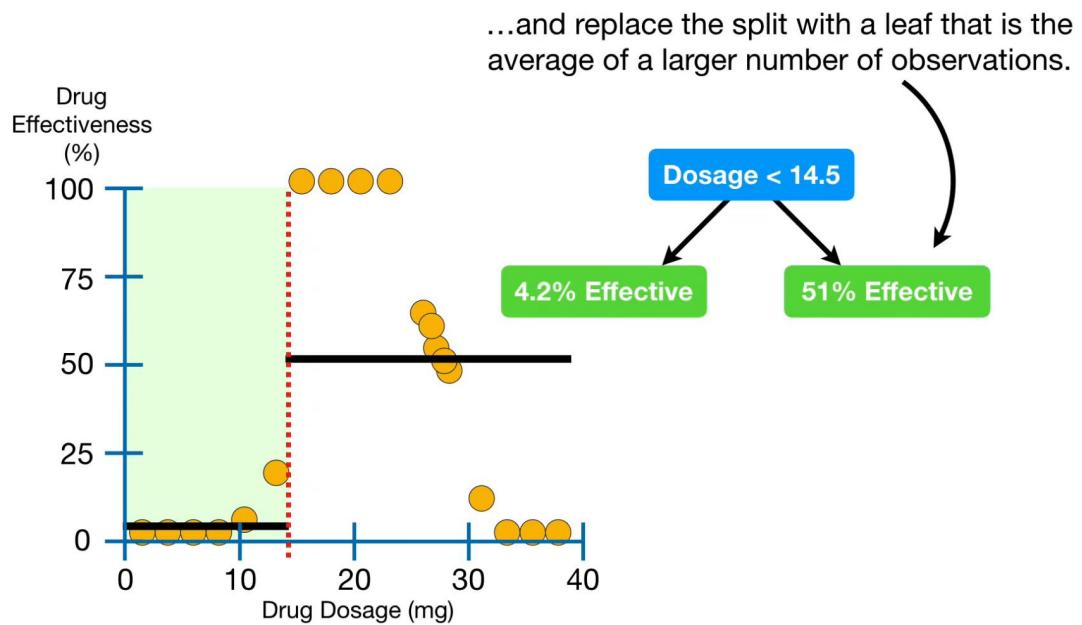
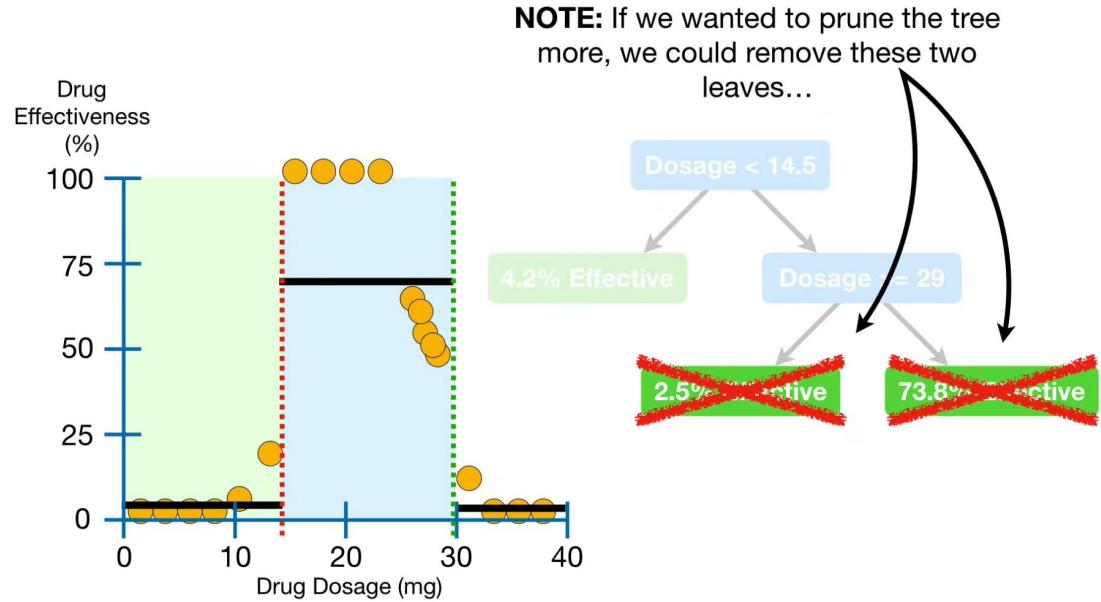


Thus, the main idea behind pruning a **Regression Tree** is to prevent overfitting the **Training Data**...

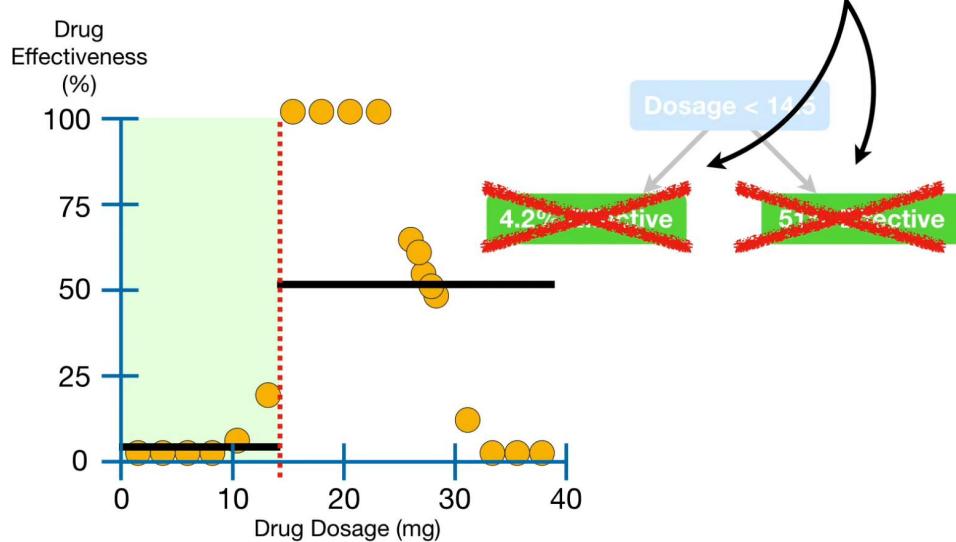


...so that the tree will do a better job with the **Testing Data**.

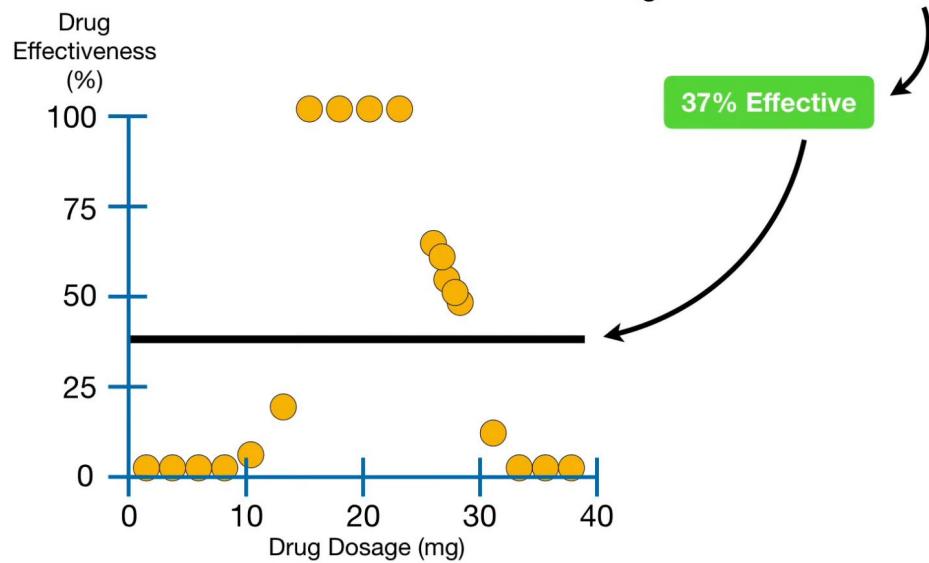




And we could then remove these two leaves...

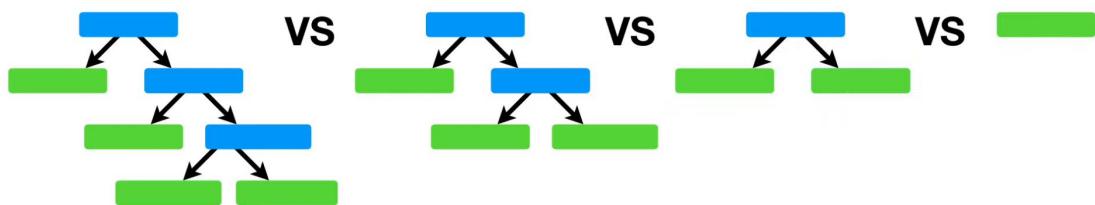


...and replace the split with a leaf that is the average of all of the observations.



So the question is:

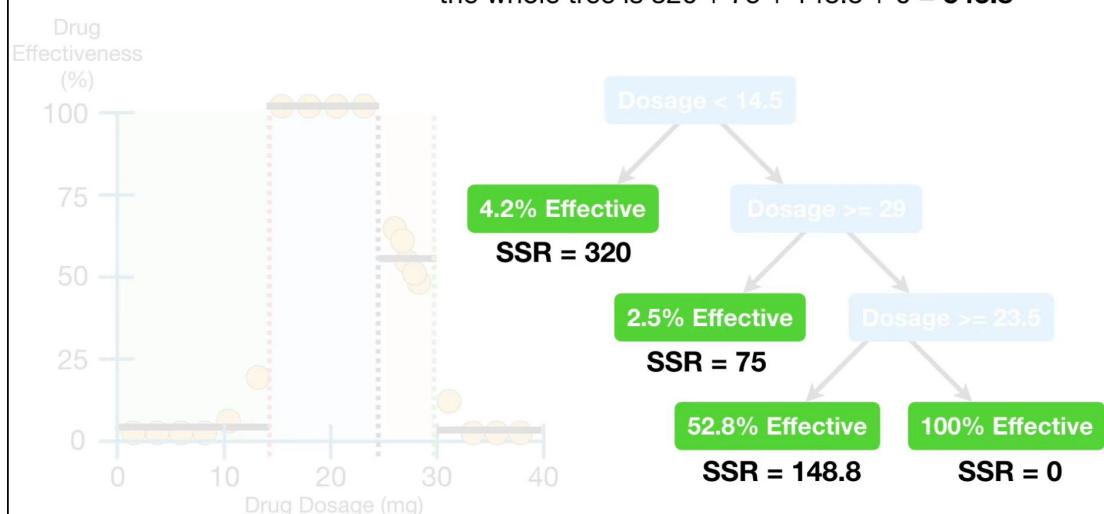
“How do we decide which tree to use?”



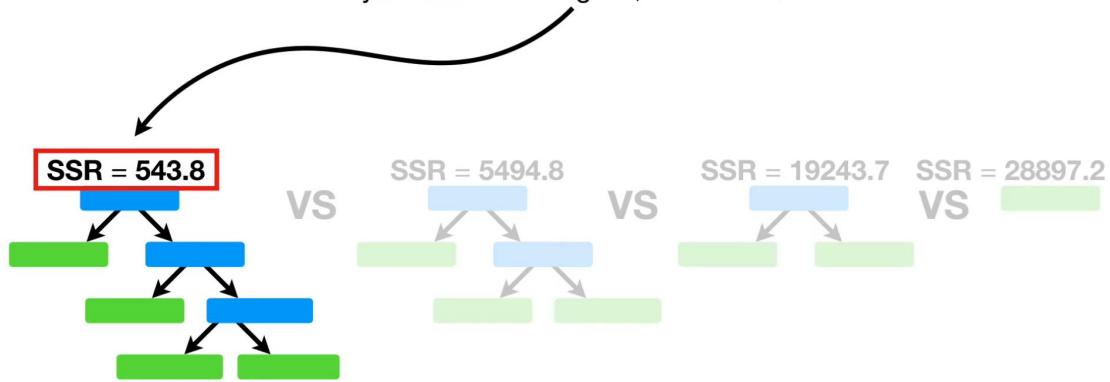
In this **StatQuest**, we will answer that question with **Cost Complexity Pruning**.

The first step in **Cost Complexity Pruning** is to calculate the **Sum of Squared Residuals** for each tree.

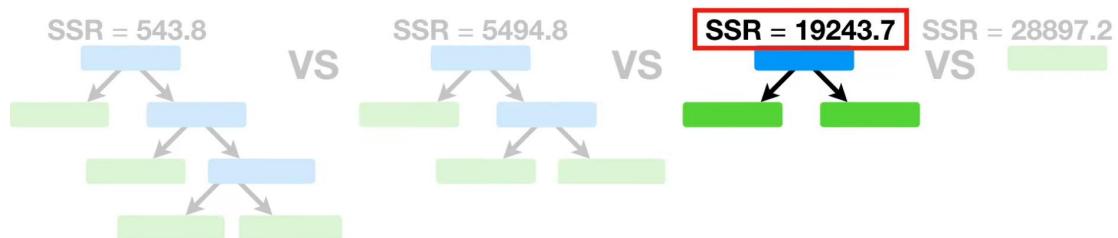
Thus, the total **Sum of Squared Residuals** (SSR) for the whole tree is $320 + 75 + 148.8 + 0 = 543.8$



NOTE: The **Sum of Squared Residuals** is relatively small for the original, full sized tree...



...but each time we remove a leaf, the **Sum of Squared Residuals** gets larger and larger.



However, we knew that was going to happen because the whole idea was for the pruned trees to **not** fit the **Training Data** as well as the full sized tree.

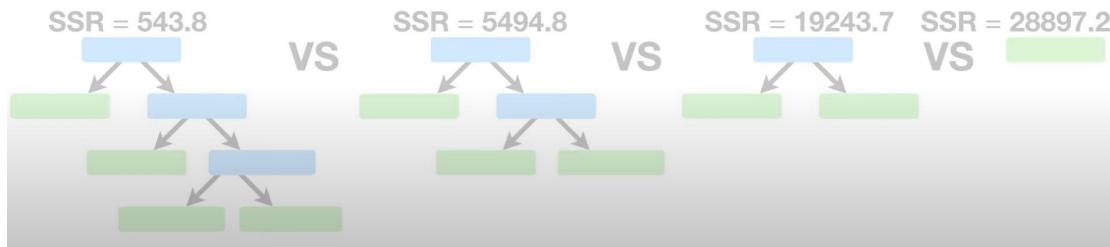
So how do we compare these trees?

Weakest Link Pruning works by calculating a **Tree Score**...

...that is based on the **Sum of Squared Residuals (SSR)** for the tree or sub-tree...

...and a **Tree Complexity Penalty** that is a function of the number of leaves, or **Terminal nodes**, in the tree or sub-tree.)

$$\text{Tree Score} = \text{SSR} + \alpha T$$



The **Tree Complexity Penalty** compensates for the difference in the number of leaves.

$$\text{Tree Score} = \text{SSR} + \alpha T$$

NOTE: α (alpha) is a tuning parameter that we find using **Cross Validation** and we'll talk more about it in a bit.

$$\text{Tree Score} = \text{SSR} + \alpha T$$

For now, let's let $a = 10,000$.

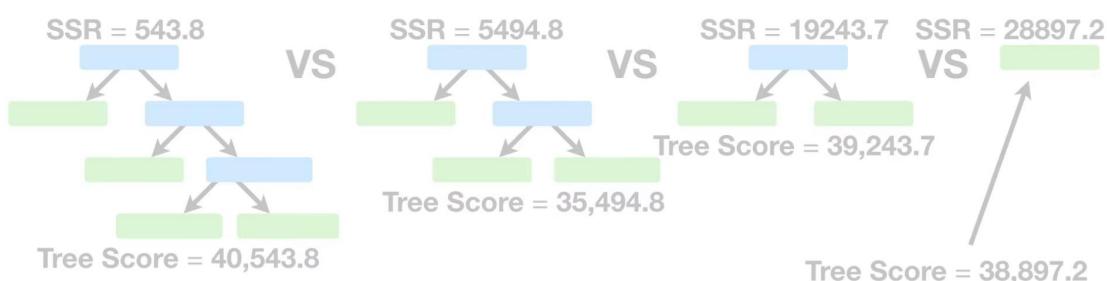


$$\text{Tree Score} = \text{SSR} + 10,000 \times T$$

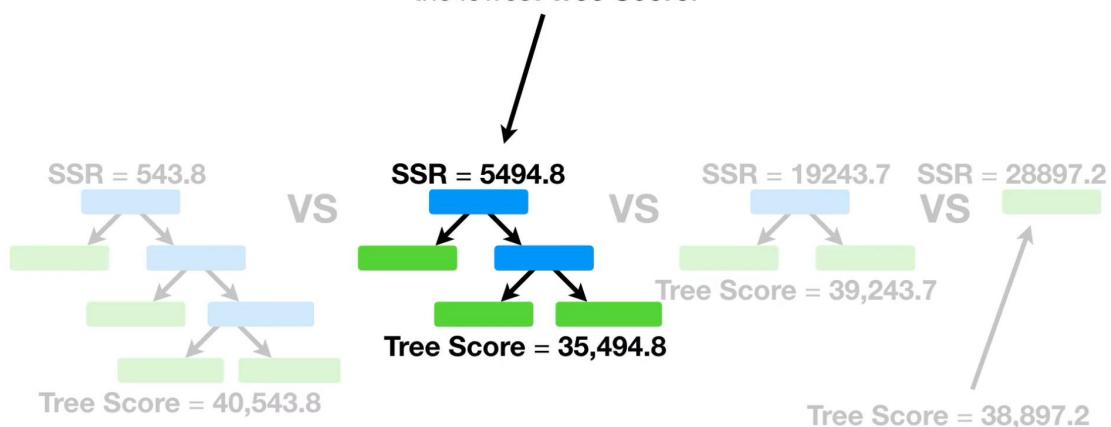
Now let's calculate the **Tree Score** for each tree.

Thus, the more leaves,
the larger the penalty.

$$\text{Tree Score} = \text{SSR} + 10,000 \times T$$

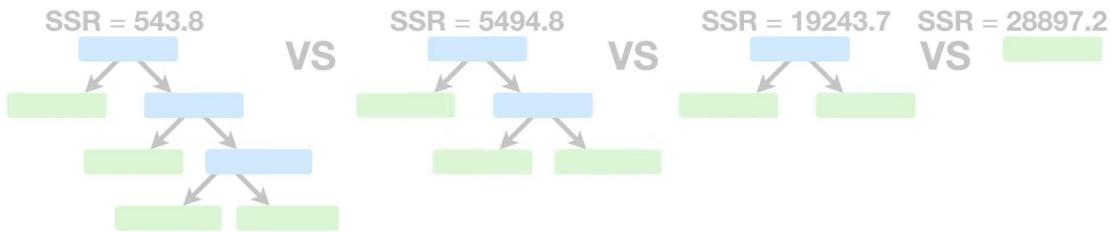


...we pick this sub-tree because it has
the lowest **Tree Score**.



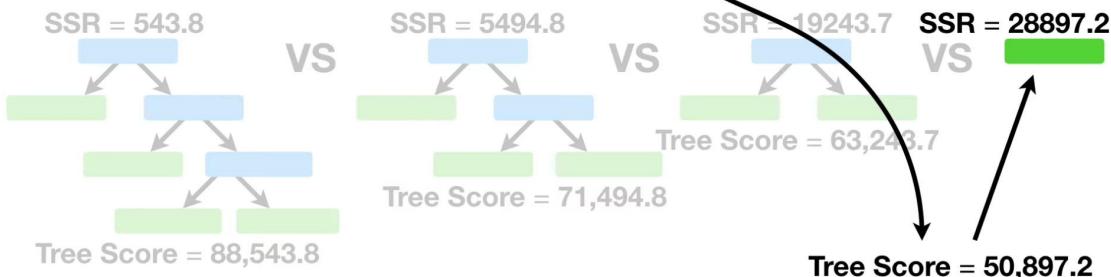
NOTE: If we set $\alpha = 22,000 \dots$

$$\text{Tree Score} = \text{SSR} + \alpha T$$



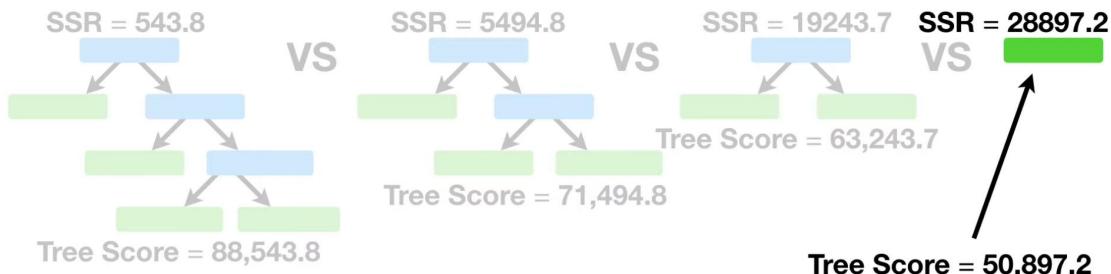
...and calculate the **Tree Scores**...

$$\text{Tree Score} = \text{SSR} + 22,000 \times T$$



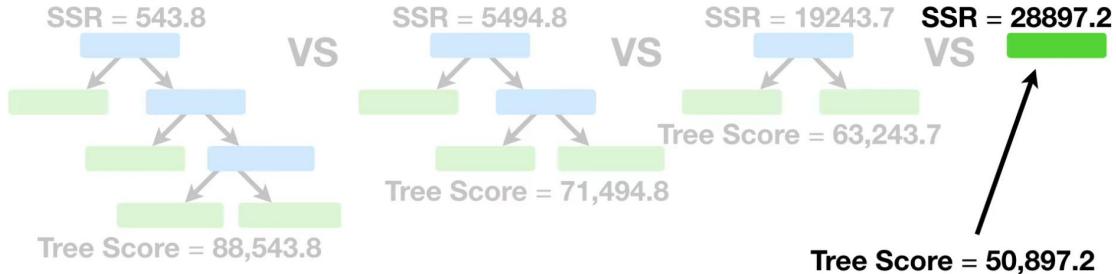
...then we would use the sub-tree with only one leaf because it has lowest **Tree Score**.

$$\text{Tree Score} = \text{SSR} + 22,000 \times T$$



Thus, the value for a makes a difference in our choice of sub-tree.

$$\text{Tree Score} = \text{SSR} + 22,000 \times T$$

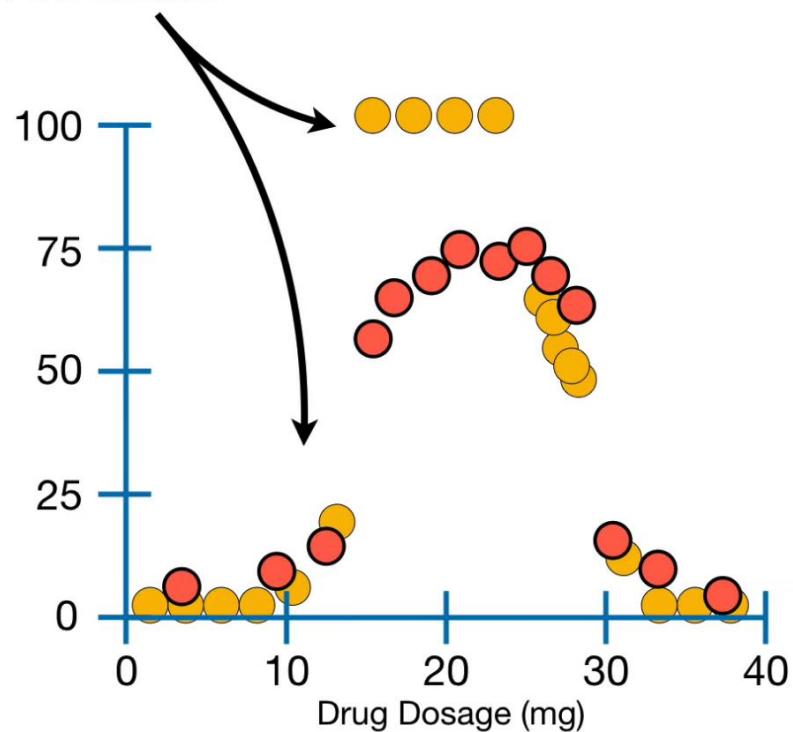


So let's talk about how build a pruned regression tree...

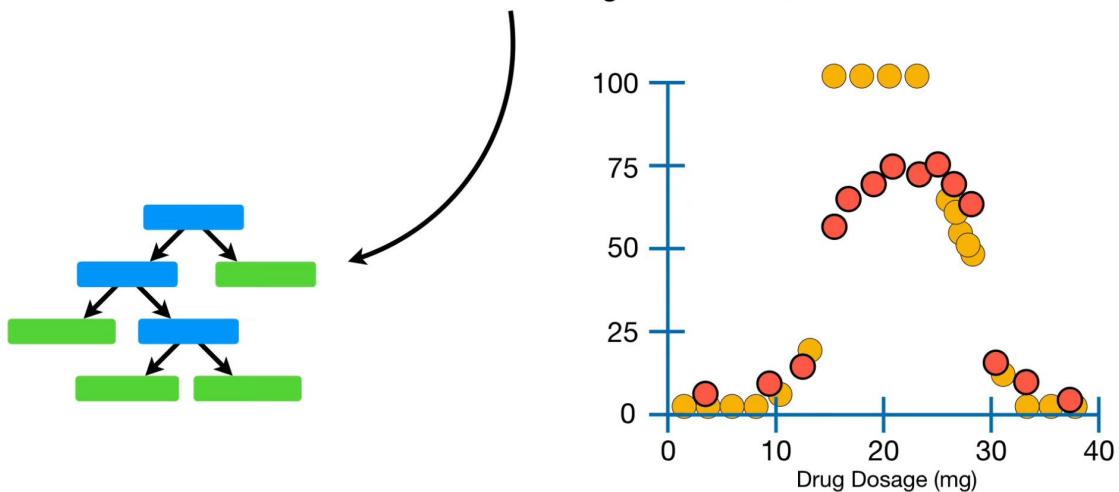
...and how to find the best value for a .

$$\text{Tree Score} = \text{SSR} + aT$$

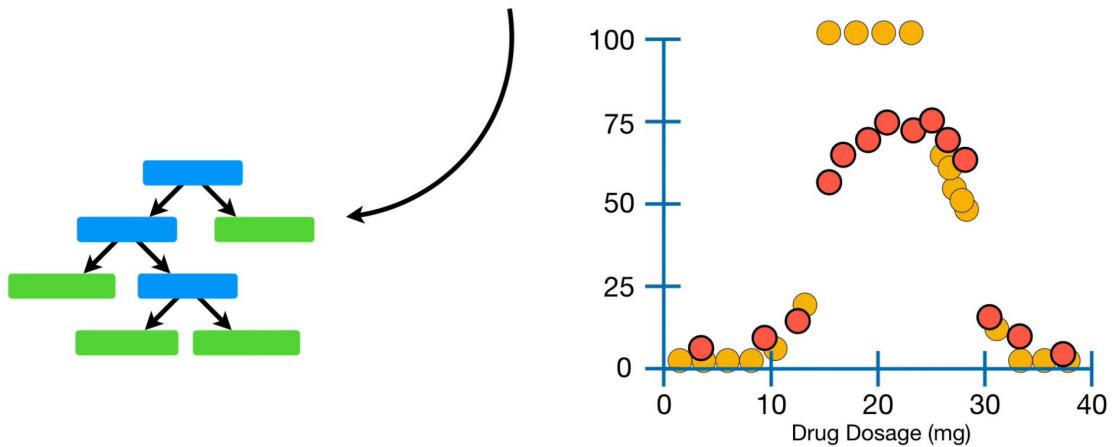
First, using ***all*** of the data...



...build a full sized **Regression Tree**.



NOTE: This full sized tree is different than before because it was fit to **all** of the data, not just the **Training Data**.



ALSO NOTE: This full sized tree has the lowest **Tree Score** when $a = 0$.

$$\text{Tree Score} = \text{SSR} + aT$$



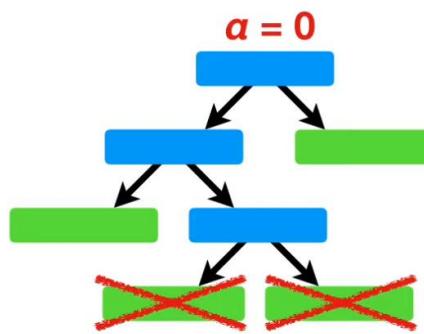
So let's put $a = 0$ here, to remind us that this tree has the lowest **Tree Score** when $a = 0$.

Tree Score = SSR



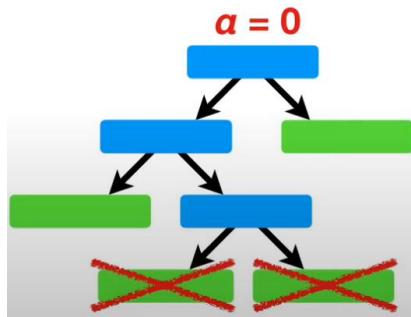
Now we increase a until pruning leaves will give us a lower **Tree Score**.

Tree Score = SSR + aT



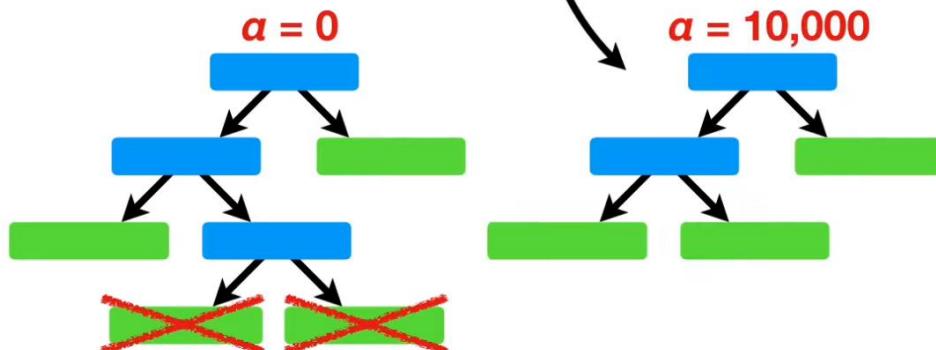
In this case, when $a = 10,000$, we'll get a lower **Tree Score** if we remove these leaves...

$$\text{Tree Score} = \text{SSR} + aT$$



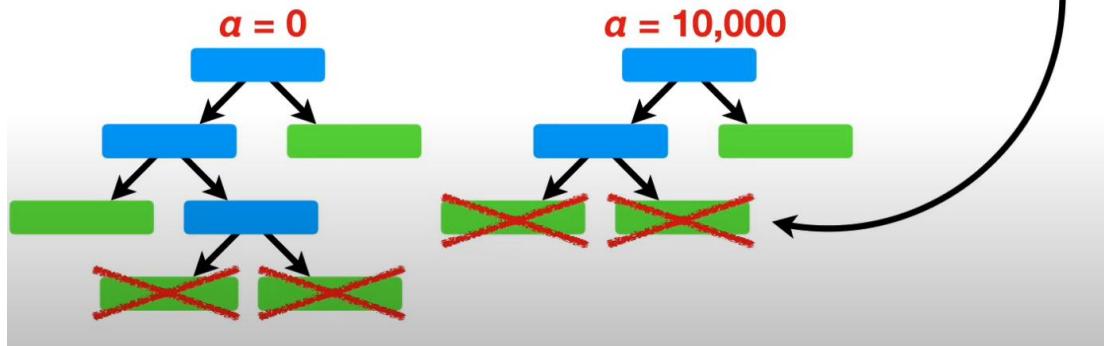
...and use this sub-tree.

$$\text{Tree Score} = \text{SSR} + aT$$



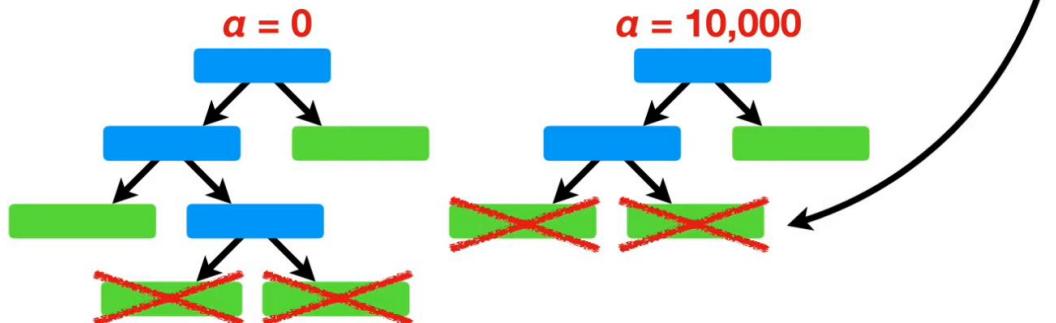
Now we increase α again until pruning leaves will give us a lower **Tree Score**.

$$\text{Tree Score} = \text{SSR} + \alpha T$$



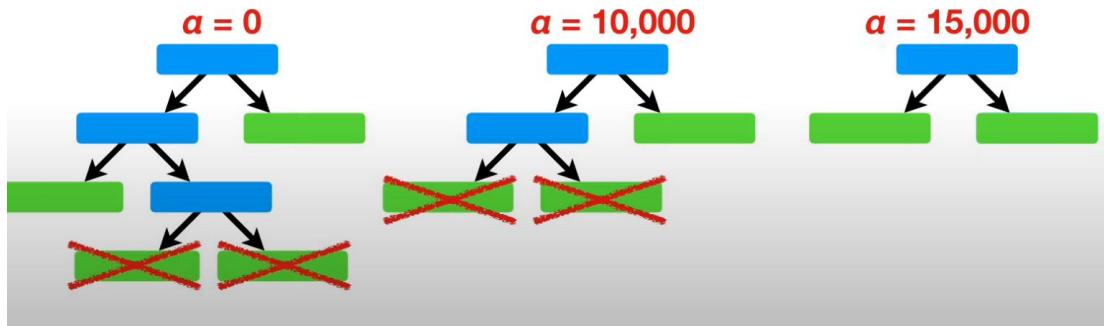
In this case when $\alpha = 15,000$, we will get a lower **Tree Score** if we remove these leaves...

$$\text{Tree Score} = \text{SSR} + \alpha T$$



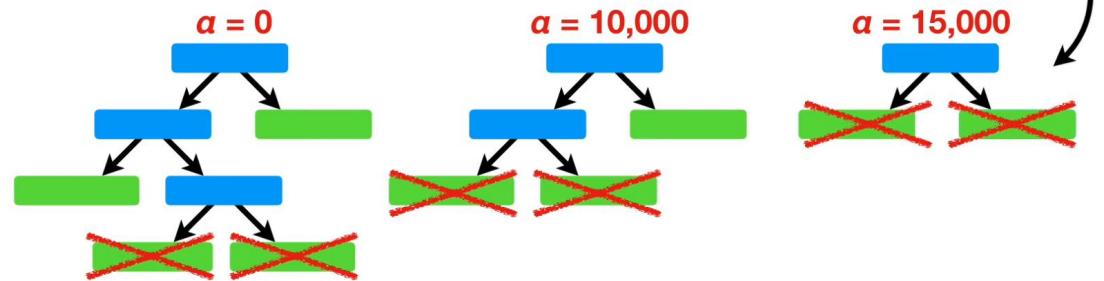
...and use this sub-tree instead.

$$\text{Tree Score} = \text{SSR} + aT$$



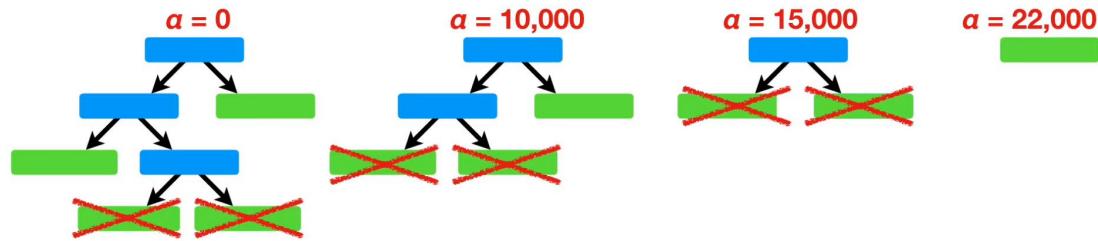
And when $a = 22,000$, we will get a lower **Tree Score** if we remove these leaves...

$$\text{Tree Score} = \text{SSR} + aT$$

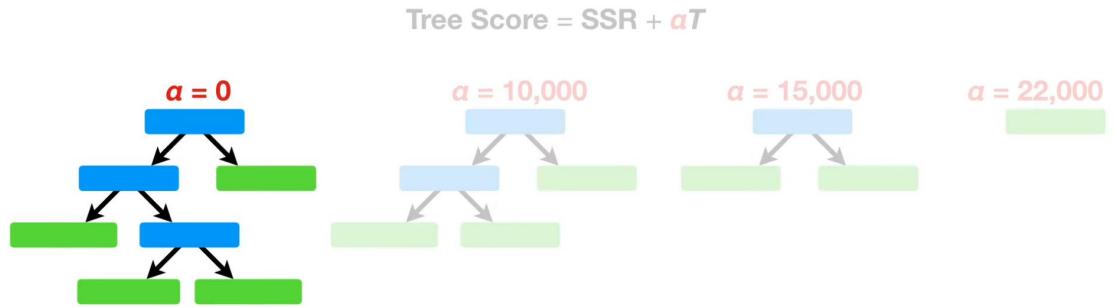


...and use this sub-tree instead.

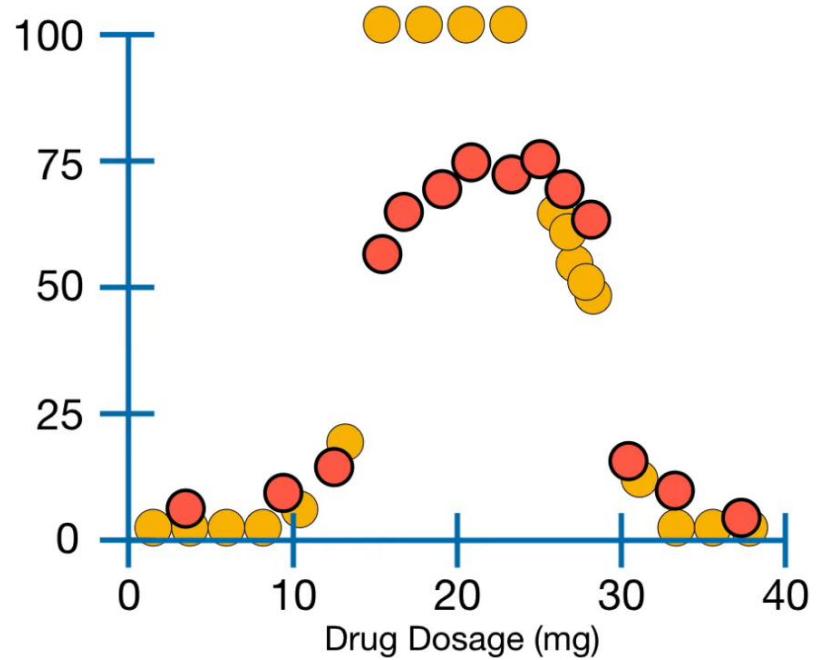
$$\text{Tree Score} = \text{SSR} + aT$$



In the end, different values for α give us a sequence of trees, from full sized to just a leaf.

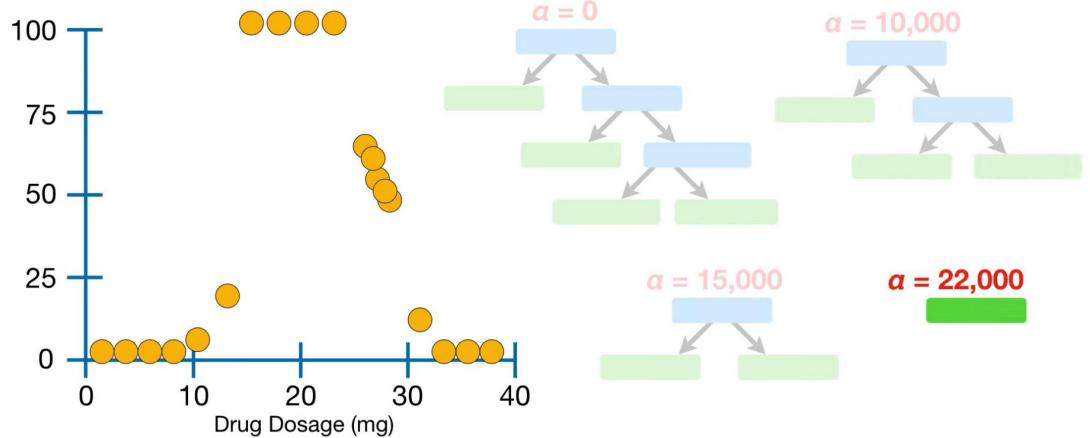


Now go back to the full dataset...



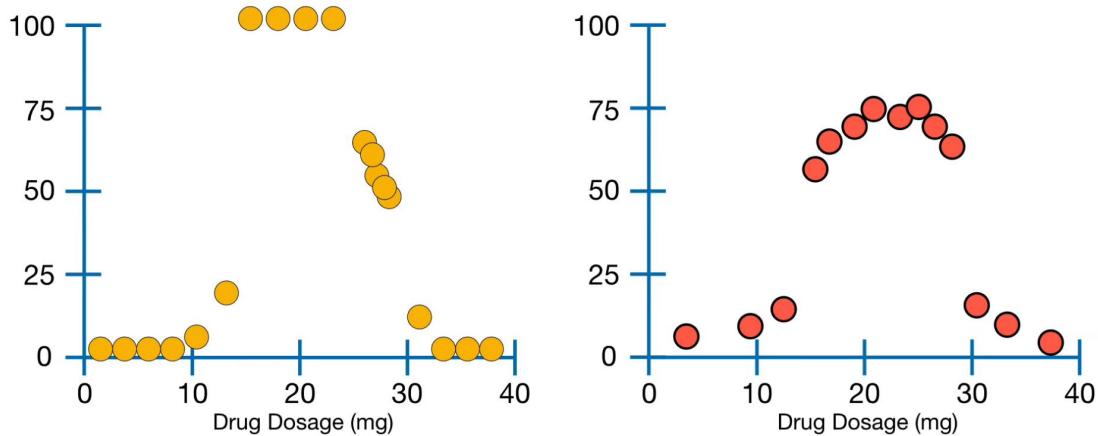
And just using the
Training Data...

...use the α values we found before to
build a full tree and a sequence of sub-
trees that minimize the **Tree Score**.



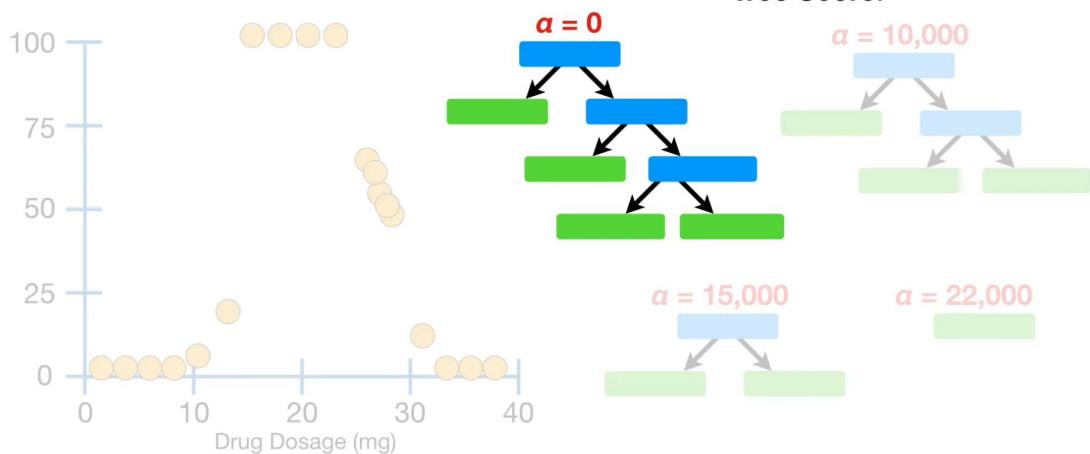
...and divide it into **Training...**

...and **Testing Datasets**.



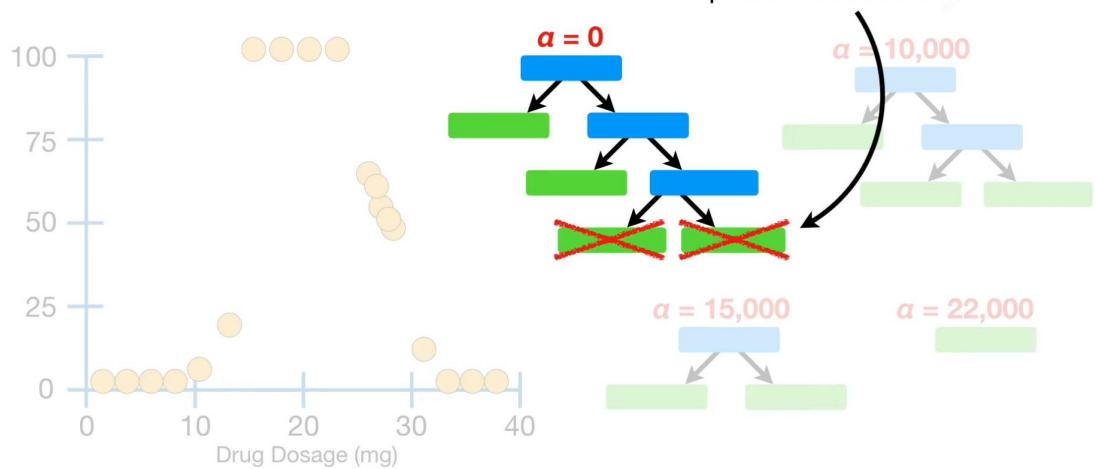
$$\text{Tree Score} = \text{SSR} + \alpha T$$

In other words, when $\alpha = 0$,
we build a full sized tree,
since it will have the lowest
Tree Score.



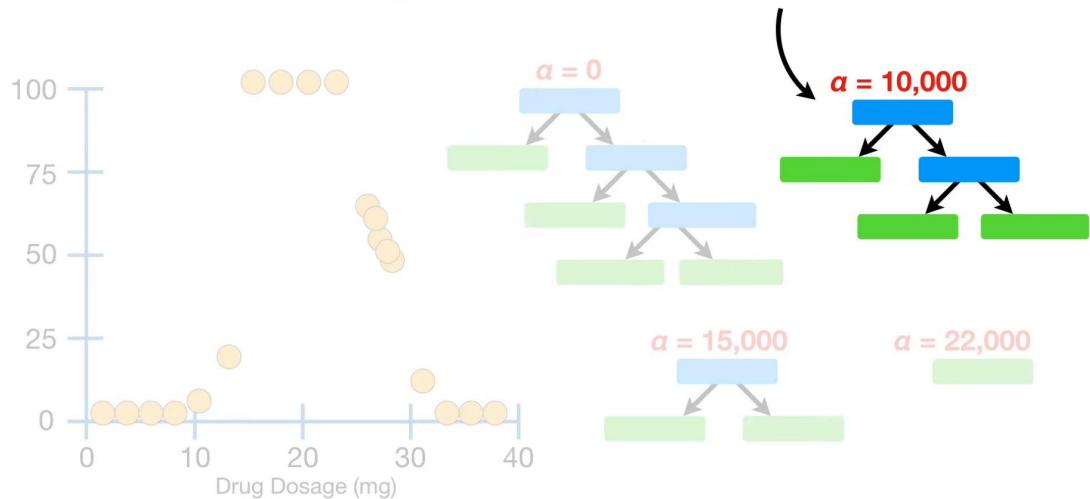
$$\text{Tree Score} = \text{SSR} + \alpha T$$

However, when $\alpha = 10,000$, we will get a lower **Tree Score** if we prune these leaves...



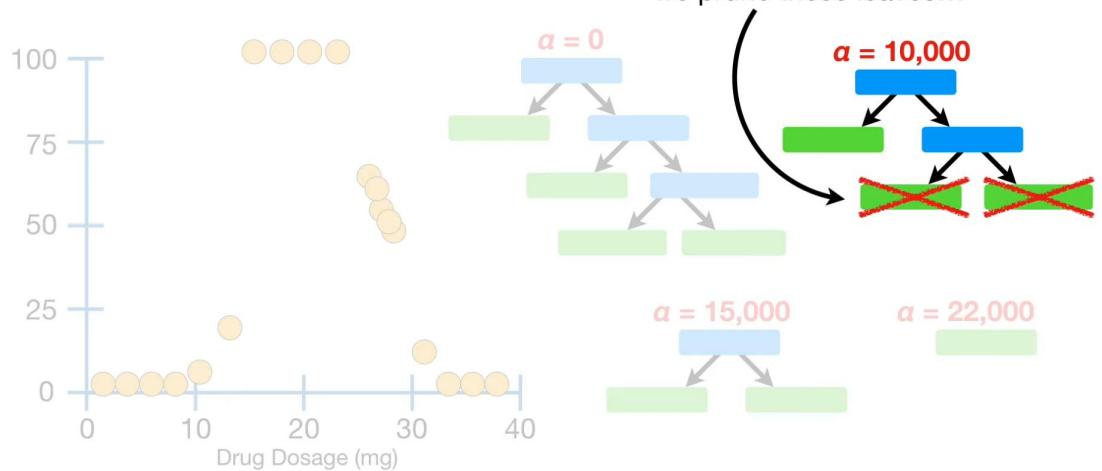
$$\text{Tree Score} = \text{SSR} + \alpha T$$

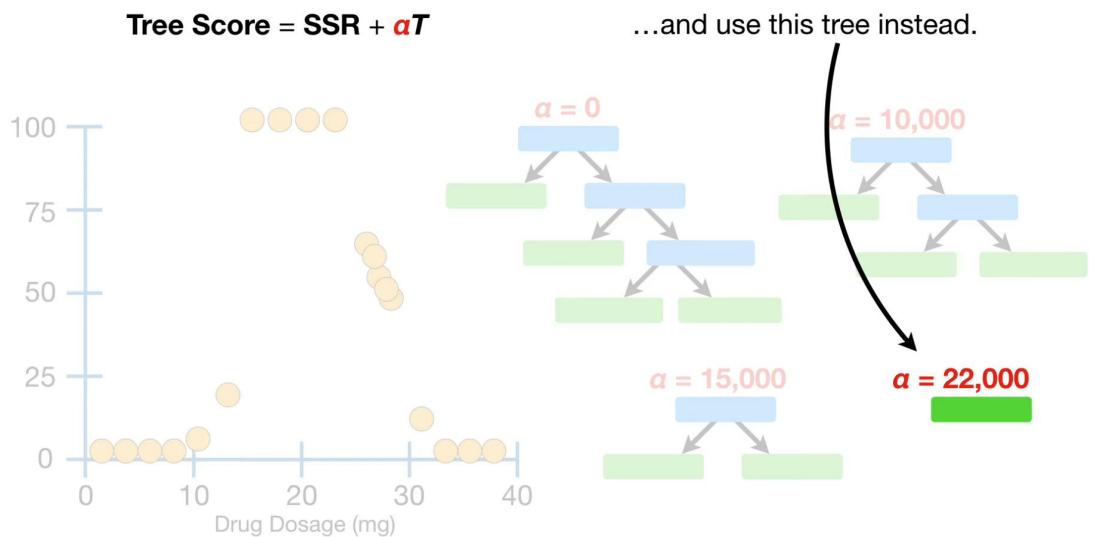
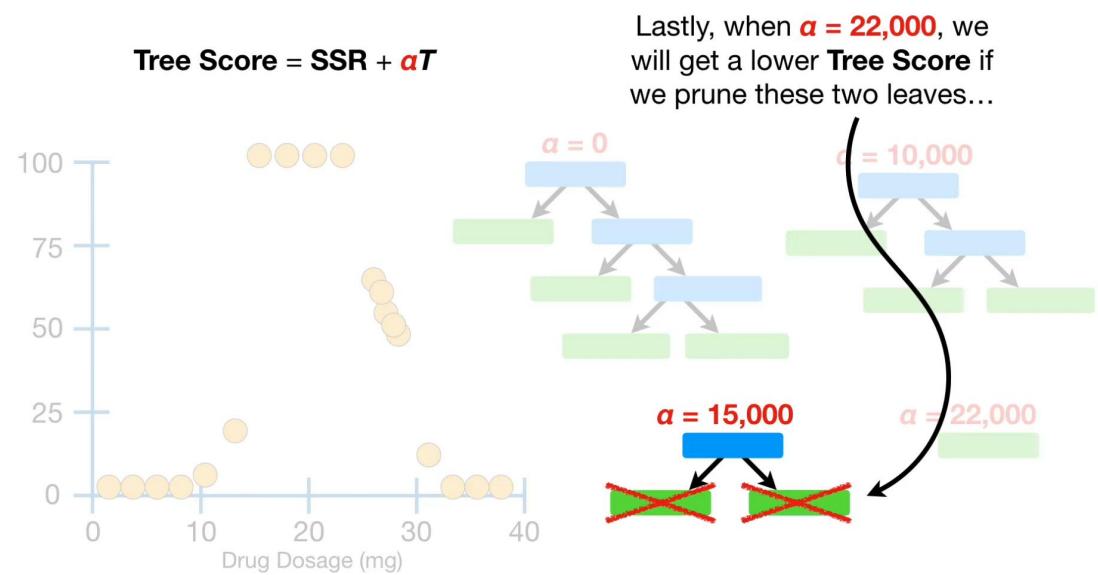
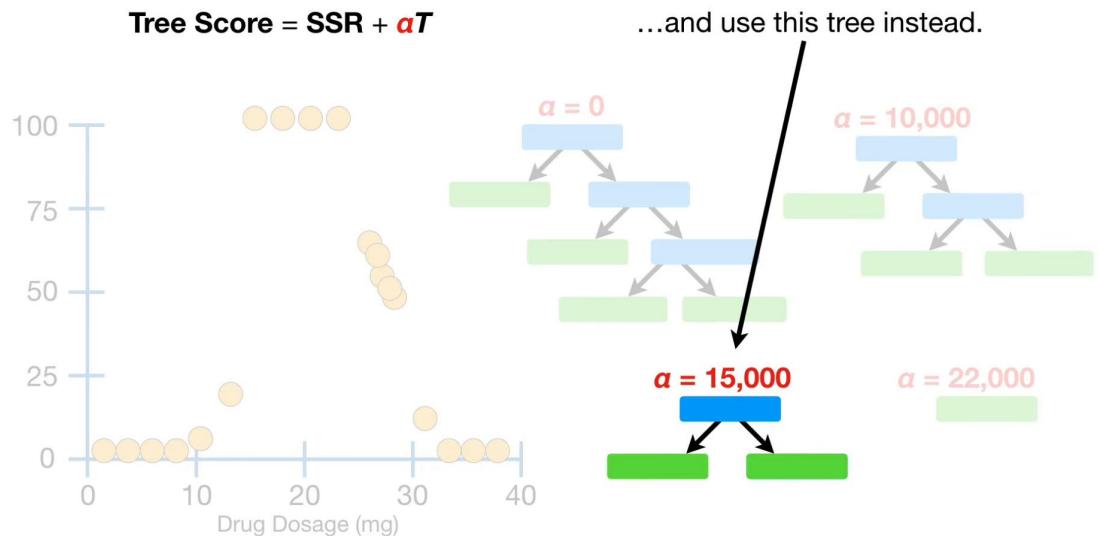
...and use this tree instead.



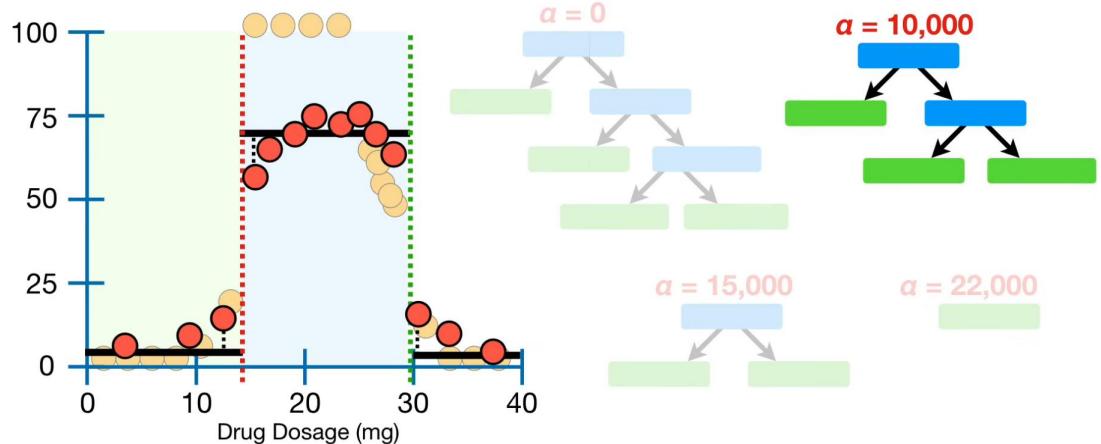
$$\text{Tree Score} = \text{SSR} + \alpha T$$

And when $\alpha = 15,000$, we will get a lower **Tree Score** if we prune these leaves...

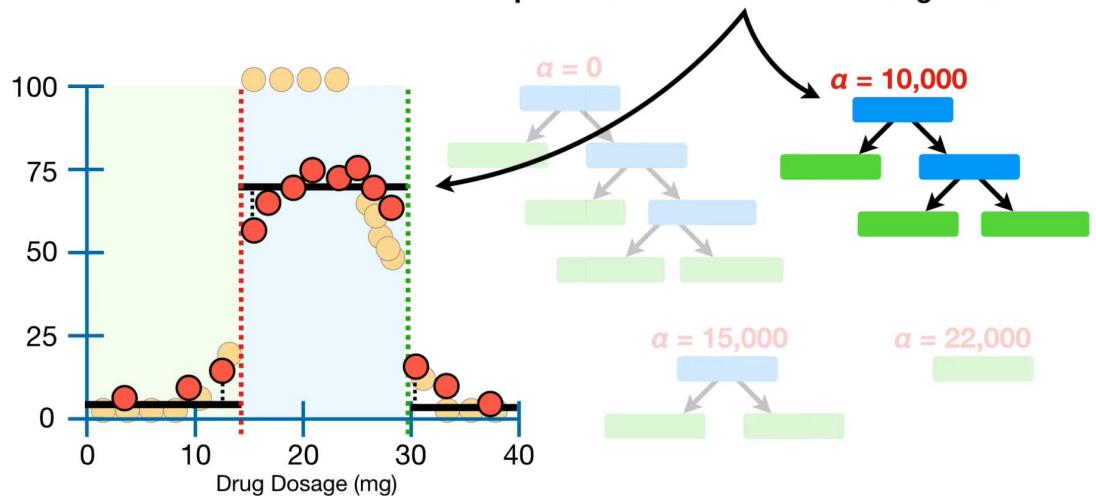




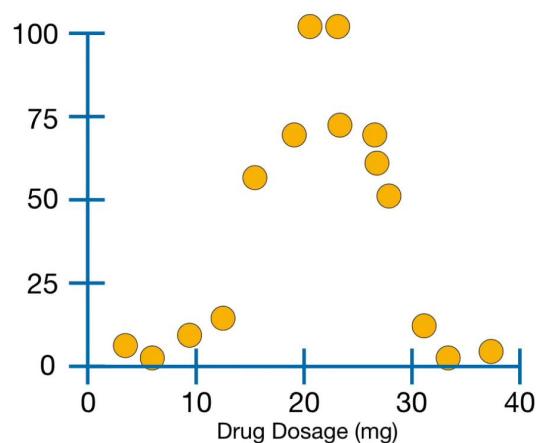
Now calculate the **Sum of Squared Residuals** for each new tree using only the **Testing Data**.



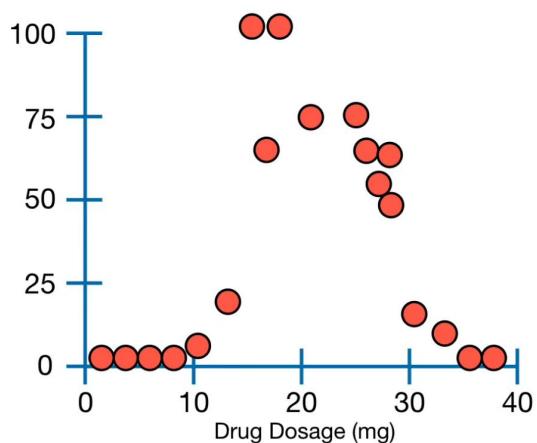
In this case, the tree with $\alpha = 10,000$ had the smallest Sum of Squared Residuals for the Testing Data.



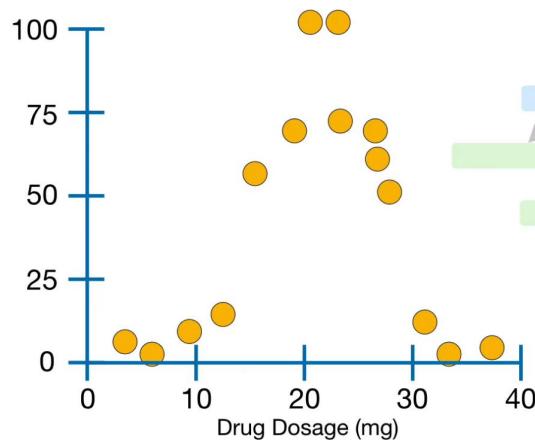
Now we go back and create new **Training Data**...



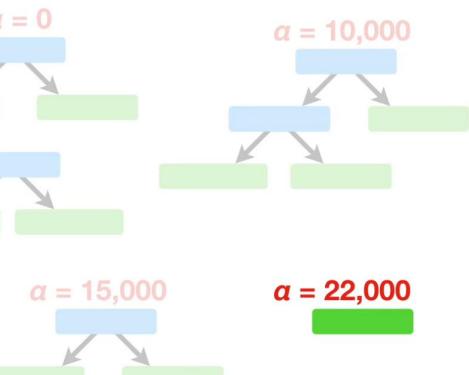
...and new **Testing Data**.



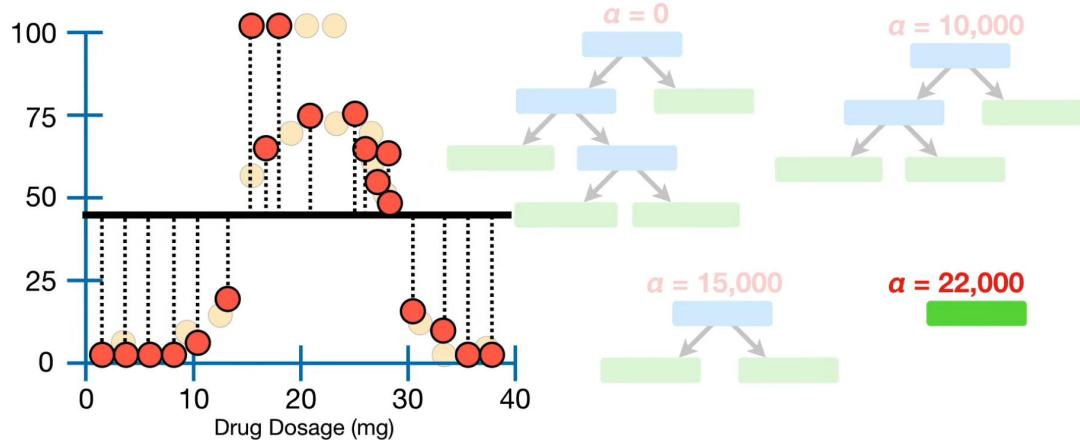
And just using the new
Training Data...



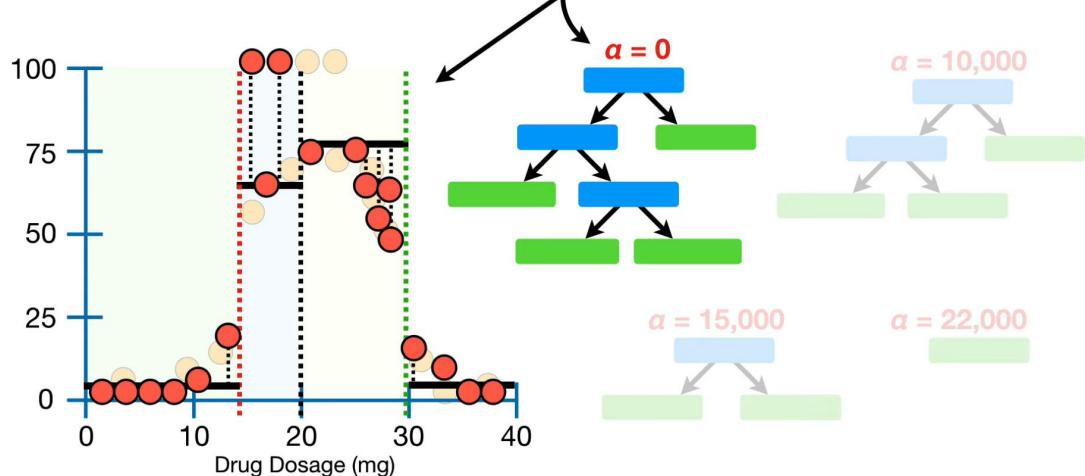
...build a new sequence of trees, from full sized to a leaf, using the α values we found before.



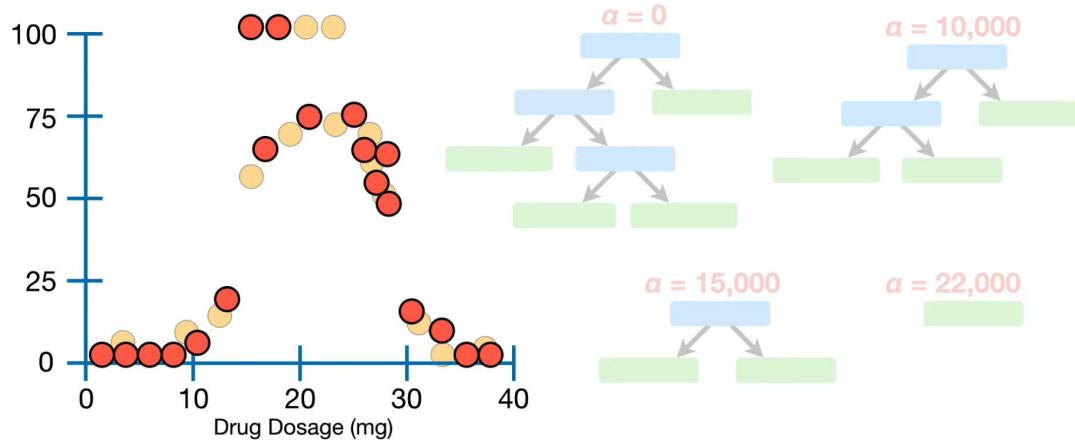
Then we calculate the **Sum of Squared Residuals** using the new **Testing Data**.



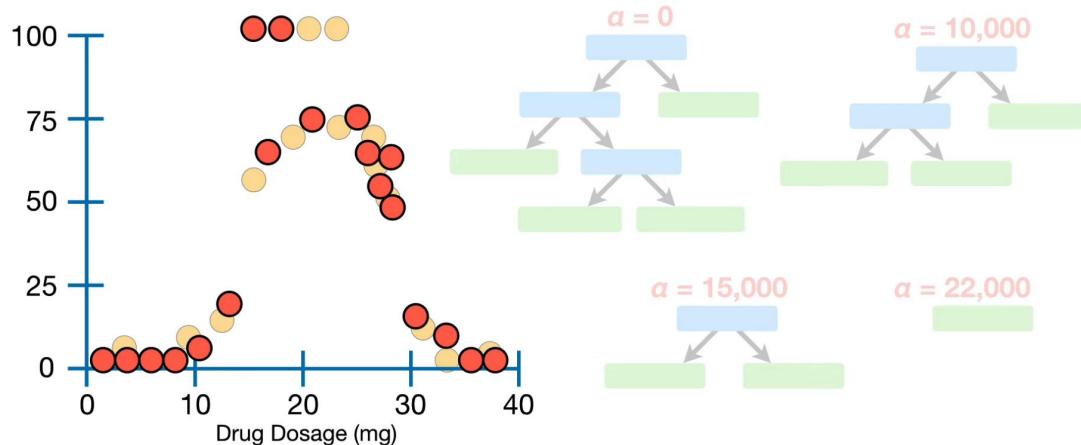
This time, the tree with $\alpha = 0$ had the lowest **Sum of Squared Residuals**.



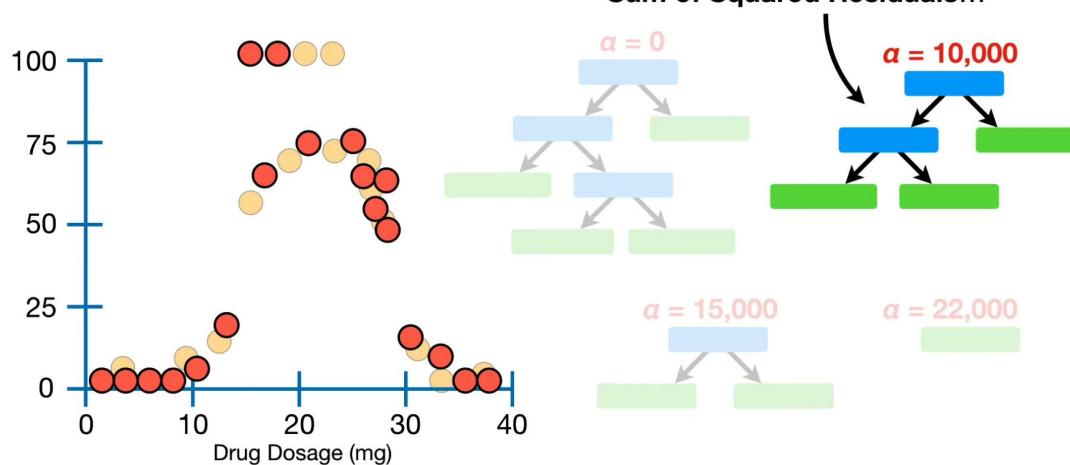
Now we just keep repeating until we have done **10-Fold Cross Validation**...



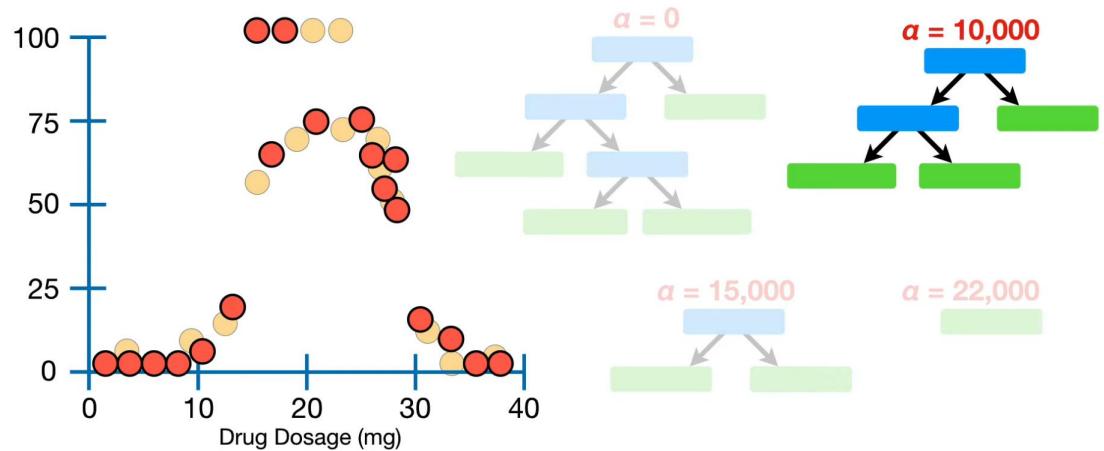
...and the value for α that, on average, gave us the lowest **Sum of Squared Residuals** with the **Testing Data**, is the final value for α .



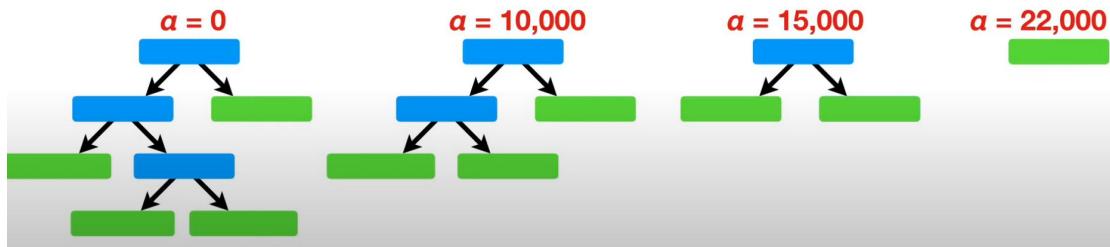
In this case, the optimal trees built with $\alpha = 10,000$ had, on average, the lowest **Sum of Squared Residuals**...



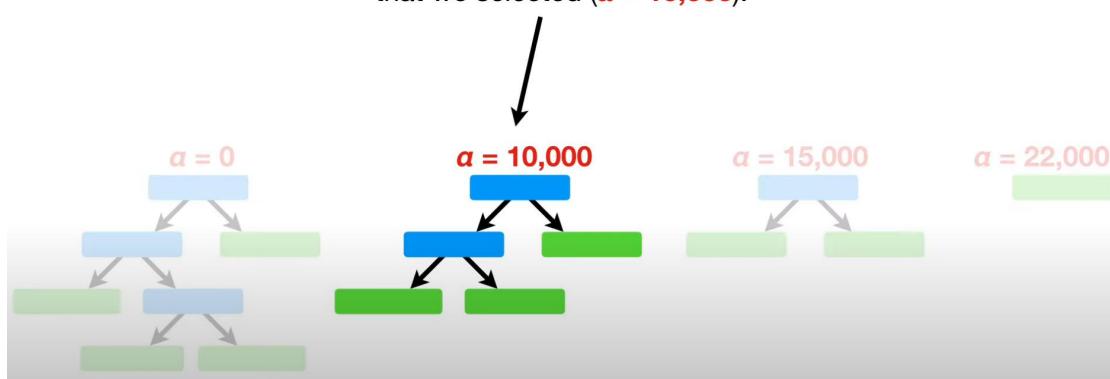
...so $\alpha = 10,000$ is our final value.



Lastly, we go back to the original trees and sub-trees made from the full data.



...and pick the tree that corresponds to the value for α that we selected ($\alpha = 10,000$).



This sub-tree will be the final,
pruned tree.

