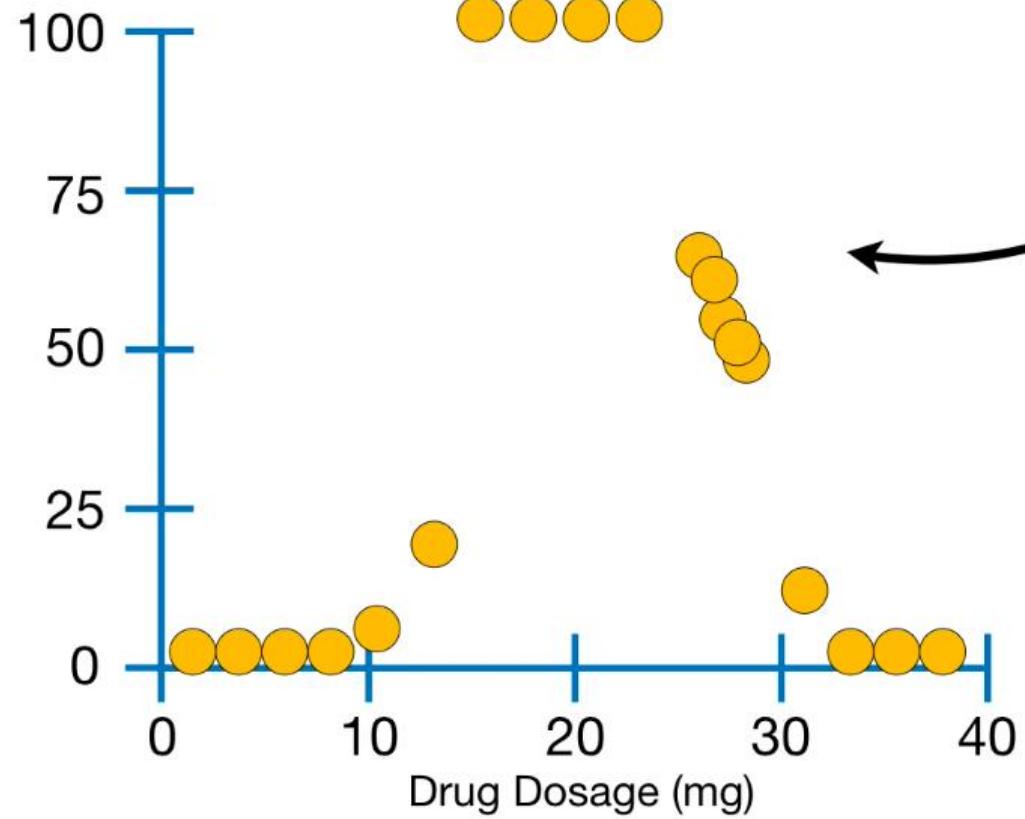
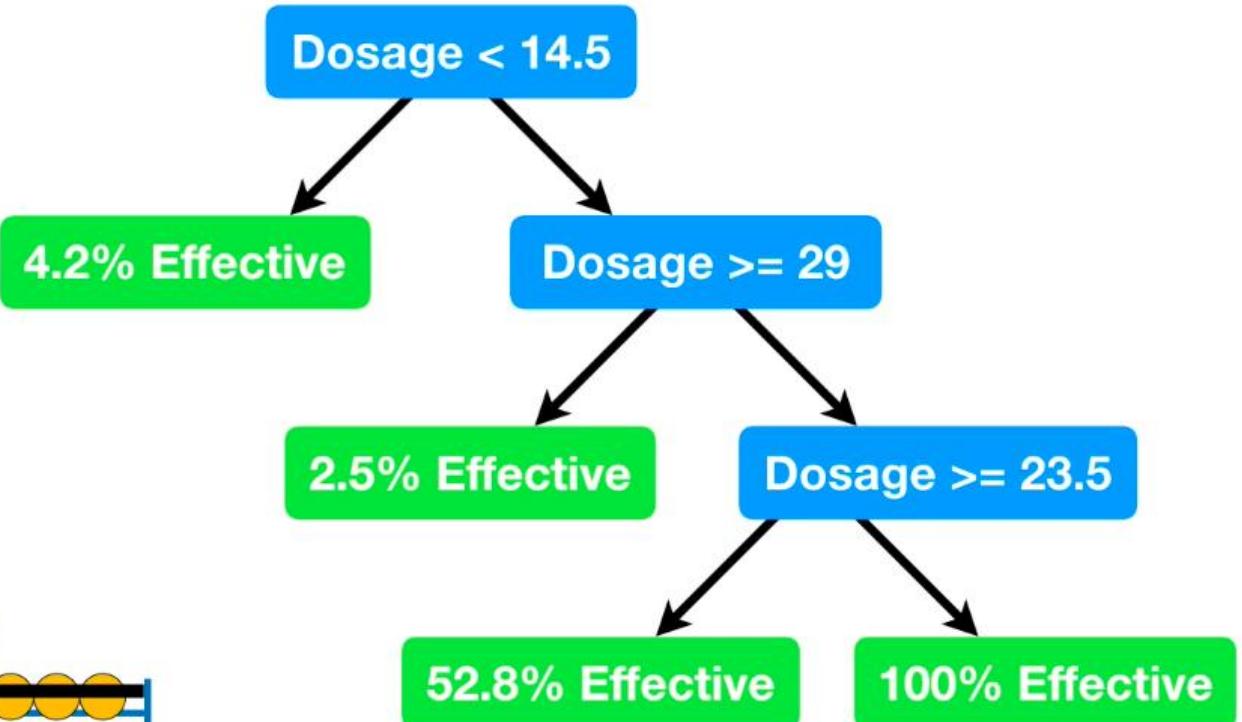
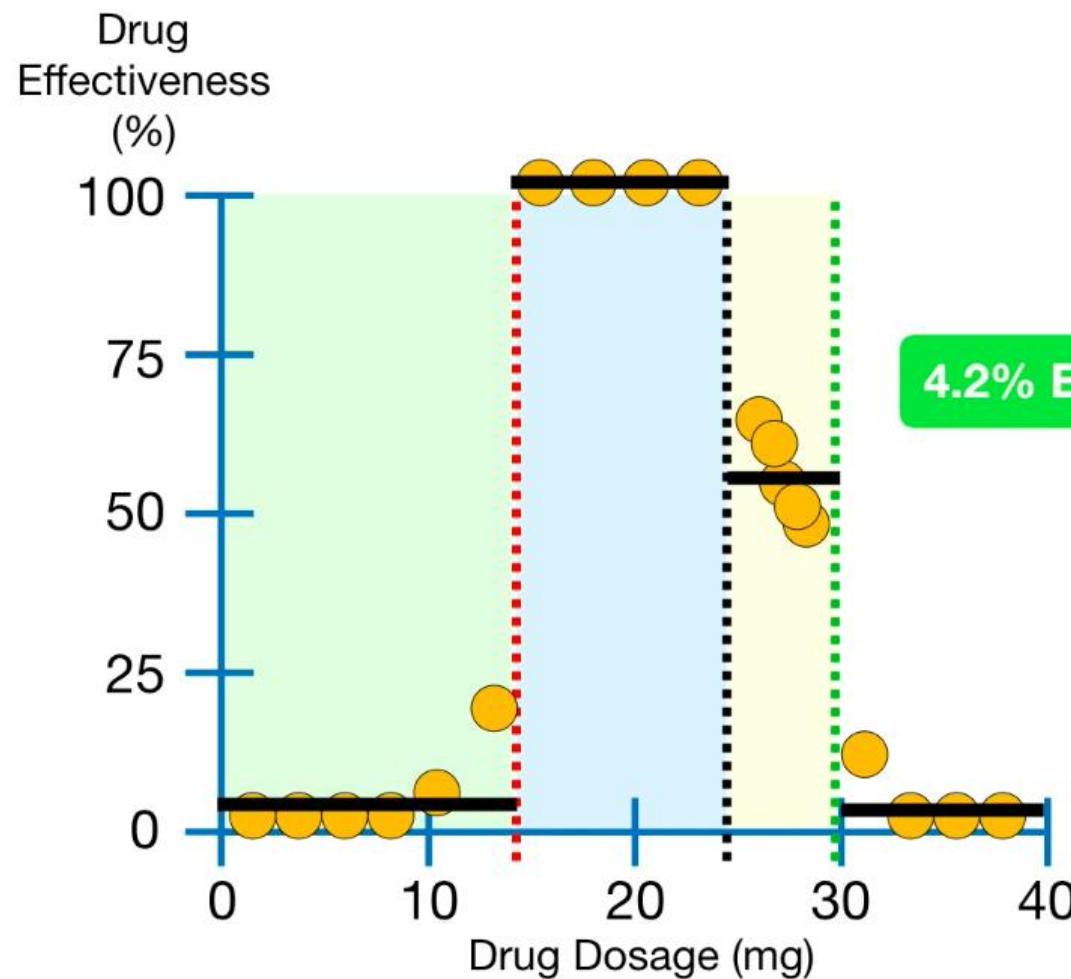


Drug
Effectiveness
(%)

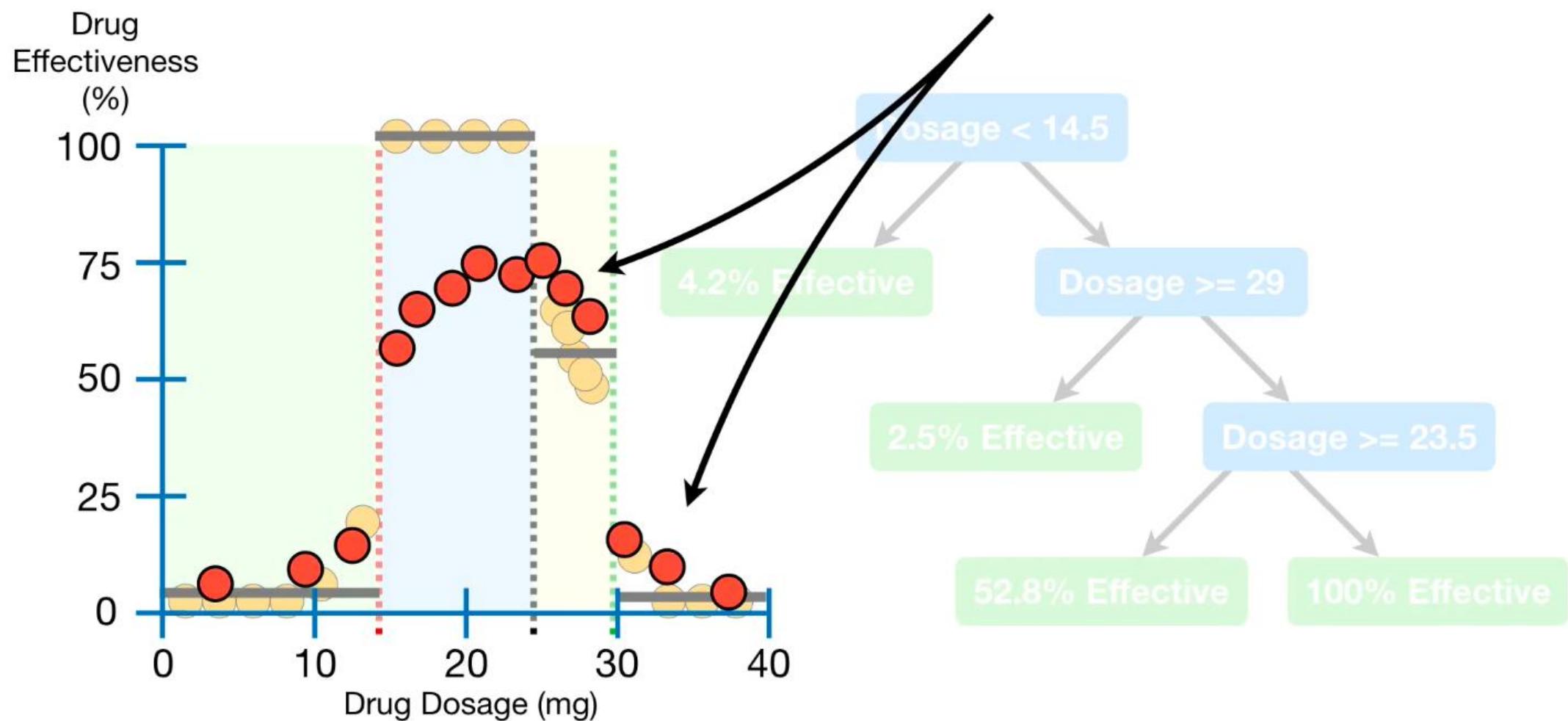


In the **StatQuest** on **Regression Trees**,
we had this data.

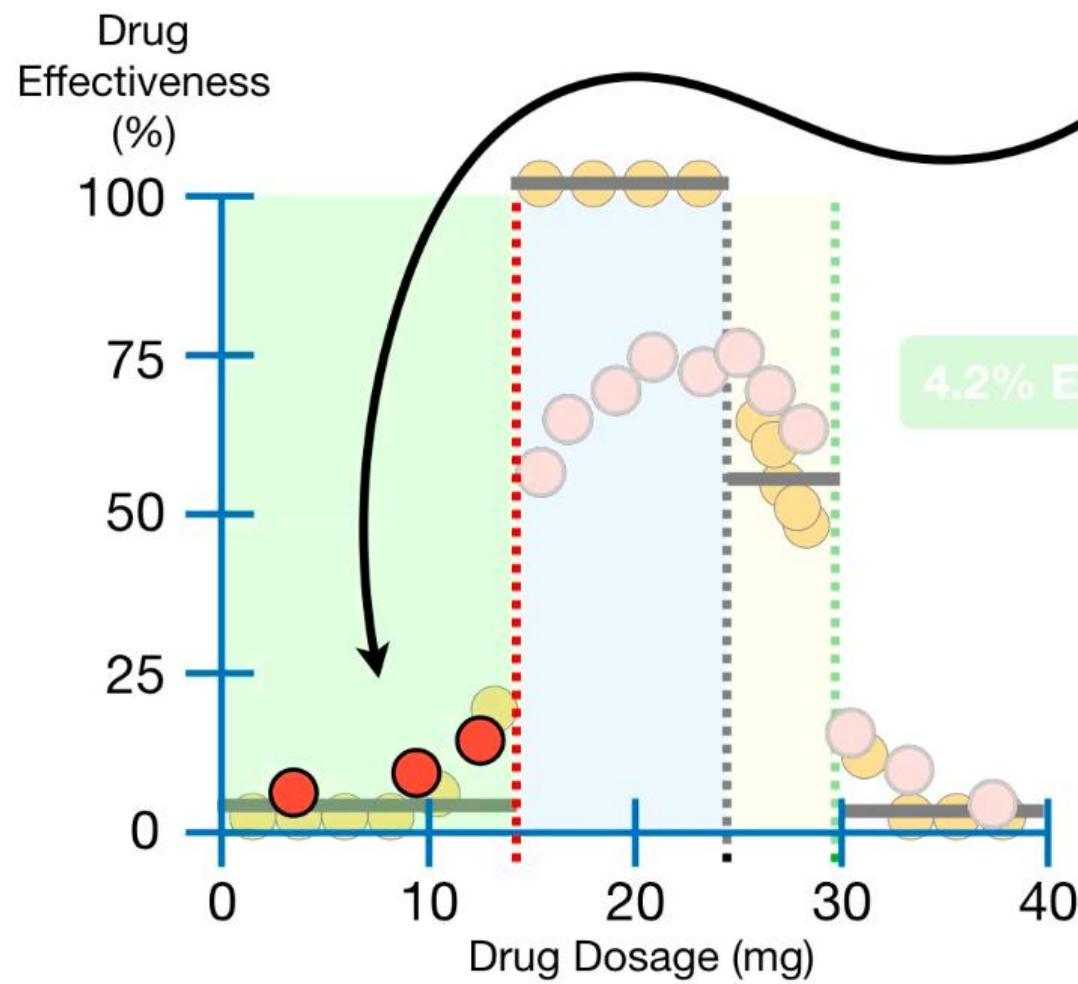
We then fit a **Regression Tree**
to the data...



However, what if these **Red Circles** were **Testing Data**?



These three observations...



Dosage < 14.5

4.2% Effective

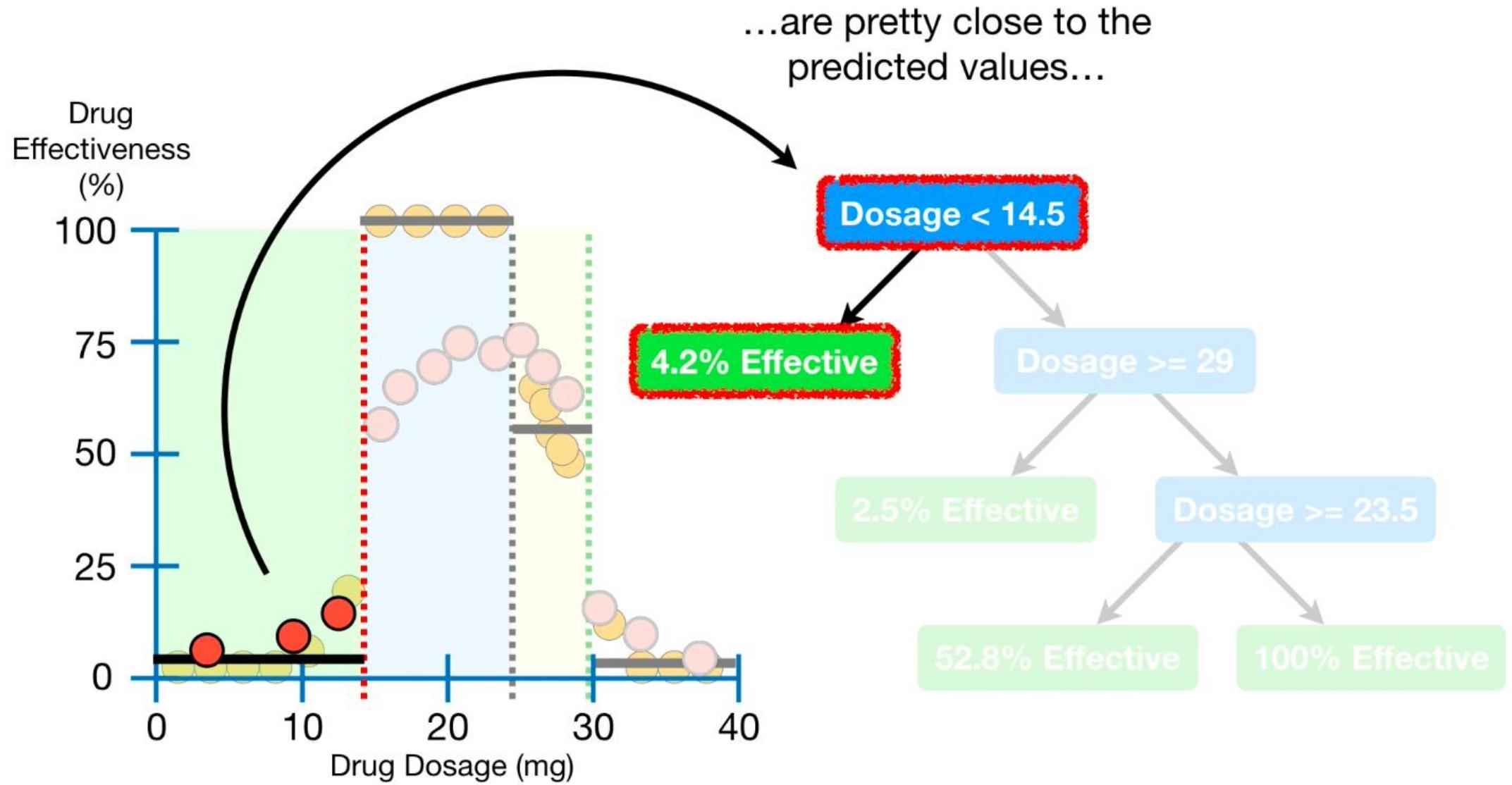
Dosage >= 29

2.5% Effective

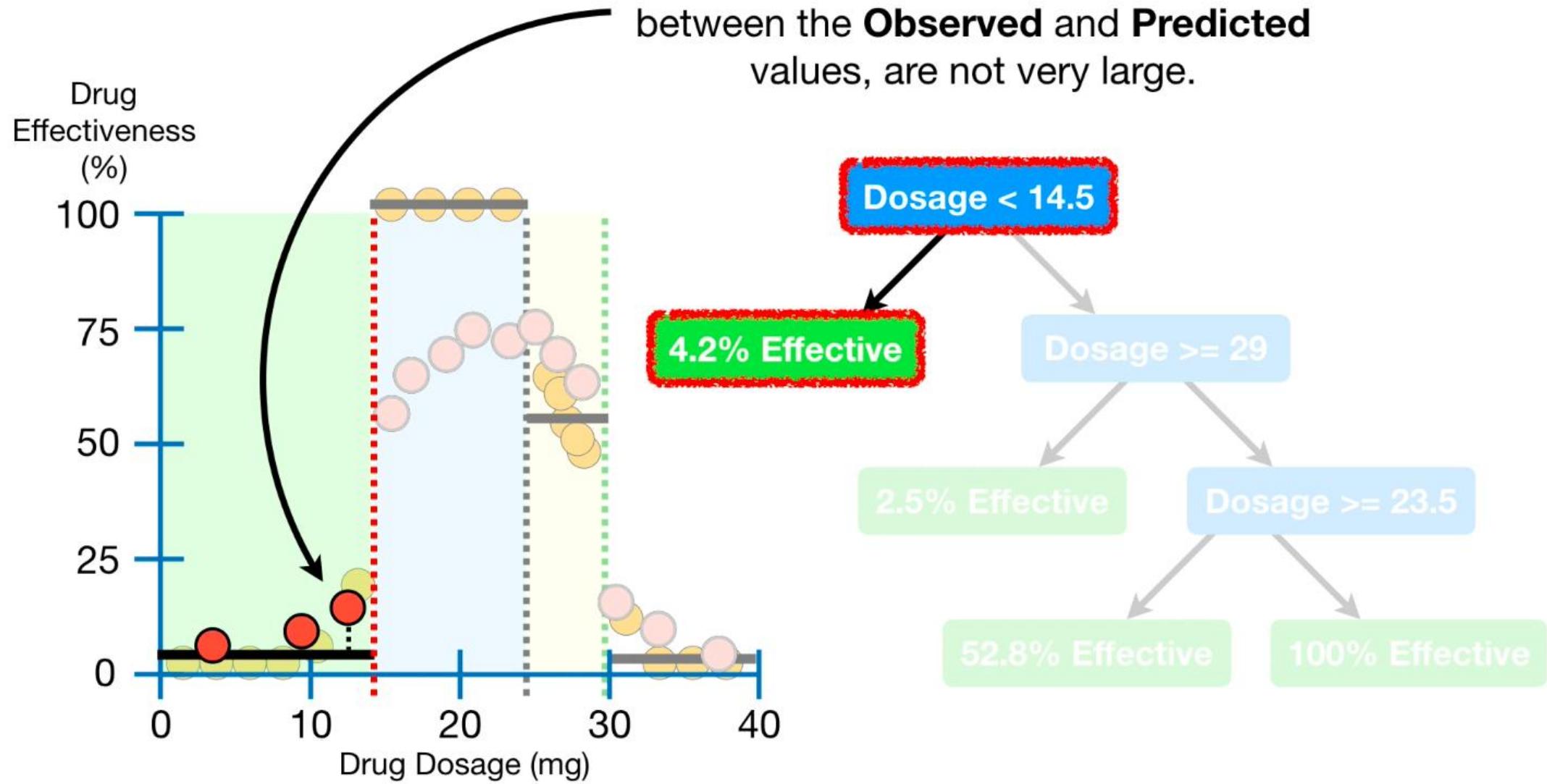
Dosage >= 23.5

52.8% Effective

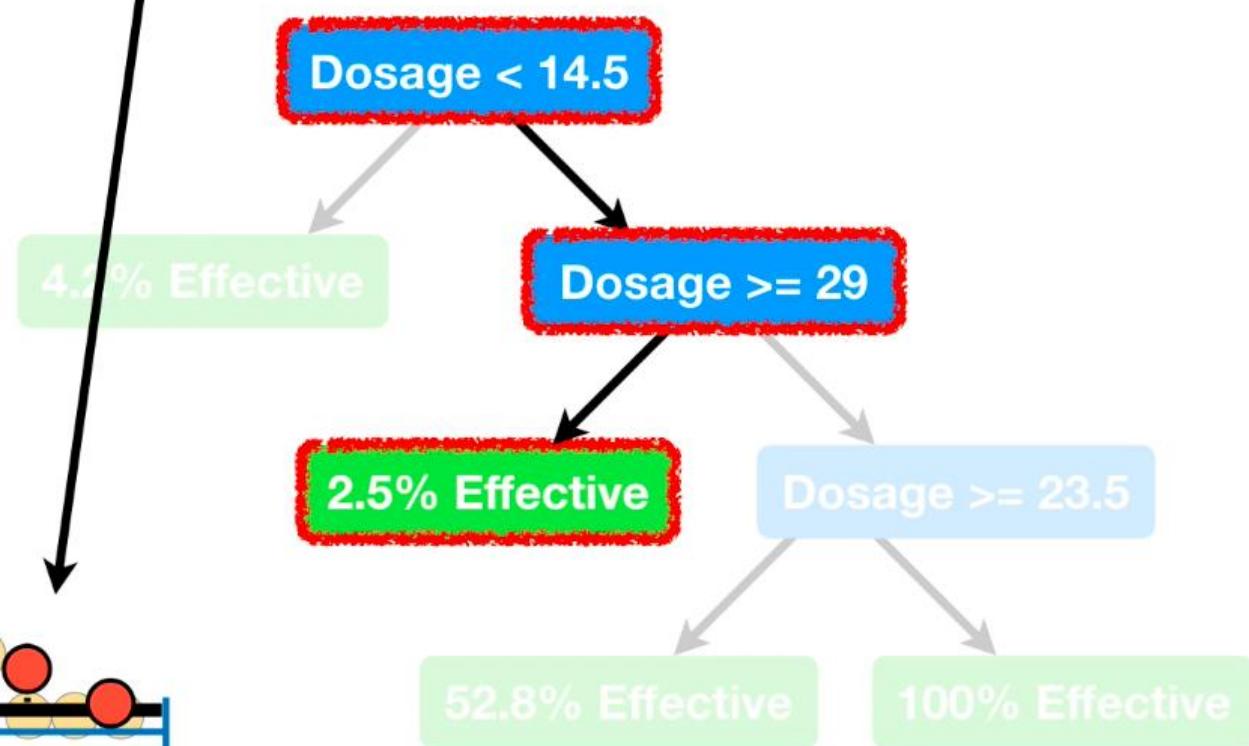
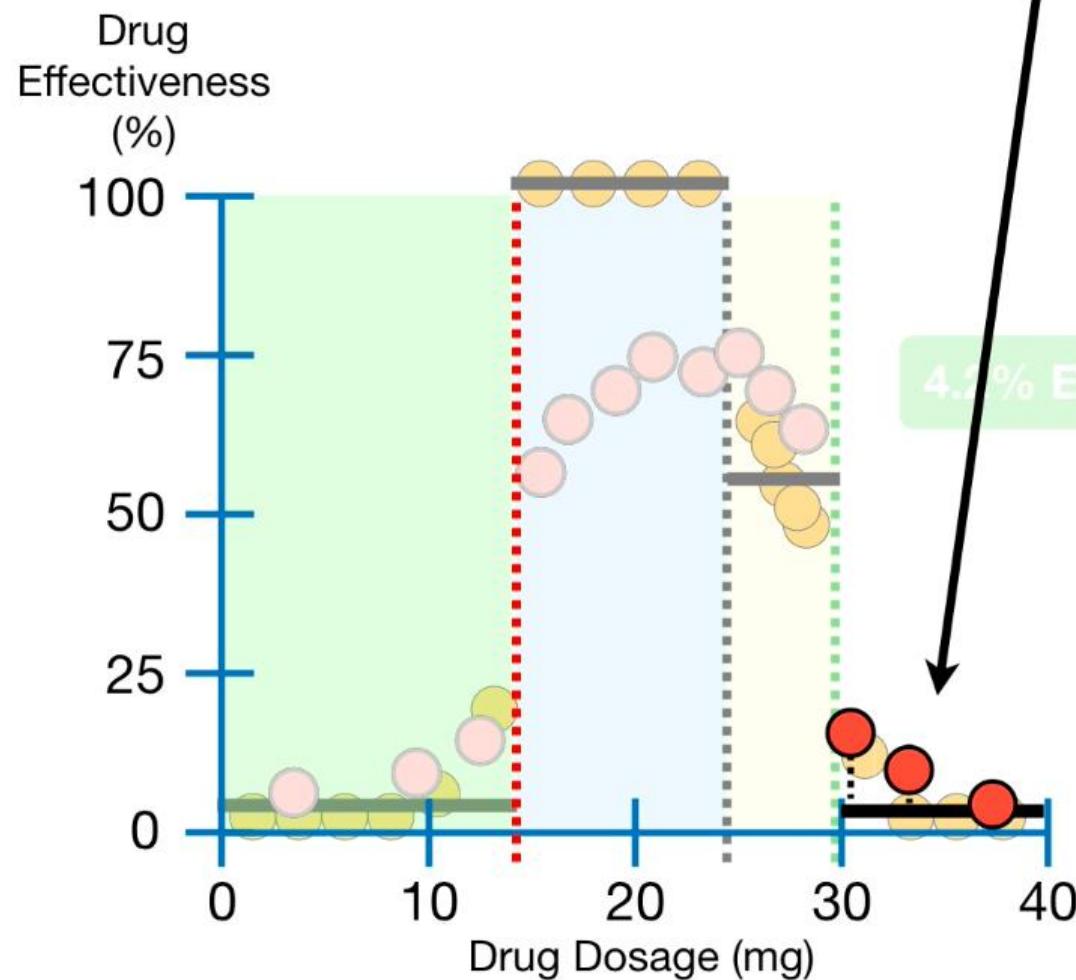
100% Effective



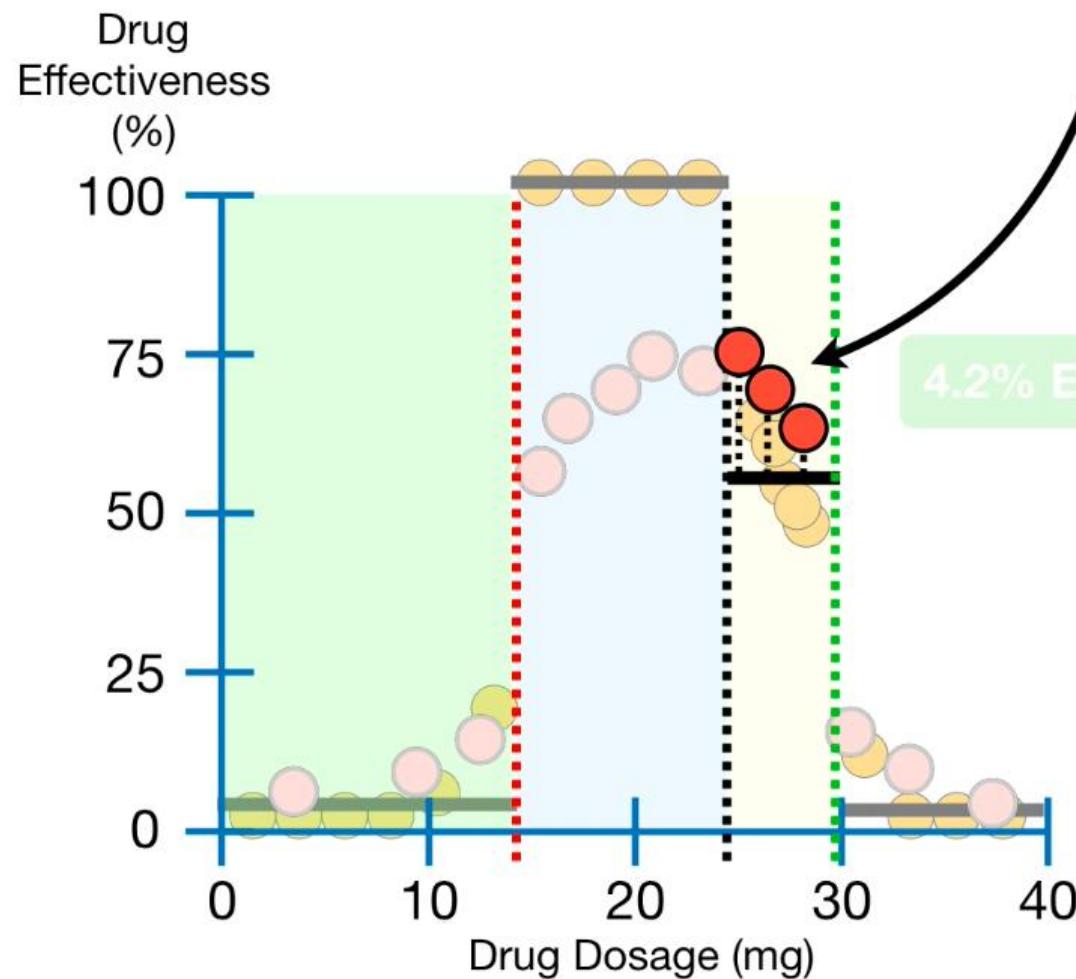
...so their **Residuals**, the difference between the **Observed** and **Predicted** values, are not very large.



Similarly, the **Residuals** for these observations in the **Testing Data** are relatively small.



However, the **Residuals** for these **Observations** are larger than before.



Dosage < 14.5

4.2% Effective

Dosage >= 29

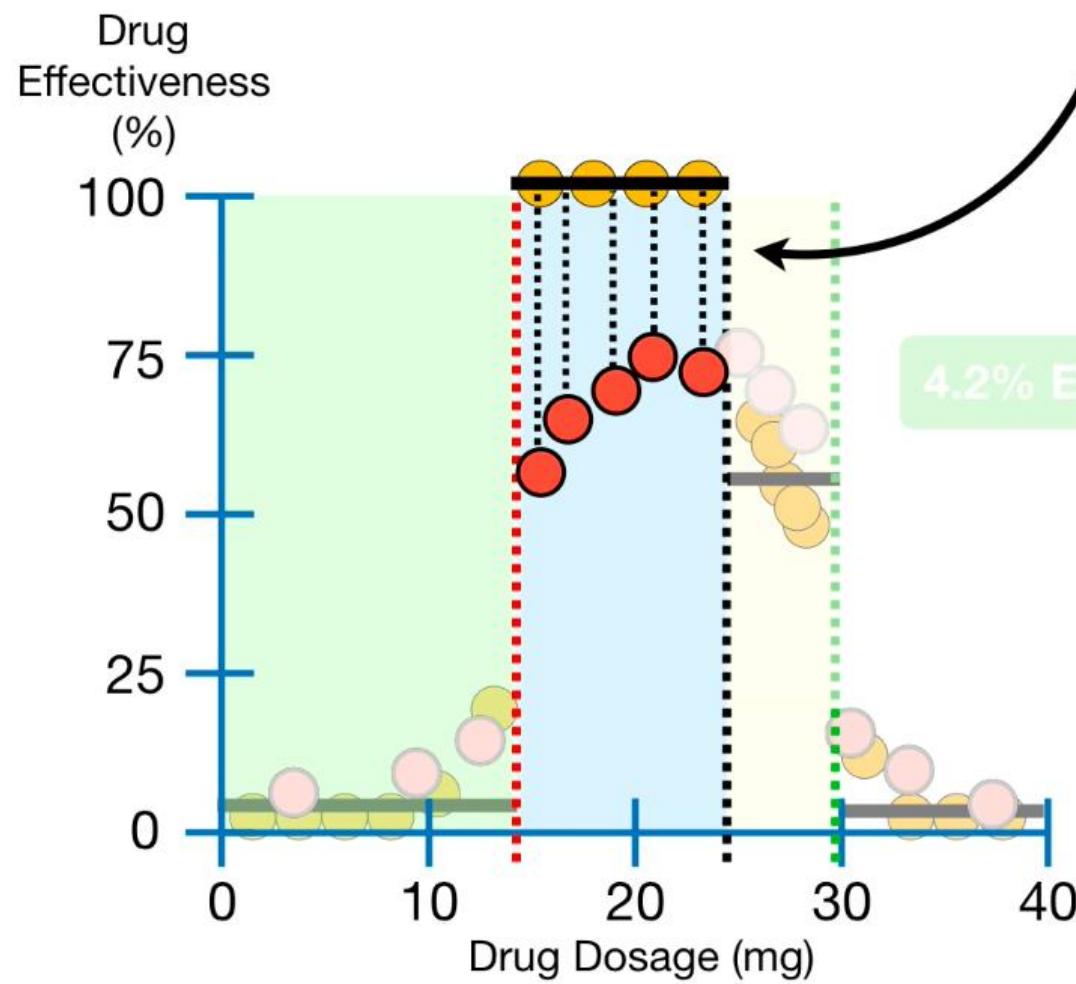
2.5% Effective

Dosage >= 23.5

52.8% Effective

100% Effective

And the **Residuals** for these **Observations** are much larger.

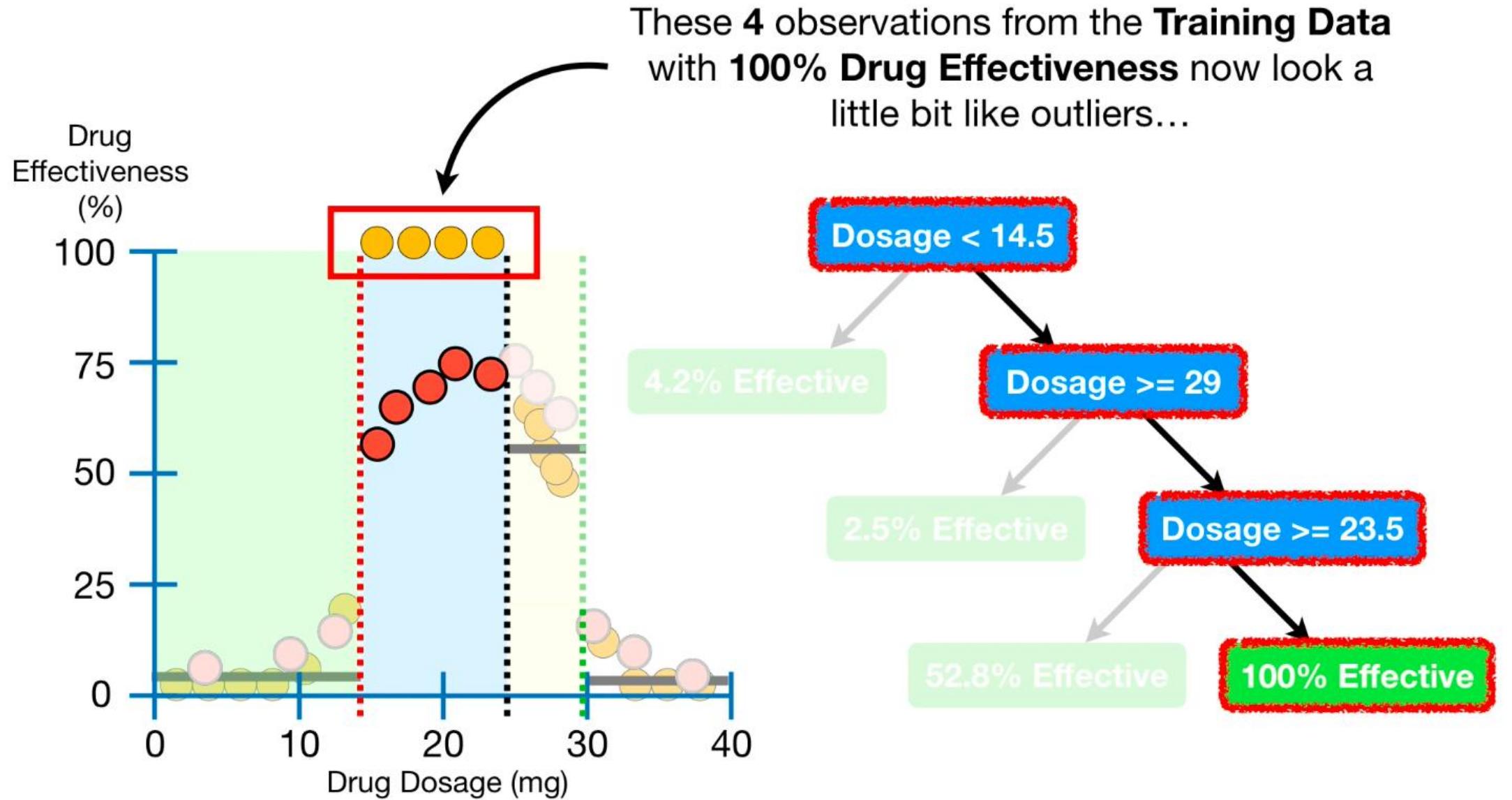


Dosage < 14.5

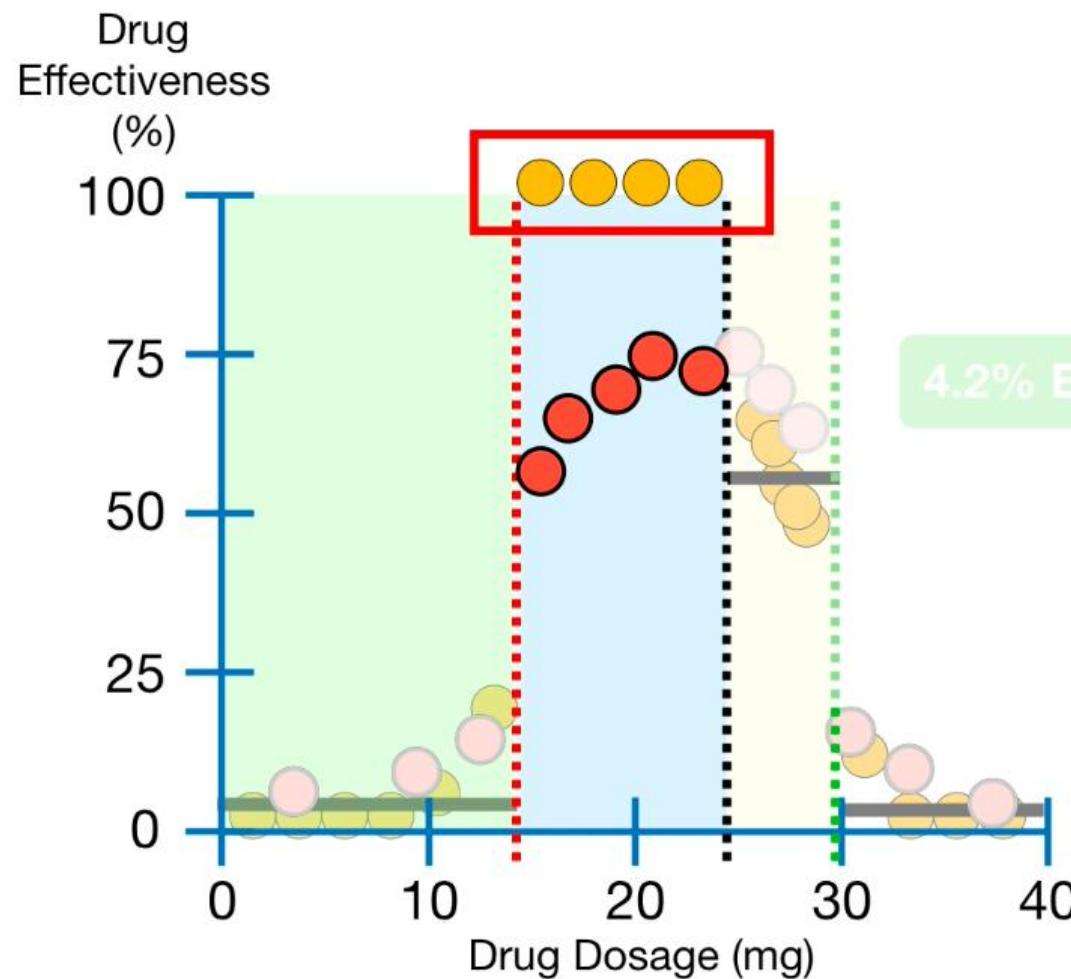
Dosage ≥ 29

Dosage ≥ 23.5

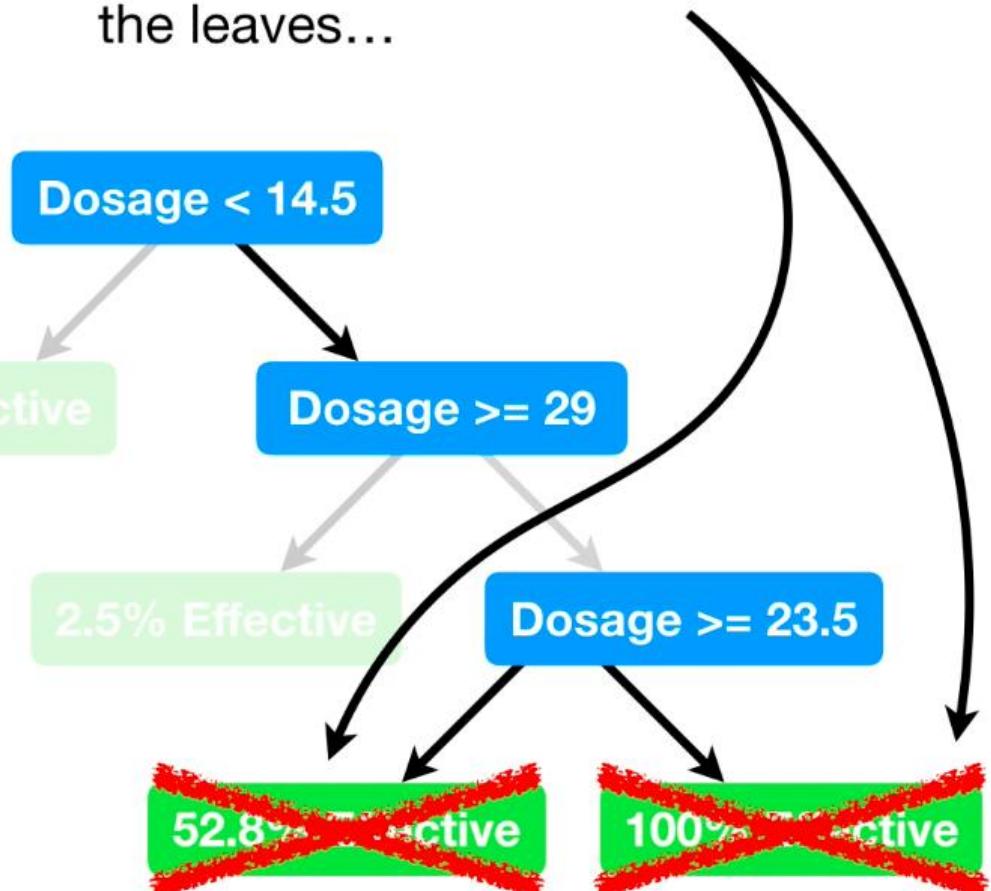
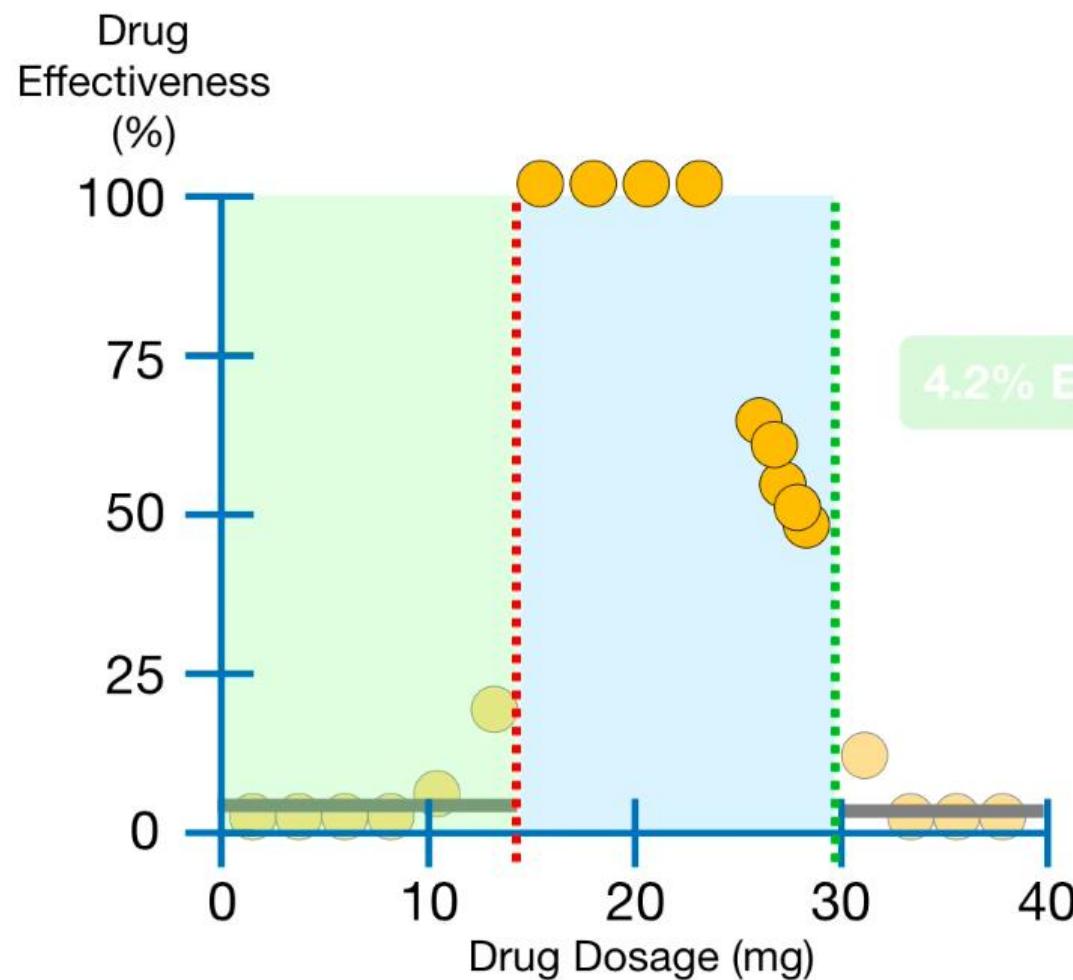
100% Effective



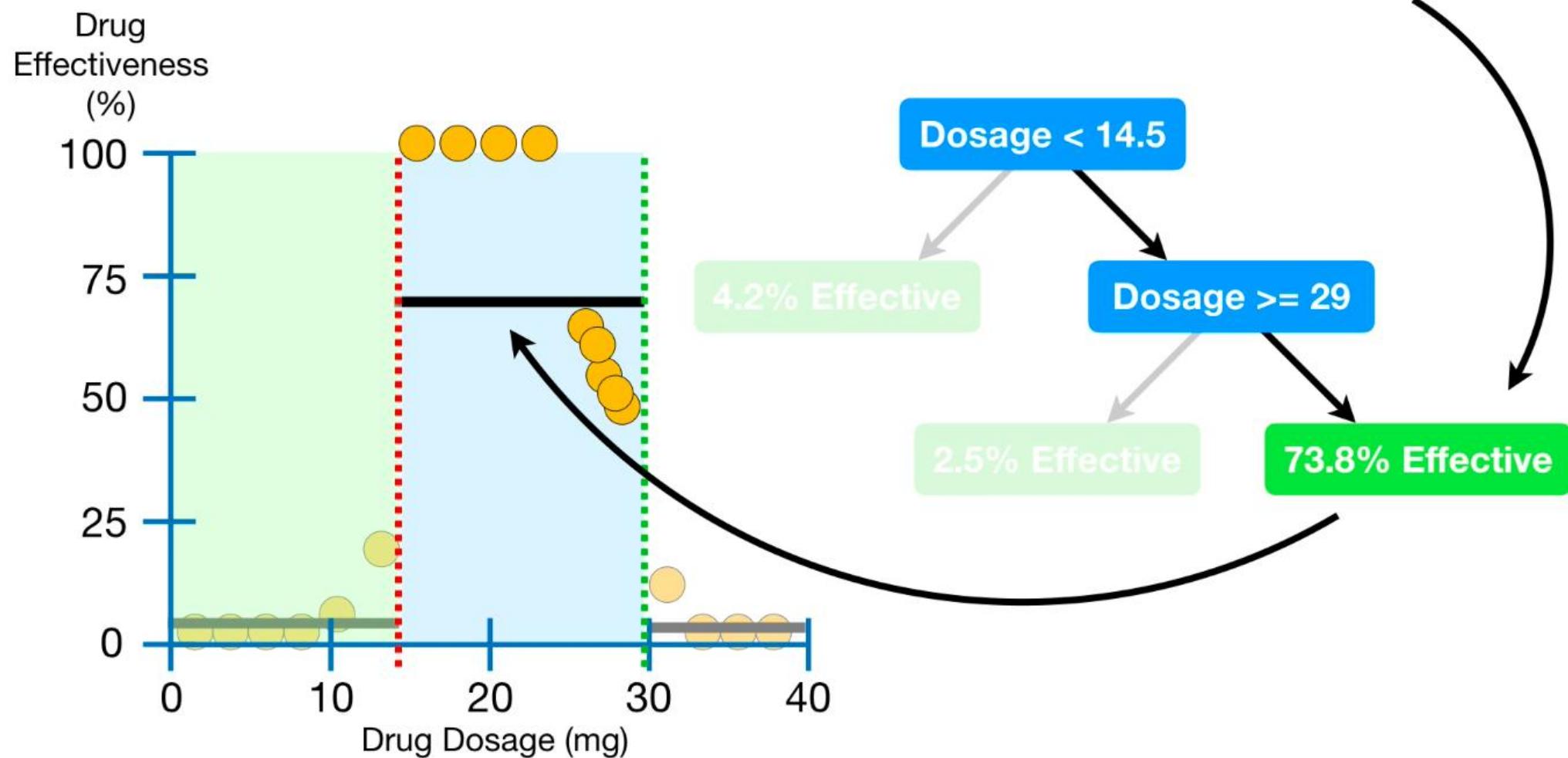
...and that means that maybe we **Overfit** the **Regression Tree** to the **Training Data**.

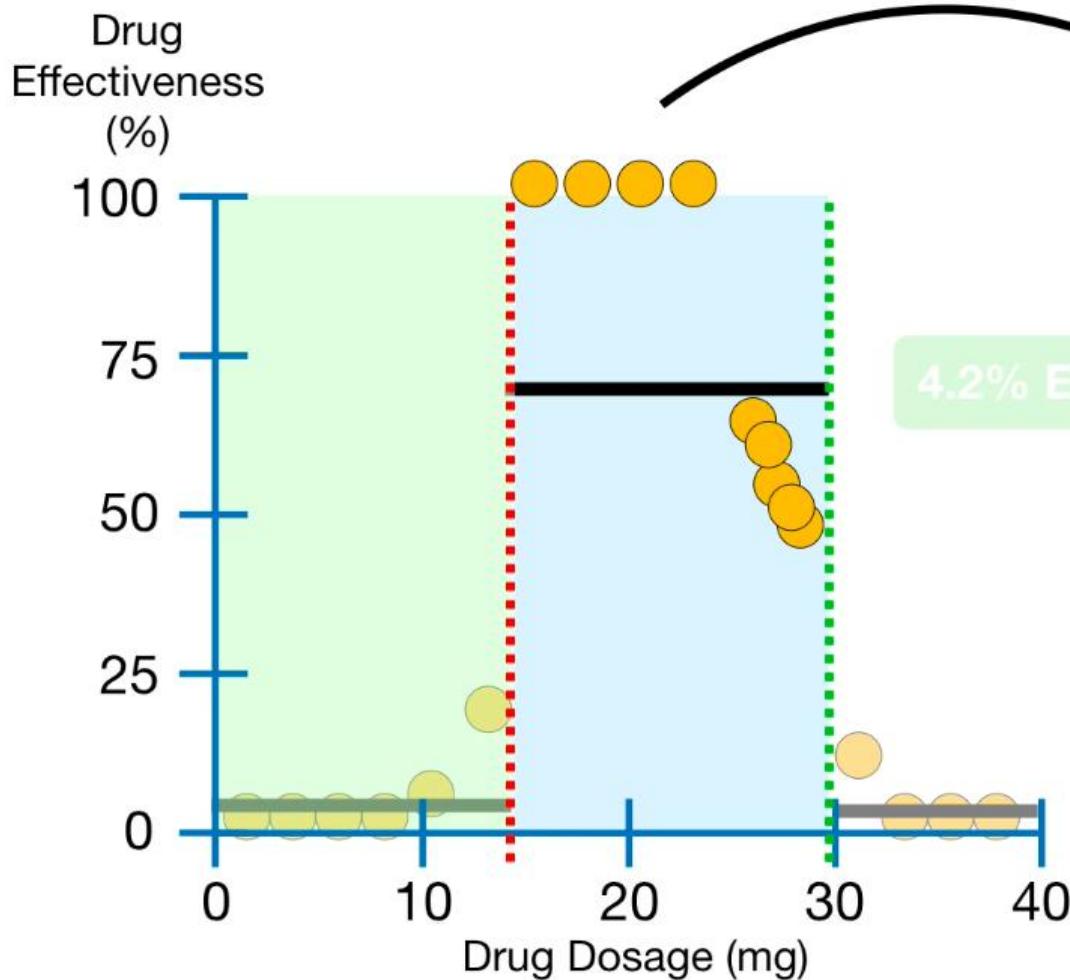


One way to prevent **Overfitting** a **Regression Tree** to the **Training Data** is to remove some of the leaves...



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





Now all of the observations between **14.5** and **29** go to the leaf on the far right.

Dosage < 14.5

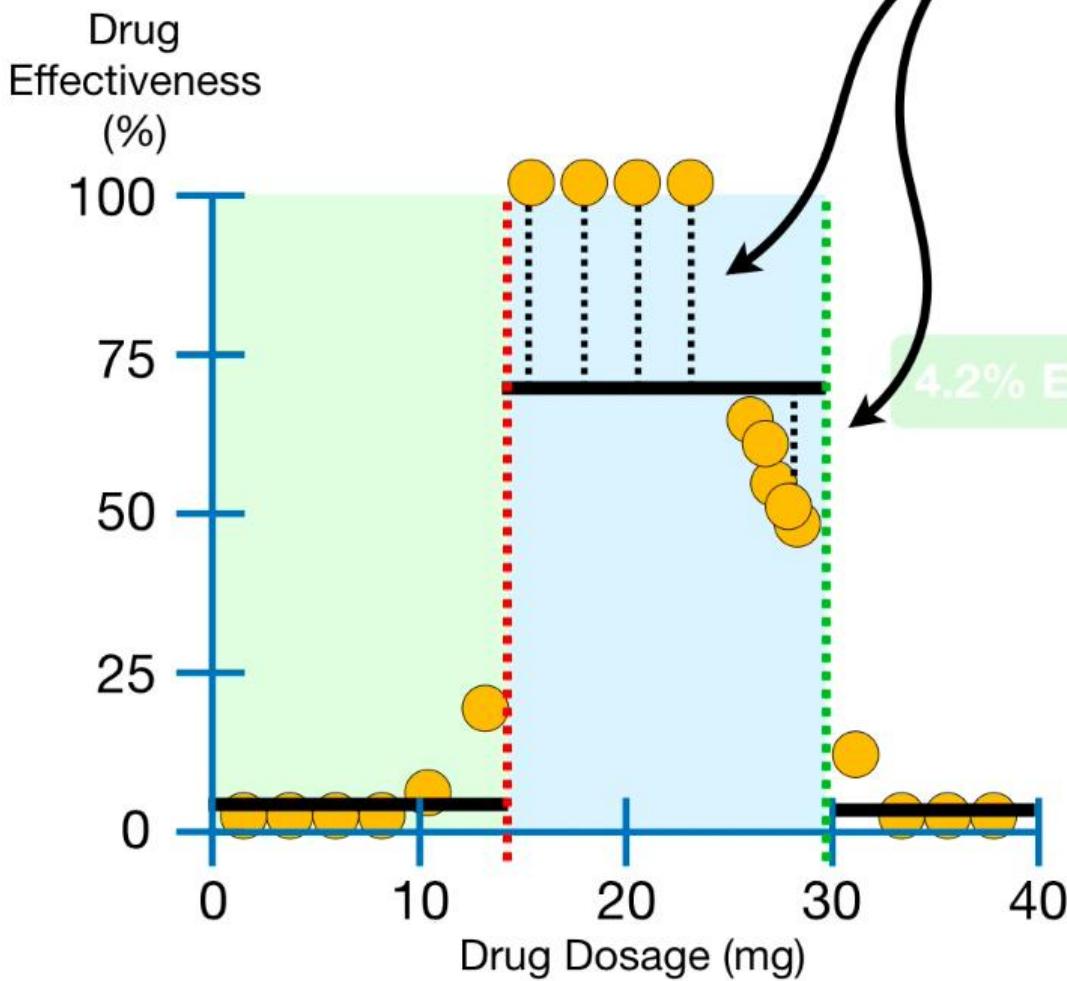
4.2% Effective

Dosage ≥ 29

2.5% Effective

73.8% Effective

The large **Residuals** tell us that the new tree doesn't fit the **Training Data** as well as before...



Dosage < 14.5

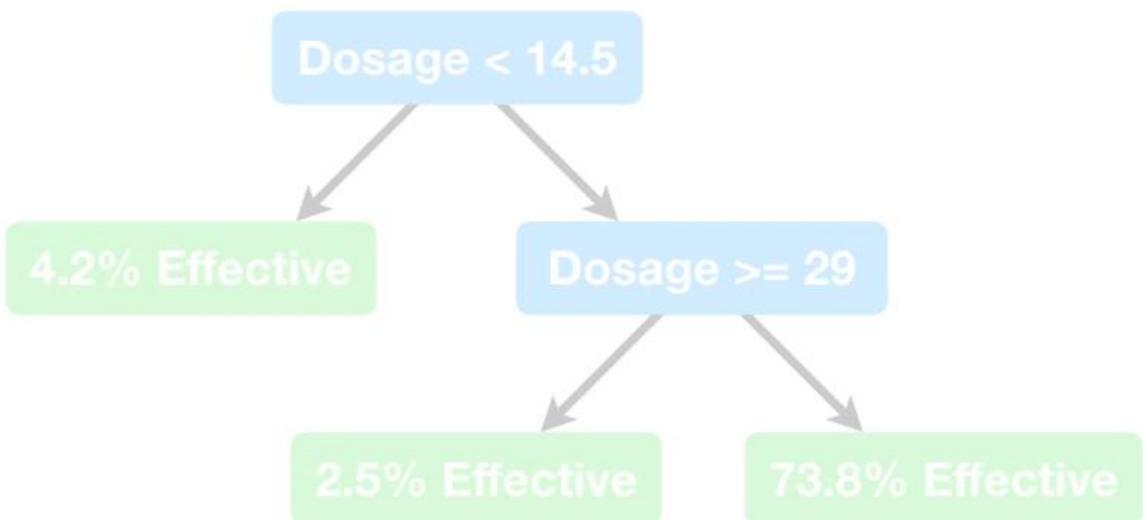
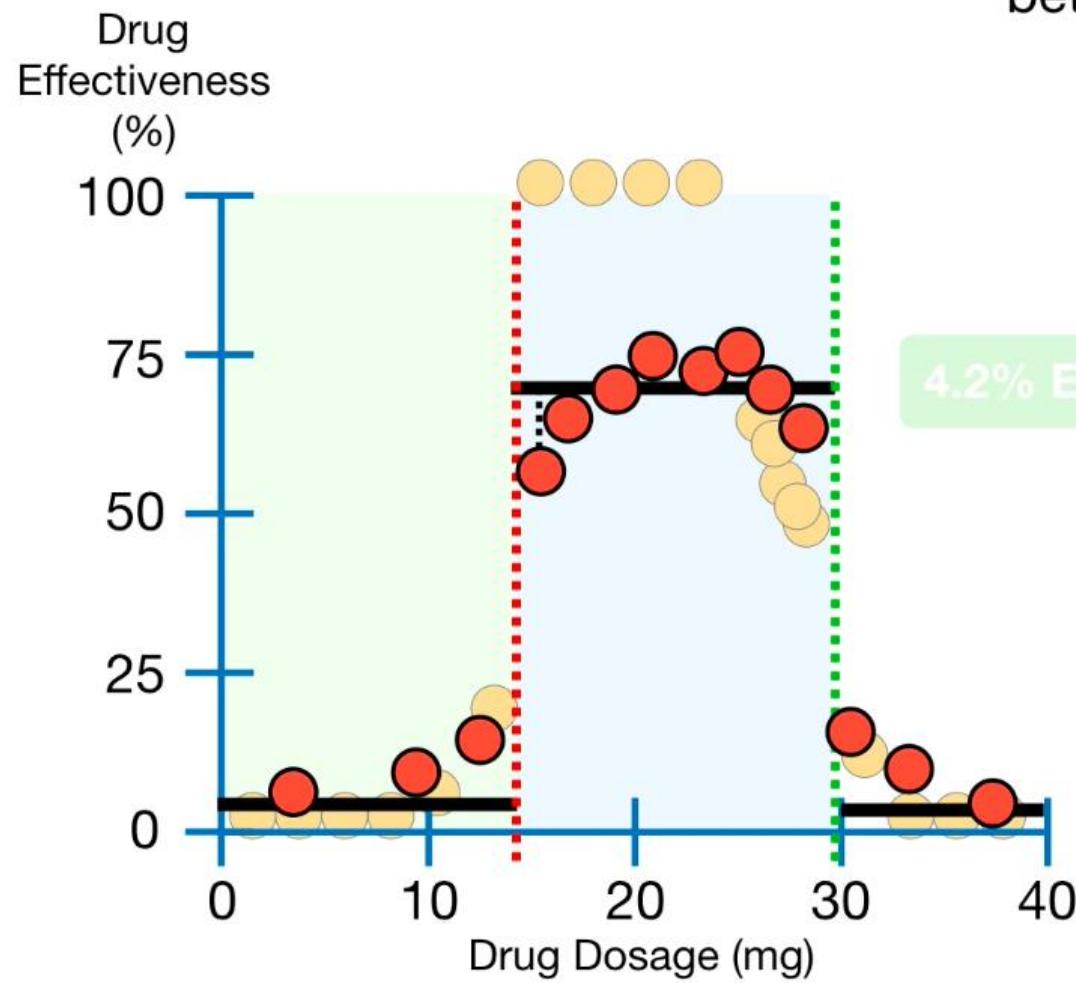
4.2% Effective

Dosage ≥ 29

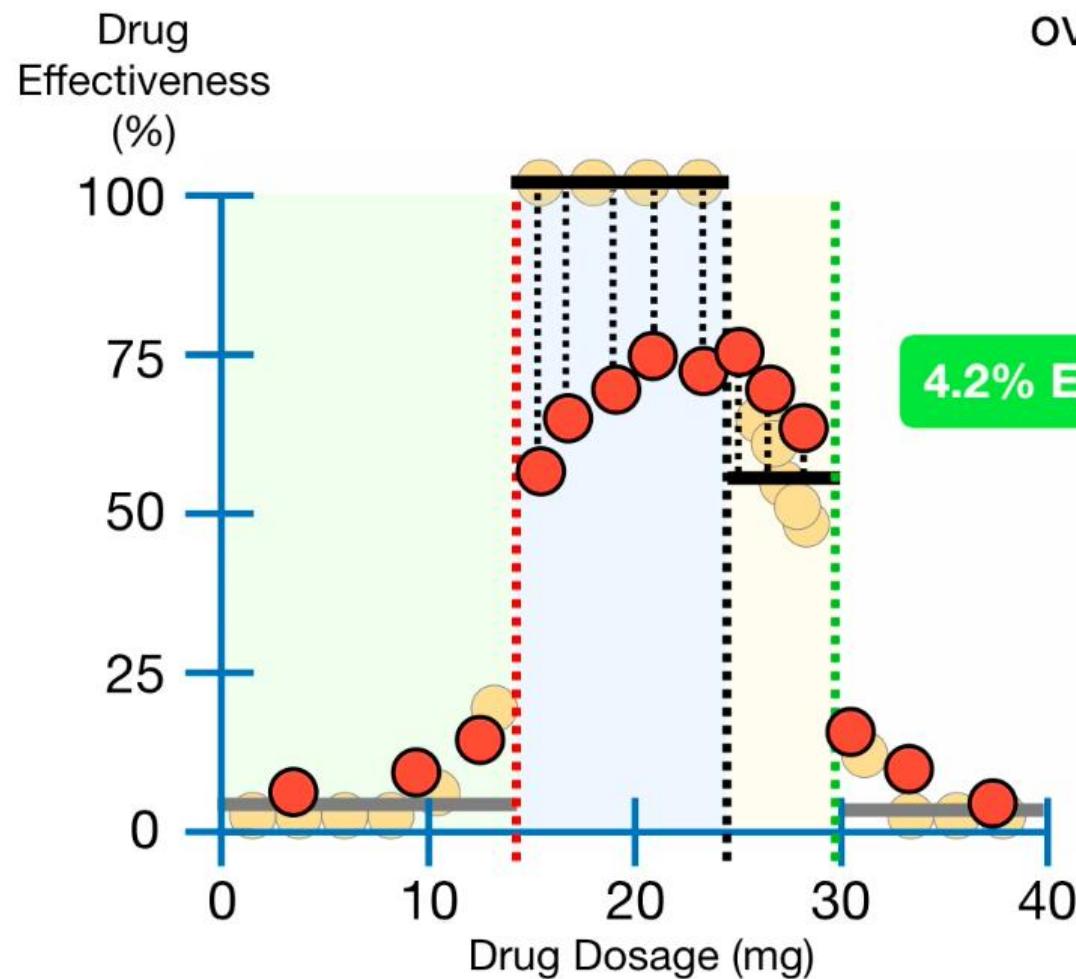
2.5% Effective

73.8% Effective

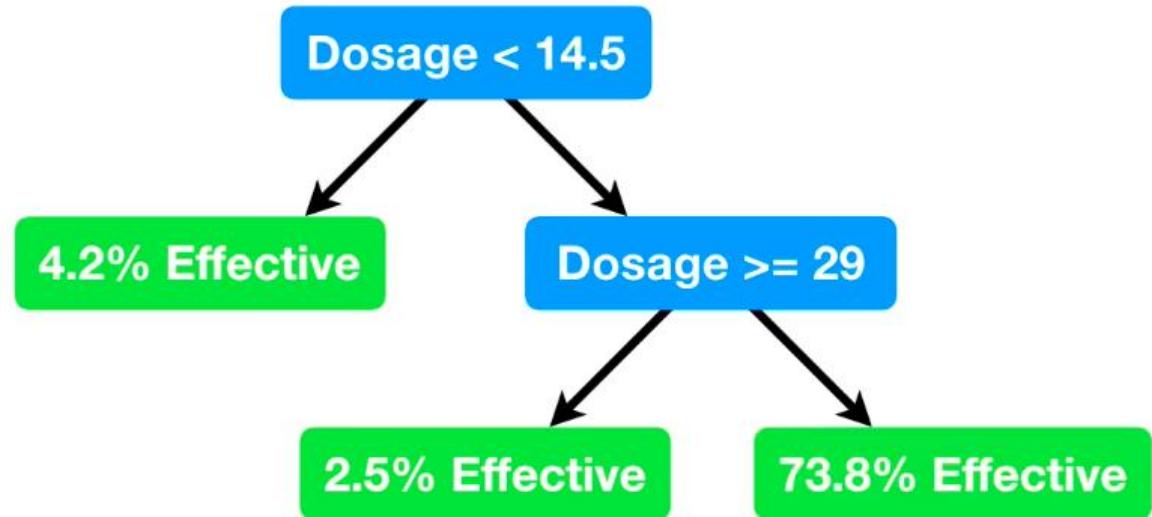
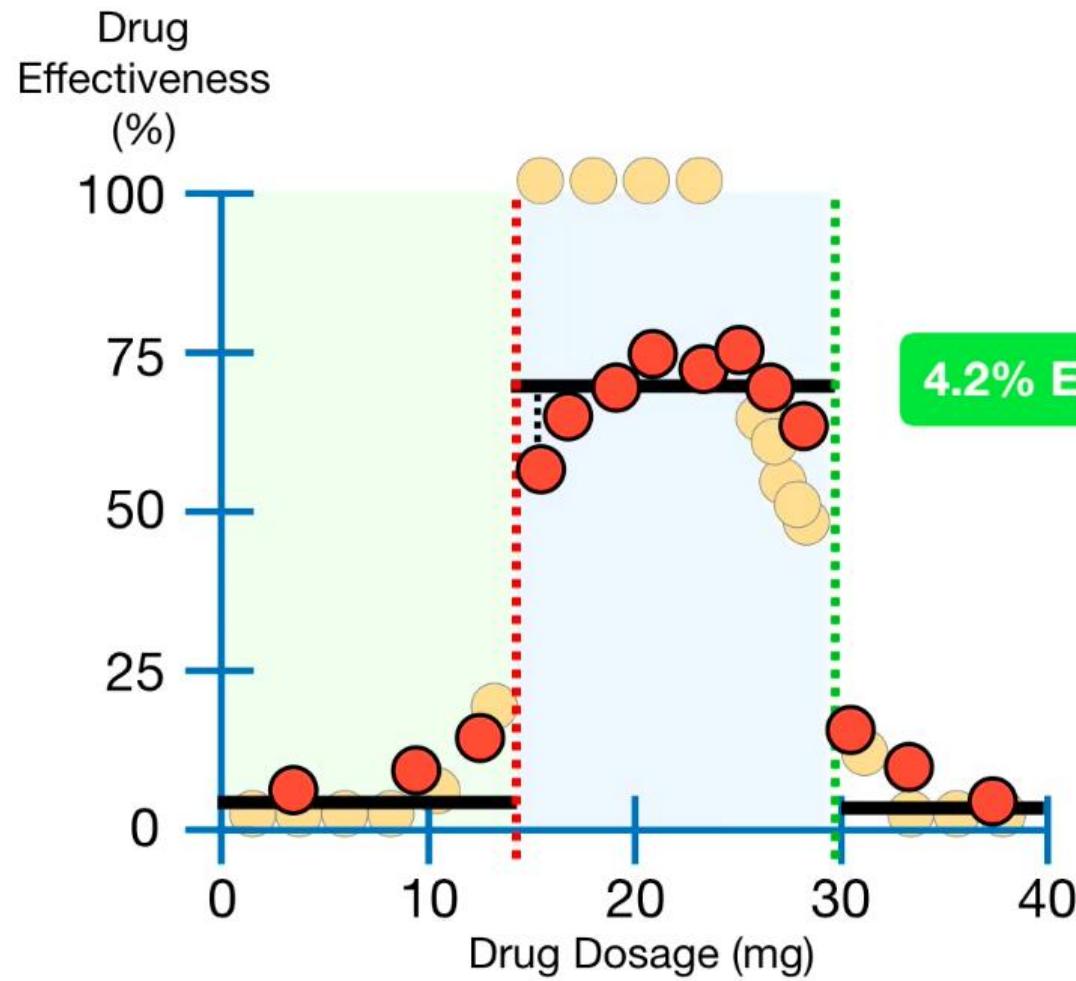
...but the new sub-tree does a much better job with the **Testing Data**.



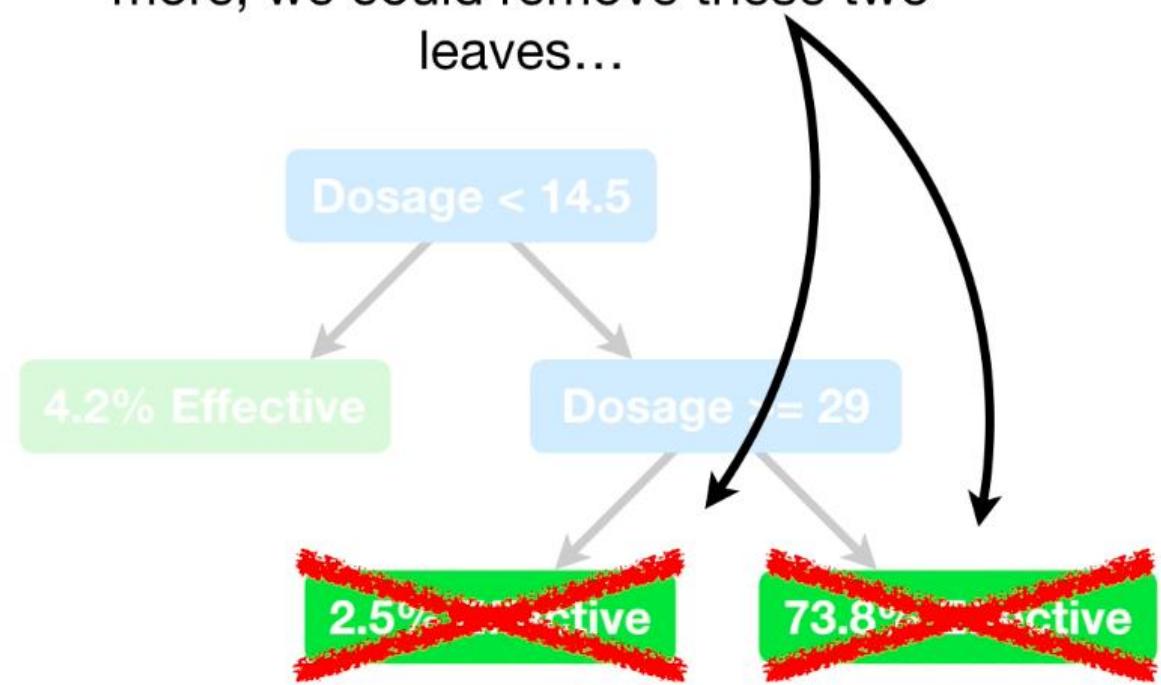
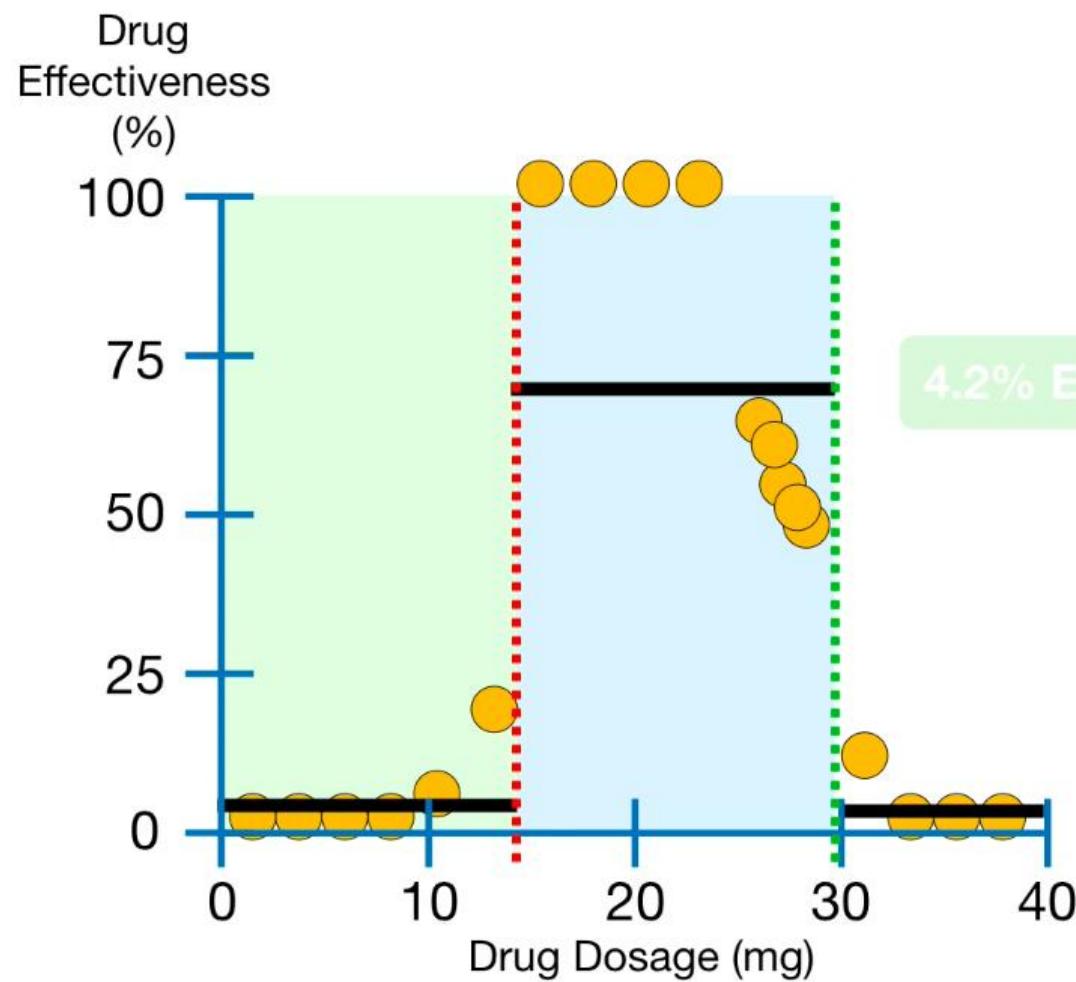
Thus, the main idea behind pruning a **Regression Tree** is to prevent overfitting the **Training Data**...



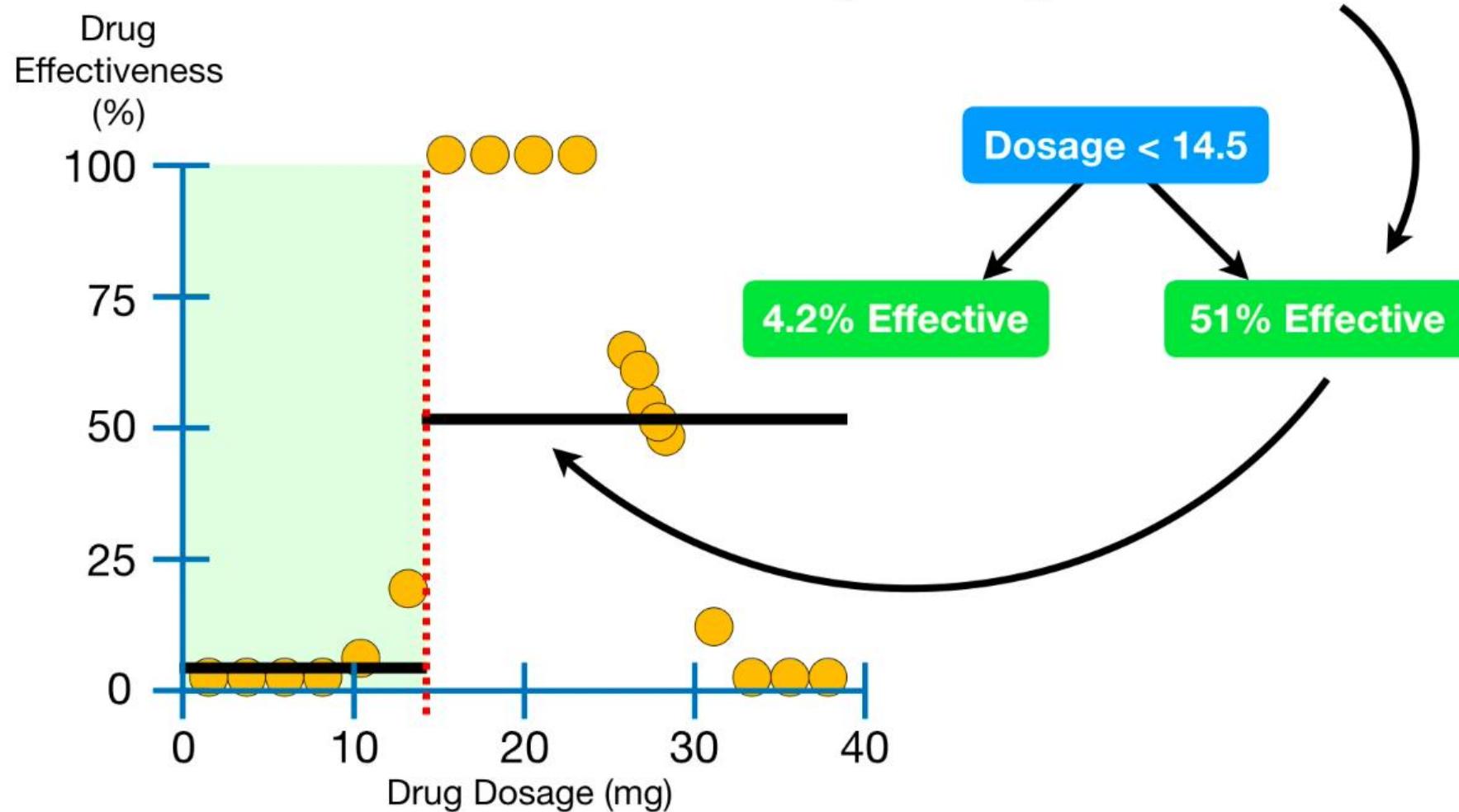
...so that the tree will do a better job
with the **Testing Data**.



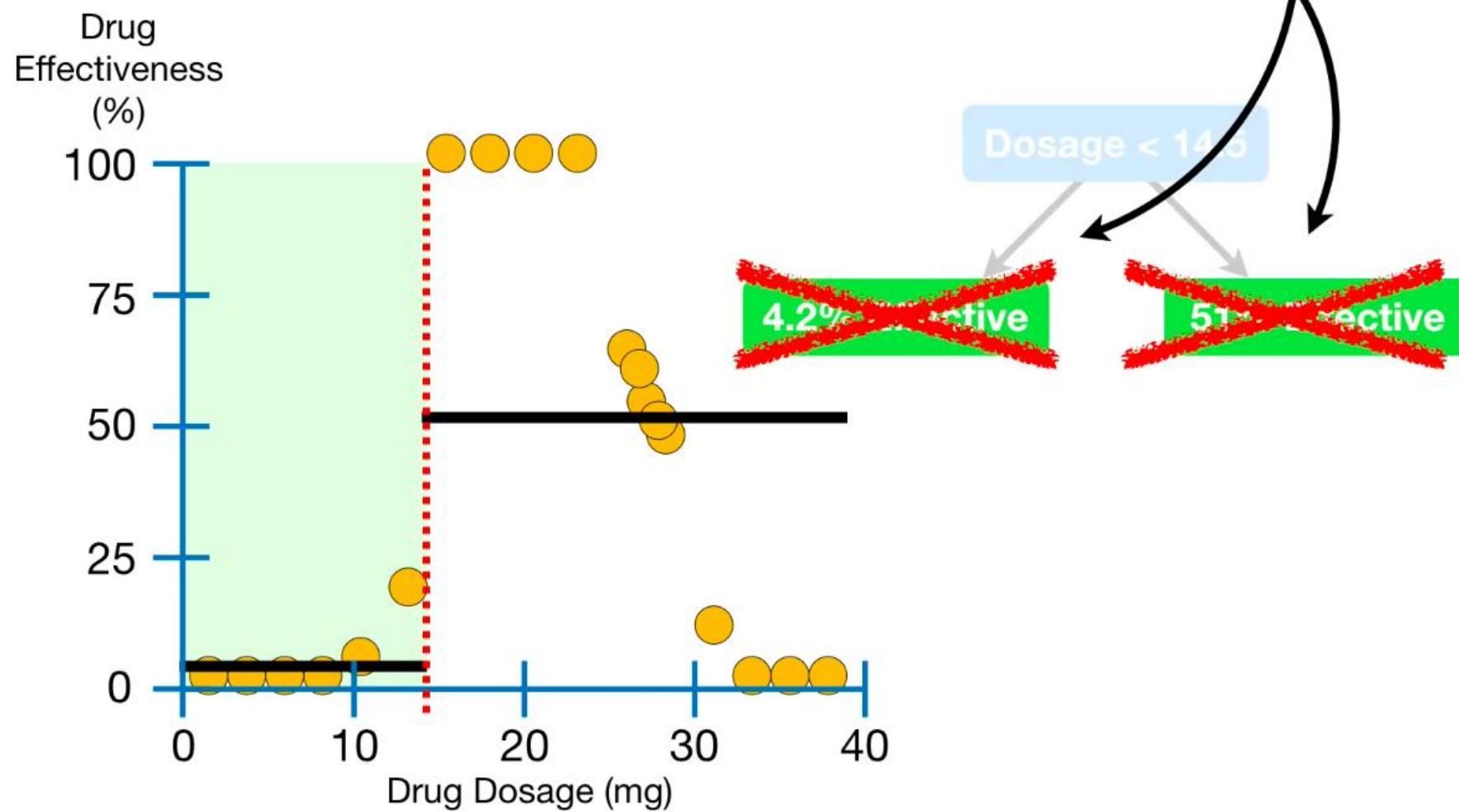
NOTE: If we wanted to prune the tree more, we could remove these two leaves...



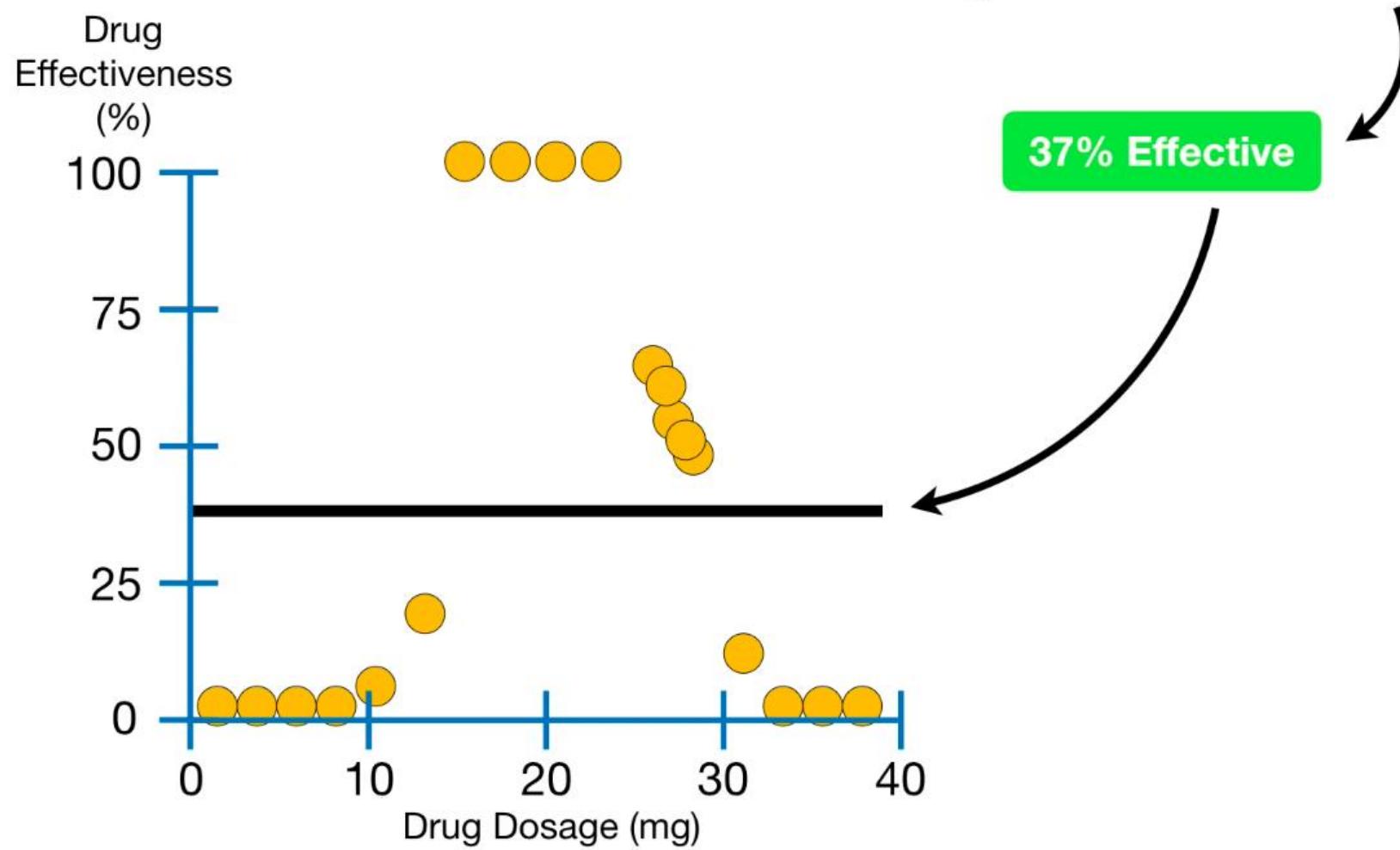
...and replace the split with a leaf that is the average of a larger number of observations.



And we could then remove these two leaves...

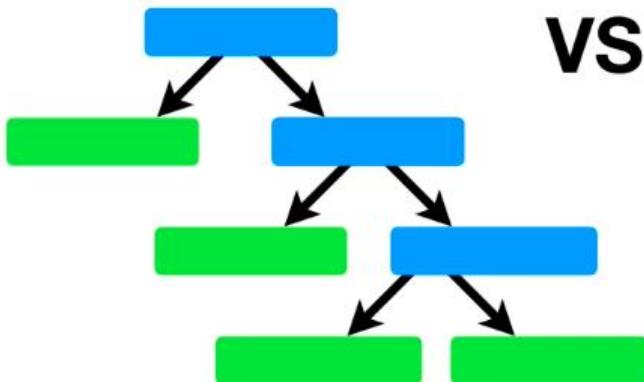


...and replace the split with a leaf that is the average of all of the observations.

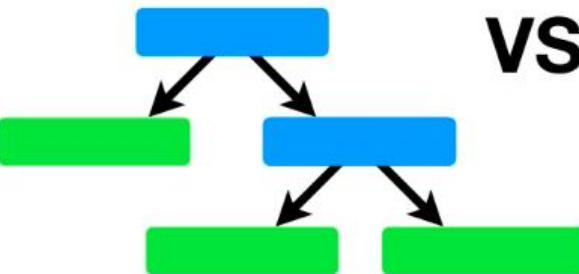


So the question is:

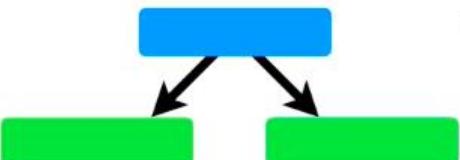
“How do we decide which tree to use?”



VS



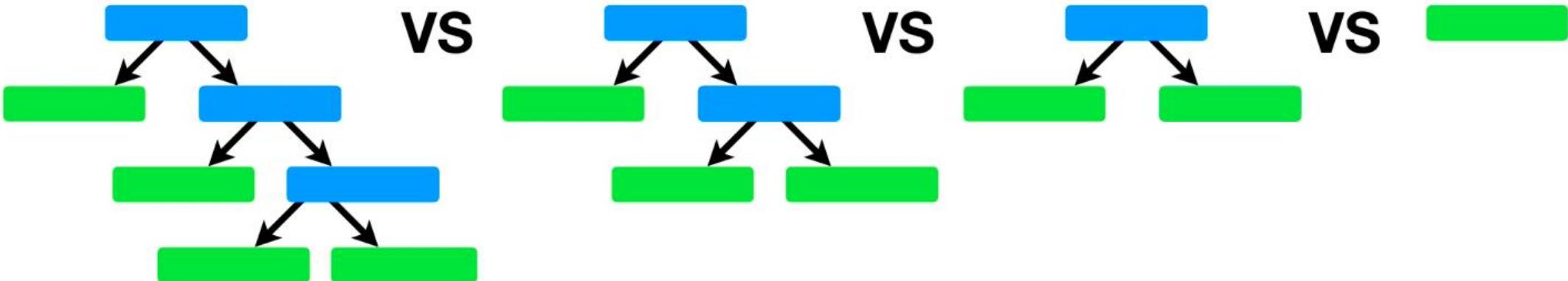
VS



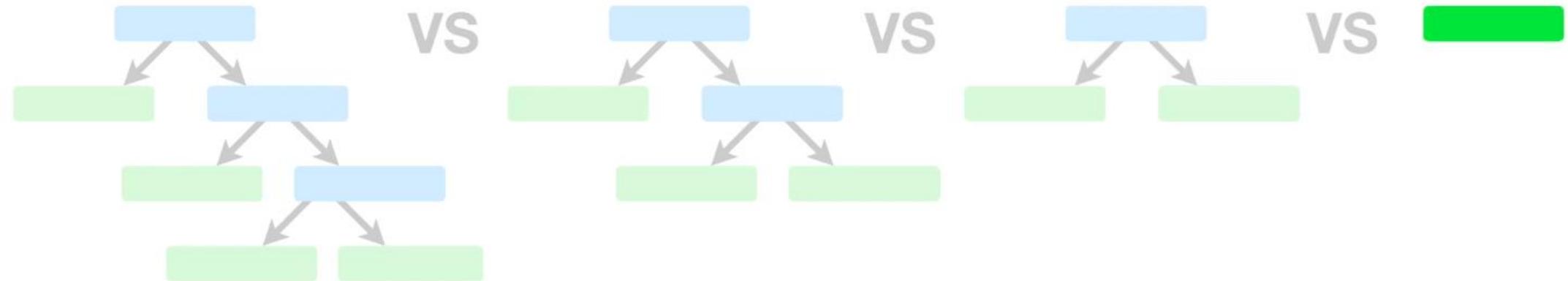
VS



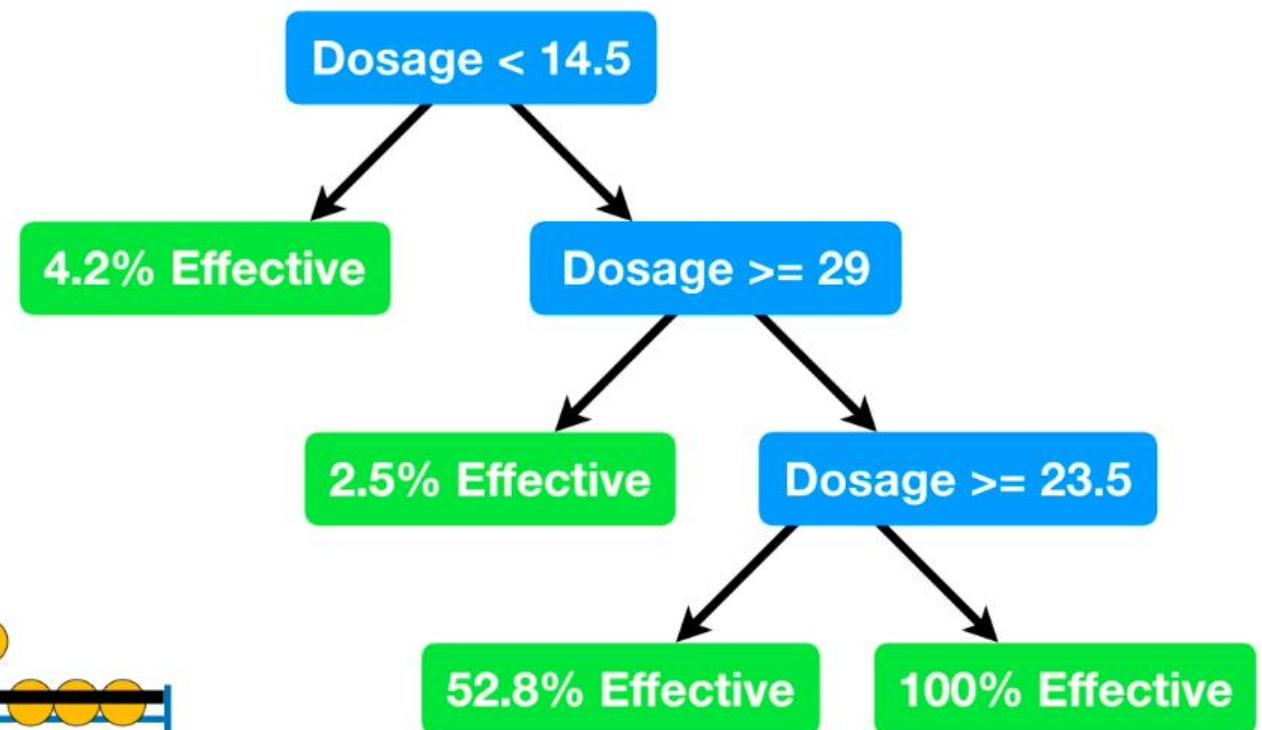
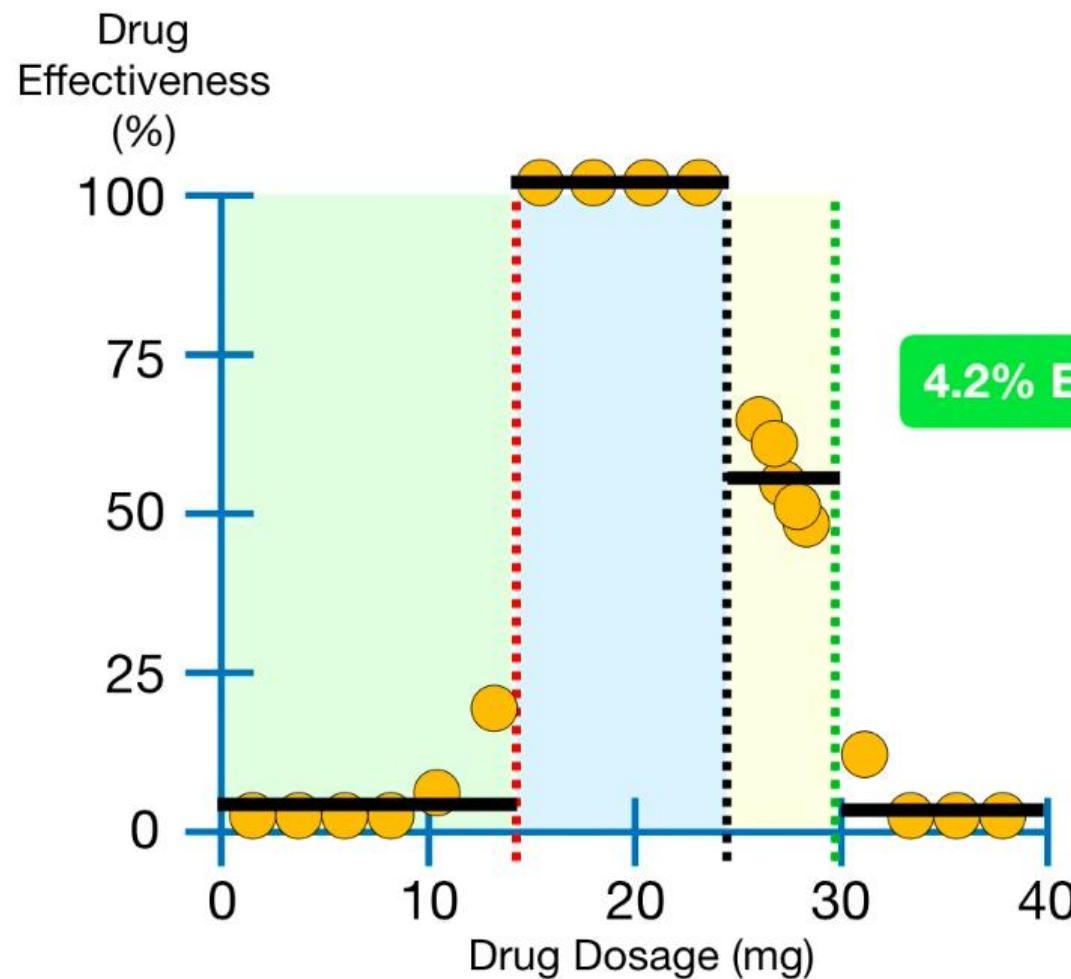
In this **StatQuest**, we will answer that question with **Cost Complexity Pruning**.



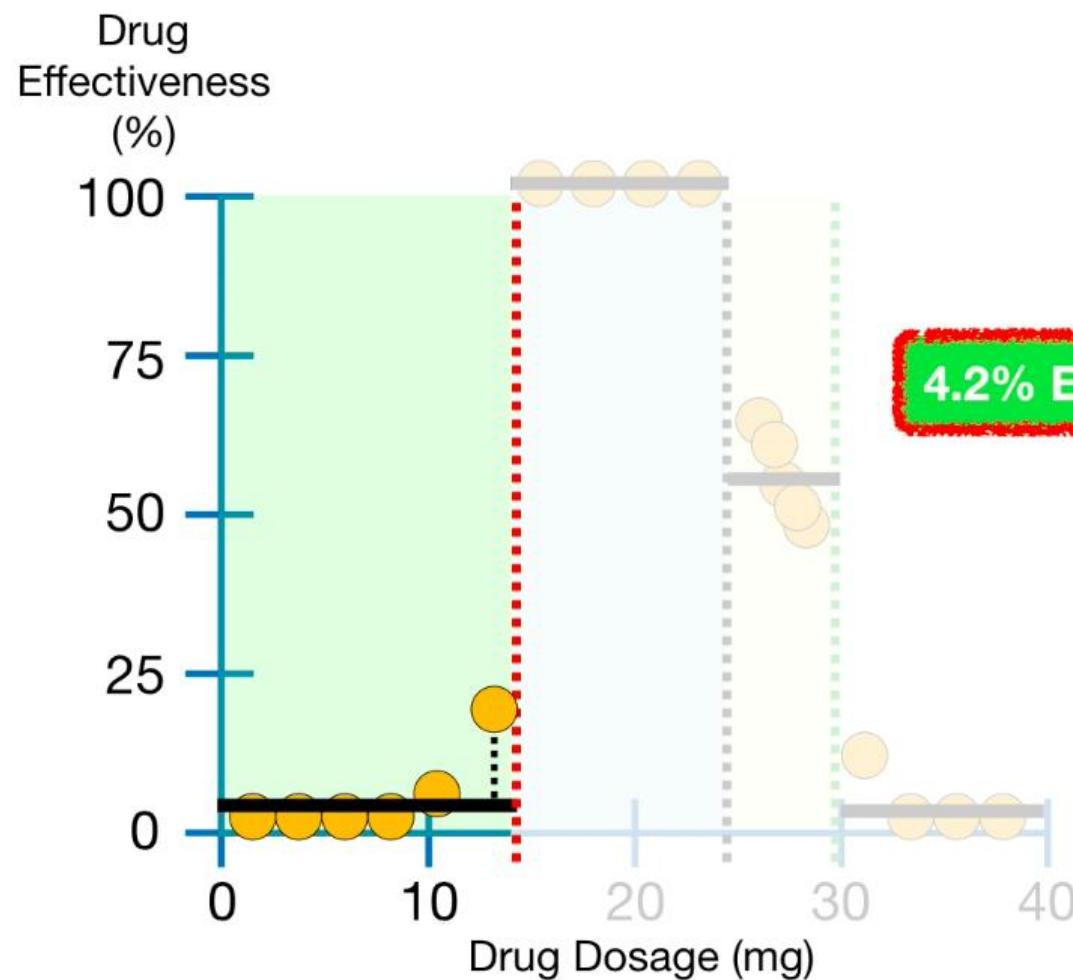
The first step in **Cost Complexity Pruning** is to calculate the **Sum of Squared Residuals** for each tree.



Here is the original, full sized tree.



The **Sum of Squared Residuals** for the **Observations with Dosages < 14.5** is...



Dosage < 14.5

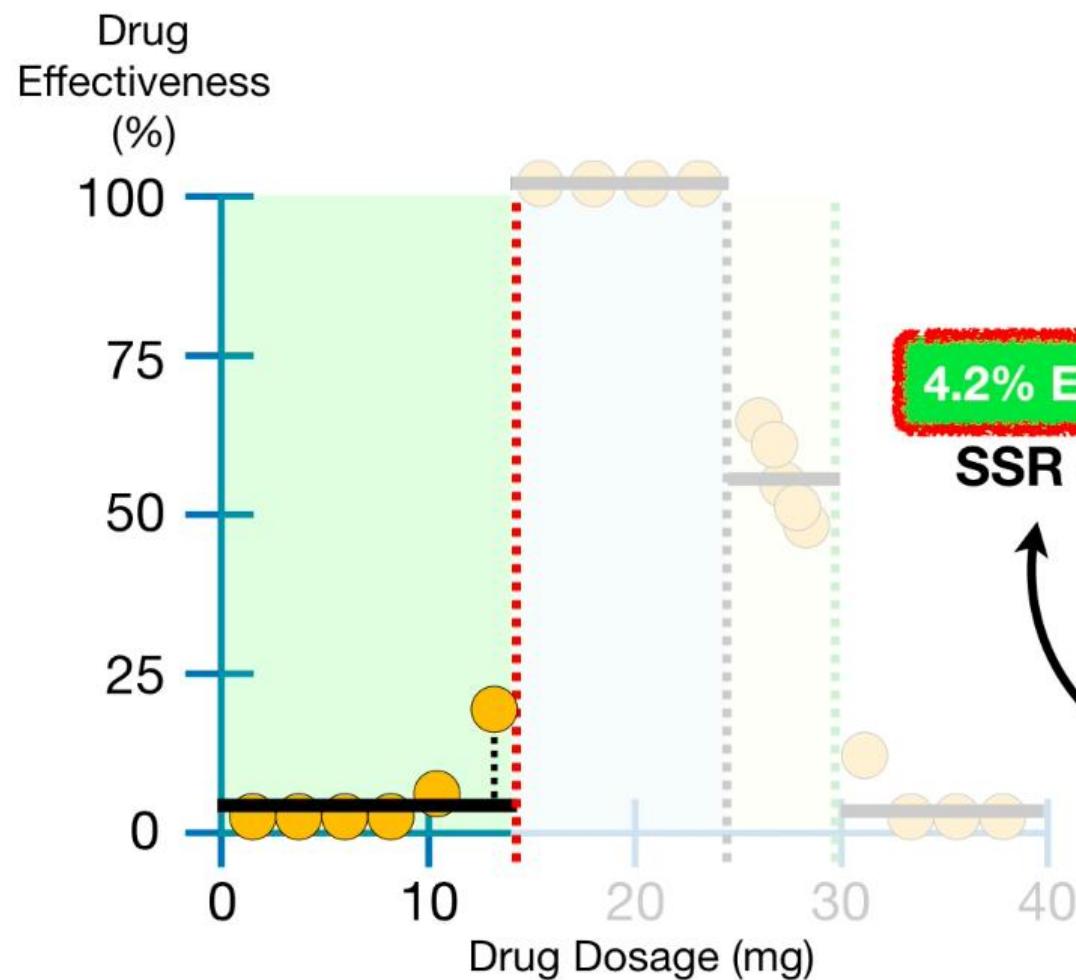
4.2% Effective

$$(0 - 4.2)^2 + (0 - 4.2)^2 + (0 - 4.2)^2 + (0 - 4.2)^2$$

$$+ (5 - 4.2)^2 + (20 - 4.2)^2$$

$$= 320.8$$

So we'll save that **Sum of Squared Residuals (SSR)** underneath the corresponding leaf.



Dosage < 14.5

4.2% Effective

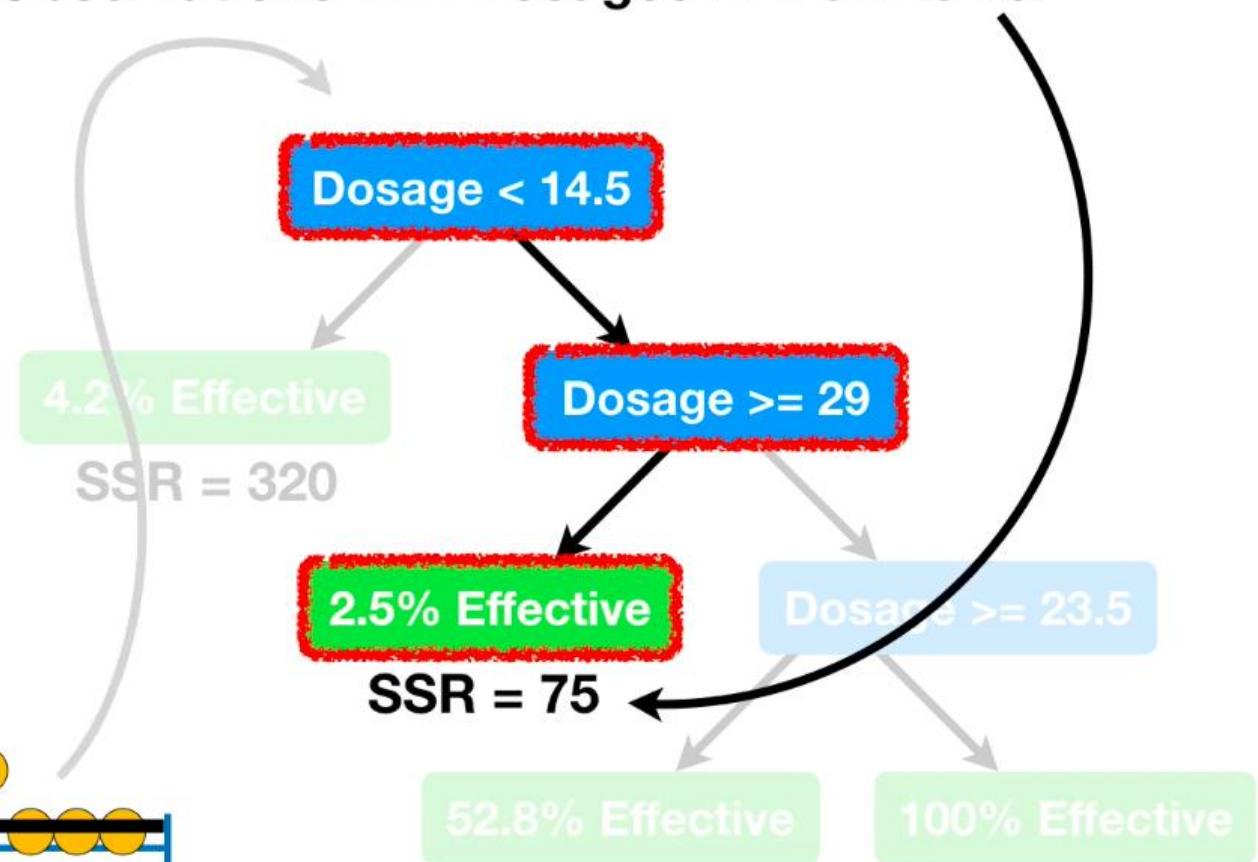
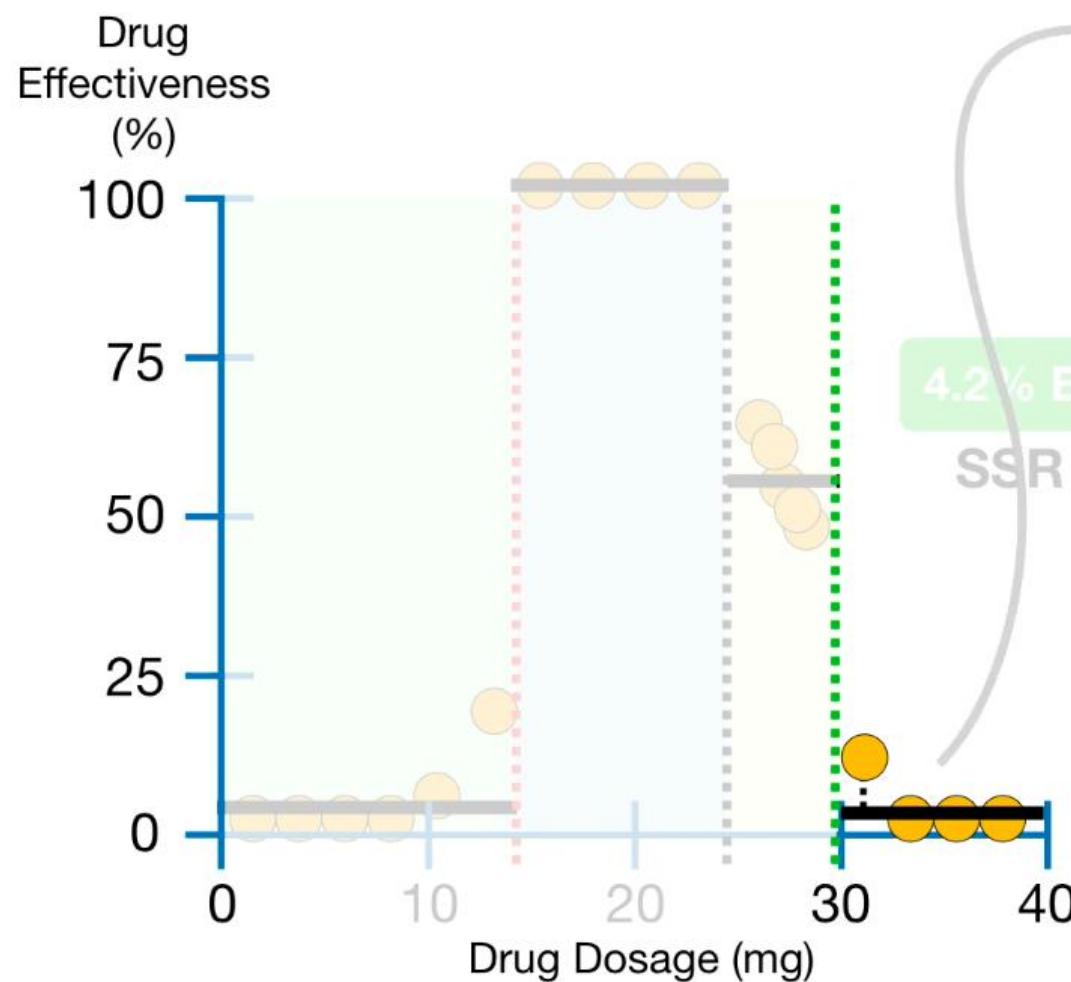
$$SSR = 320$$

$$(0 - 4.2)^2 + (0 - 4.2)^2 + (0 - 4.2)^2 + (0 - 4.2)^2$$

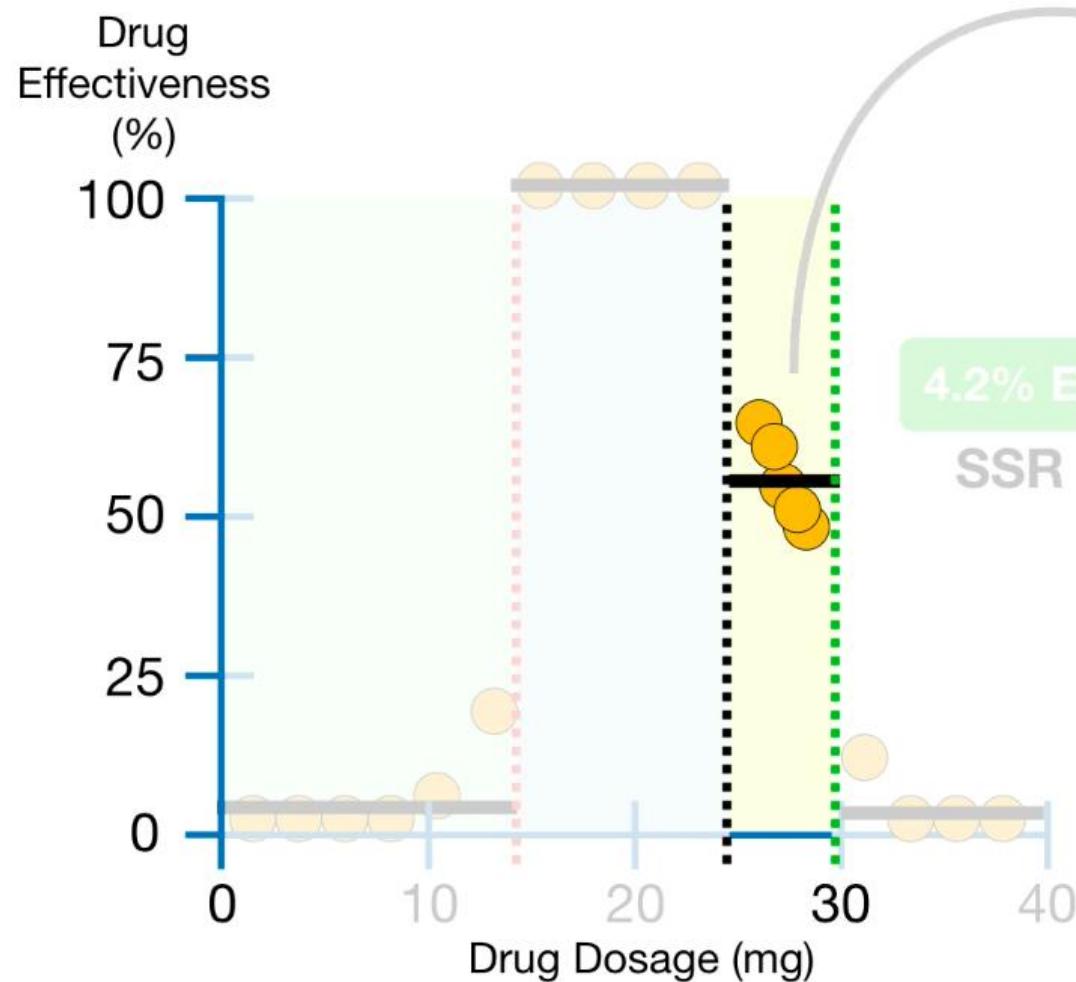
$$+ (5 - 4.2)^2 + (20 - 4.2)^2$$

$$= 320.8$$

The Sum of Squared Residuals for Observations with Dosages ≥ 29 ... is 75.



The Sum of Squared Residuals for Observations with Dosages ≥ 23.5 and < 29 ... is 148.8.



Dosage < 14.5

4.2% Effective
SSR = 320

Dosage ≥ 29

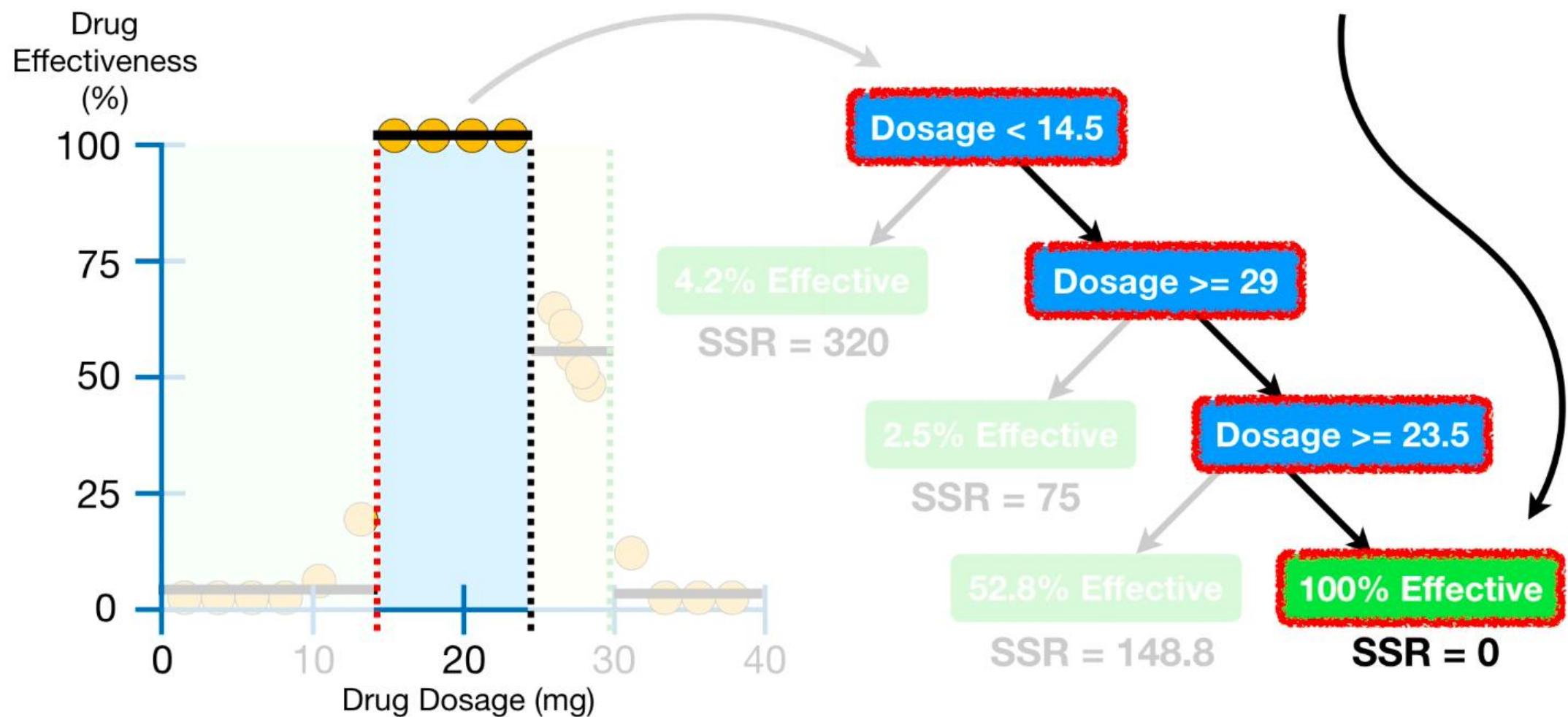
2.5% Effective
SSR = 75

Dosage ≥ 23.5

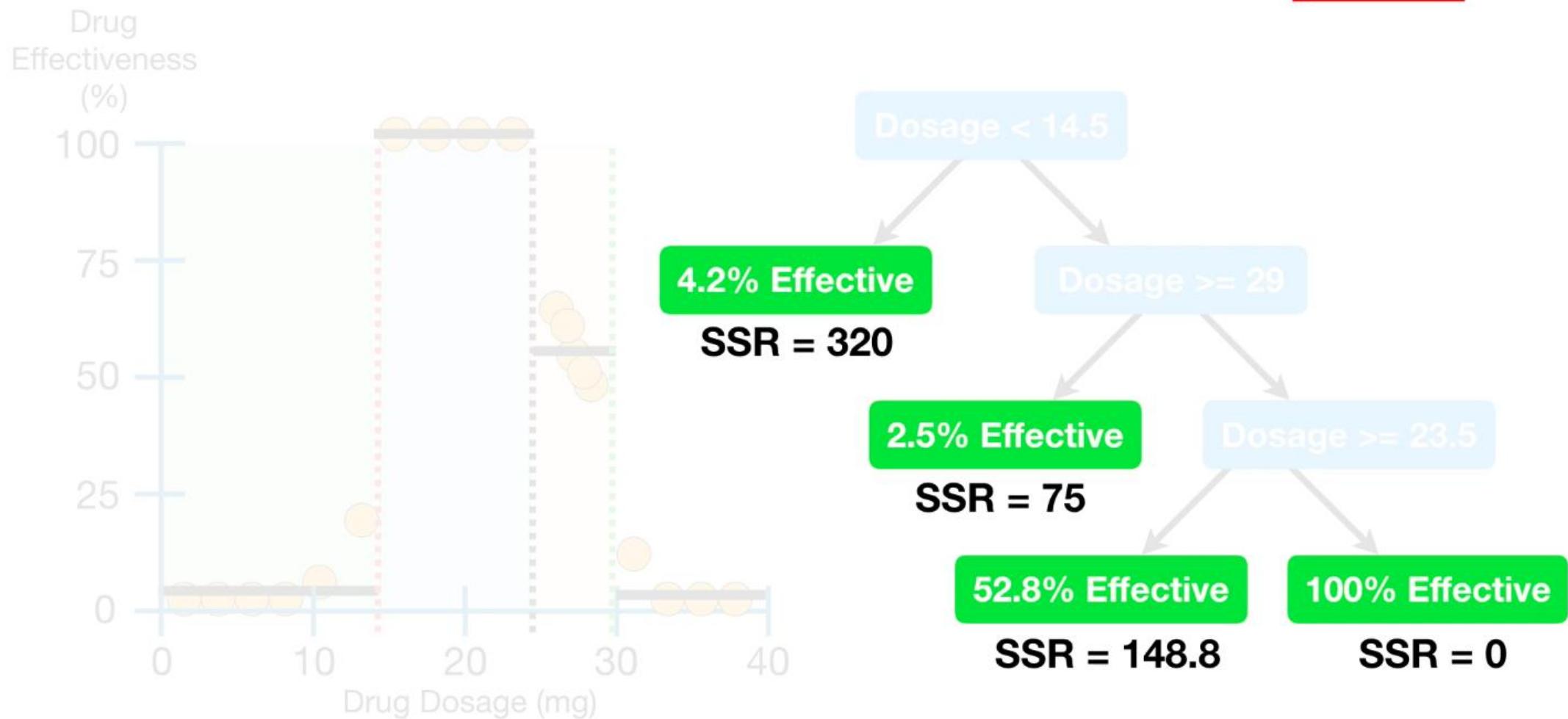
52.8% Effective
SSR = 148.8

100% Effective

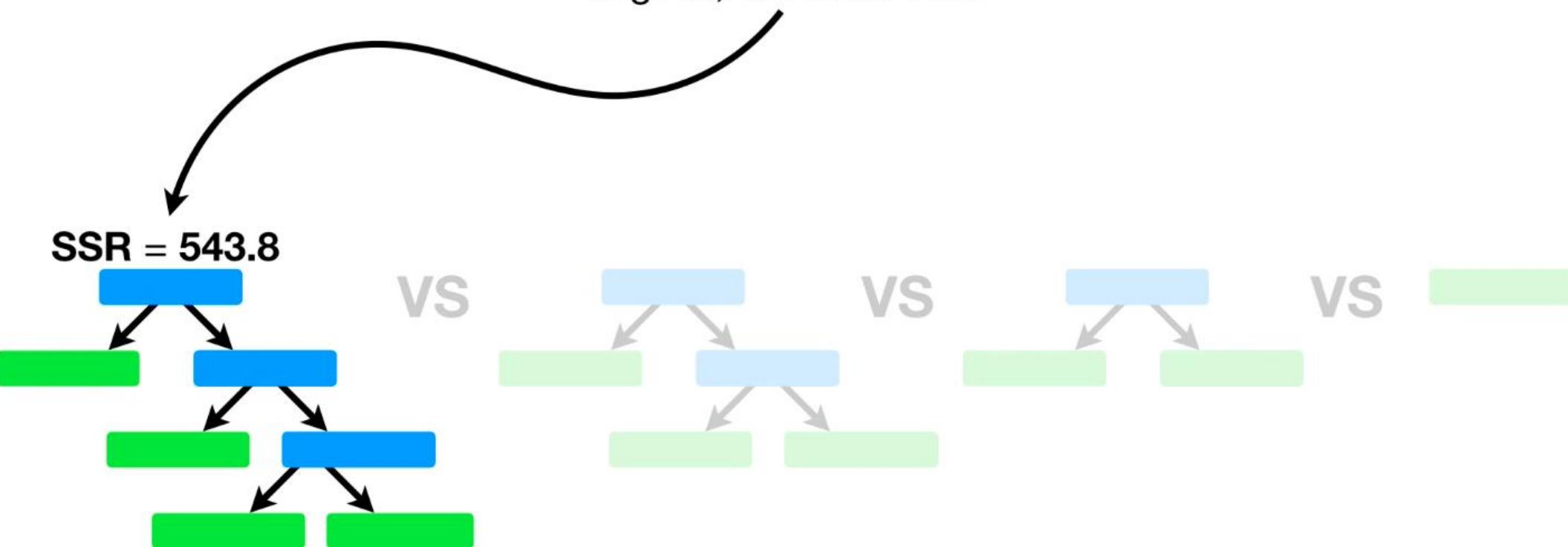
And the **Sum of Squared Residuals** for Observations with **Dosages ≥ 14.5** and $< 23.5\dots$ is 0.



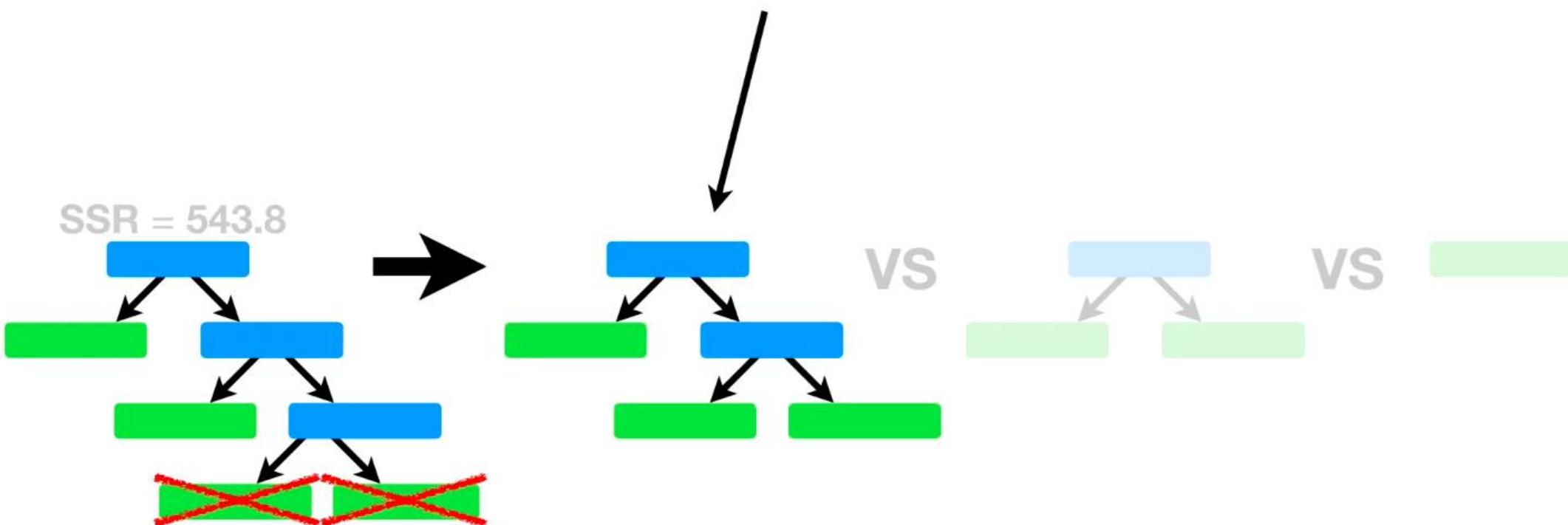
Thus, the total **Sum of Squared Residuals** (SSR) for the whole tree is $320 + 75 + 148.8 + 0 = \boxed{543.8}$



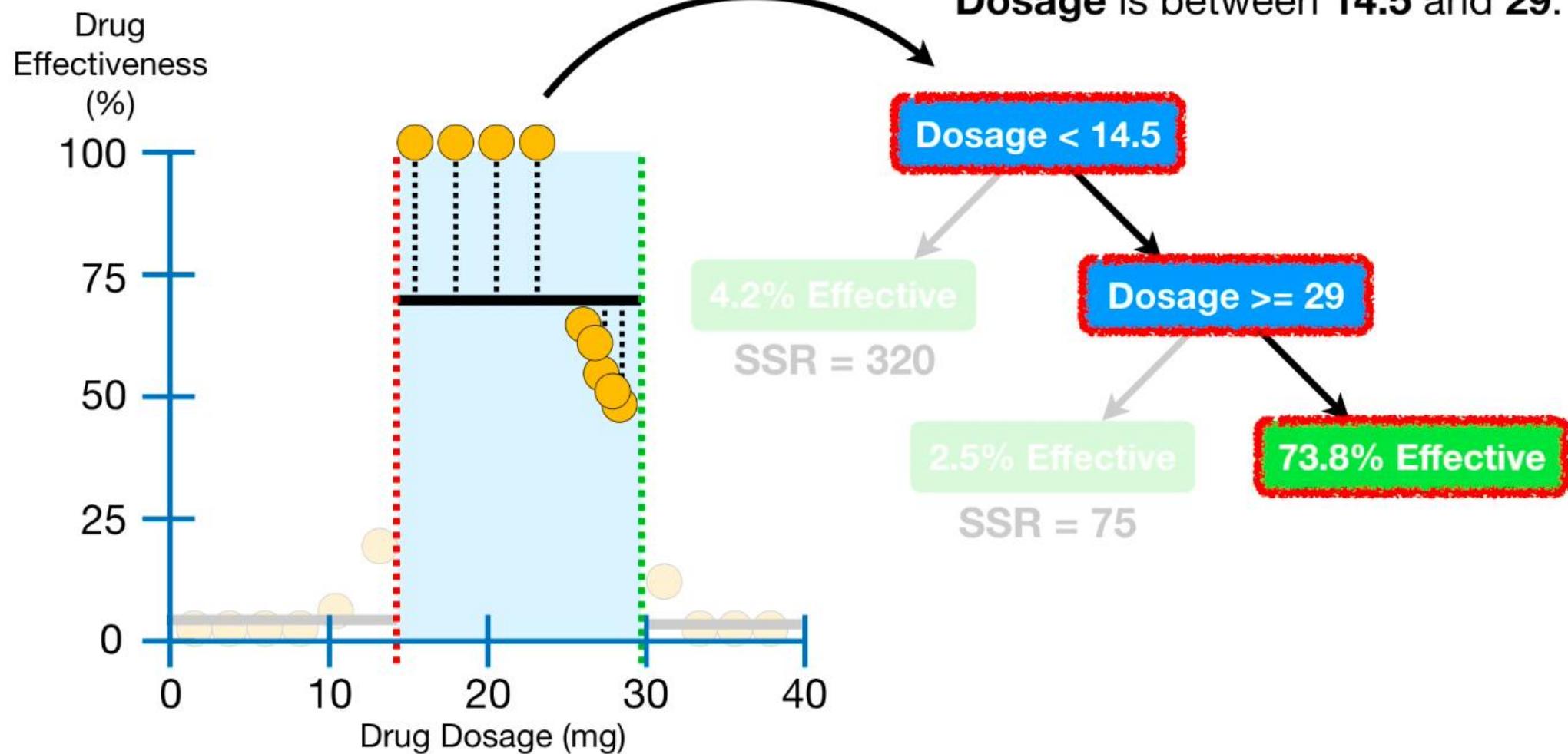
So let's put **SSR = 543.8** on top of the original, full sized tree.



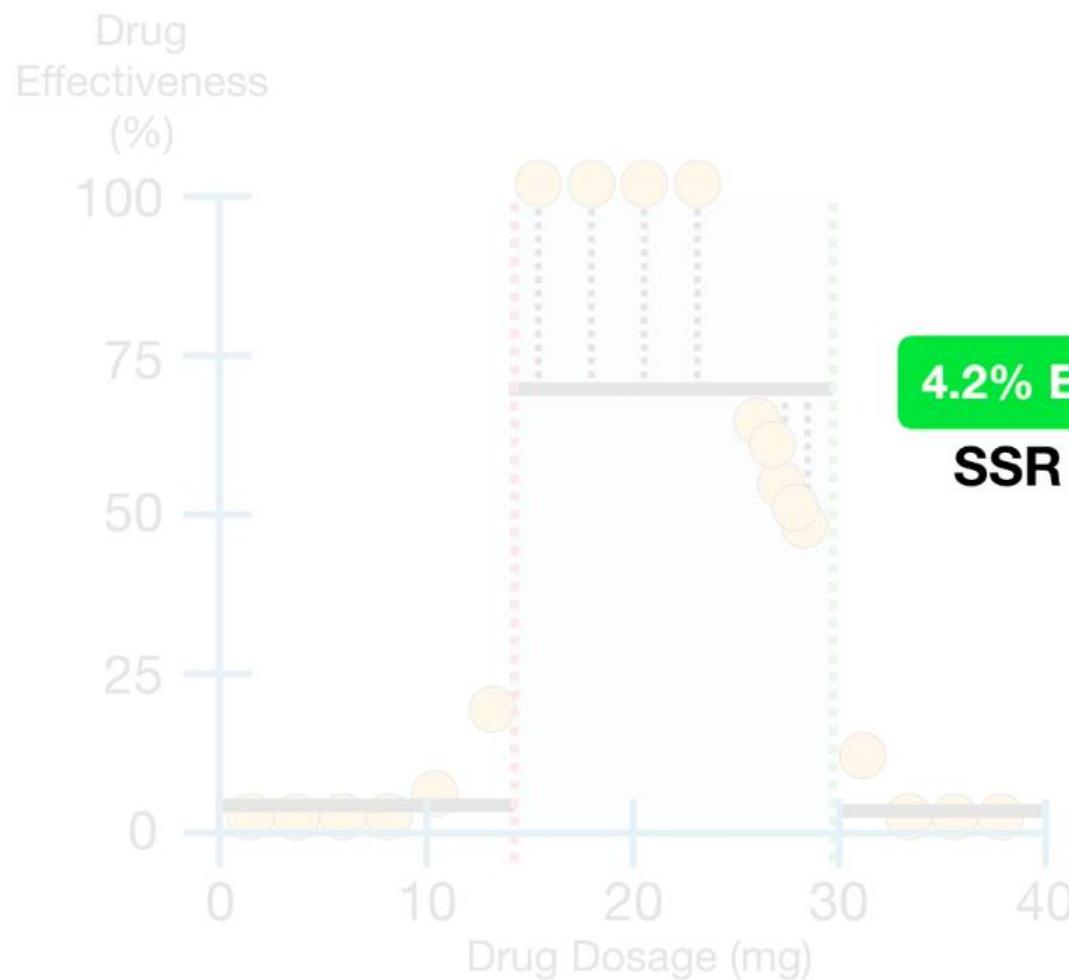
Now let's calculate the **Sum of Squared Residuals** for the sub-tree with one fewer leaf.



...but we have to calculate a new **Sum of Squared Residuals** for when the **Dosage** is between **14.5** and **29**.



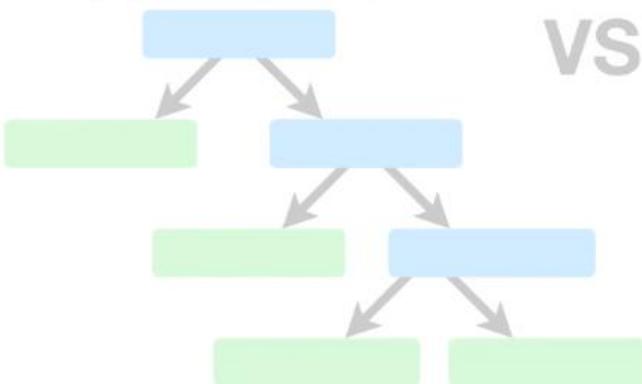
Thus, the total **Sum of Squared Residuals** for this tree is $320 + 75 + 5099.8 = \boxed{5494.8}$



So let's put **SSR = 5494.8** on top
of the sub-tree with **3** leaves.

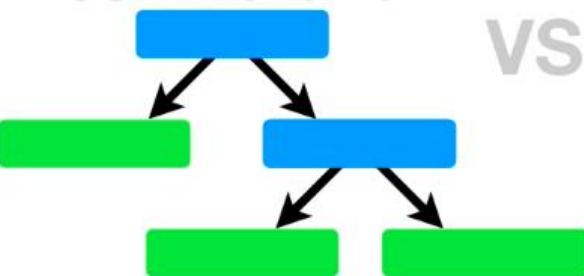


SSR = 543.8



VS

SSR = 5494.8

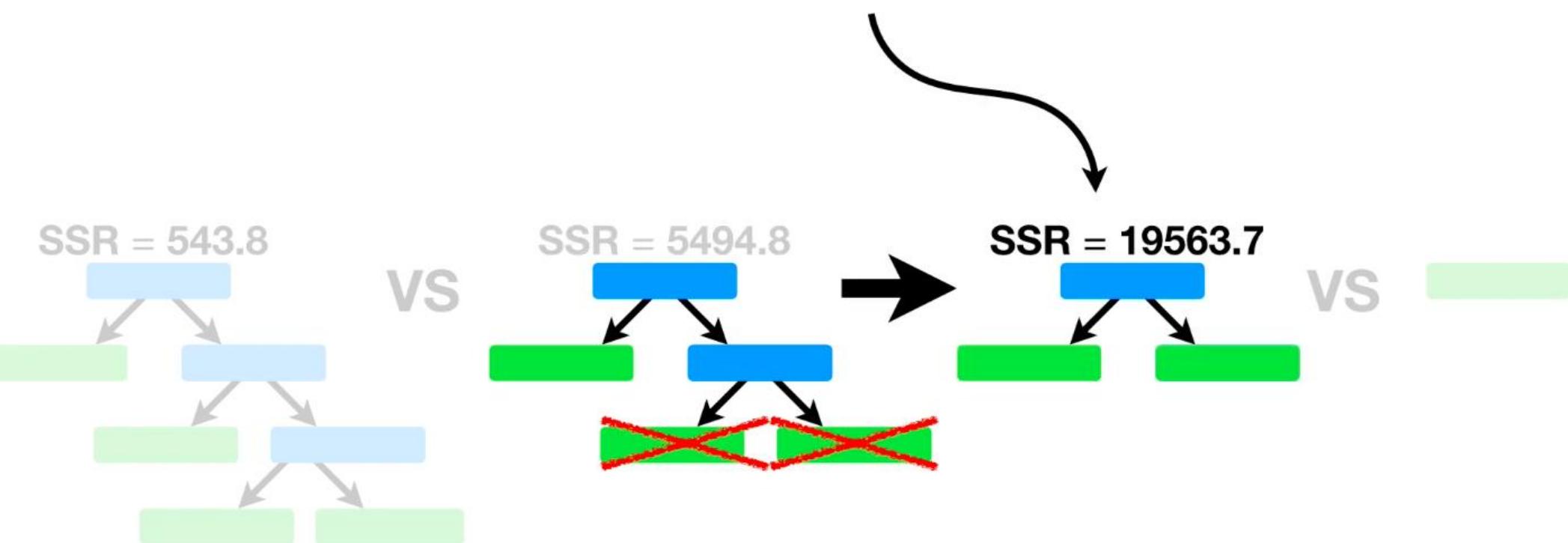


VS



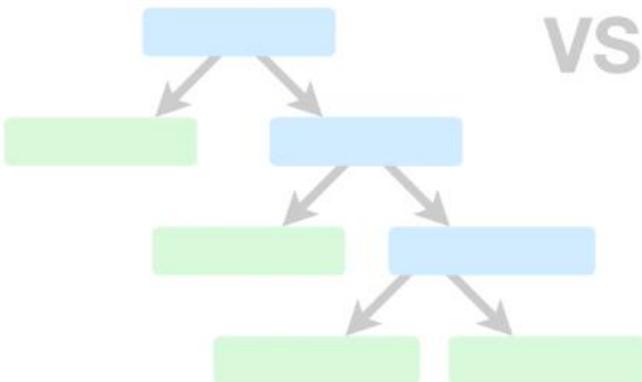
VS

Similarly, the **Sum of Squared Residuals** for
the sub-tree with **2 leaves**...



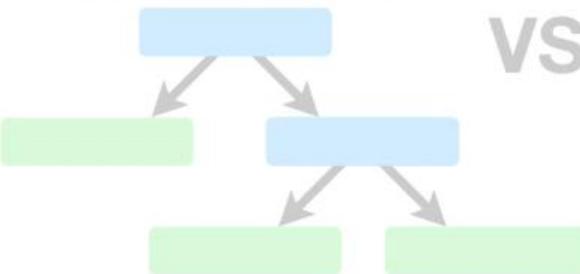
Lastly, the **Sum of Squared Residuals** for the
sub-tree with only one leaf...

$SSR = 543.8$



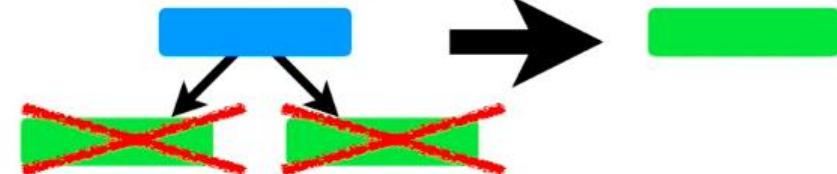
VS

$SSR = 5494.8$



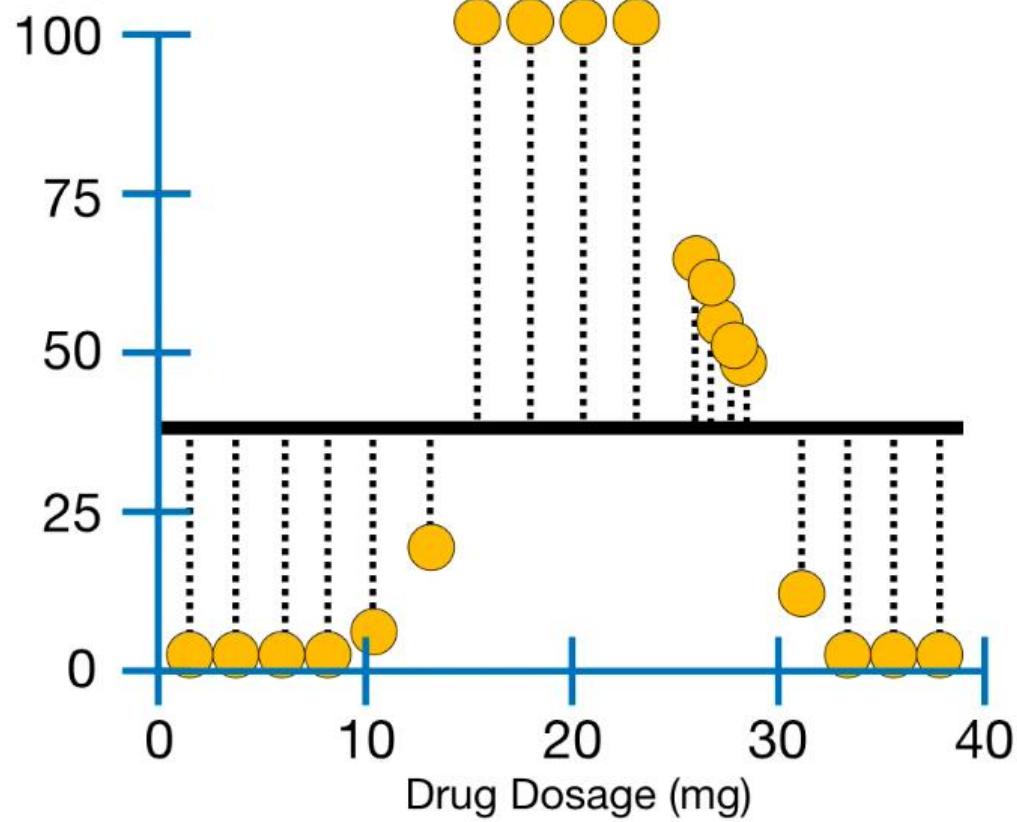
VS

$SSR = 19243.7$



...is **28,897.2**.

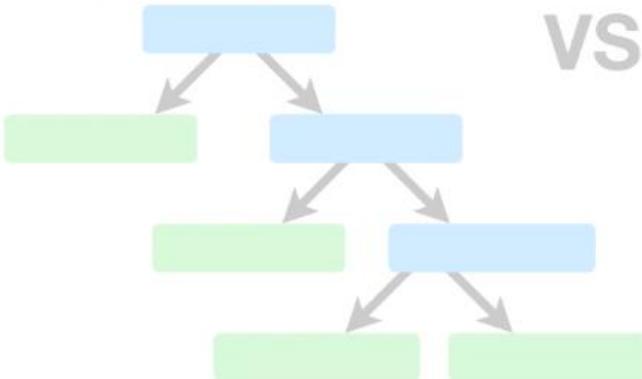
Drug
Effectiveness
(%)



37% Effective

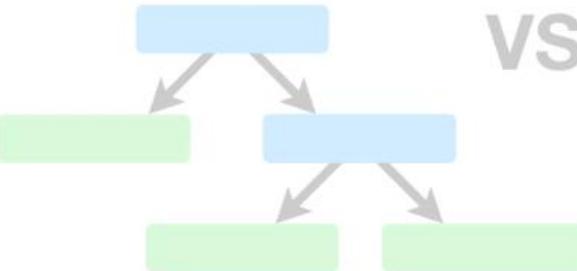
...so let's put **SSR = 28897.2** on top of the
sub-tree with **1 leaf**.

SSR = 543.8



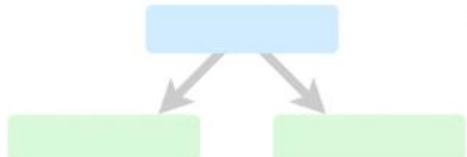
VS

SSR = 5494.8



VS

SSR = 19243.7

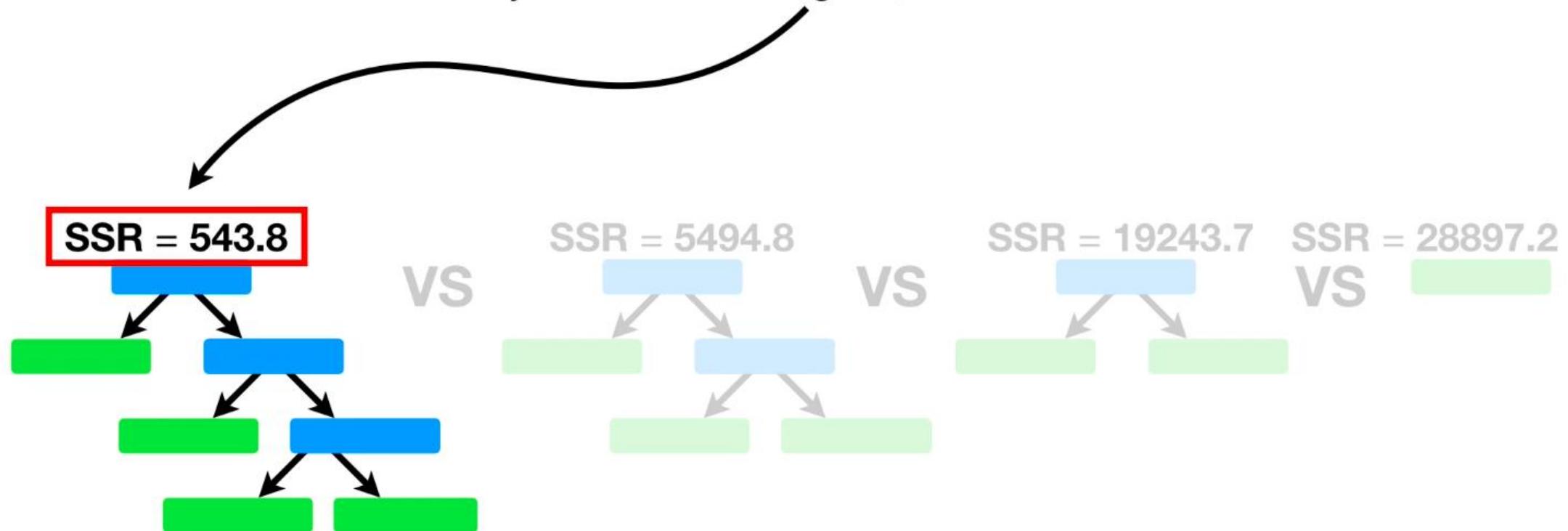


SSR = 28897.2

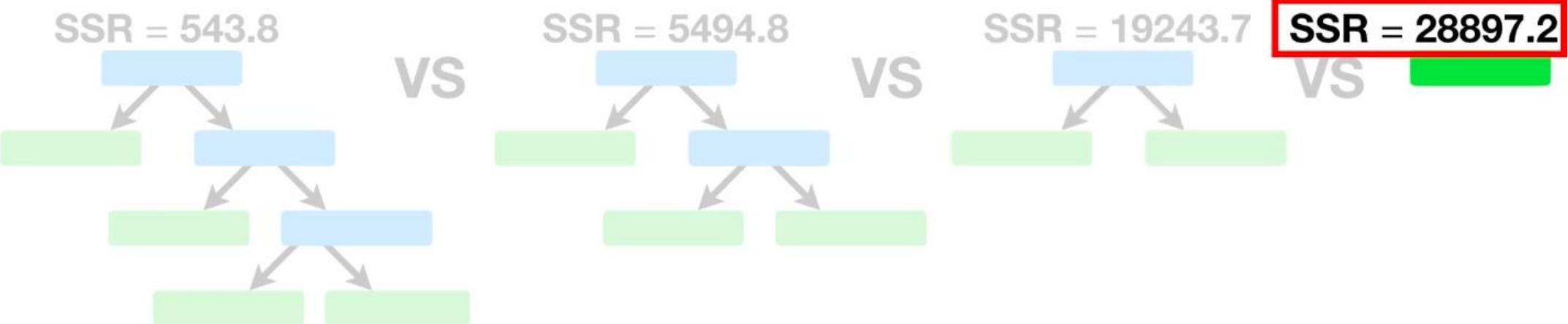
VS



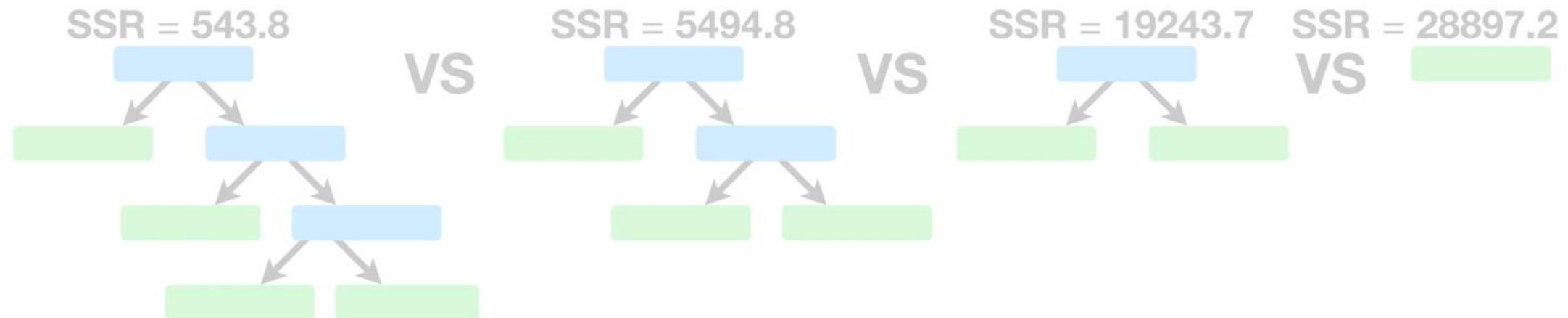
NOTE: The Sum of Squared Residuals is
relatively small for the original, full sized tree...



...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.

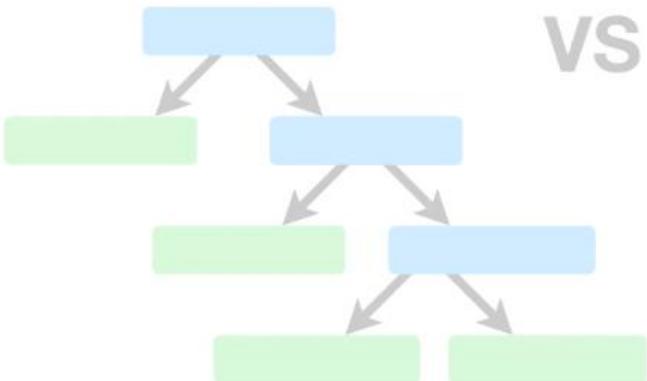


Weakest Link Pruning works by calculating a
Tree Score...



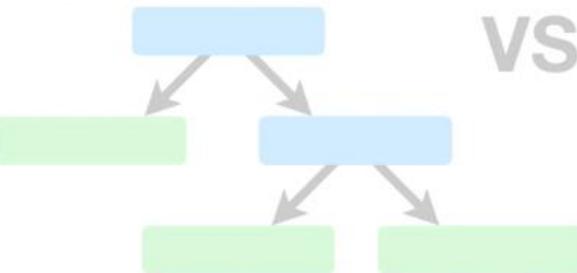
Tree Score

$SSR = 543.8$



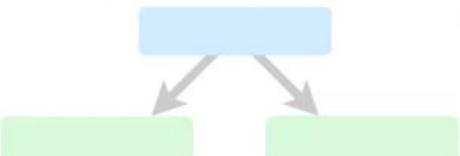
VS

$SSR = 5494.8$



VS

$SSR = 19243.7$



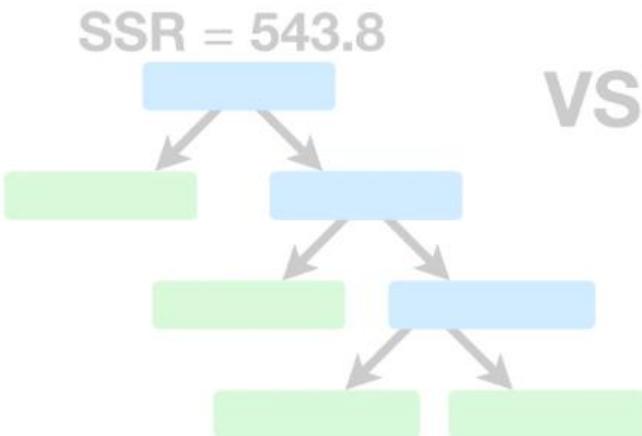
$SSR = 28897.2$

VS

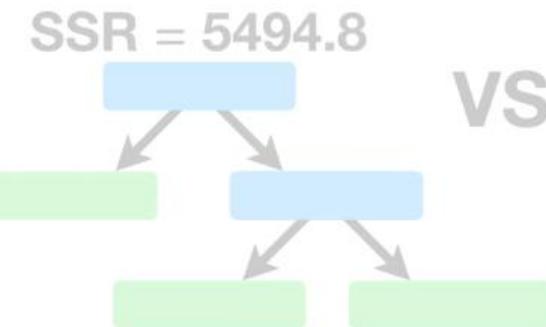
...that is based on the **Sum of Squared Residuals (SSR)** for the tree or sub-tree...



Tree Score = SSR



VS



VS



$SSR = 28897.2$

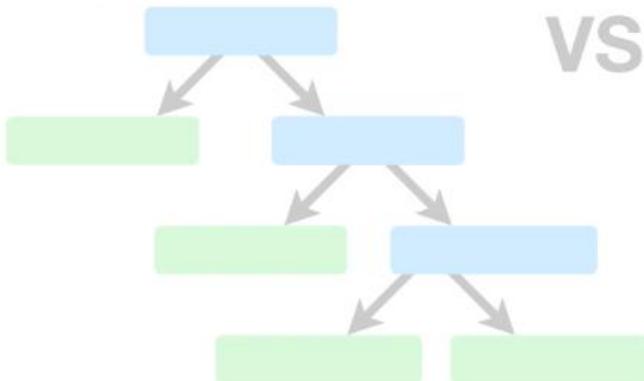


```
graph TD; Root[SSR = 28897.2] --> Node1[ ]; Node1 --> Leaf1[ ]; Node1 --> Node2[ ]; Node2 --> Leaf2[ ]; Node2 --> Node3[ ]; Node3 --> Leaf3[ ]; Node3 --> Leaf4[ ]
```

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

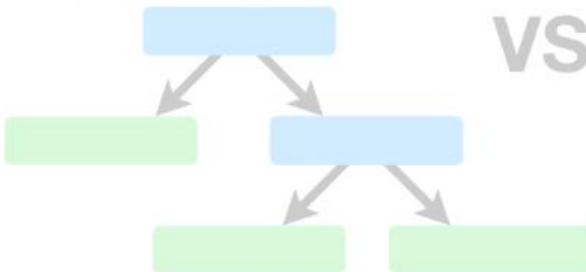
$$\text{Tree Score} = \text{SSR} + aT$$

$$\text{SSR} = 543.8$$



VS

$$\text{SSR} = 5494.8$$



VS

$$\text{SSR} = 19243.7$$



$$\text{SSR} = 28897.2$$

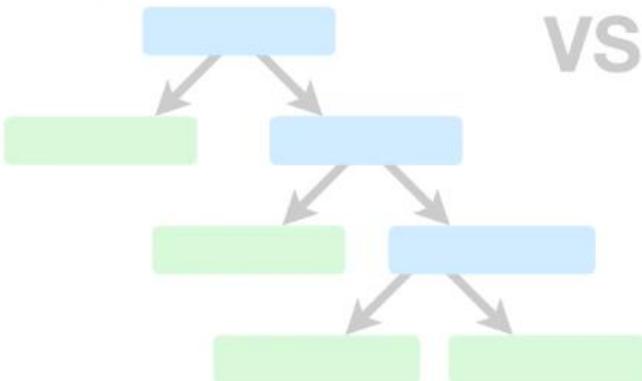
VS

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



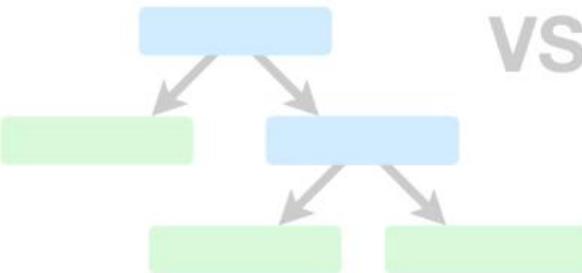
$$\text{Tree Score} = \text{SSR} + \boxed{aT}$$

$\text{SSR} = 543.8$



VS

$\text{SSR} = 5494.8$



VS

$\text{SSR} = 19243.7$

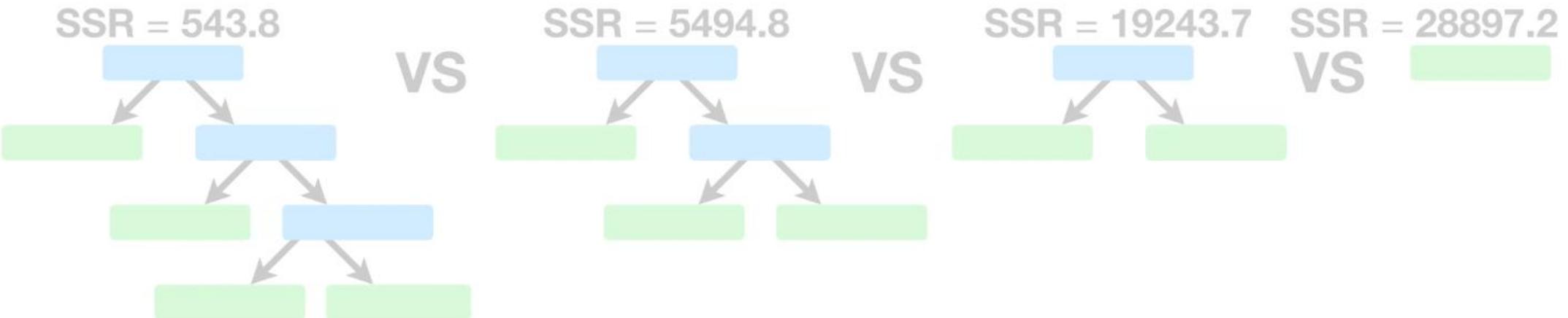


$\text{SSR} = 28897.2$

VS

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

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

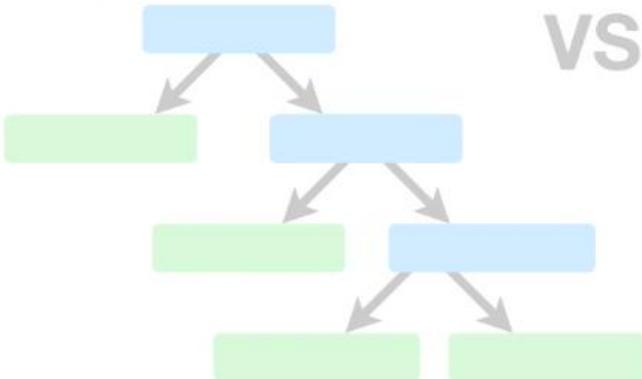


For now, let's let $a = 10,000$.



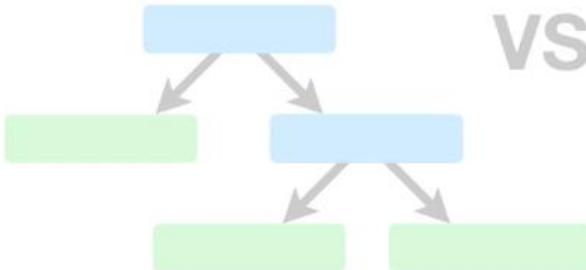
$$\text{Tree Score} = \text{SSR} + 10,000 \times T$$

$$\text{SSR} = 543.8$$



VS

$$\text{SSR} = 5494.8$$



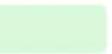
VS

$$\text{SSR} = 19243.7$$



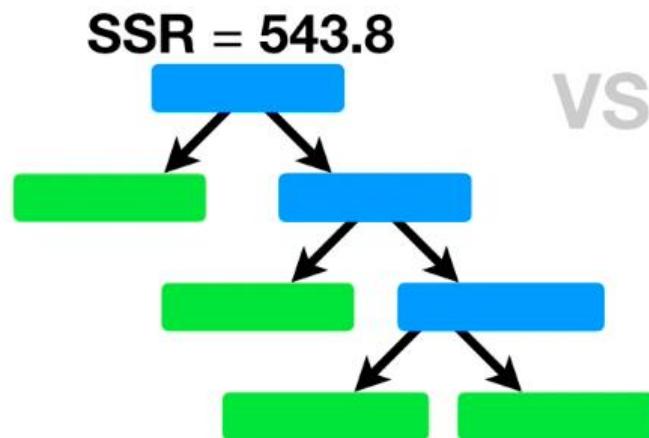
$$\text{SSR} = 28897.2$$

VS

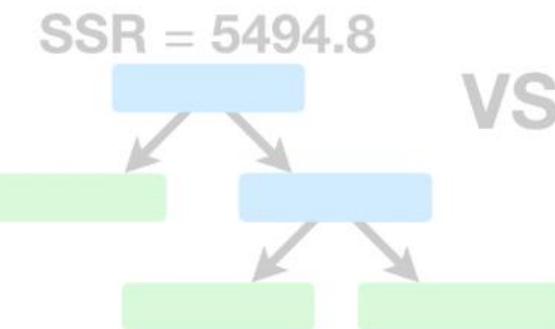


So the **Tree Score** for the original,
full sized tree is **40,543.8**.

$$\text{Tree Score} = 543.8 + 10,000 \times 4 = 40,543.8$$



VS



VS



SSR = 28897.2

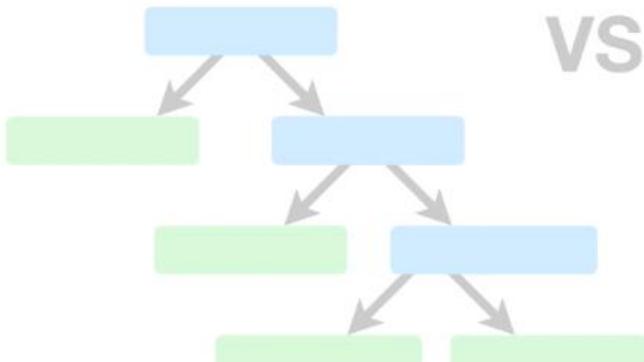
VS

```
graph TD; Root[SSR = 28897.2] --> Node1[ ]; Node1 --> Leaf1[ ]; Node1 --> Node2[ ]; Node2 --> Leaf2[ ]
```

...and the total
Tree Score = 35,494.8.

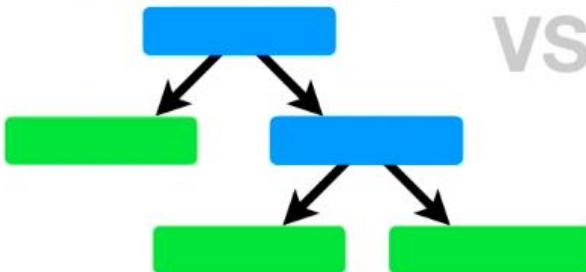
$$\text{Tree Score} = 5494.8 + \mathbf{10,000 \times 3} = 35494.8$$

$$\text{SSR} = 543.8$$

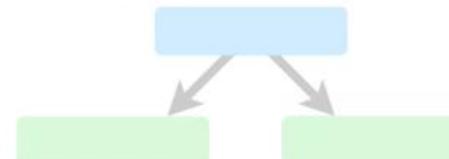


$$\text{Tree Score} = 40,543.8$$

$$\text{SSR} = 5494.8$$



$$\text{SSR} = 19243.7$$



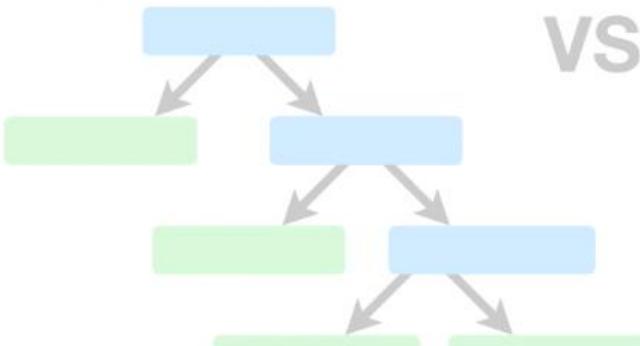
$$\text{SSR} = 28897.2$$



...39,243.7.

$$\text{Tree Score} = 19243.7 + 10,000 \times 2 = 39243.7$$

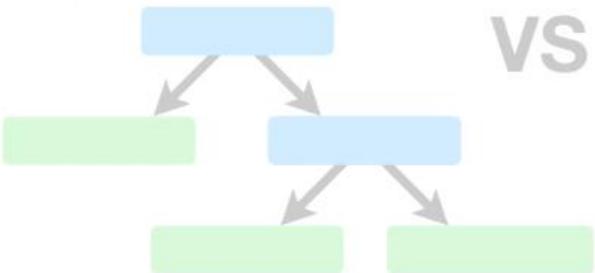
SSR = 543.8



Tree Score = 40,543.8

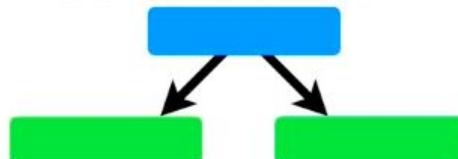
VS

SSR = 5494.8



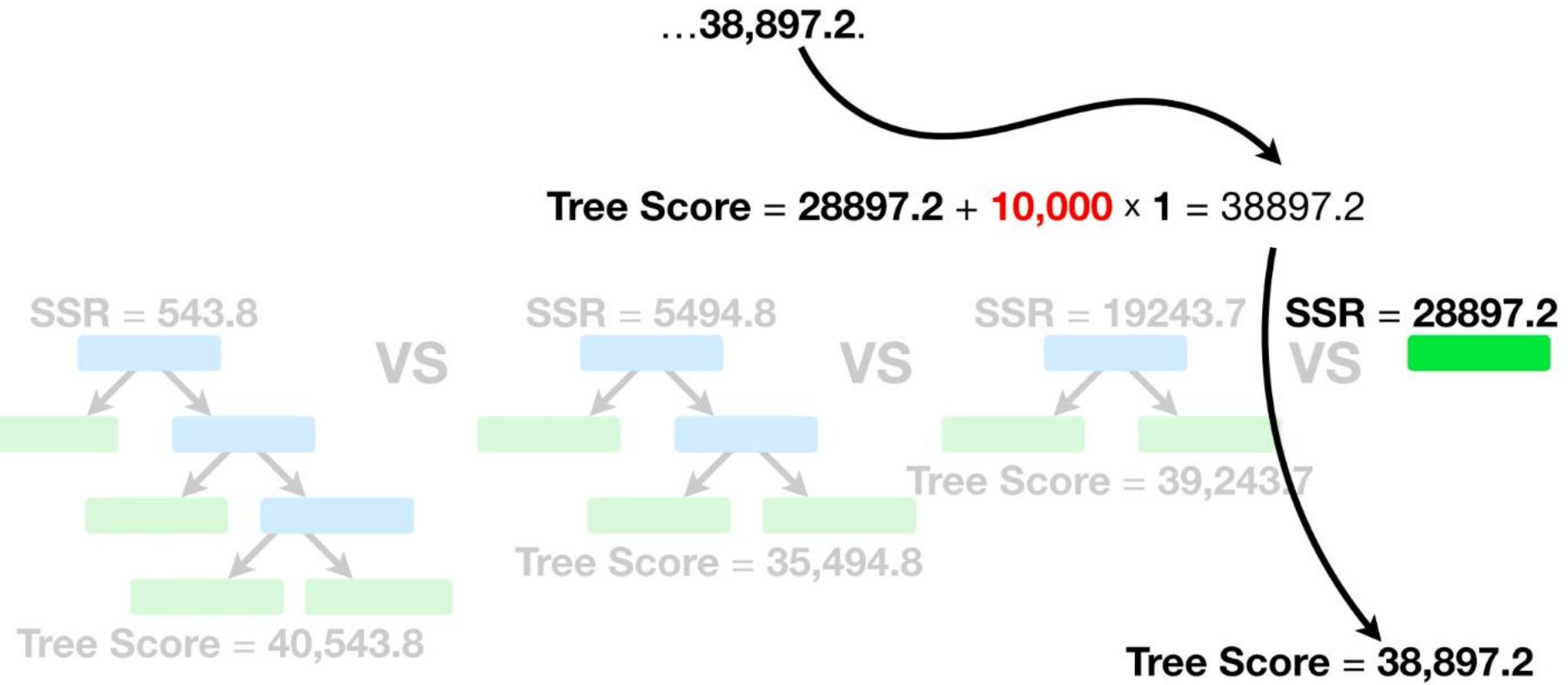
Tree Score = 35,494.8

SSR = 19243.7



SSR = 28897.2

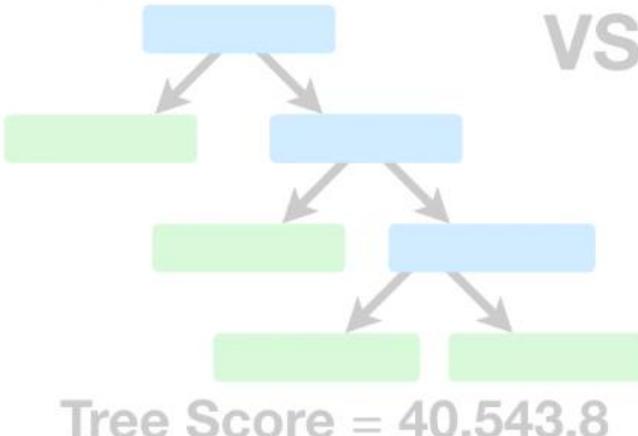
VS



NOTE: Because $a = 10,000$, the **Tree Complexity Penalty** for the tree with 1 leaf was **10,000....**

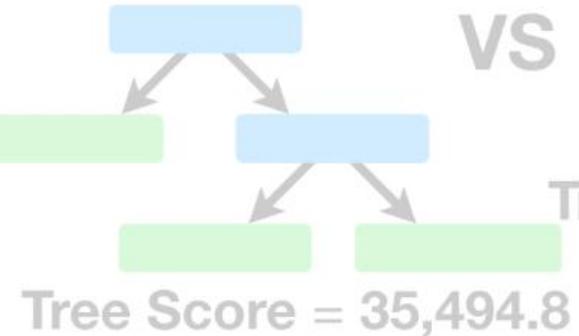
$$\text{Tree Score} = \text{SSR} + 10,000 \times T$$

$$\text{SSR} = 543.8$$



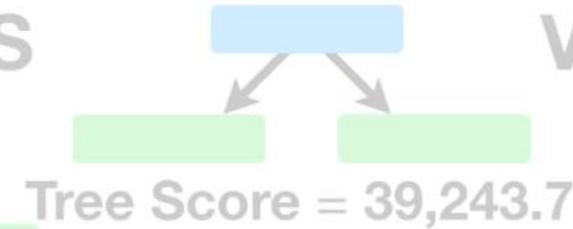
VS

$$\text{SSR} = 5494.8$$

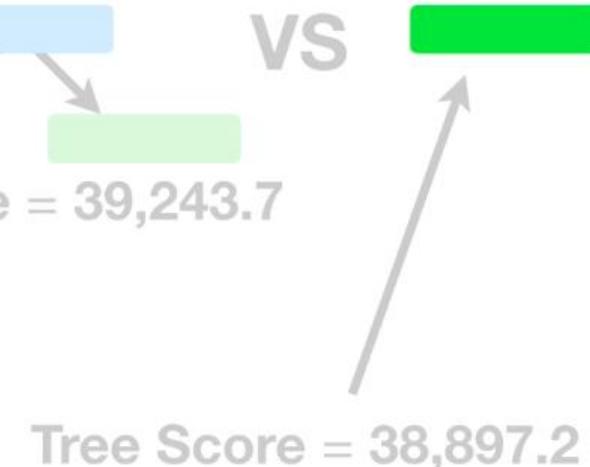


VS

$$\text{SSR} = 19243.7$$



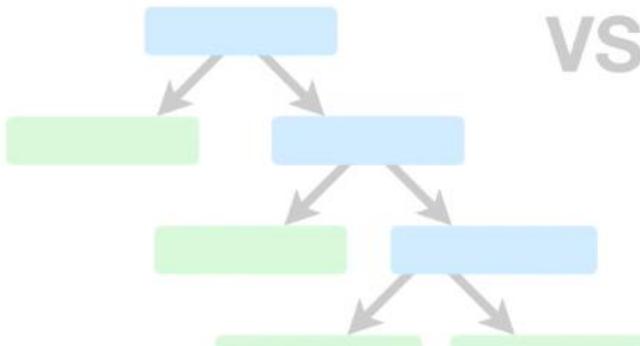
$$\text{SSR} = 28897.2$$



...and the **Tree Complexity Penalty** for the tree
with **2 leaves** was **20,000...**

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

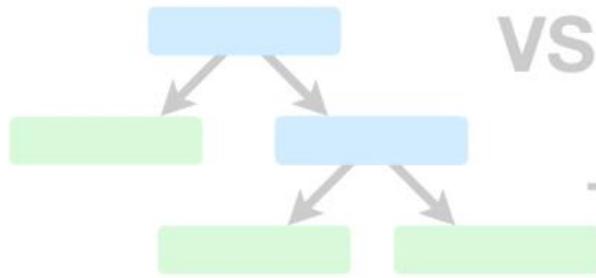
$$\text{SSR} = 543.8$$



$$\text{Tree Score} = 40,543.8$$

VS

$$\text{SSR} = 5494.8$$



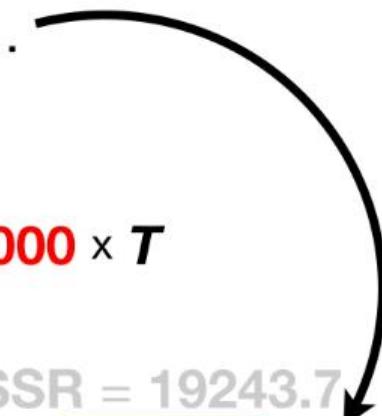
$$\text{Tree Score} = 35,494.8$$

$$\text{SSR} = 19243.7$$



$$\text{Tree Score} = 39,243.7$$

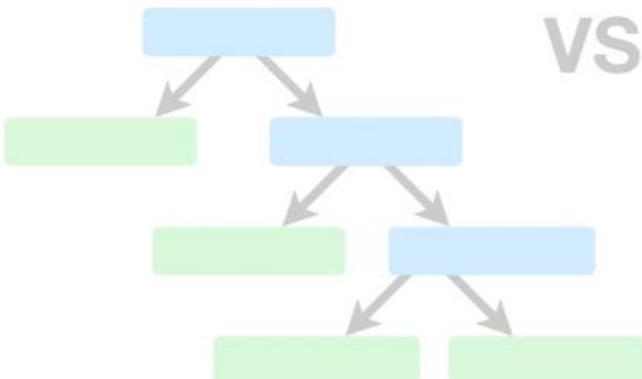
$$\text{Tree Score} = 38,897.2$$



...and the **Tree Complexity Penalty** for the tree
with 3 leaves was **30,000...**

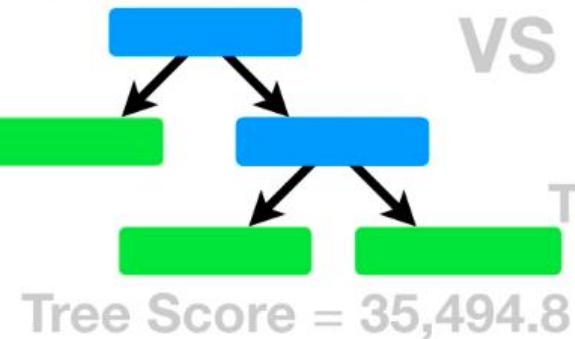
$$\text{Tree Score} = \text{SSR} + 10,000 \times T$$

$$\text{SSR} = 543.8$$



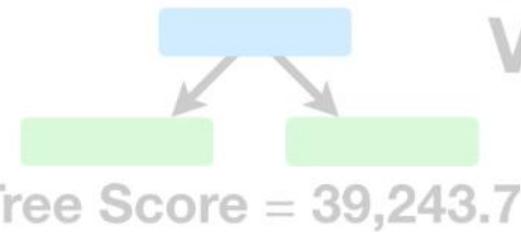
$$\text{Tree Score} = 40,543.8$$

$$\text{SSR} = 5494.8$$



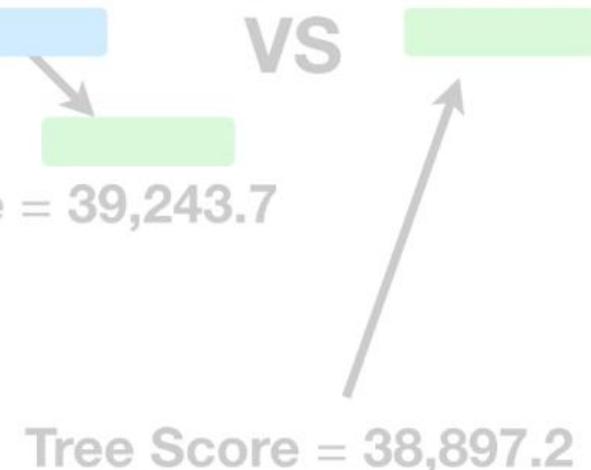
$$\text{Tree Score} = 35,494.8$$

$$\text{SSR} = 19243.7$$



$$\text{Tree Score} = 39,243.7$$

$$\text{SSR} = 28897.2$$



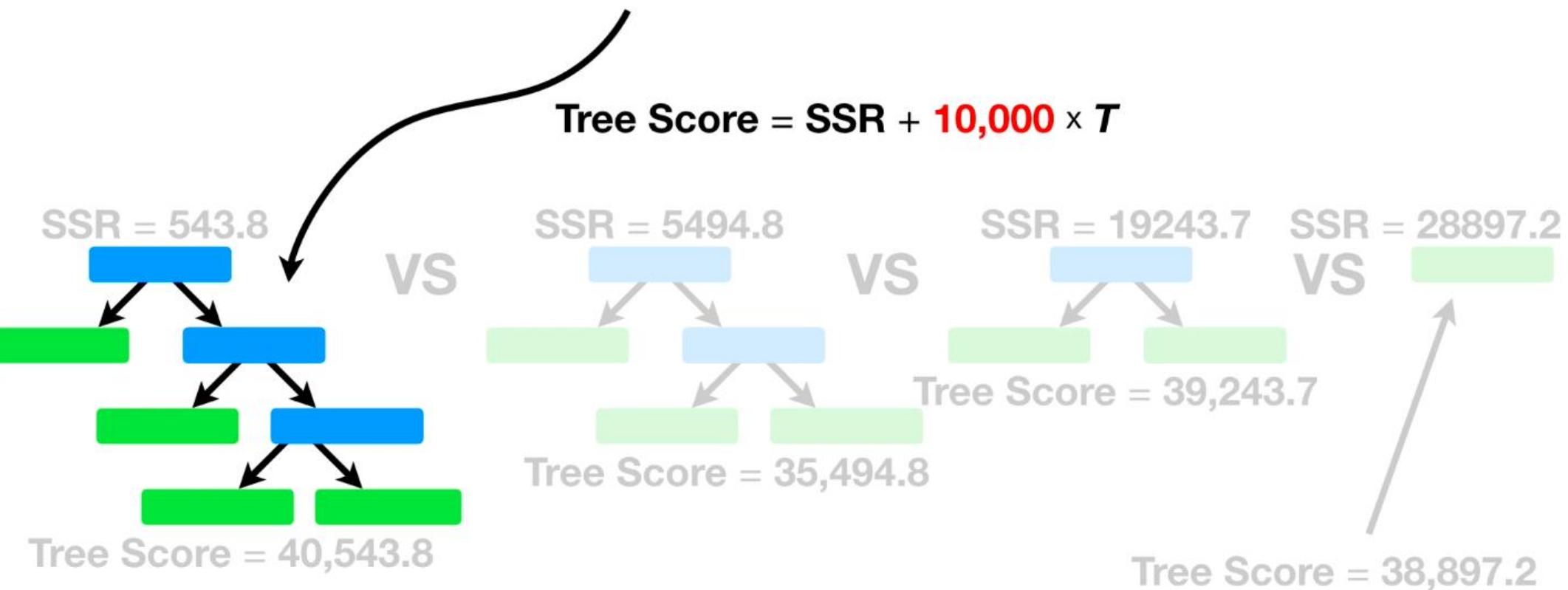
$$\text{Tree Score} = 38,897.2$$

VS

VS



...and the **Tree Complexity Penalty** for the original, full sized tree with 4 leaves was **40,000**.

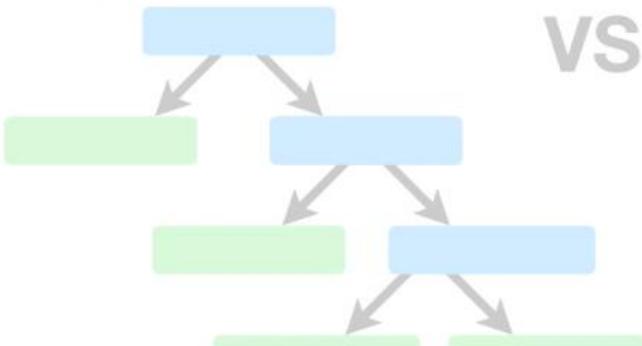


Thus, the more leaves,
the larger the penalty.



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

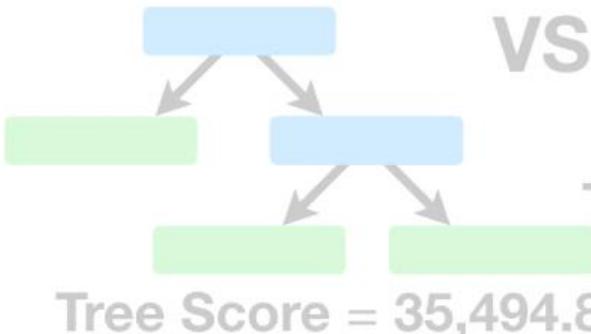
$$\text{SSR} = 543.8$$



$$\text{Tree Score} = 40,543.8$$

VS

$$\text{SSR} = 5494.8$$



$$\text{Tree Score} = 35,494.8$$

$$\text{SSR} = 19243.7$$

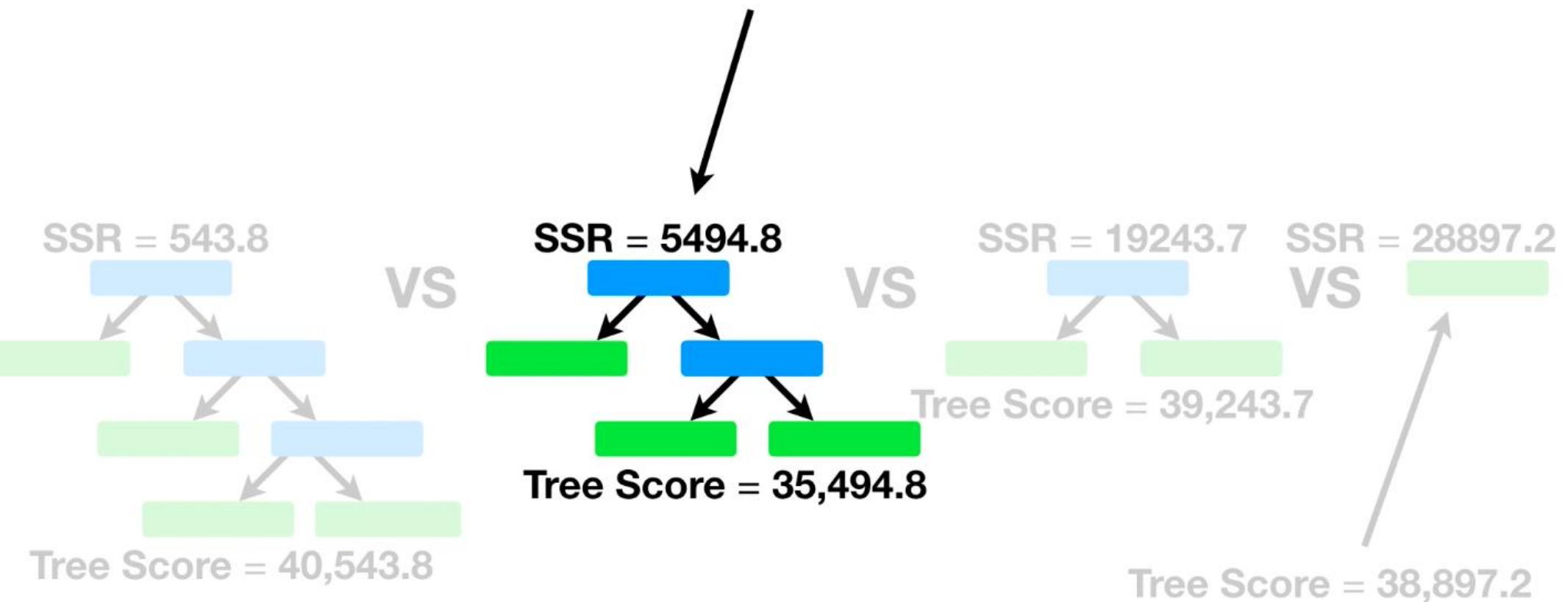
VS

$$\text{Tree Score} = 39,243.7$$

$$\text{SSR} = 28897.2$$

$$\text{Tree Score} = 38,897.2$$

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

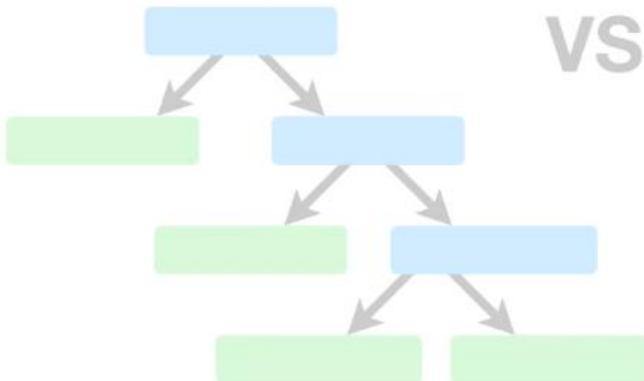


NOTE: If we set $\alpha = 22,000 \dots$



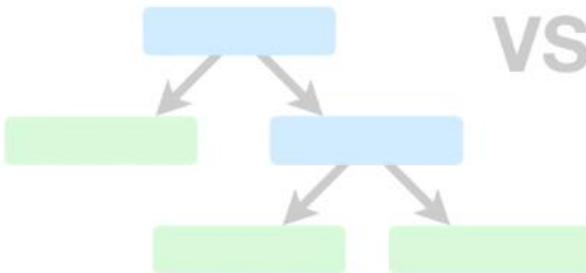
$$\text{Tree Score} = \text{SSR} + 22,000 \times T$$

$$\text{SSR} = 543.8$$



VS

$$\text{SSR} = 5494.8$$



VS

$$\text{SSR} = 19243.7$$



$$\text{SSR} = 28897.2$$

VS

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

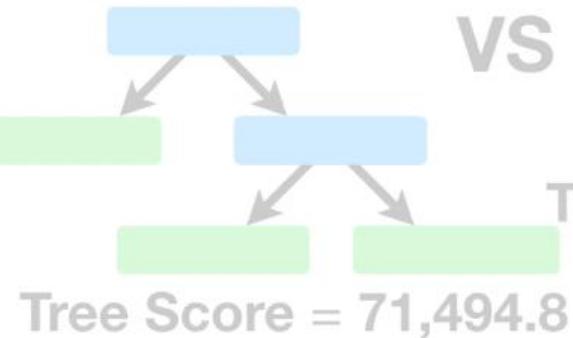
$$\text{Tree Score} = \text{SSR} + 22,000 \times T$$

$$\text{SSR} = 543.8$$



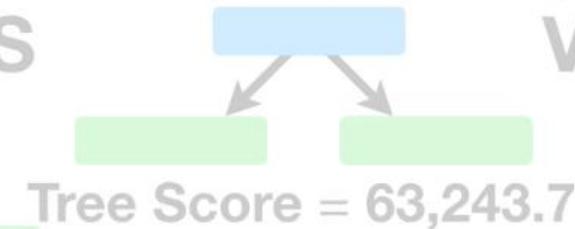
VS

$$\text{SSR} = 5494.8$$

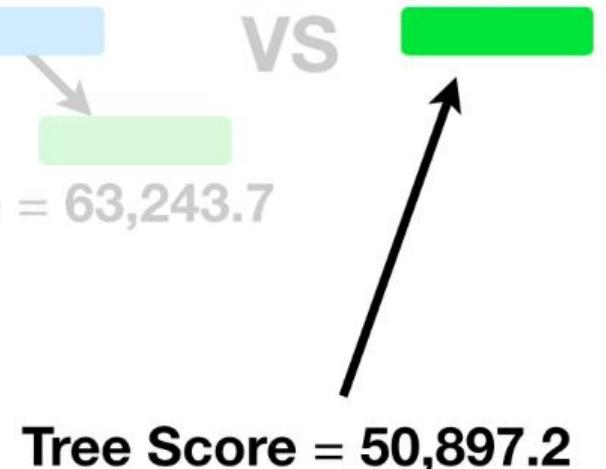


VS

$$\text{SSR} = 19243.7$$



$$\text{SSR} = 28897.2$$



$$\text{Tree Score} = 88,543.8$$

$$\text{Tree Score} = 71,494.8$$

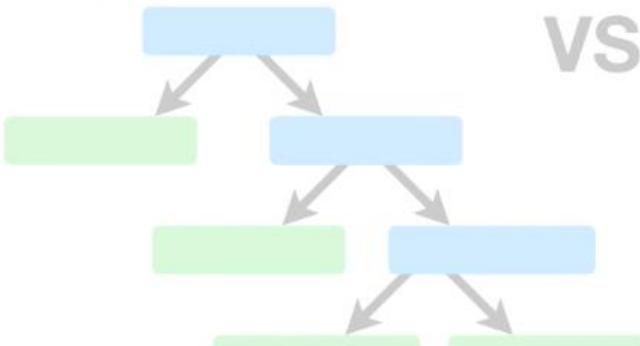
$$\text{Tree Score} = 50,897.2$$

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



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

$$\text{SSR} = 543.8$$



$$\text{Tree Score} = 88,543.8$$

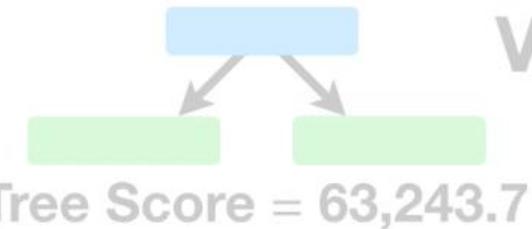
VS

$$\text{SSR} = 5494.8$$



$$\text{Tree Score} = 71,494.8$$

$$\text{SSR} = 19243.7$$



$$\text{Tree Score} = 63,243.7$$

$$\text{SSR} = 28897.2$$

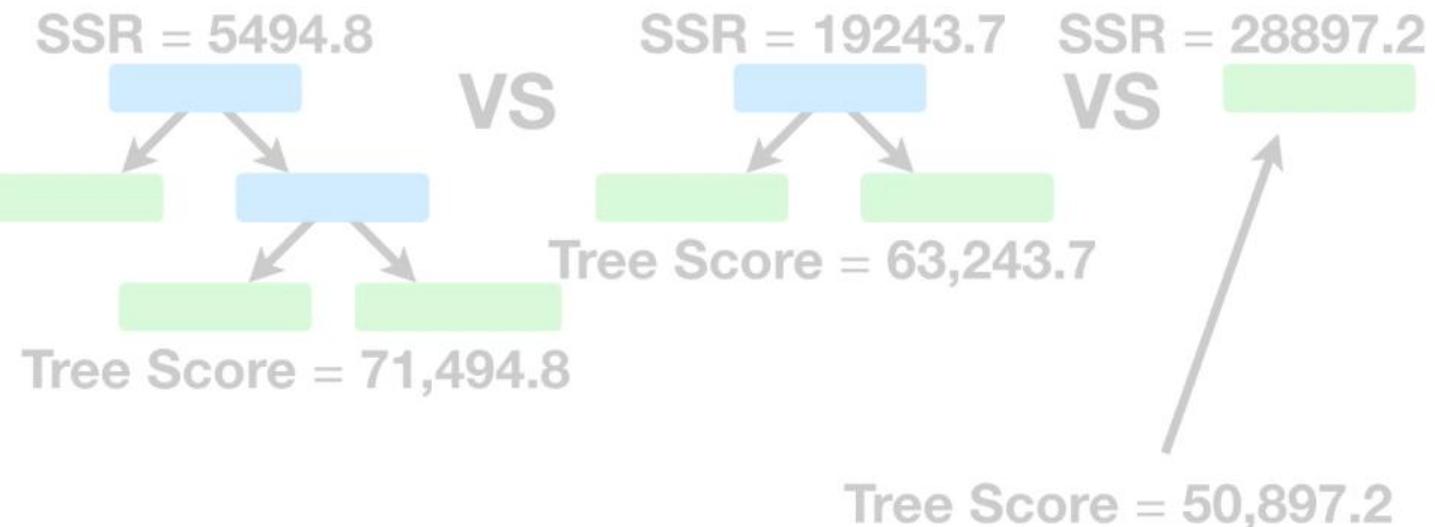
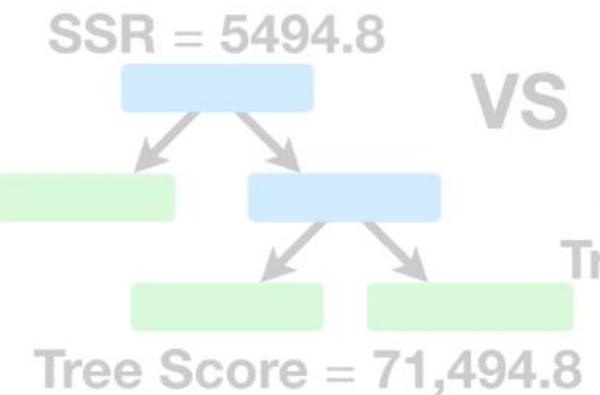
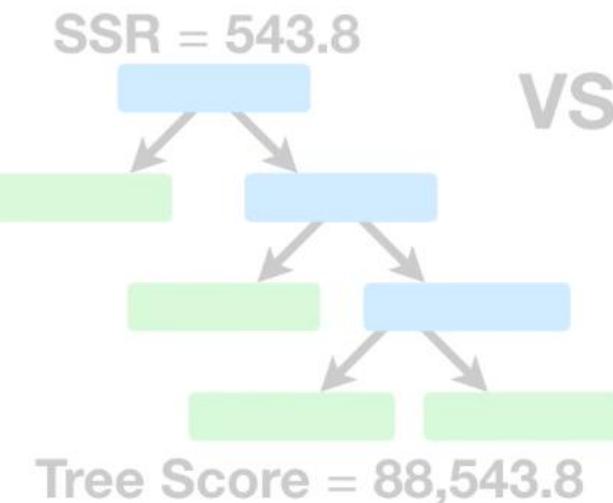
VS

$$\text{Tree Score} = 50,897.2$$



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

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

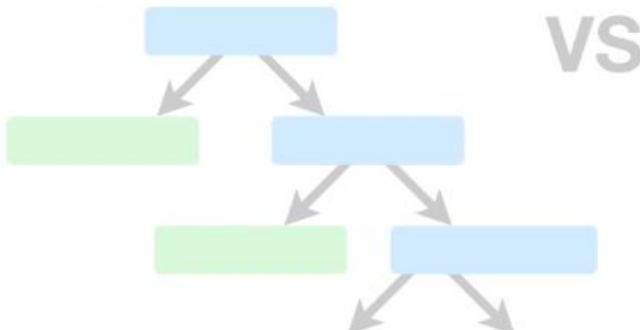


...and how to find the best value for a .



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

$$\mathbf{SSR} = 543.8$$



$$\text{Tree Score} = 88,543.8$$

VS

$$\mathbf{SSR} = 5494.8$$



$$\text{Tree Score} = 71,494.8$$

VS

$$\mathbf{SSR} = 19243.7$$



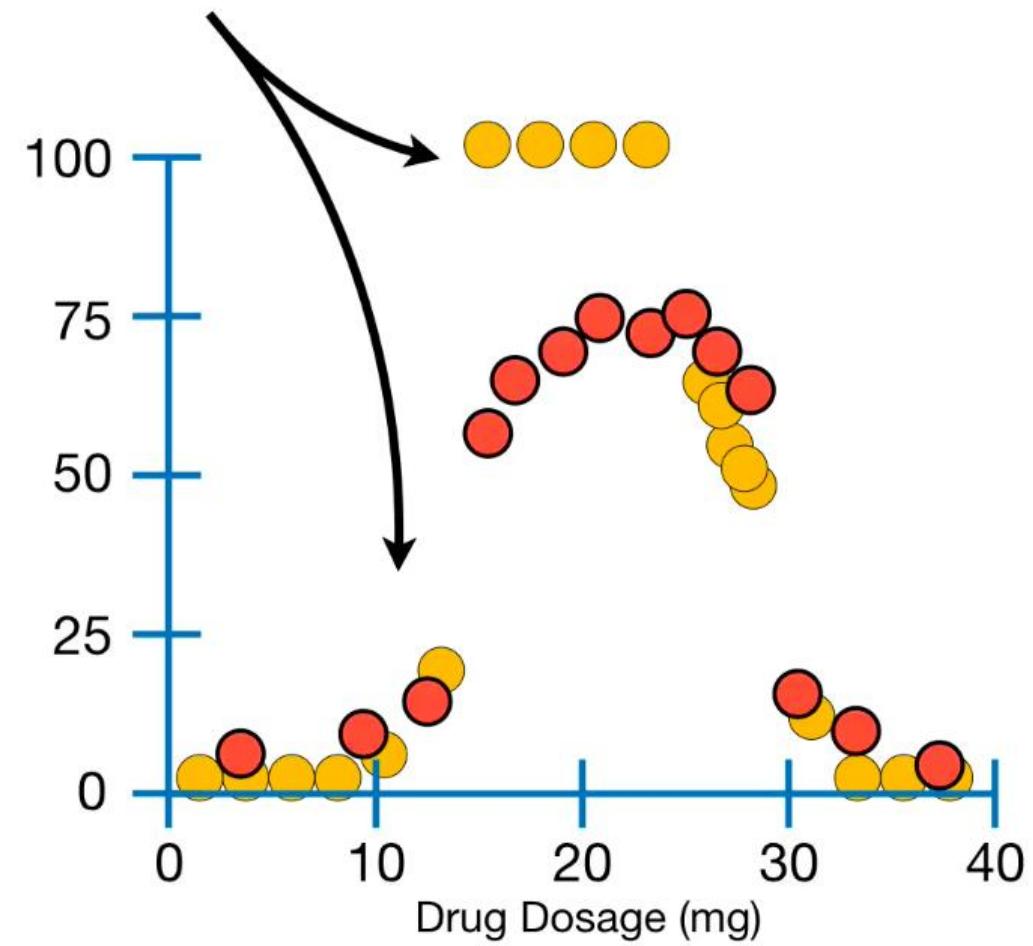
$$\text{Tree Score} = 63,243.7$$

$$\mathbf{SSR} = 28897.2$$

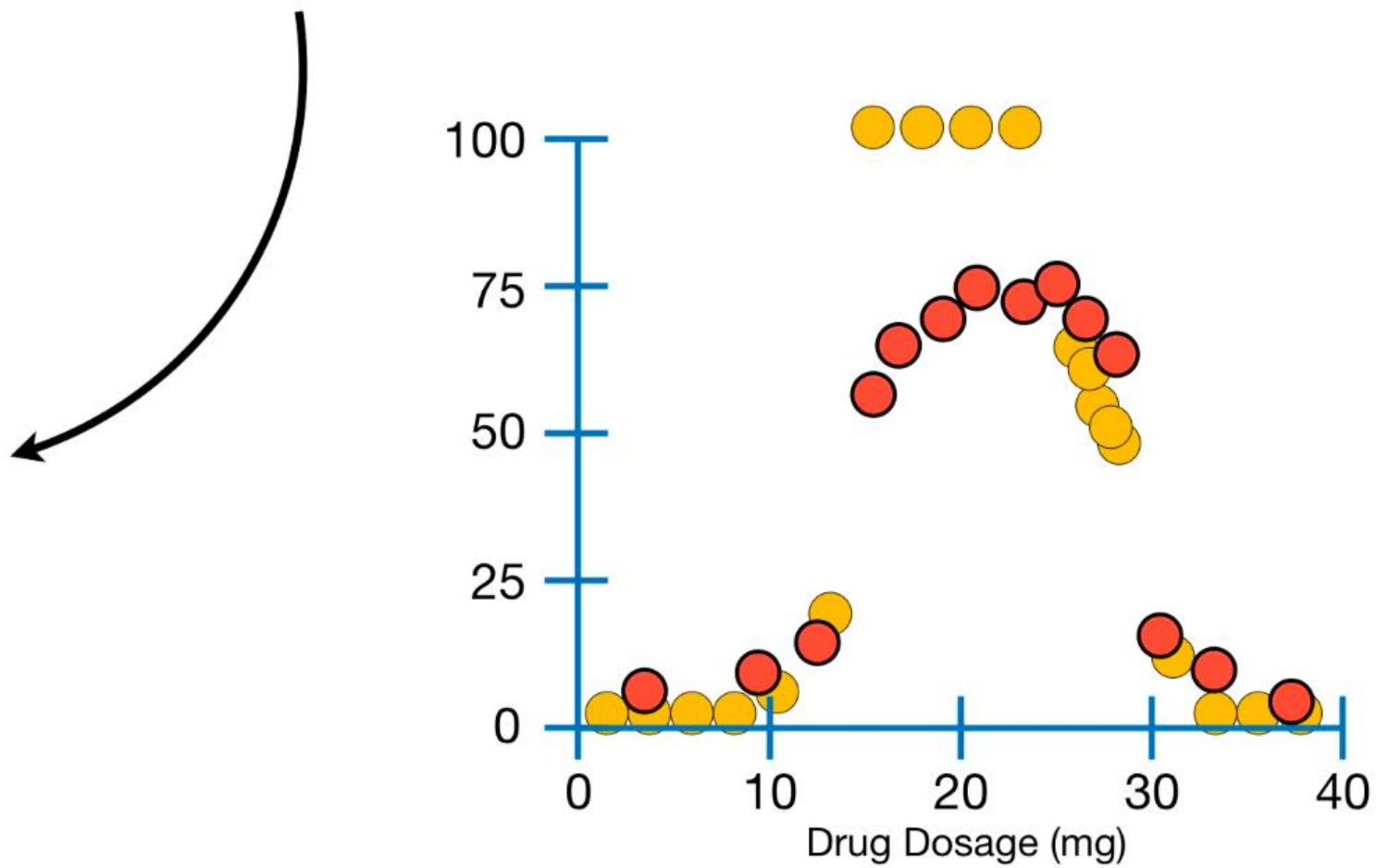
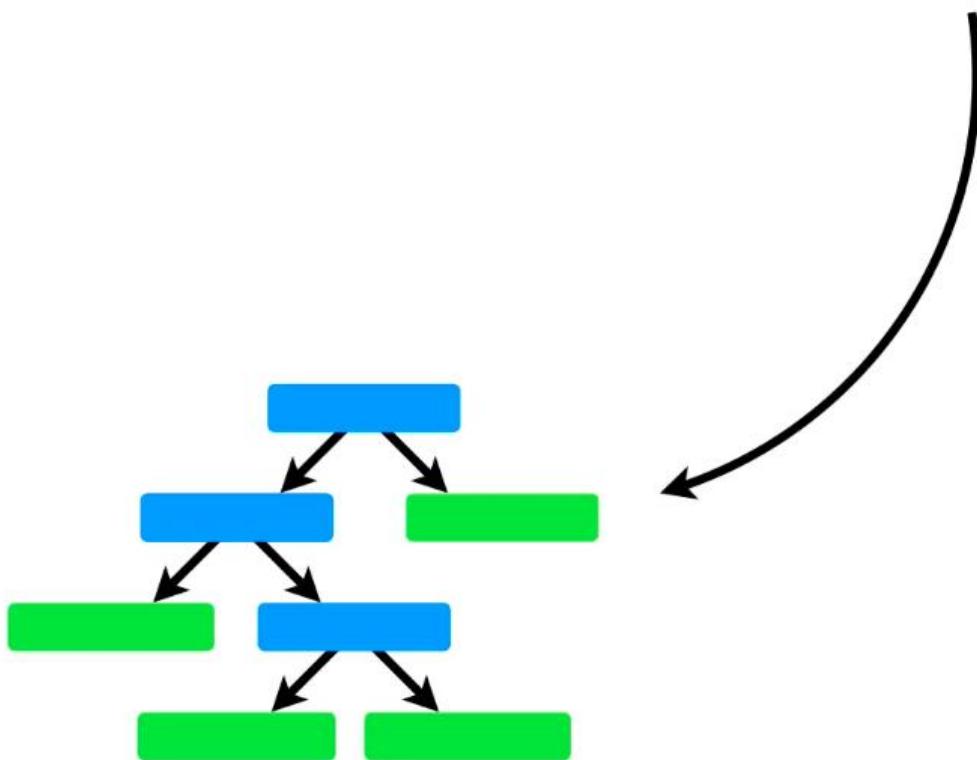
VS

$$\text{Tree Score} = 50,897.2$$

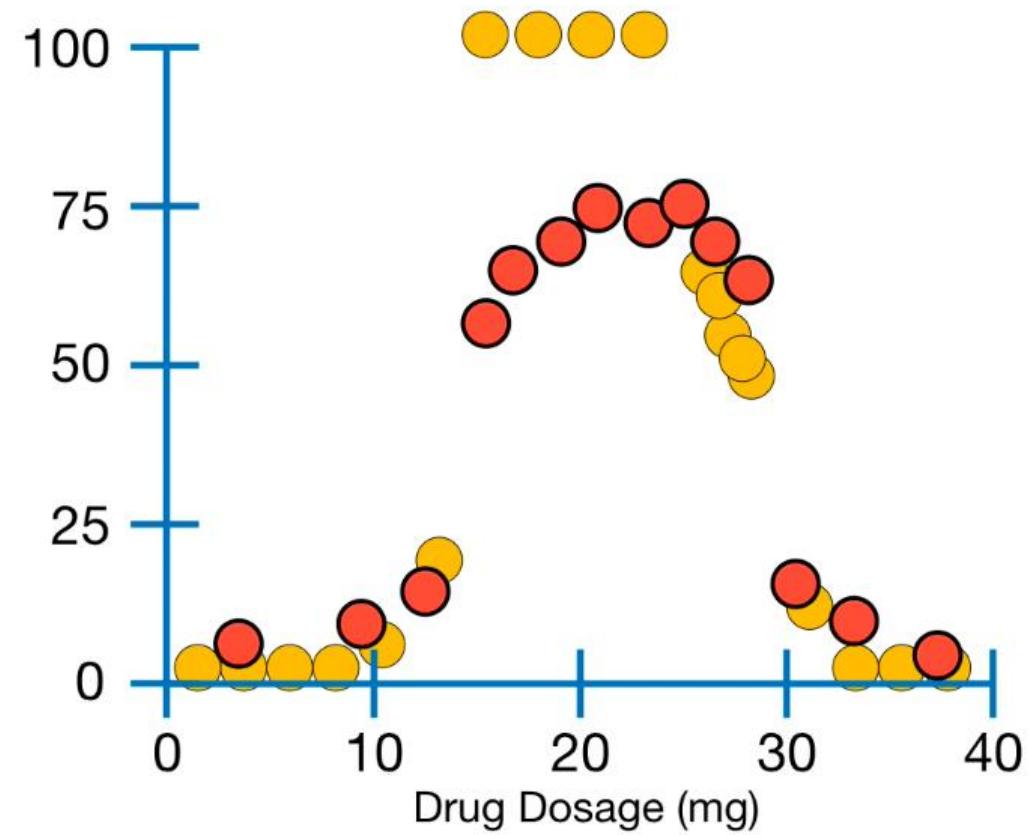
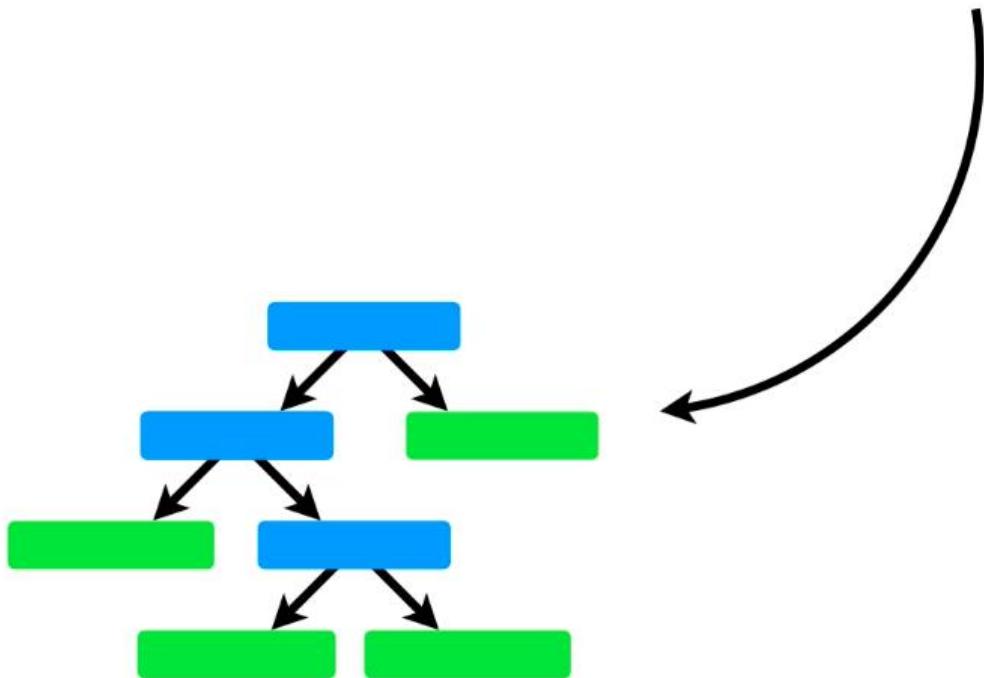
First, using *all* of the data...



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

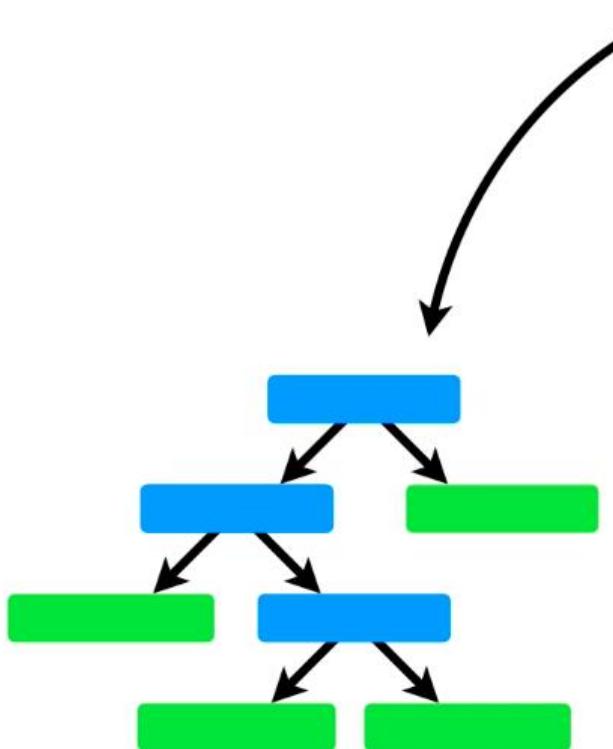


NOTE: This full sized tree is different than before because it was fit to *all* of the data, not just the **Training Data**.



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

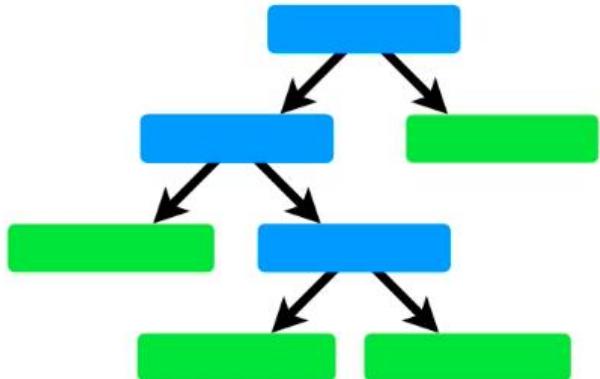
$$\text{Tree Score} = \text{SSR} + aT$$



...and the **Tree Score** is just the **Sum of the Squared Residuals**...

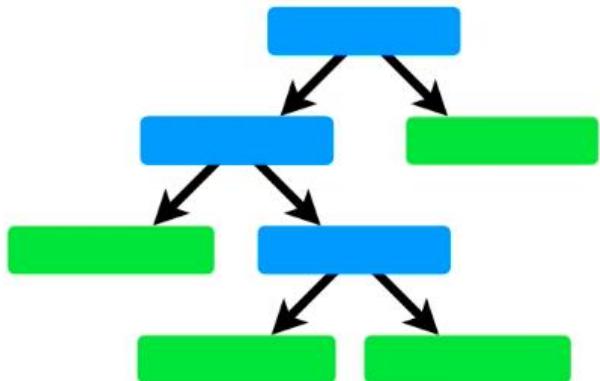


$$\text{Tree Score} = \text{SSR}$$



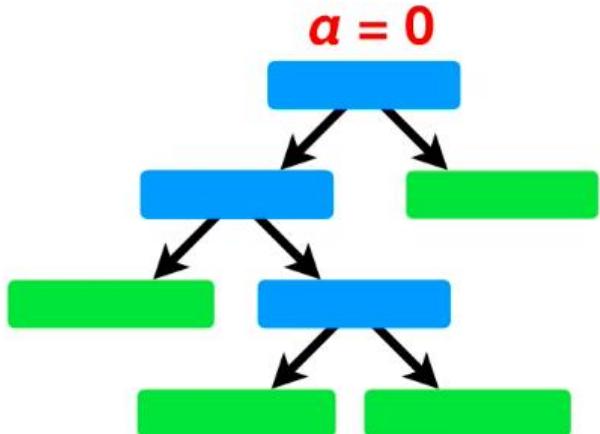
...and, as we saw earlier, all of the sub-trees will have larger **Sum of Squared Residuals**.

Tree Score = SSR



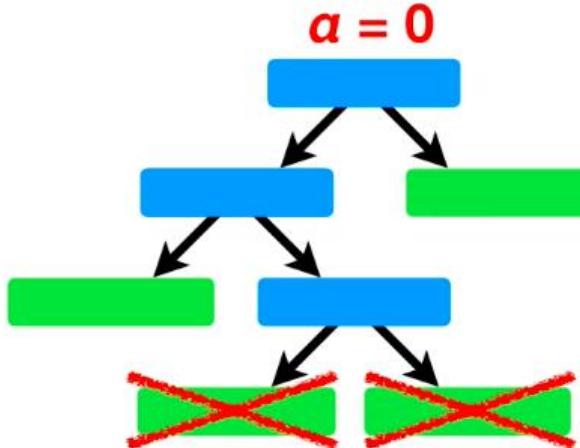
So let's put $a = 0$ here, to remind us
that this tree has the lowest **Tree Score**
when $a = 0$.

Tree Score = SSR



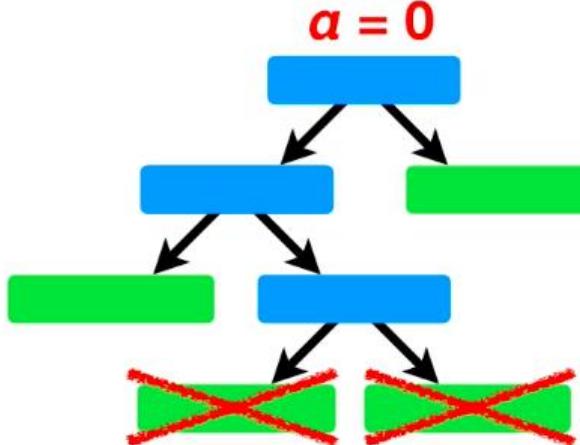
Now we increase a until pruning leaves will give us a lower **Tree Score**.

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



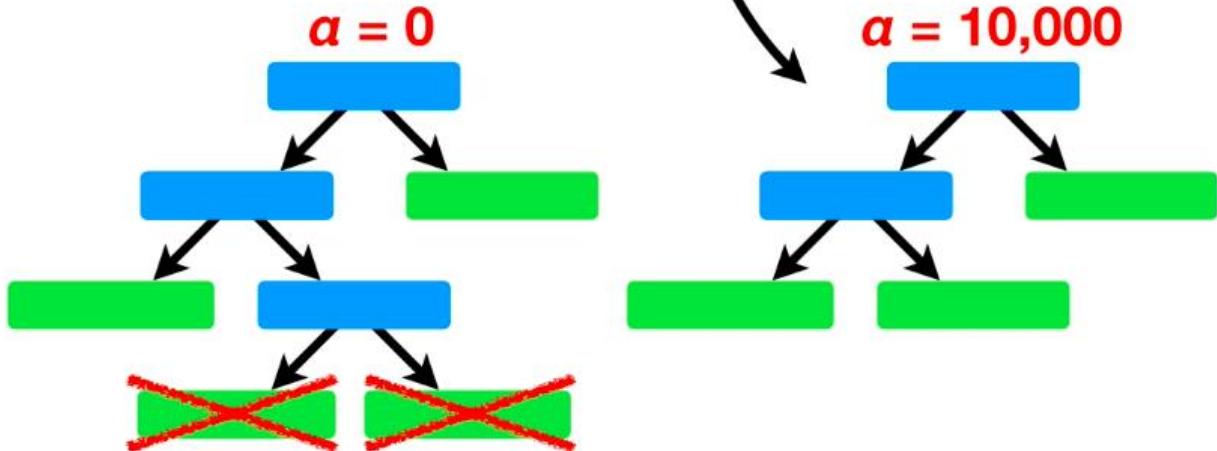
In this case, when $a = 10,000$, we'll get a lower **Tree Score** if we remove these leaves...

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



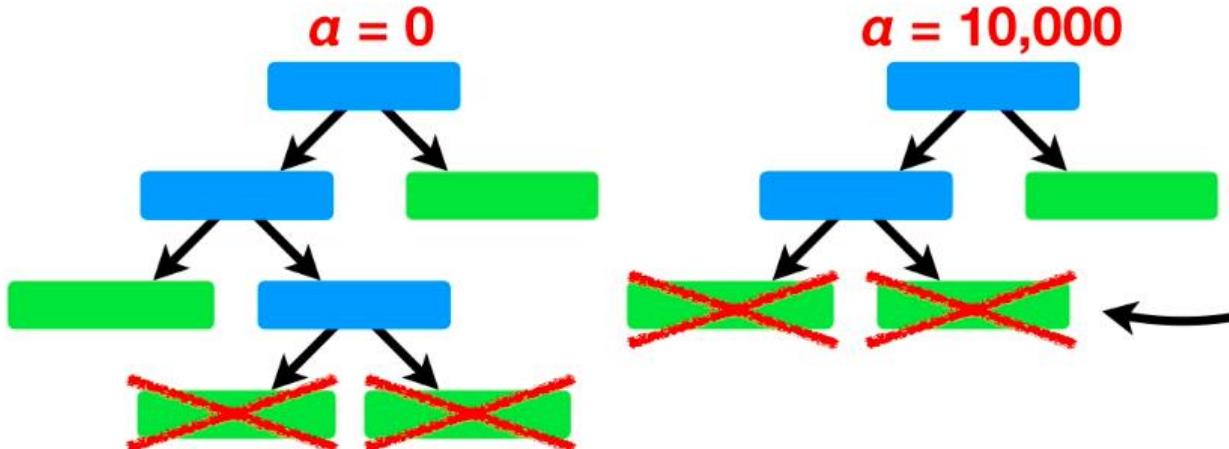
...and use this sub-tree.

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



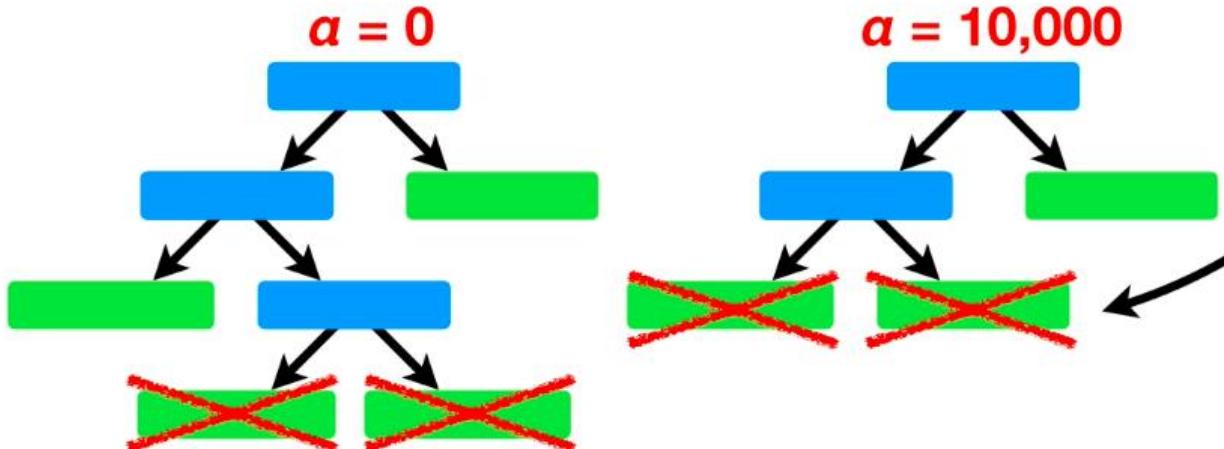
Now we increase a again until pruning leaves will give us a lower **Tree Score**.

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



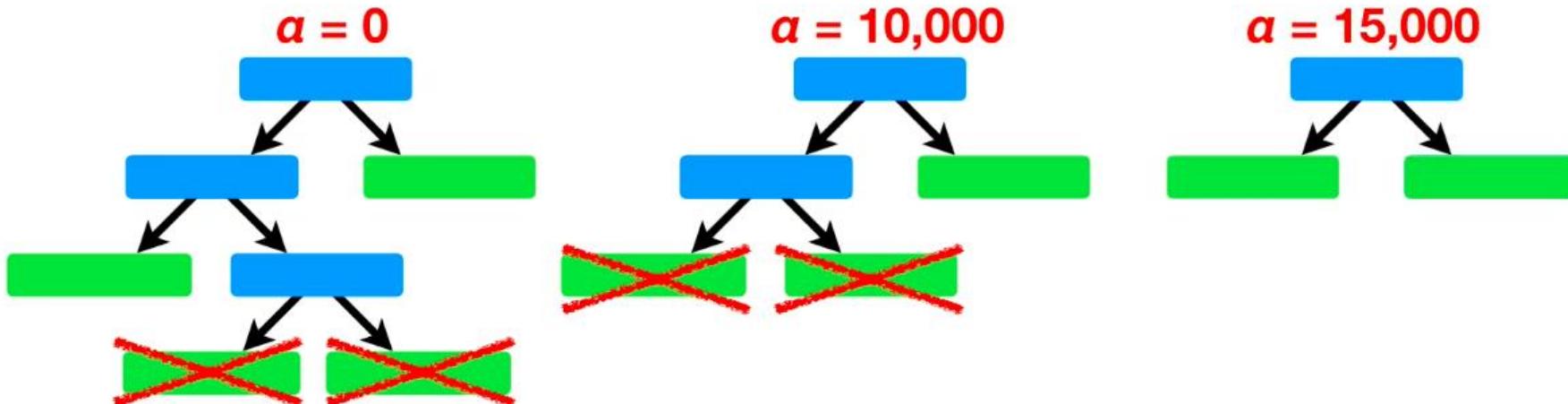
In this case when $a = 15,000$,
we will get a lower **Tree Score**
if we remove these leaves...

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



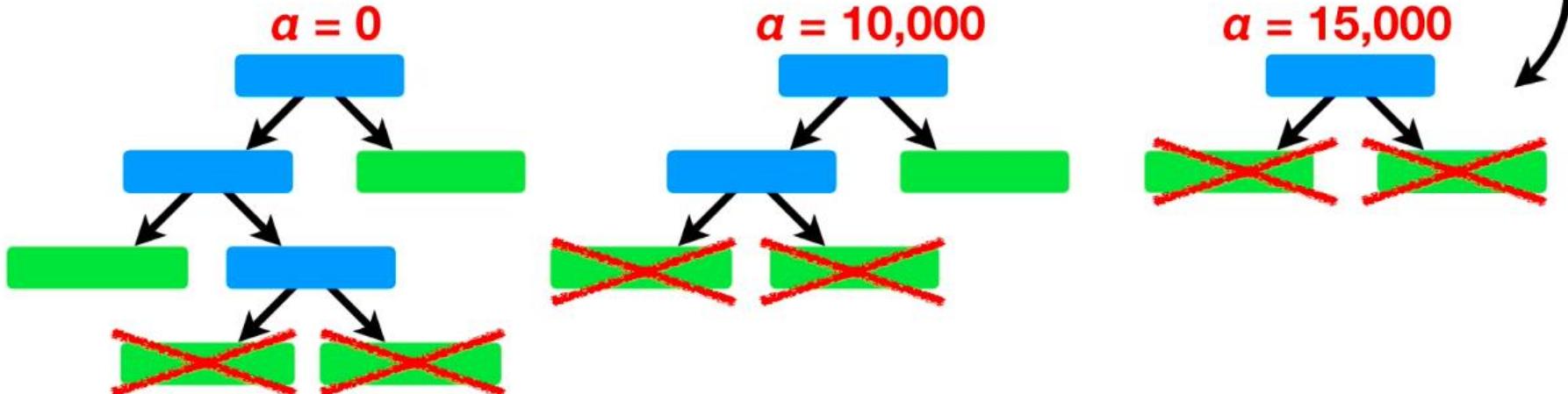
...and use this sub-tree
instead.

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

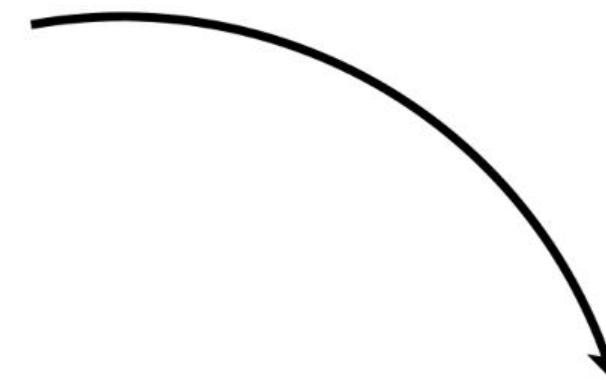


And when $a = 22,000$, we will
get a lower **Tree Score** if we
remove these leaves...

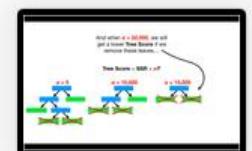
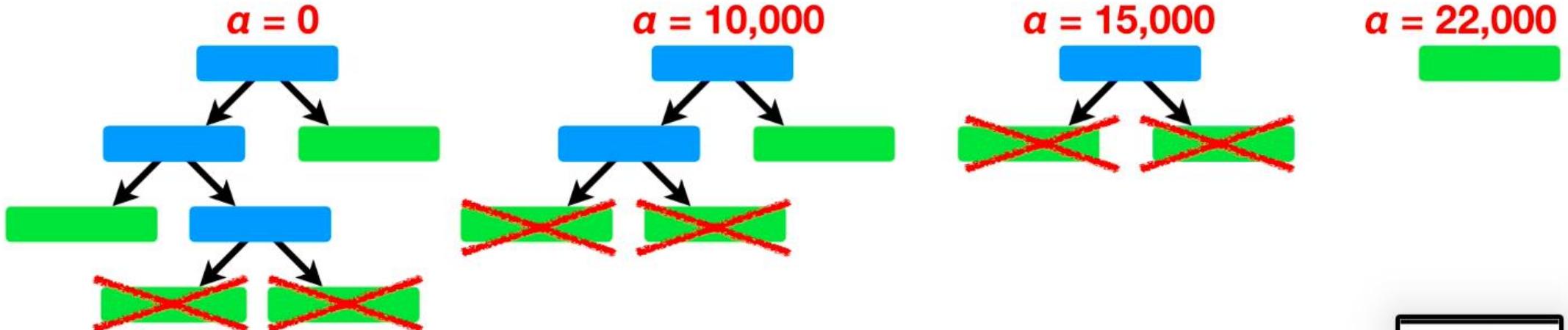
$$\text{Tree Score} = \text{SSR} + aT$$



...and use this sub-tree
instead.

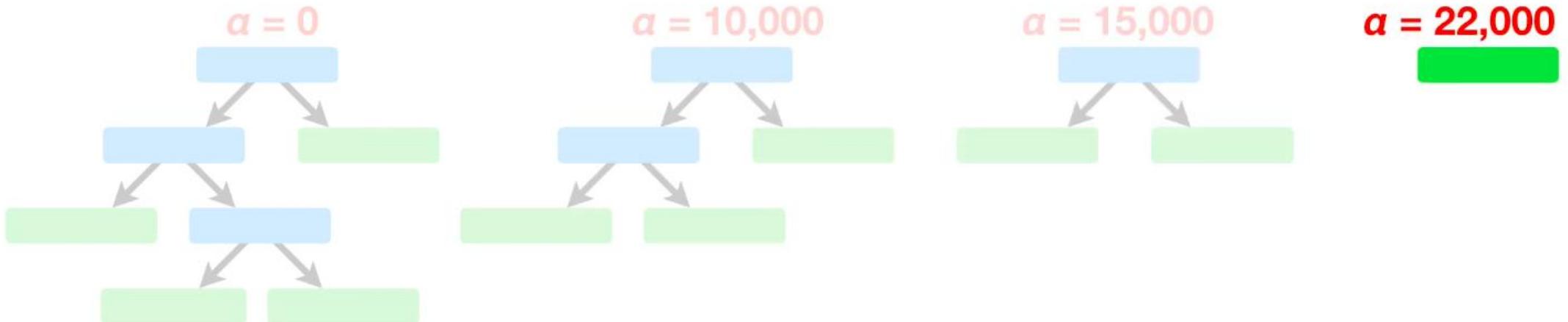


$$\text{Tree Score} = \text{SSR} + aT$$

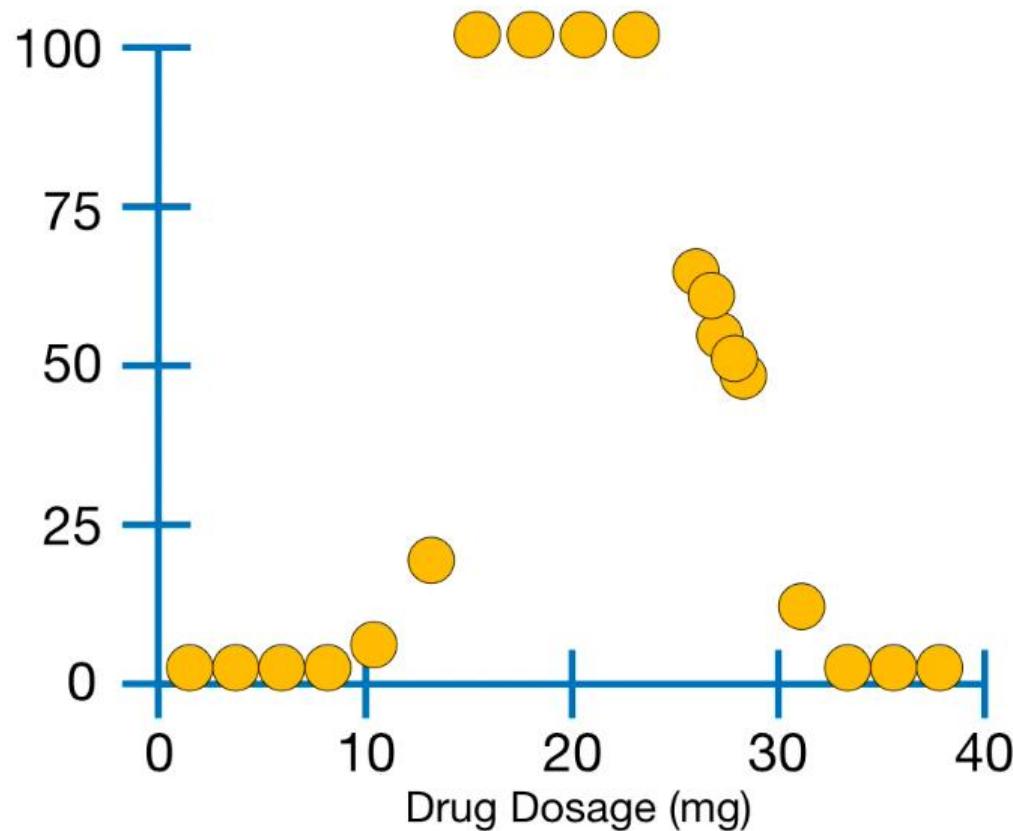


In the end, different values for a give us a sequence of trees, from full sized to just a leaf.

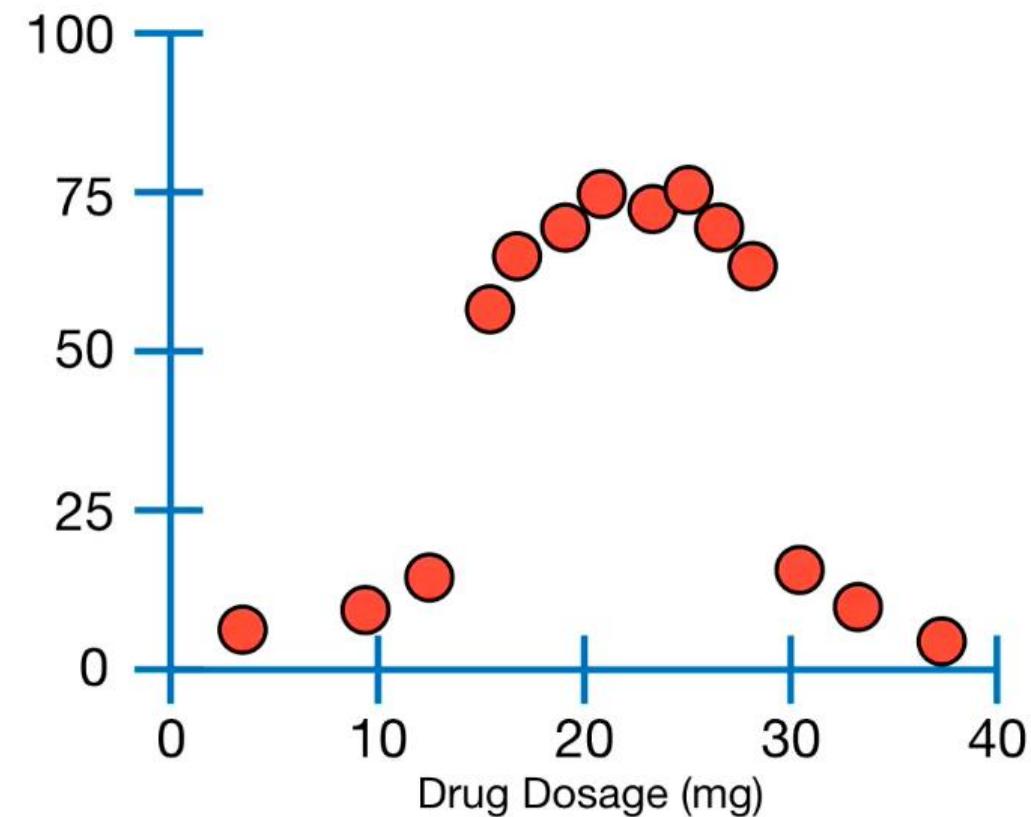
$$\text{Tree Score} = \text{SSR} + aT$$



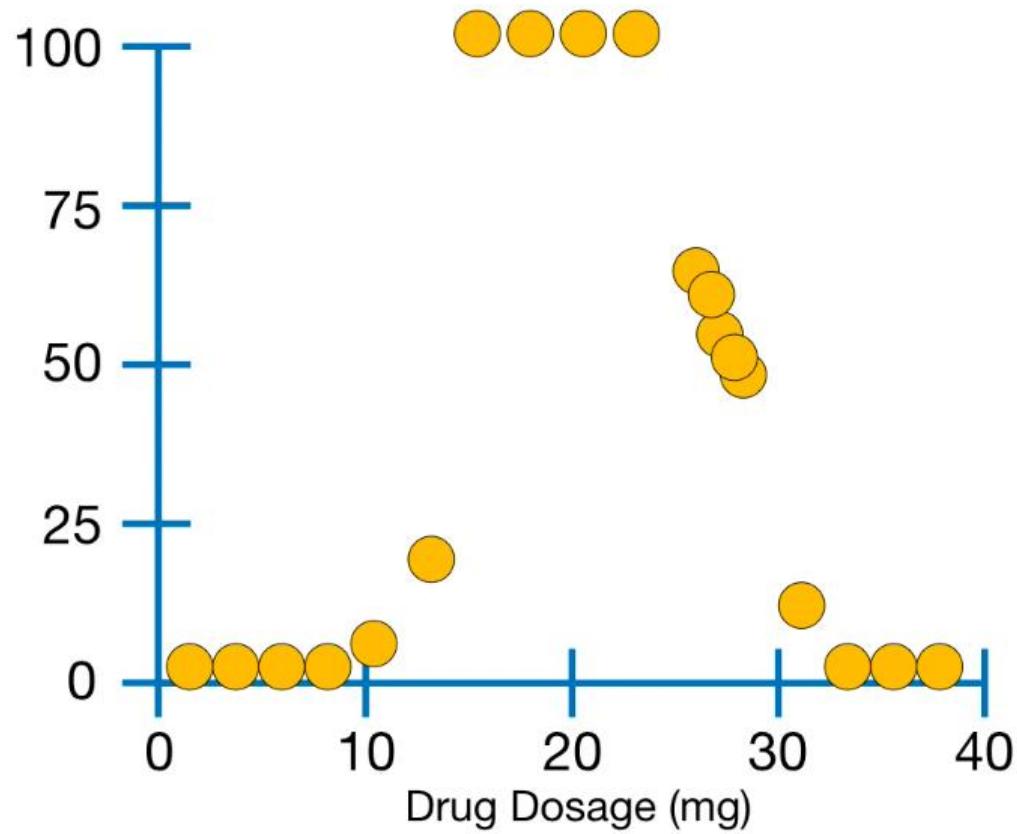
...and divide it into **Training**...



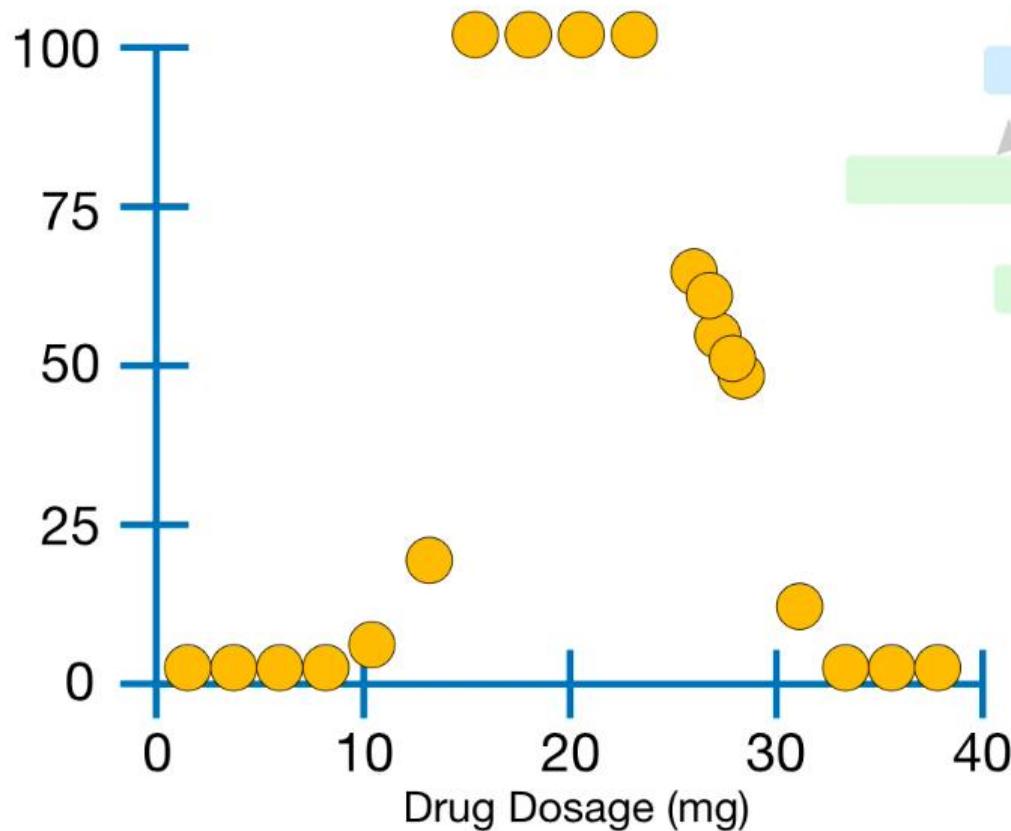
...and **Testing Datasets**.



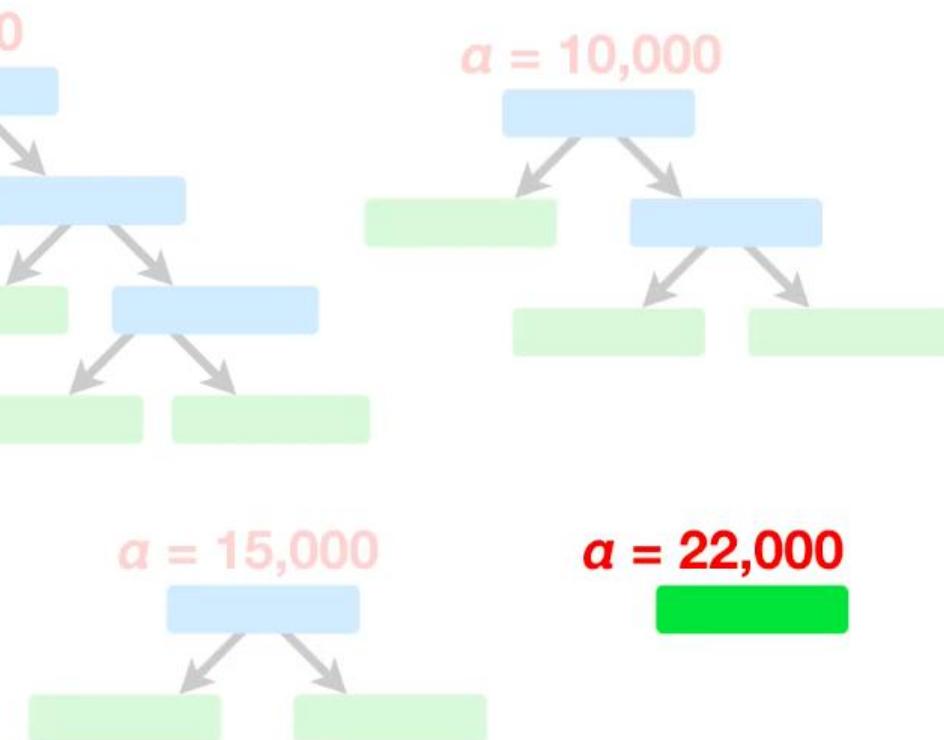
And just using the
Training Data...



And just using the
Training Data...

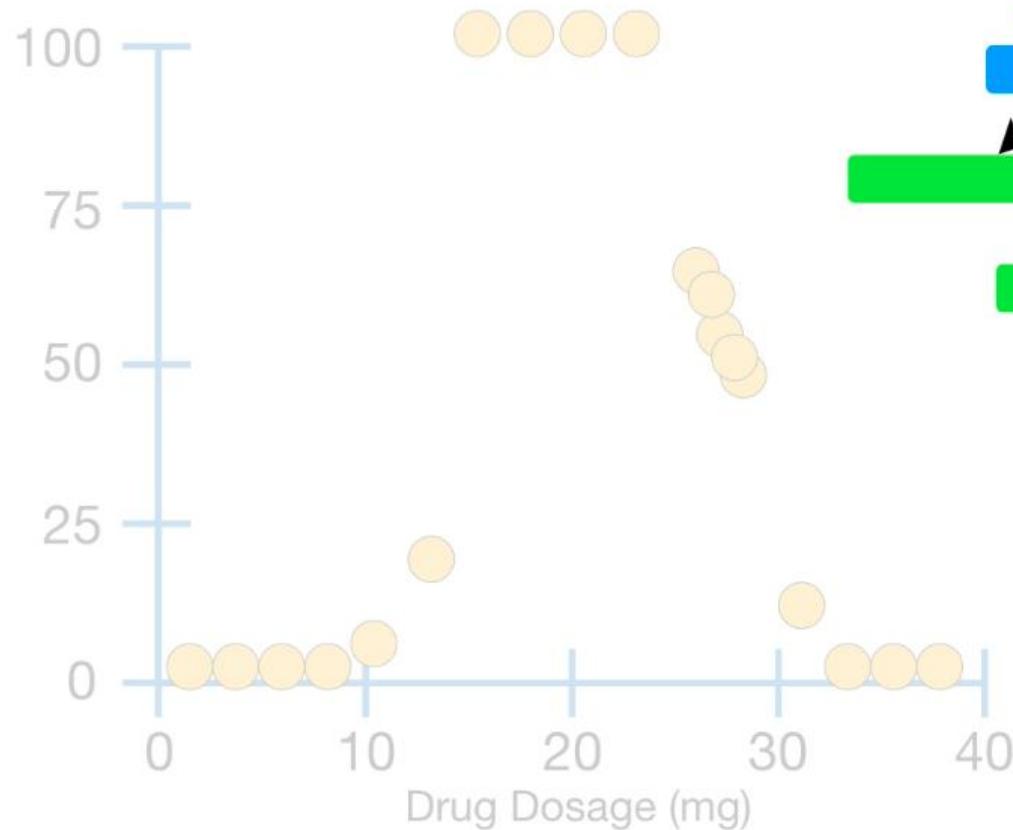


...use the α values we found before to build a full tree and a sequence of sub-trees that minimize the **Tree Score**.



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

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

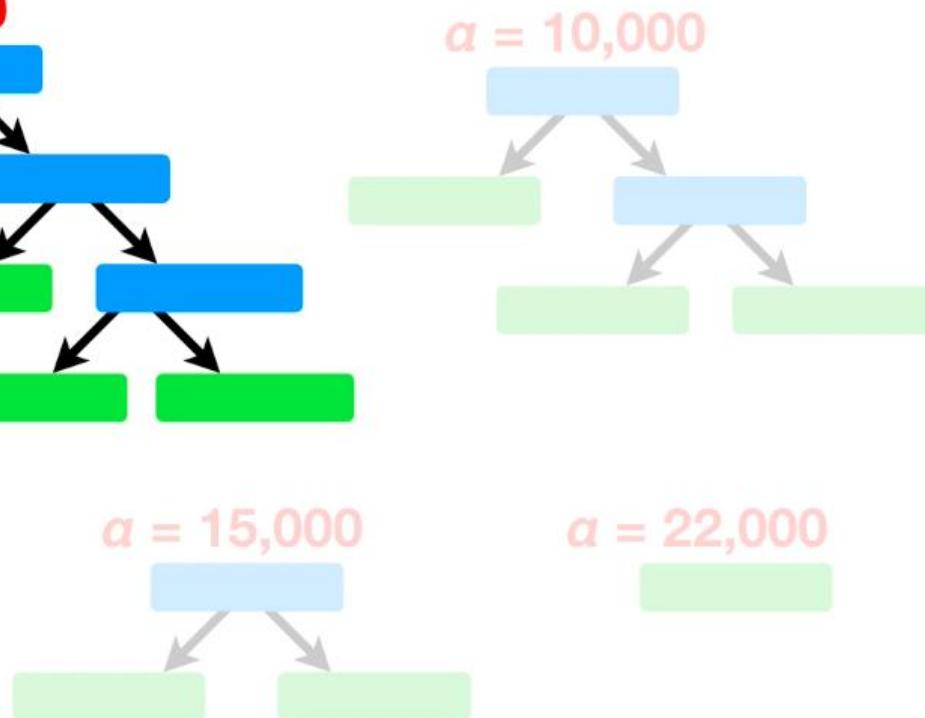


$\alpha = 0$

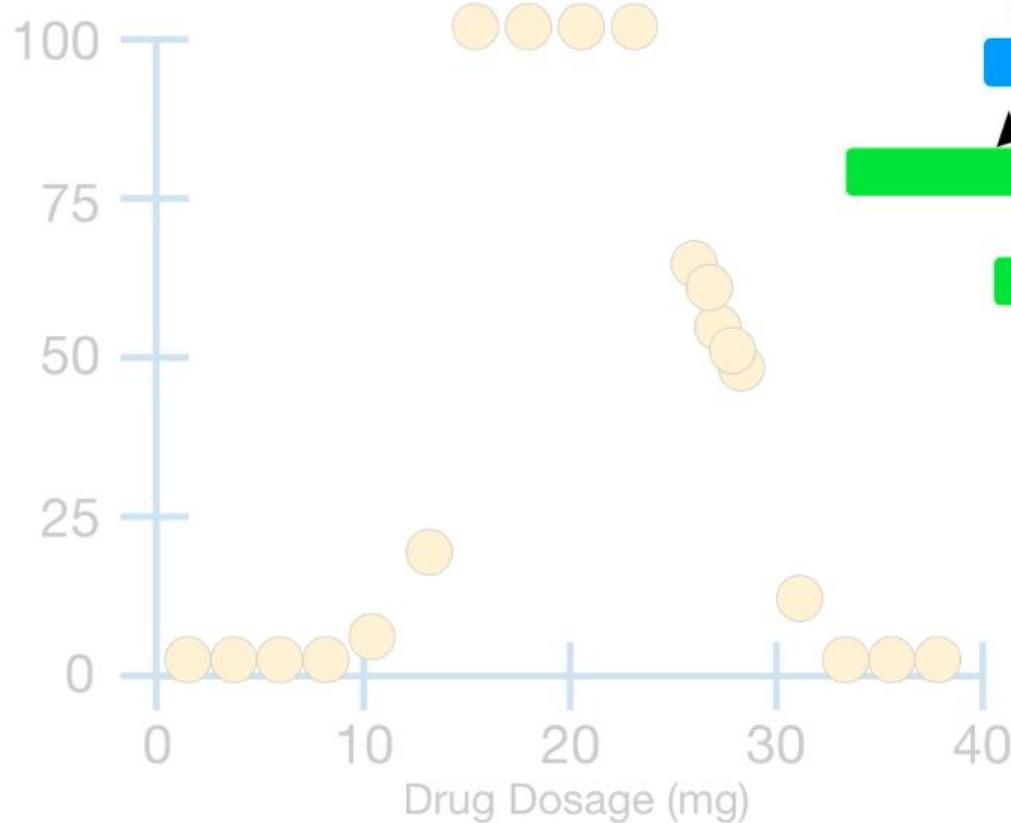
$\alpha = 10,000$

$\alpha = 15,000$

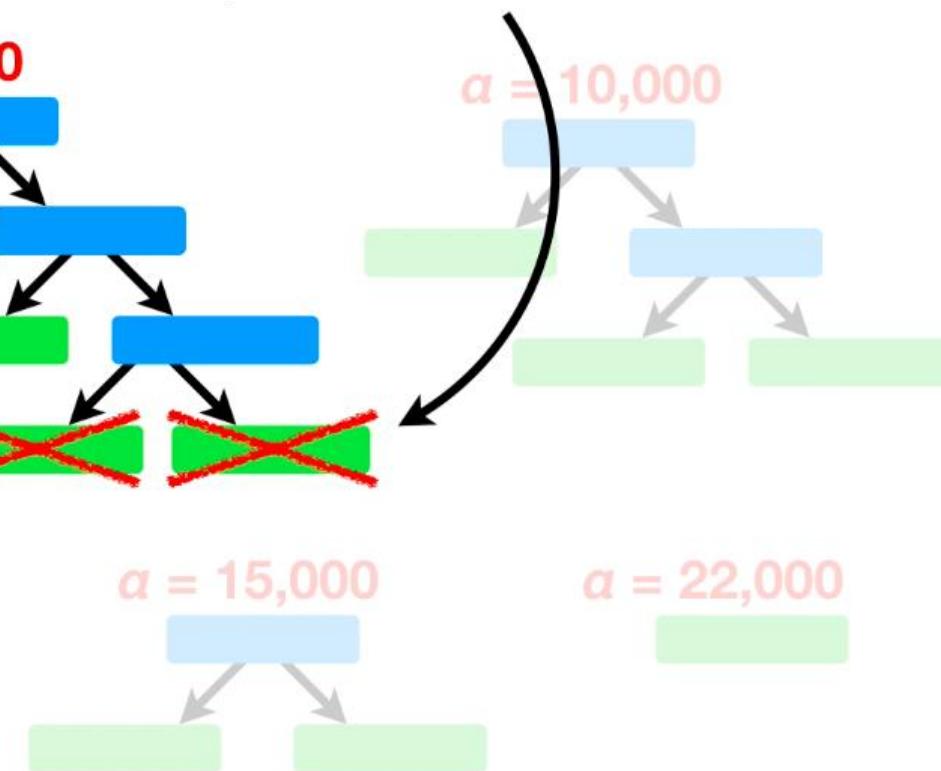
$\alpha = 22,000$



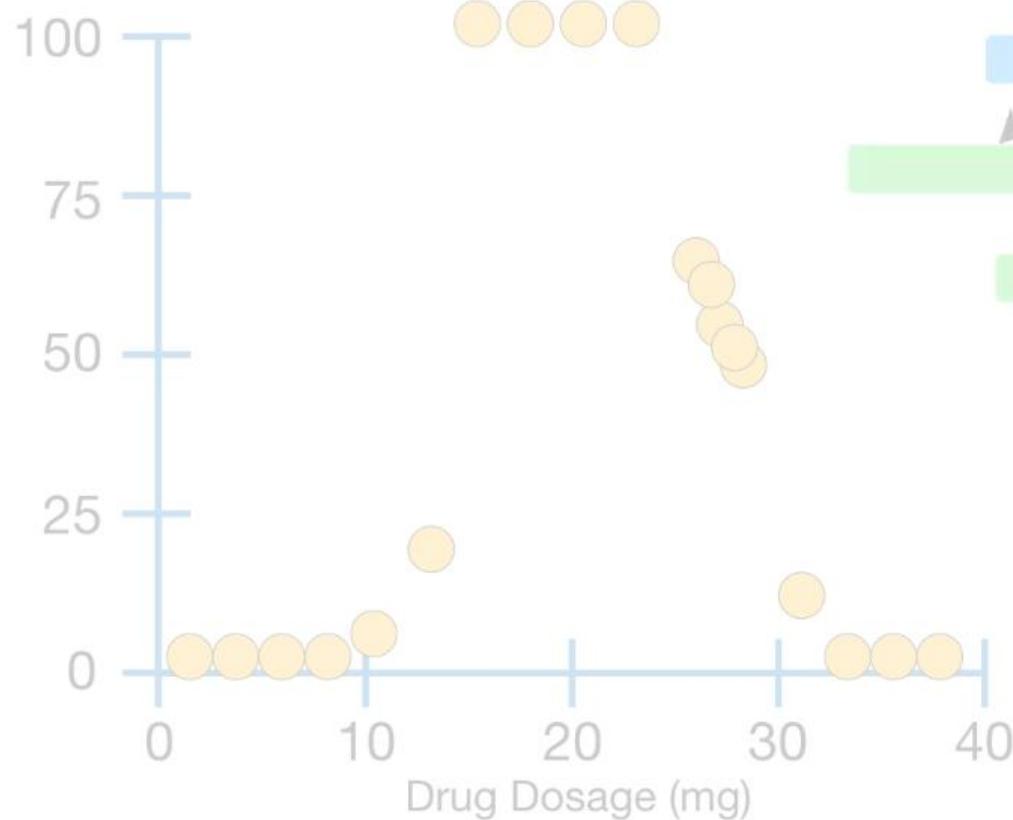
$$\text{Tree Score} = \text{SSR} + aT$$



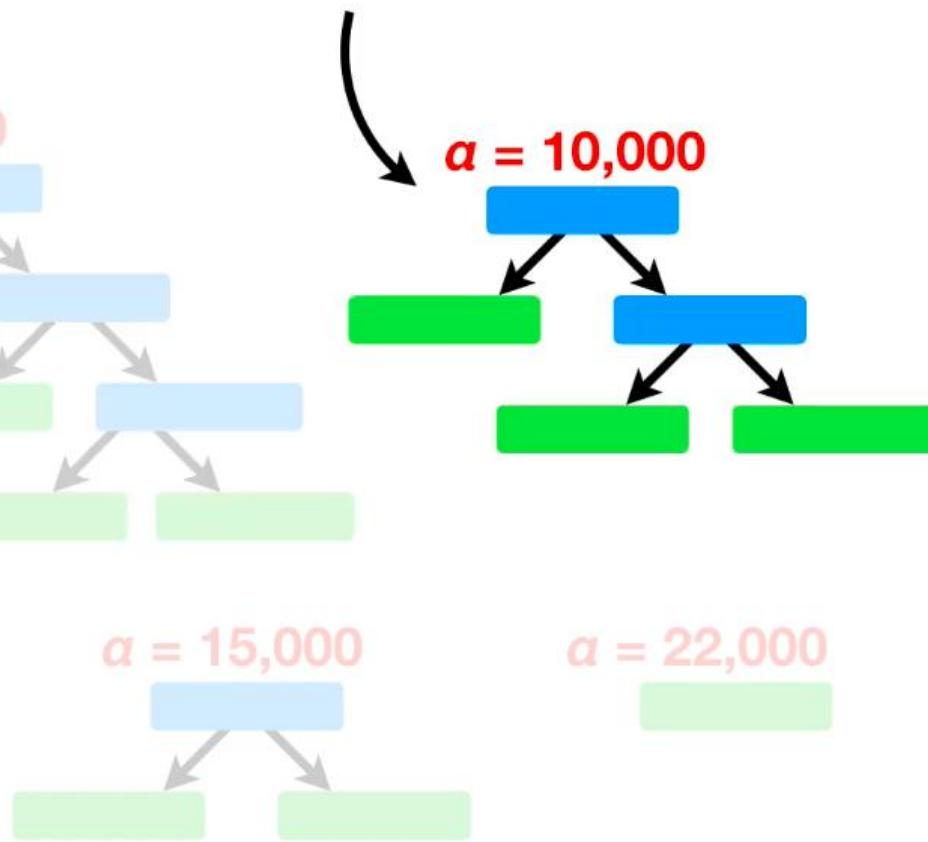
However, when $a = 10,000$,
we will get a lower **Tree Score**
if we prune these leaves...



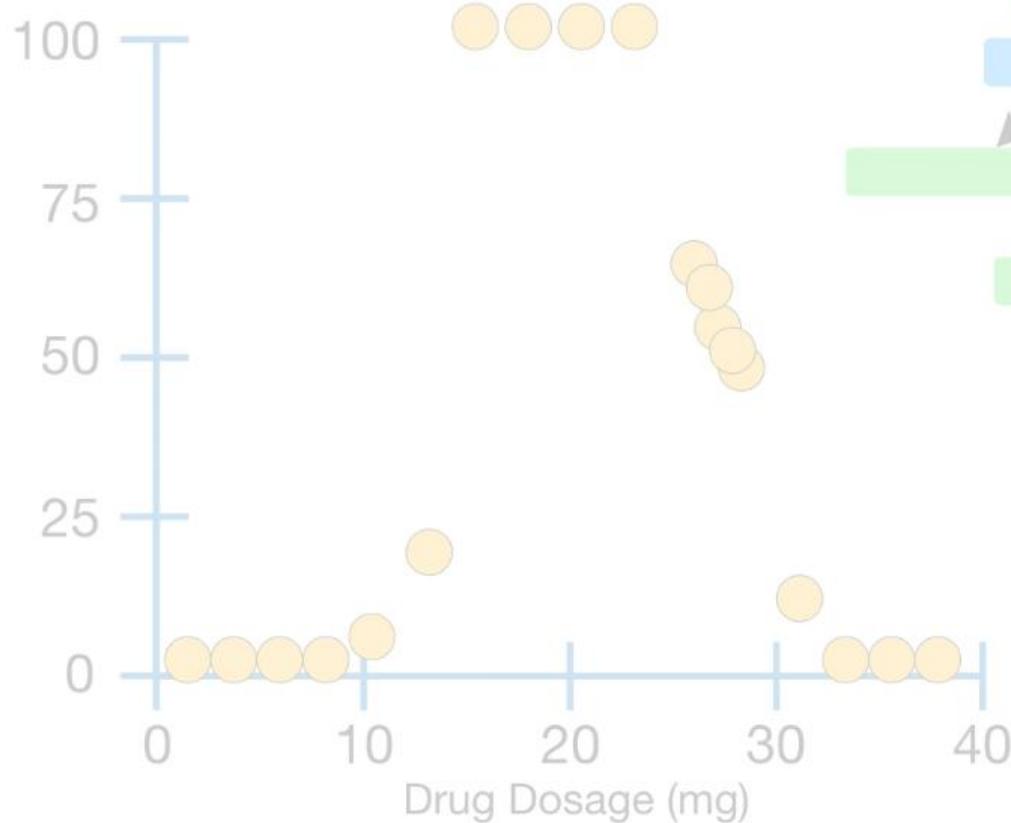
Tree Score = SSR + αT



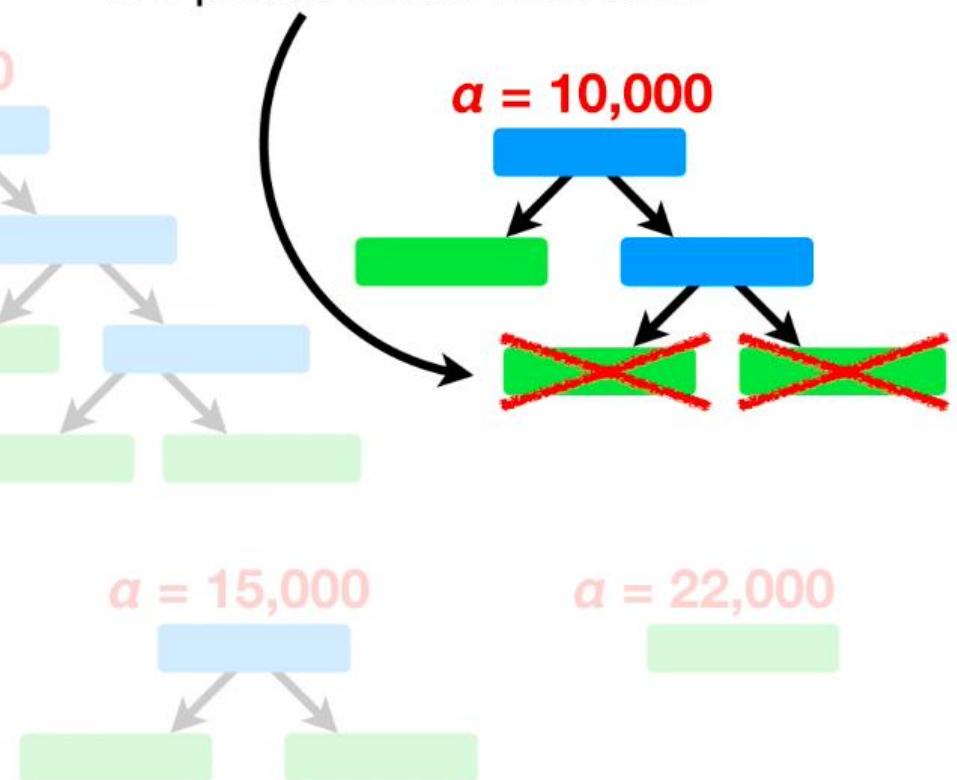
...and use this tree instead.



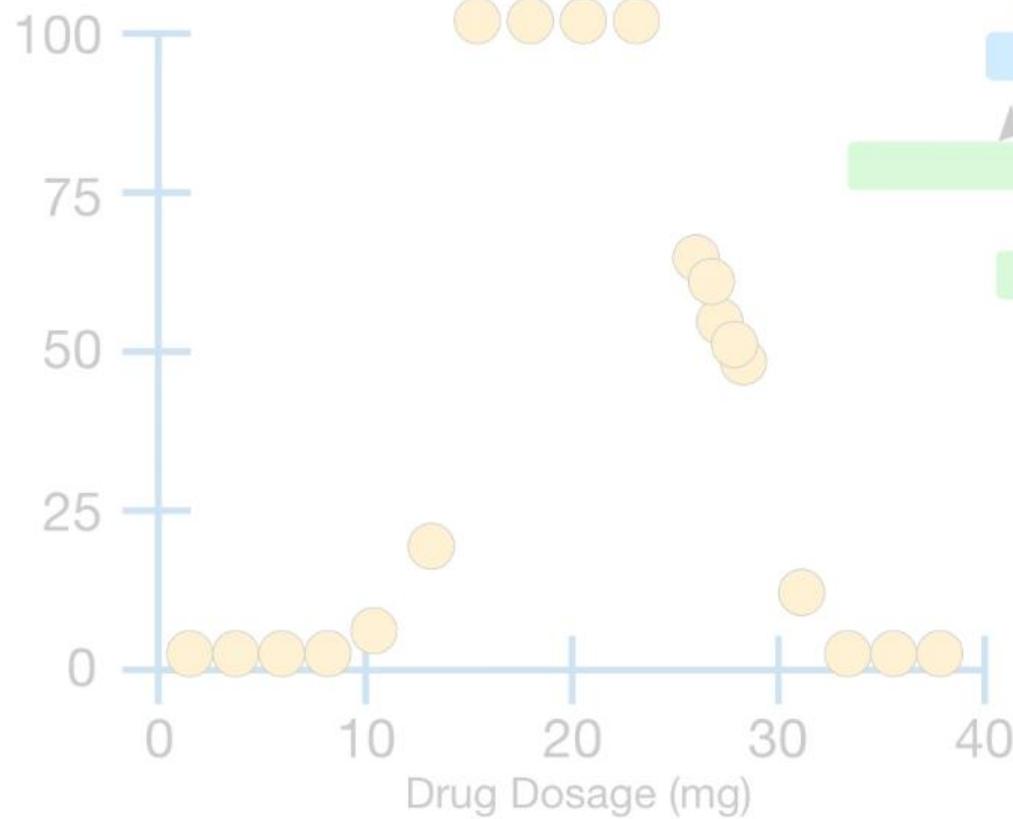
$$\text{Tree Score} = \text{SSR} + \alpha T$$



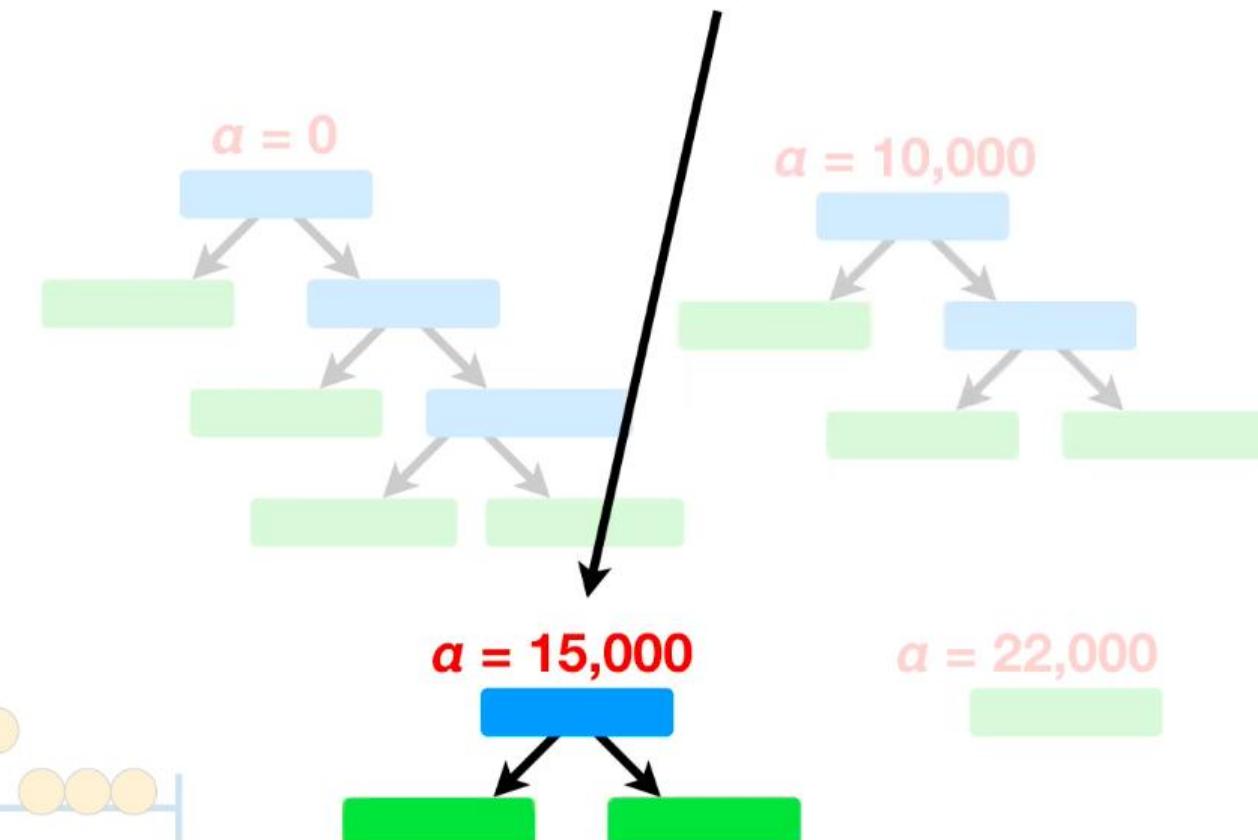
And when $\alpha = 15,000$, we will get a lower **Tree Score** if we prune these leaves...



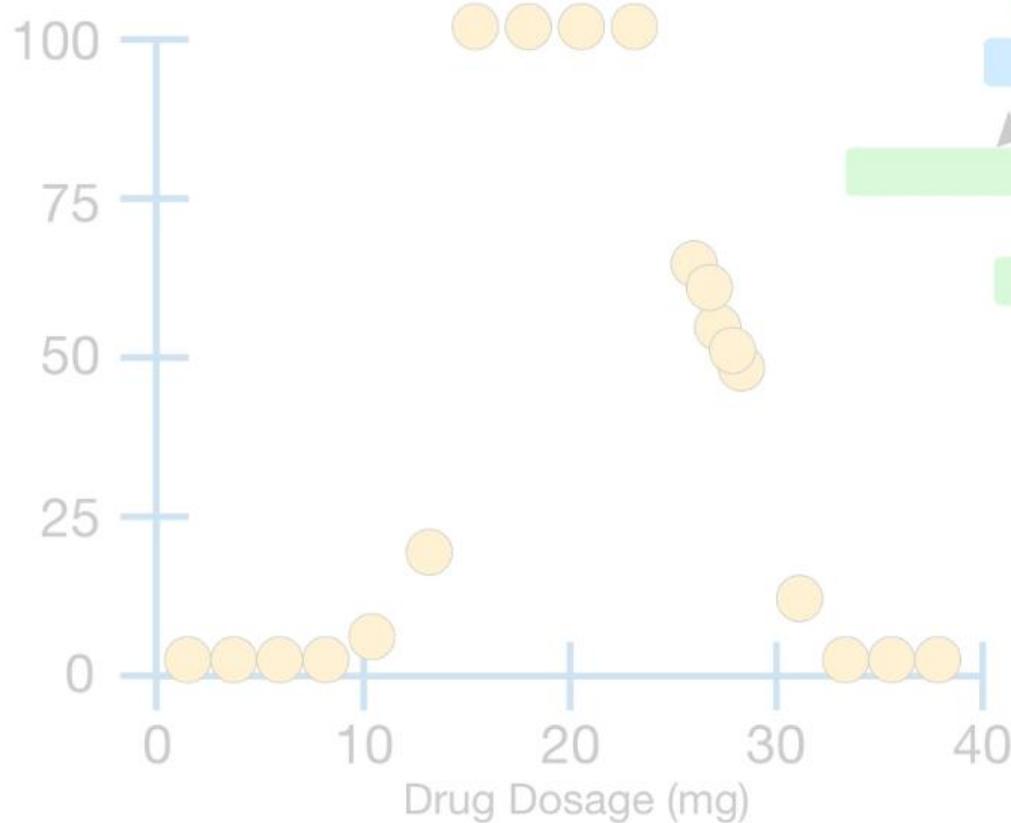
Tree Score = SSR + αT



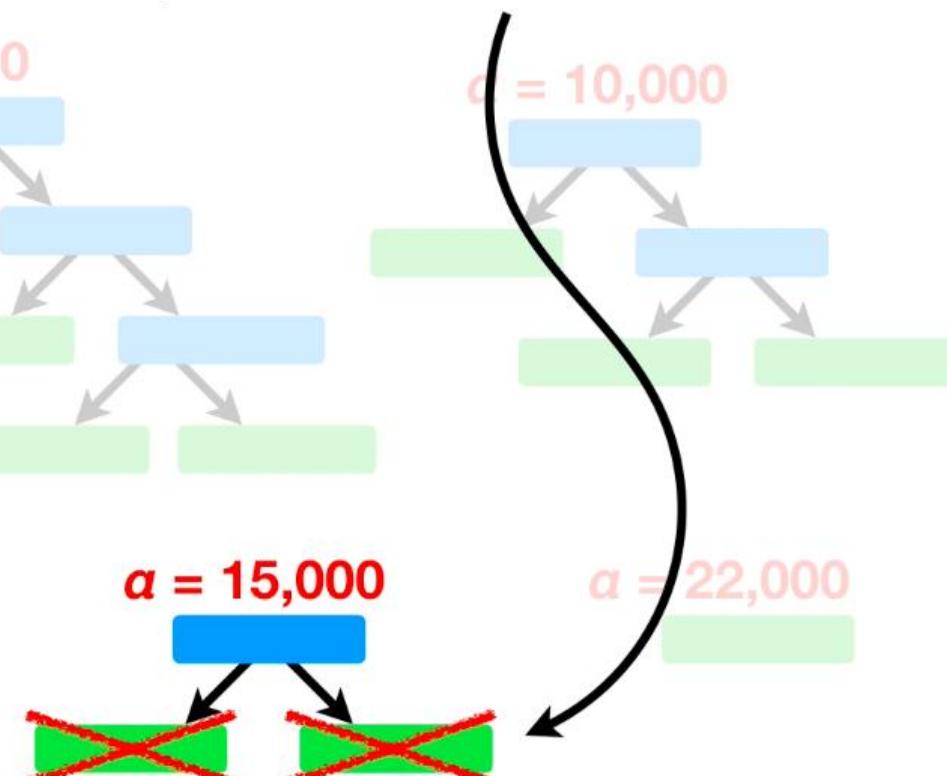
...and use this tree instead.



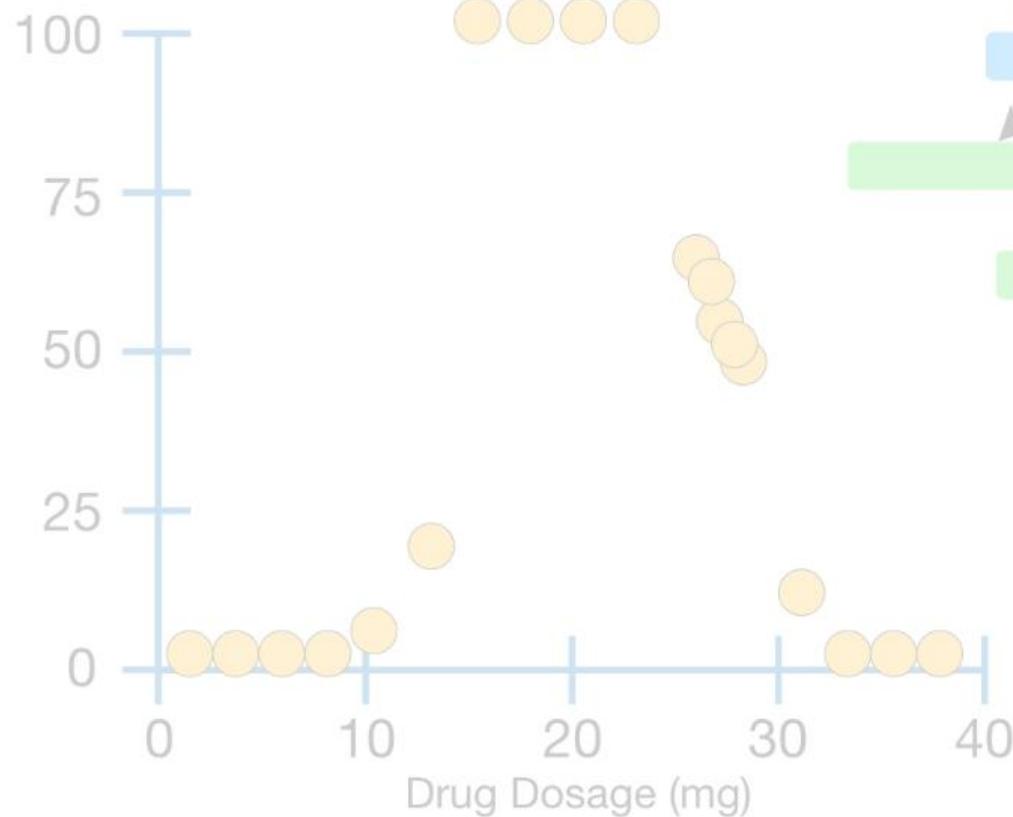
$$\text{Tree Score} = \text{SSR} + \alpha T$$



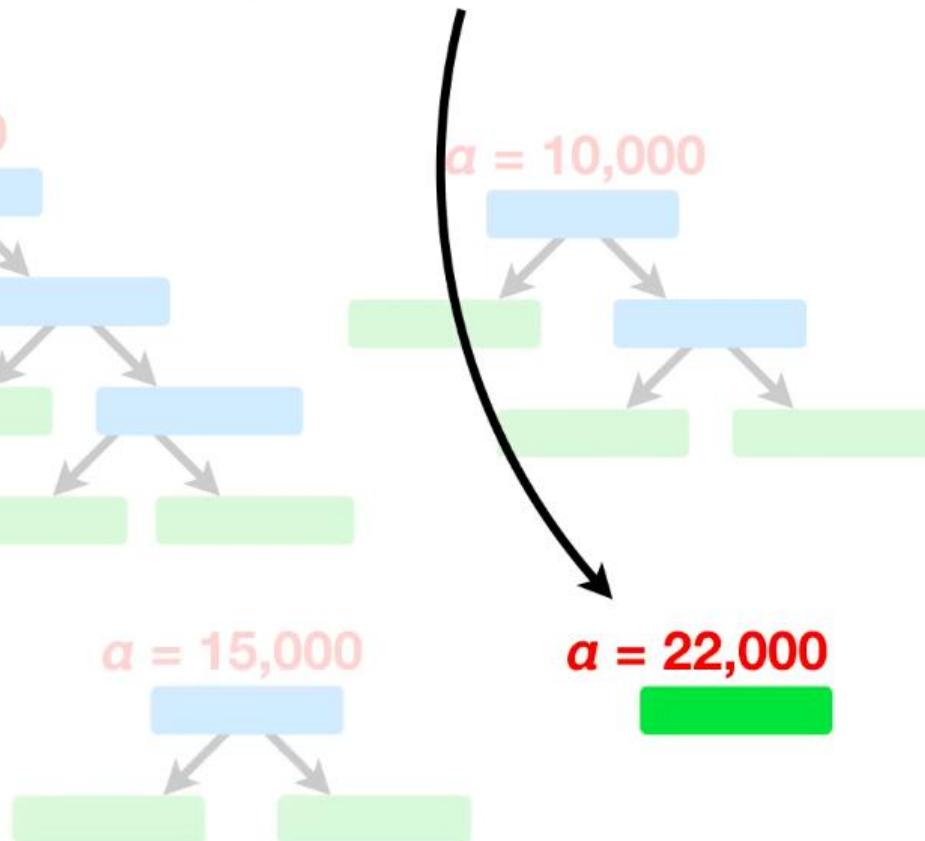
Lastly, when $\alpha = 22,000$, we will get a lower **Tree Score** if we prune these two leaves...



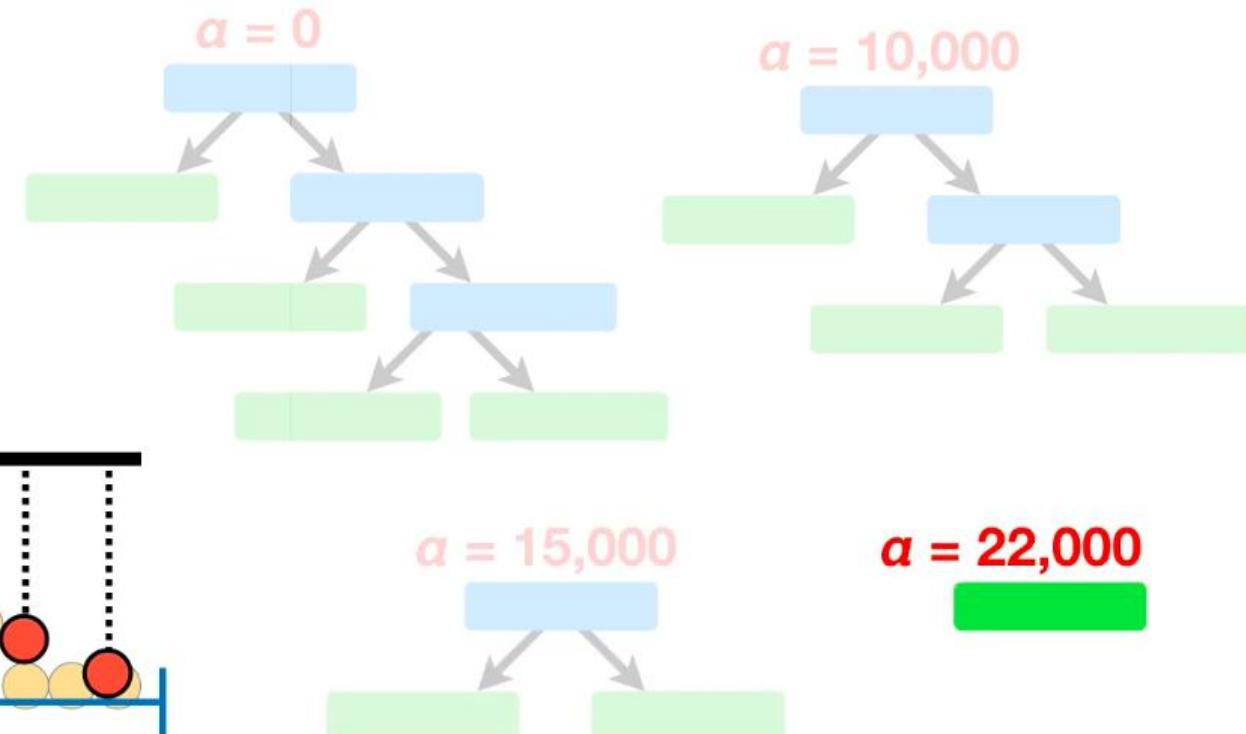
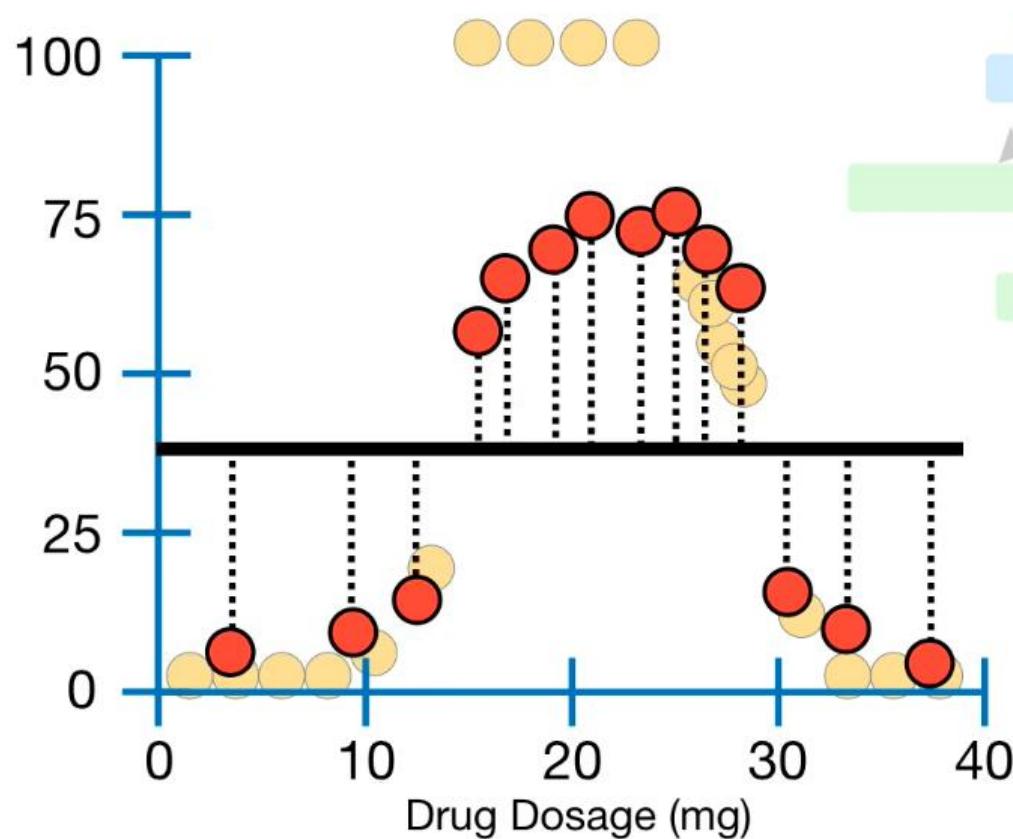
Tree Score = SSR + αT



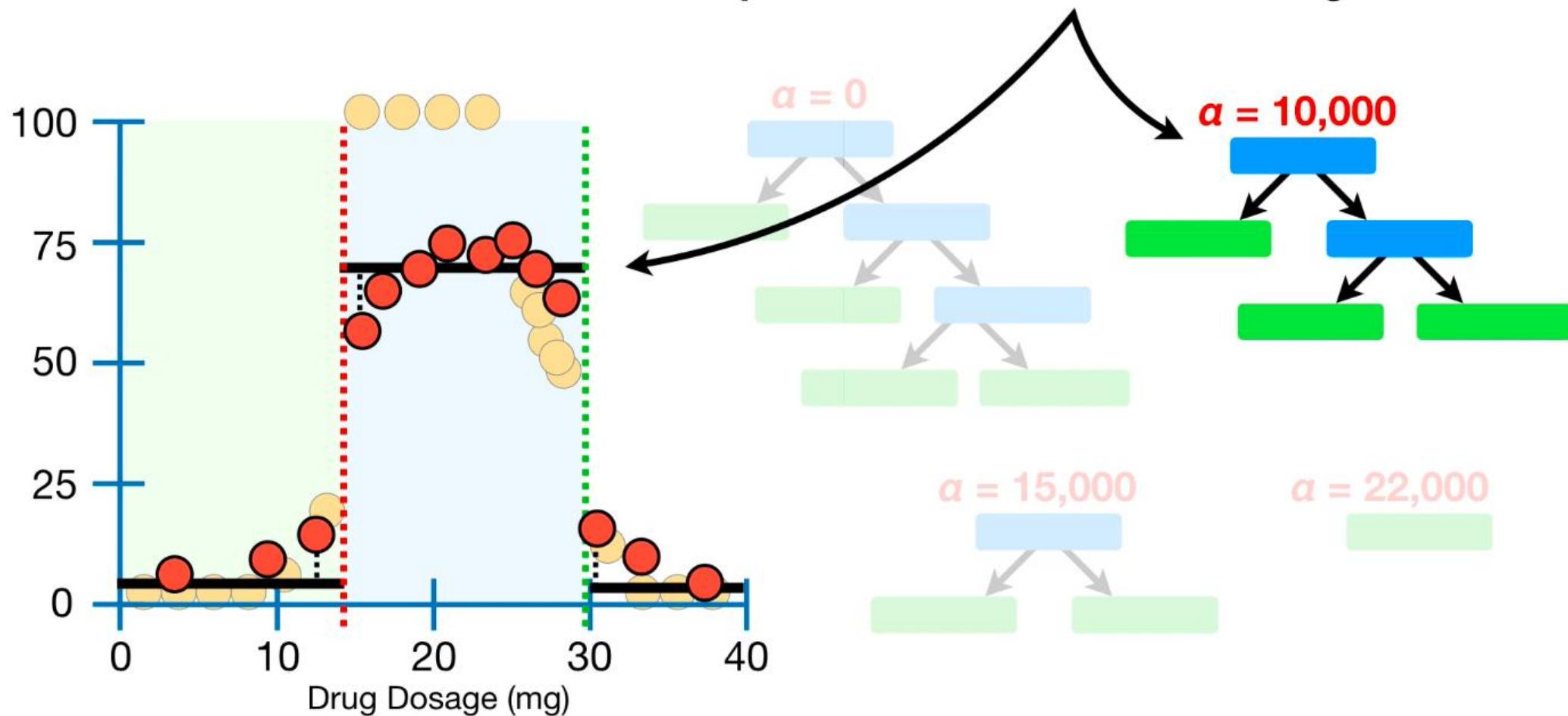
...and use this tree instead.



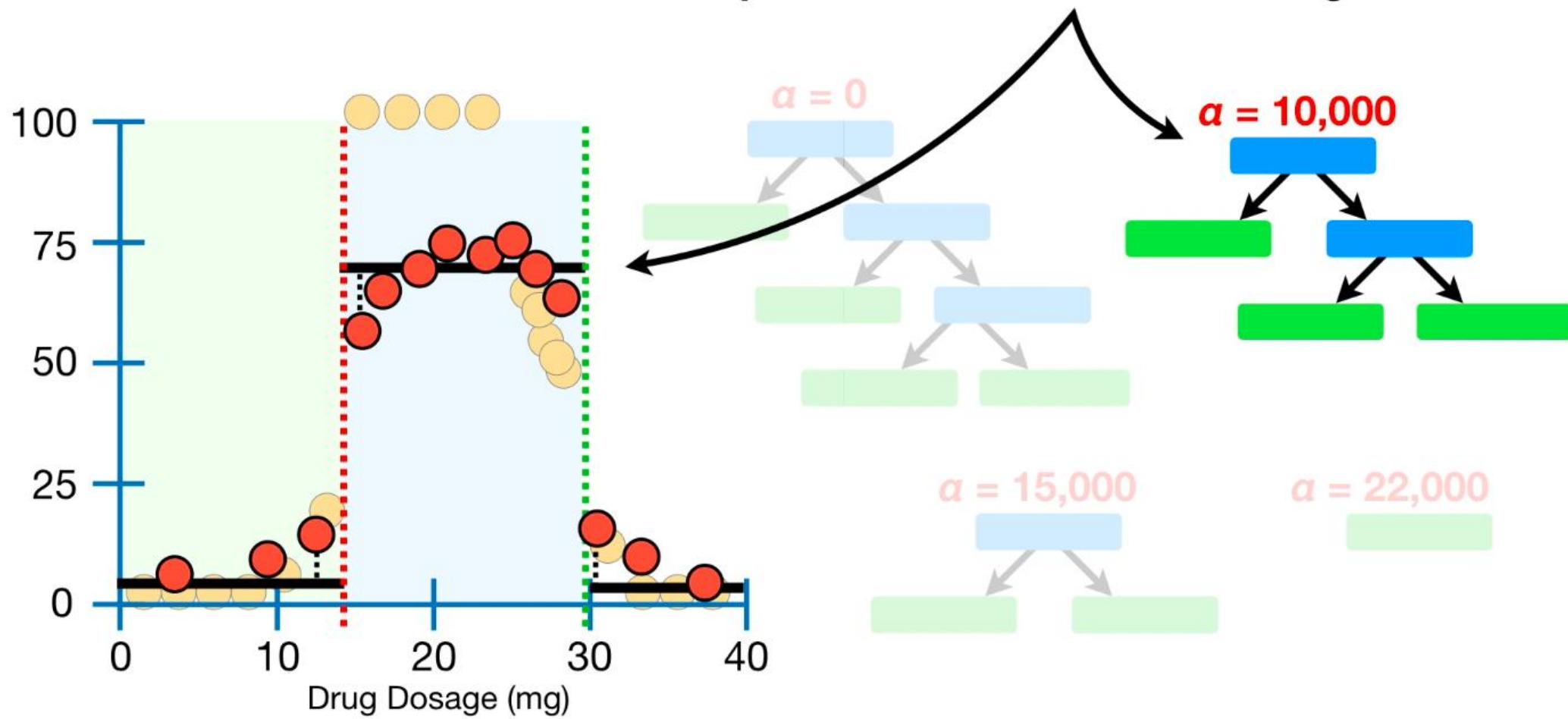
Now calculate the **Sum of Squared Residuals** for each new tree using only the **Testing Data**.



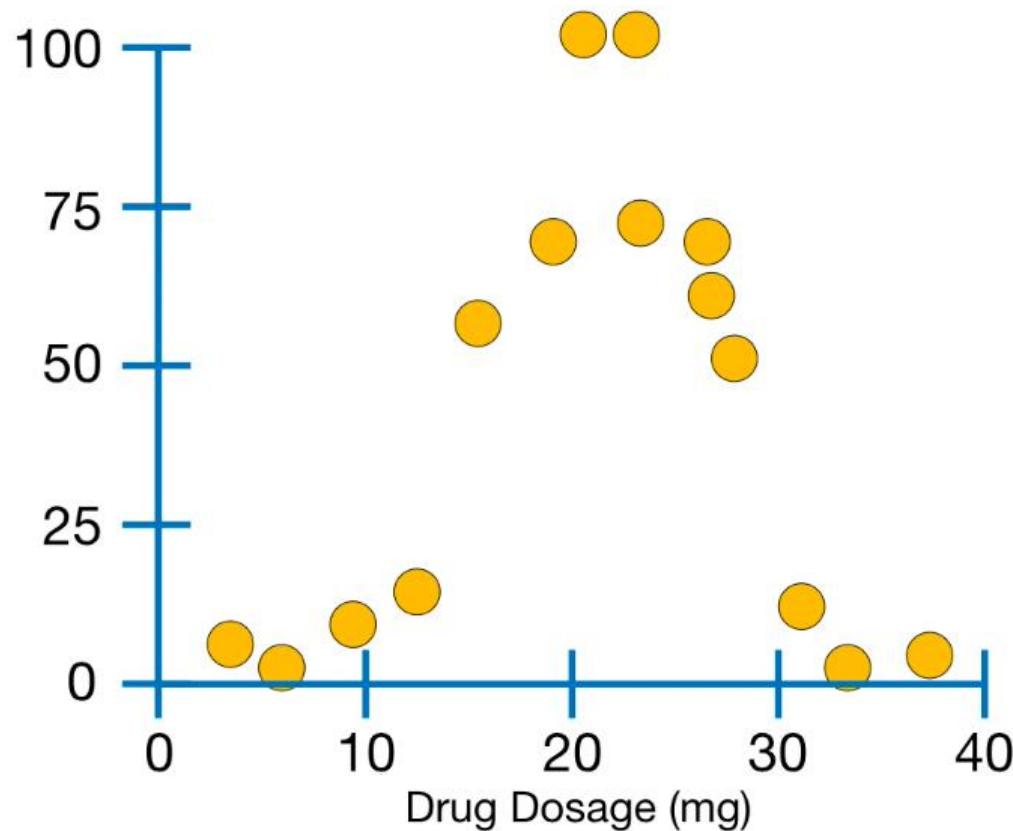
In this case, the tree with $\alpha = 10,000$ had the smallest
Sum of Squared Residuals for the **Testing Data**.



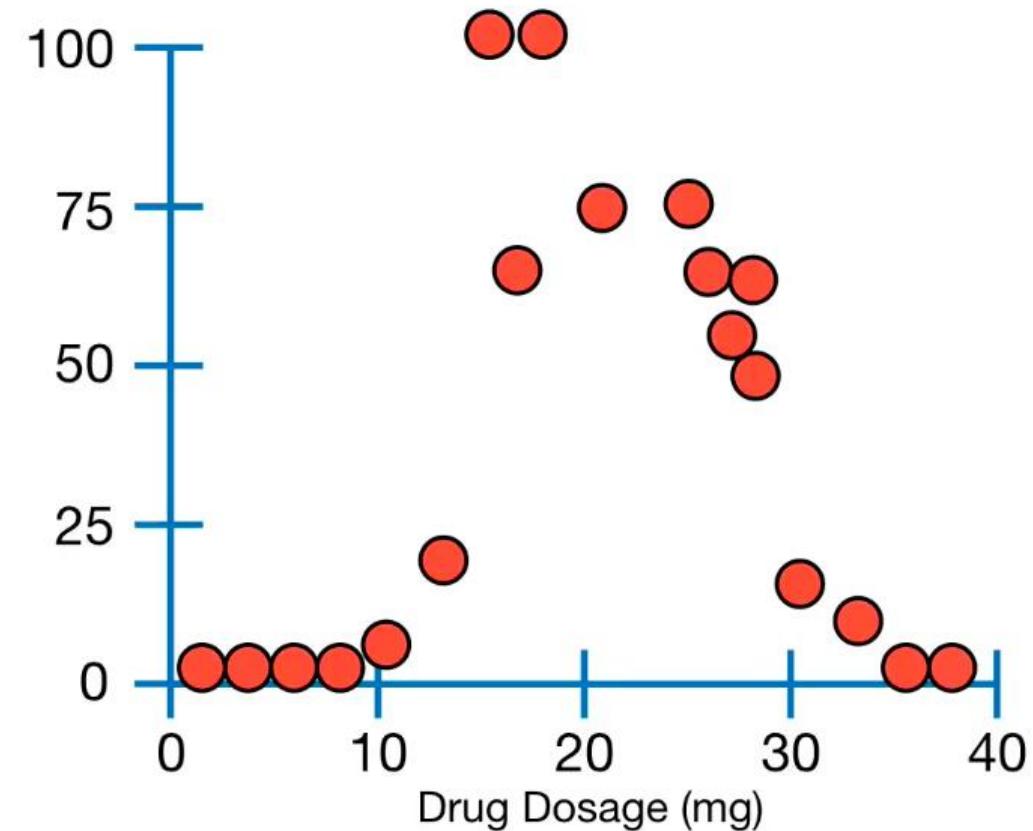
In this case, the tree with $\alpha = 10,000$ had the smallest
Sum of Squared Residuals for the **Testing Data**.



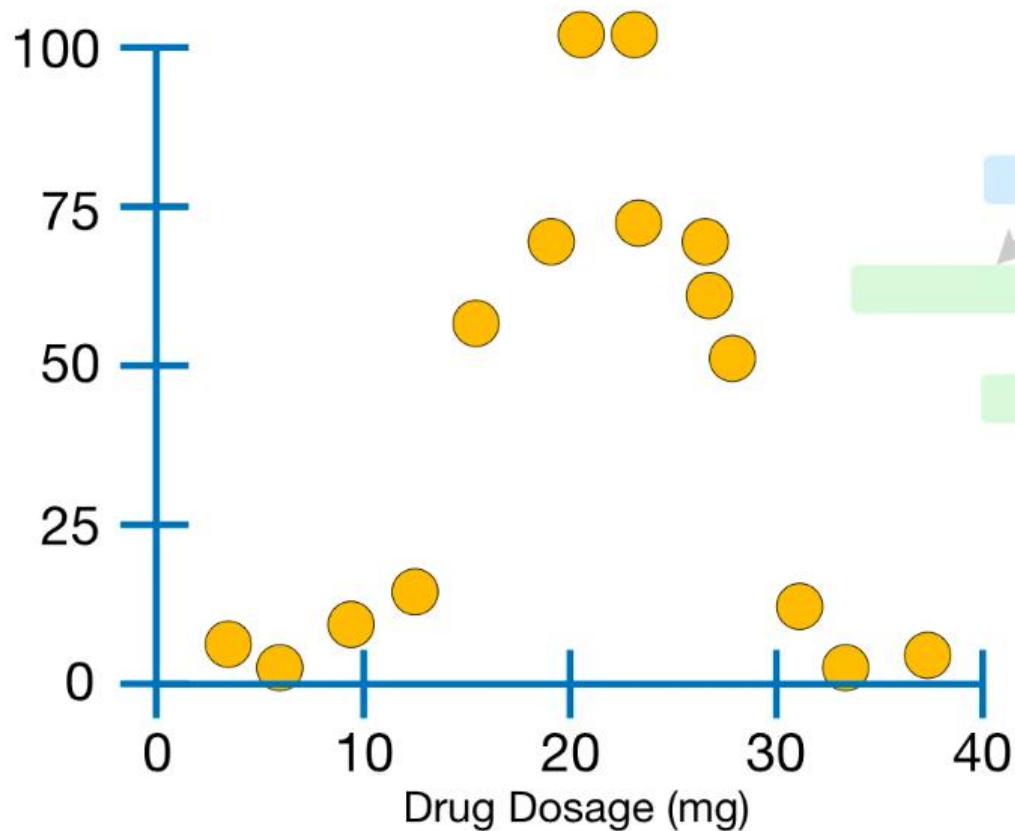
Now we go back and create
new **Training Data**...



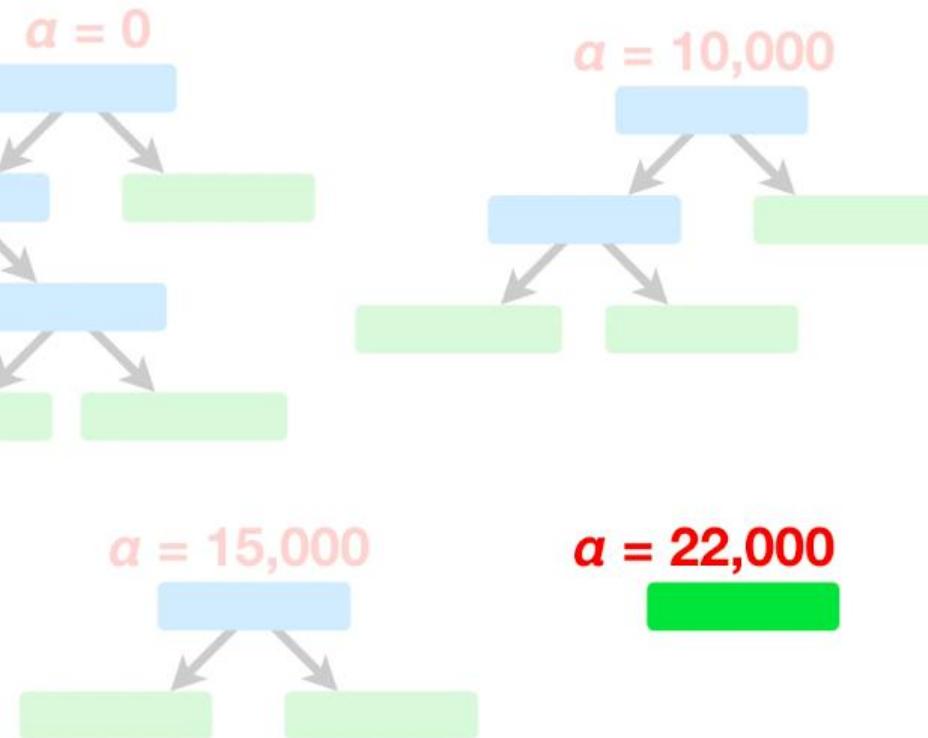
...and new **Testing Data**.



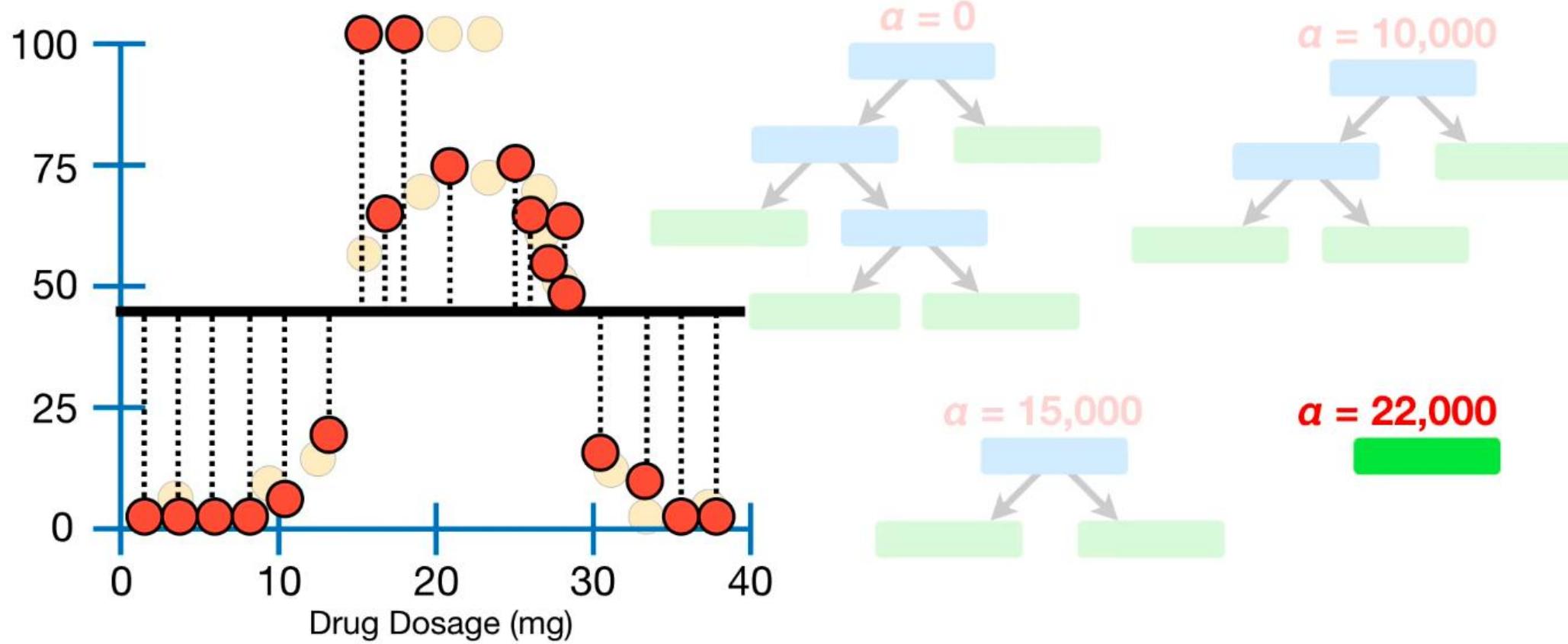
And just using the new
Training Data...

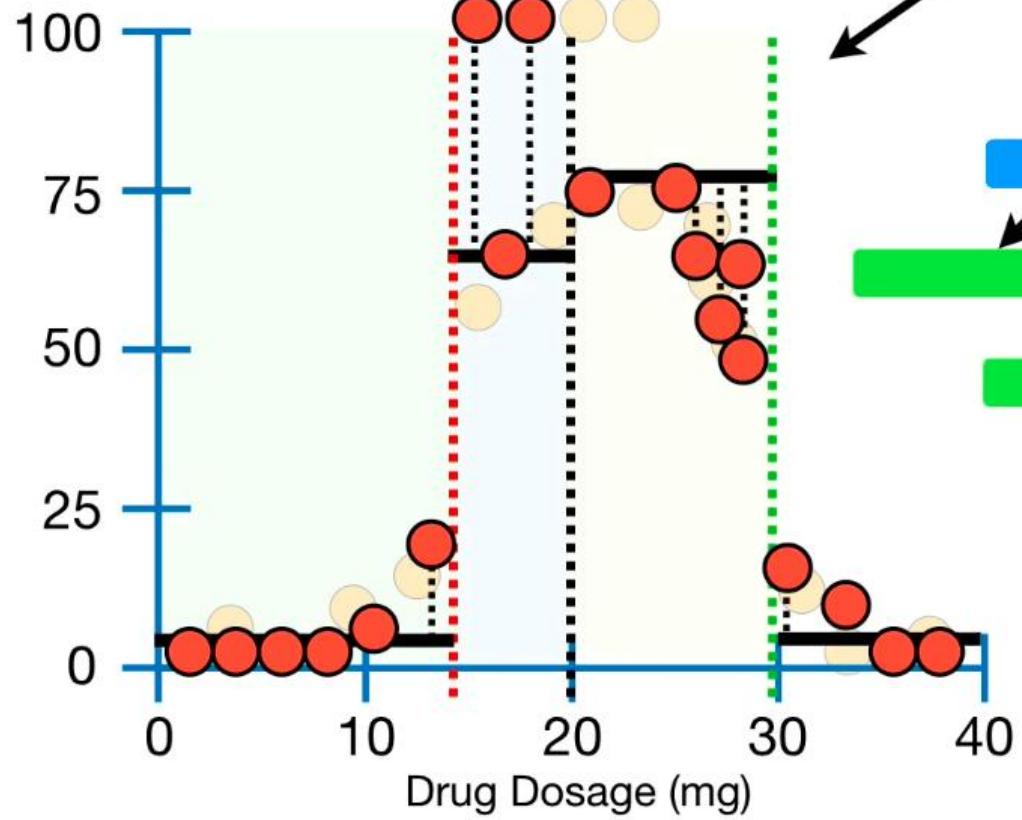


...build a new sequence of trees, from full sized to a leaf, using the α values we found before.



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





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

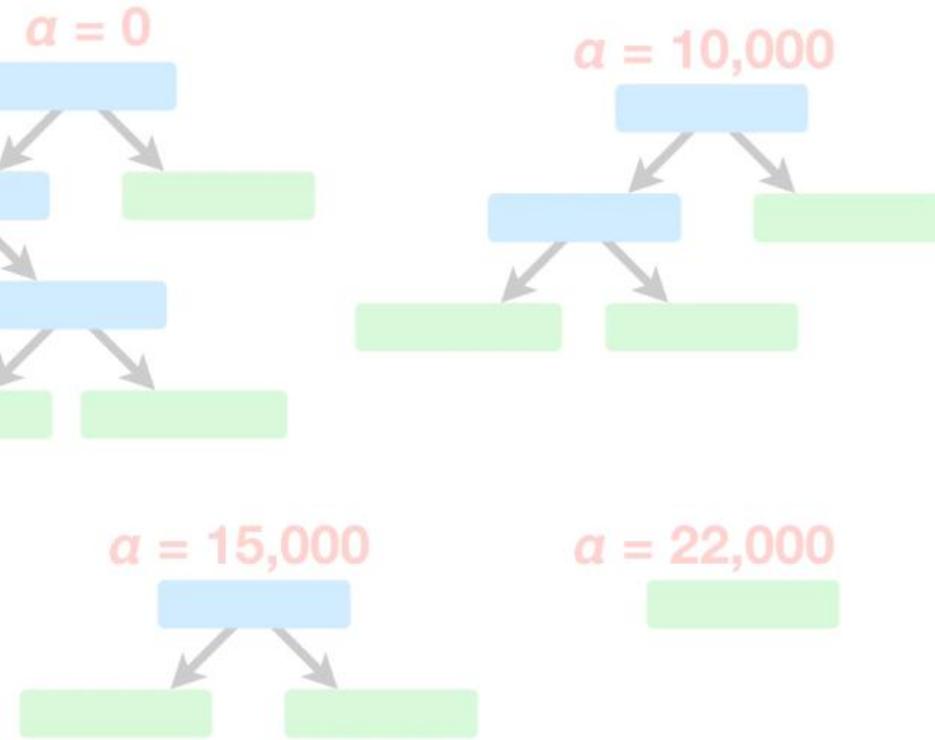
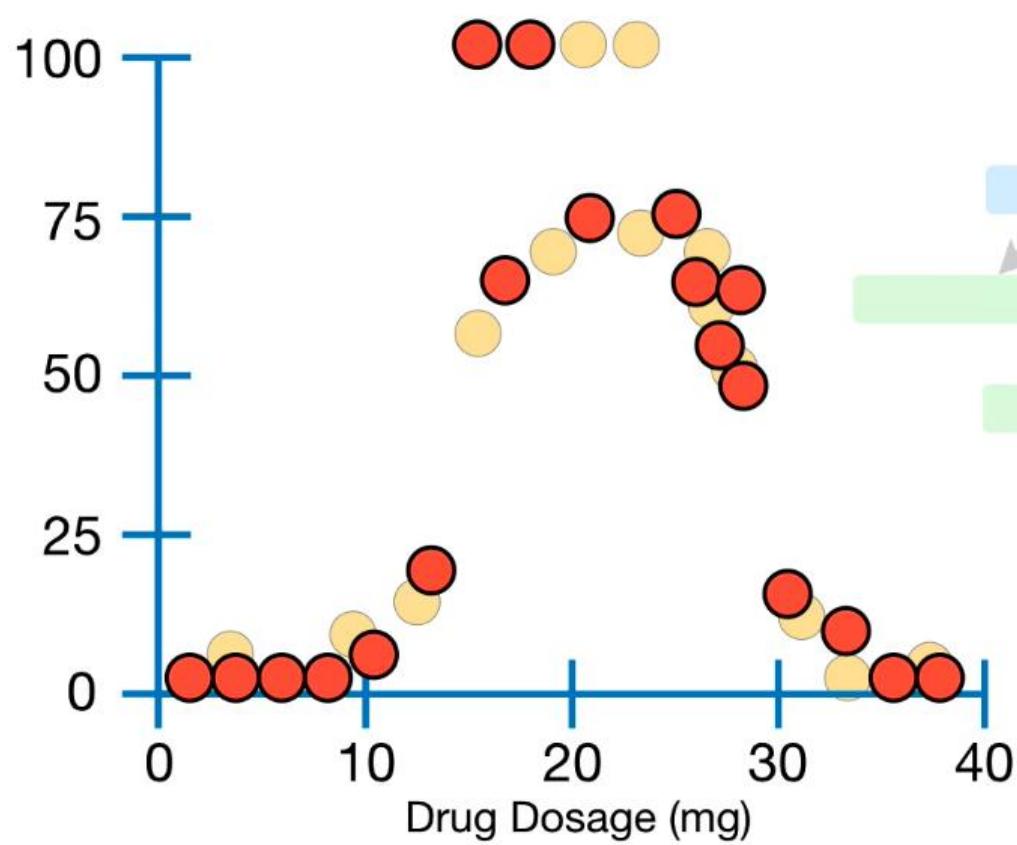
$a = 0$

$a = 10,000$

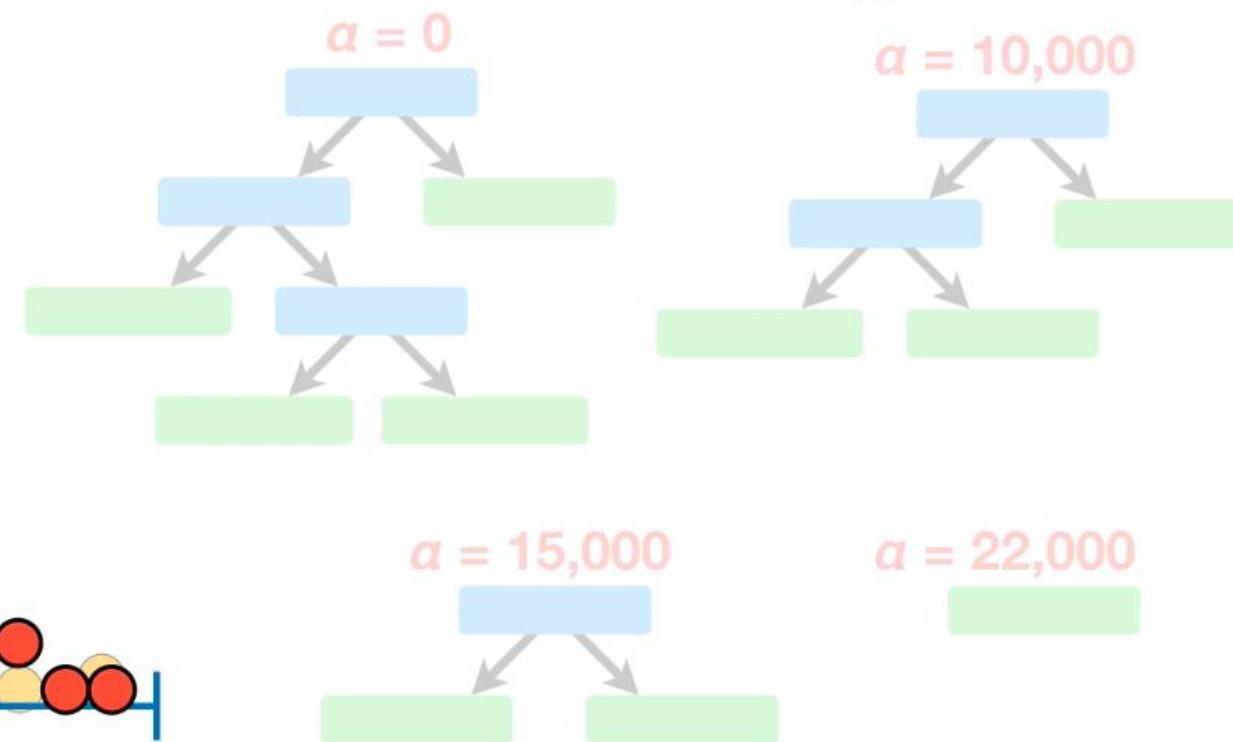
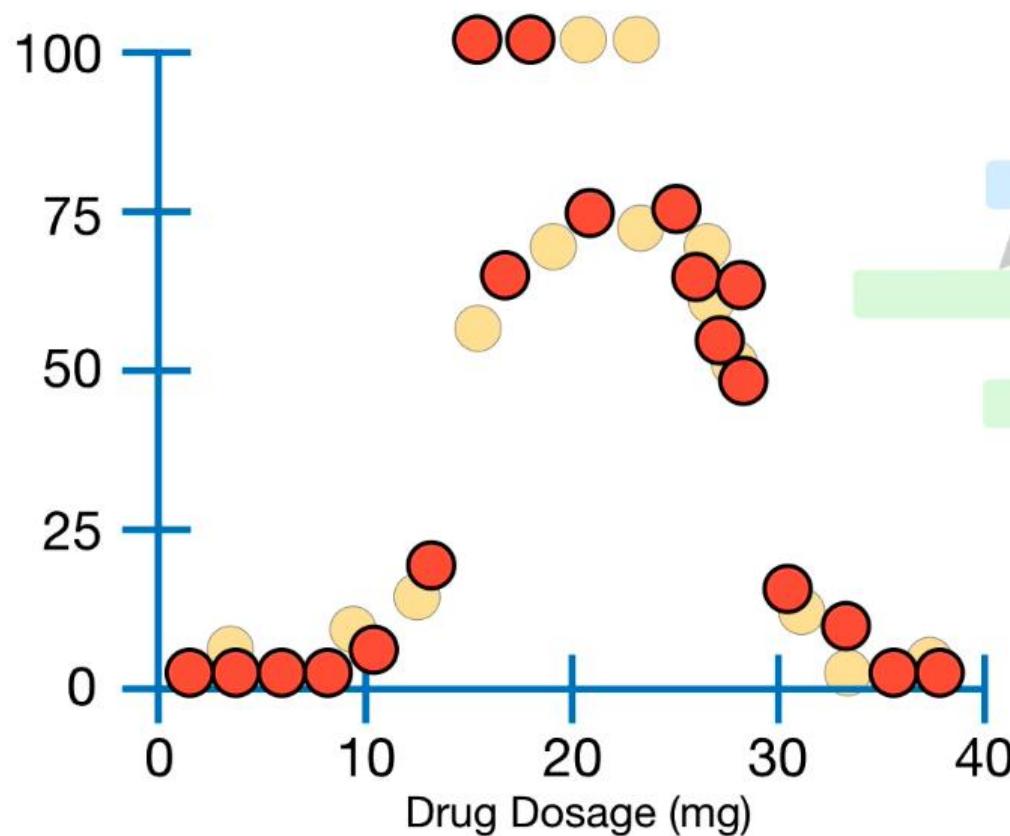
$a = 15,000$

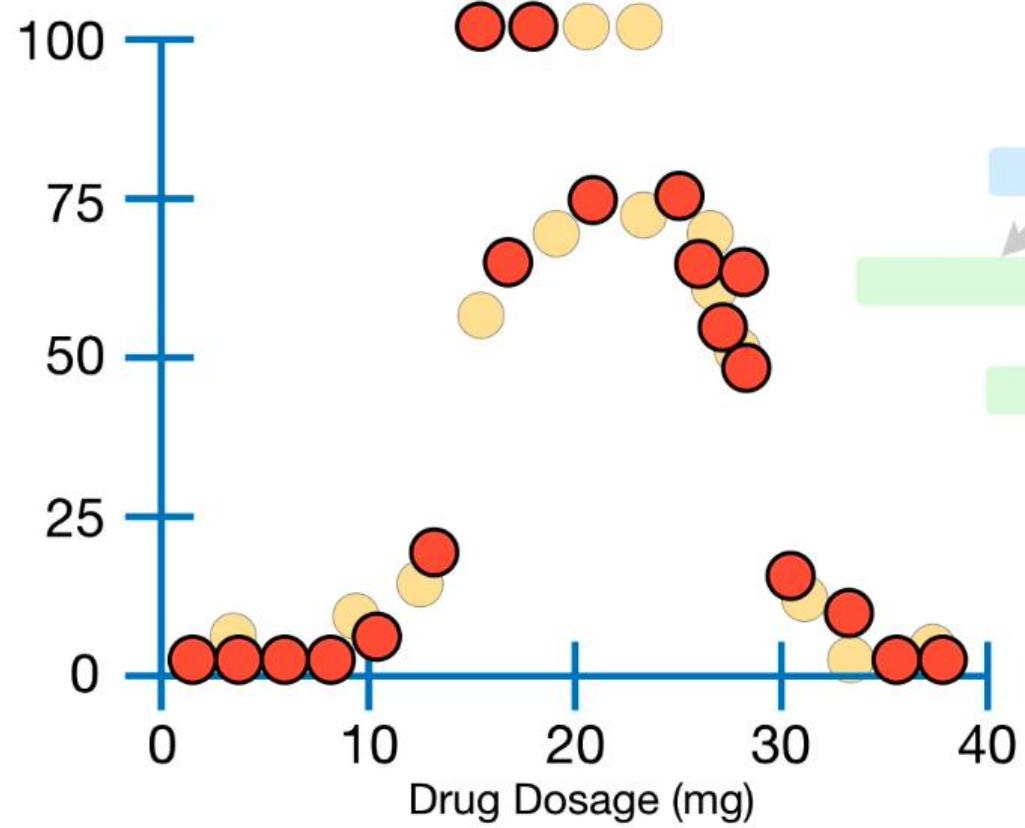
$a = 22,000$

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

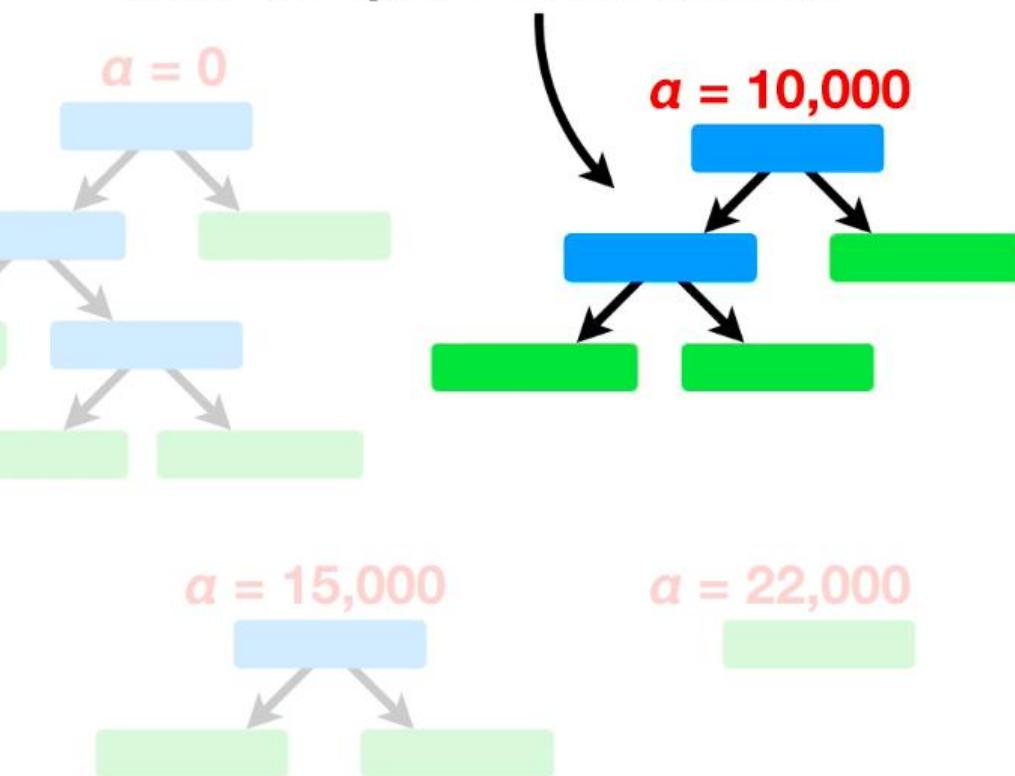


...and the value for a that, on average, gave us the lowest **Sum of Squared Residuals** with the **Testing Data**, is the final value for a .

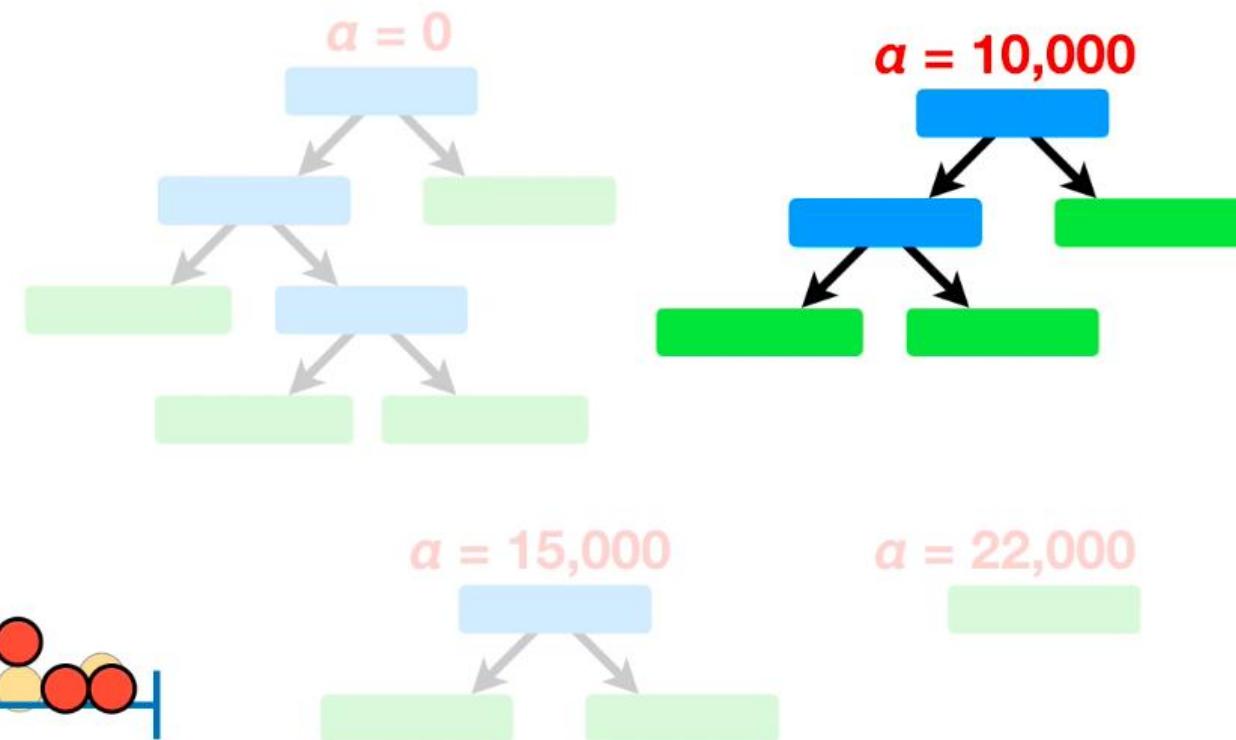
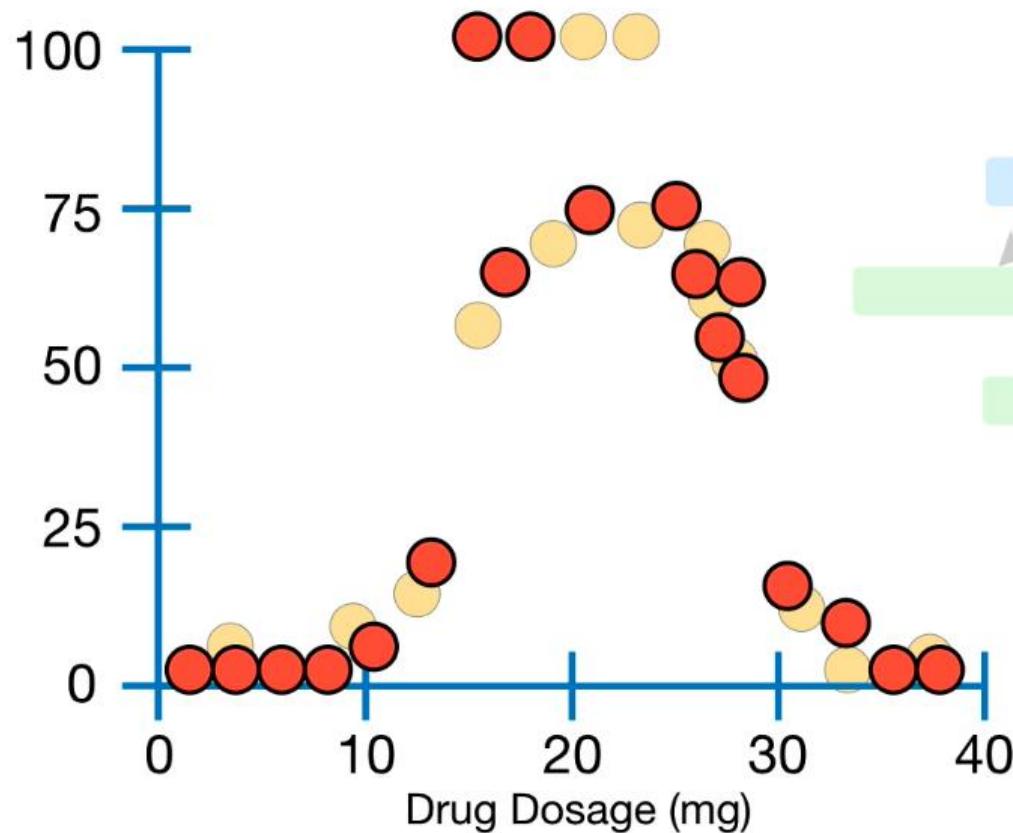




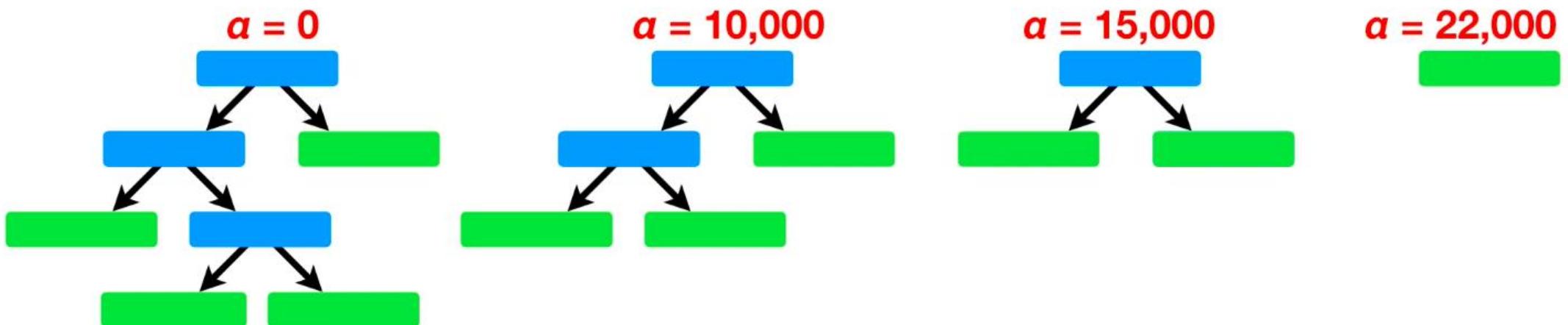
In this case, the optimal trees built with $\alpha = 10,000$ had, on average, the lowest **Sum of Squared Residuals...**



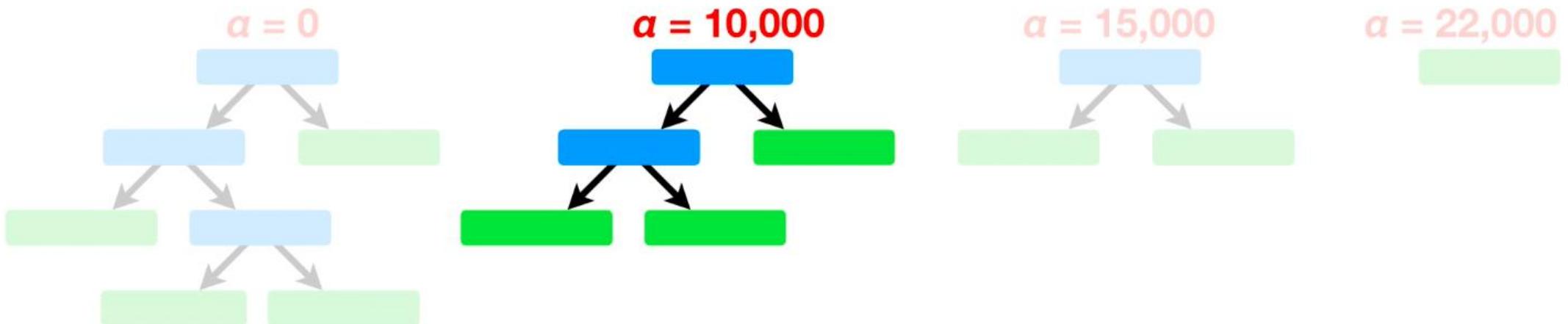
...so $\alpha = 10,000$ is our final value.



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 a that we selected ($a = 10,000$).



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

