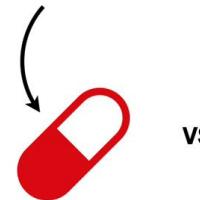


Imagine we developed  
a new drug...



Imagine we developed  
a new drug...



...to cure the  
common cold.



vs.

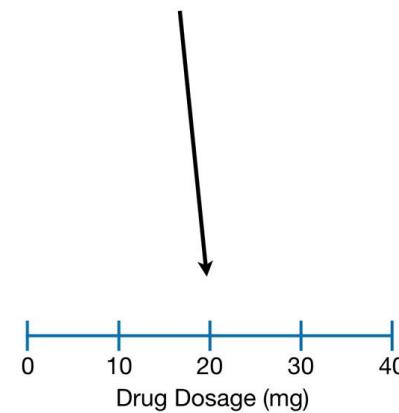
However, we don't know the optimal  
dosage to give to patients.



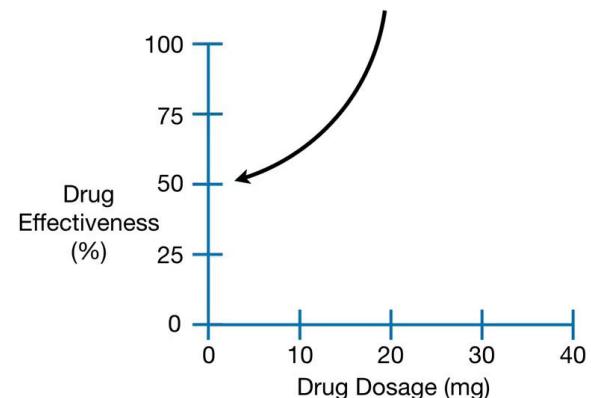
vs.



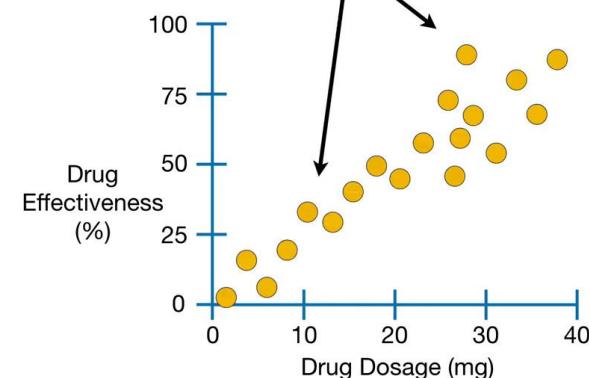
So we do a clinical trial with  
different dosages...



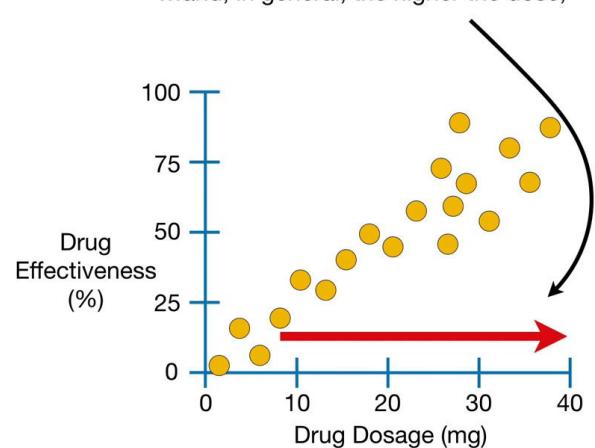
...and measure how effective each dosage is.



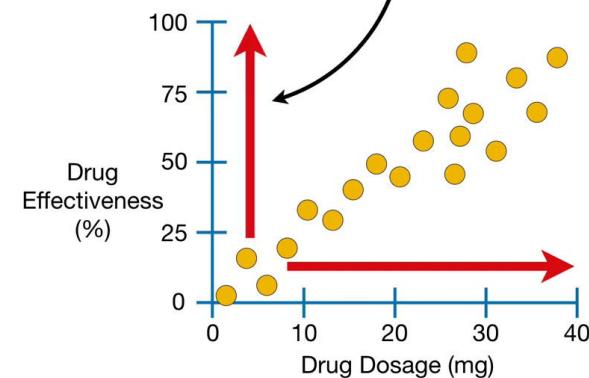
If the data looked like this...



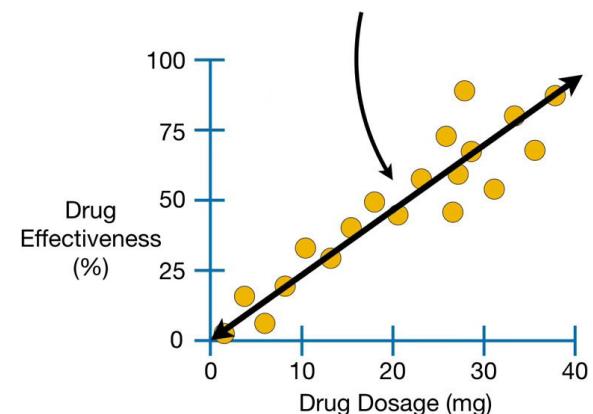
...and, in general, the higher the dose,



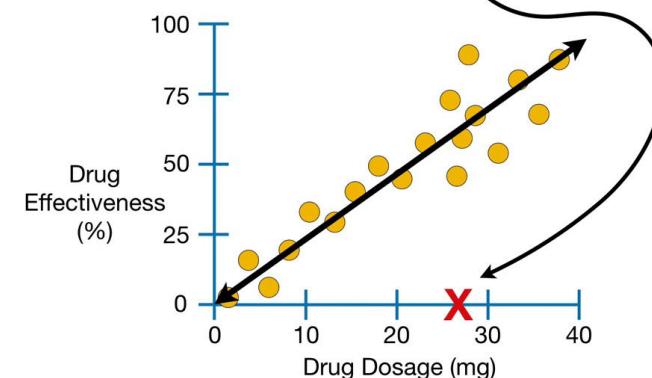
...and, in general, the higher the dose,  
the more effective the drug...



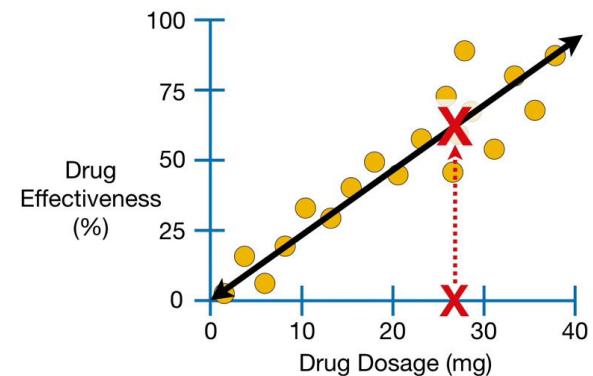
...then we could easily fit a line to the data...



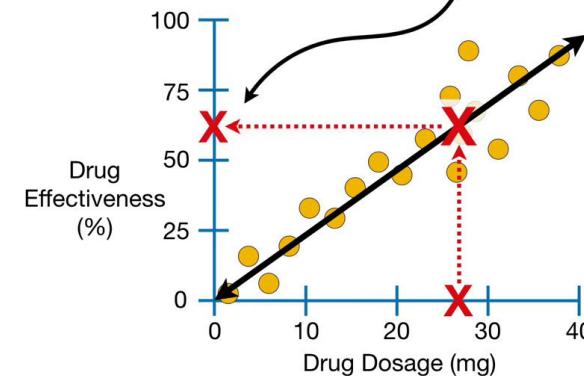
...and if someone told us they were taking a **27 mg Dose**...



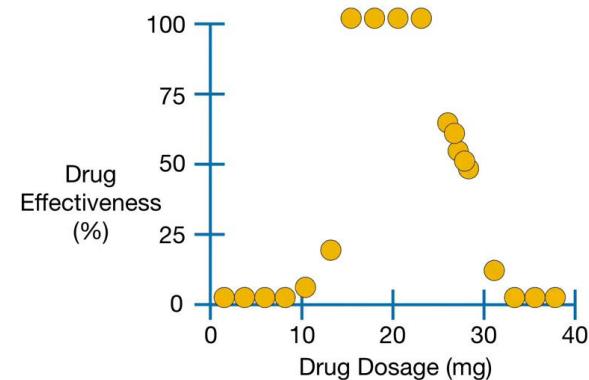
...we could use the line to predict that a **27 mg Dose** should be **62% Effective**.



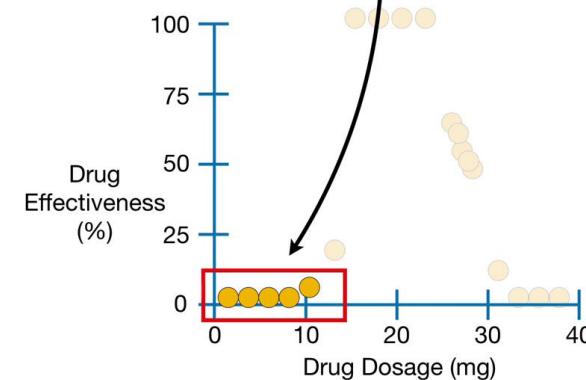
...we could use the line to predict that a **27 mg Dose** should be **62% Effective**.



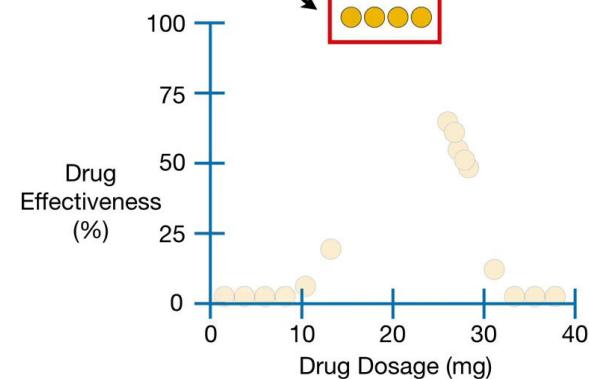
However, what if the data looked like this?



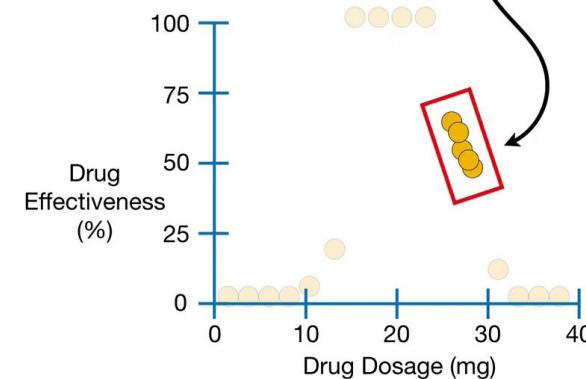
Low dosages are not effective...



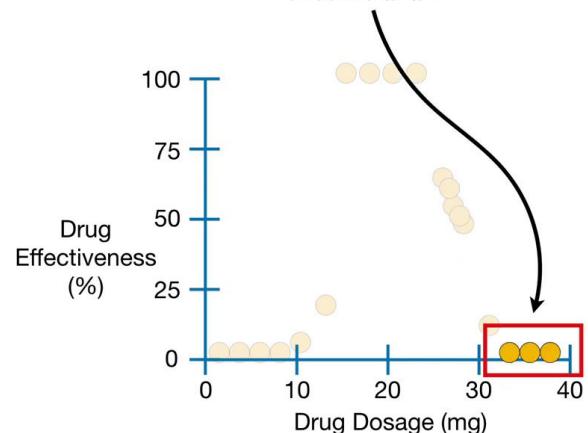
...moderate dosages work really well...



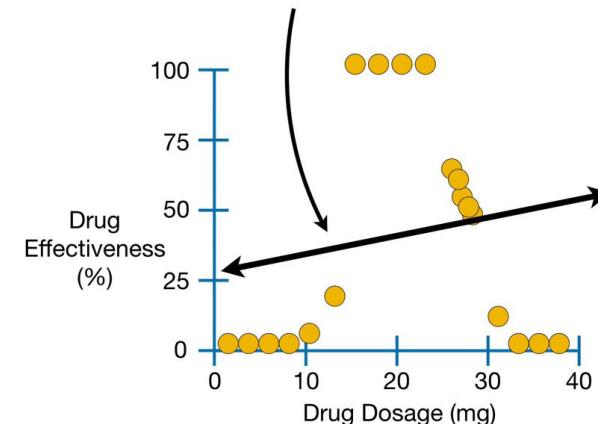
...somewhat higher dosages work at about 50% effectiveness...



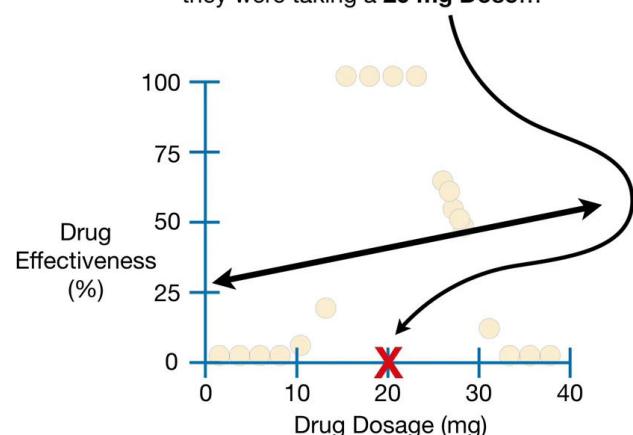
...and high dosages are not effective at all.



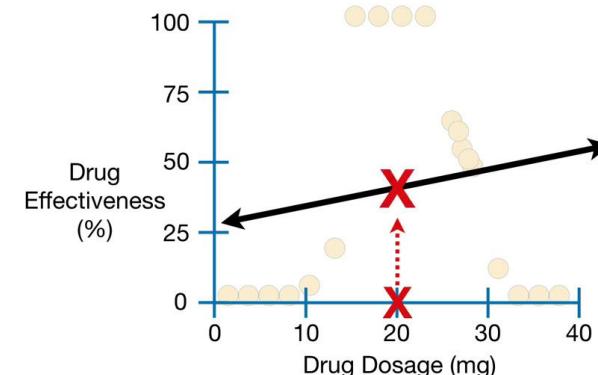
In this case, fitting a straight line to the data will not be very useful.



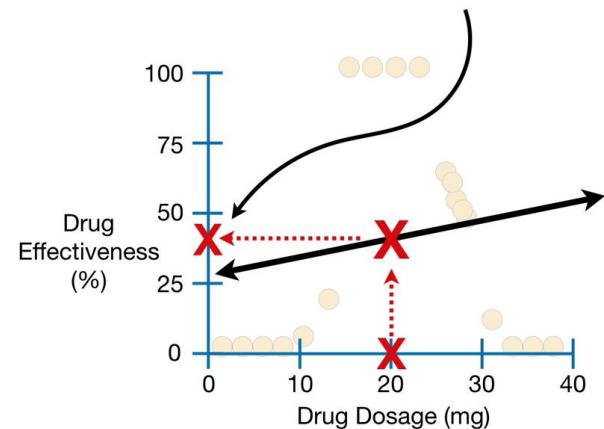
For example, if someone told us they were taking a **20 mg Dose**...



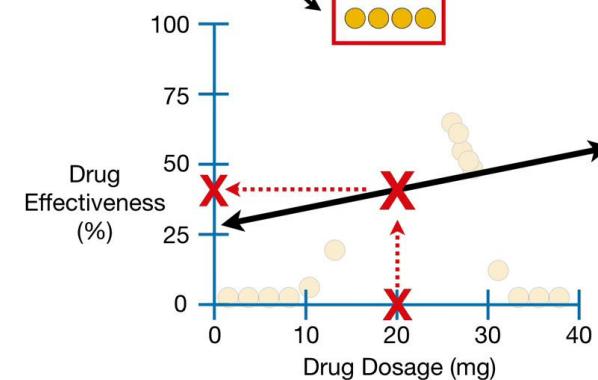
...then we would predict that a **20 mg Dose** should be **45% Effective**...



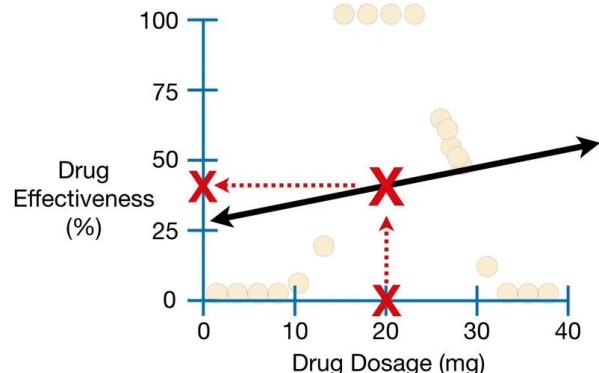
...then we would predict that a **20 mg Dose** should be **45% Effective**...



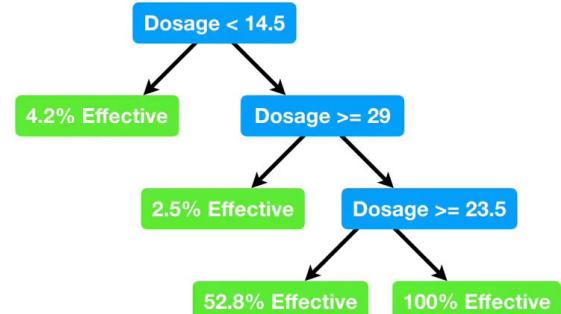
...even though the observed data says that it should be **100% Effective**.



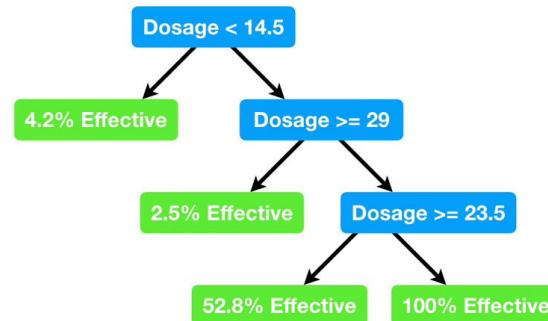
So we need to use something other than a straight line to make predictions.



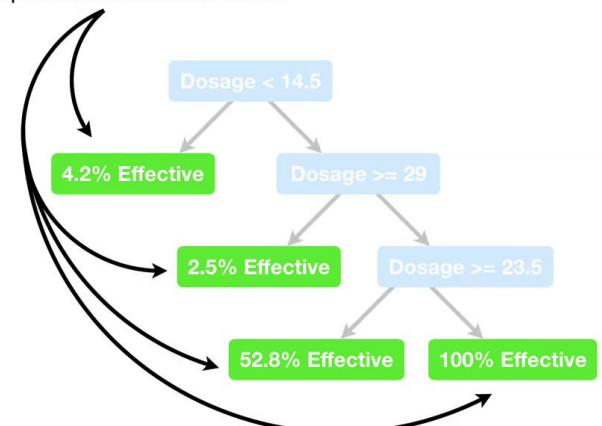
One option is to use a **Regression Tree**.



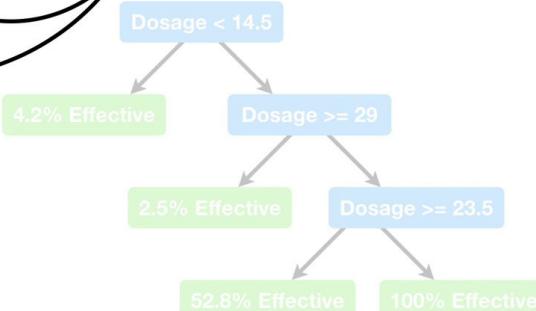
**Regression Trees** are a type of Decision Tree.



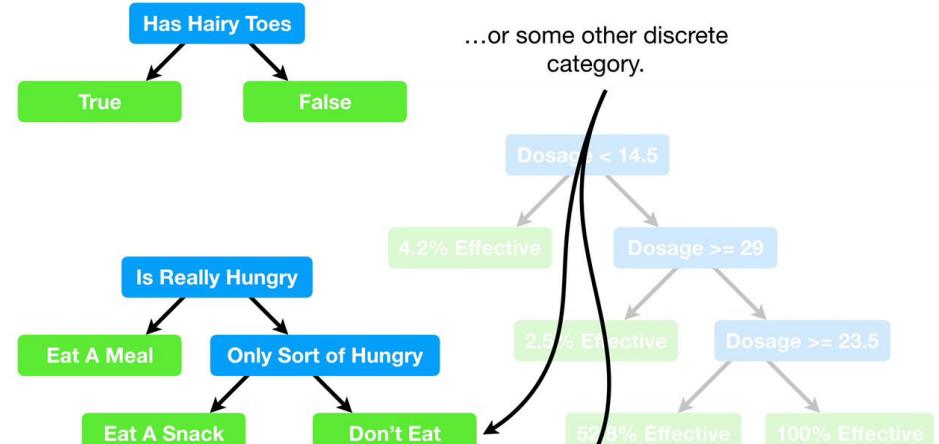
In a **Regression Tree**, each leaf represents a numeric value.



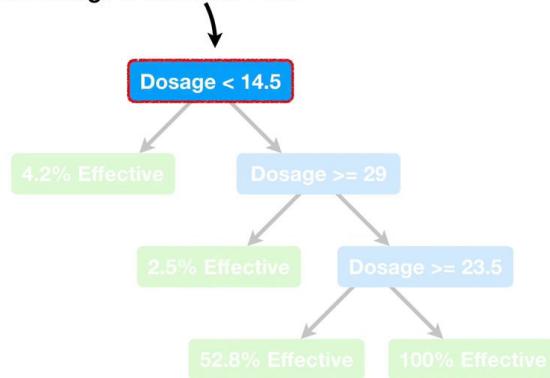
In contrast, **Classification Trees** have True or False in their leaves...



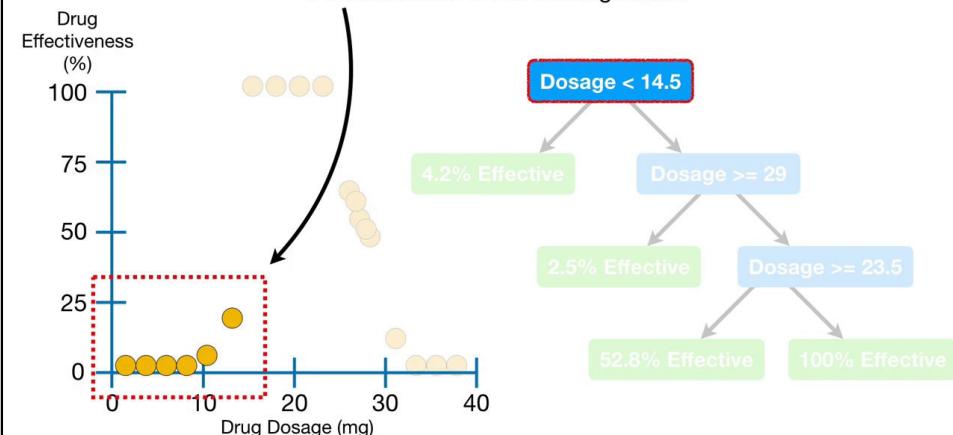
...or some other discrete category.



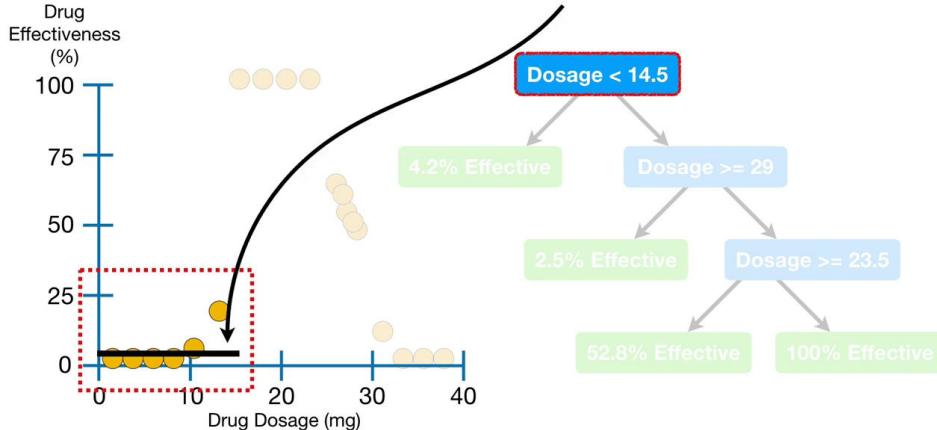
With this **Regression Tree**, we start by asking if the **Dosage** is less than **14.5**.



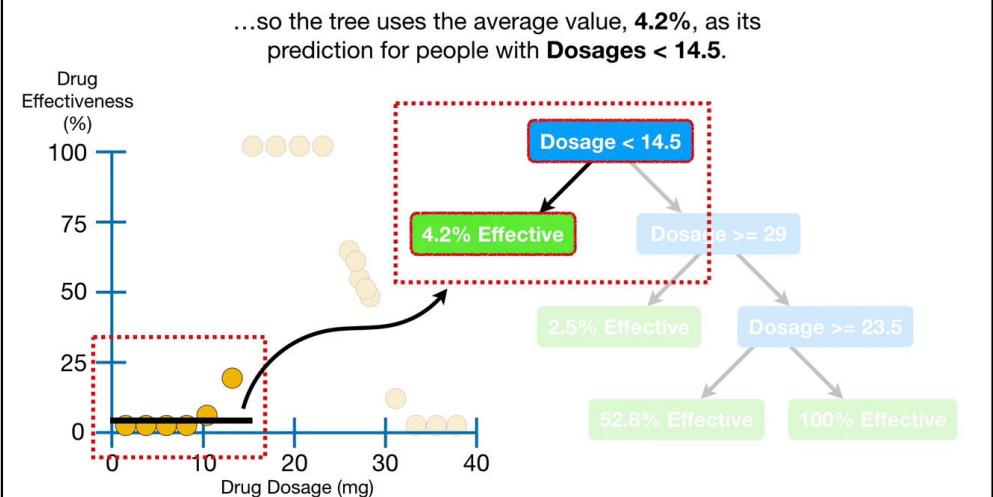
...if so, then we are talking about these 6 observations in the training data...

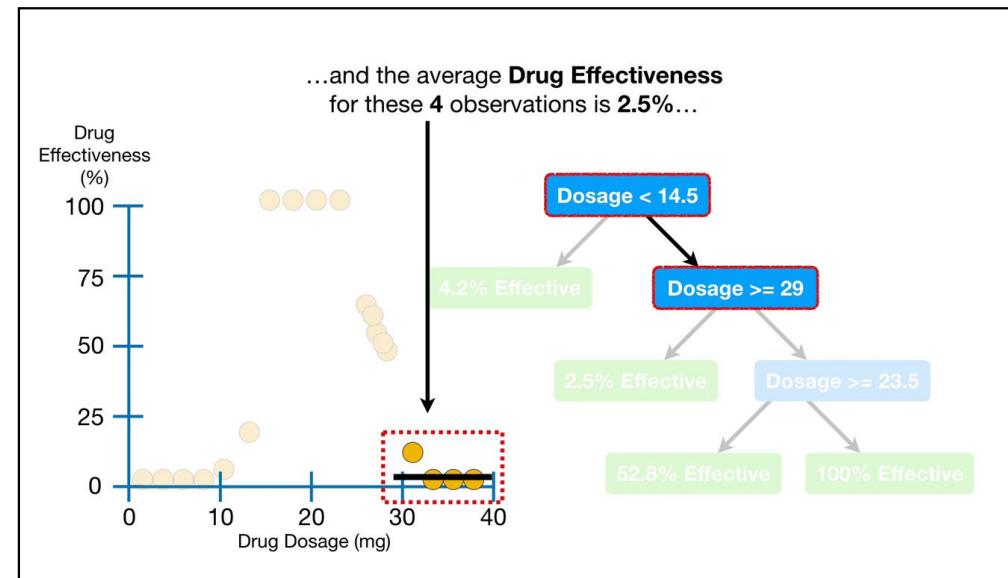
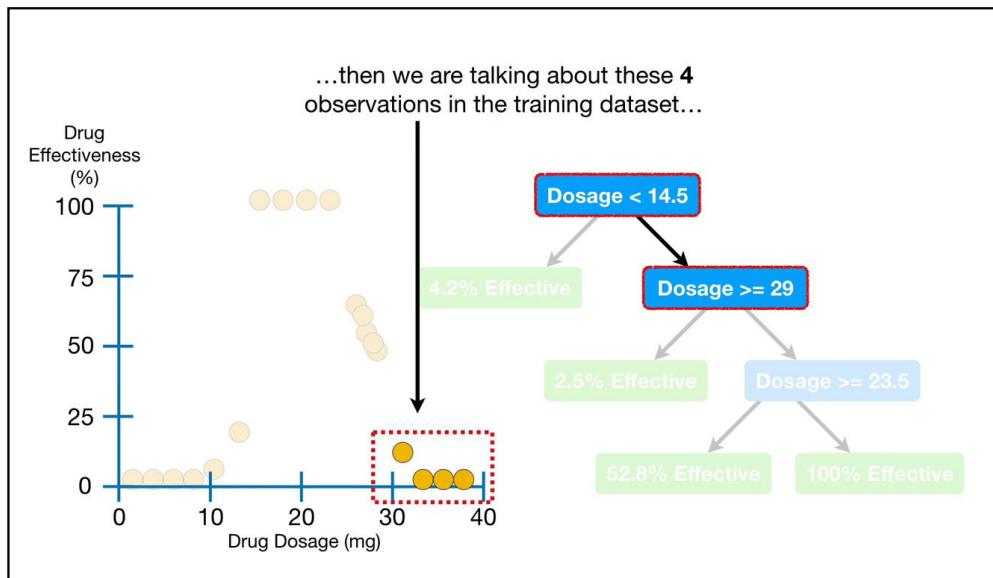
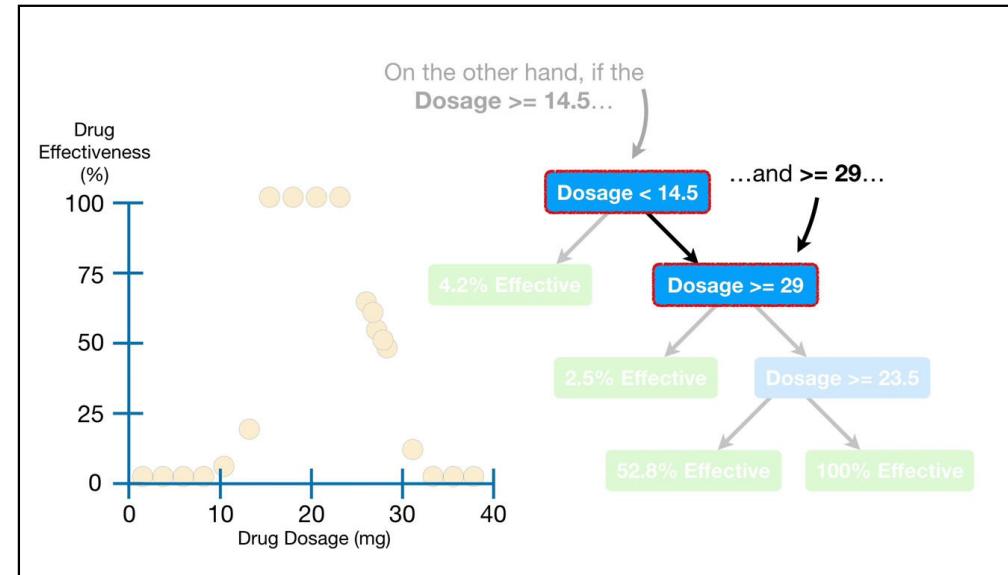
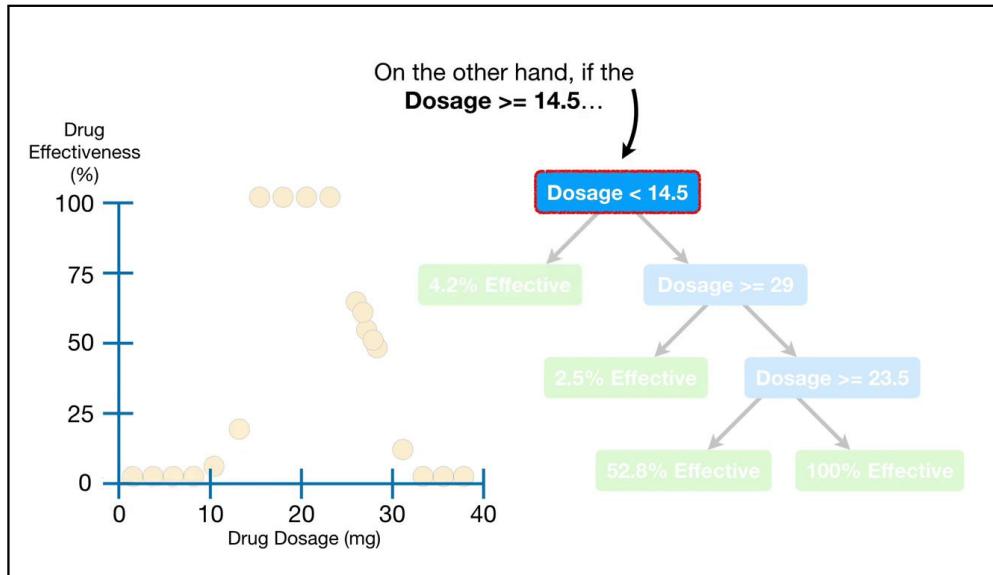


...and the average **Drug Effectiveness** for these 6 observations is **4.2%**...

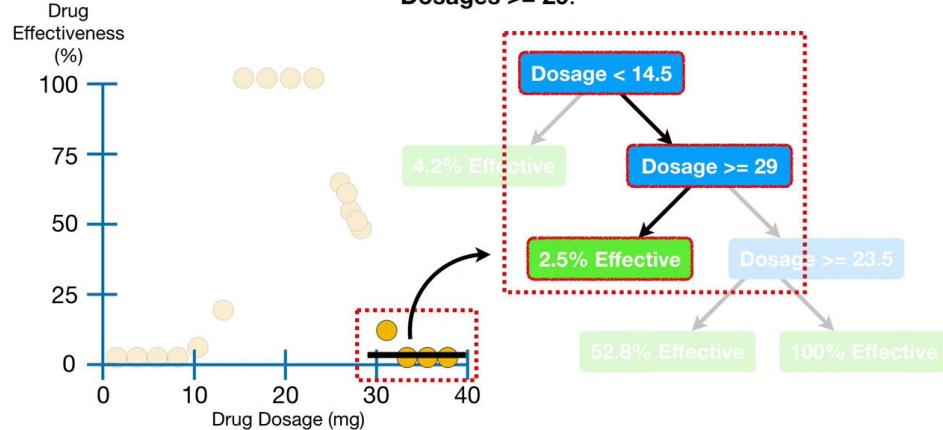


...so the tree uses the average value, **4.2%**, as its prediction for people with **Dosages < 14.5**.

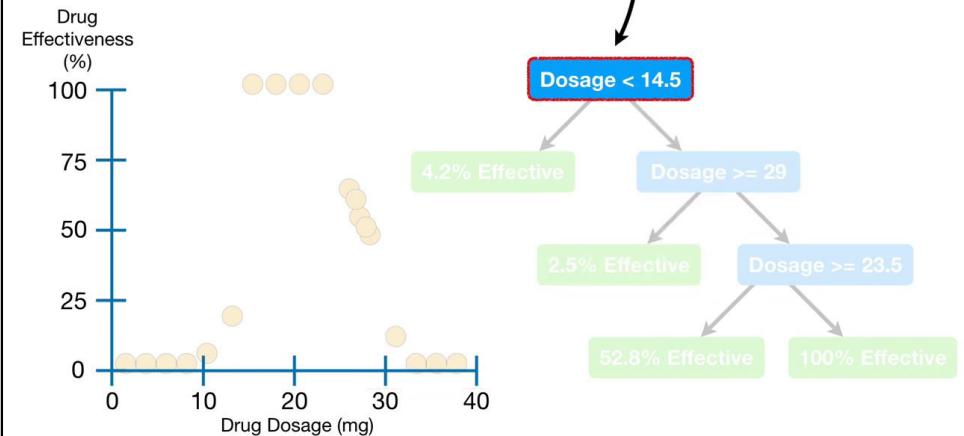




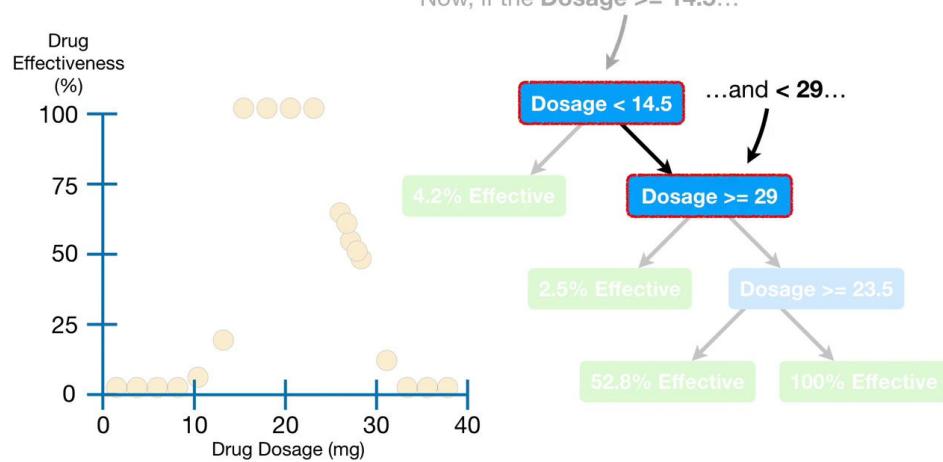
...so the tree uses the average value, 2.5%, as its prediction for people with Dosages  $\geq 29$ .



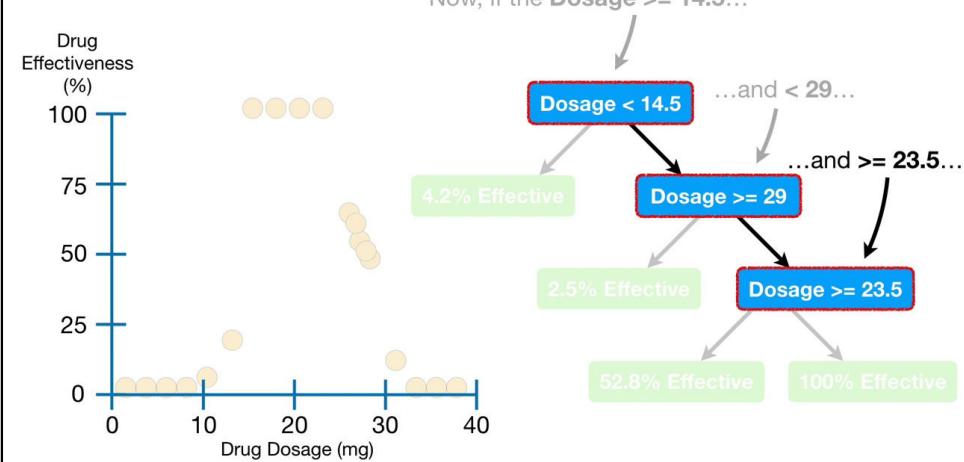
Now, if the Dosage  $\geq 14.5$ ...

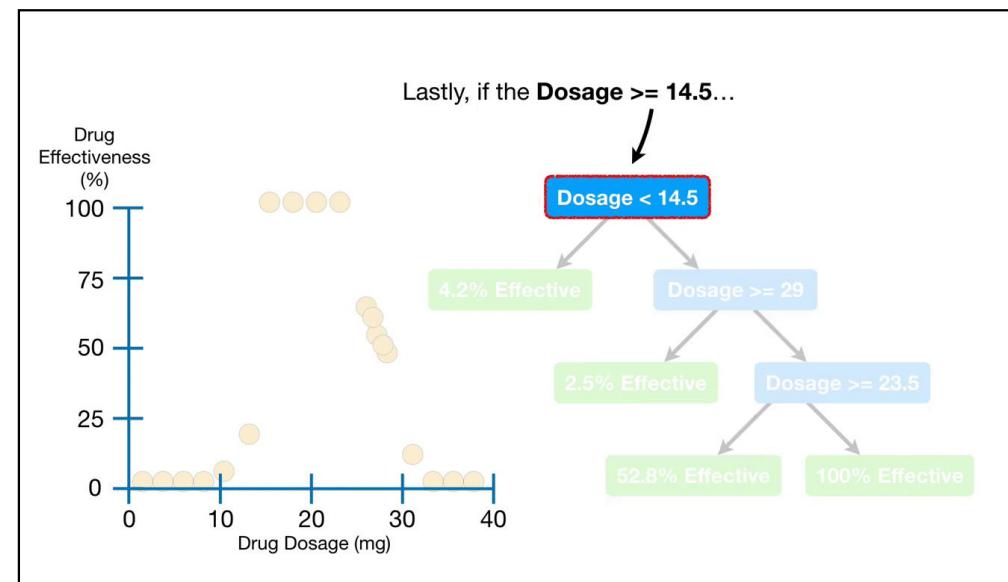
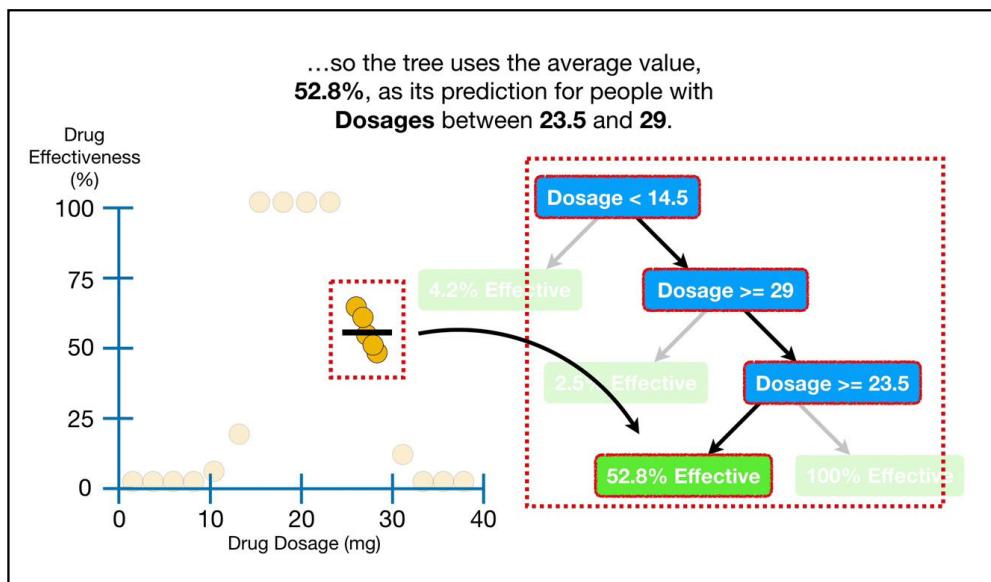
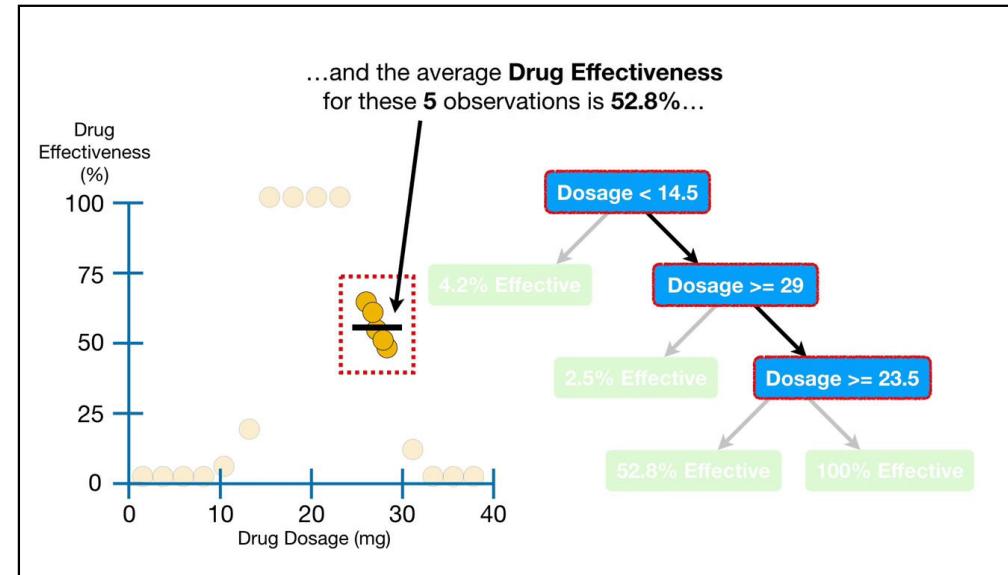
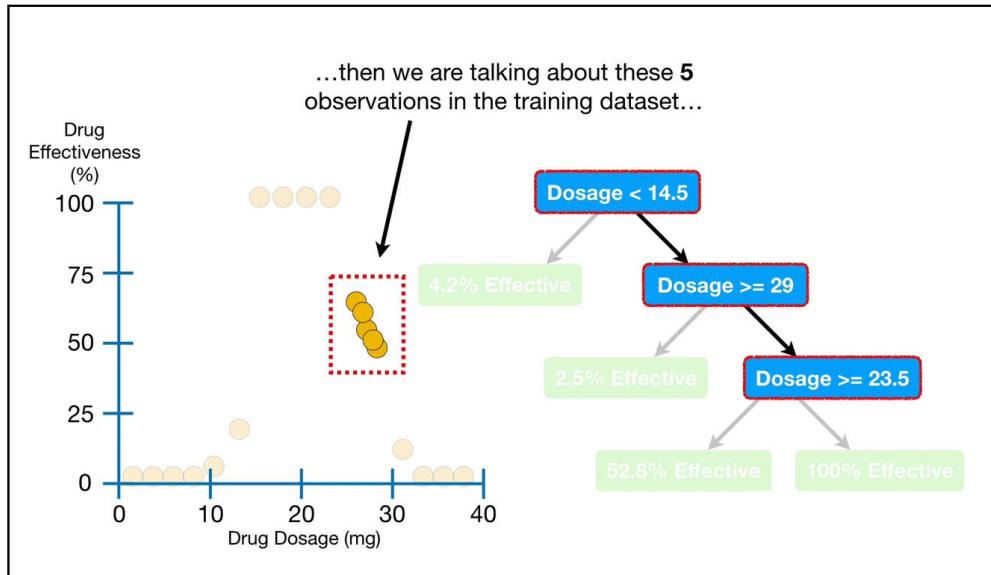


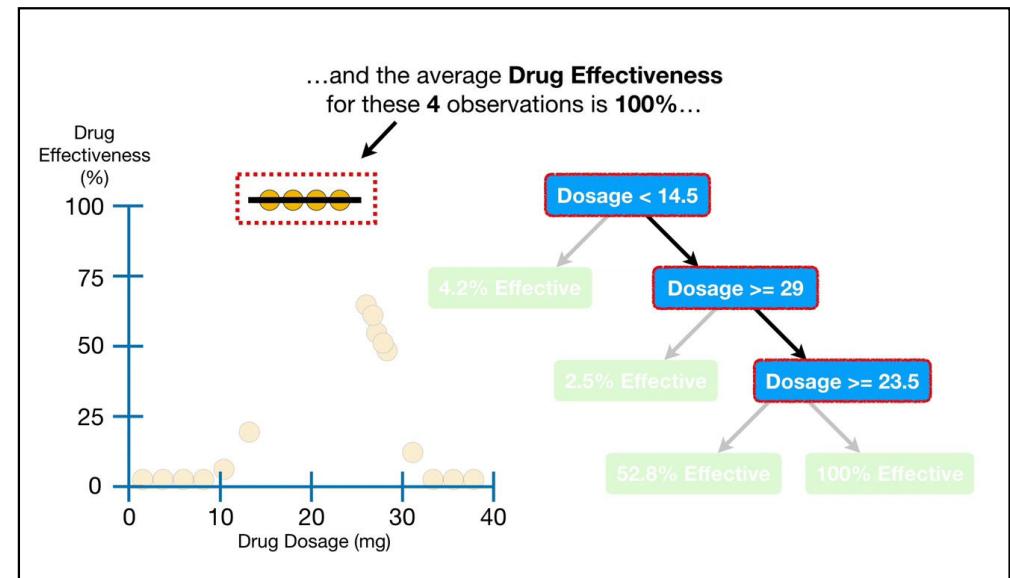
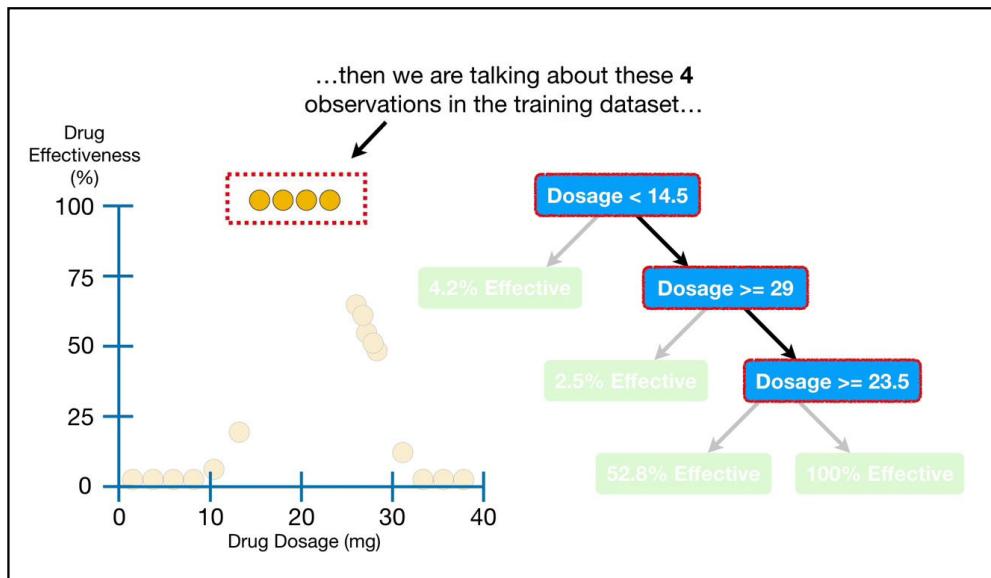
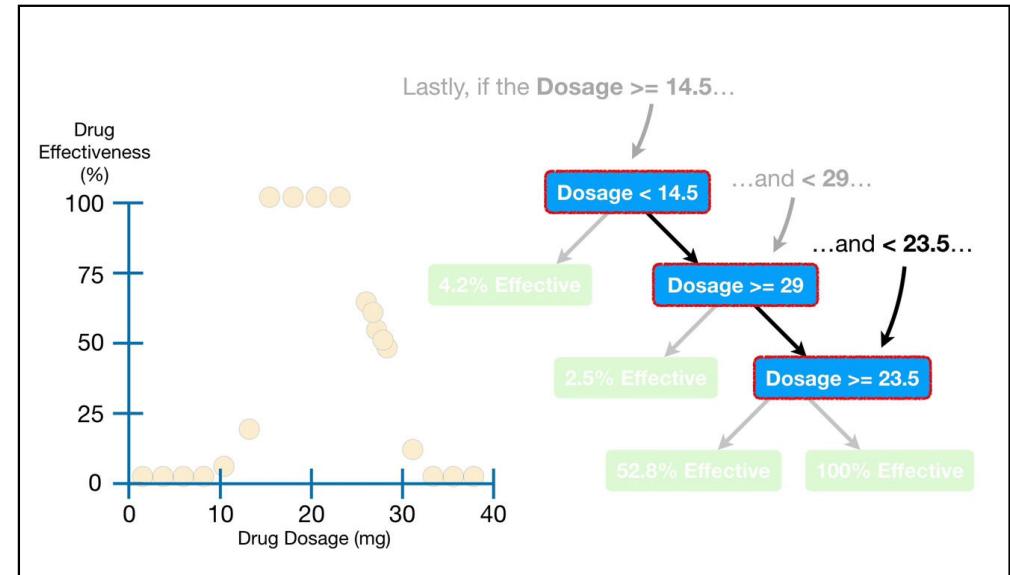
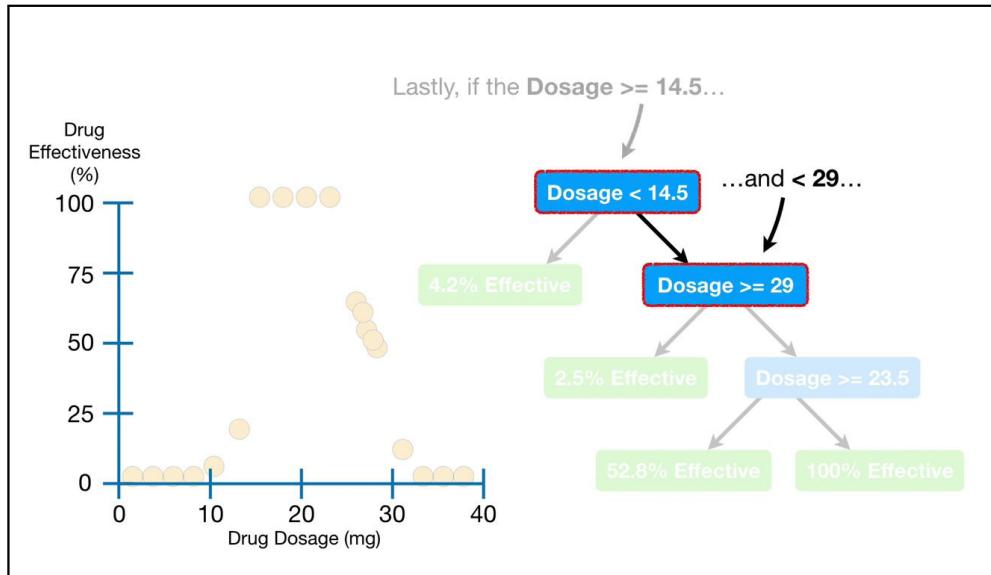
Now, if the Dosage  $\geq 14.5$ ...



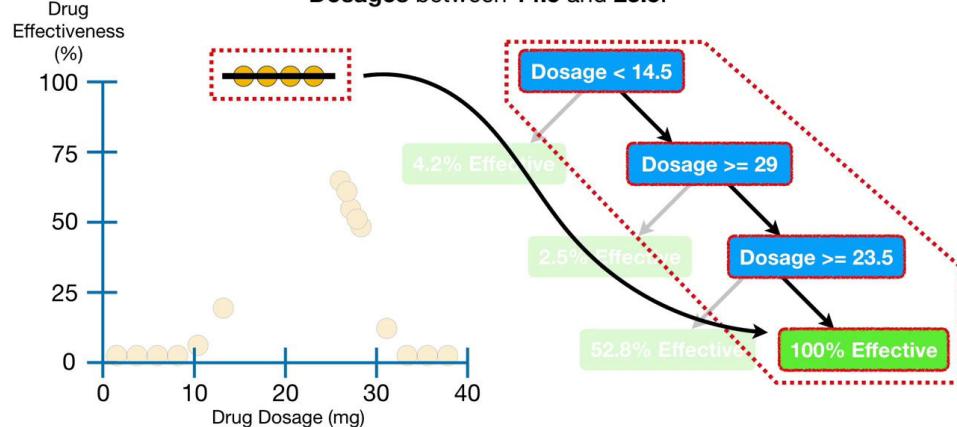
Now, if the Dosage  $\geq 14.5$ ...



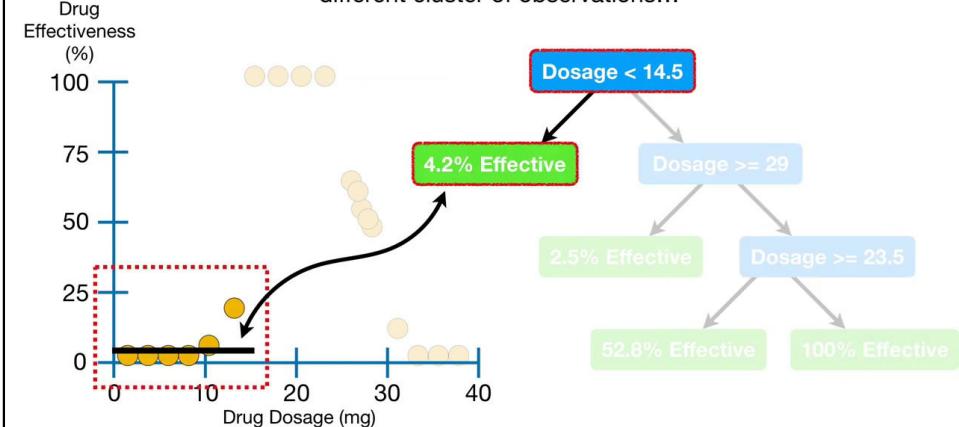




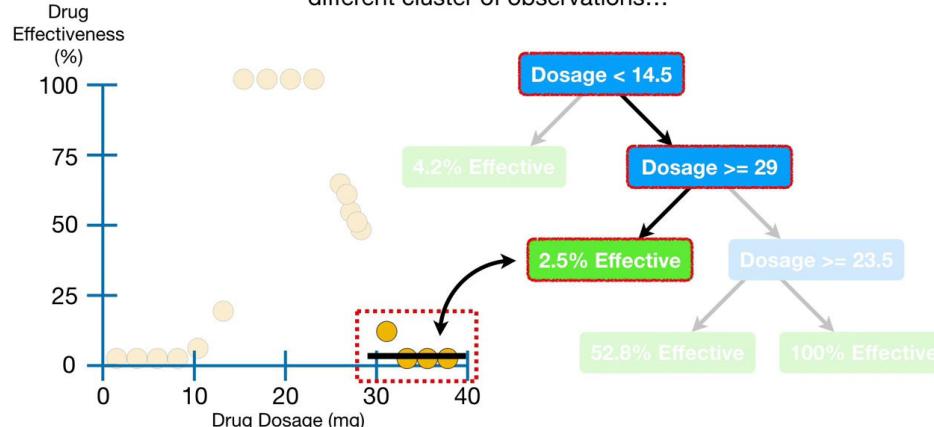
...so the tree uses the average value, 100%, as its prediction for people with **Dosages** between 14.5 and 23.5.



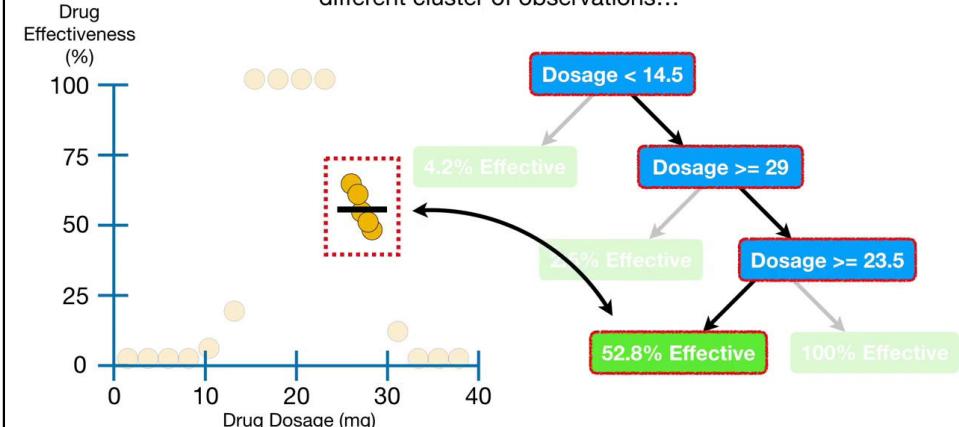
Since each leaf corresponds to the average **Drug Effectiveness** in a different cluster of observations...

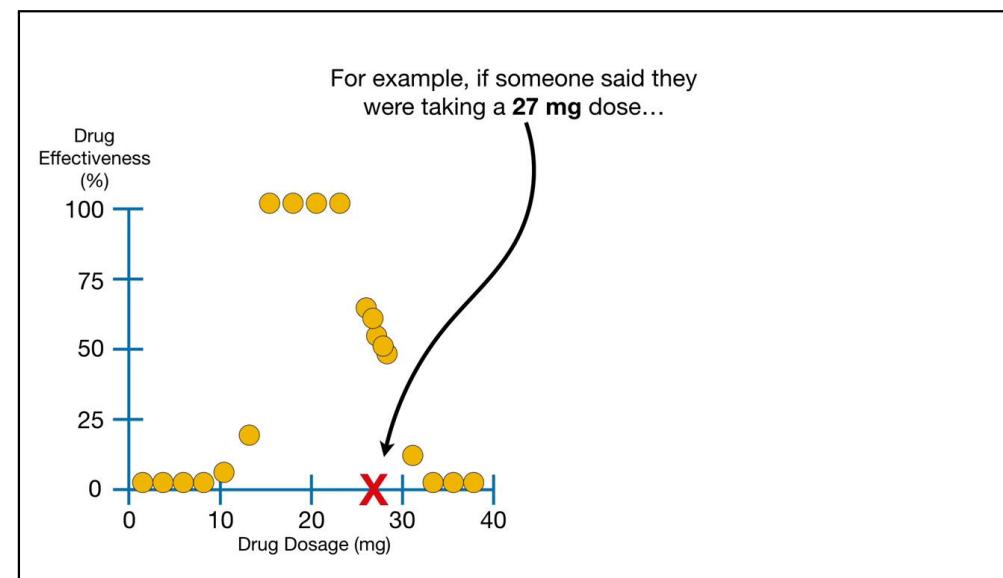
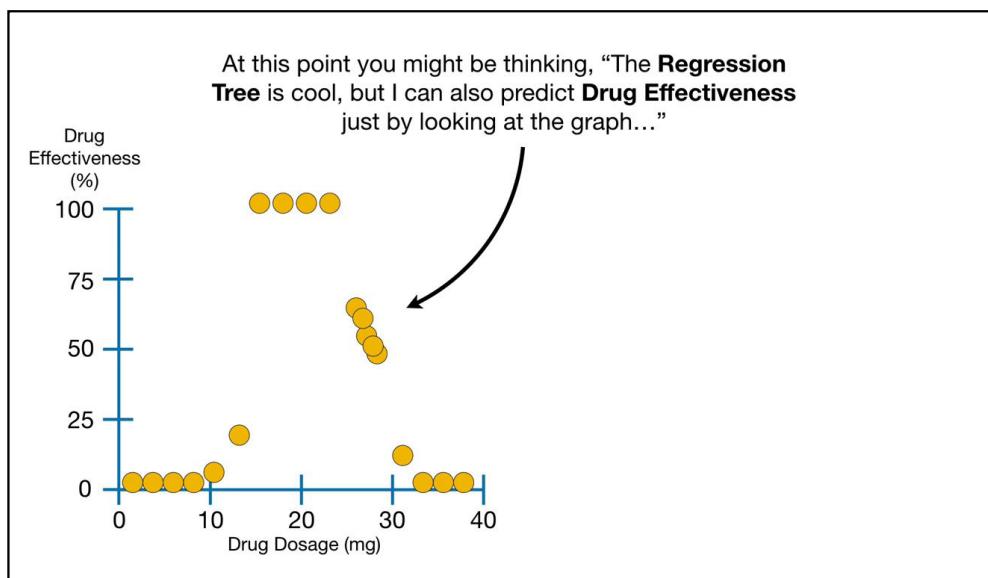
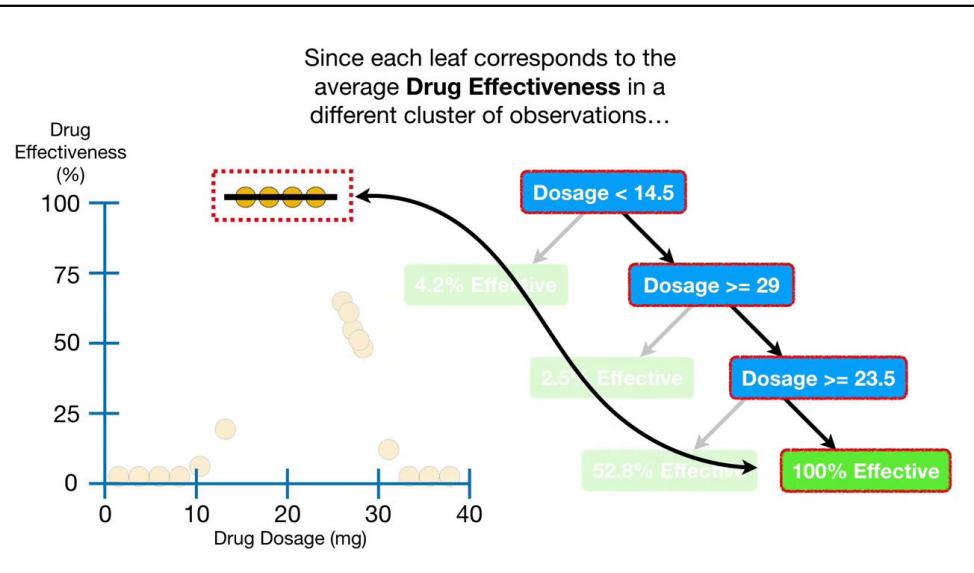


Since each leaf corresponds to the average **Drug Effectiveness** in a different cluster of observations...

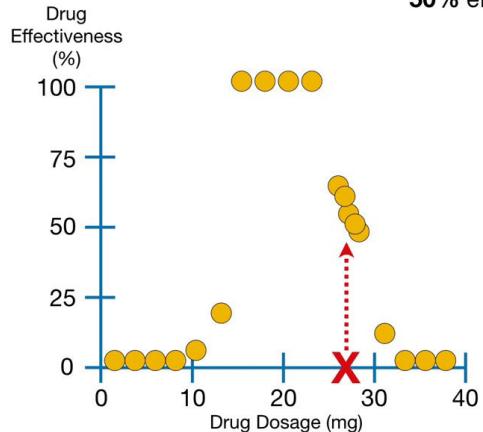


Since each leaf corresponds to the average **Drug Effectiveness** in a different cluster of observations...

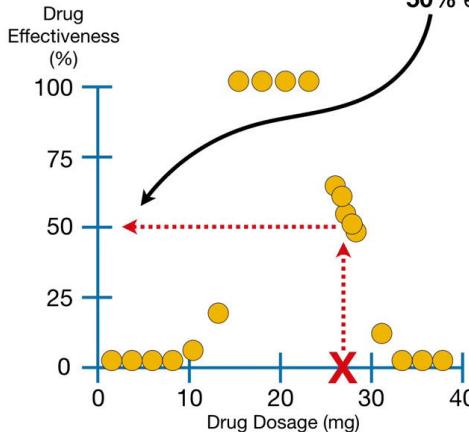




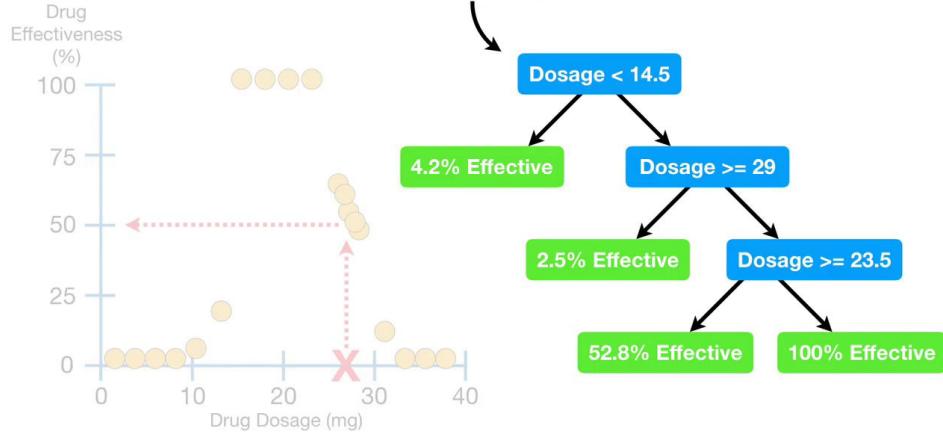
...then, just by looking at the graph,  
I can tell that the drug will be about  
**50% effective.**



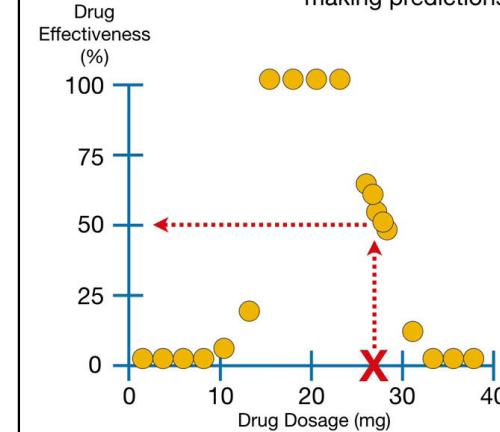
...then, just by looking at the graph,  
I can tell that the drug will be about  
**50% effective.**



So why make a big deal about the  
**Regression Tree?**



When the data are super simple and we are only using  
one predictor, **Dosage**, to predict **Drug Effectiveness**,  
making predictions by eye isn't terrible.

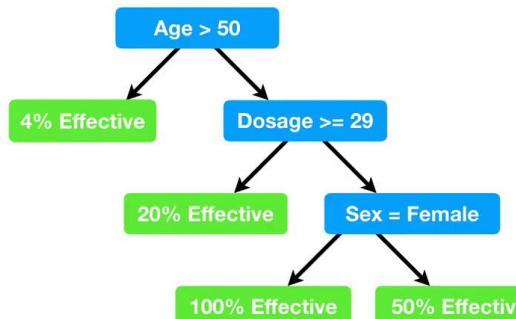


Dosage	Drug Effect.
10	58
20	60
35	57
5	44
etc...	etc...

But when we have **3** or more predictors, like **Dosage**, **Age** and **Sex**, to predict **Drug Effectiveness**, drawing a graph is very difficult, if not impossible.

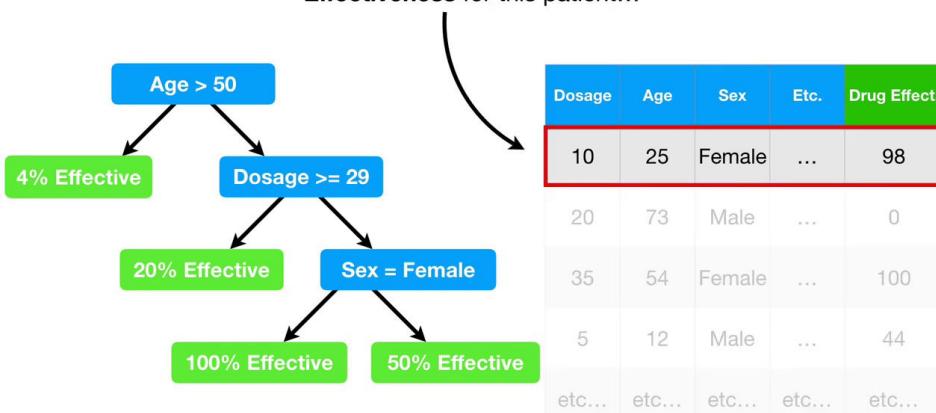
Dosage	Age	Sex	Etc.	Drug Effect.
10	25	Female	...	98
20	73	Male	...	0
35	54	Female	...	100
5	12	Male	...	44
etc...	etc...	etc...	etc...	etc...

In contrast, a **Regression Tree** easily accommodates the additional predictors.

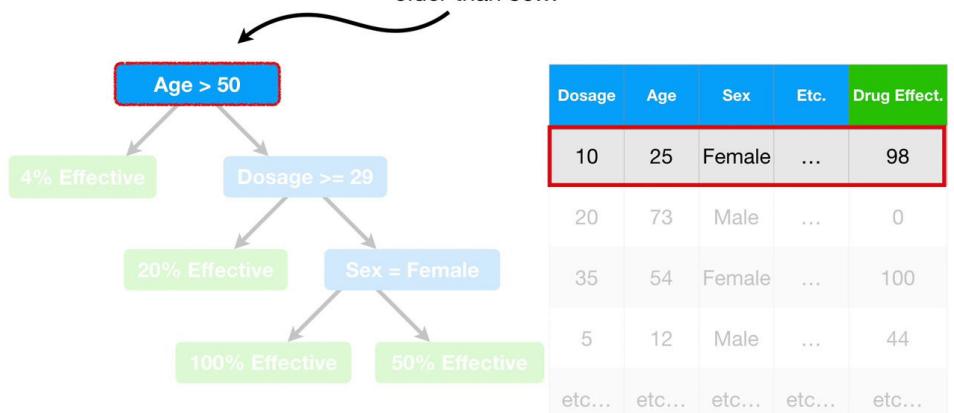


Dosage	Age	Sex	Etc.	Drug Effect.
10	25	Female	...	98
20	73	Male	...	0
35	54	Female	...	100
5	12	Male	...	44
etc...	etc...	etc...	etc...	etc...

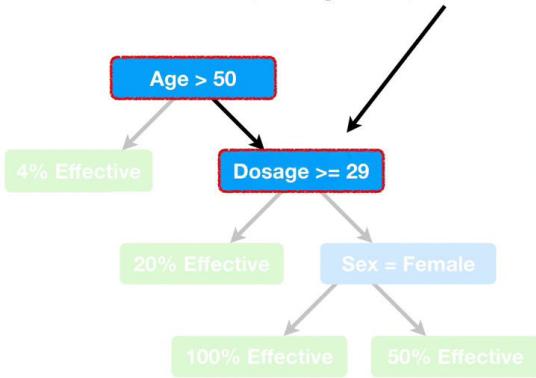
For example, if we wanted to predict the **Drug Effectiveness** for this patient...



...we start by asking if they are older than 50...

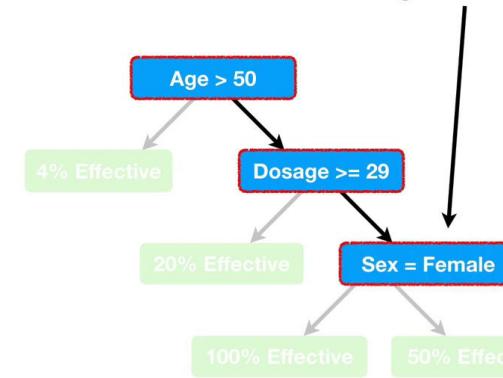


...and since they *not* over 50, we follow the branch on the *right* and ask if their **Dosage  $\geq 29$** ...



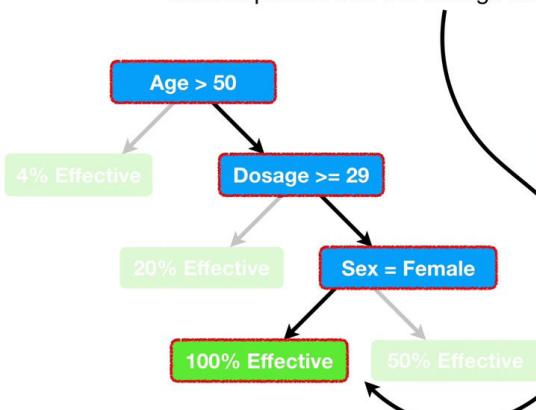
Dosage	Age	Sex	Etc.	Drug Effect.
10	25	Female	...	98
20	73	Male	...	0
35	54	Female	...	100
5	12	Male	...	44
etc...	etc...	etc...	etc...	etc...

...and since their dosage is *not*  $\geq 29$ , we follow the branch on the *right* and ask if they are **Female**...



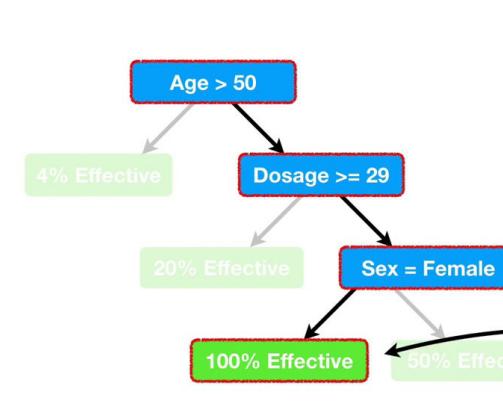
Dosage	Age	Sex	Etc.	Drug Effect.
10	25	Female	...	98
20	73	Male	...	0
35	54	Female	...	100
5	12	Male	...	44
etc...	etc...	etc...	etc...	etc...

...and since they *are* **Female**, we follow the branch on the *left* and predict that the dosage will be **100% Effective**...



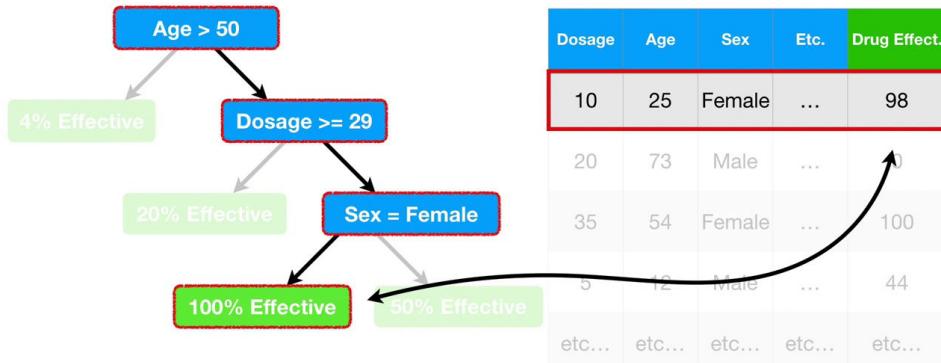
Dosage	Age	Sex	Etc.	Drug Effect.
10	25	Female	...	98
20	73	Male	...	0
35	54	Female	...	100
5	12	Male	...	44
etc...	etc...	etc...	etc...	etc...

...and that's not too far off from the truth, **98%**.



Dosage	Age	Sex	Etc.	Drug Effect.
10	25	Female	...	98
20	73	Male	...	0
35	54	Female	...	100
5	12	Male	...	44
etc...	etc...	etc...	etc...	etc...

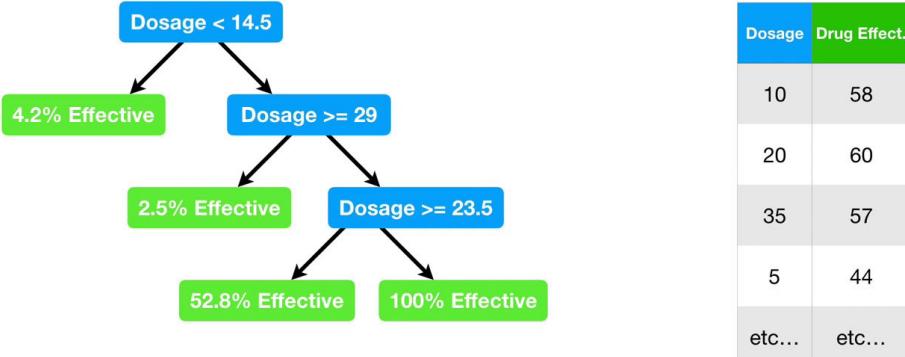
OK, now that we know that **Regression Trees** can easily handle complicated data...



...let's go back to the original data, with just one predictor, **Dosage**...

Dosage	Drug Effect.
10	58
20	60
35	57
5	44
etc...	etc...

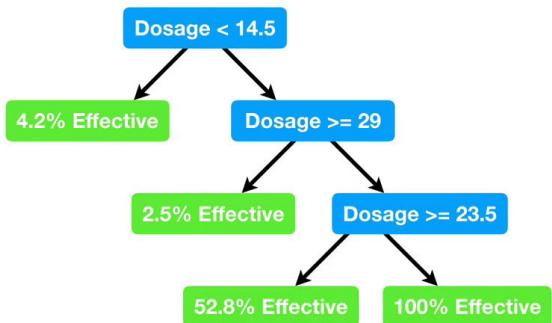
...and talk about how to build this **Regression Tree** from scratch...



...and since **Regression Trees** are built from the top down...

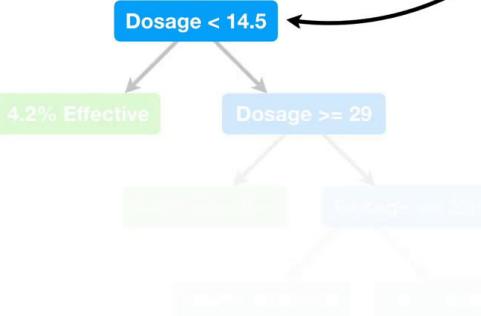
Dosage	Drug Effect.
10	58
20	60
35	57
5	44
etc...	etc...

...and since **Regression Trees** are built from the top down...

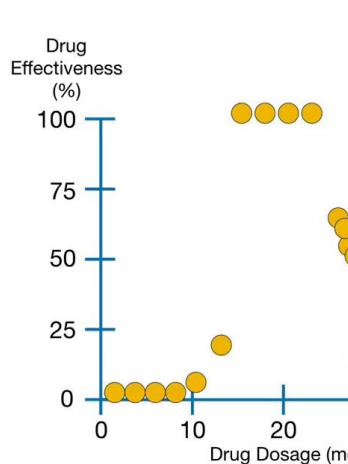


Dosage	Drug Effect.
10	58
20	60
35	57
5	44
etc...	etc...

...the first thing we do is figure out why we start by asking if **Dosage < 14.5**.

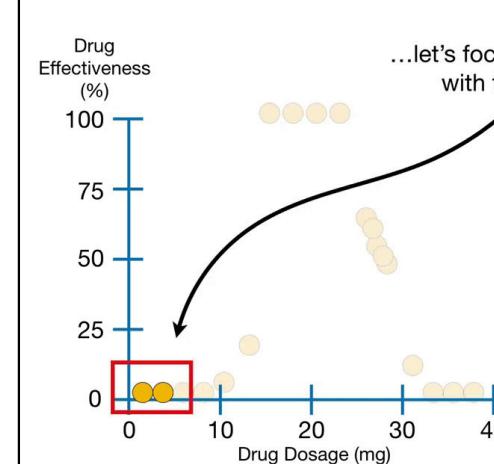


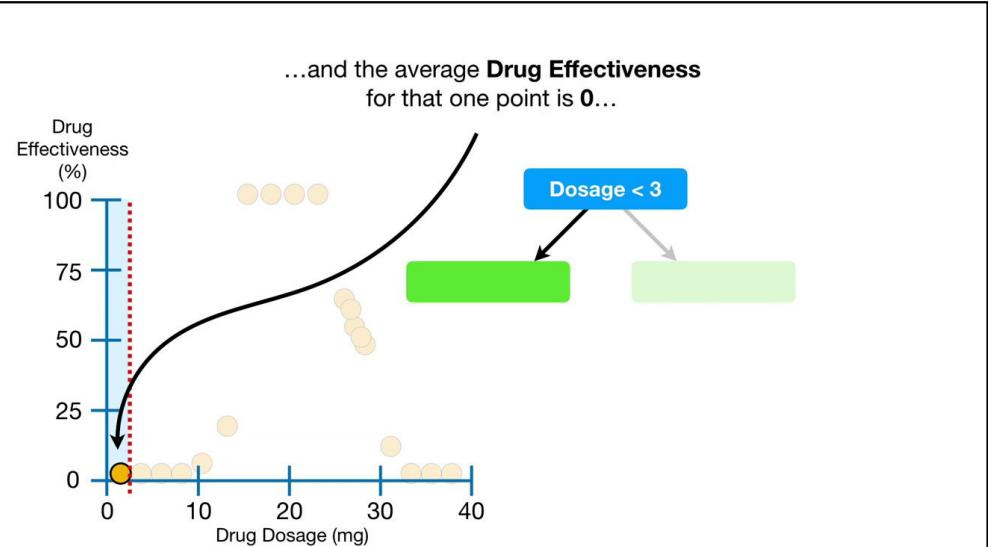
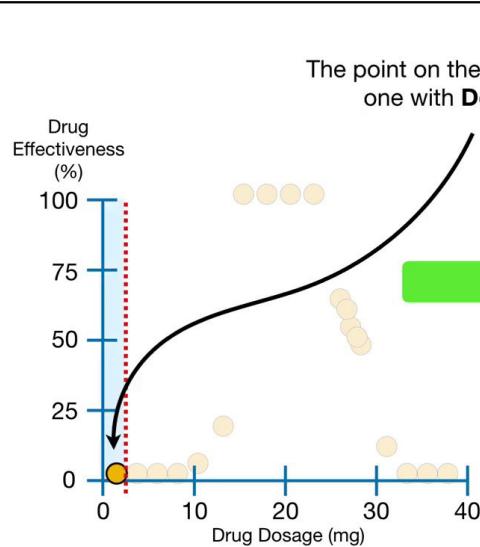
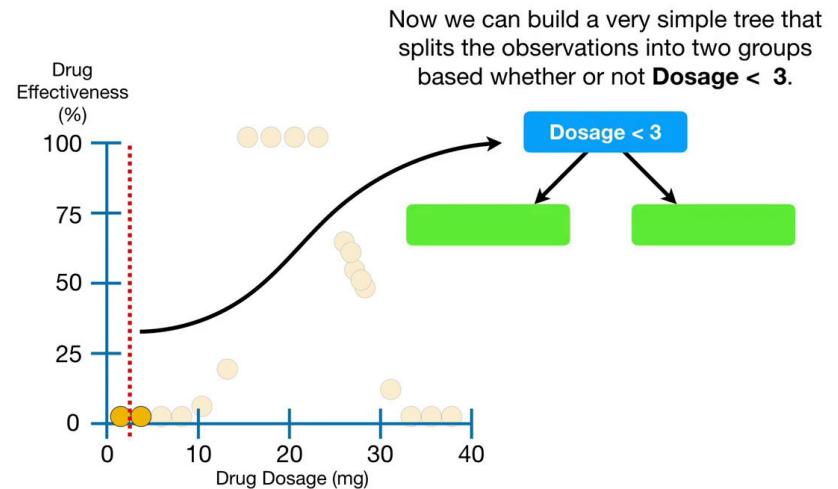
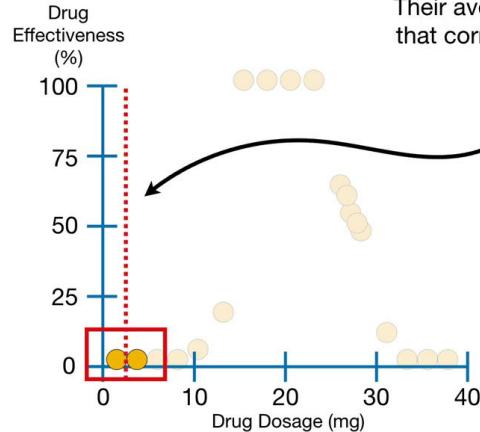
Dosage	Drug Effect.
10	58
20	60
35	57
5	44
etc...	etc...



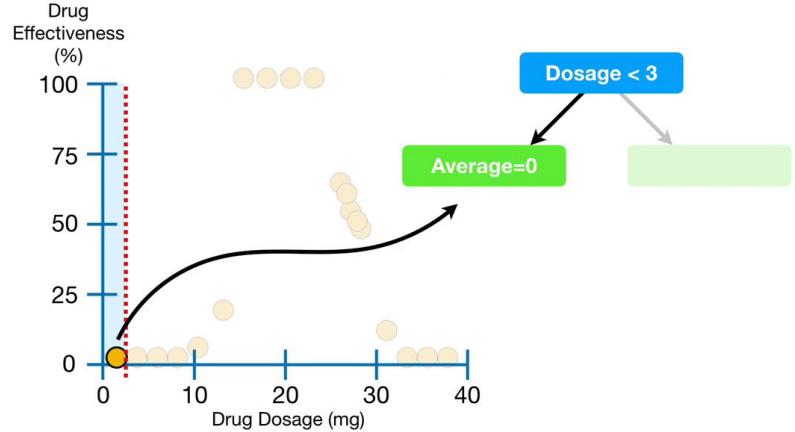
Dosage	Drug Effect.
10	58
20	60
35	57
5	44
etc...	etc...

...let's focus on the two observations with the smallest **Dosages**.

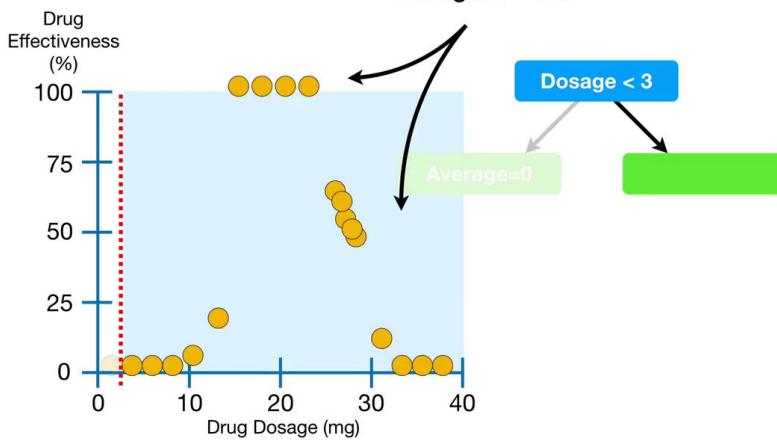




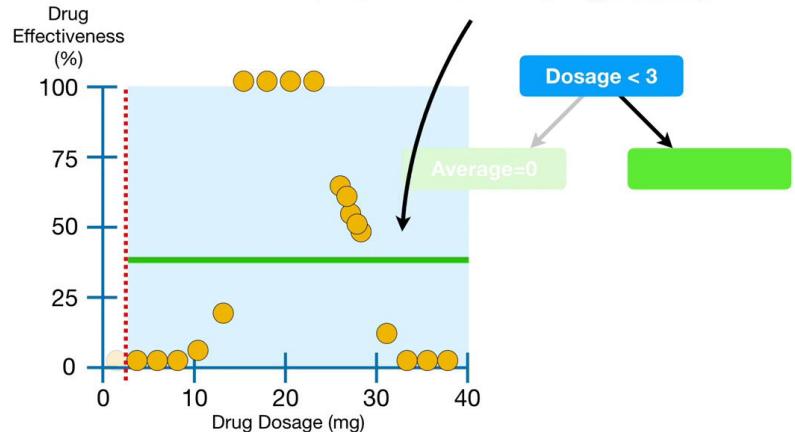
...so we put **0** in the leaf on the left side, for when **Dosage < 3**.



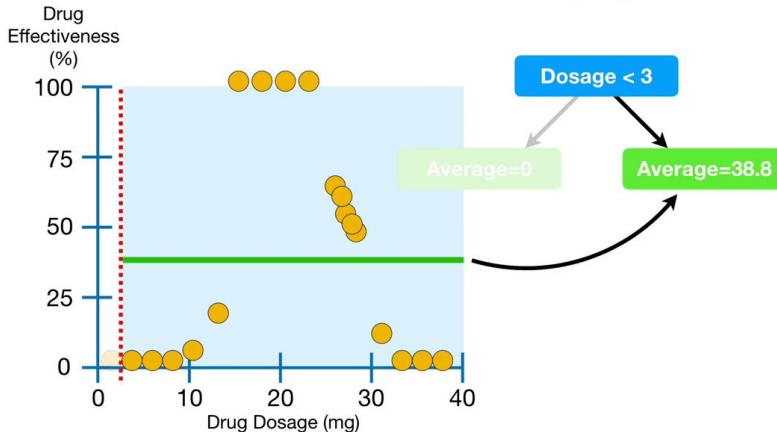
All of the other points have **Dosages  $\geq 3$** ...

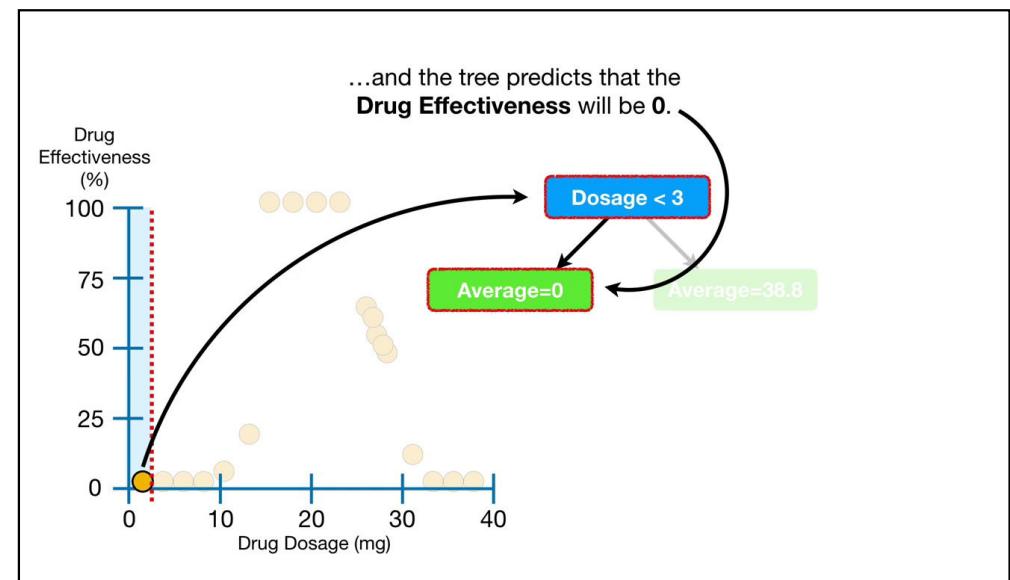
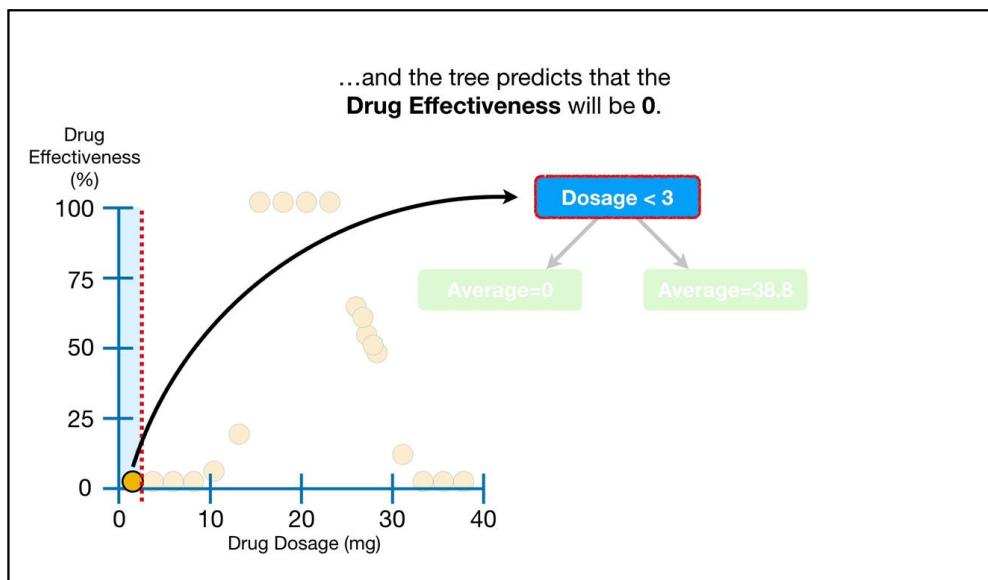
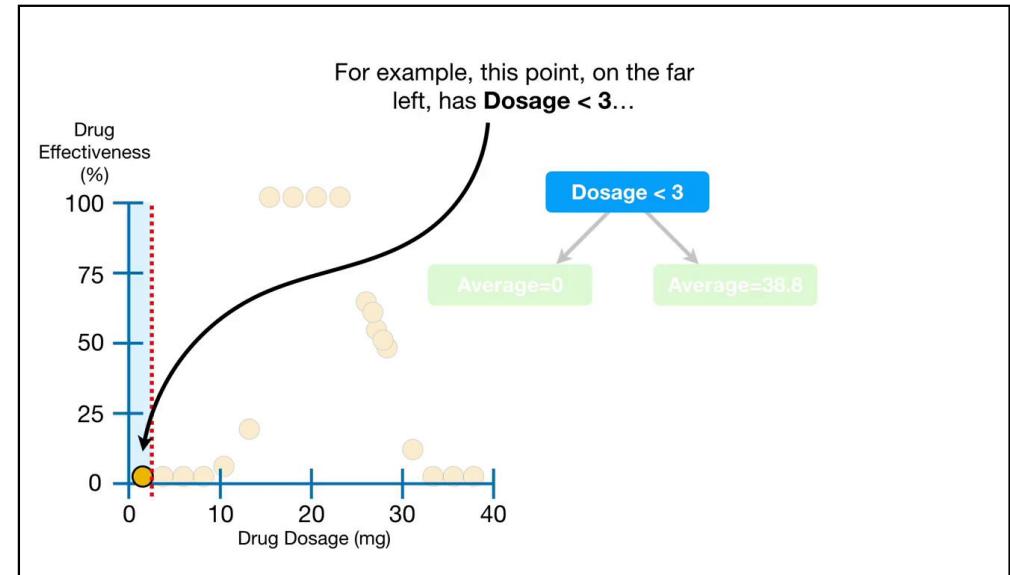
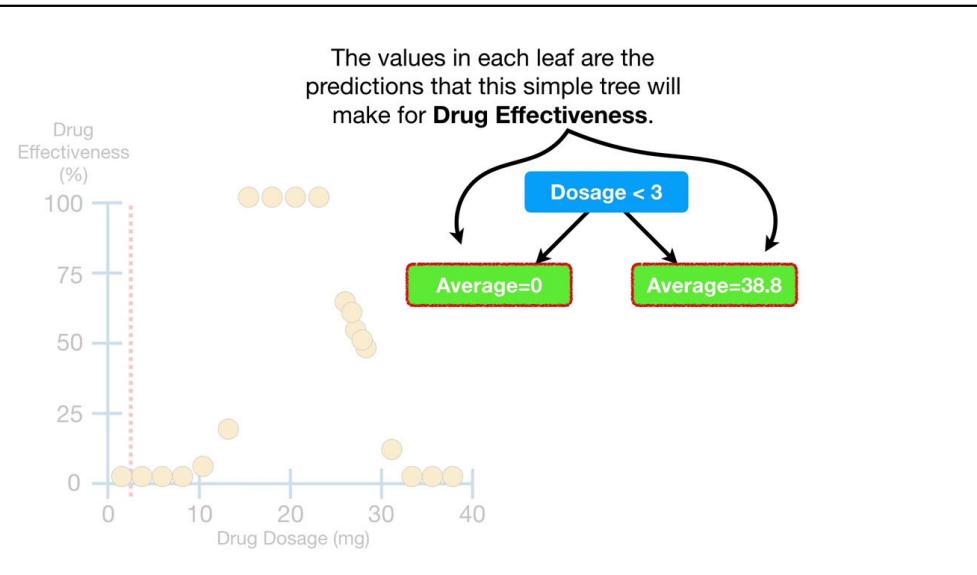


...and the average **Drug Effectiveness** for all of the points with **Dosages  $\geq 3$**  is **38.8**, (the **green line**)...

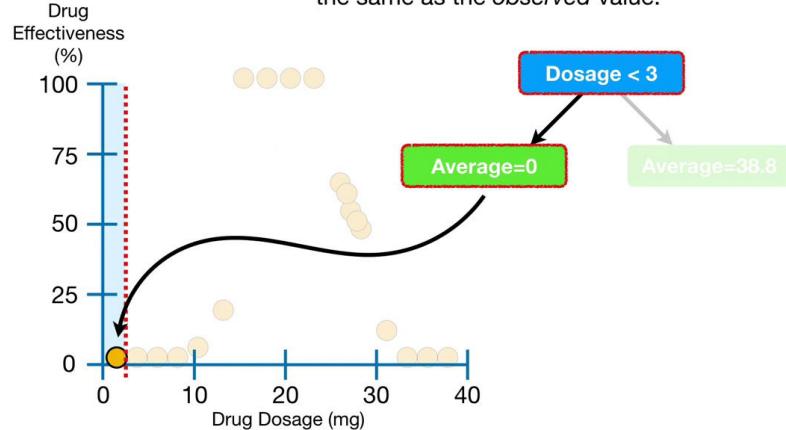


...so we put **38.8** in the leaf on the right side, for when the **Dosage  $\geq 3$** .

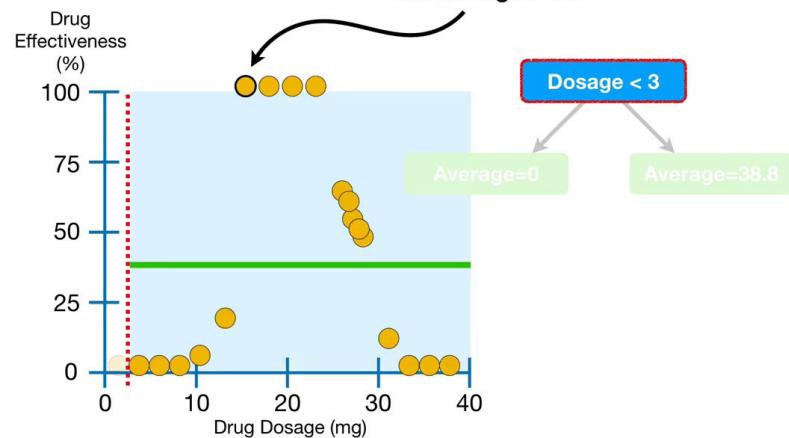




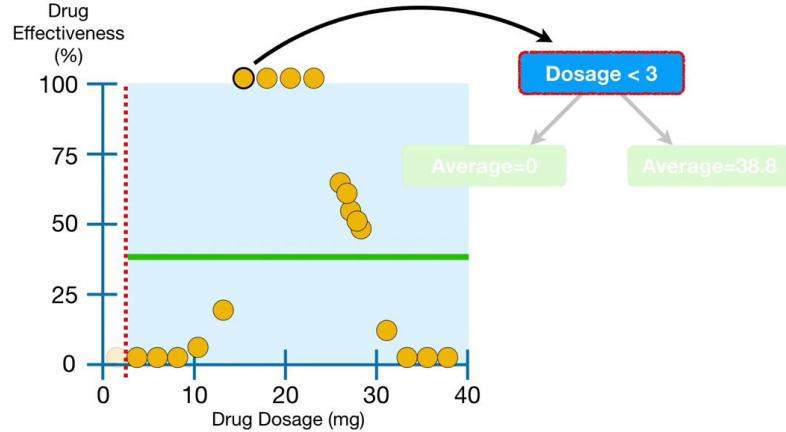
The prediction for this point, **Drug Effectiveness = 0**, is pretty good since it is the same as the observed value.



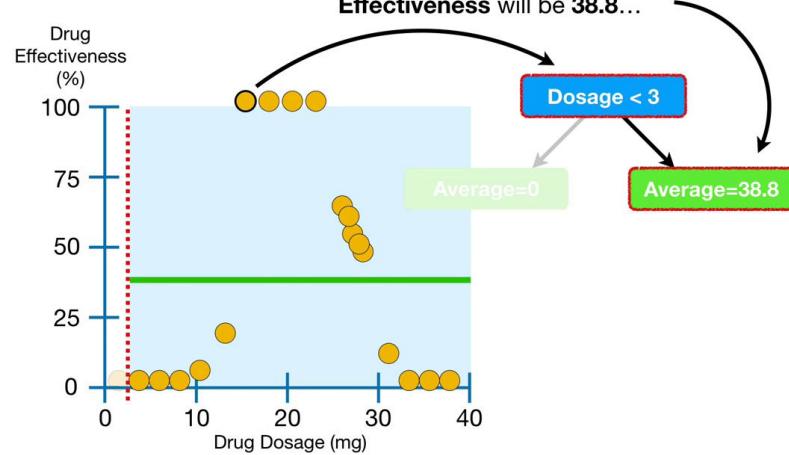
In contrast, for this point, which has **Dosage > 3...**

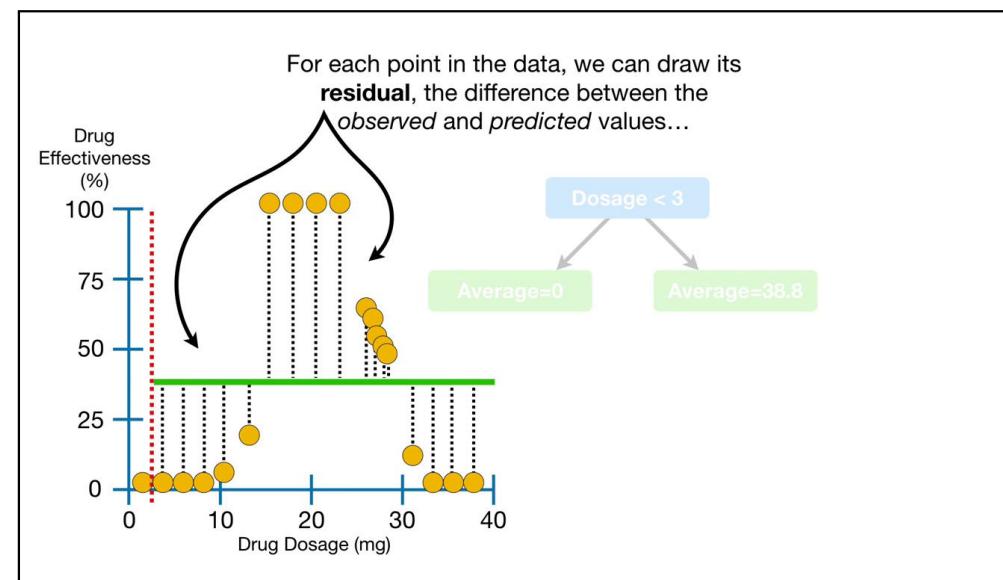
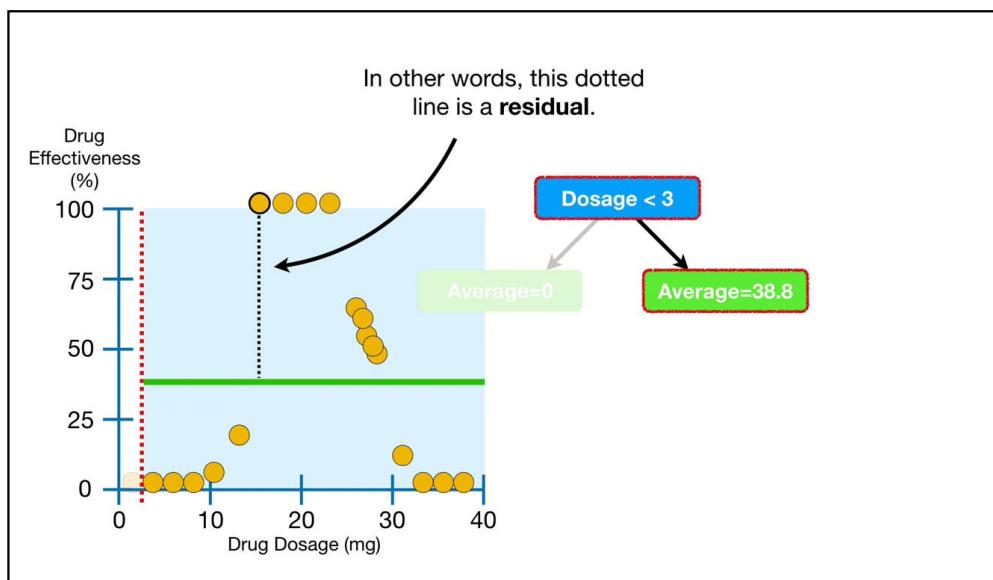
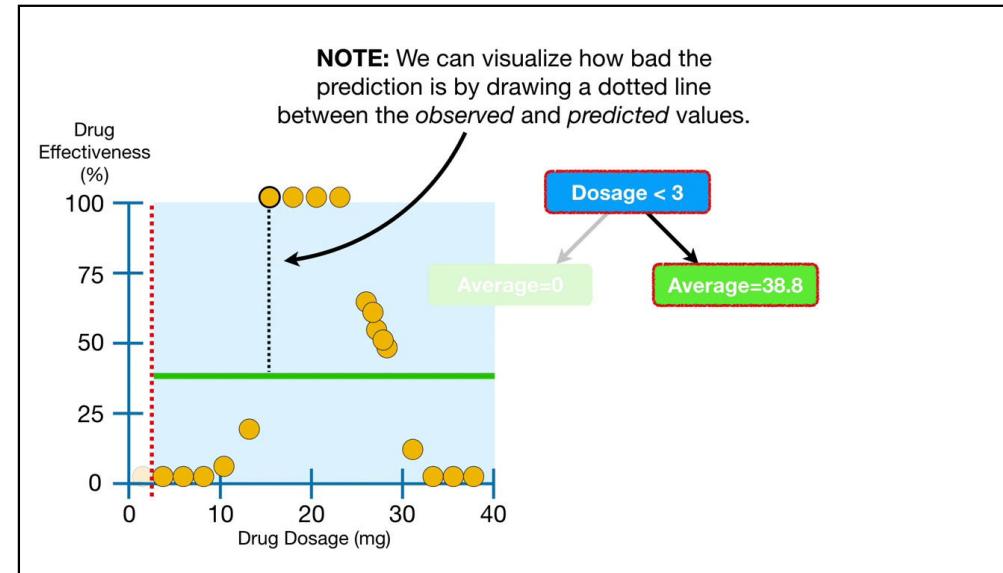
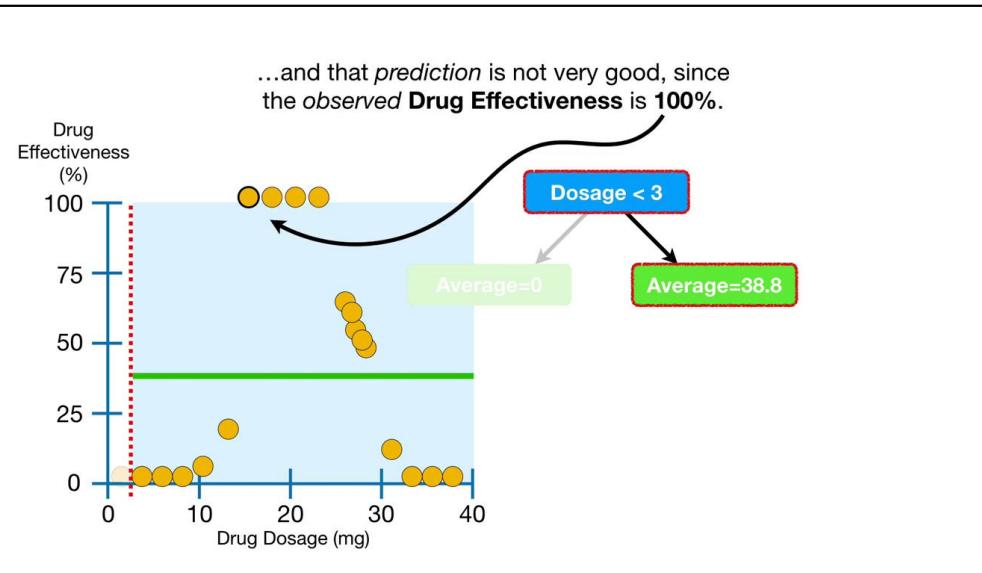


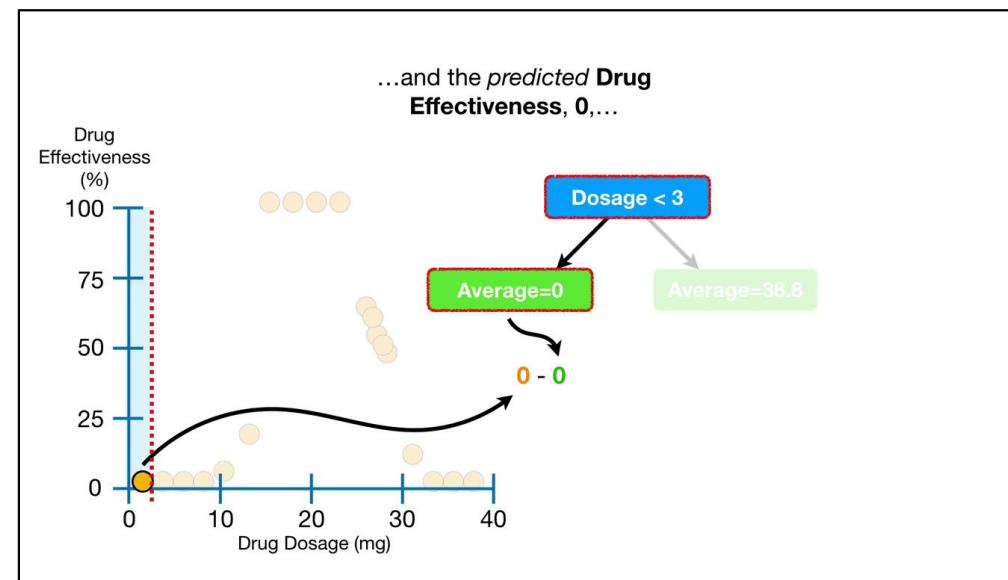
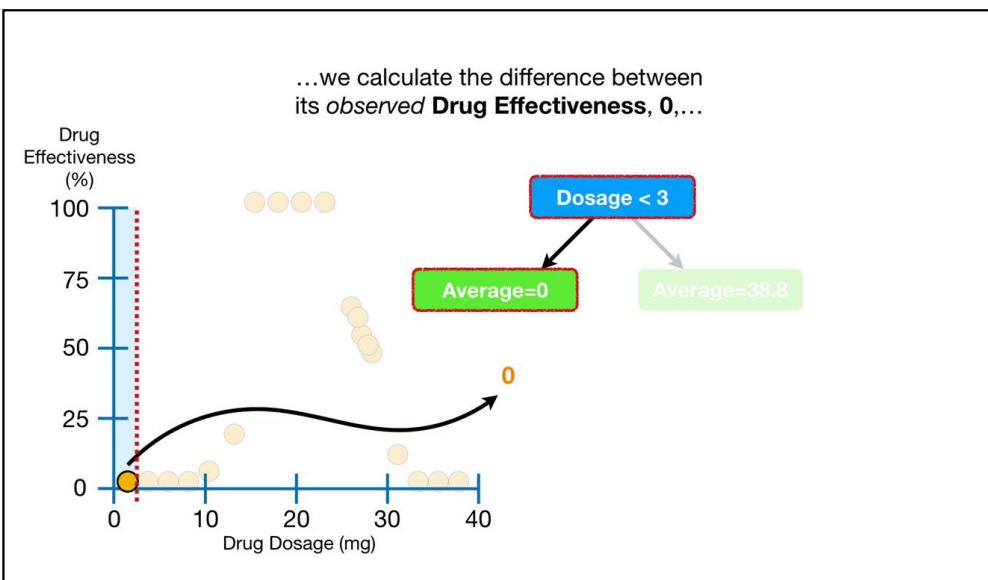
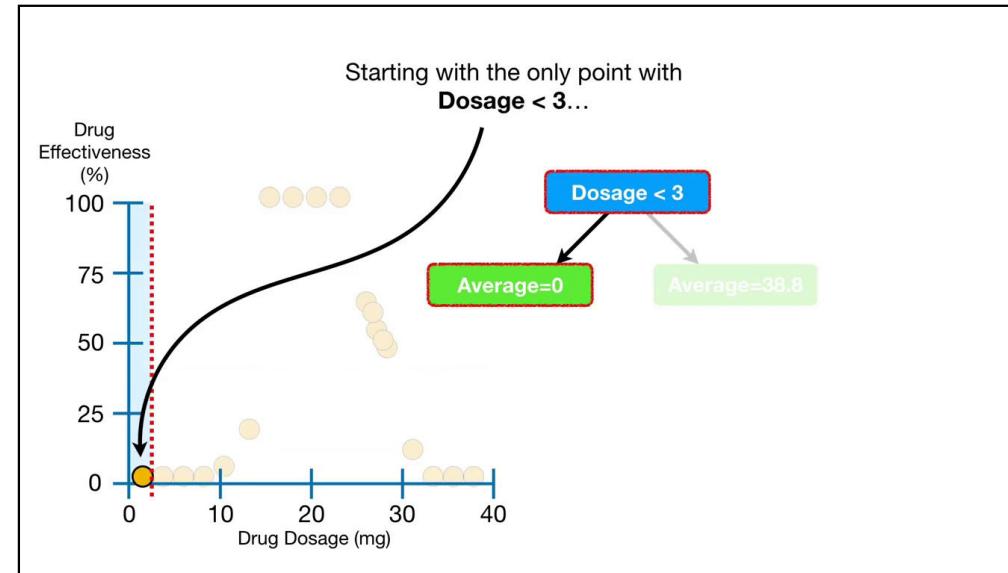
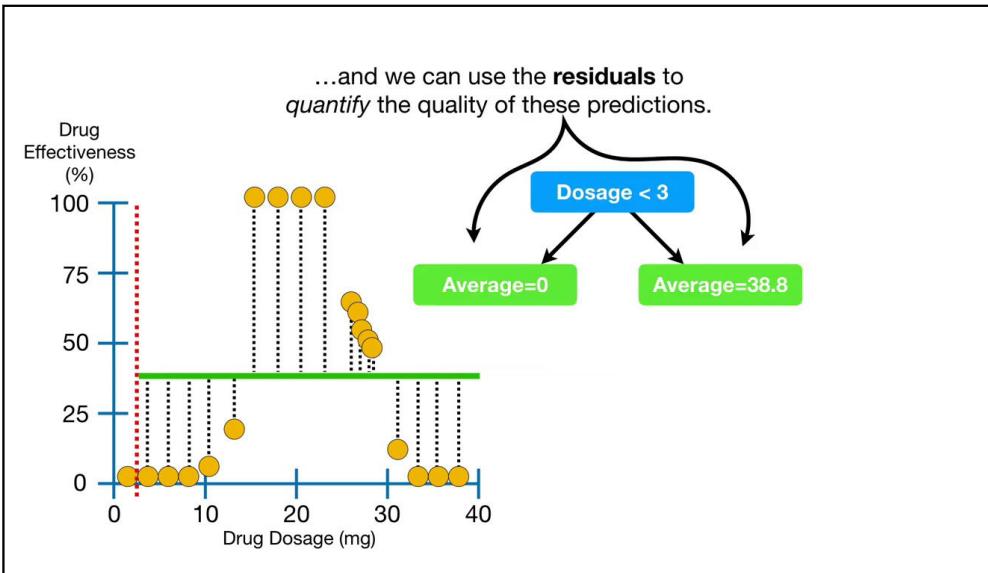
...the tree predicts that the **Drug Effectiveness** will be **38.8...**

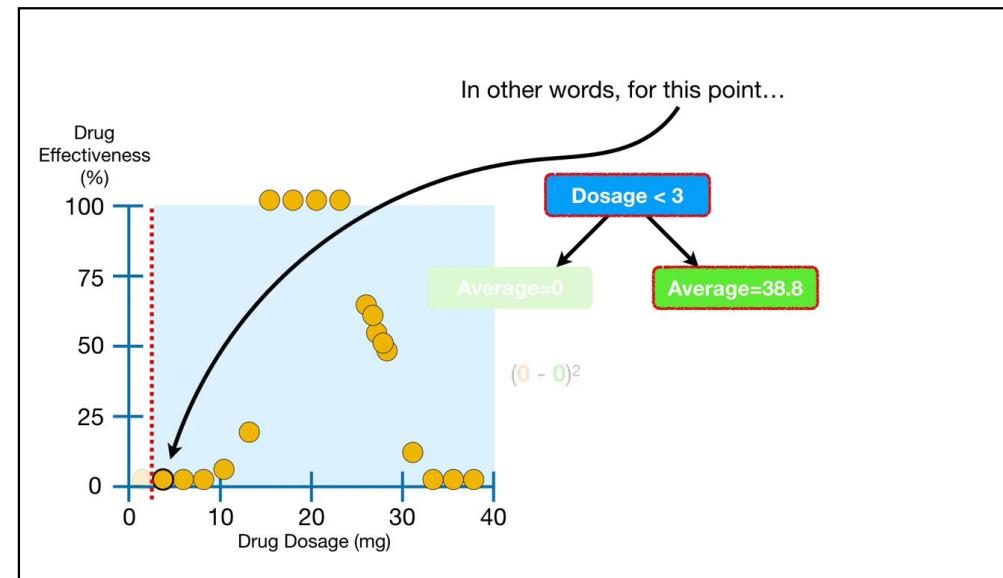
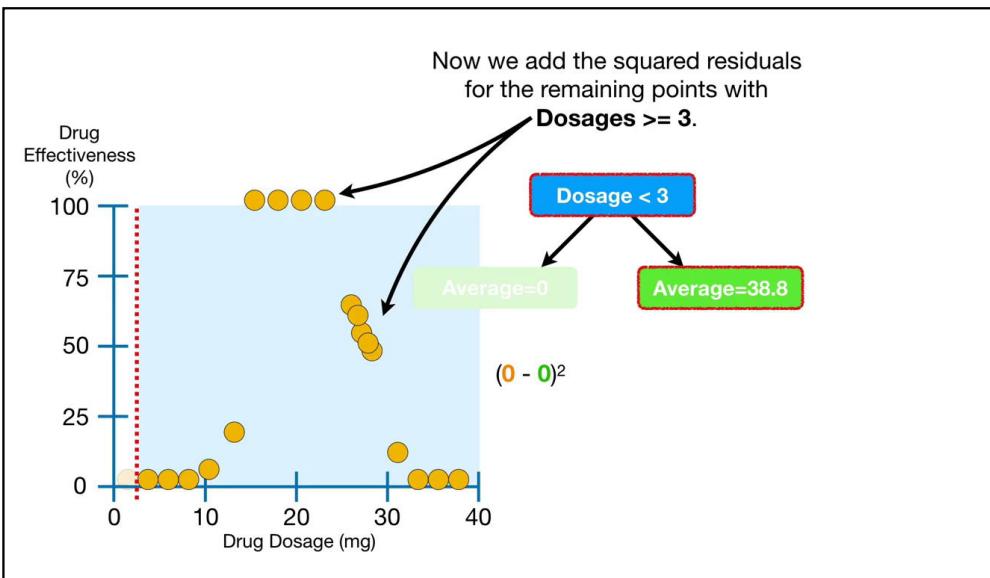
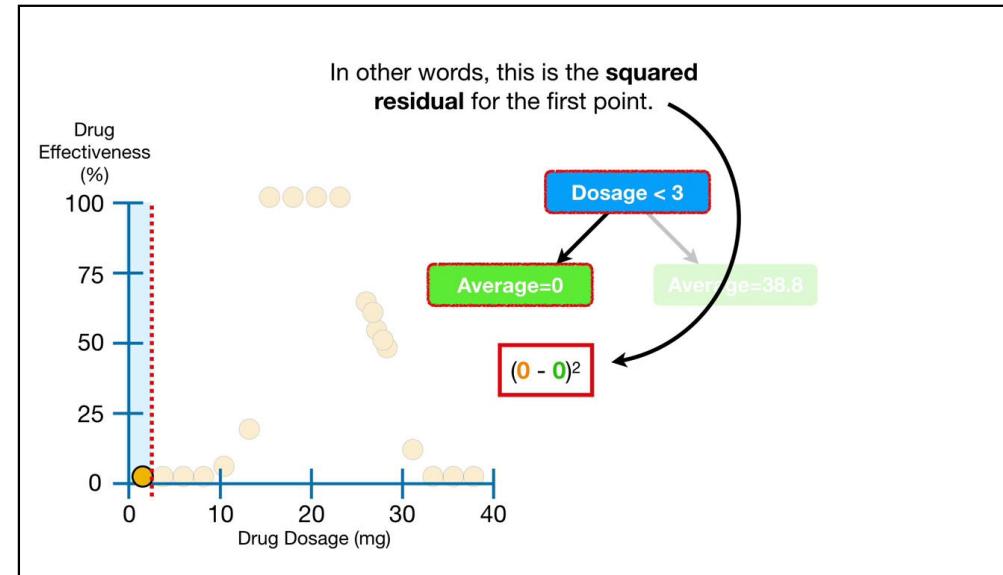
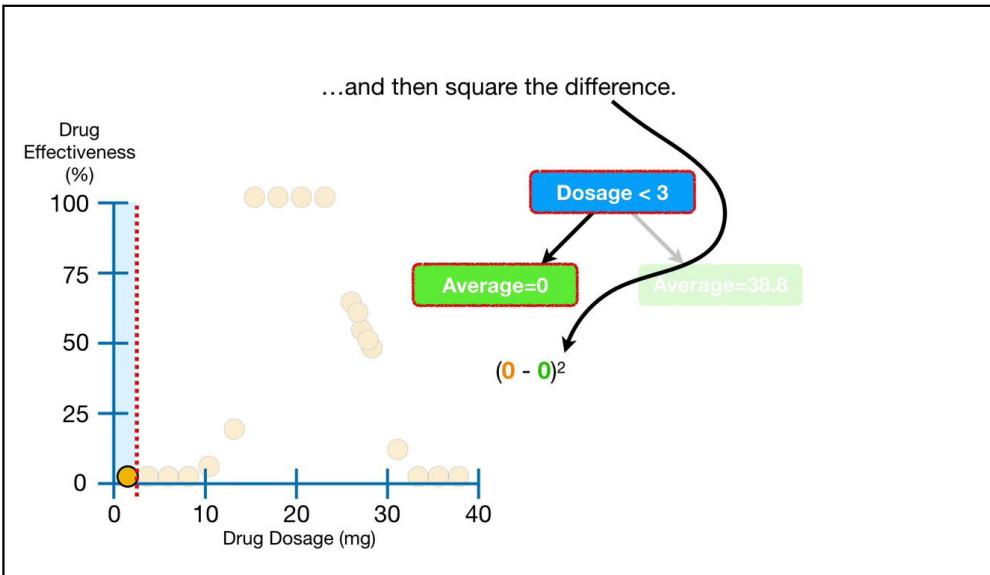


...the tree predicts that the **Drug Effectiveness** will be **38.8...**

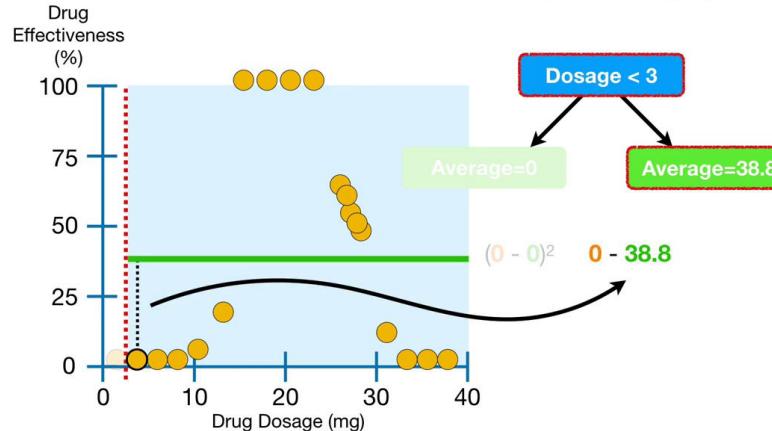




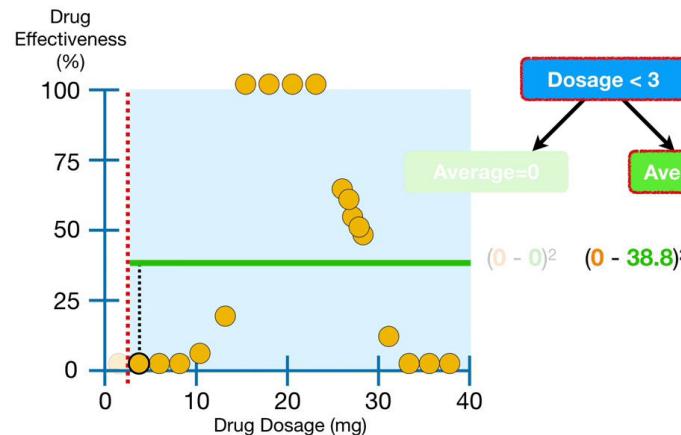




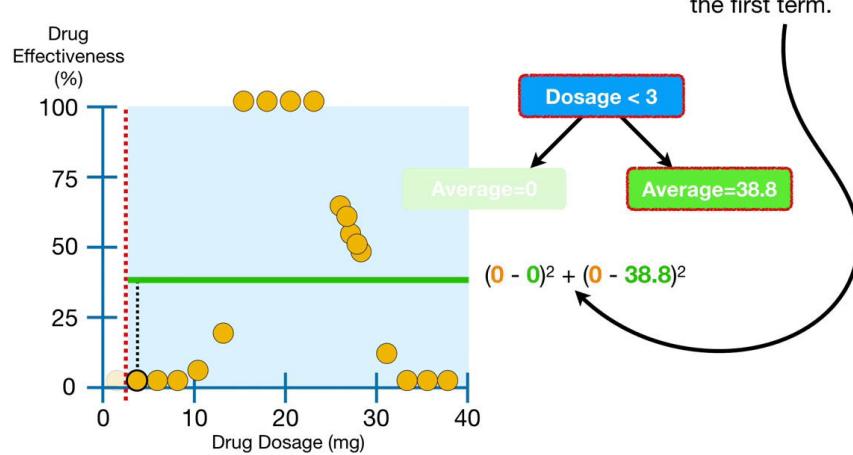
...we calculate the the difference between the observed and *predicted* values...



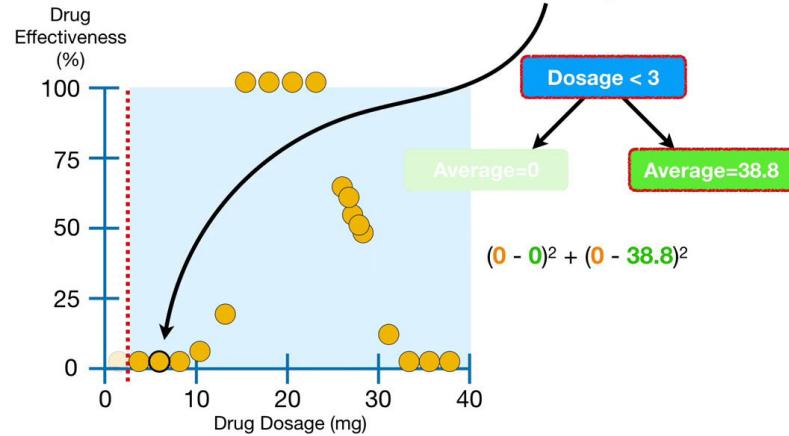
...and square it..

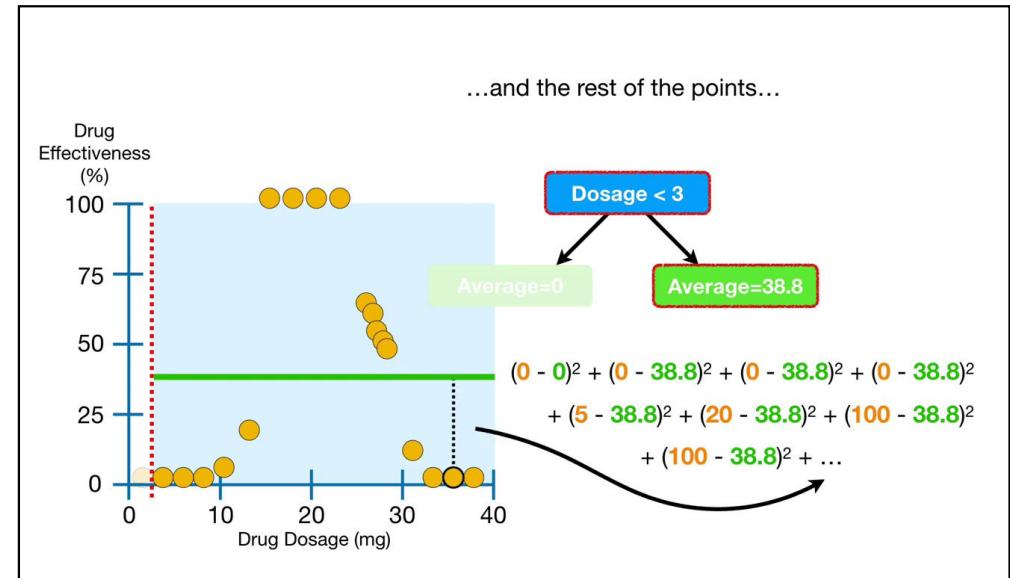
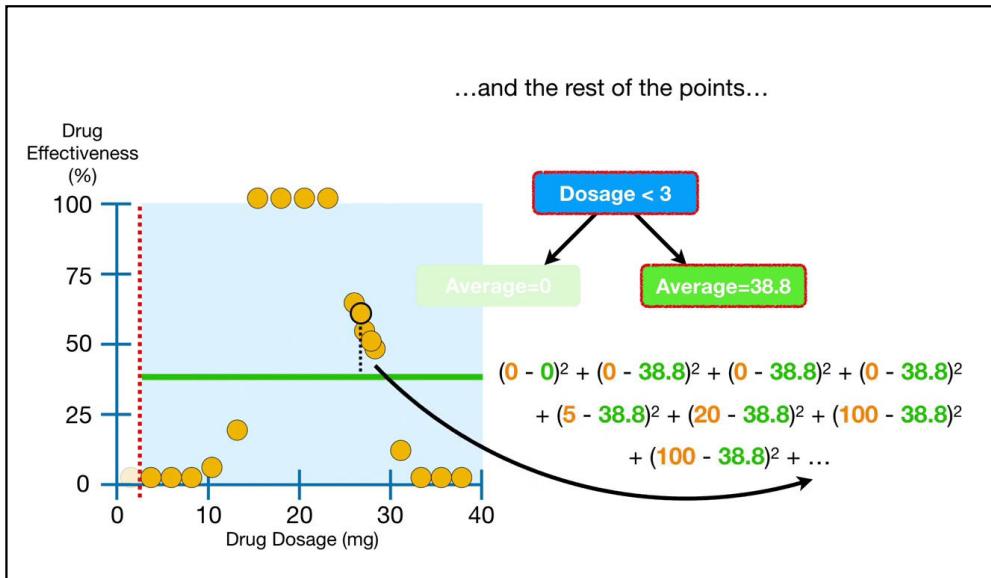
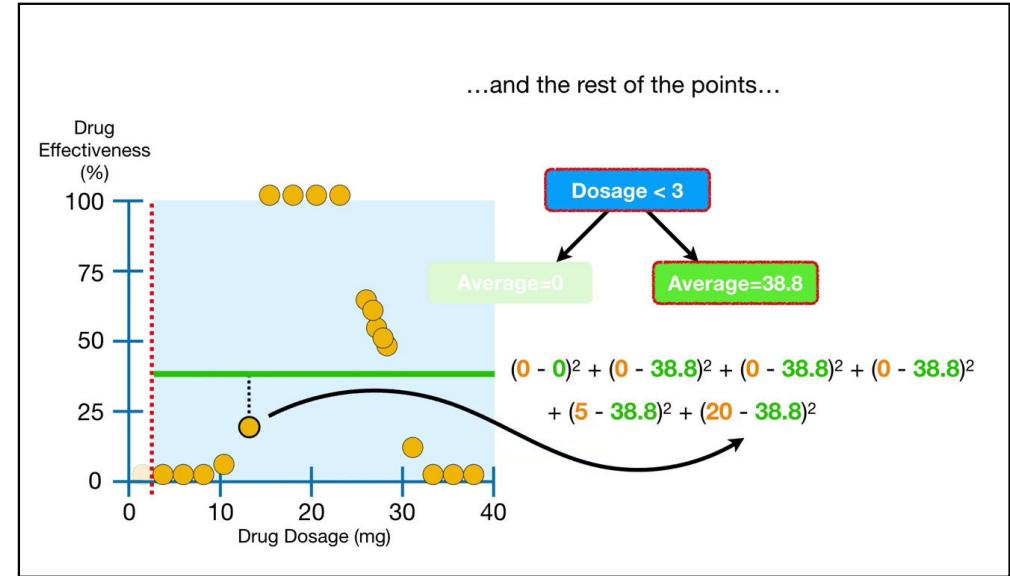
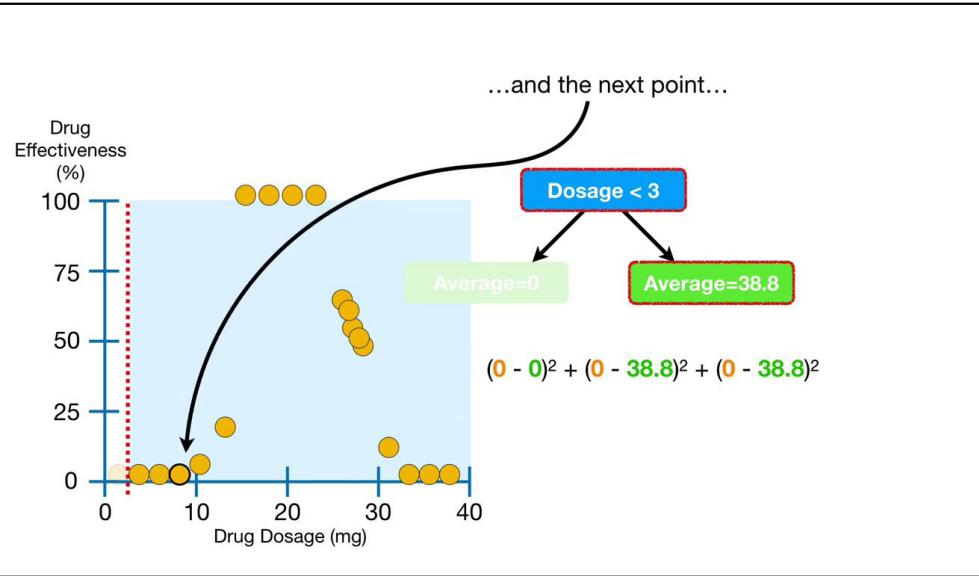


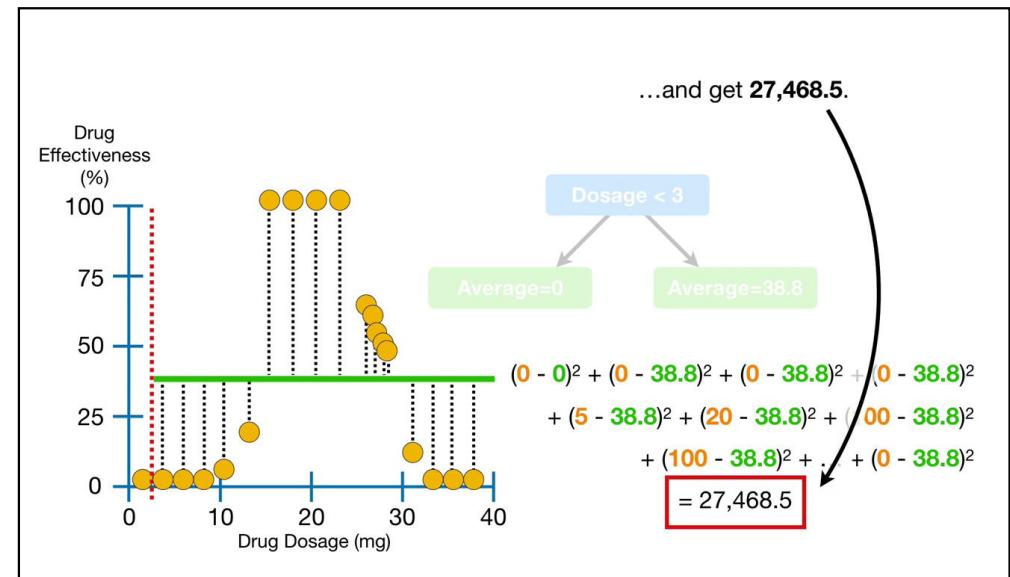
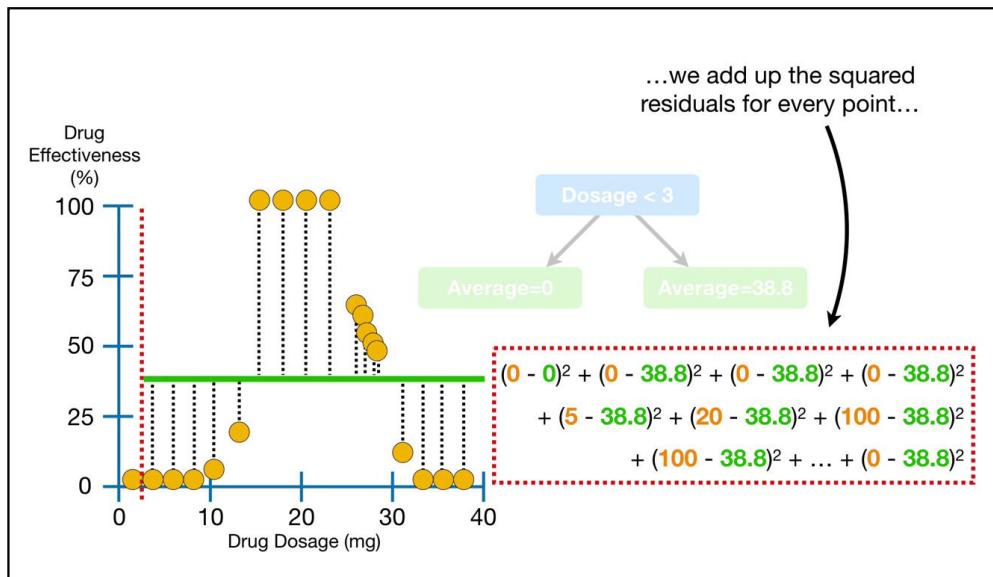
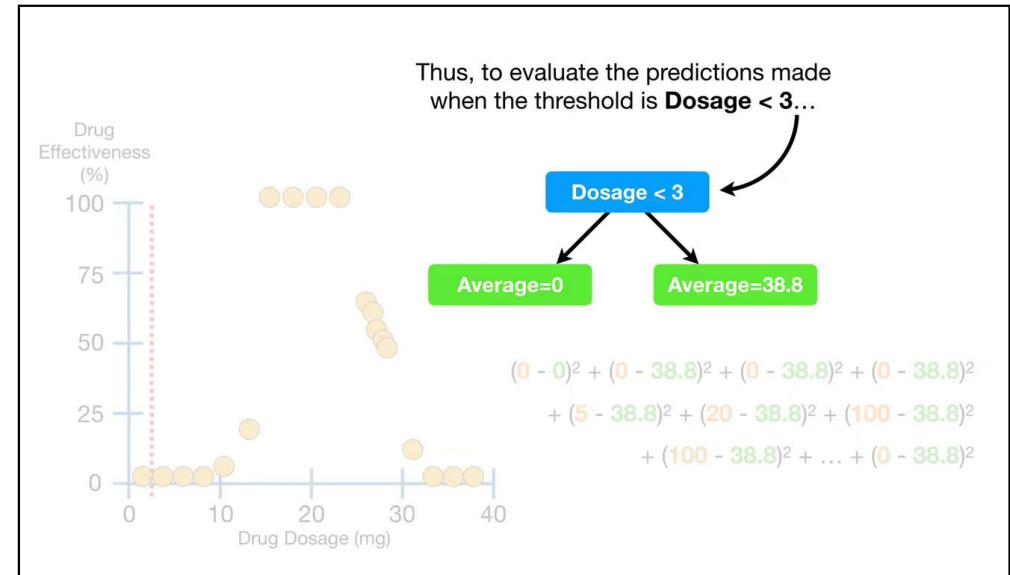
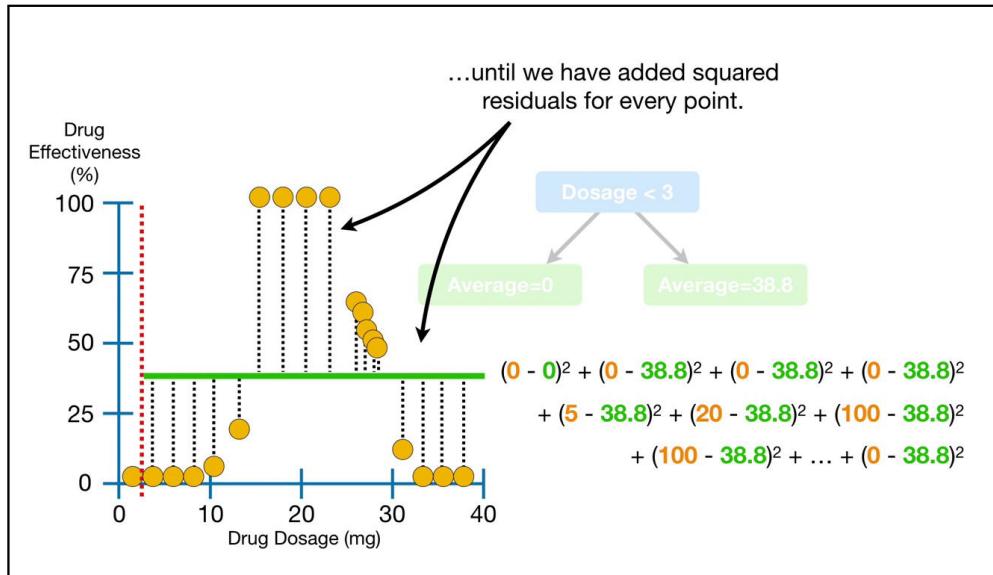
...and then add it to the first term.

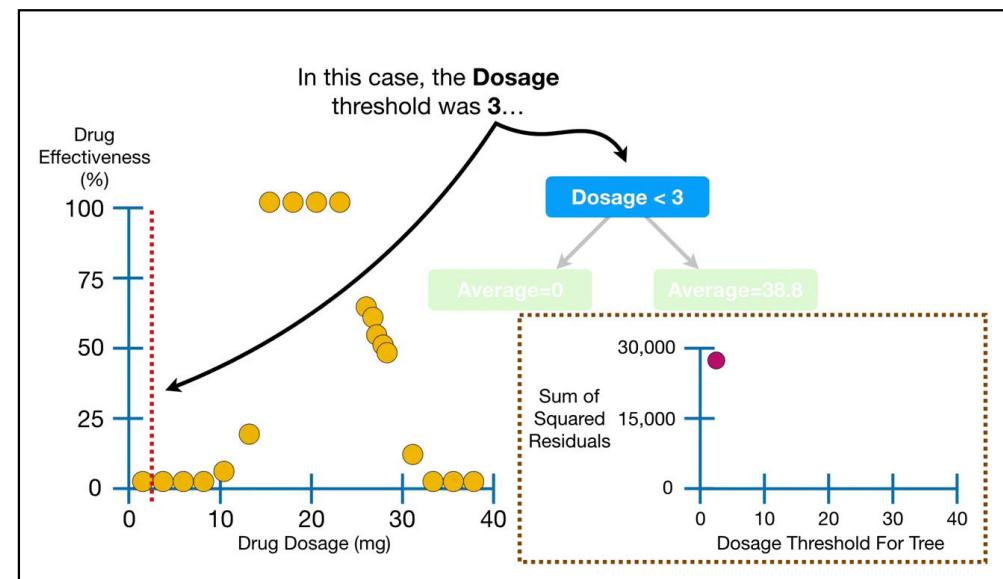
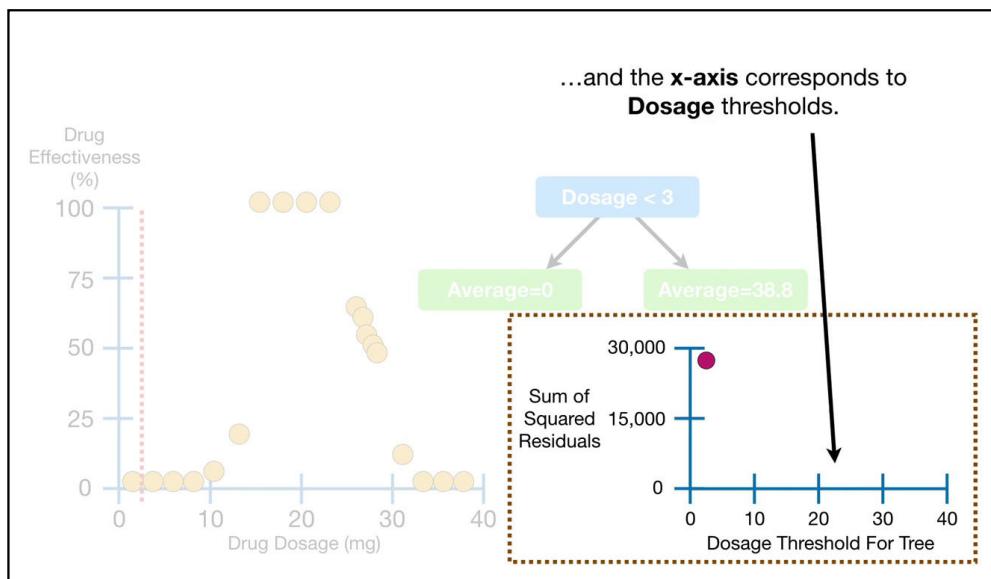
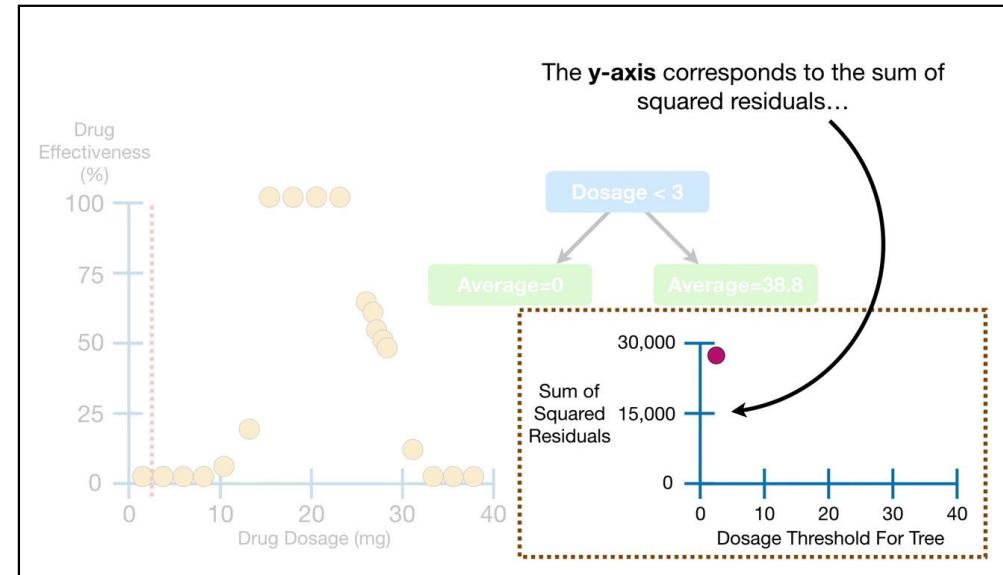
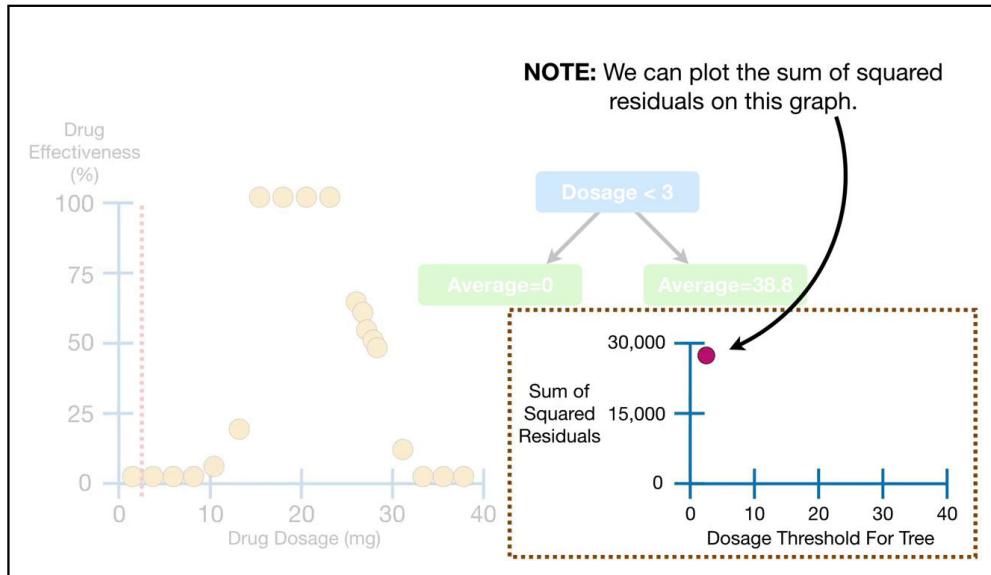


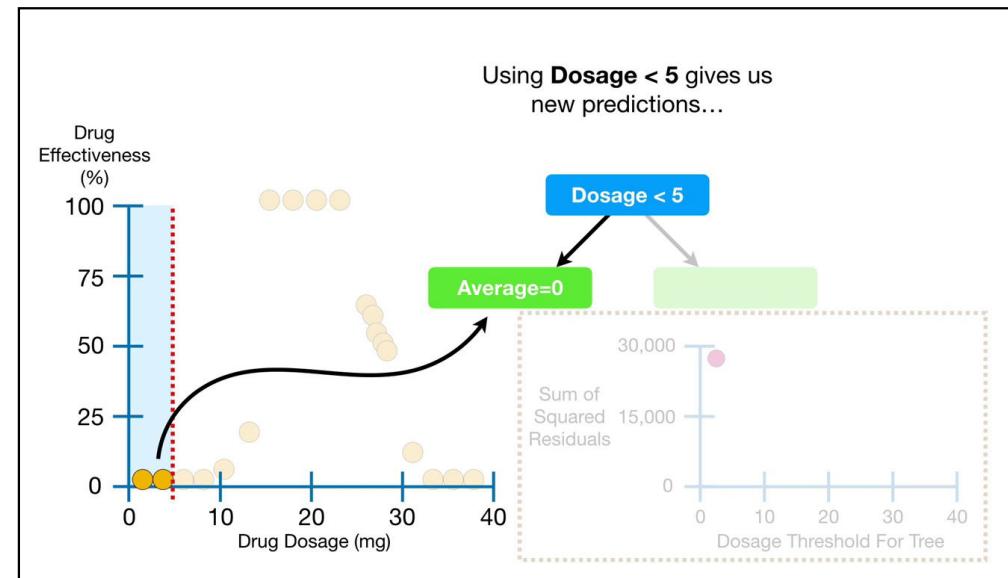
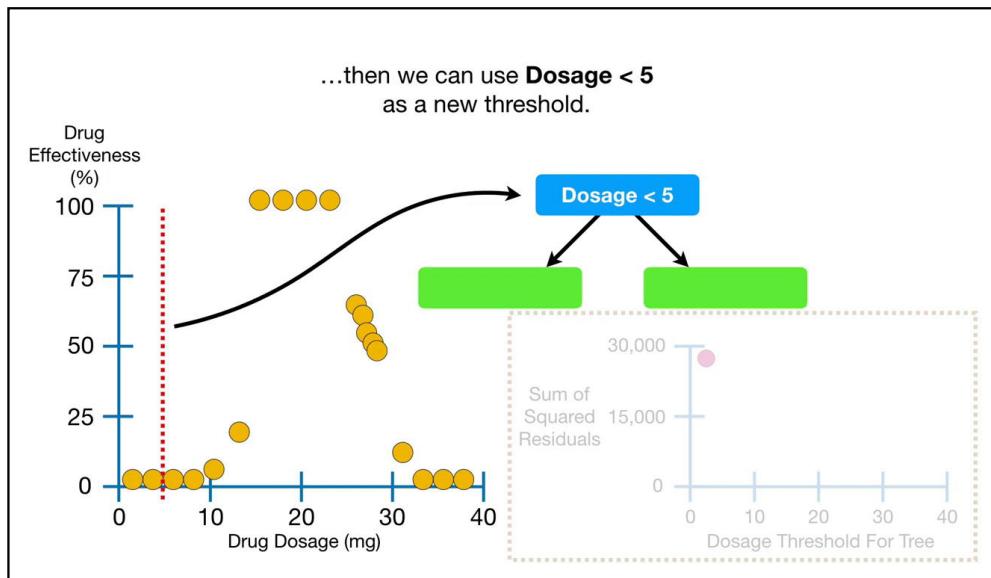
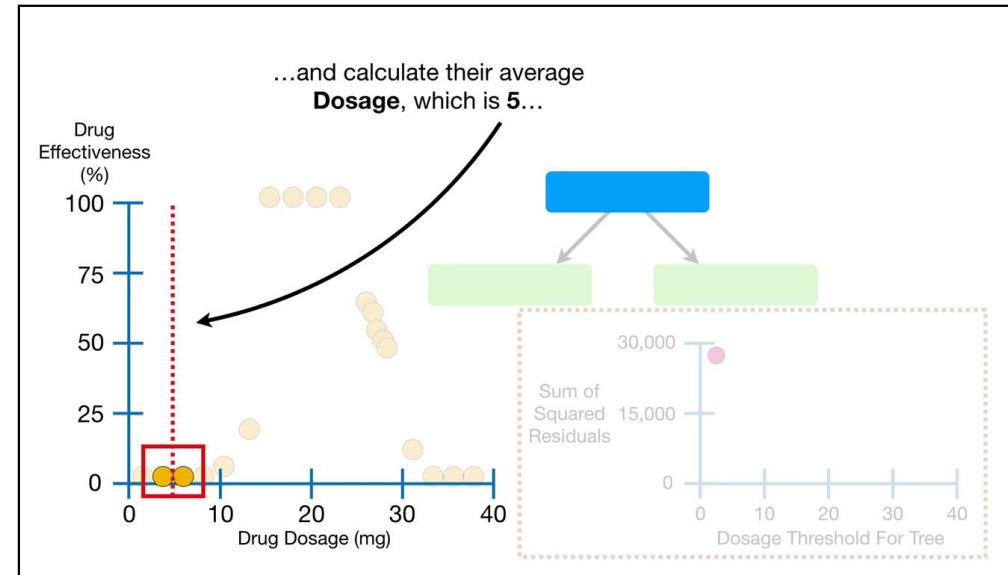
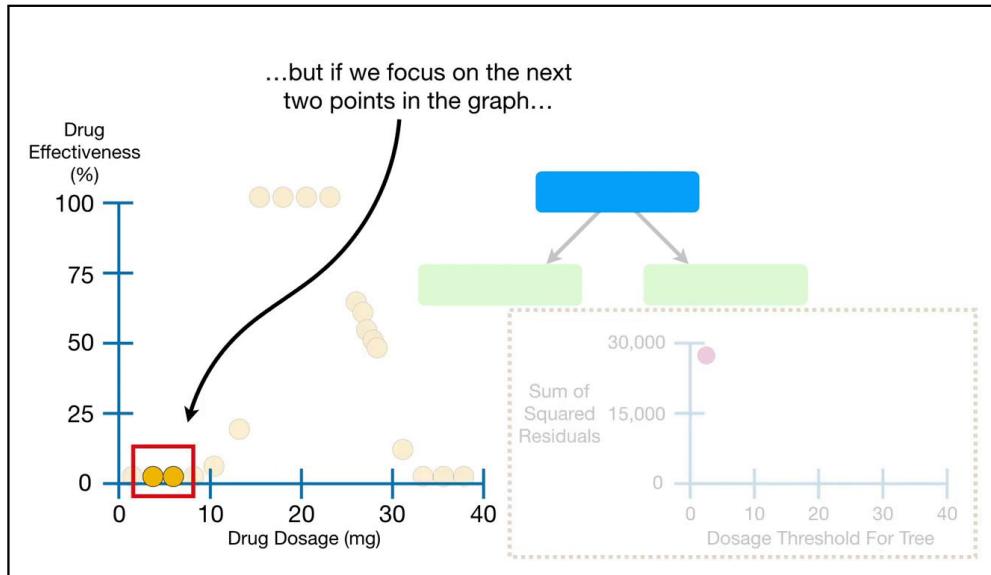
Then we do the same thing for the next point...

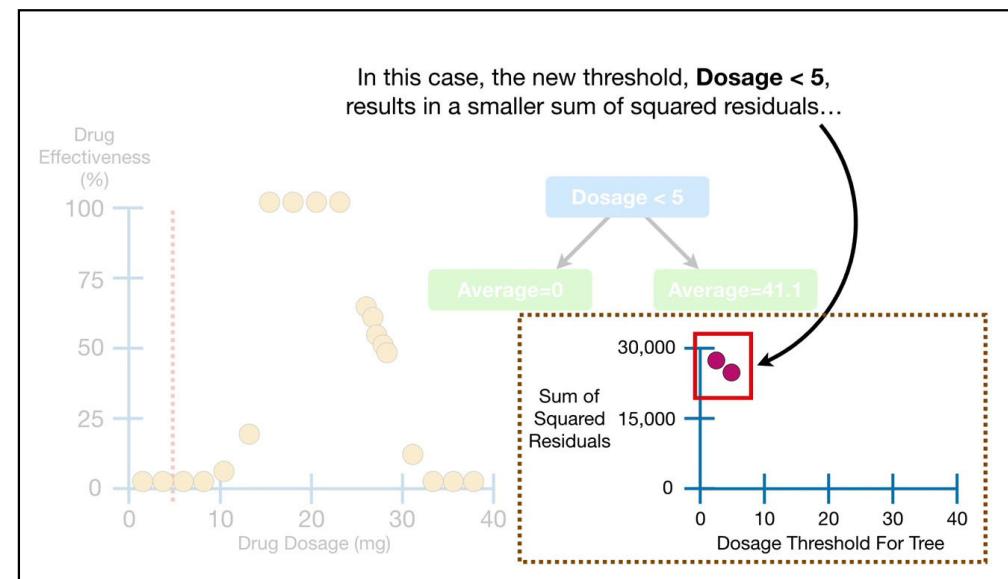
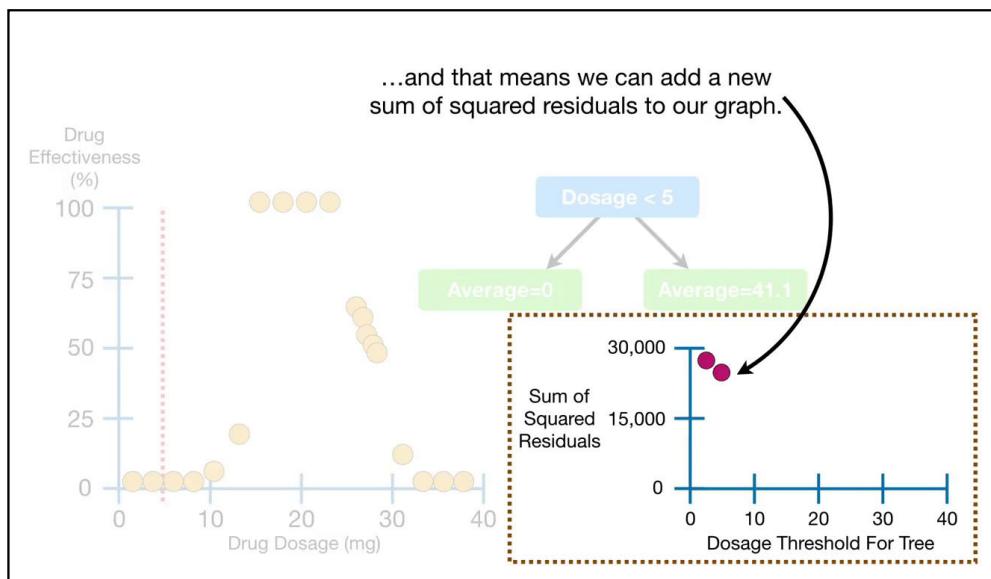
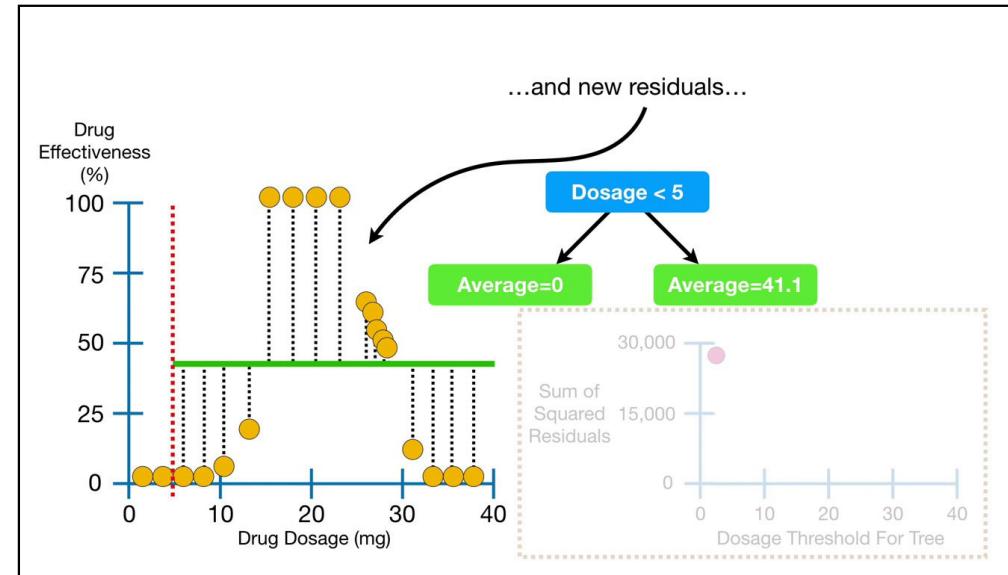
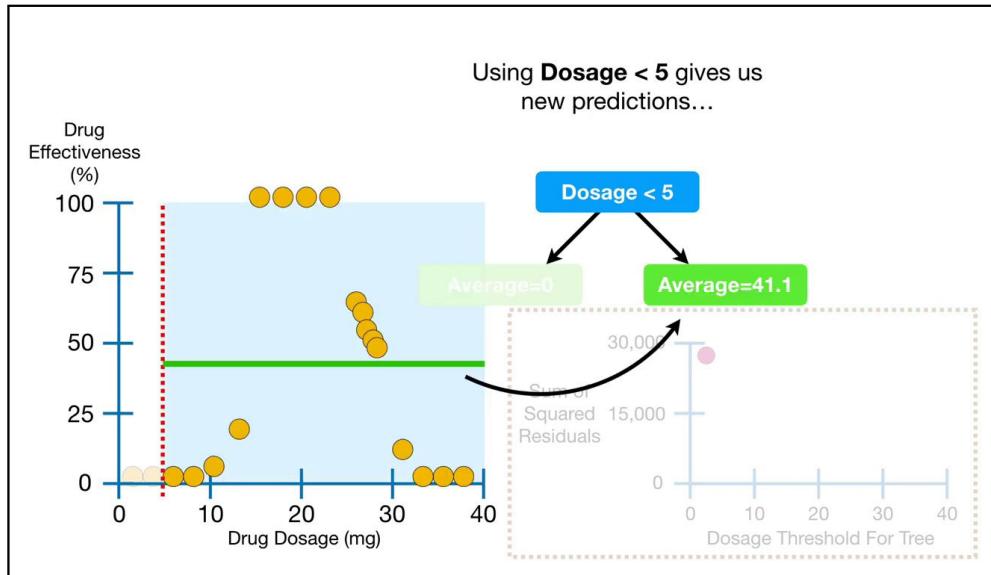


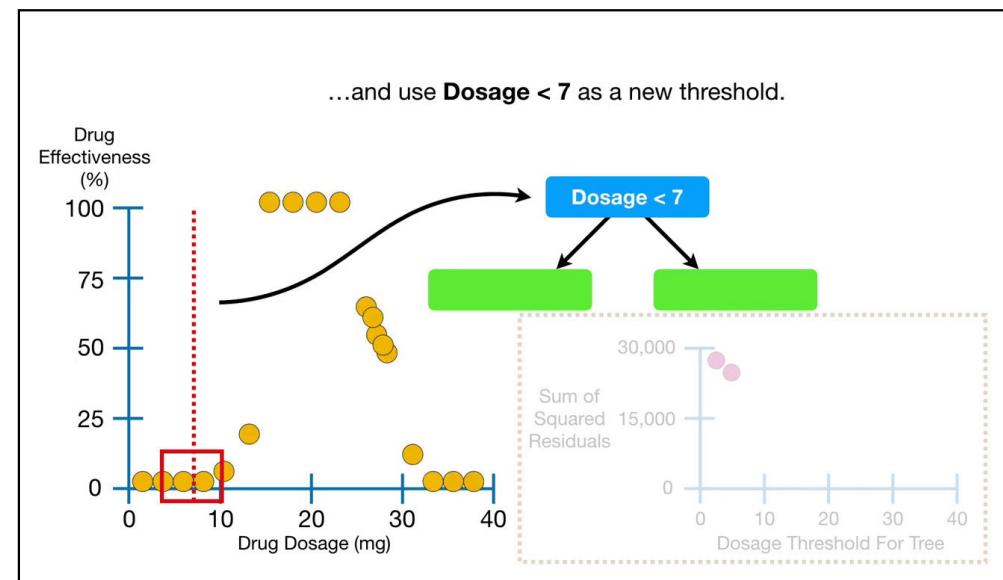
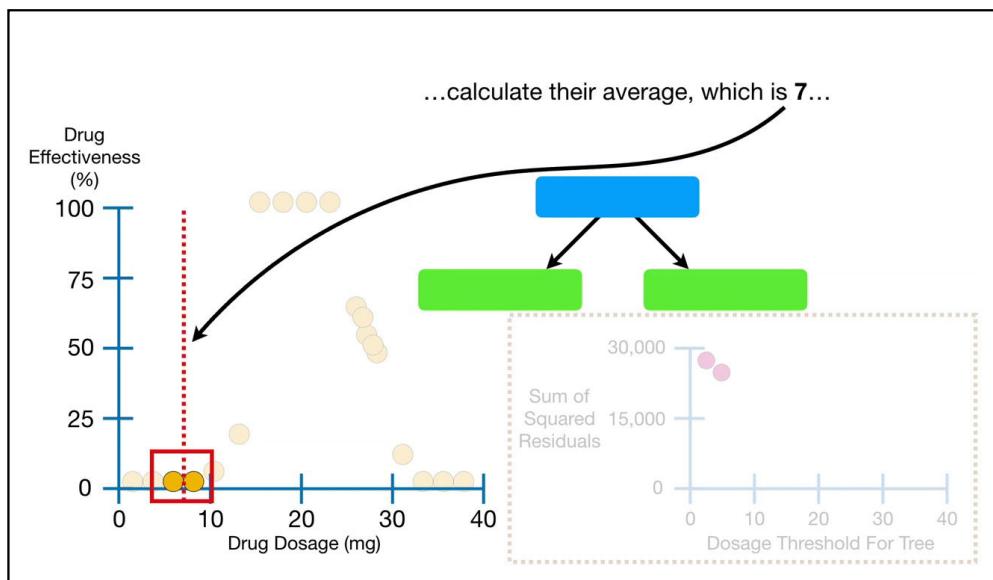
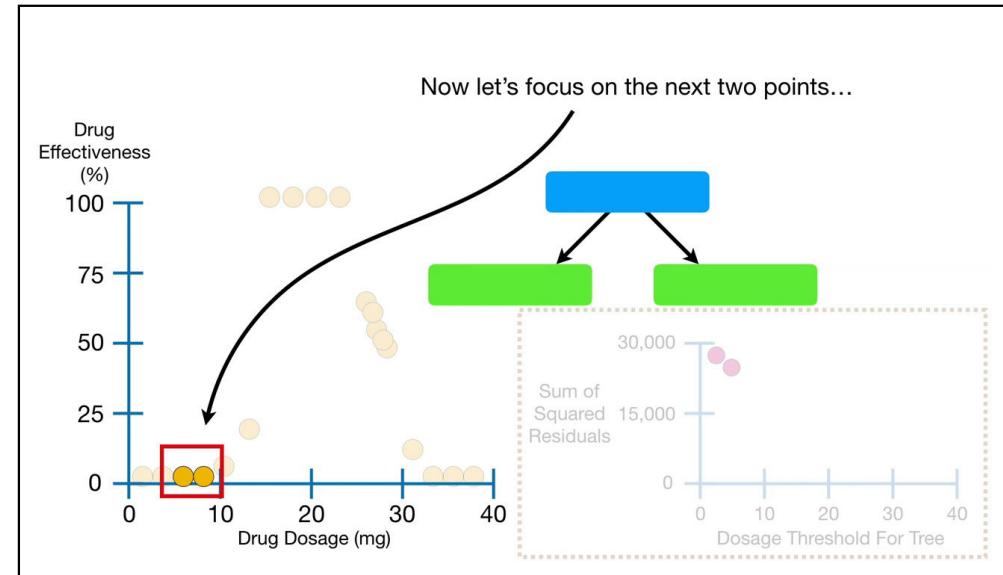
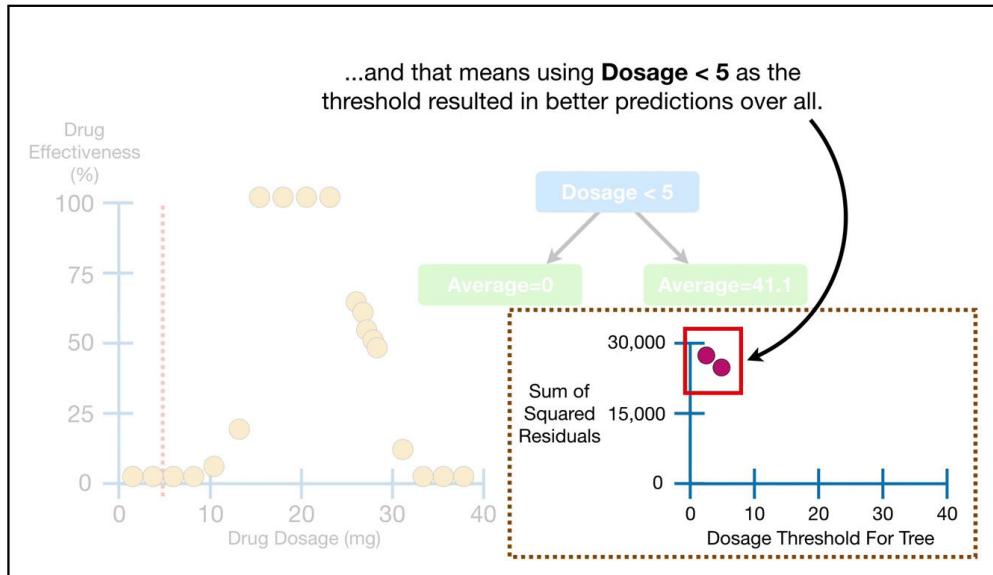


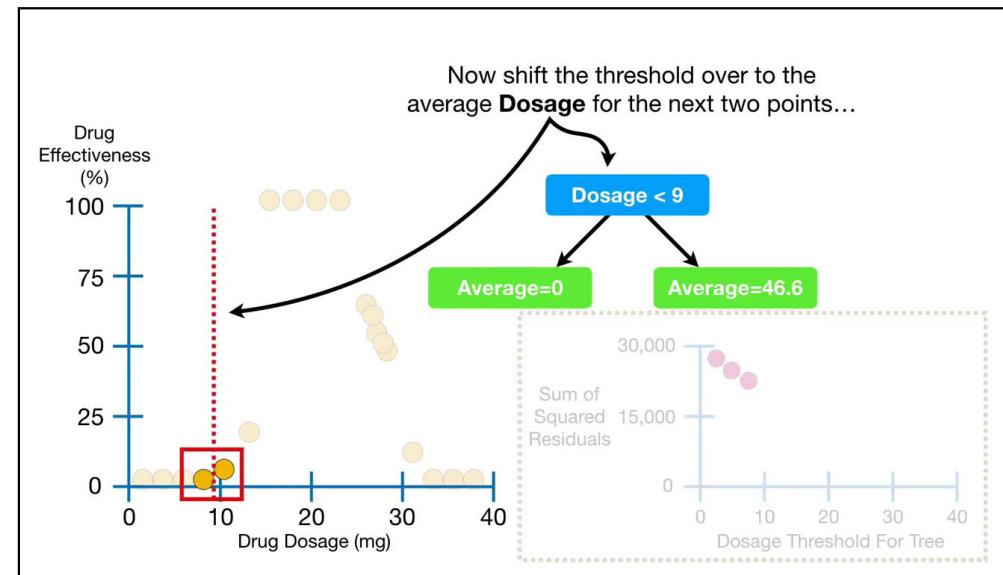
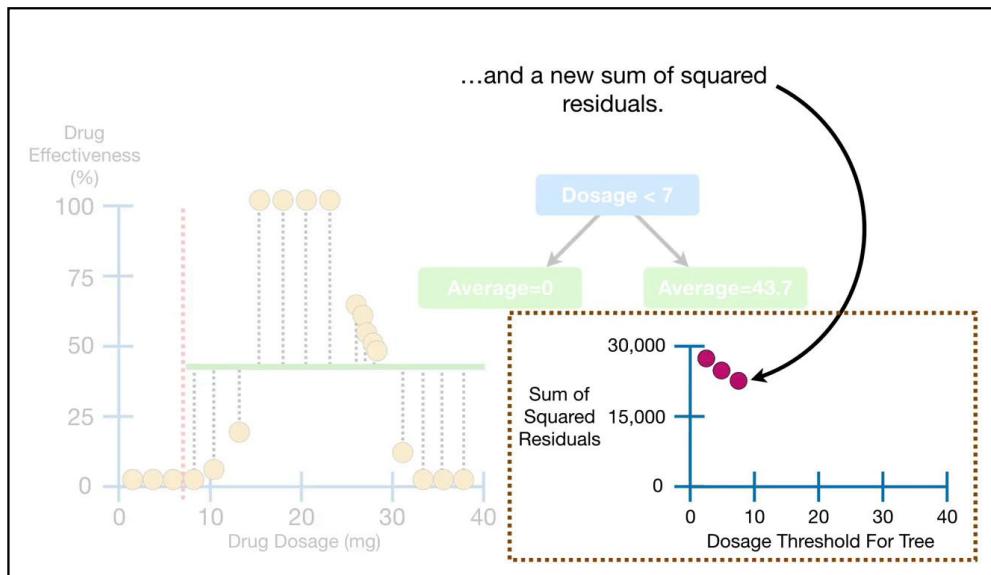
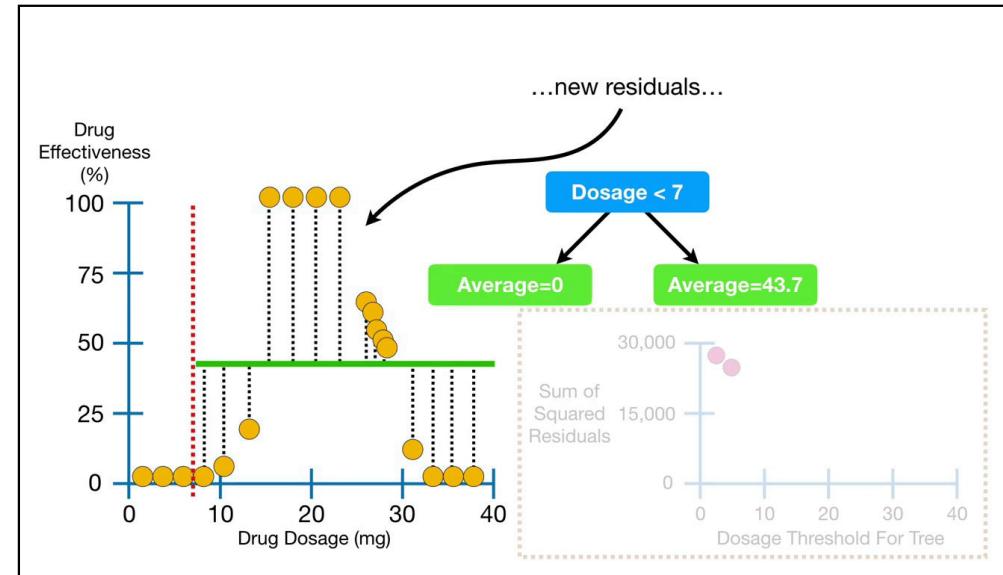
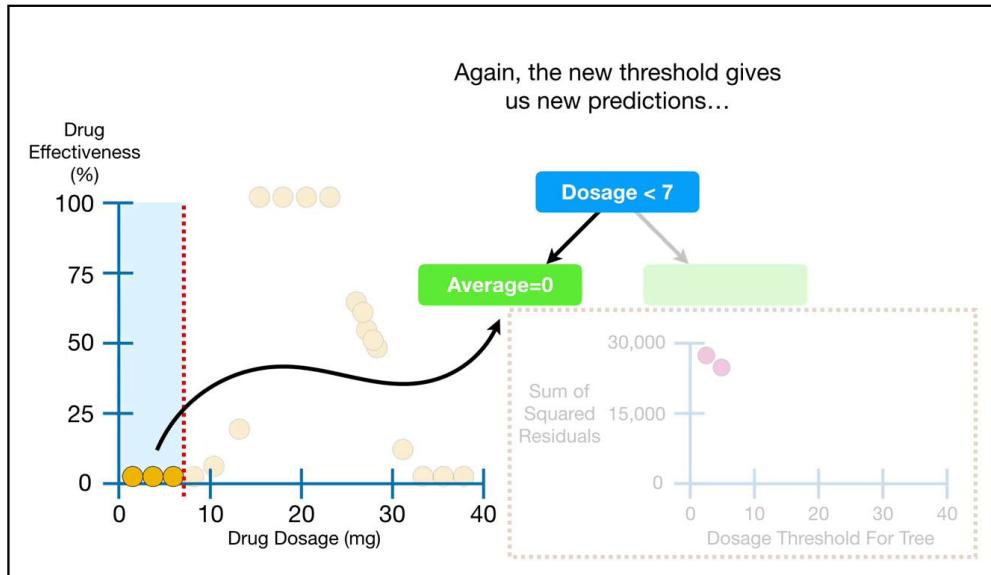


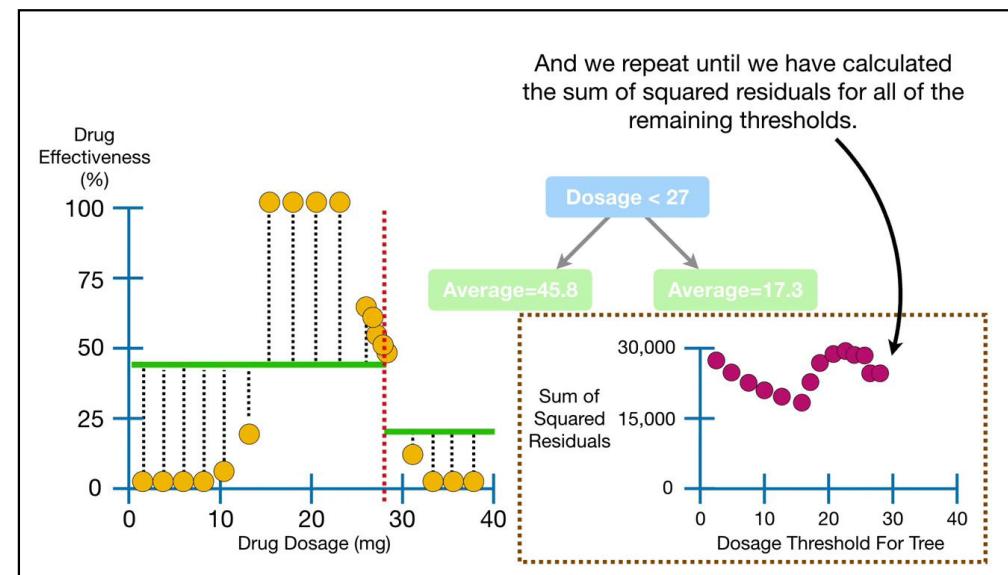
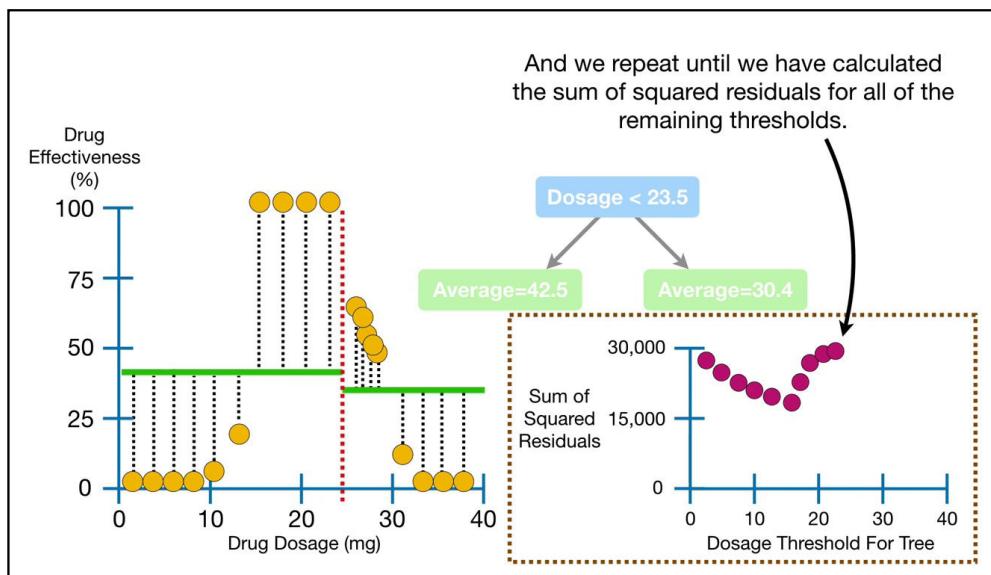
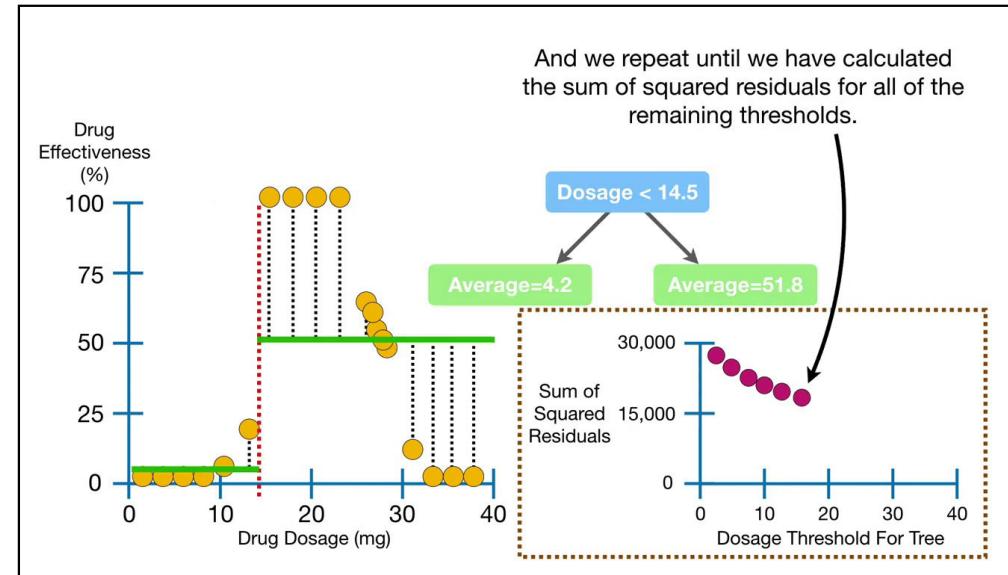
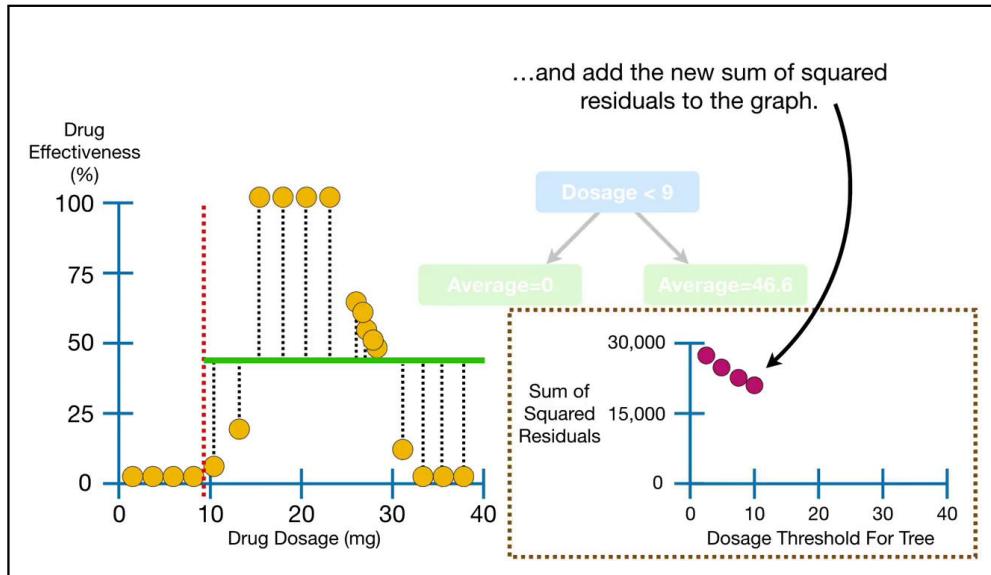


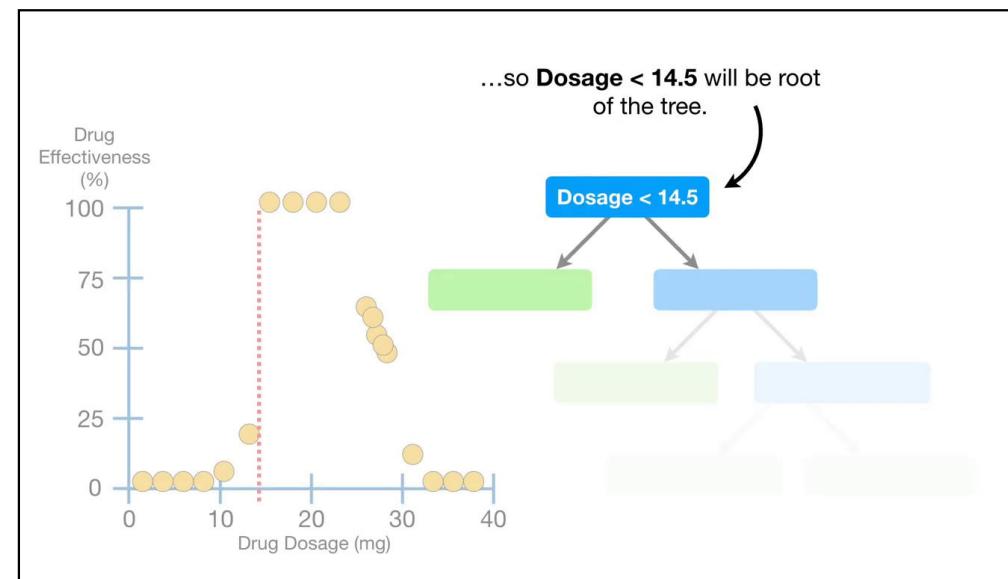
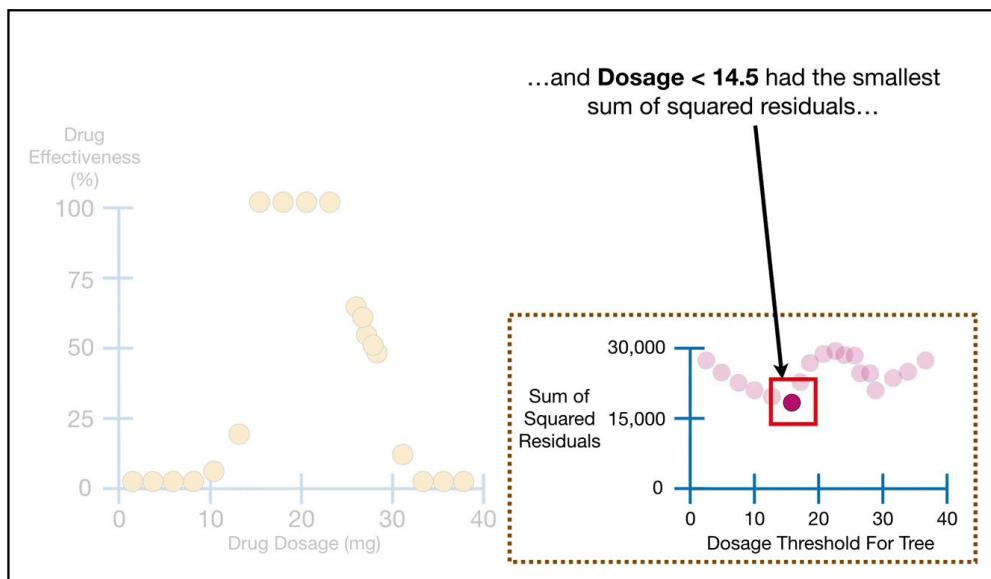
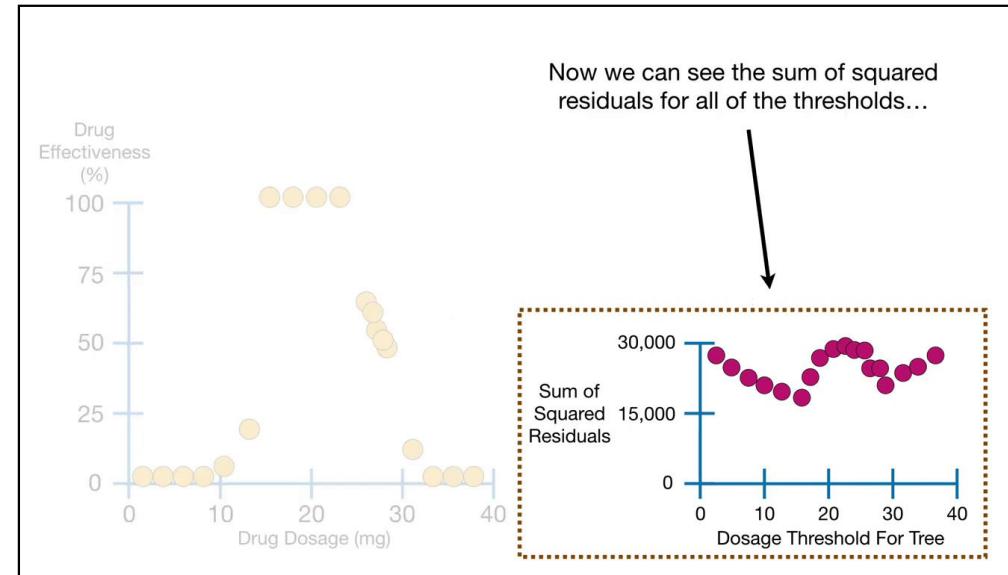
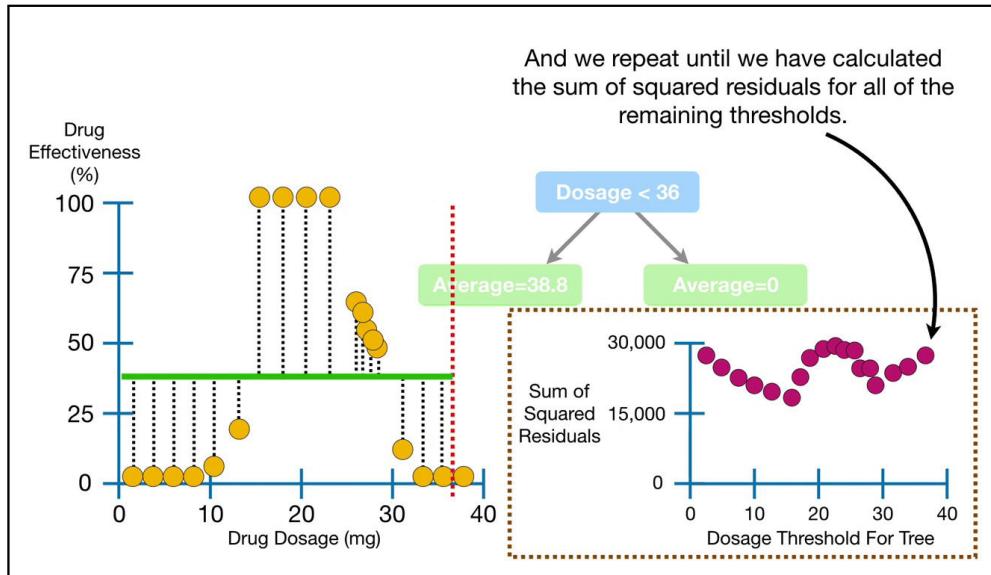


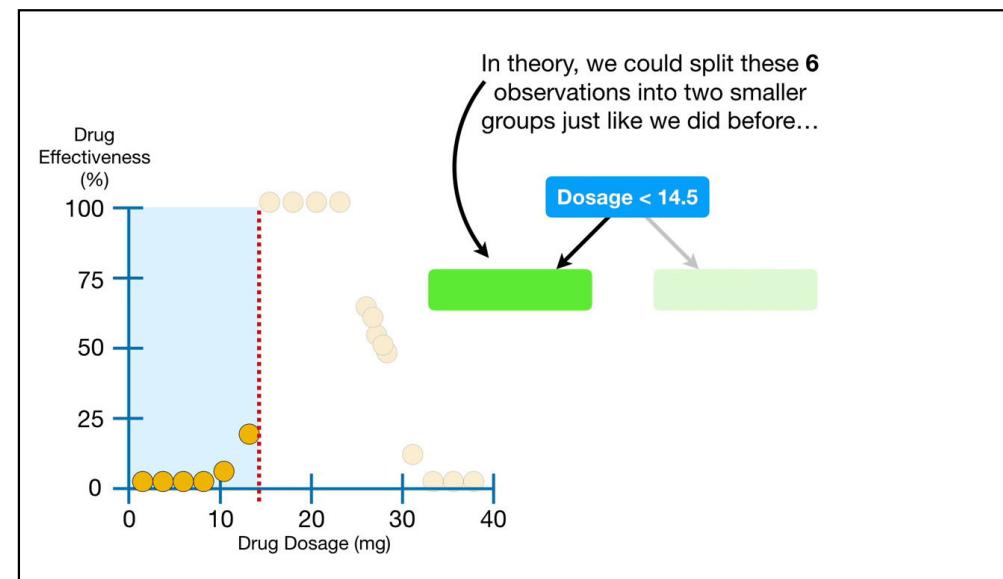
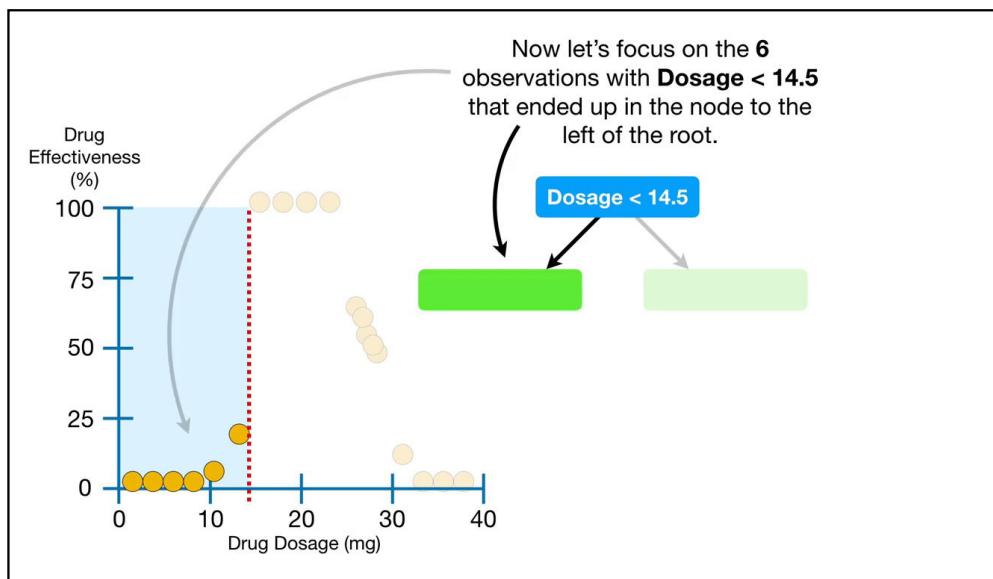
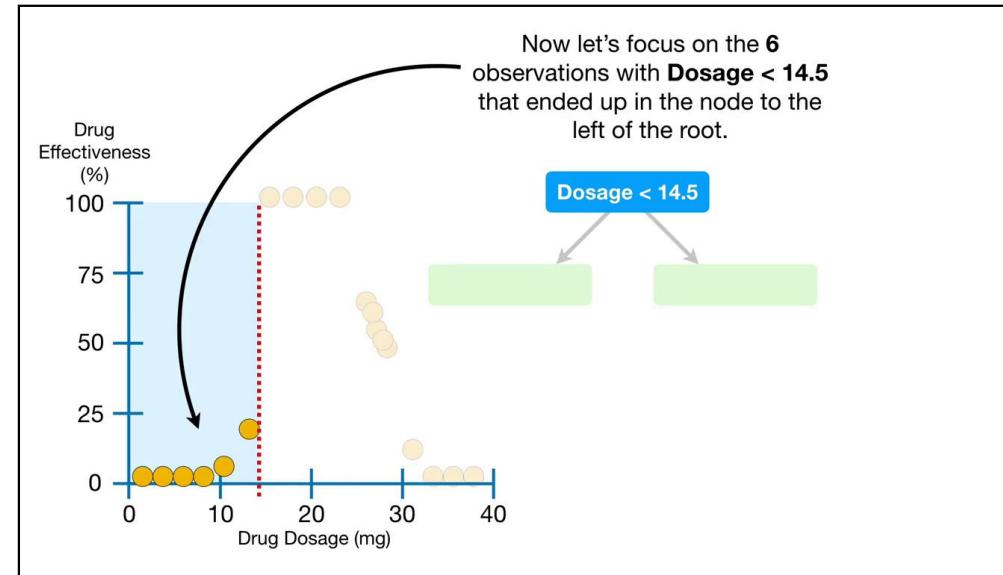
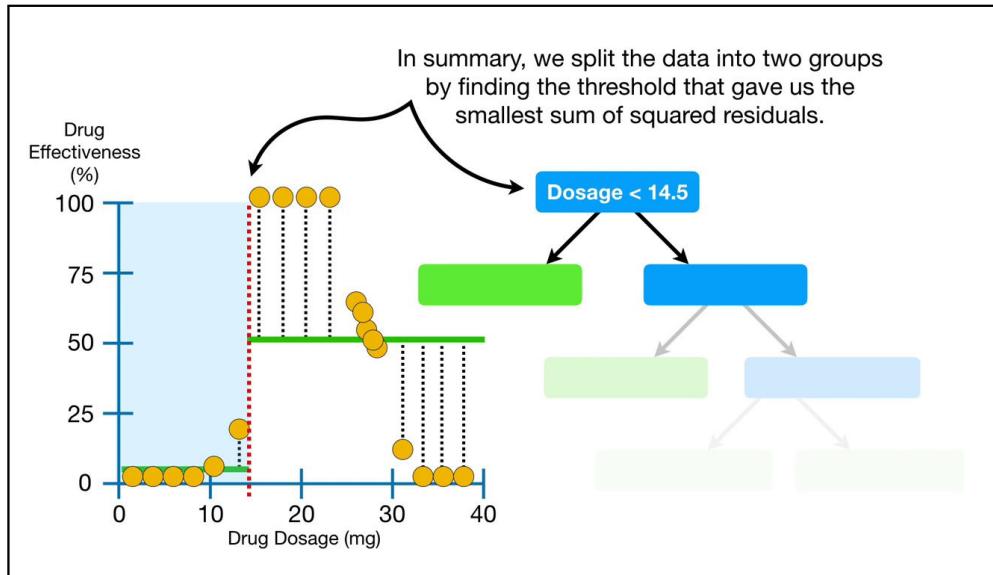


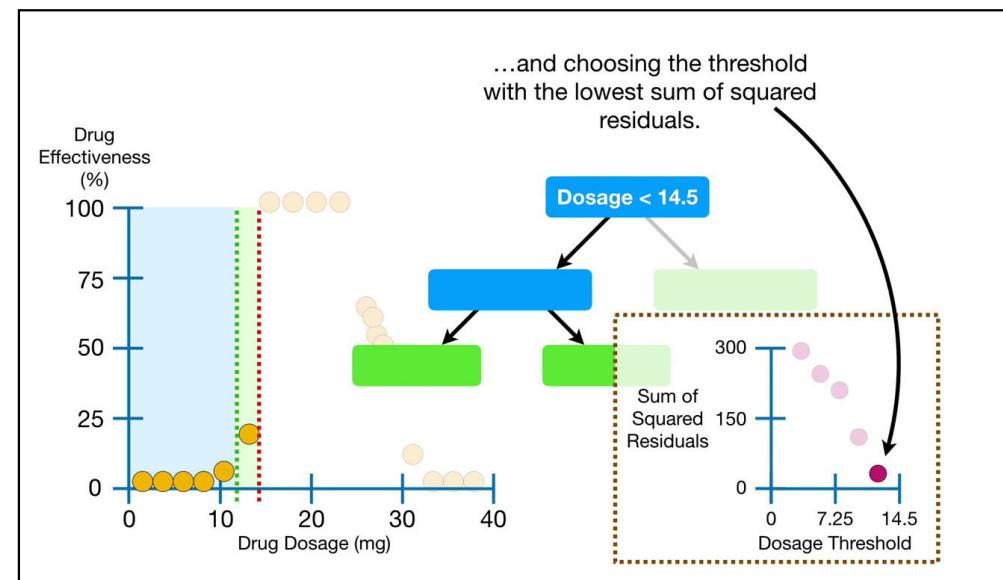
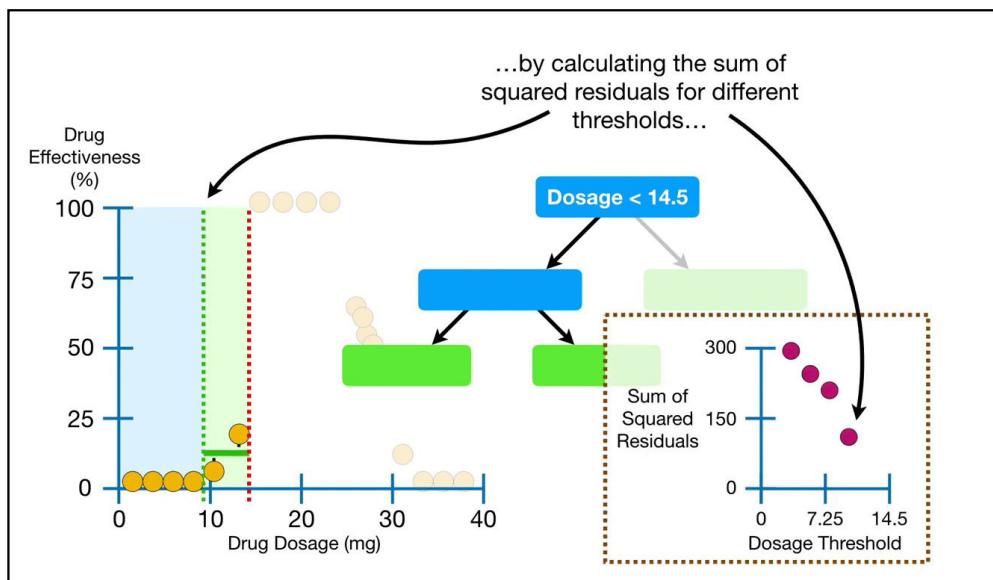
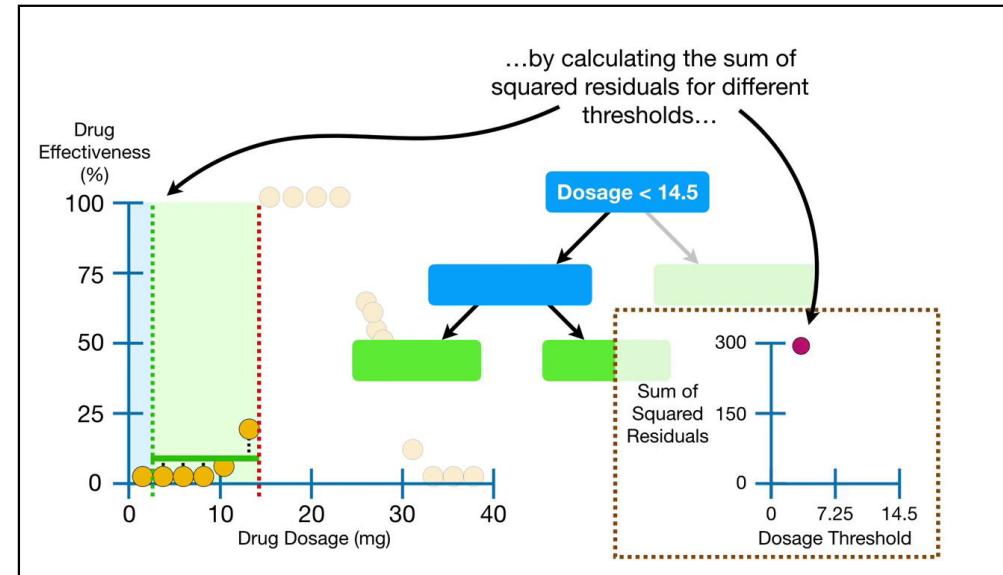
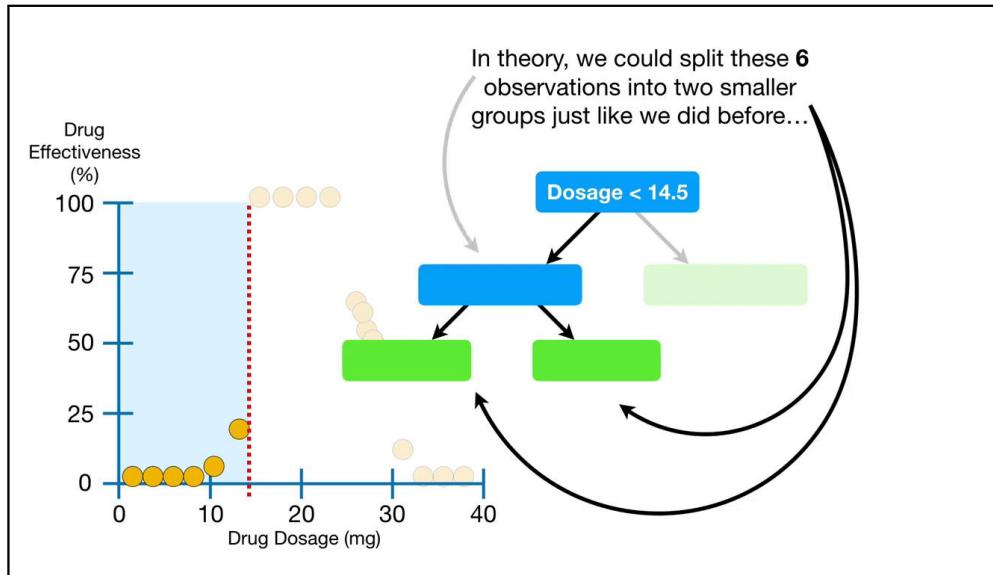


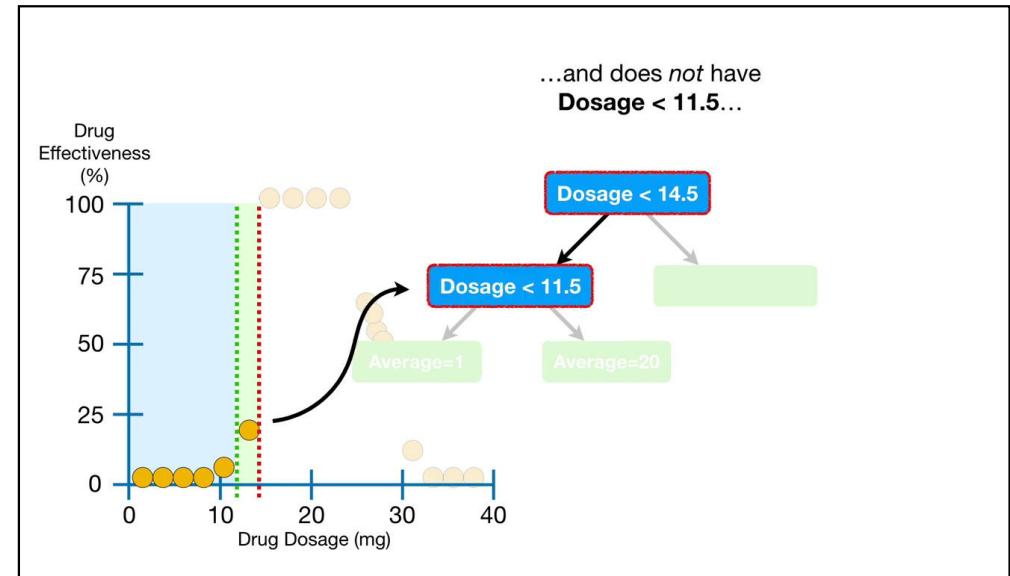
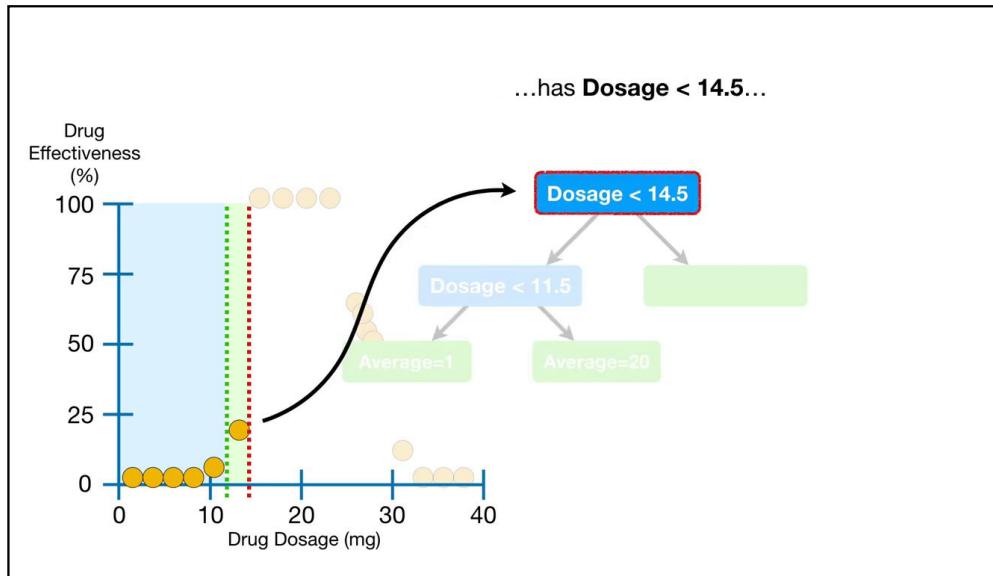
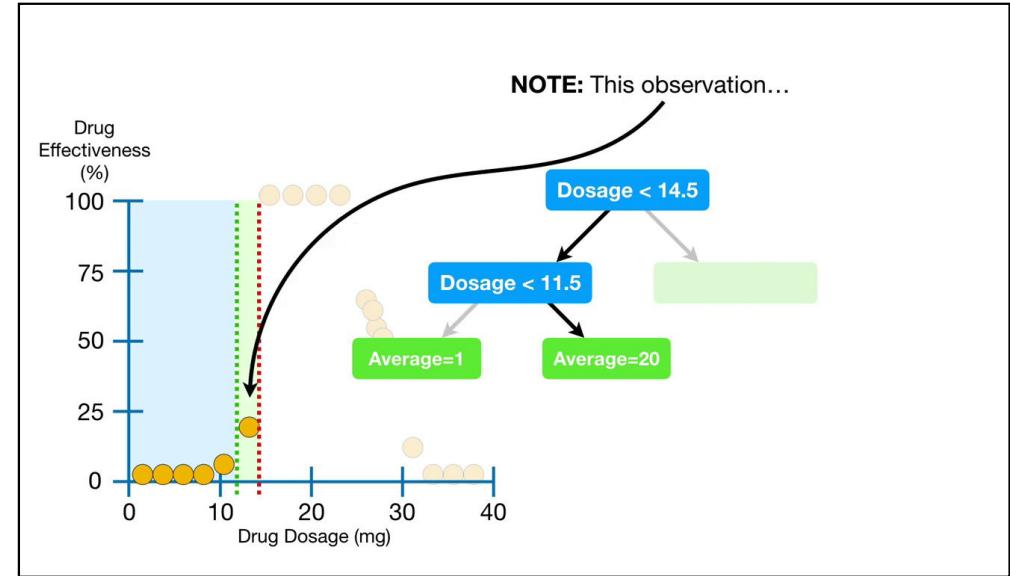
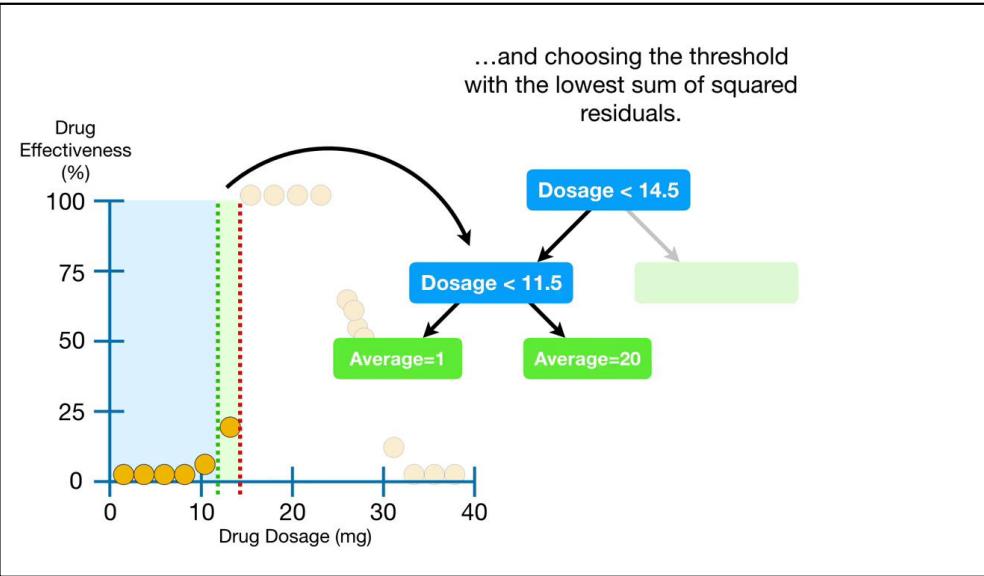


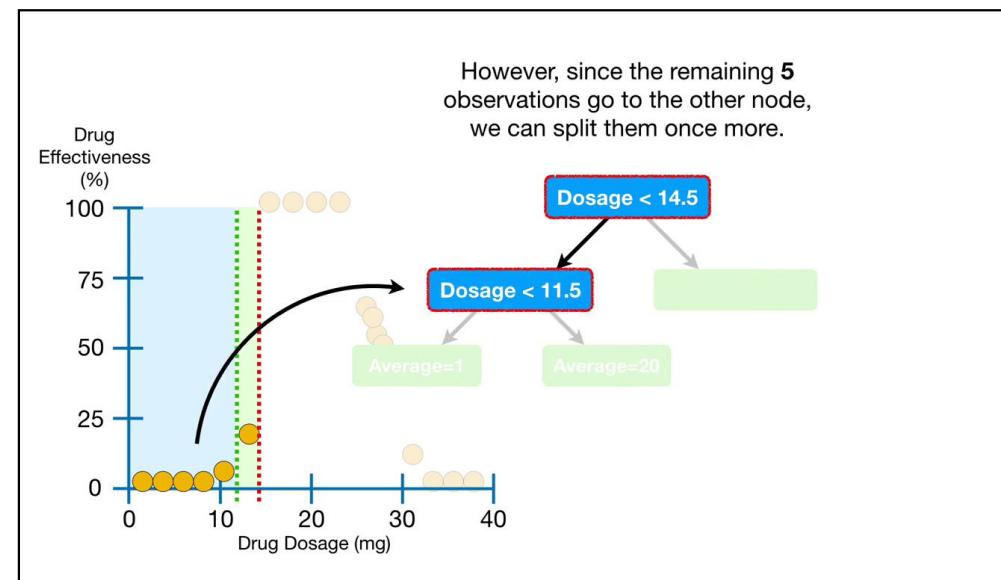
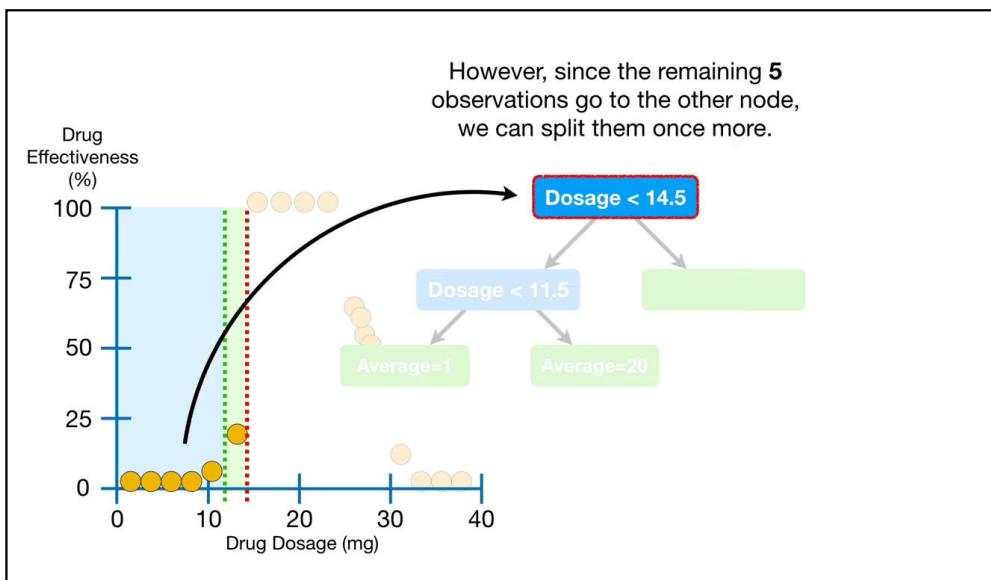
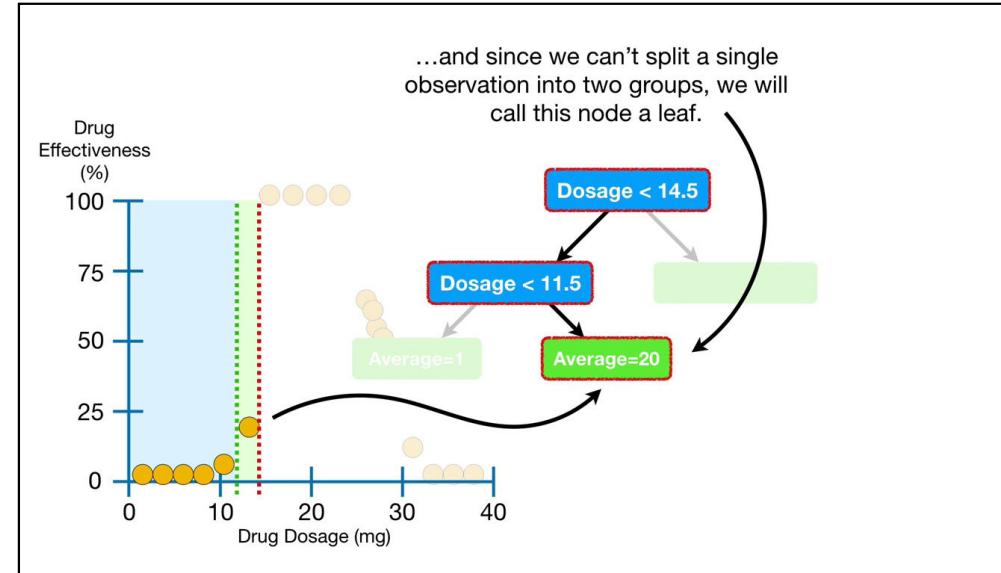
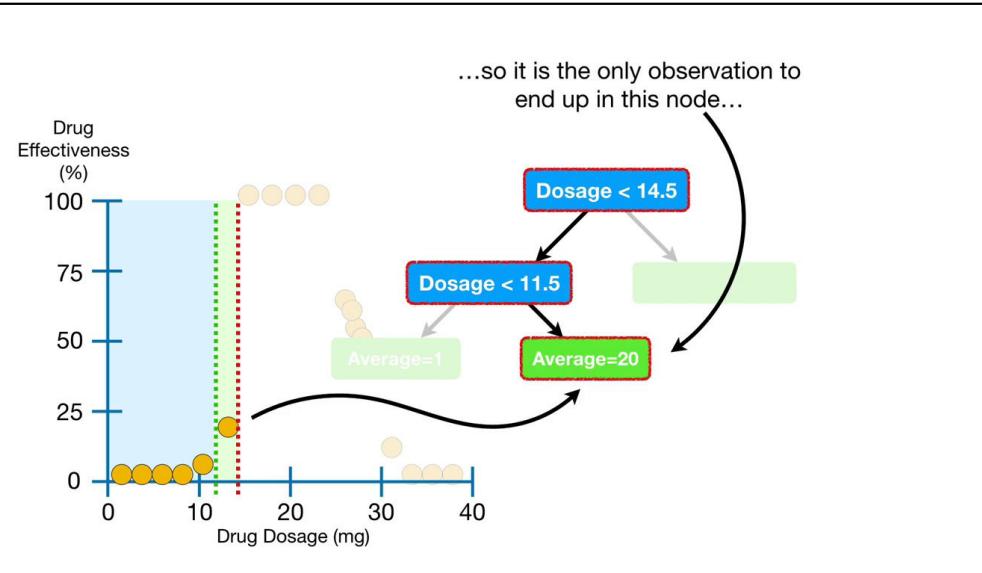


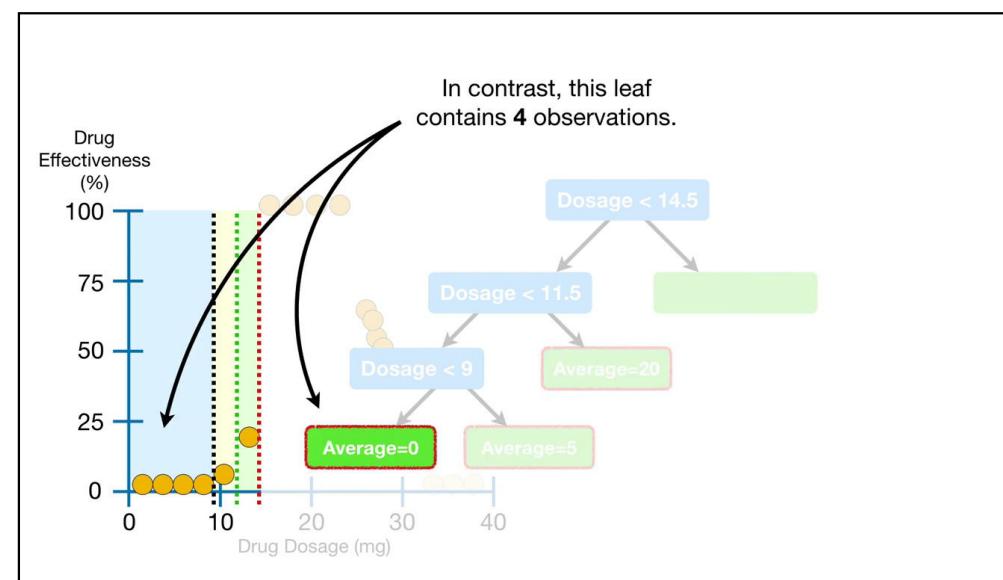
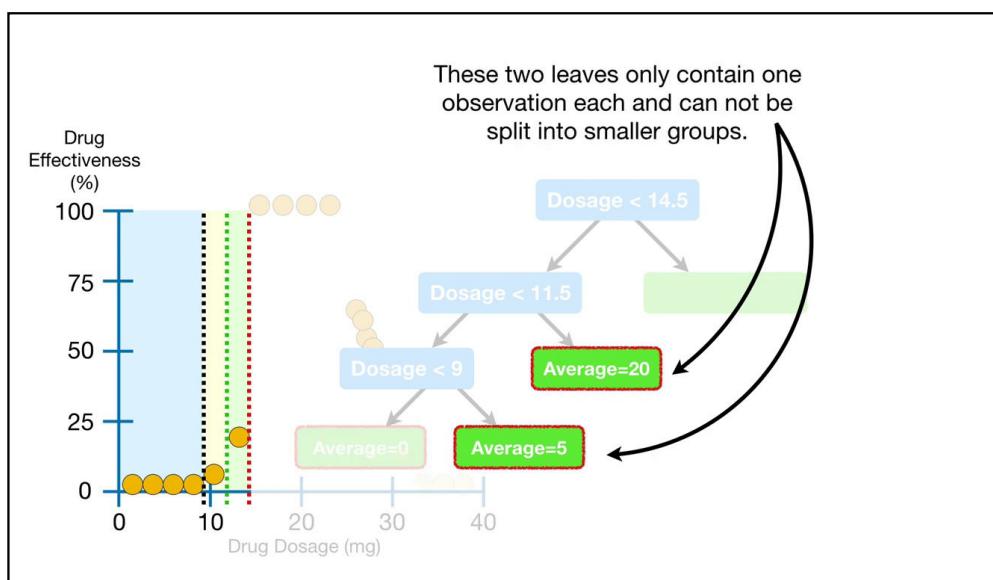
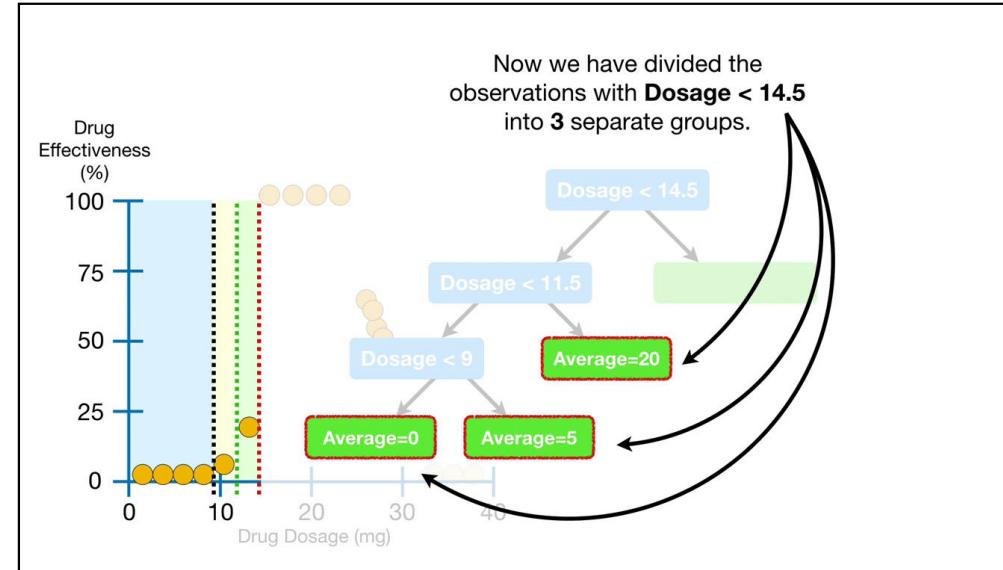
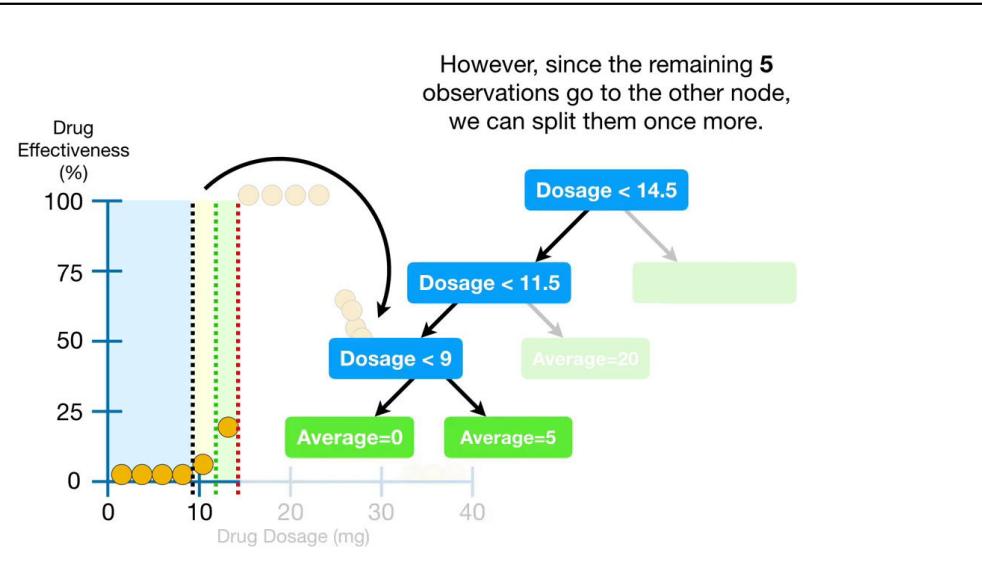


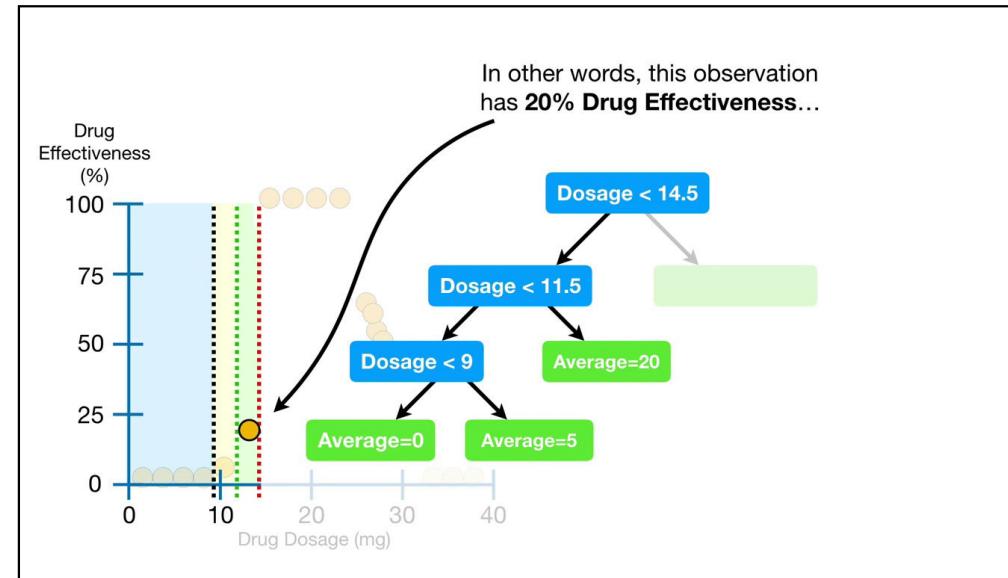
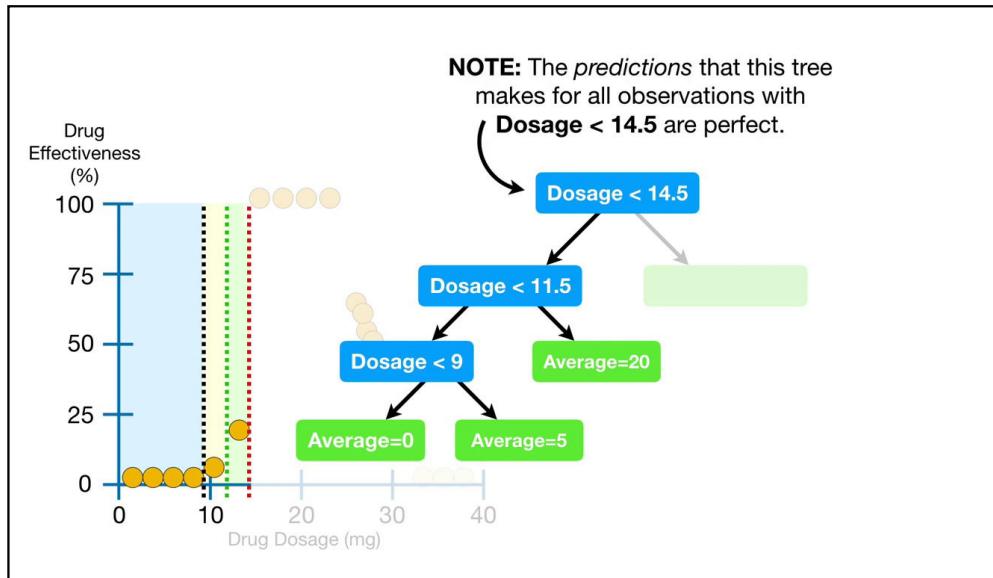
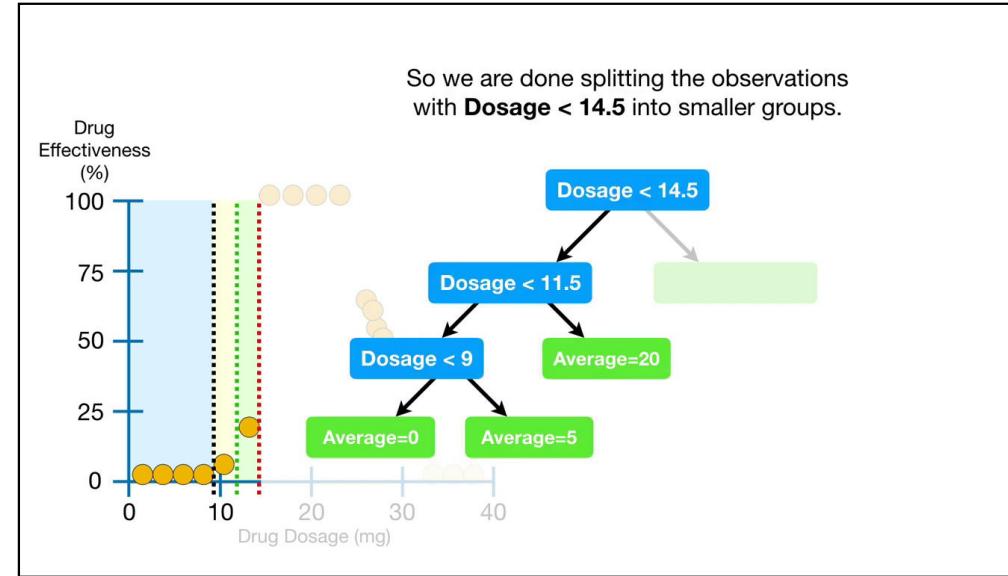
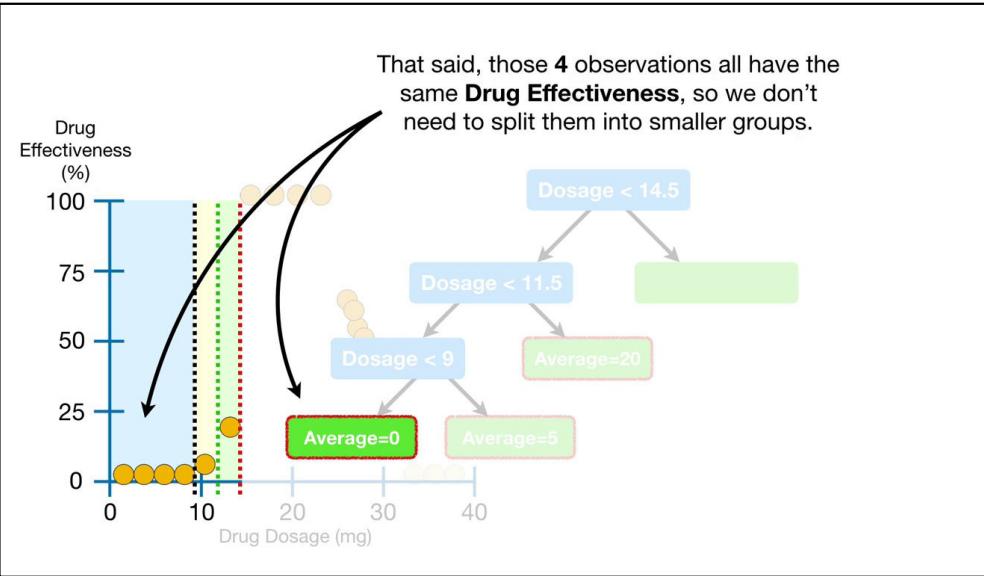


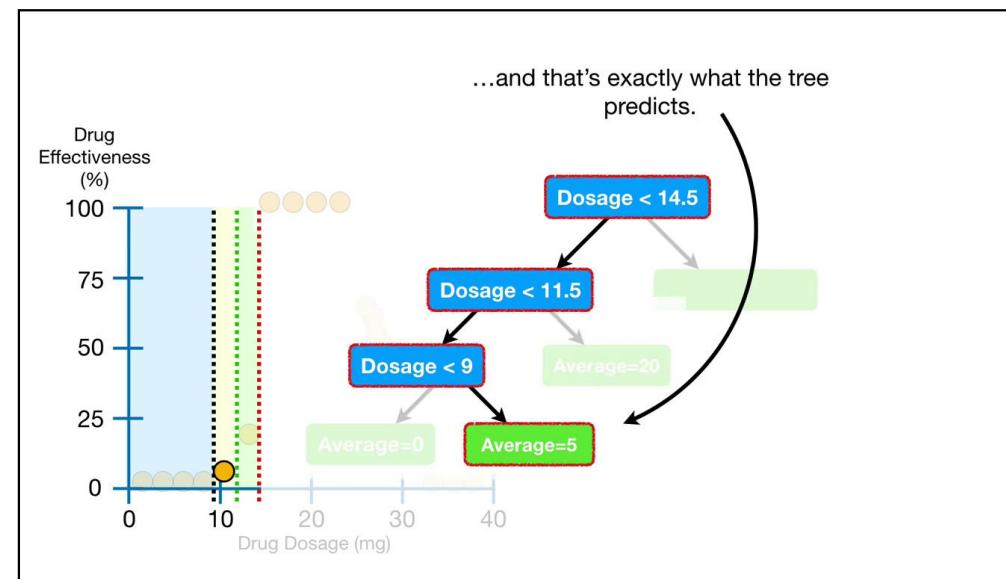
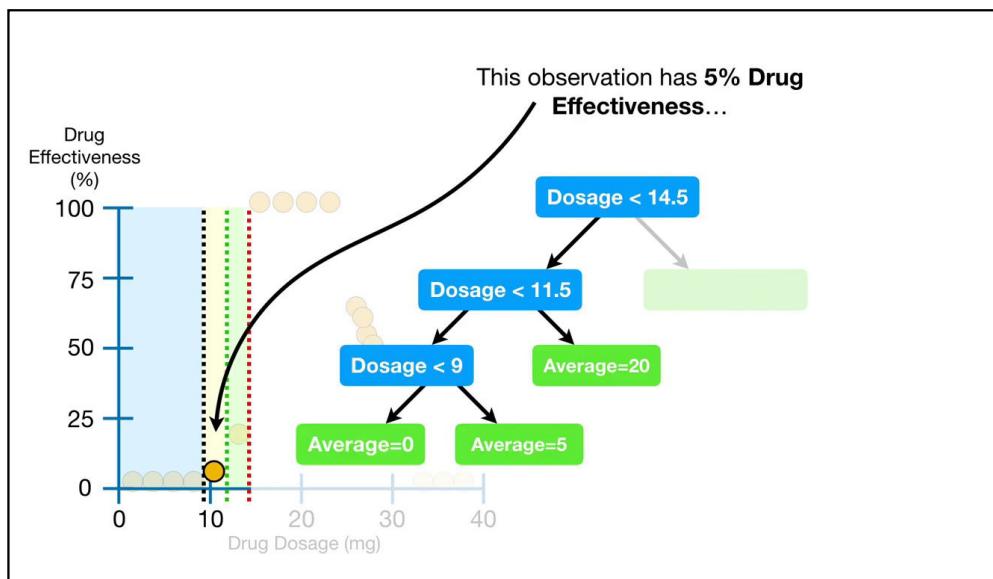
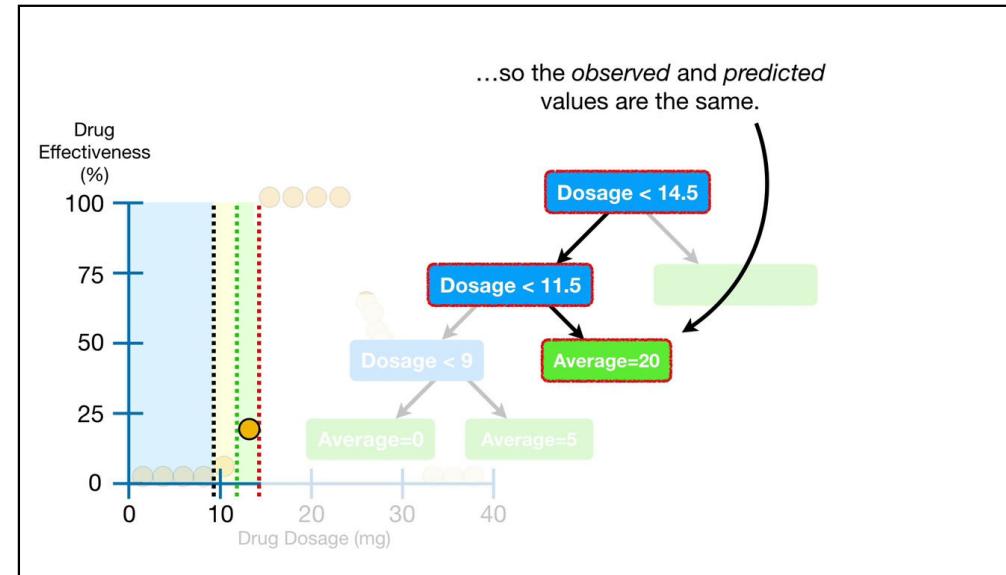
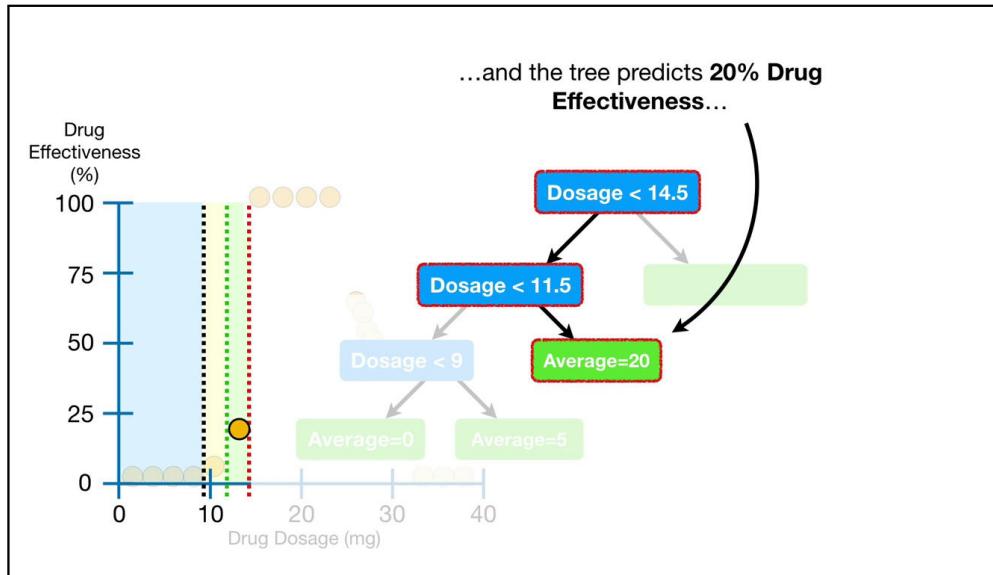


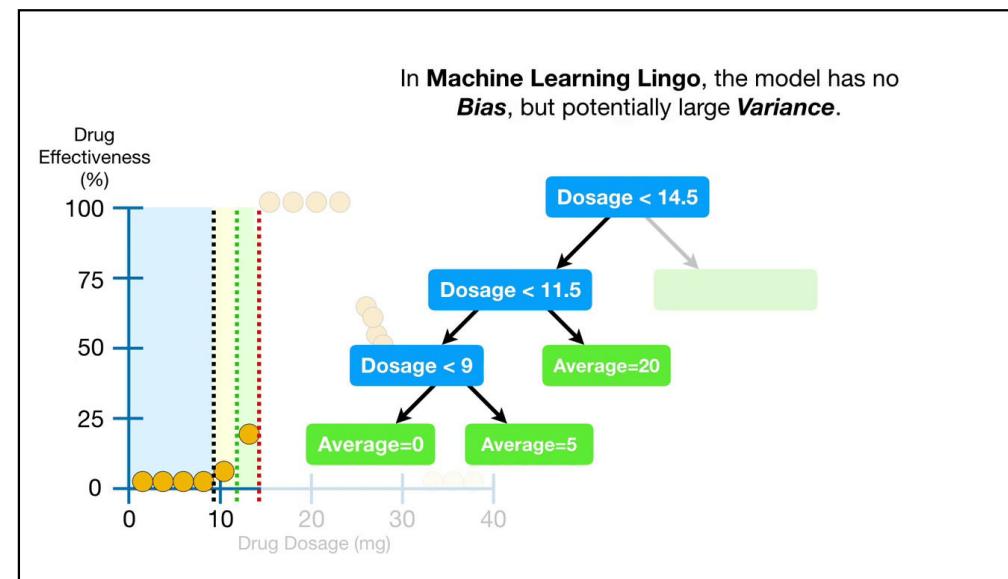
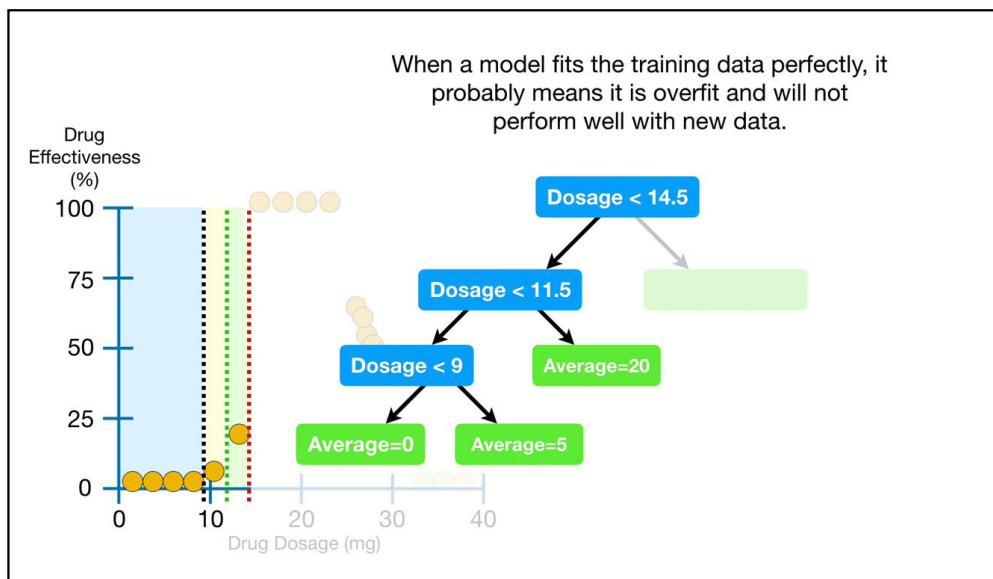
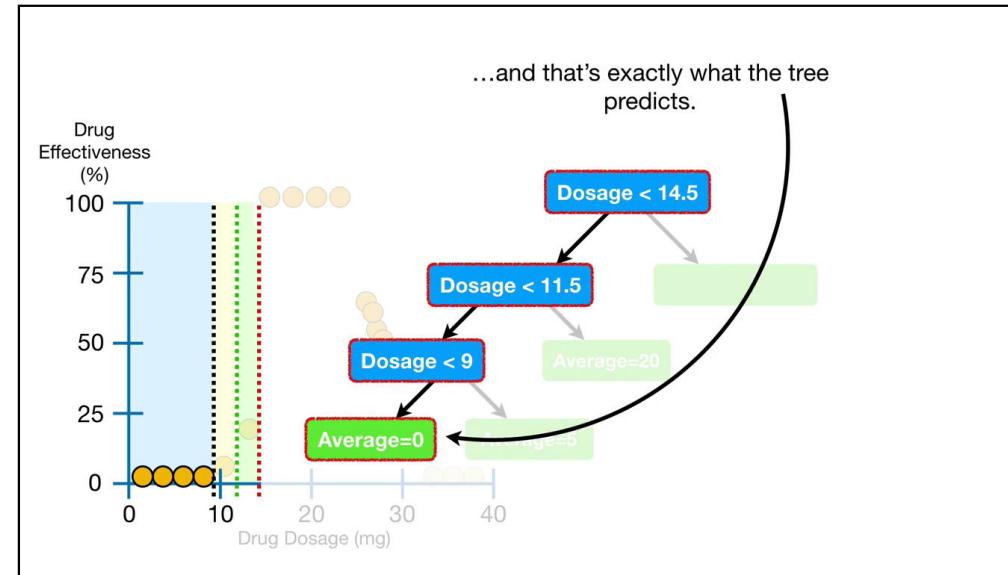
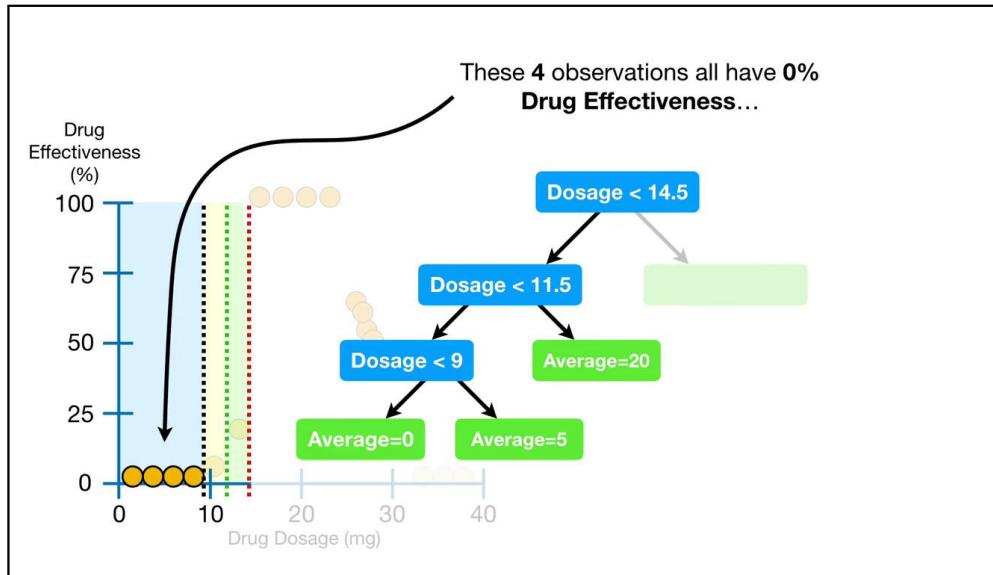




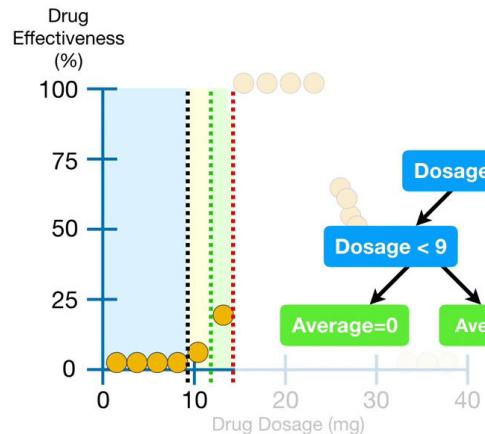




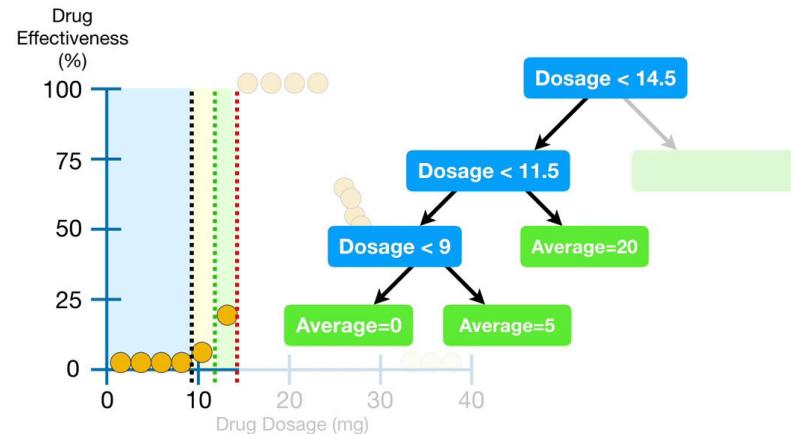




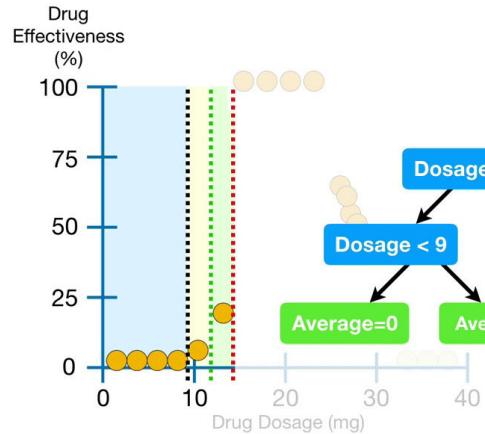
Is there a way to prevent our tree from overfitting the training data?



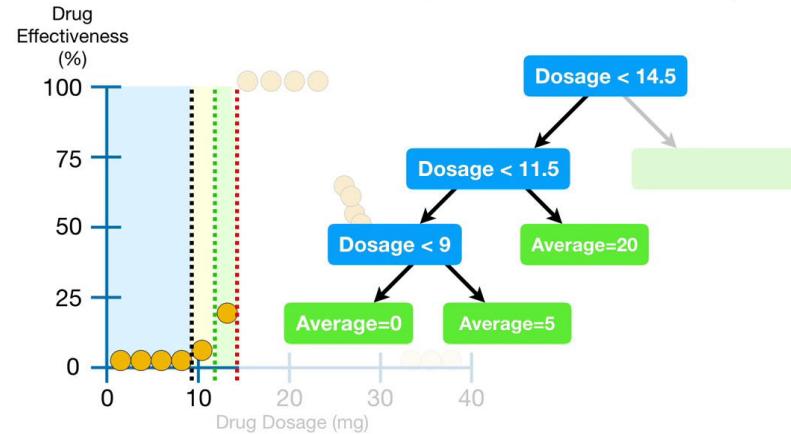
Yes, there are a bunch of techniques.

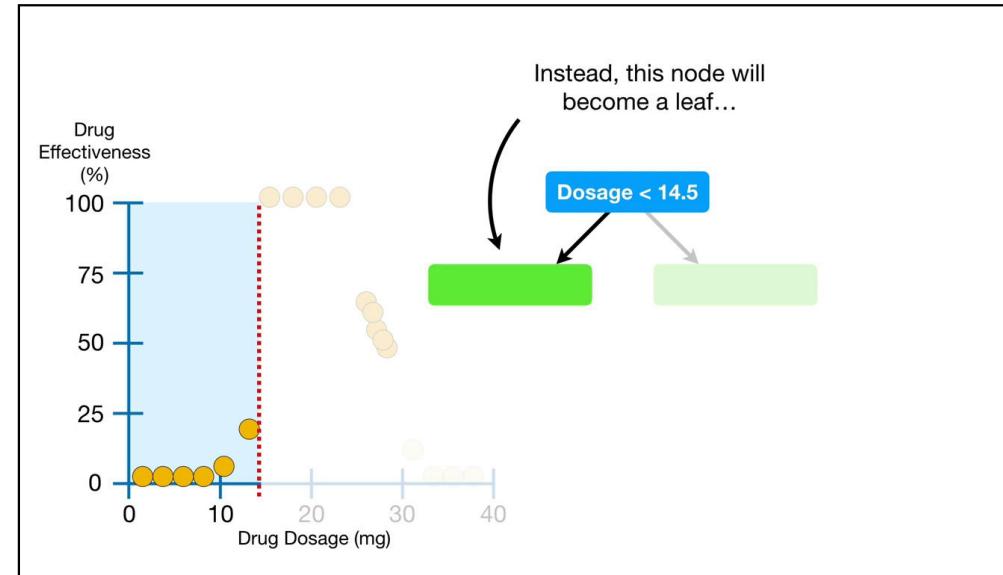
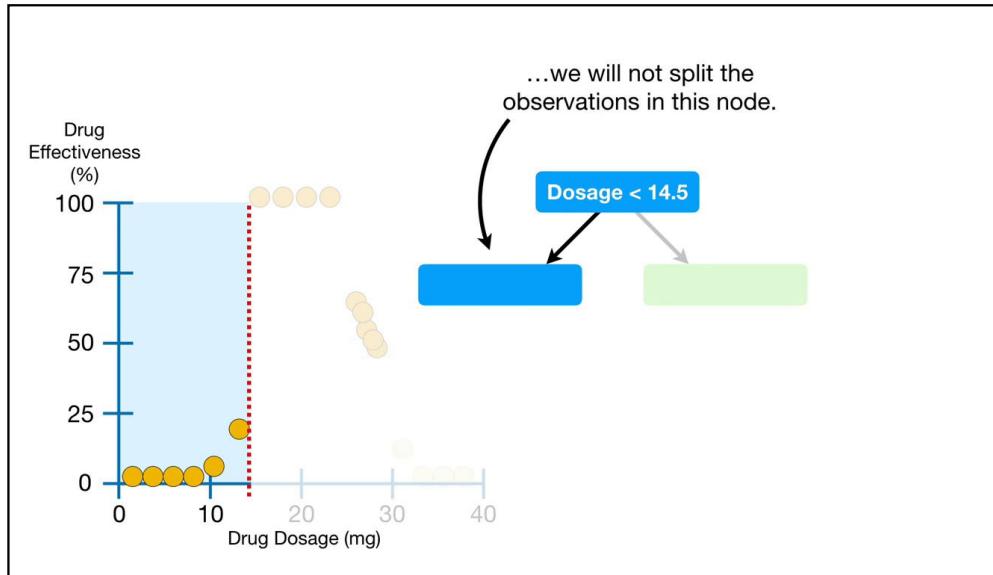
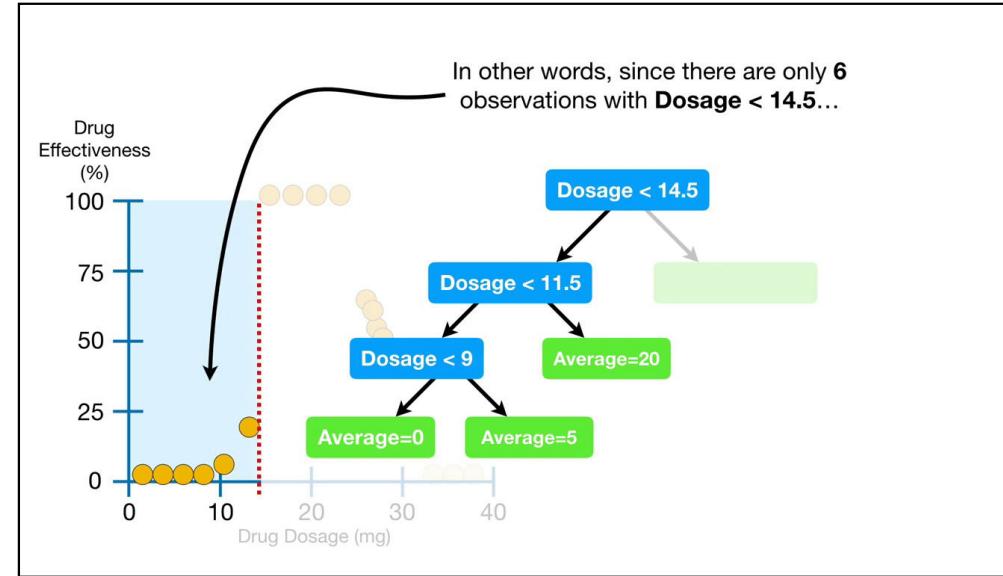
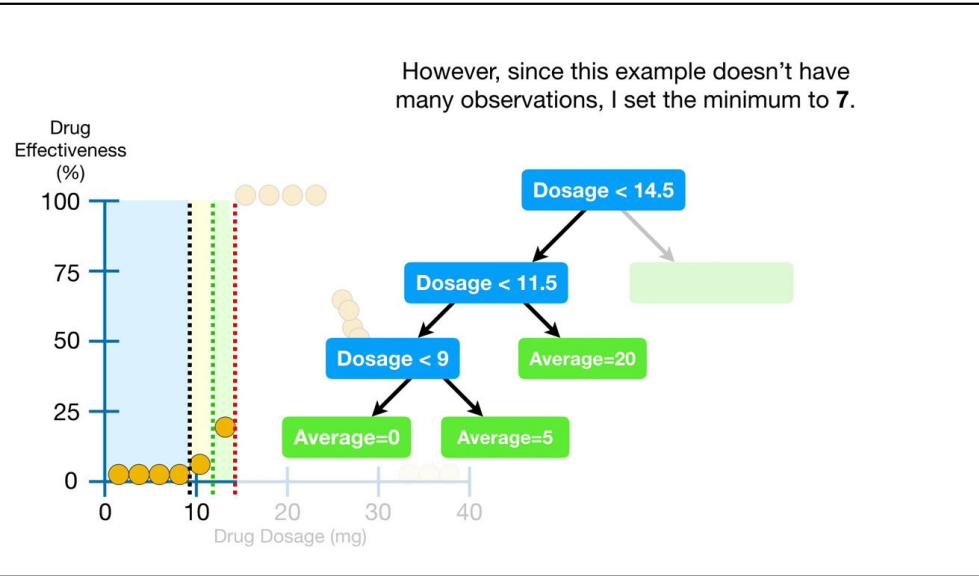


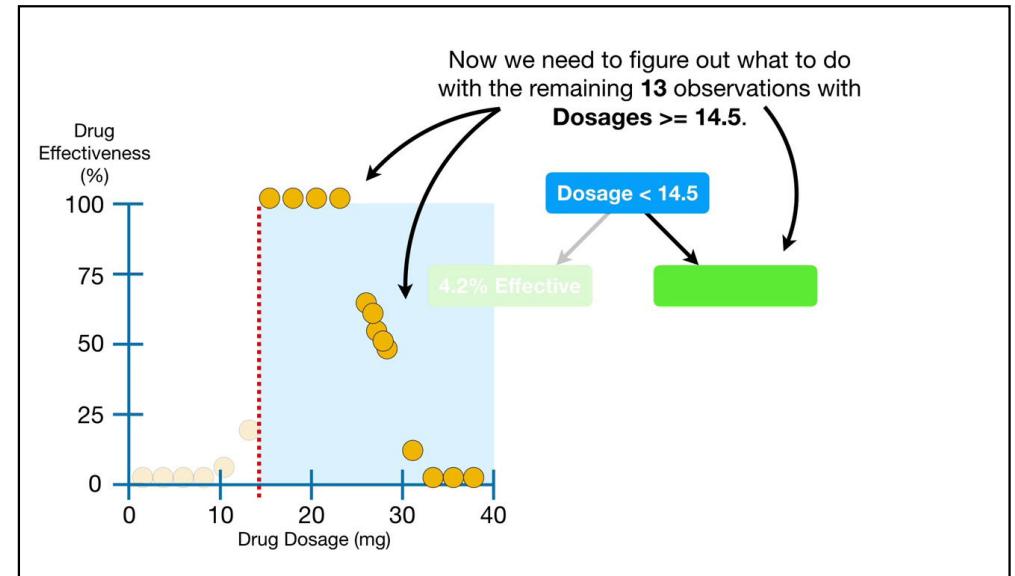
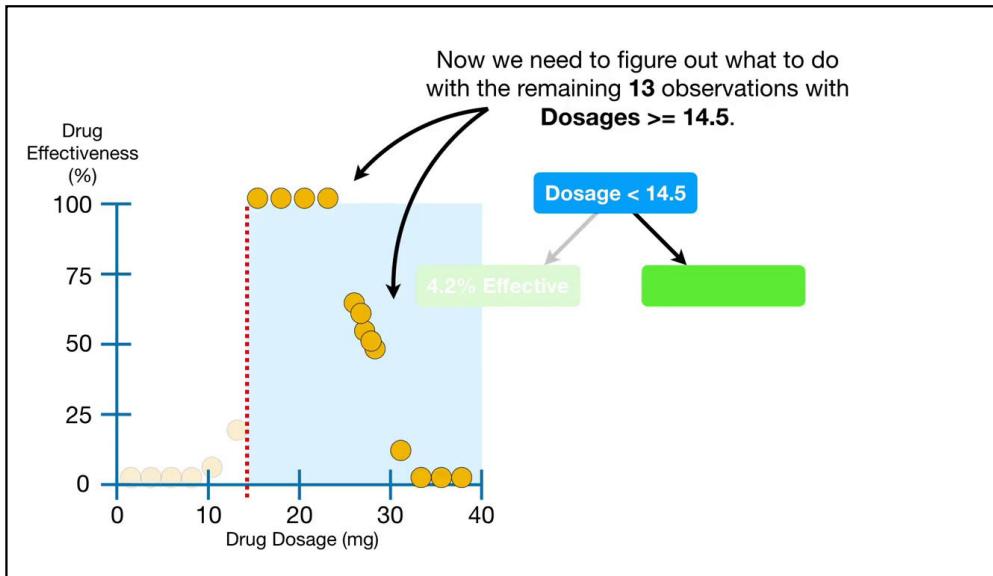
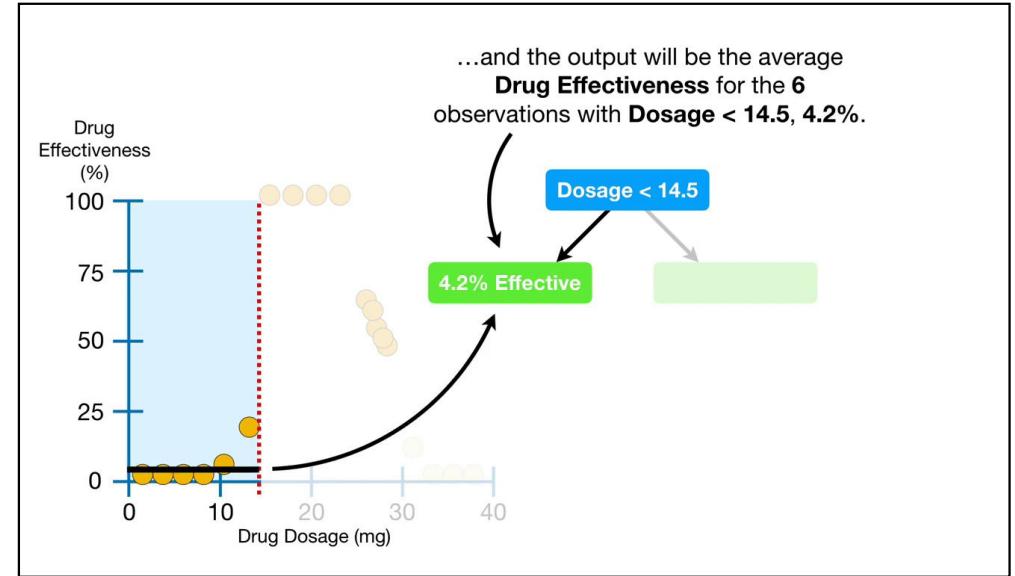
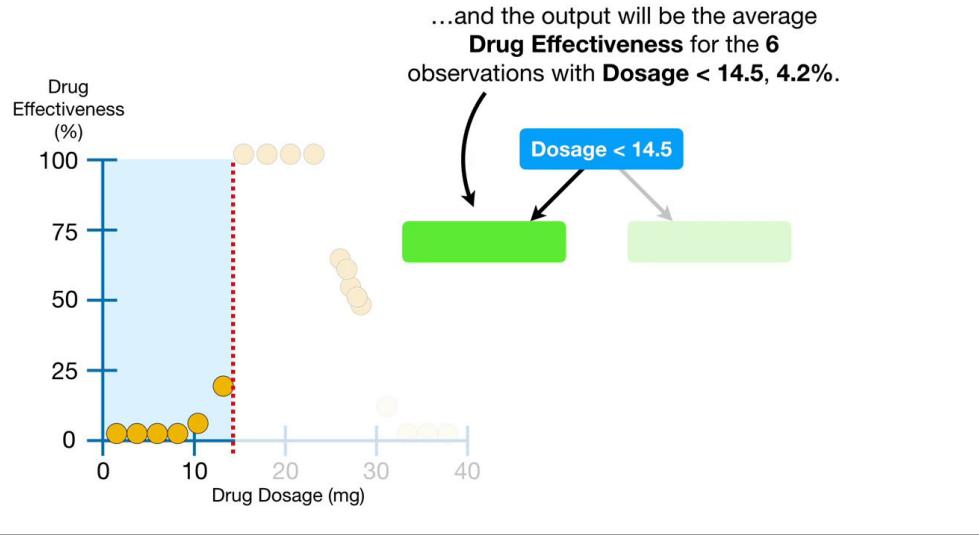
The simplest is to only split observations when there are more than some minimum number.

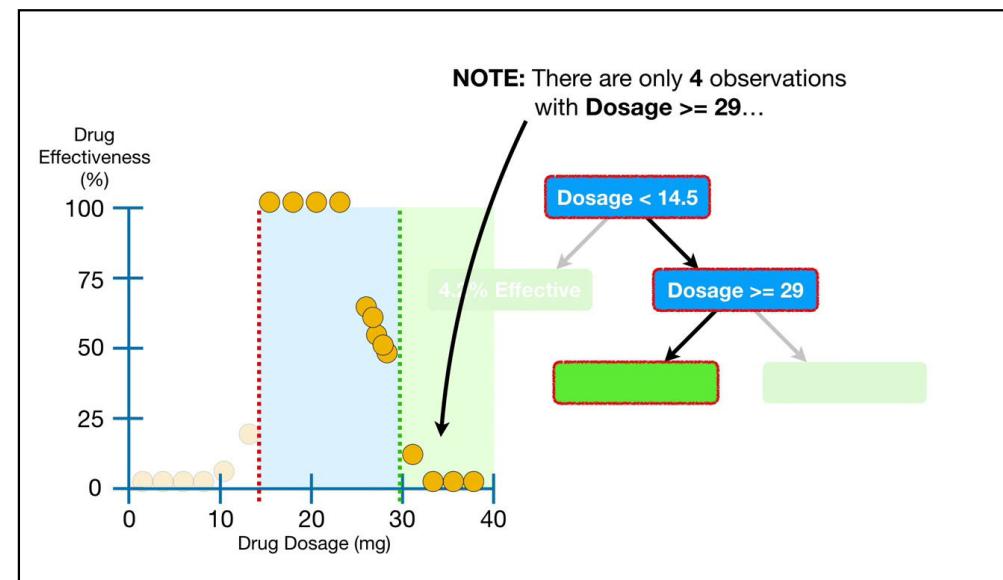
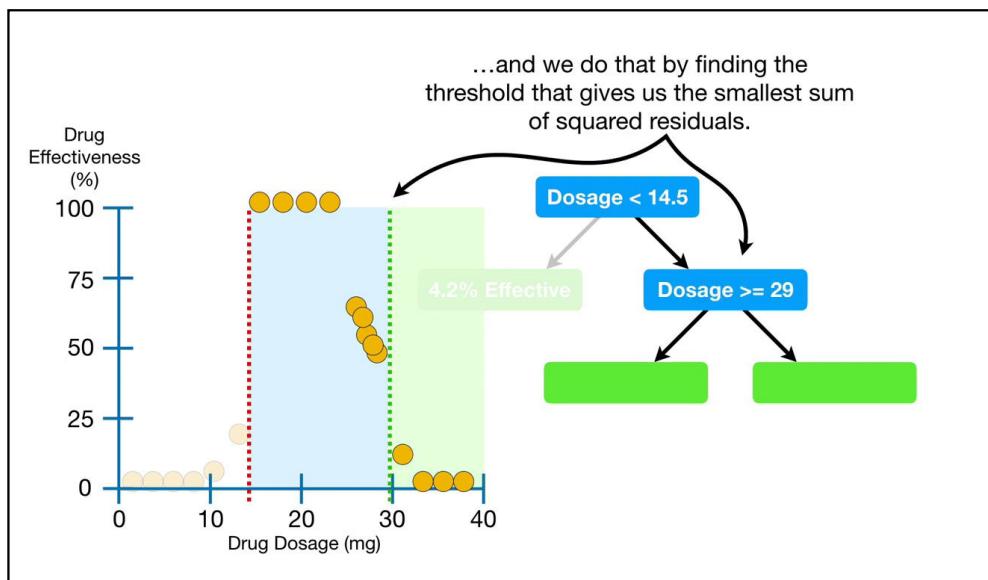
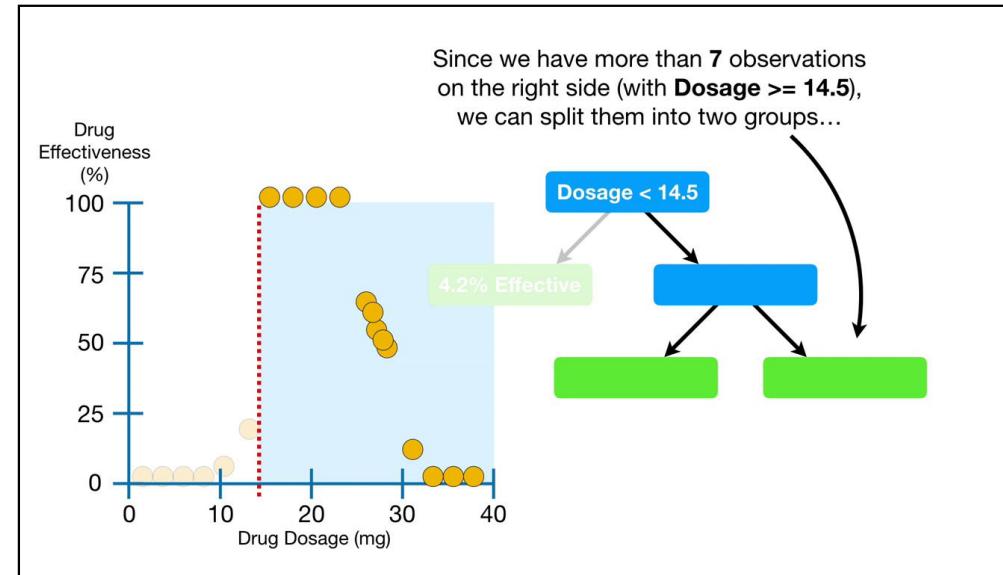
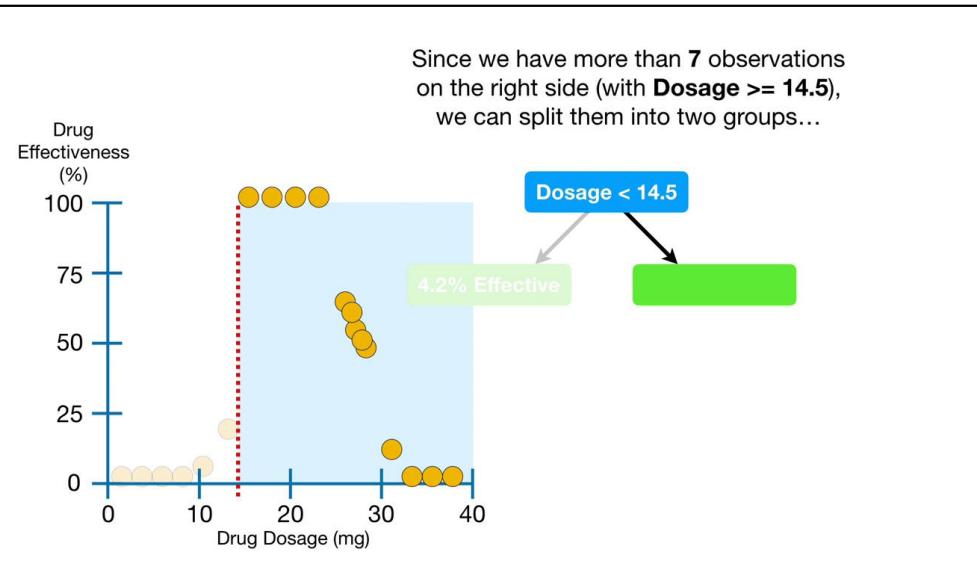


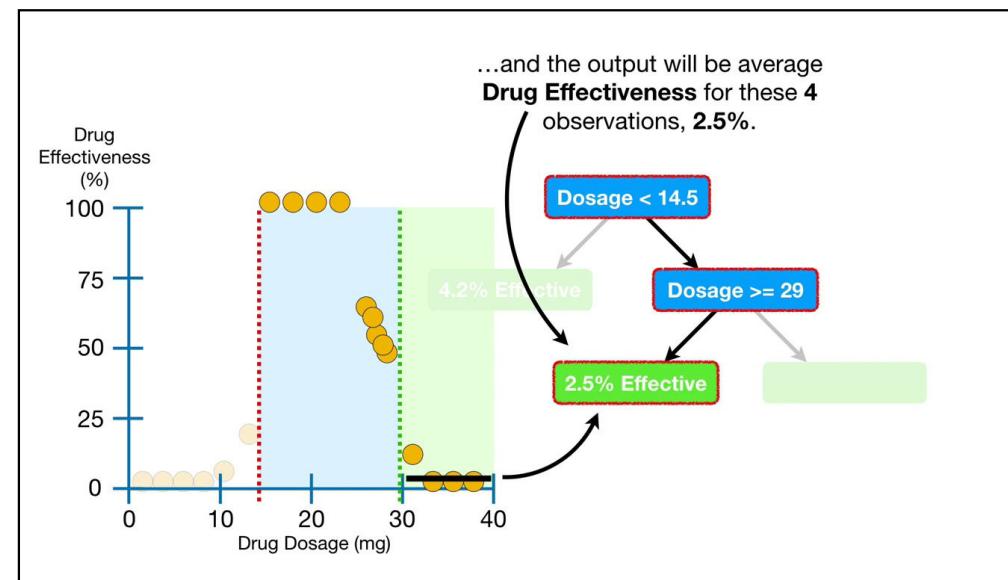
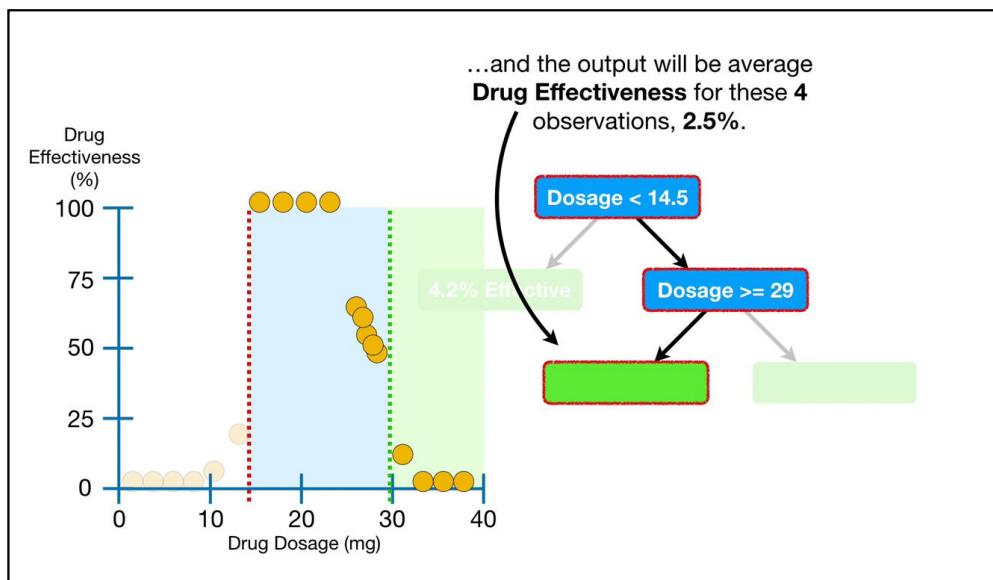
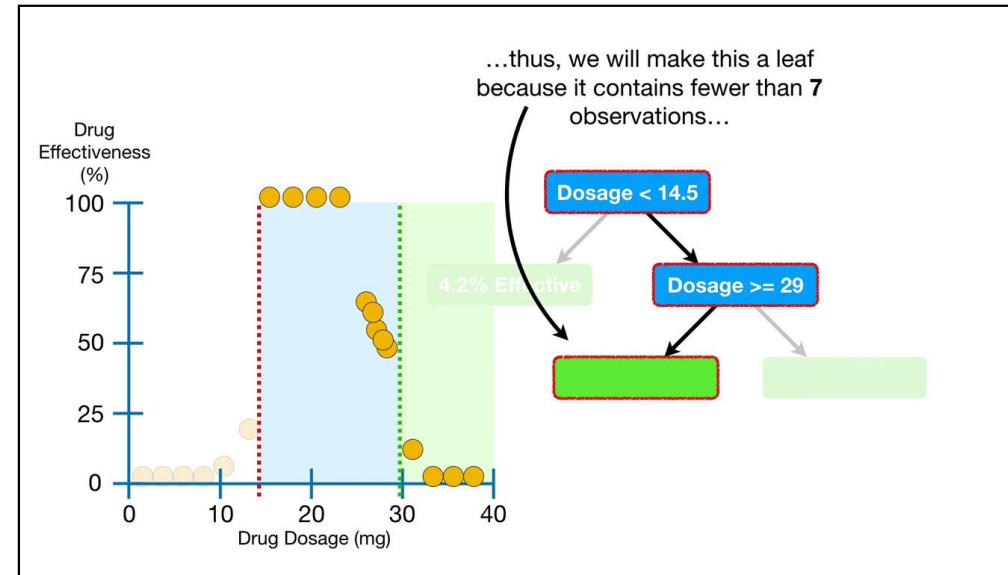
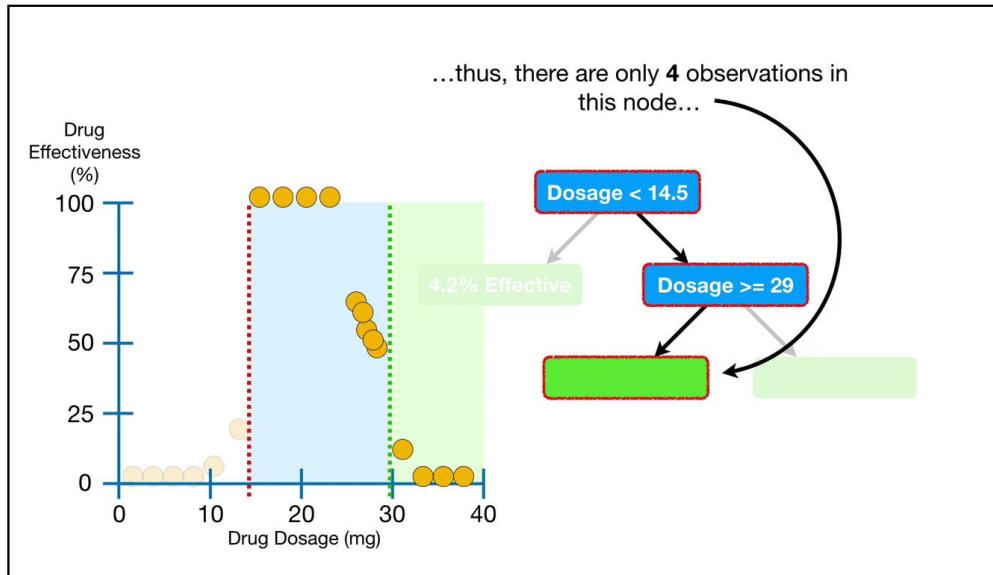
Typically, the minimum number of observations to allow for a split is **20**.

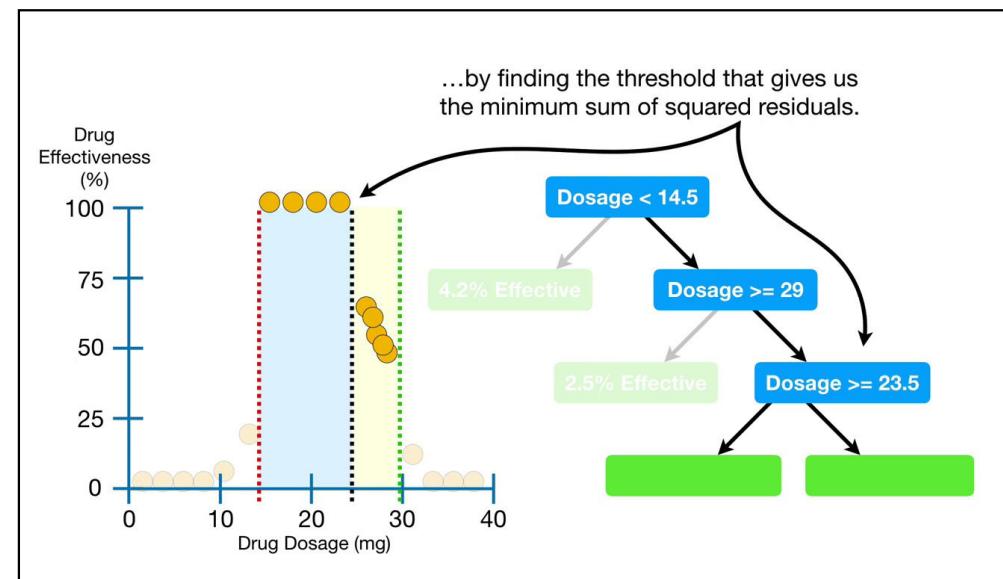
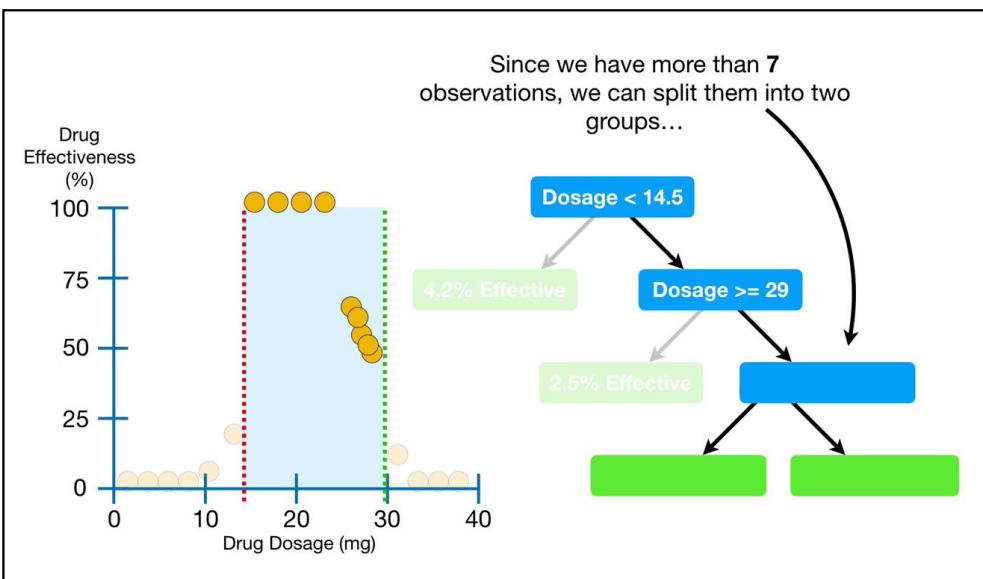
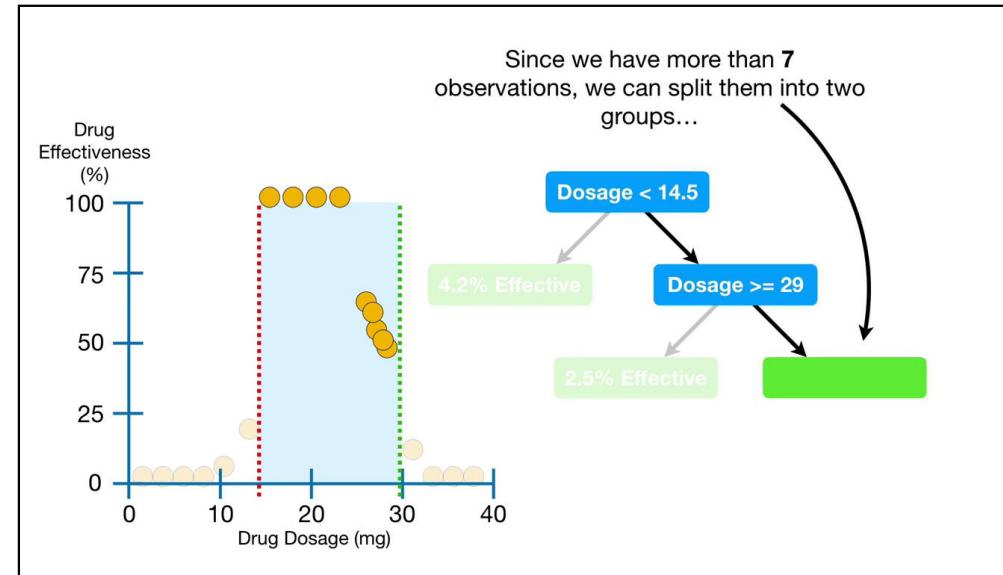
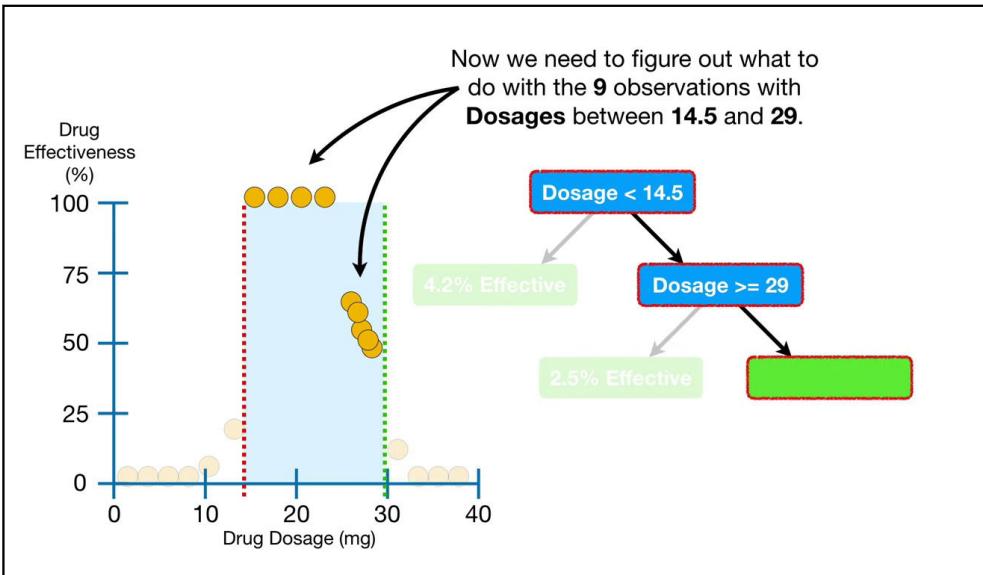


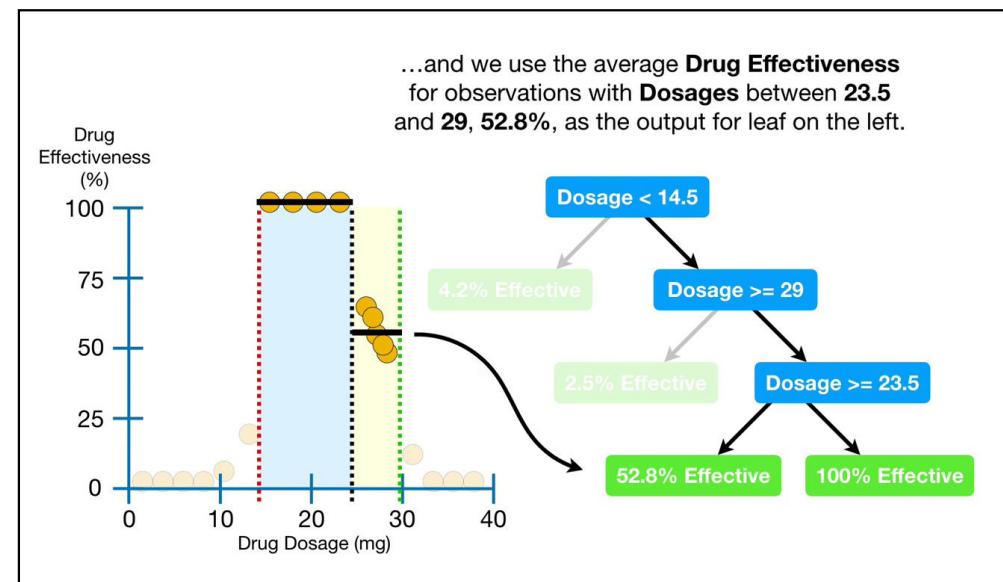
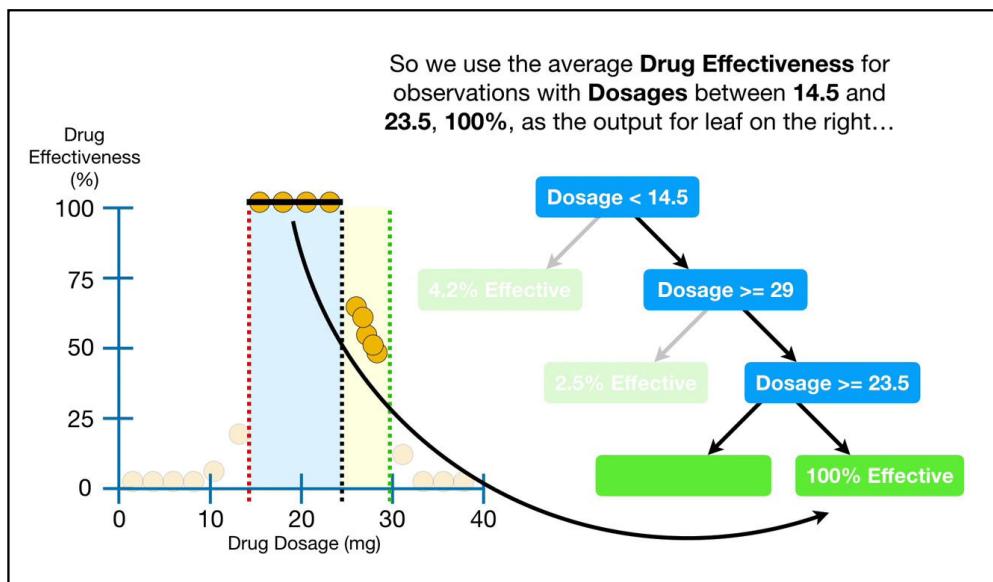
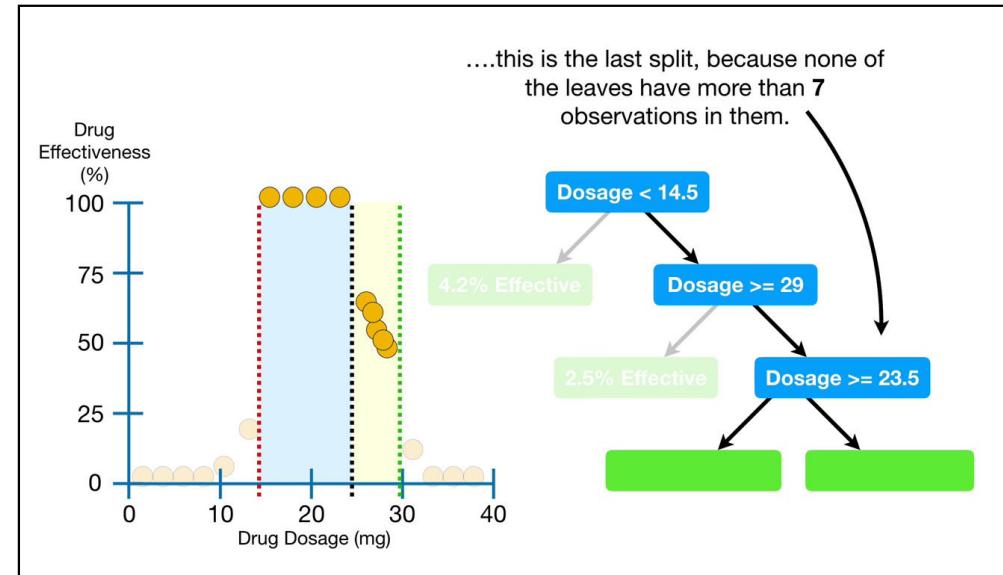
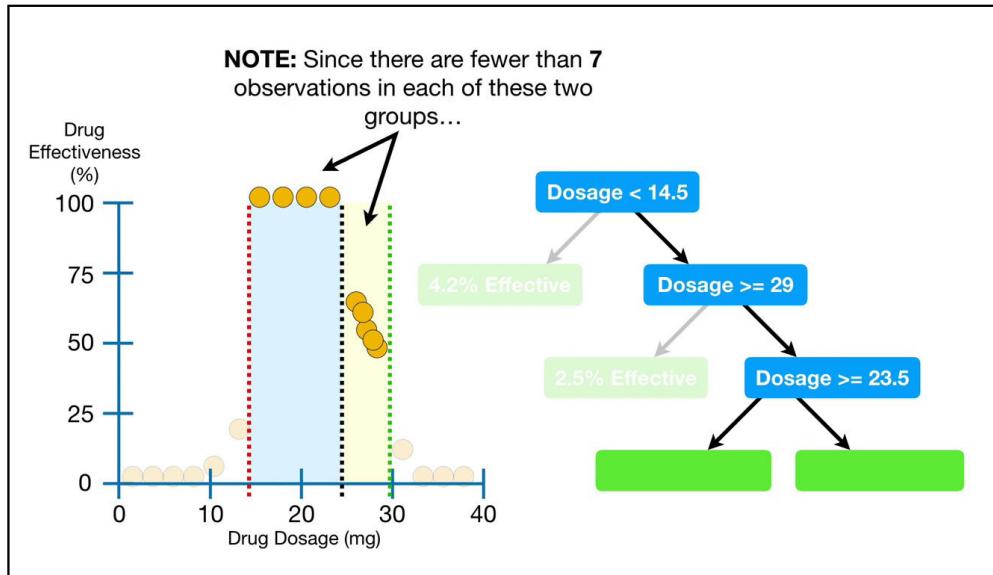


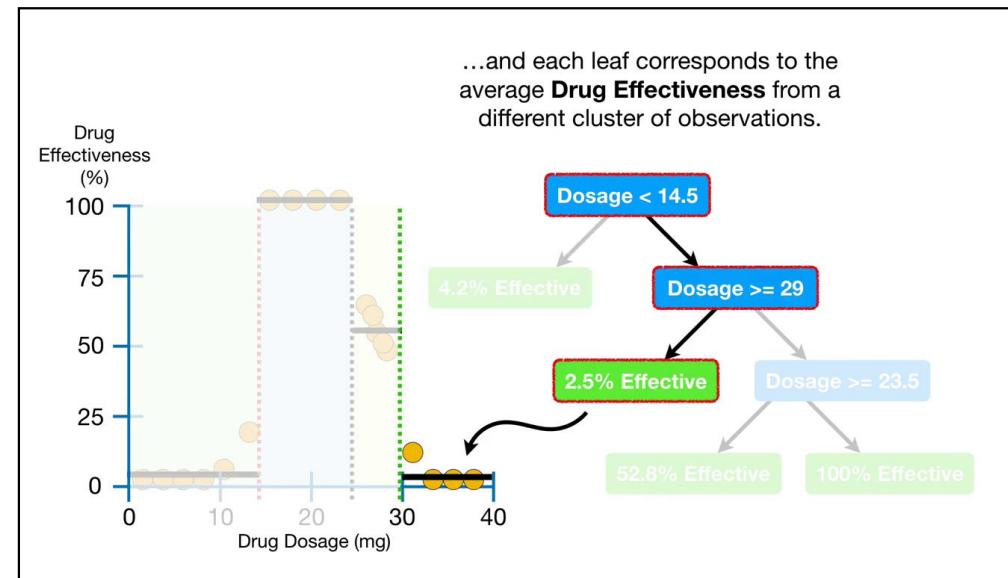
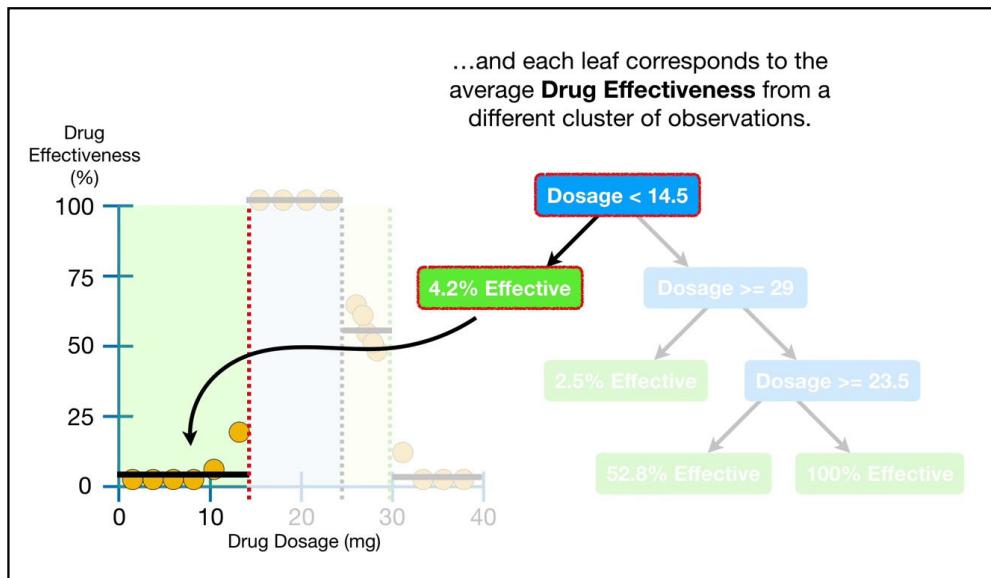
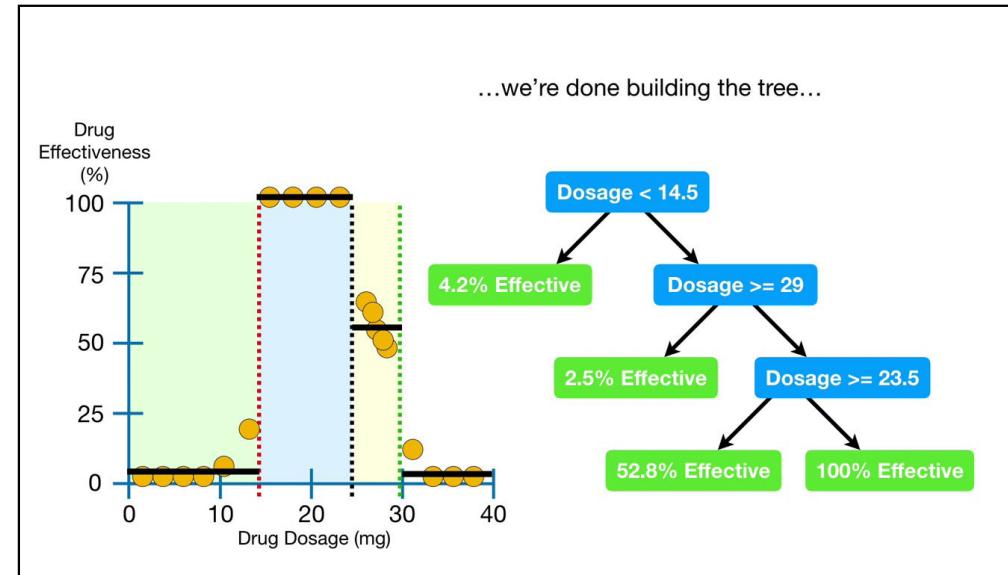
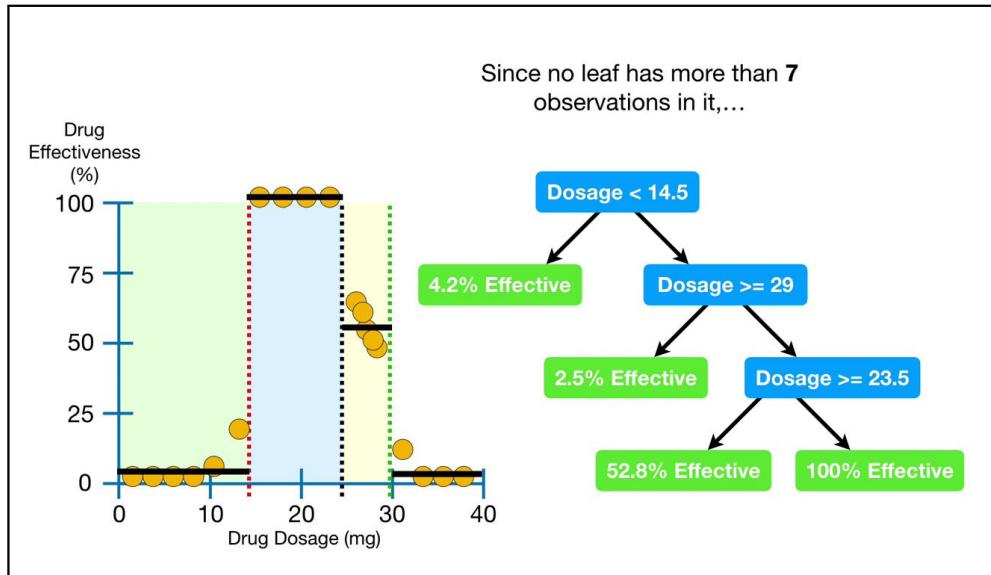


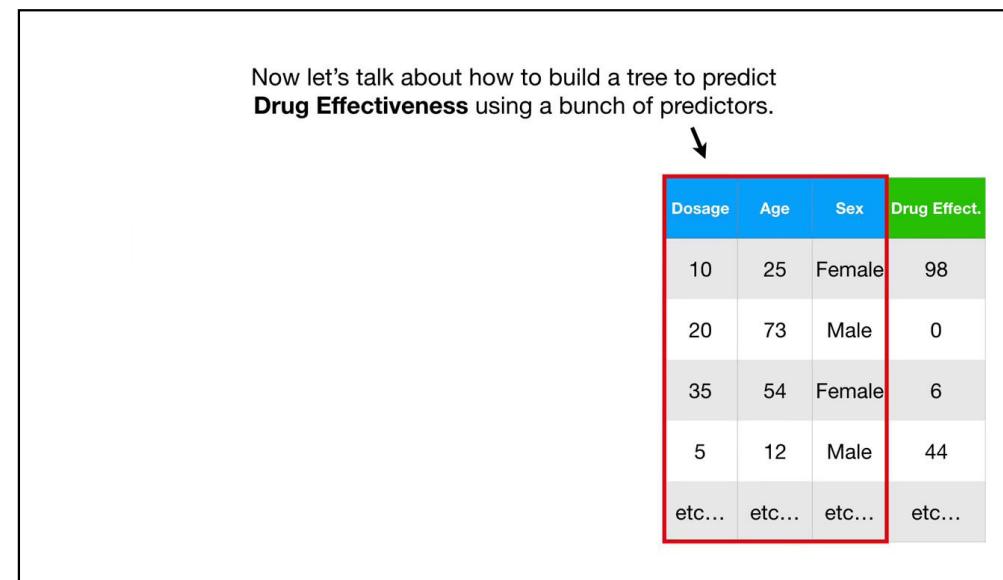
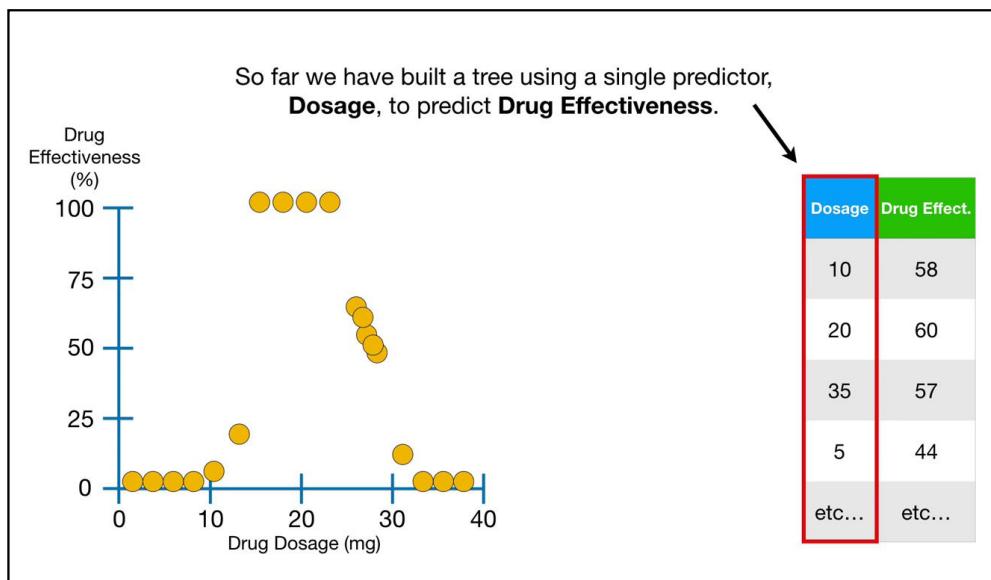
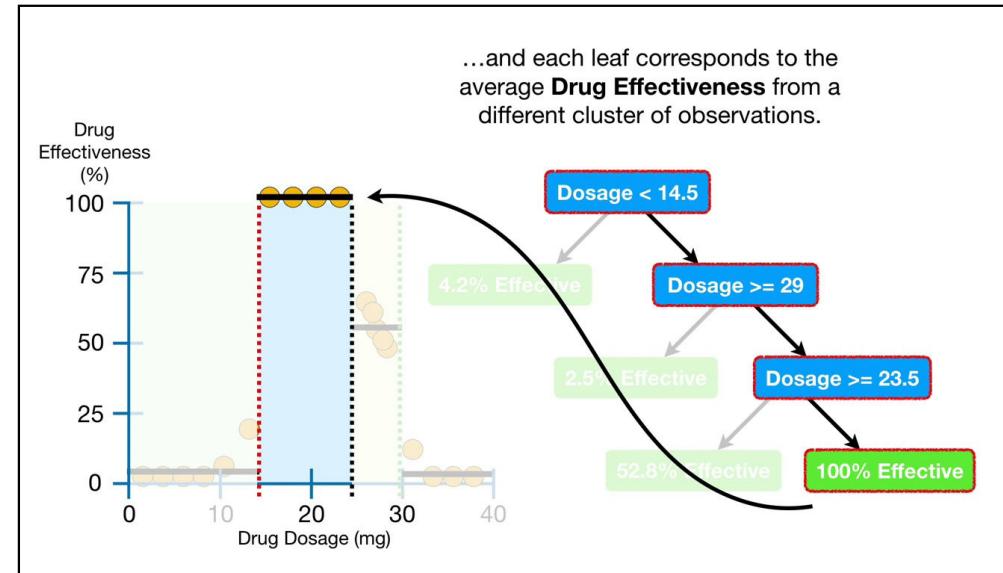
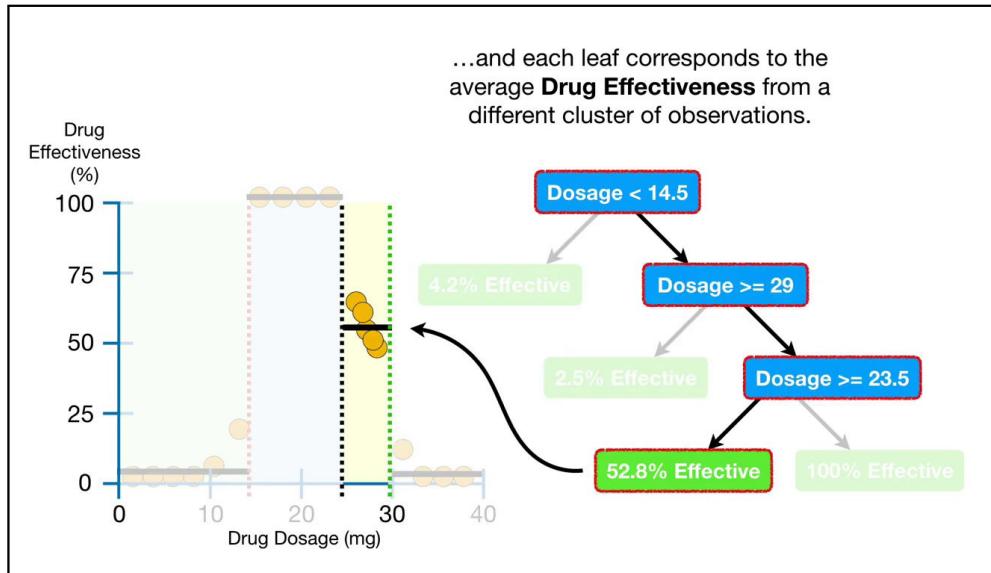












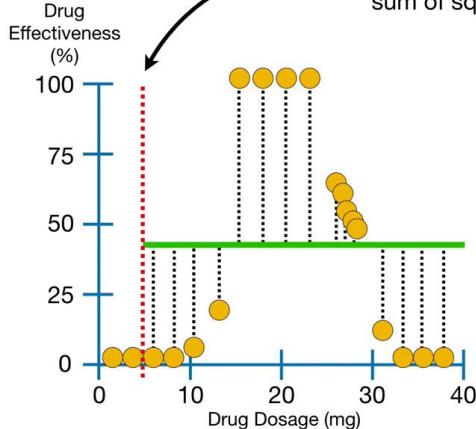
Just like before, we will start by using **Dosage** to predict **Drug Effectiveness**.

Dosage	Age	Sex	Drug Effect.
10	25	Female	98
20	73	Male	0
35	54	Female	6
5	12	Male	44
etc...	etc...	etc...	etc...

Just like before, we will start by using **Dosage** to predict **Drug Effectiveness**.

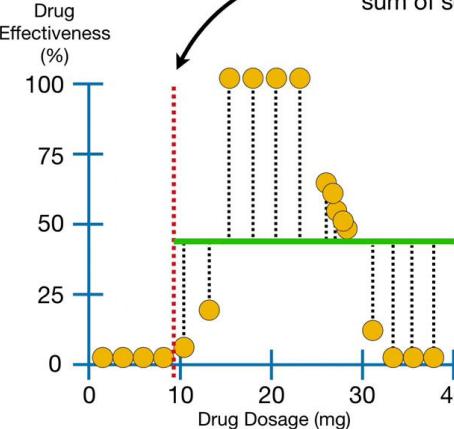
Dosage	Age	Sex	Drug Effect.
10	25	Female	98
20	73	Male	0
35	54	Female	6
5	12	Male	44
etc...	etc...	etc...	etc...

Thus, just like before, we will try different thresholds for **Dosage** and calculate the sum of squared residuals at each step...

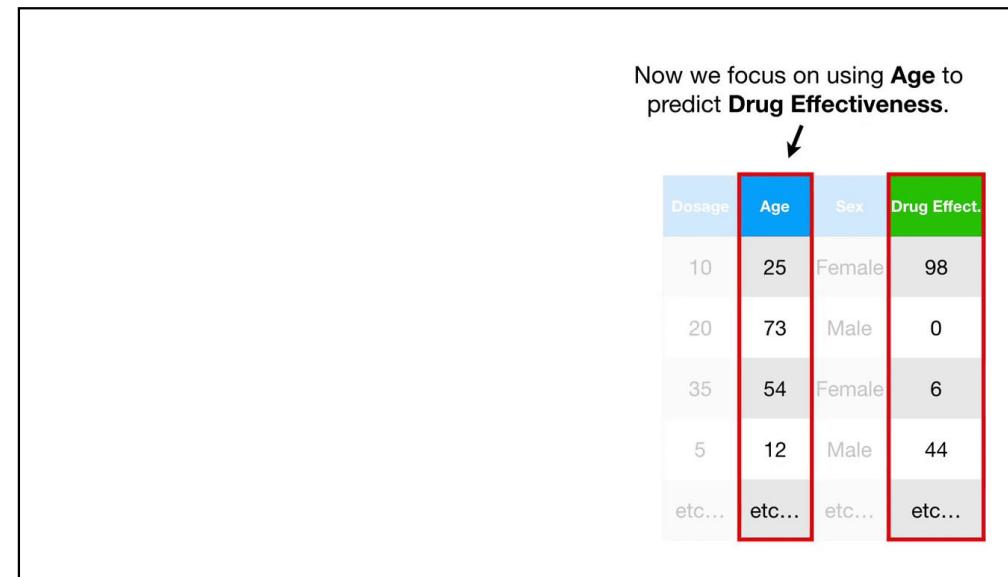
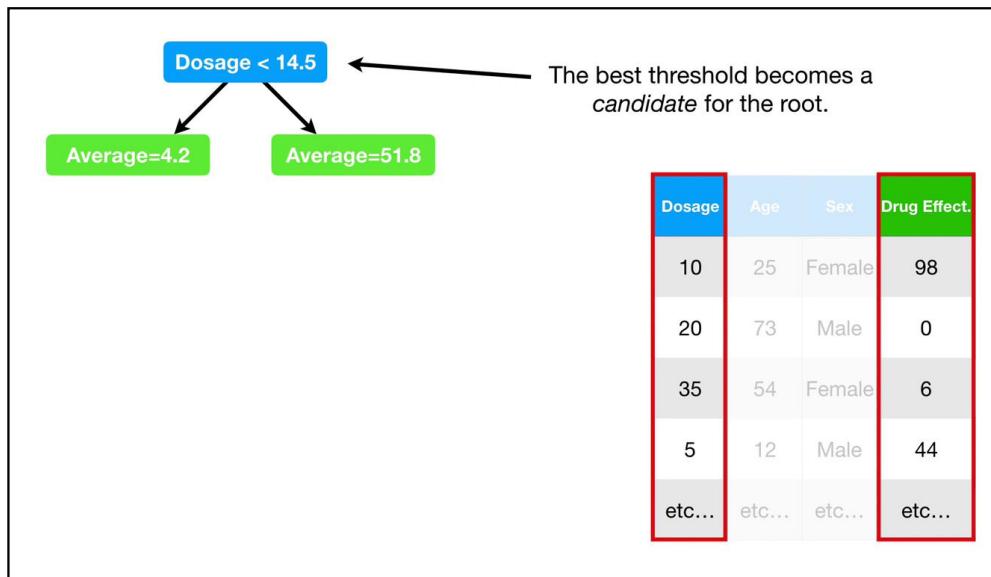
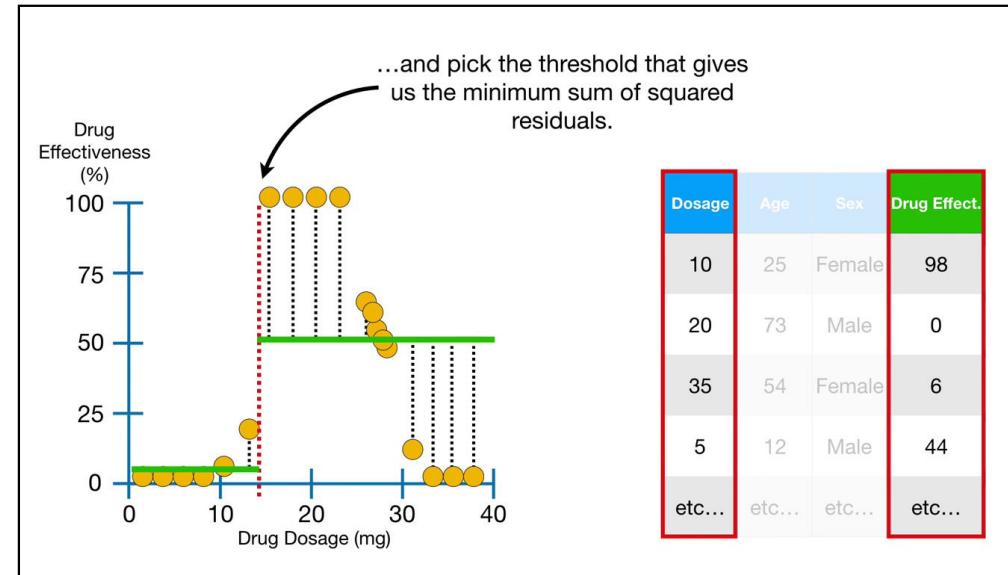
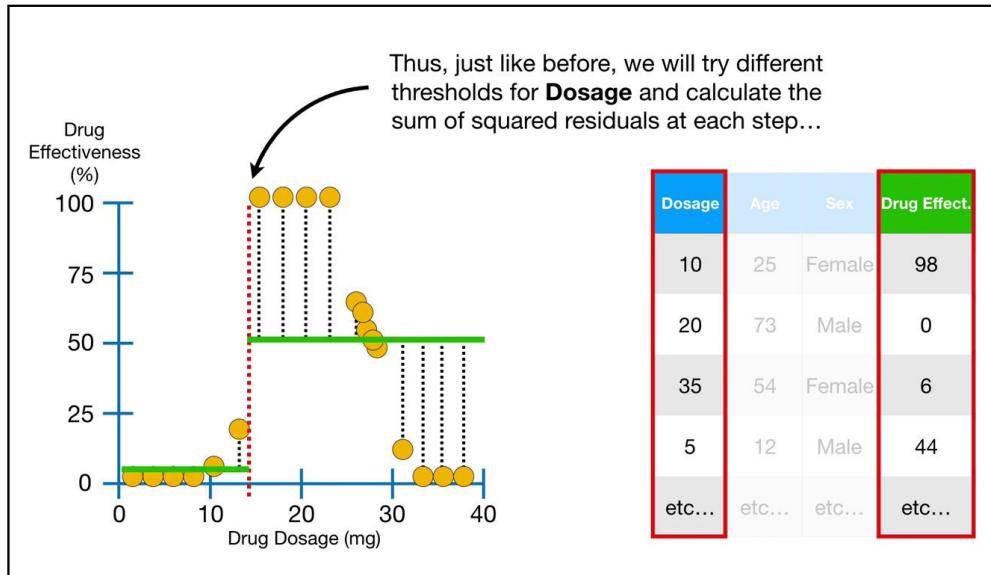


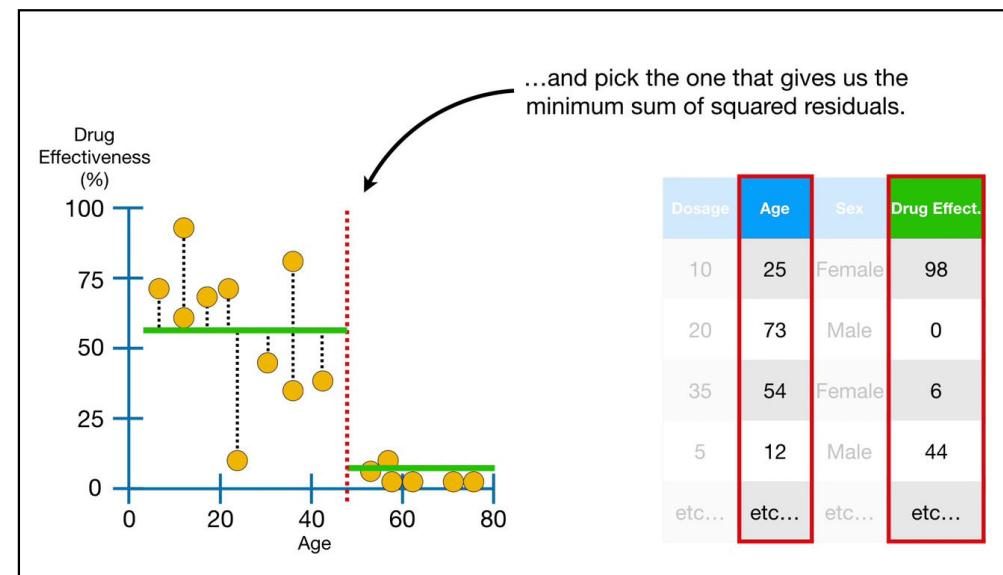
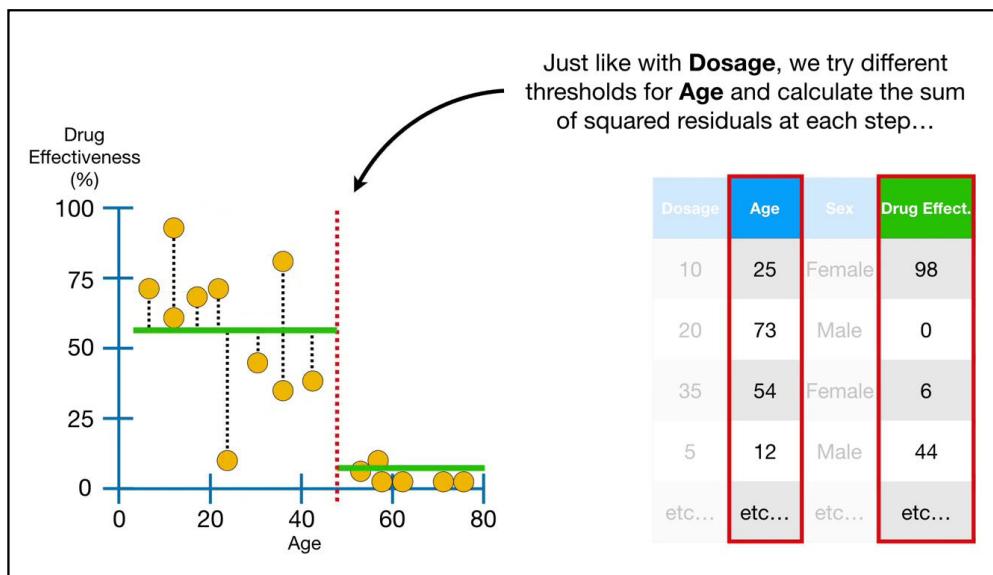
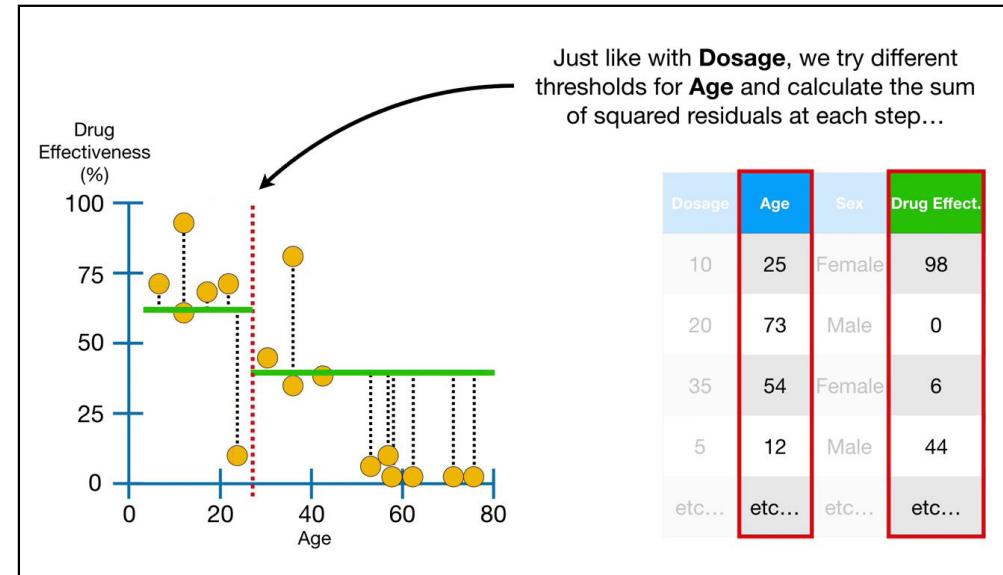
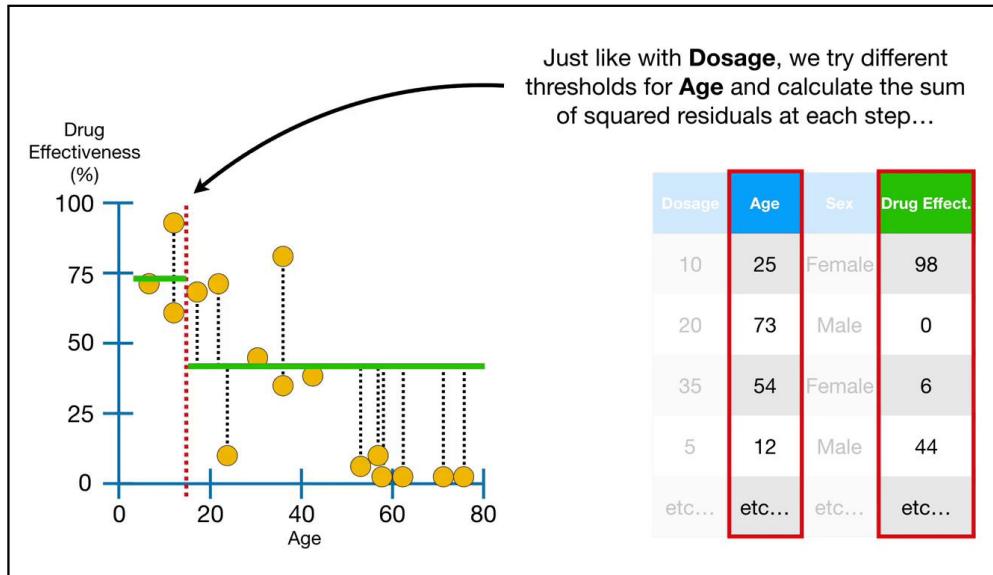
Dosage	Age	Sex	Drug Effect.
10	25	Female	98
20	73	Male	0
35	54	Female	6
5	12	Male	44
etc...	etc...	etc...	etc...

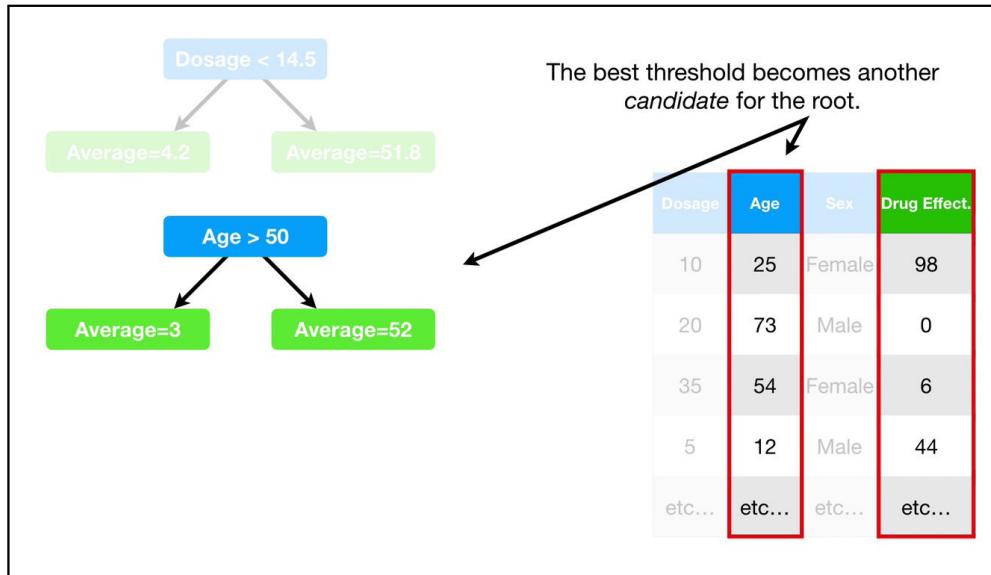
Thus, just like before, we will try different thresholds for **Dosage** and calculate the sum of squared residuals at each step...



Dosage	Age	Sex	Drug Effect.
10	25	Female	98
20	73	Male	0
35	54	Female	6
5	12	Male	44
etc...	etc...	etc...	etc...

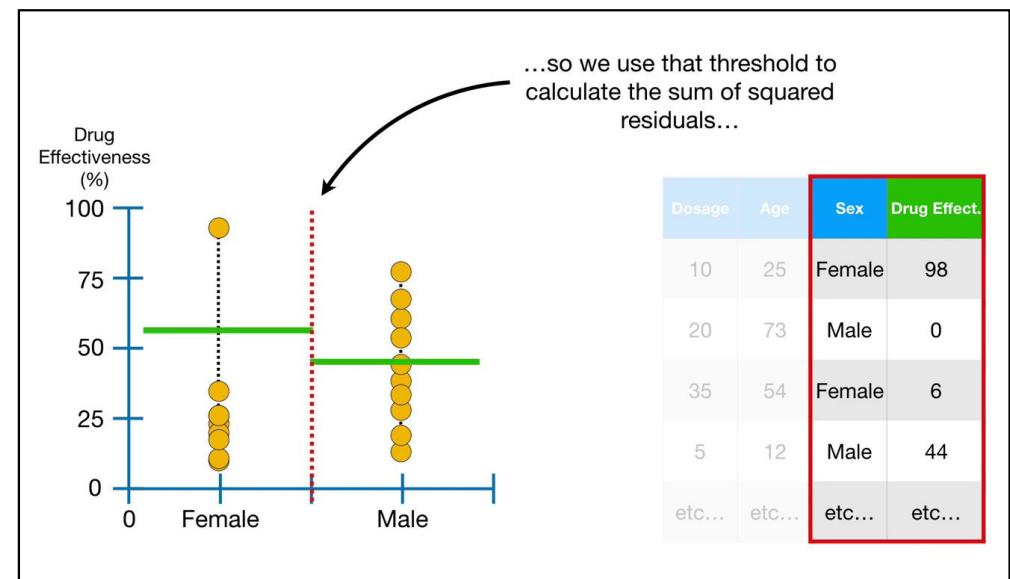
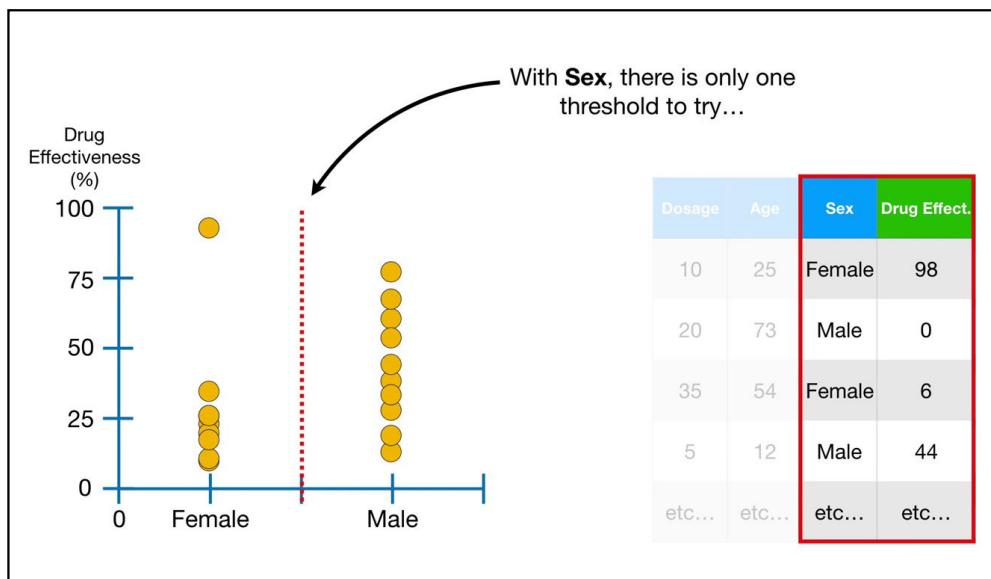


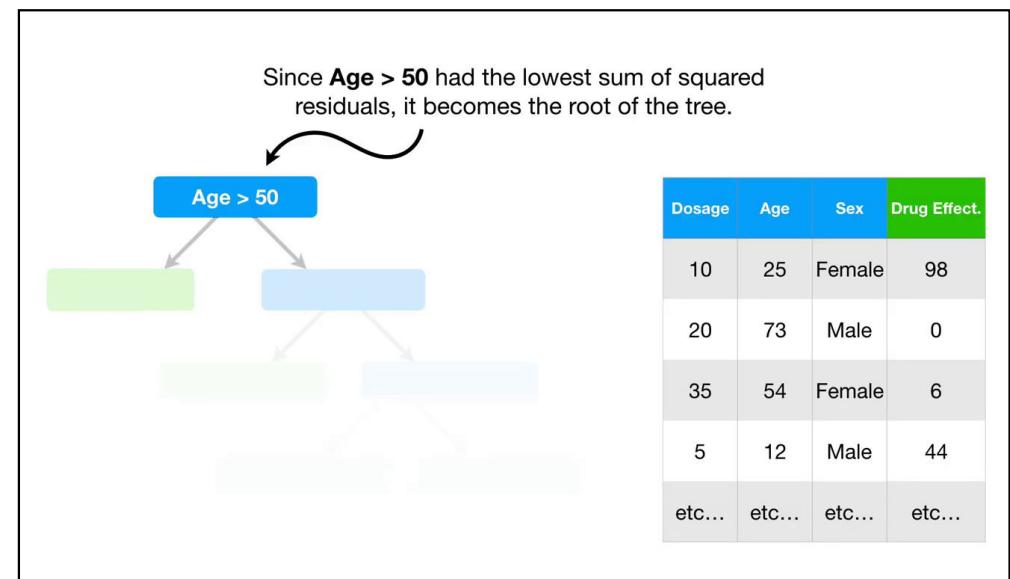
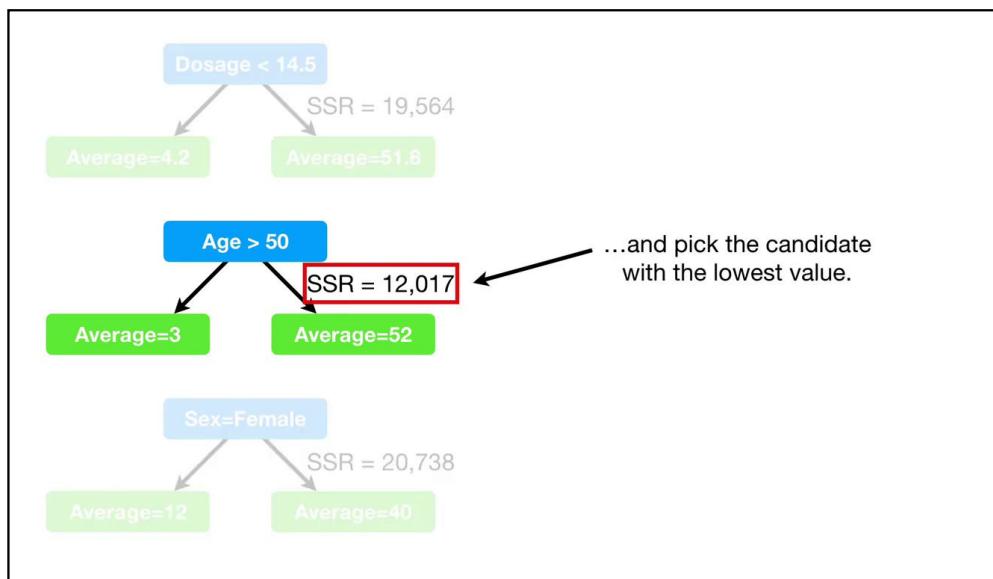
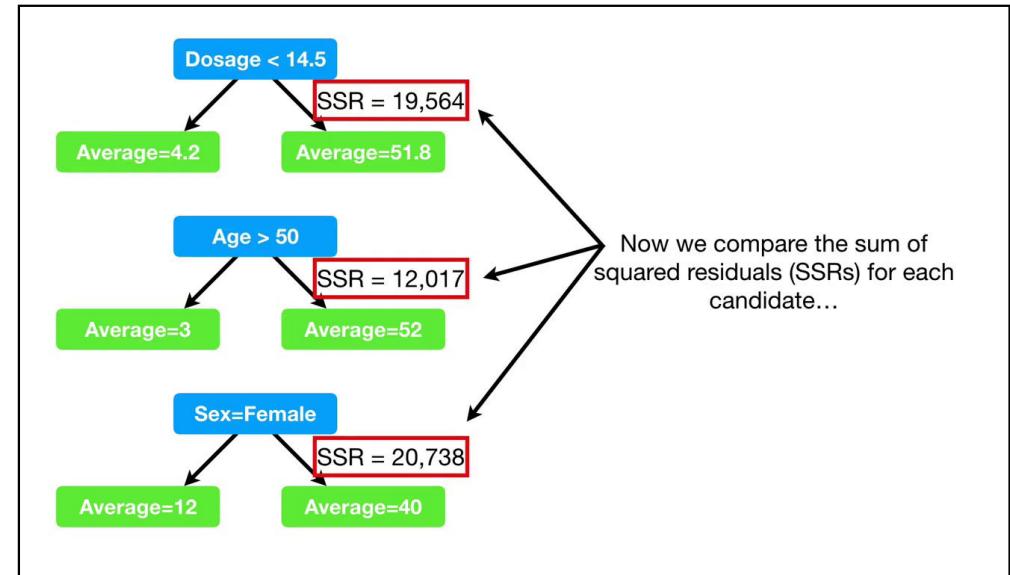
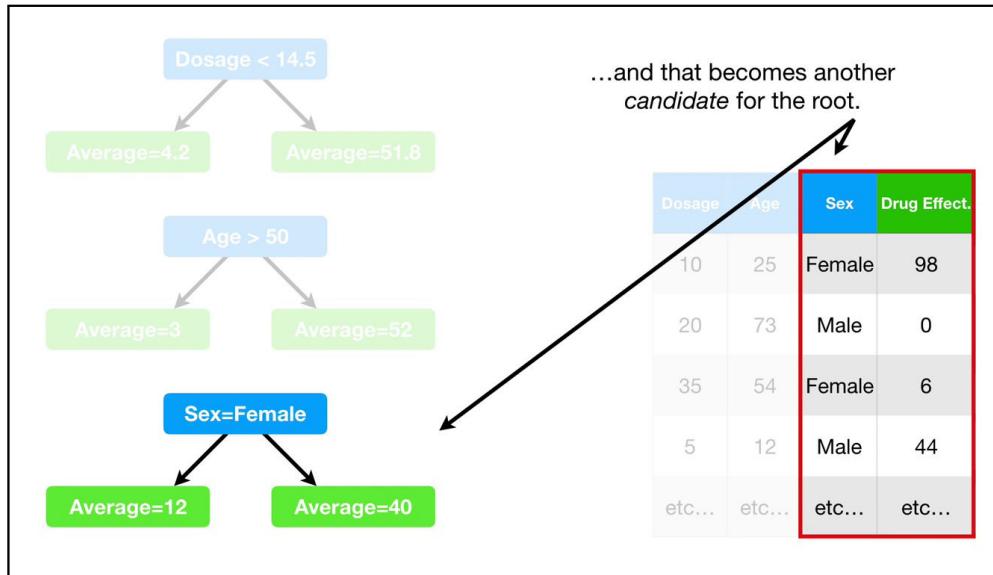




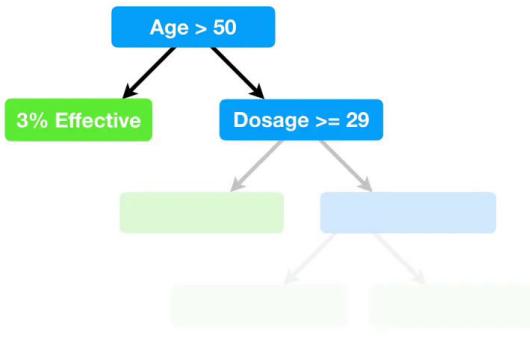
Now we focus on using **Sex** to predict **Drug Effectiveness**.

Dosage	Age	Sex	Drug Effect.
10	25	Female	98
20	73	Male	0
35	54	Female	6
5	12	Male	44
etc...	etc...	etc...	etc...



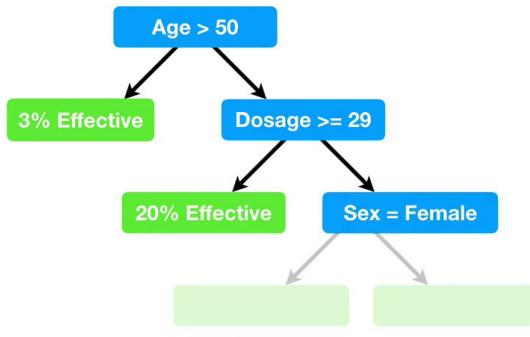


Then we grow the tree just like before, except now we compare the lowest sum of squared residuals from each predictor.



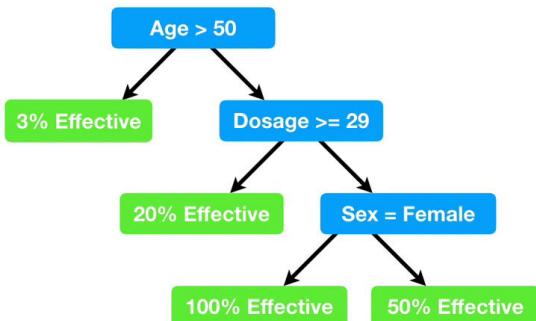
Dosage	Age	Sex	Drug Effect.
10	25	Female	98
20	73	Male	0
35	54	Female	6
5	12	Male	44
etc...	etc...	etc...	etc...

Then we grow the tree just like before, except now we compare the lowest sum of squared residuals from each predictor.



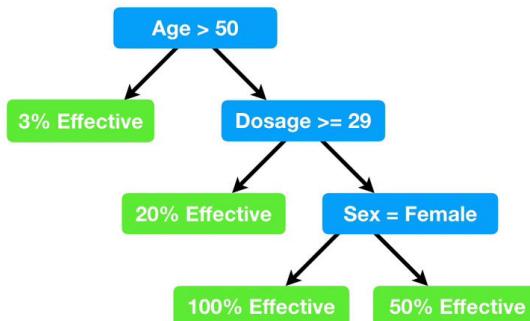
Dosage	Age	Sex	Drug Effect.
10	25	Female	98
20	73	Male	0
35	54	Female	6
5	12	Male	44
etc...	etc...	etc...	etc...

Then we grow the tree just like before, except now we compare the lowest sum of squared residuals from each predictor.



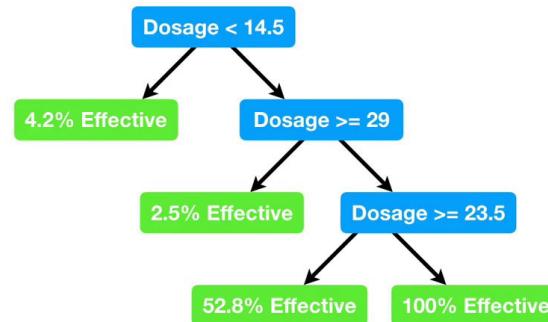
Dosage	Age	Sex	Drug Effect.
10	25	Female	98
20	73	Male	0
35	54	Female	6
5	12	Male	44
etc...	etc...	etc...	etc...

And just like before, when a leaf has less than a minimum number of observations, which is usually **20**, but we are using **7**, we stop trying to divide them.

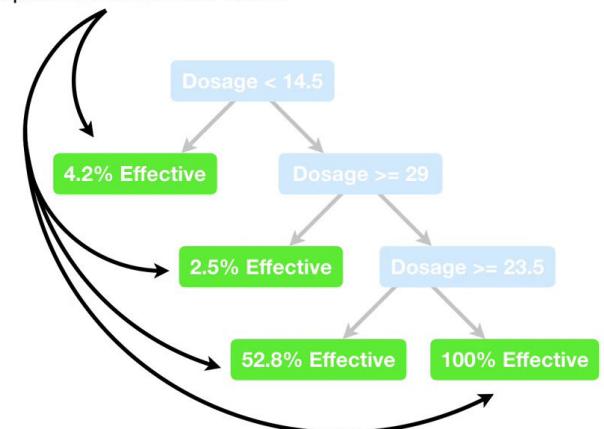


Dosage	Age	Sex	Drug Effect.
10	25	Female	98
20	73	Male	0
35	54	Female	6
5	12	Male	44
etc...	etc...	etc...	etc...

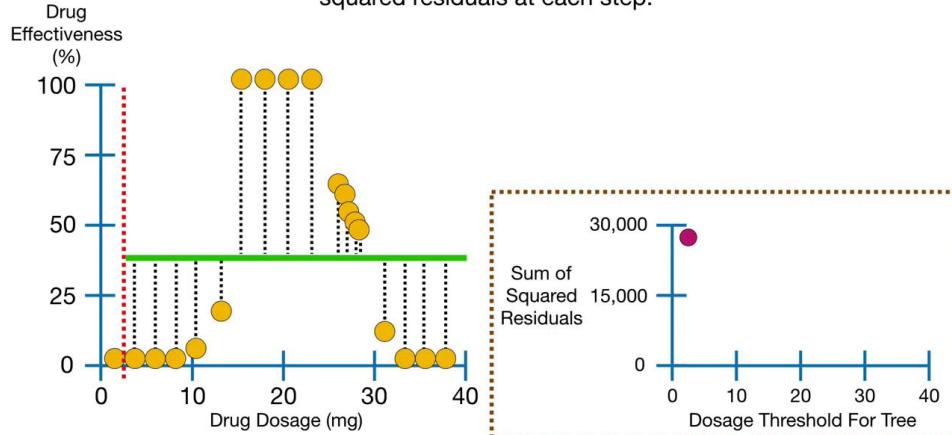
**Regression Trees** are a type of Decision Tree.



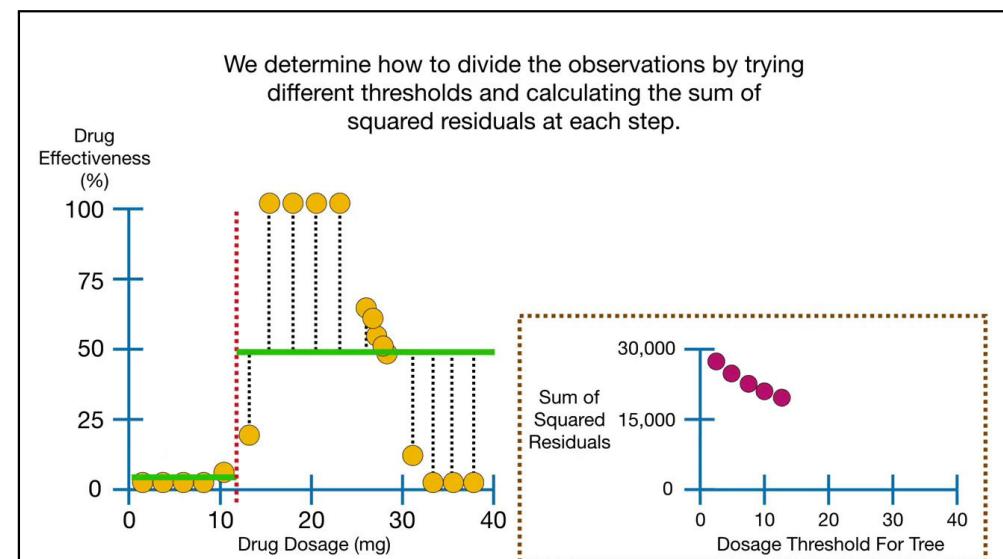
In a **Regression Tree**, each leaf represents a numeric value.



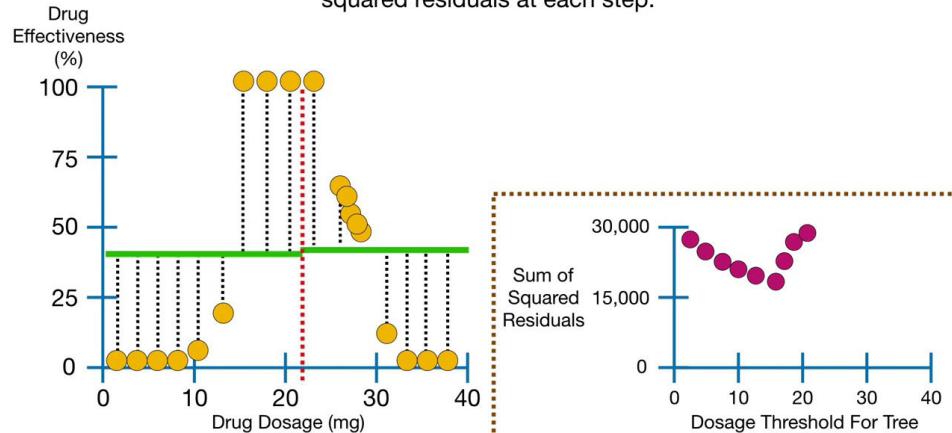
We determine how to divide the observations by trying different thresholds and calculating the sum of squared residuals at each step.



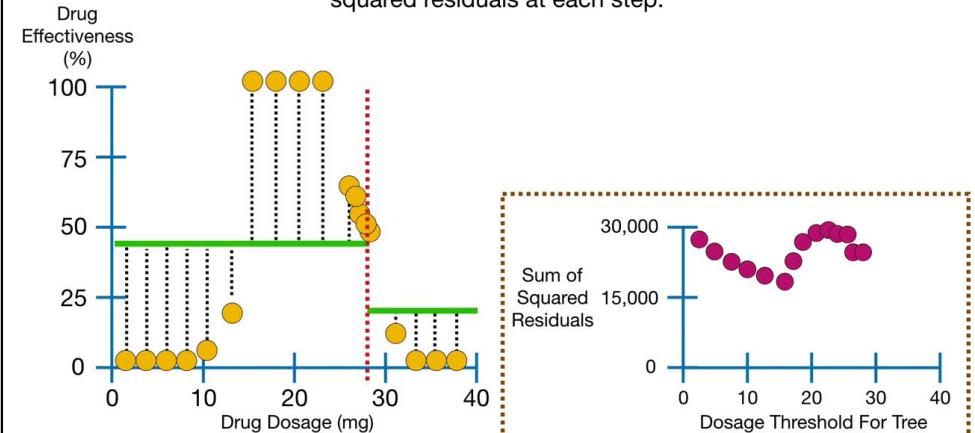
We determine how to divide the observations by trying different thresholds and calculating the sum of squared residuals at each step.



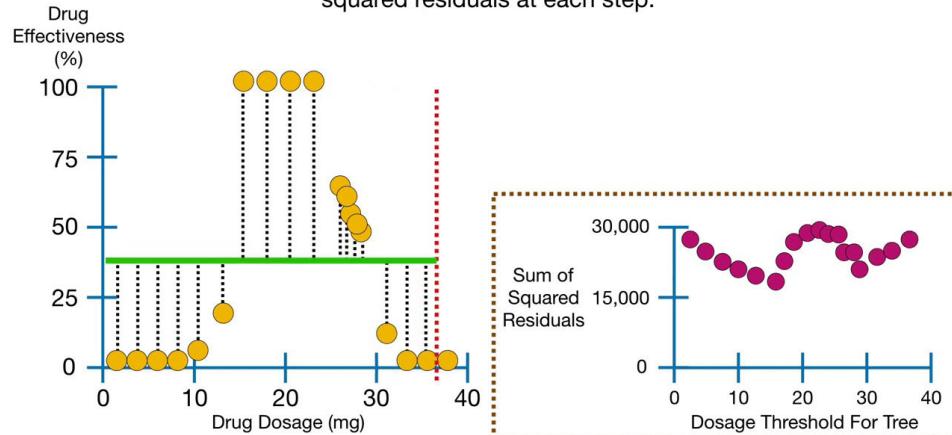
We determine how to divide the observations by trying different thresholds and calculating the sum of squared residuals at each step.



We determine how to divide the observations by trying different thresholds and calculating the sum of squared residuals at each step.



We determine how to divide the observations by trying different thresholds and calculating the sum of squared residuals at each step.



The threshold with the smallest sum of squared residuals...

