

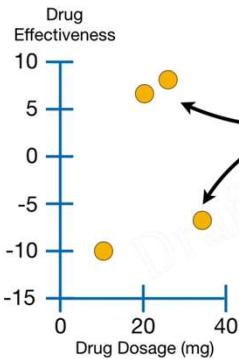
BOOSTING

- AdaBoost
- Gradient Boost
- XGBoost

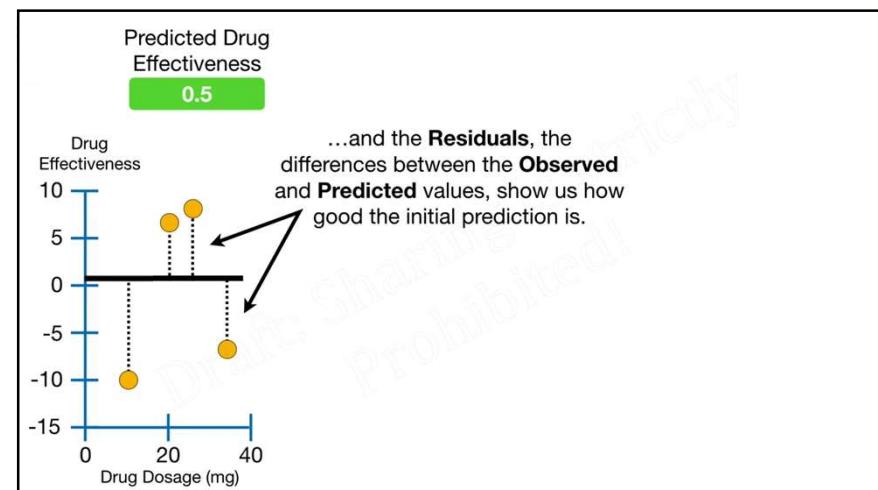
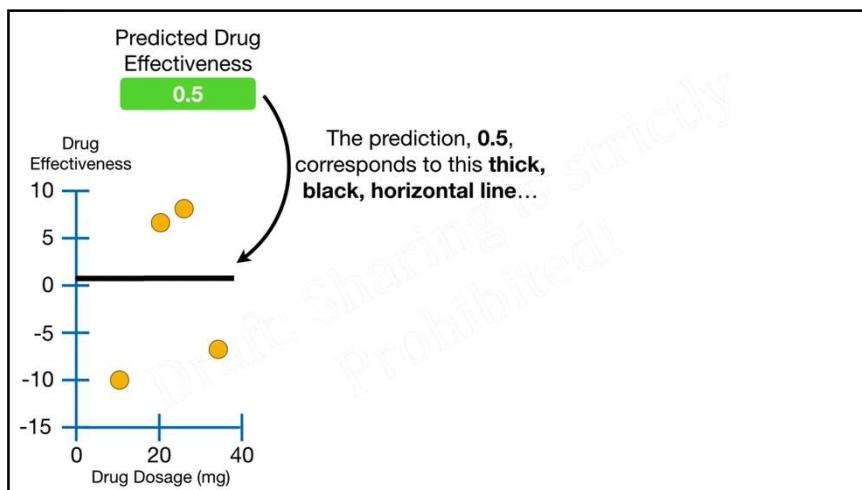
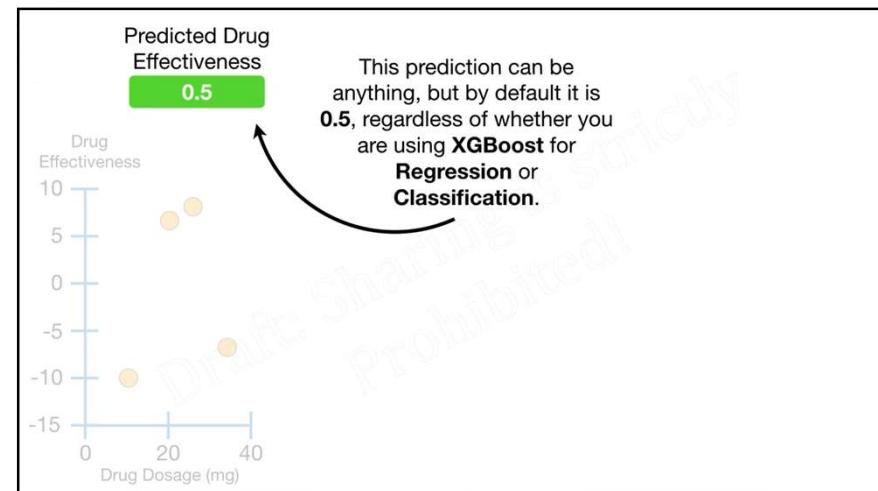
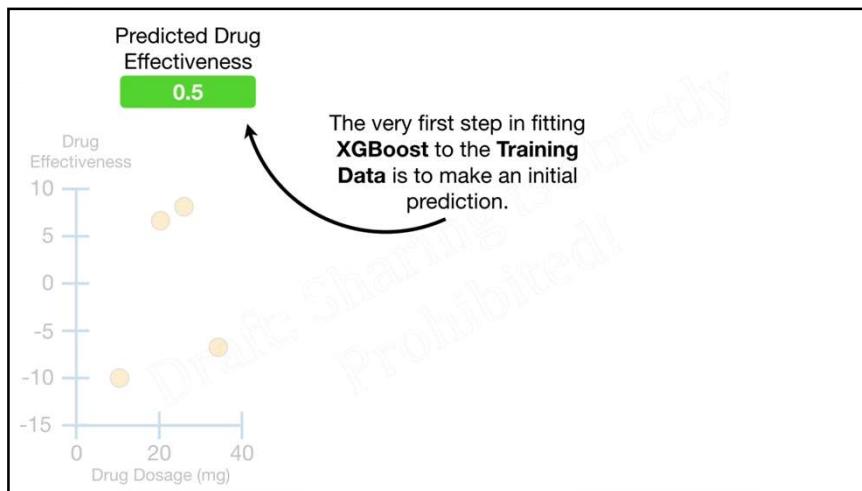
XGBOOST

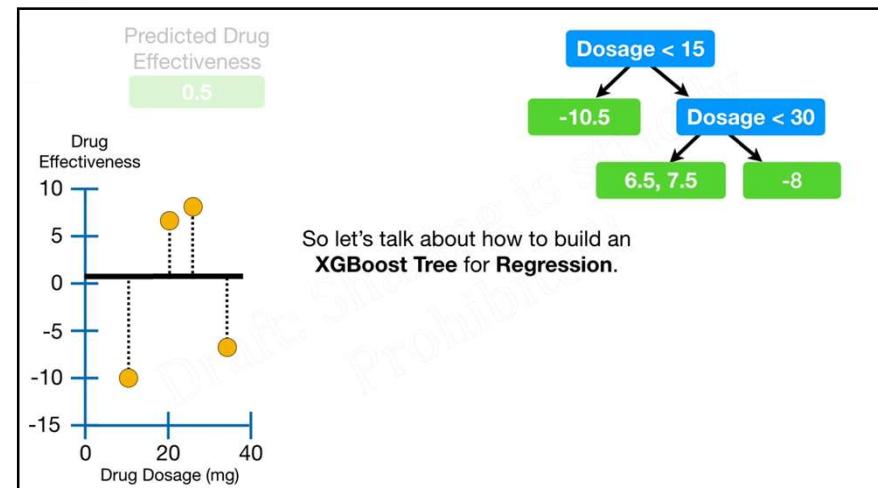
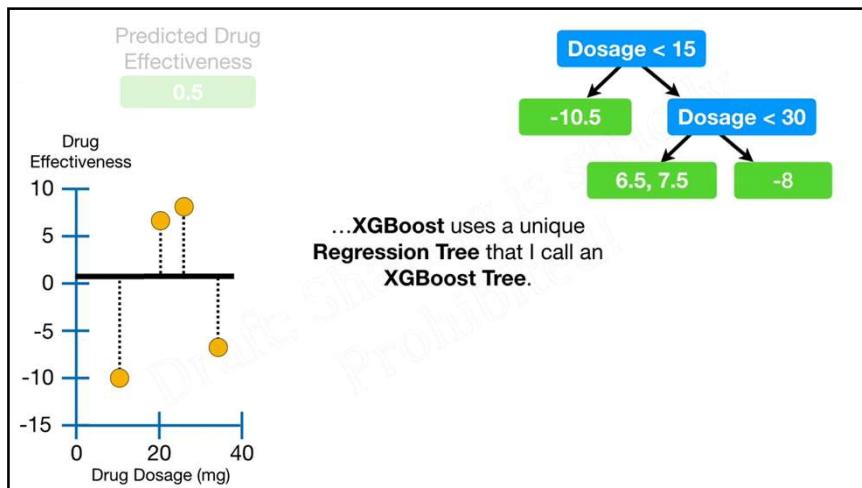
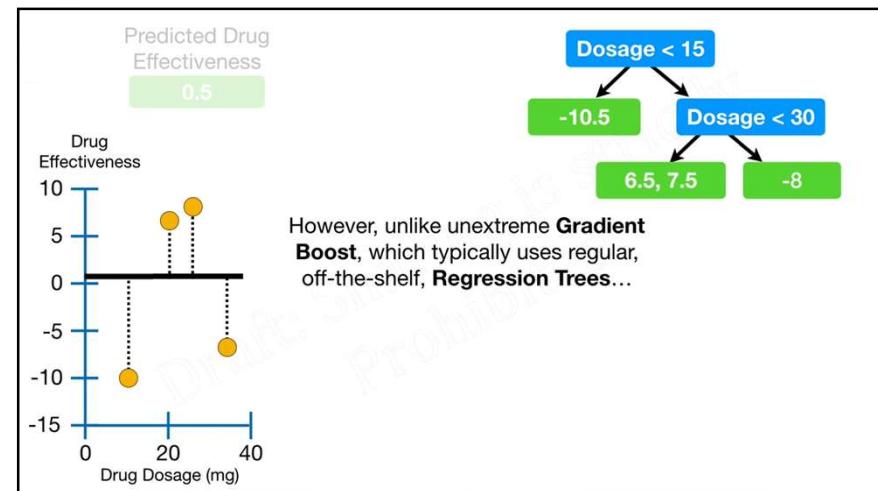
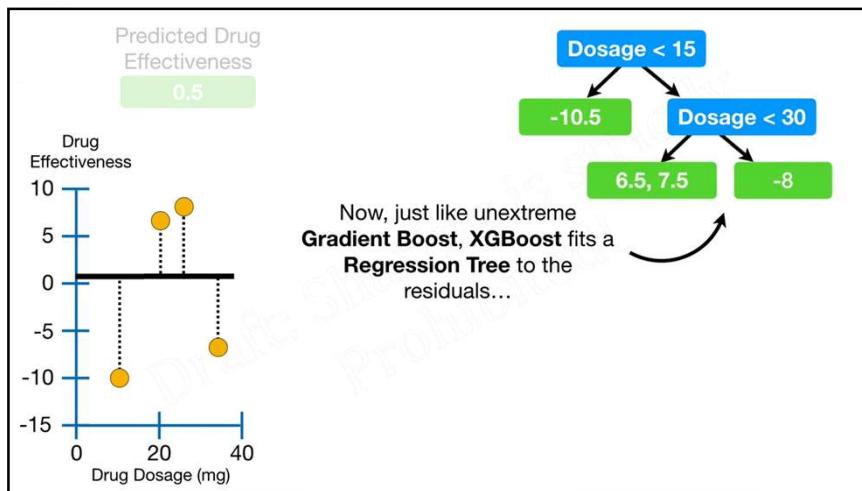
XGBoost was designed to be used with large, complicated data sets.

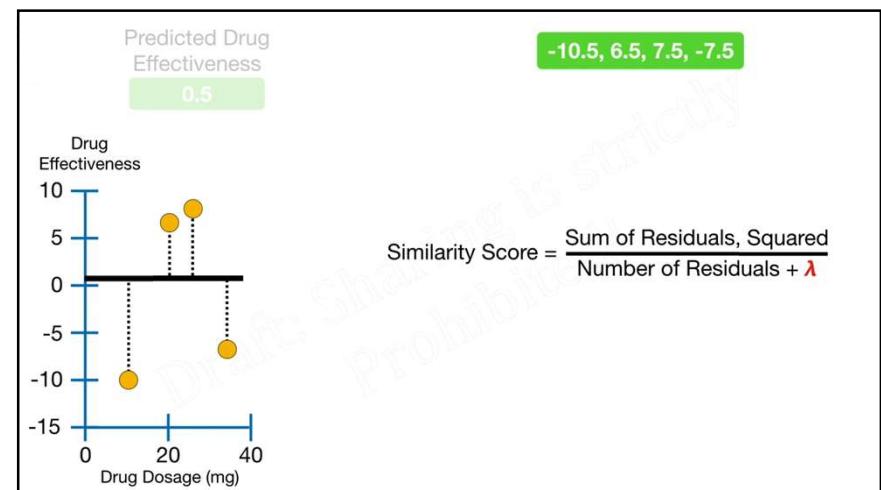
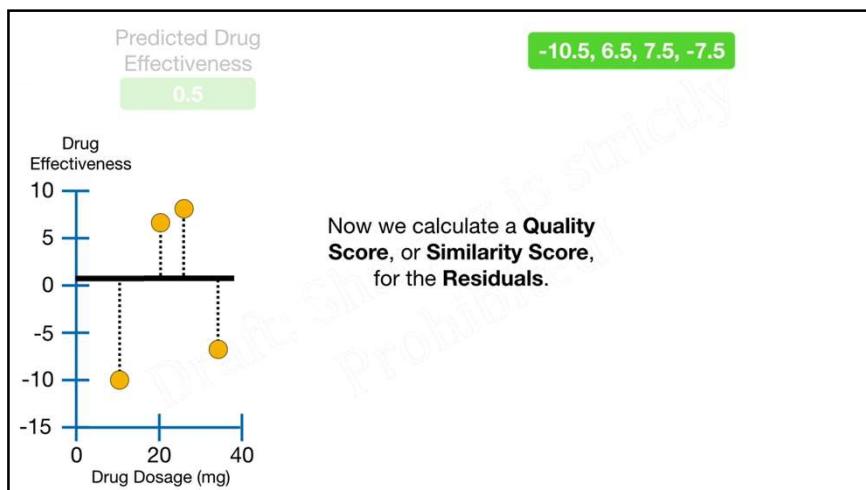
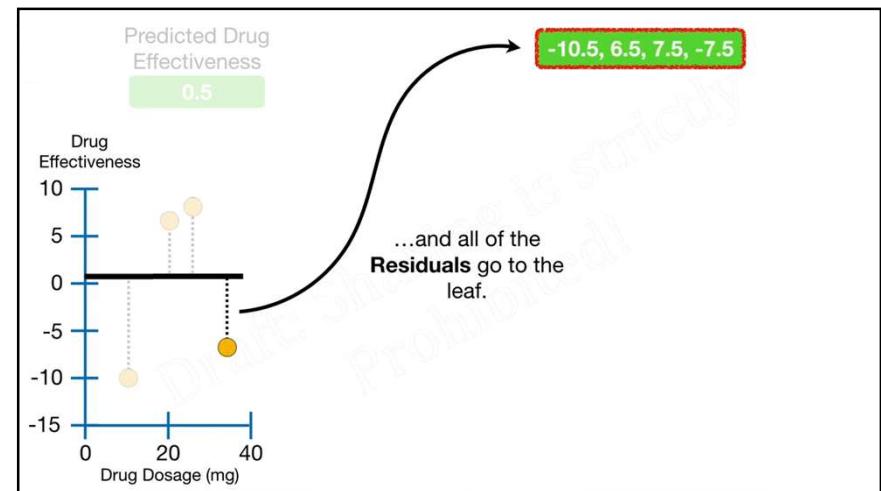
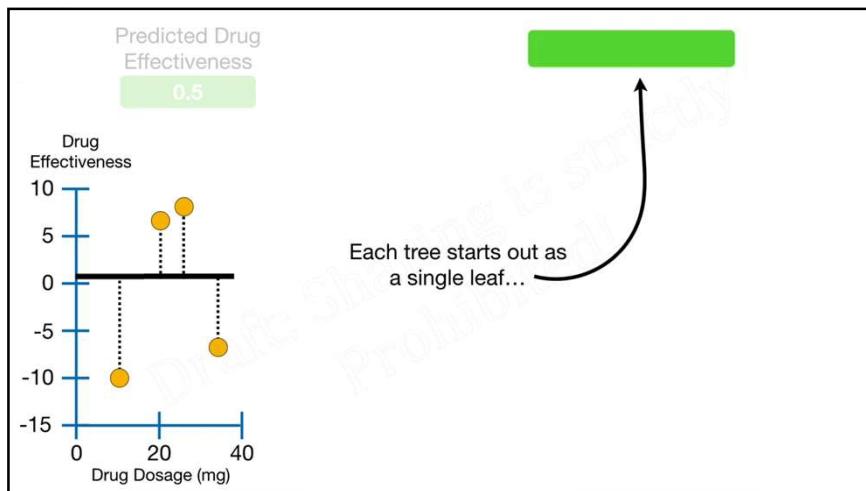
Mass	Age	BP	Color	Movie	Car	Hair	etc.
120	23	102	Brown	T2	Ford	Long	...
150	25	98	Brown	Frozen	Kia	Short	...
165	22	130	Black	Spiderman	Ford	Short	...
123	45	98	Red	T2	Kia	Long	...
156	33	78	Brown	Frozen 2	Ford	Long	...
...

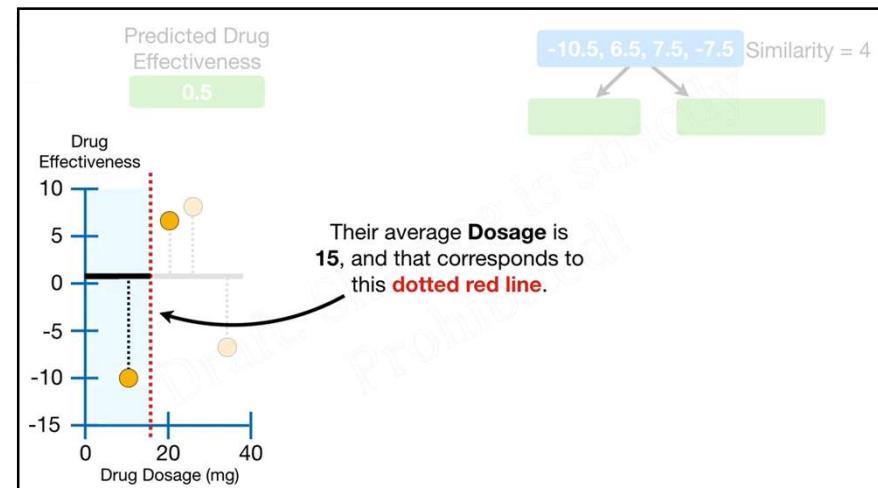
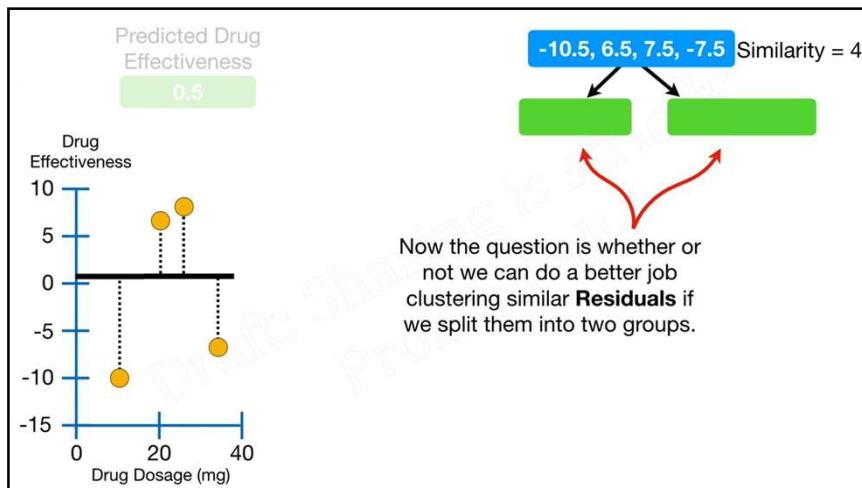
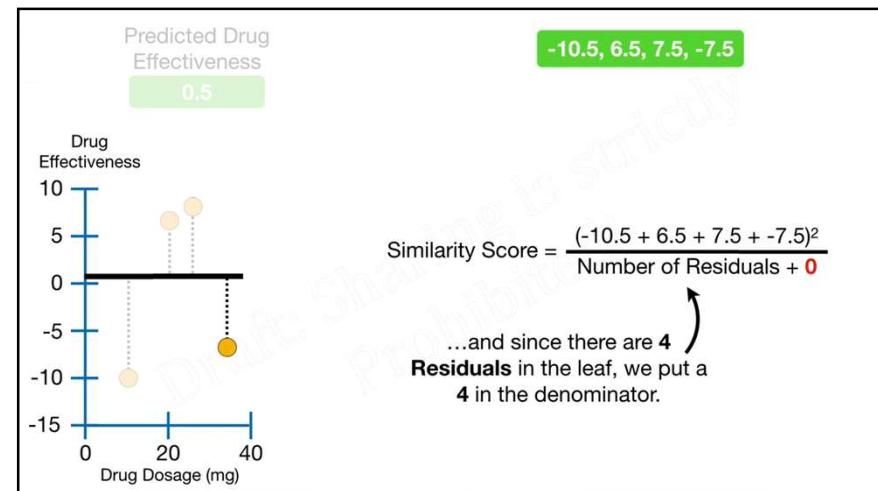
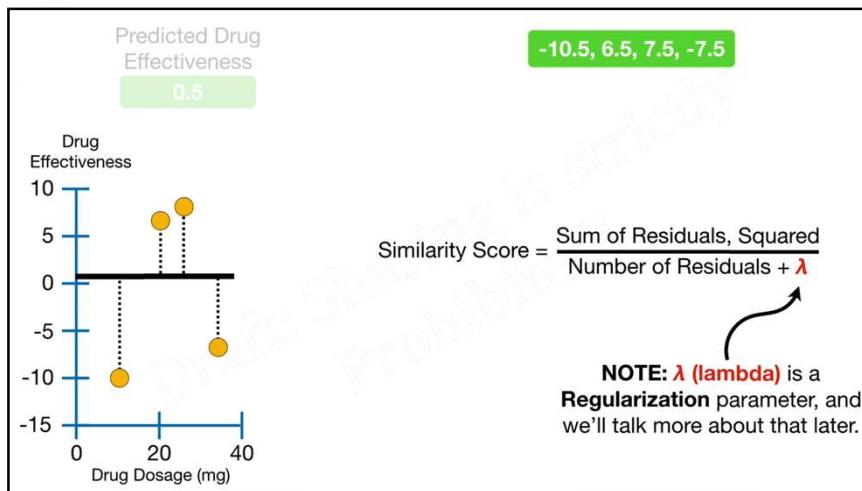


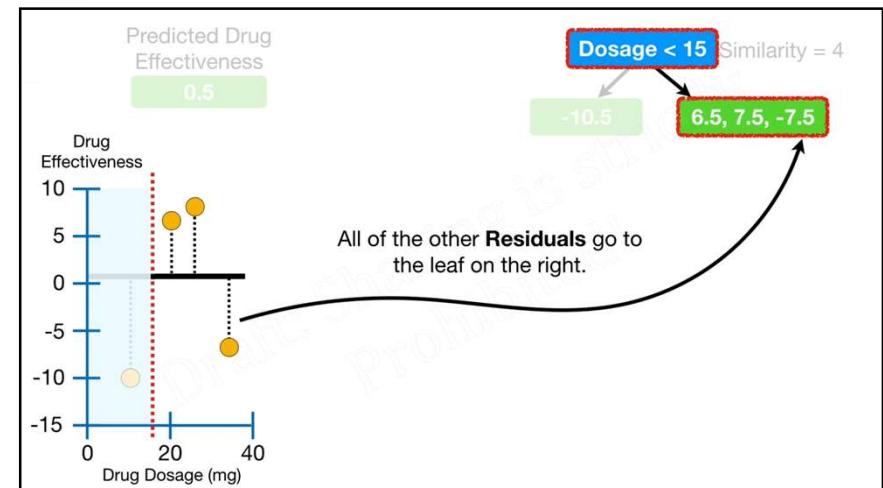
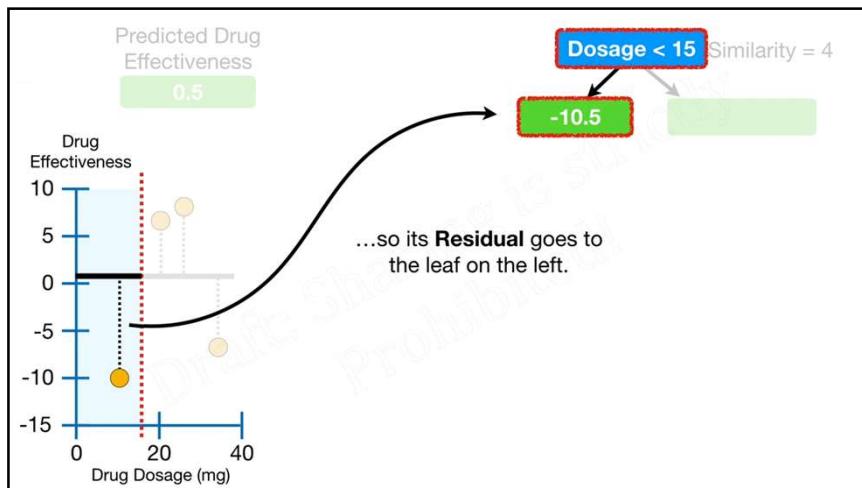
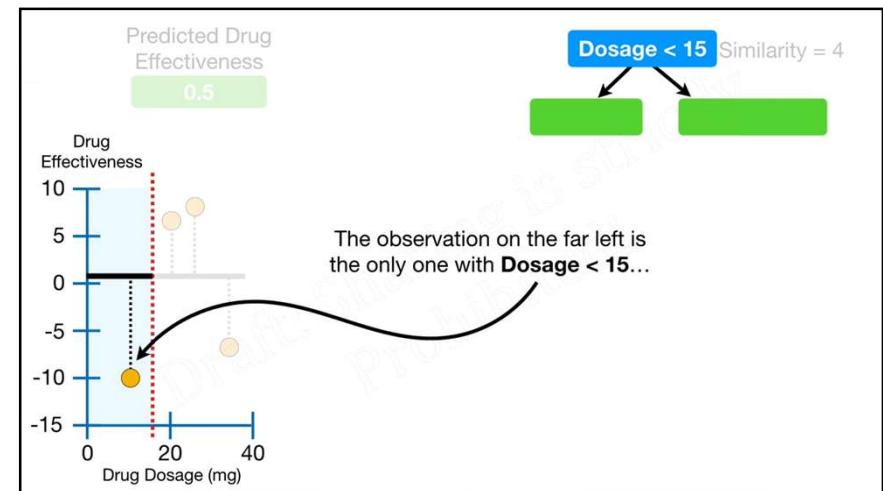
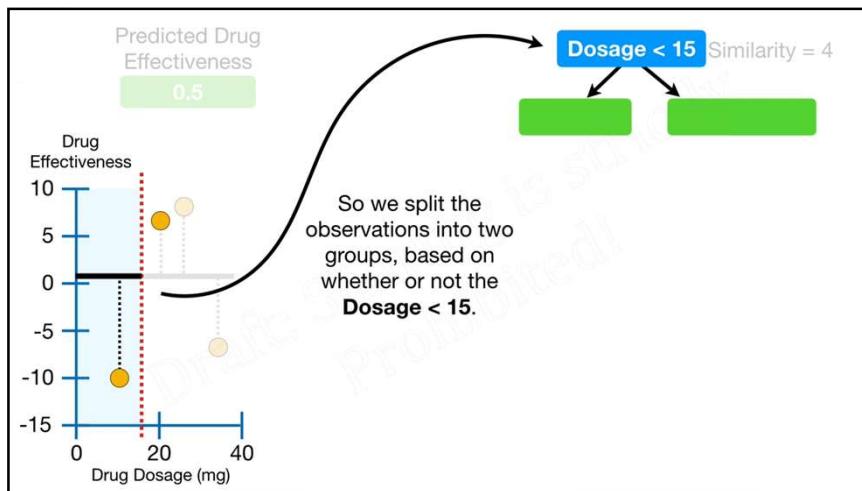
However, to keep the examples from getting out of hand, we will use this super simple **Training Data**.

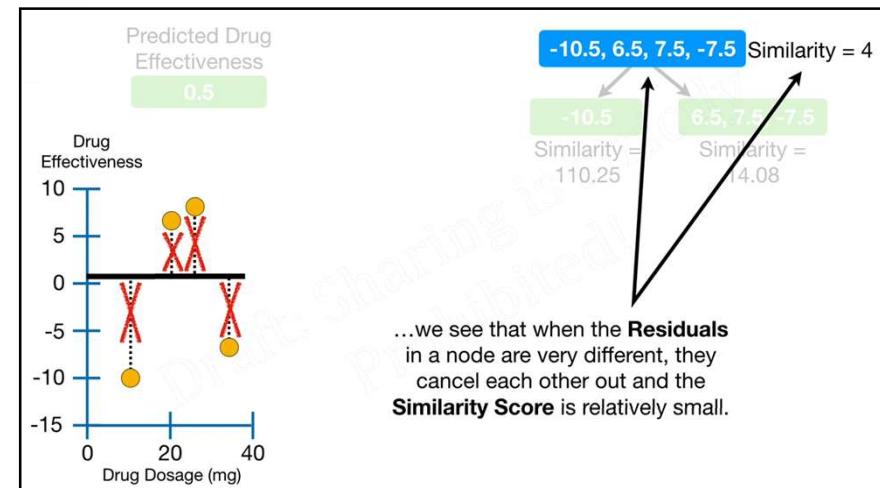
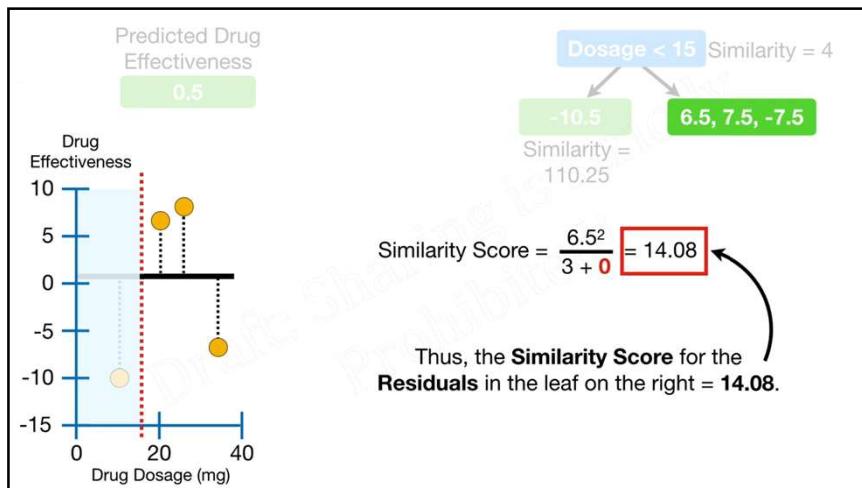
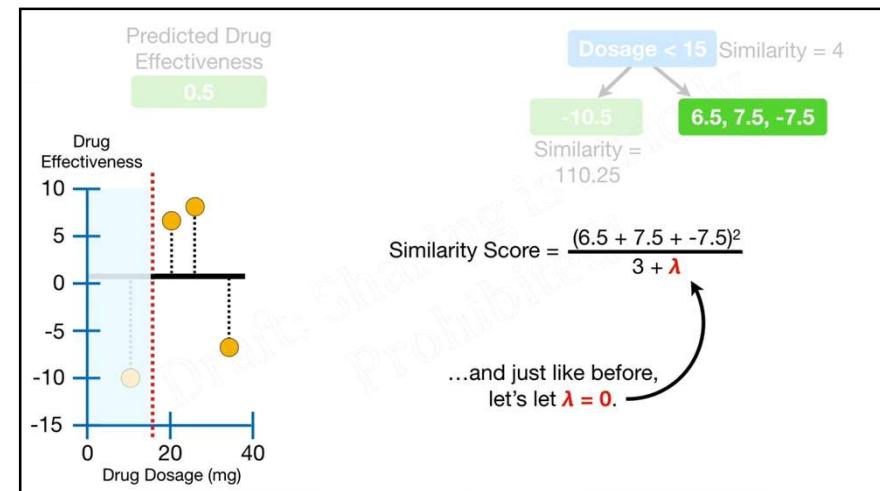
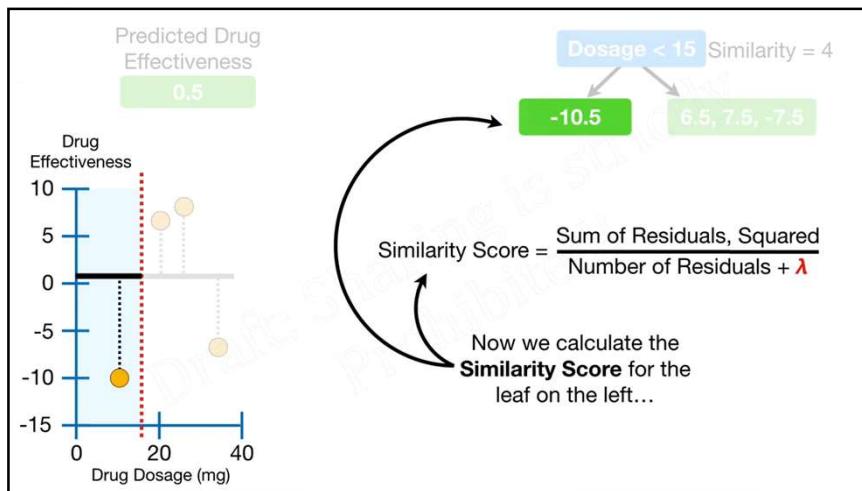


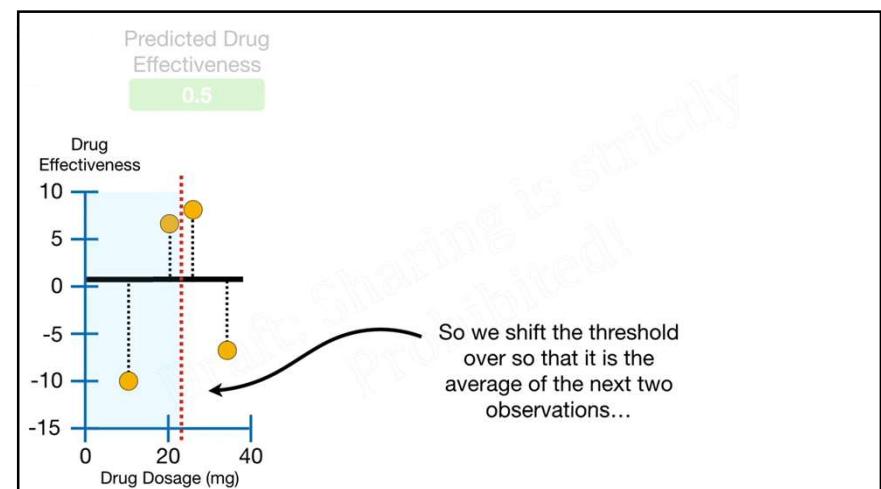
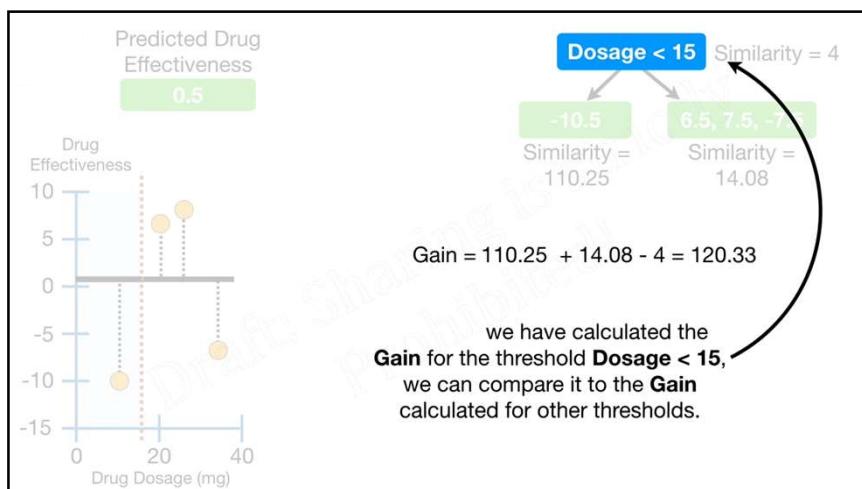
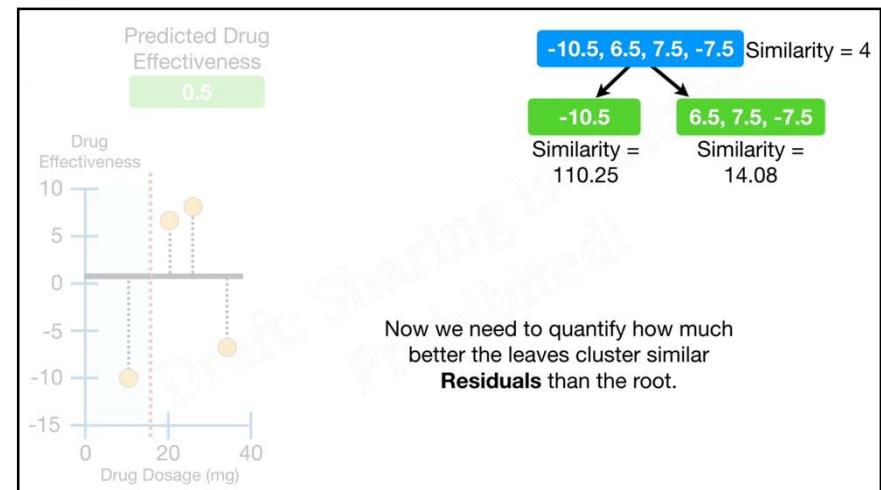
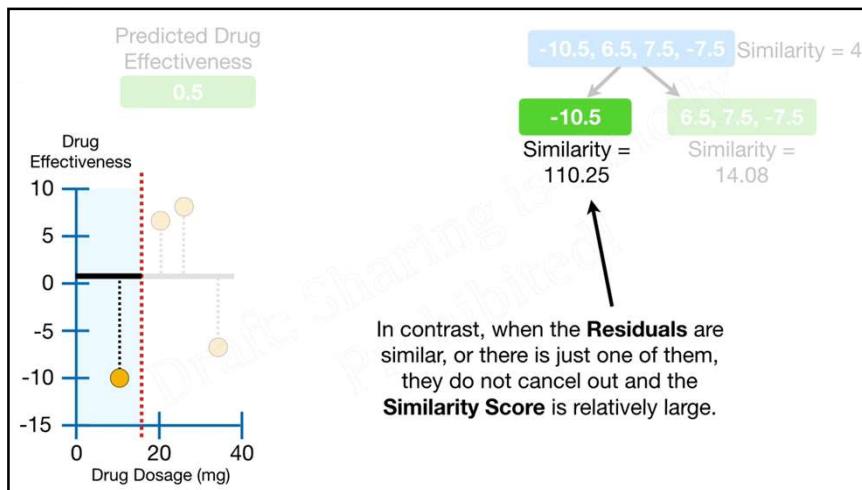


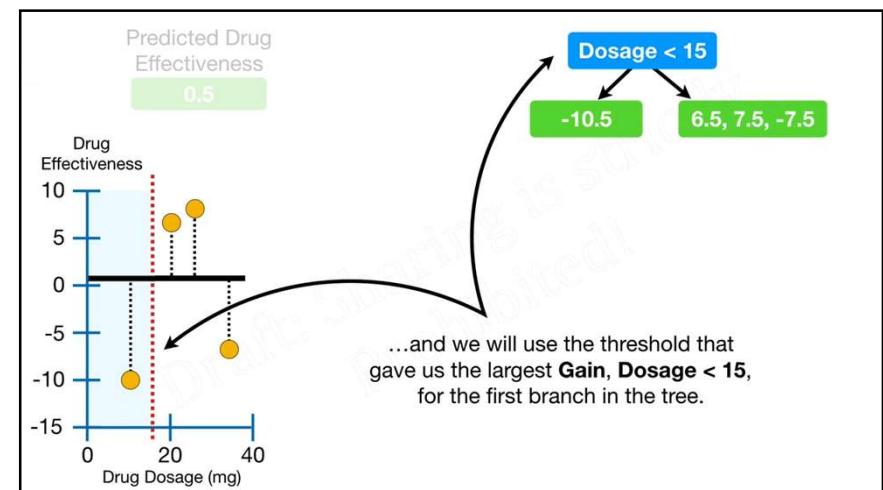
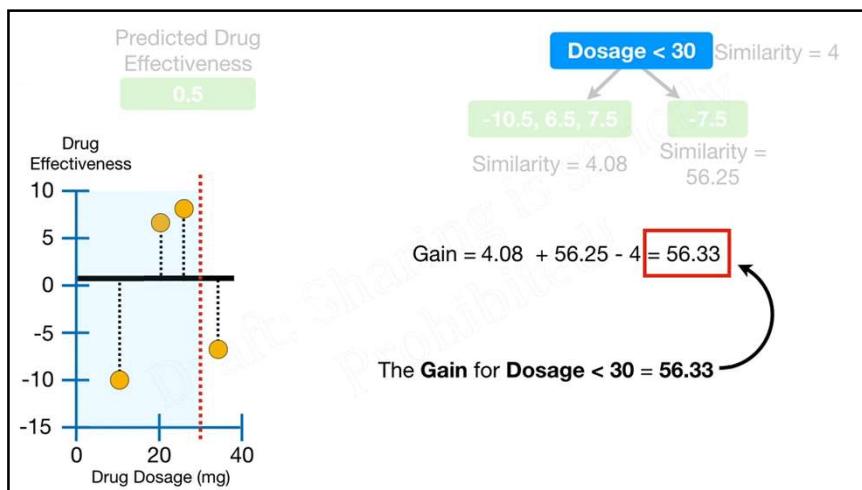
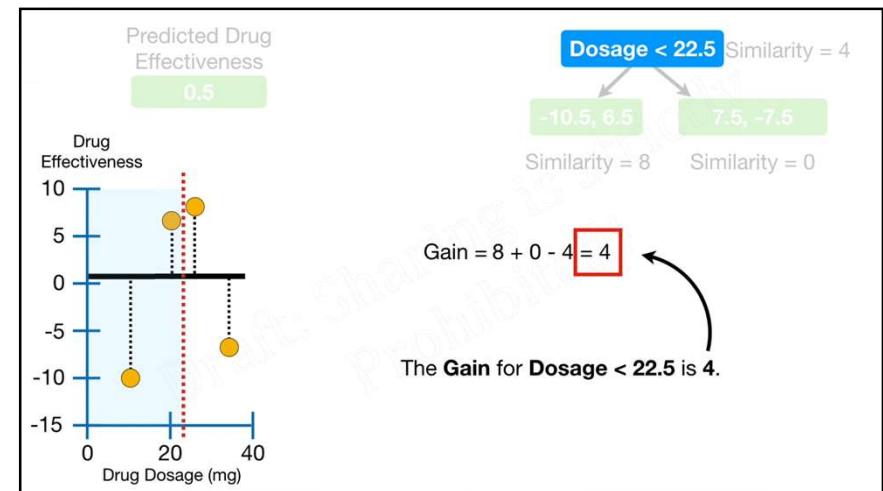
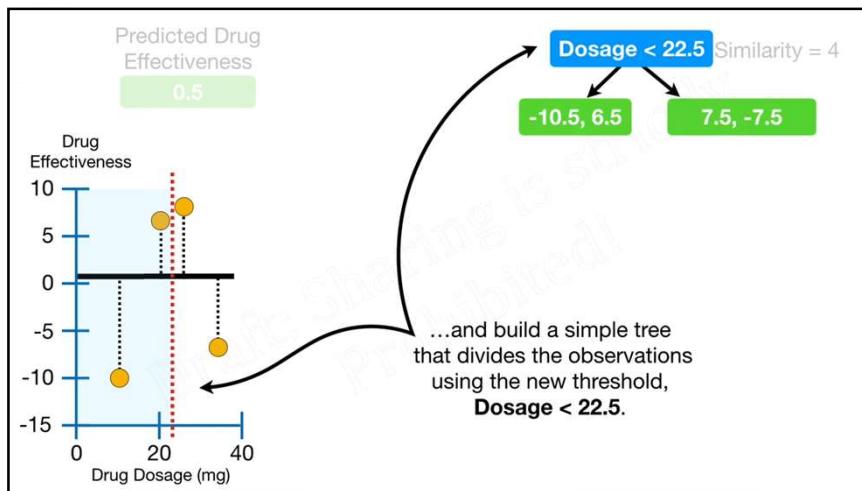


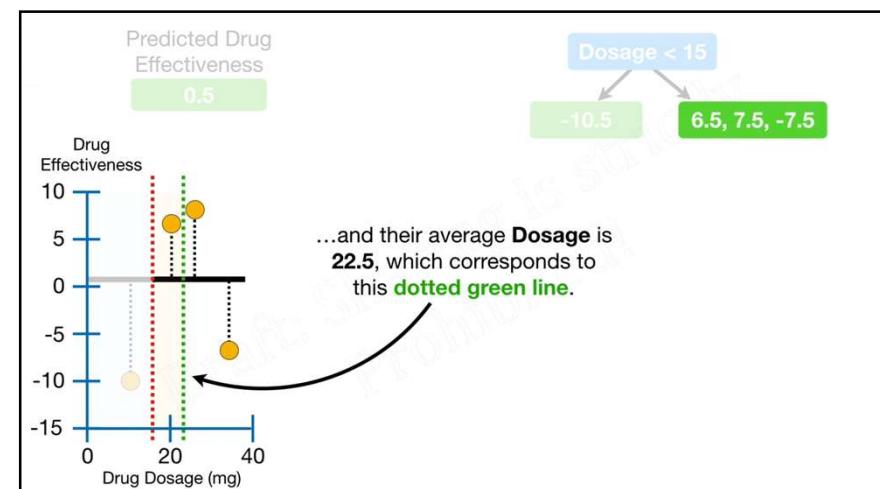
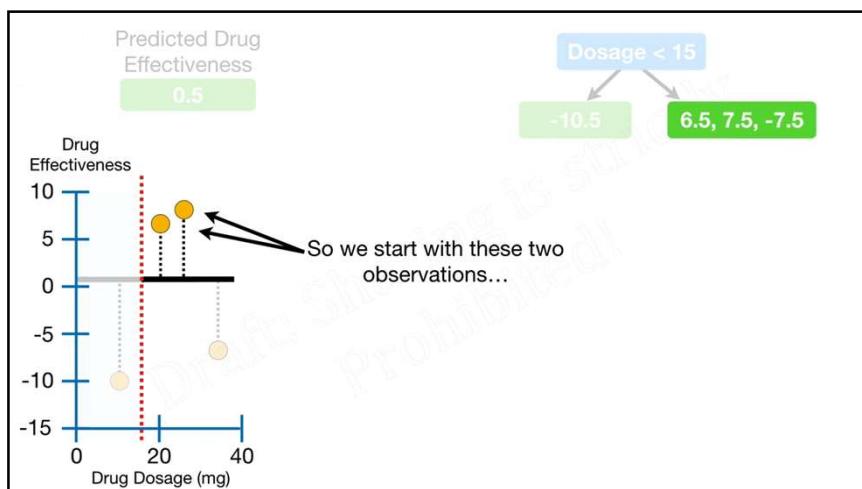
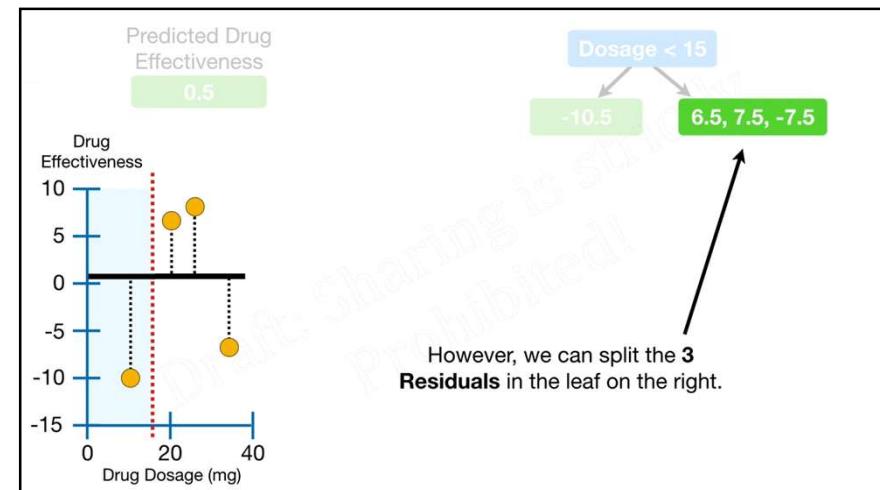
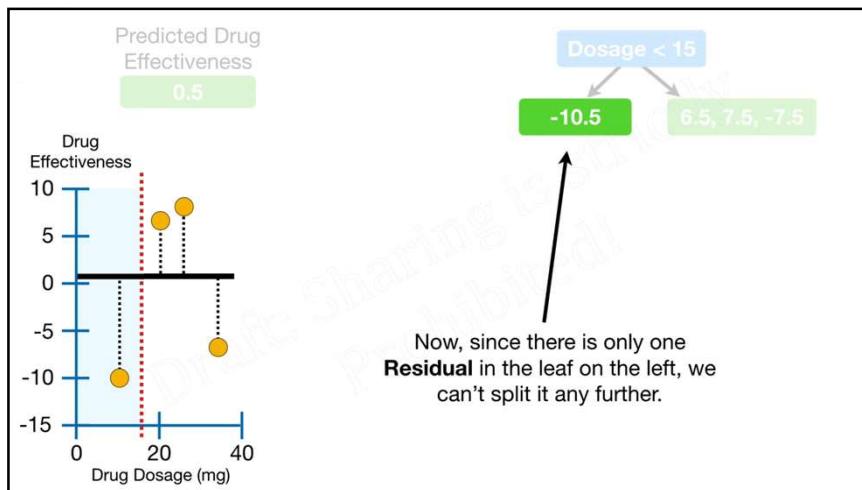


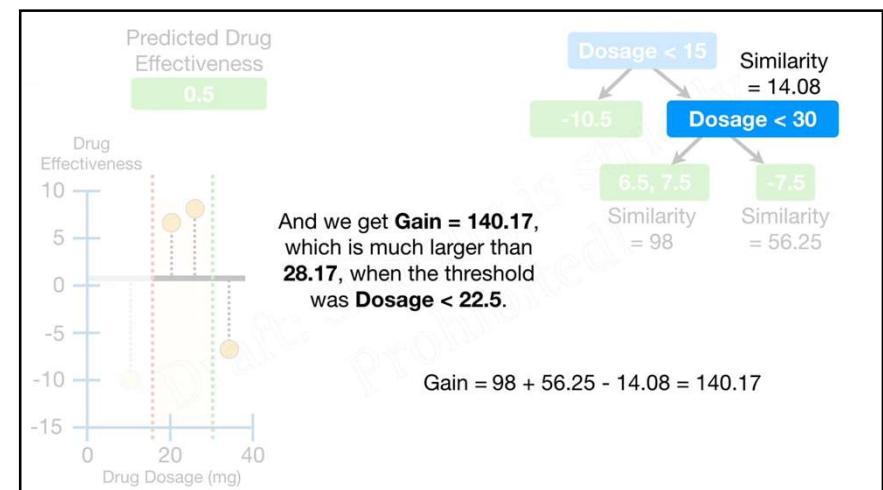
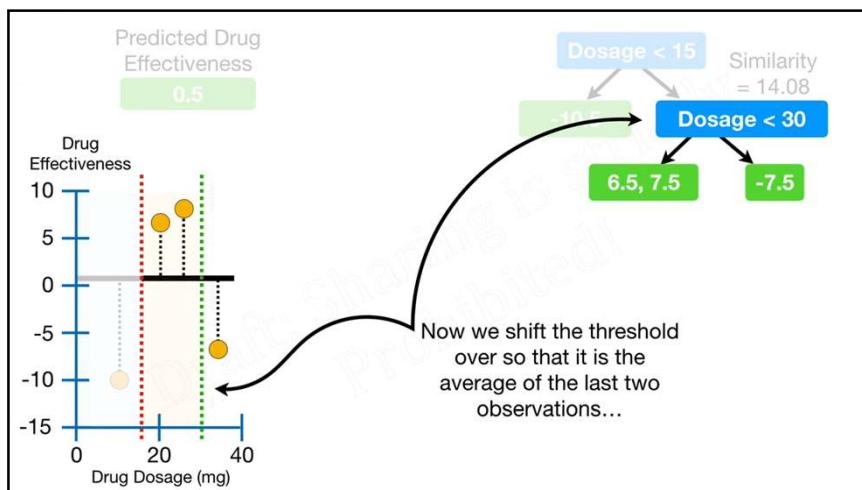
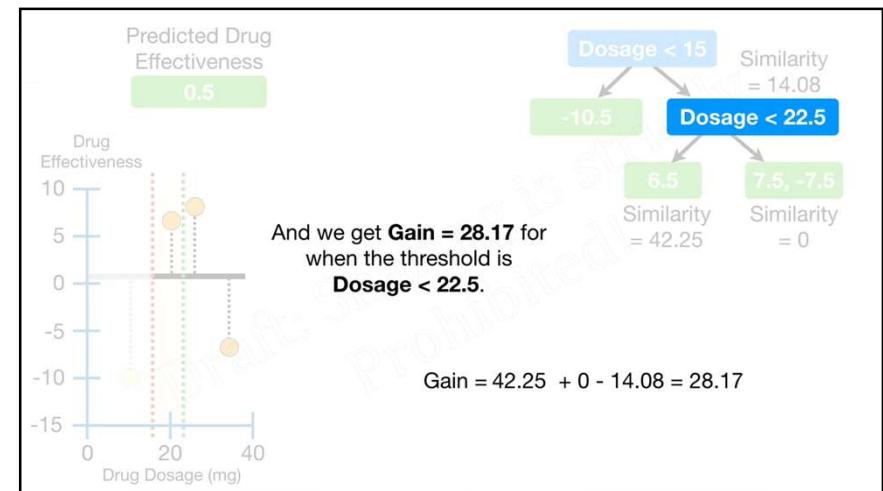
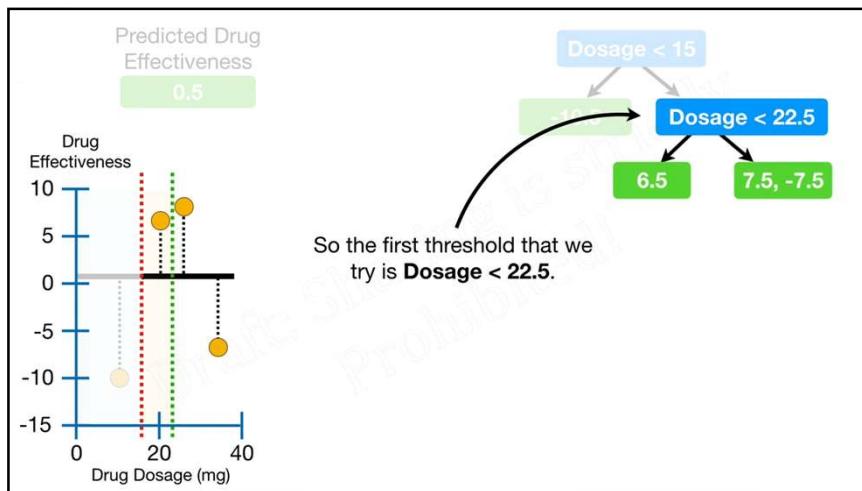


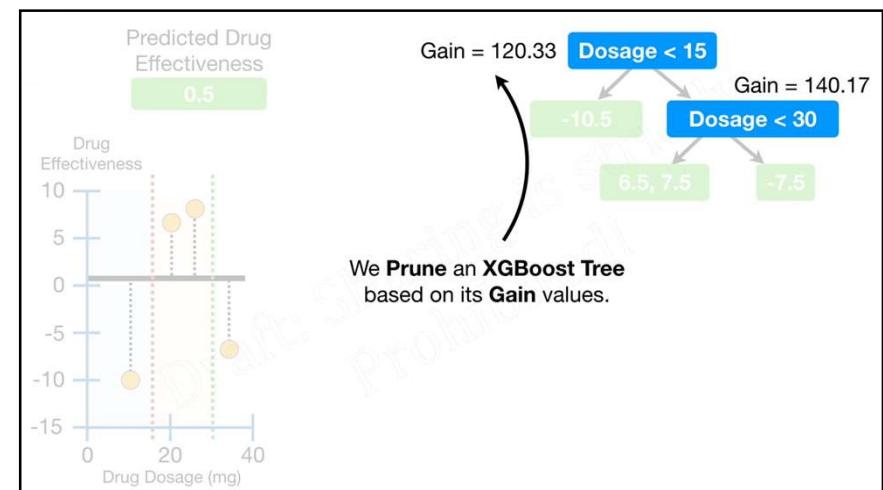
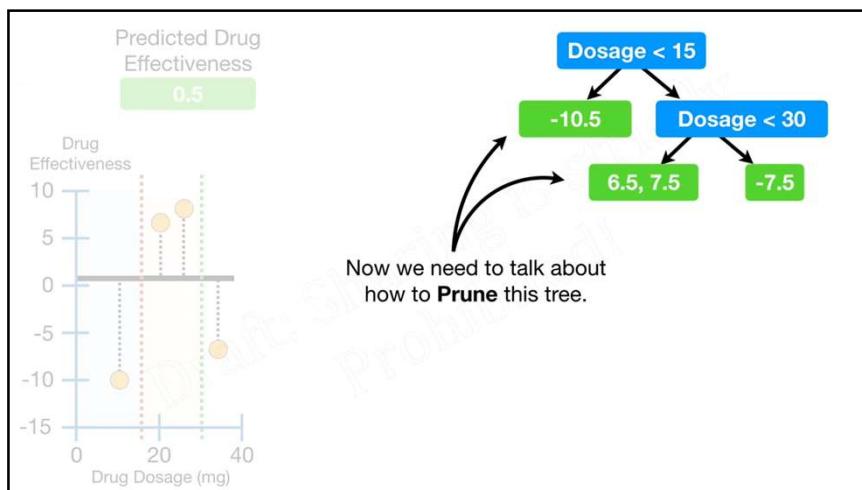
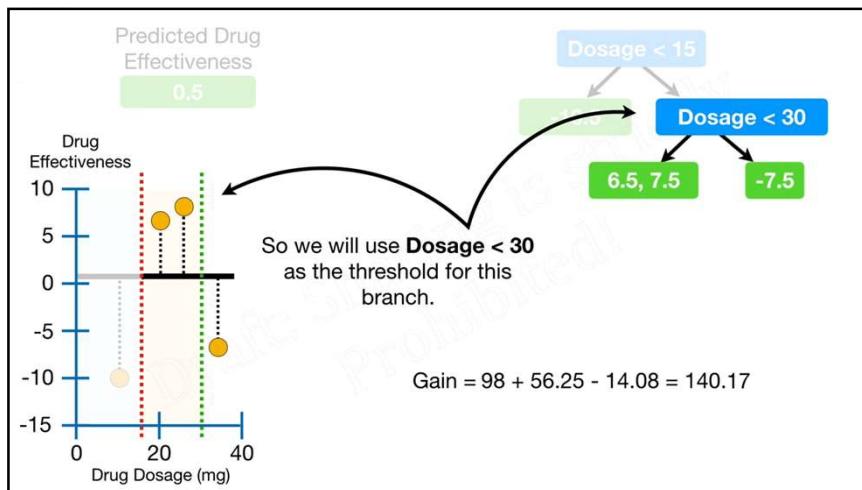


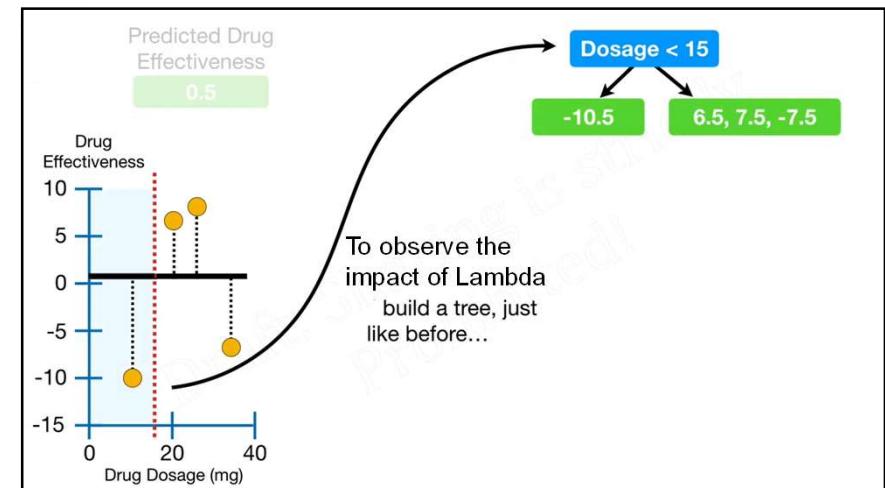
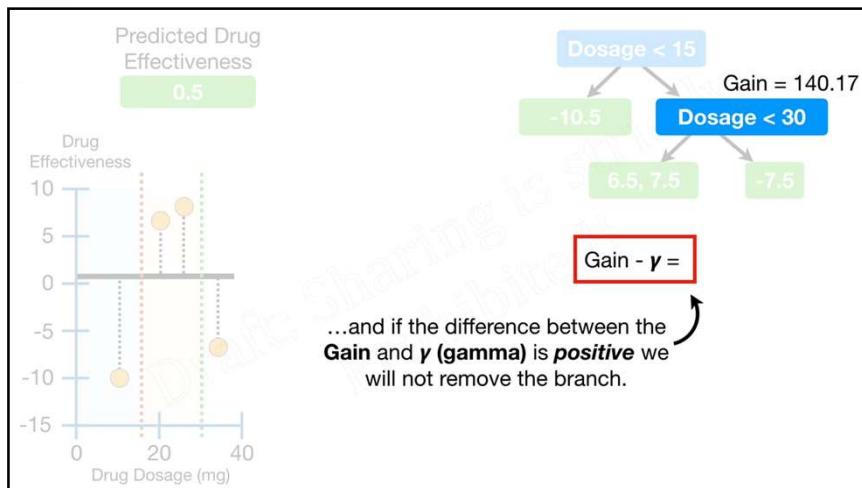
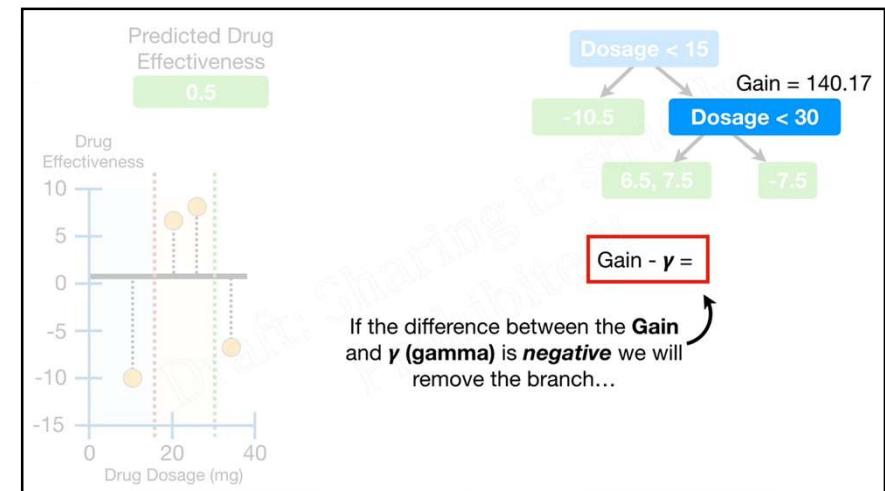


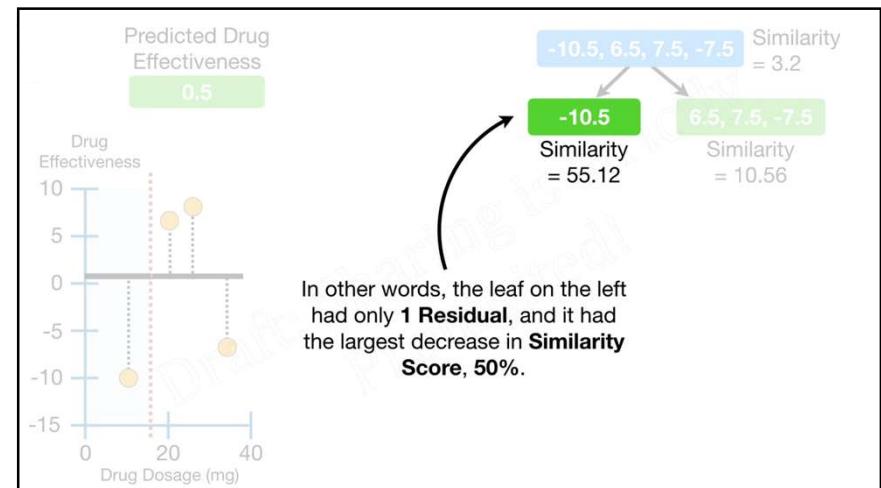
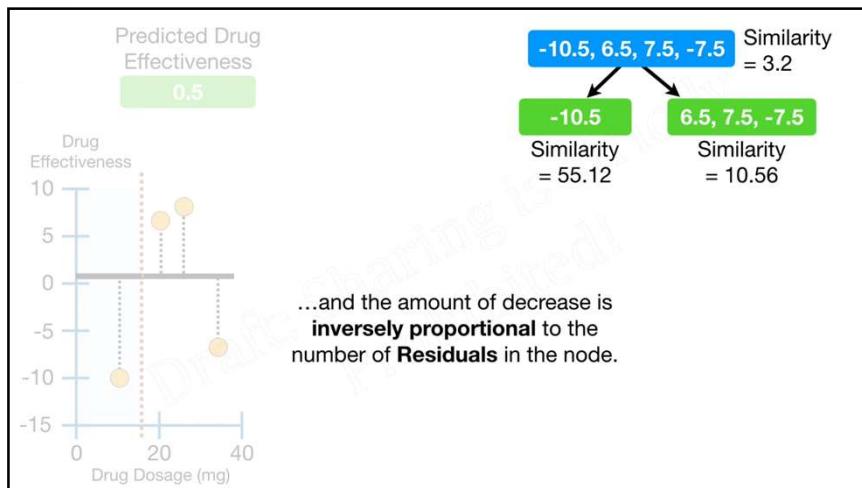
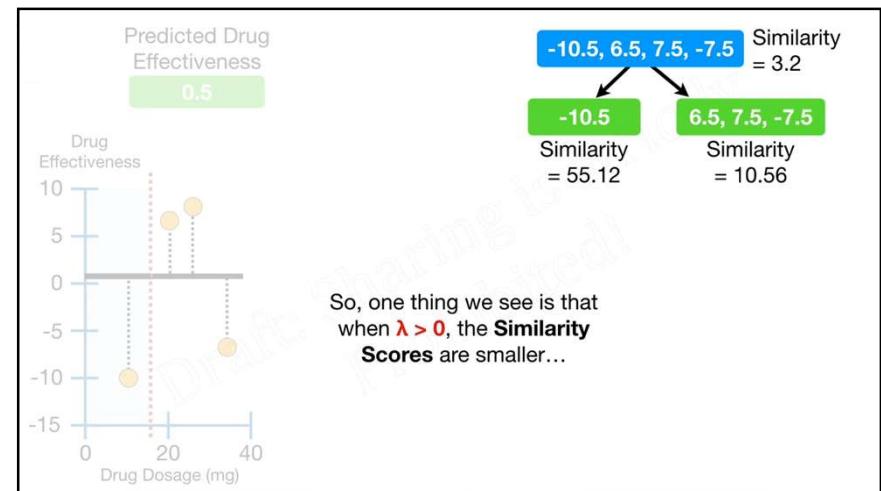
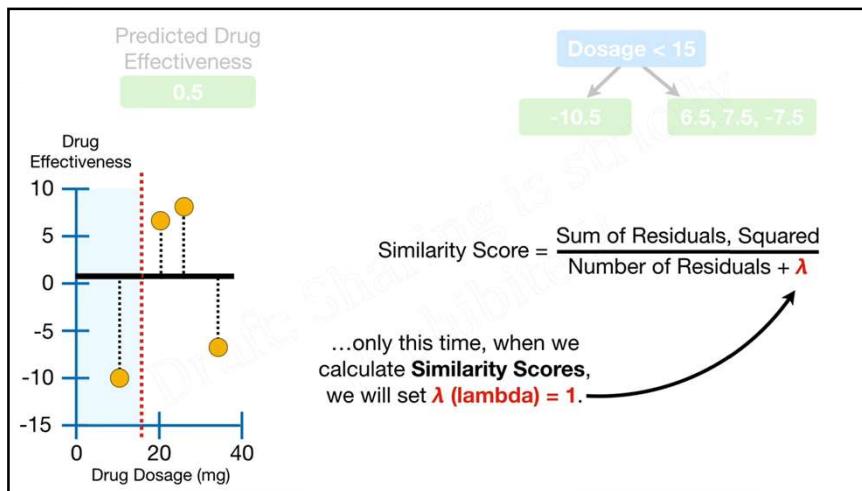


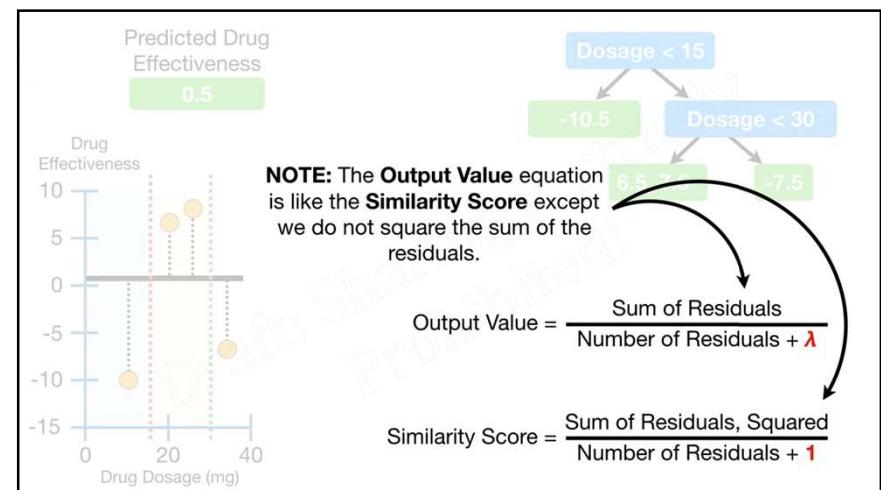
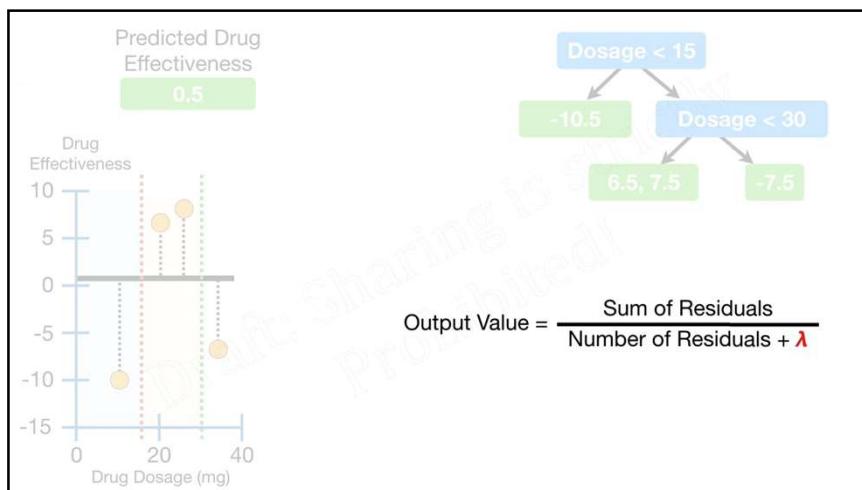
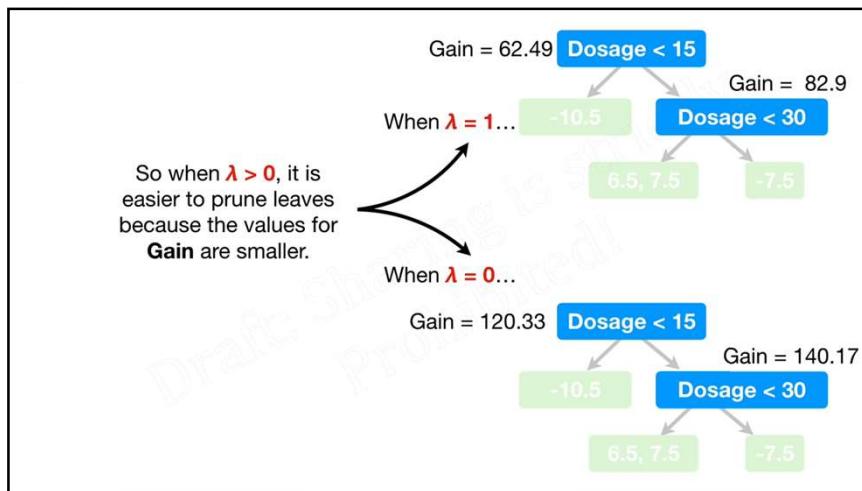


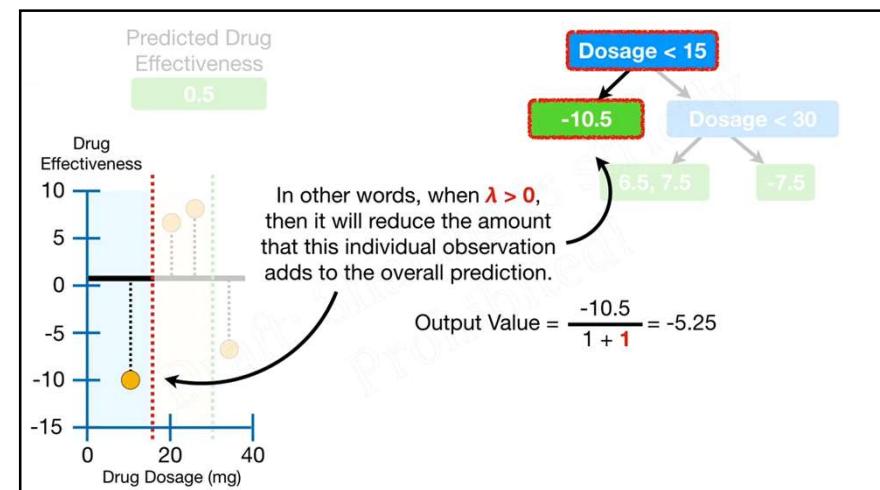
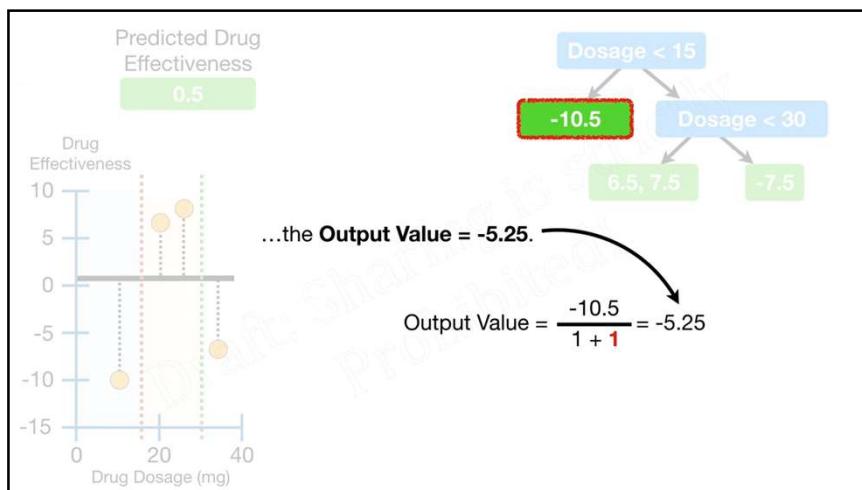
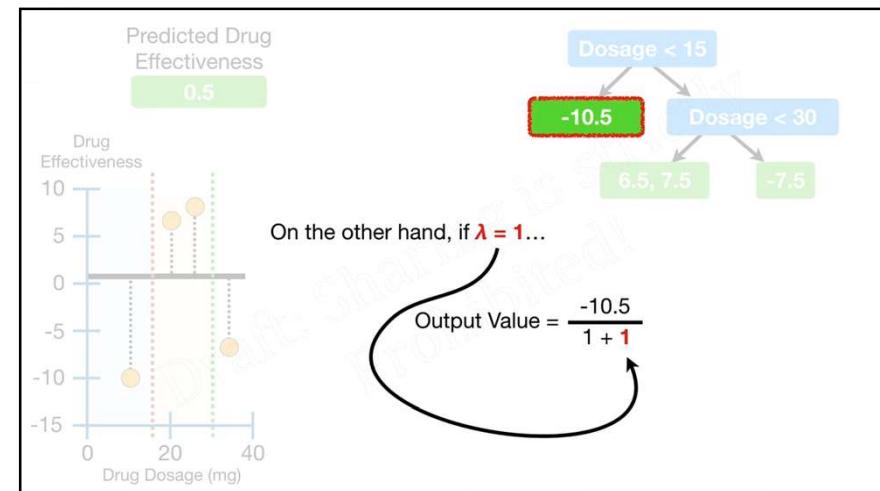
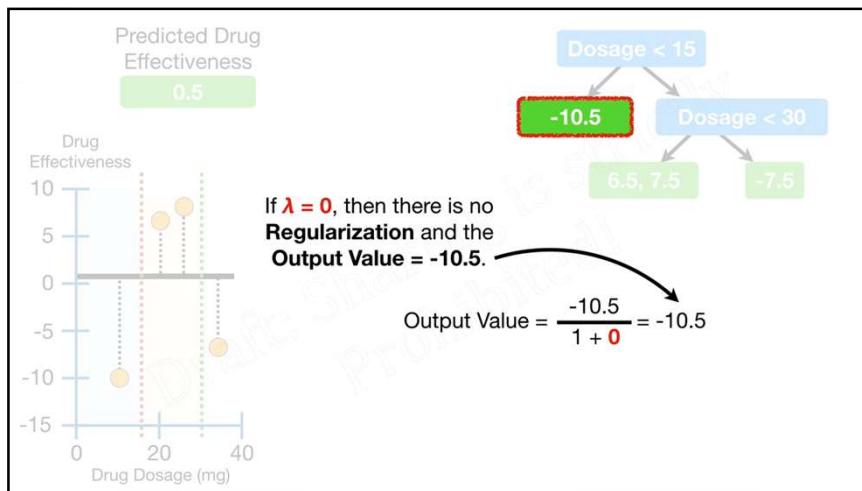


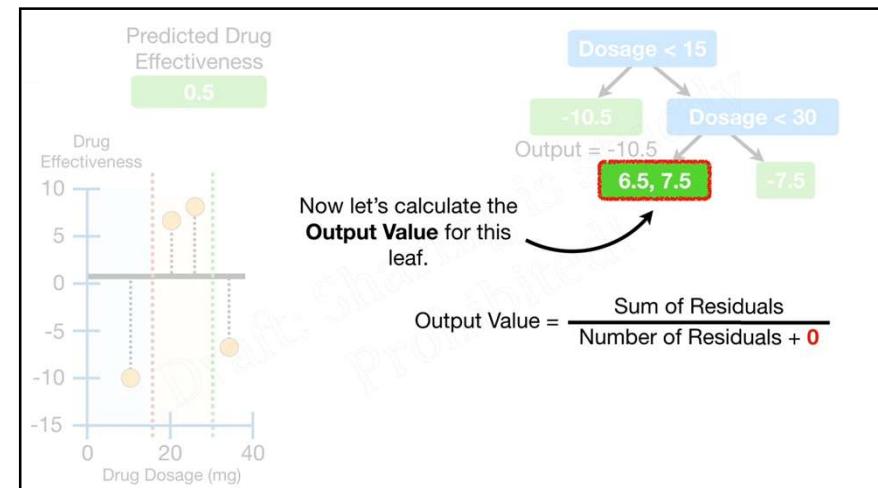
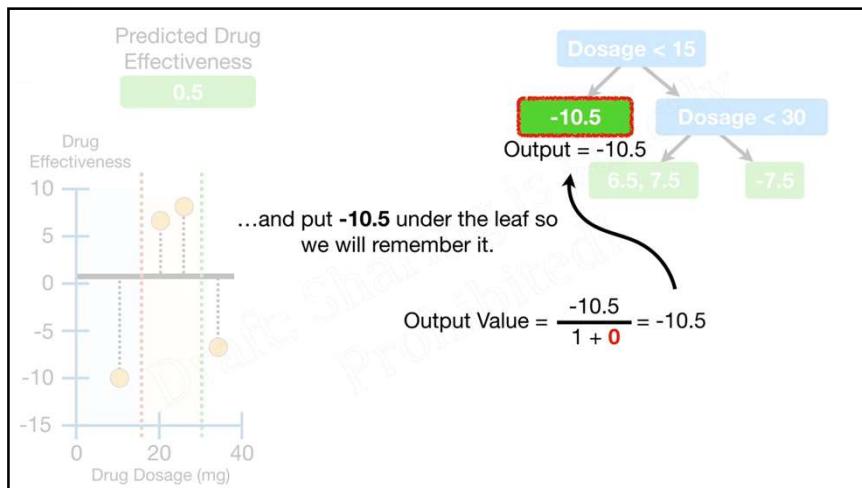
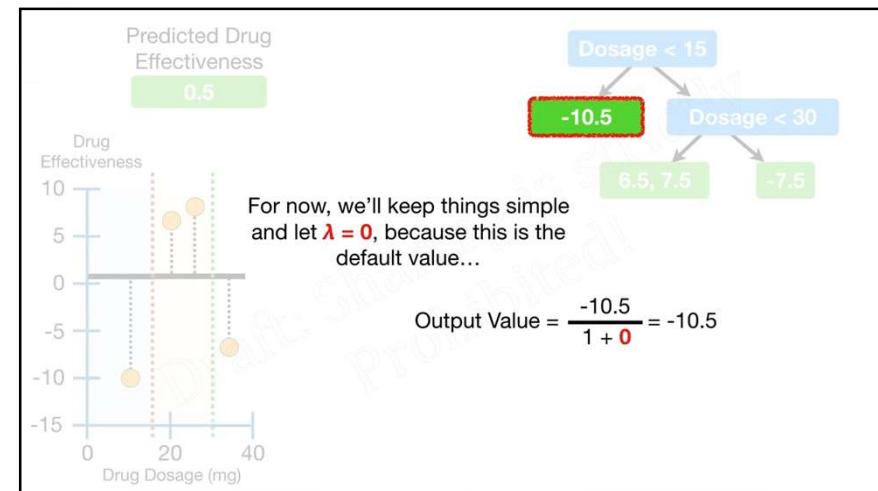
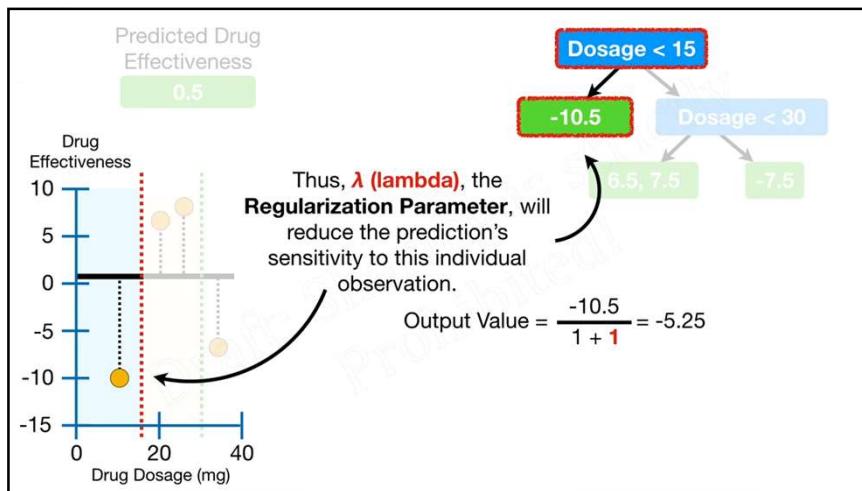


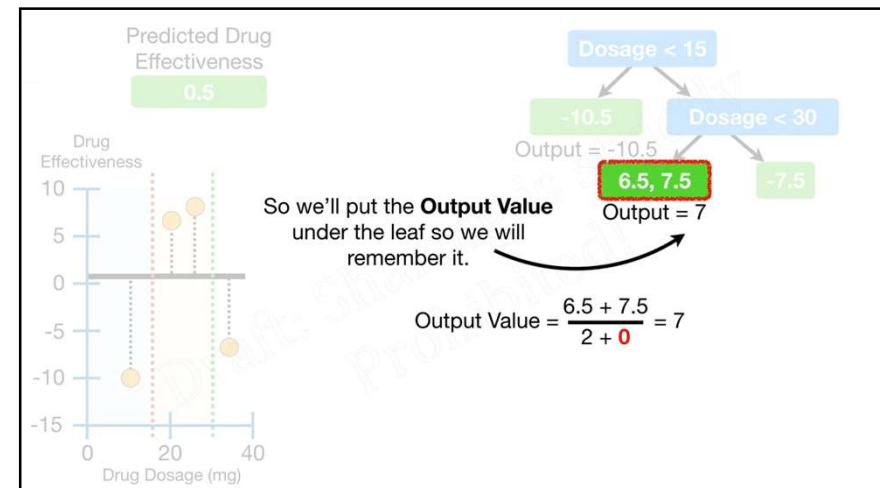
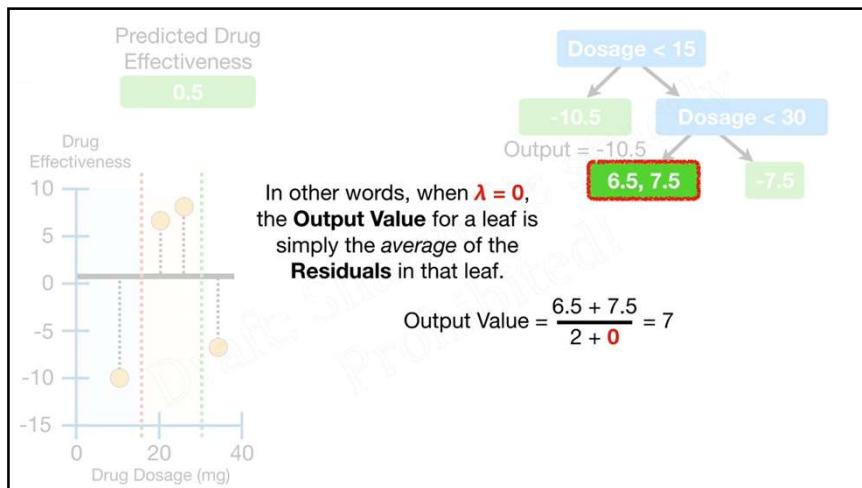
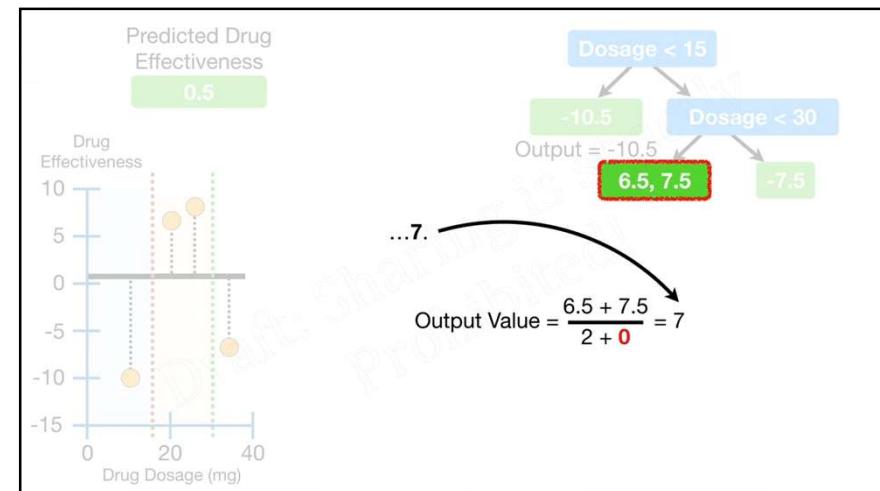
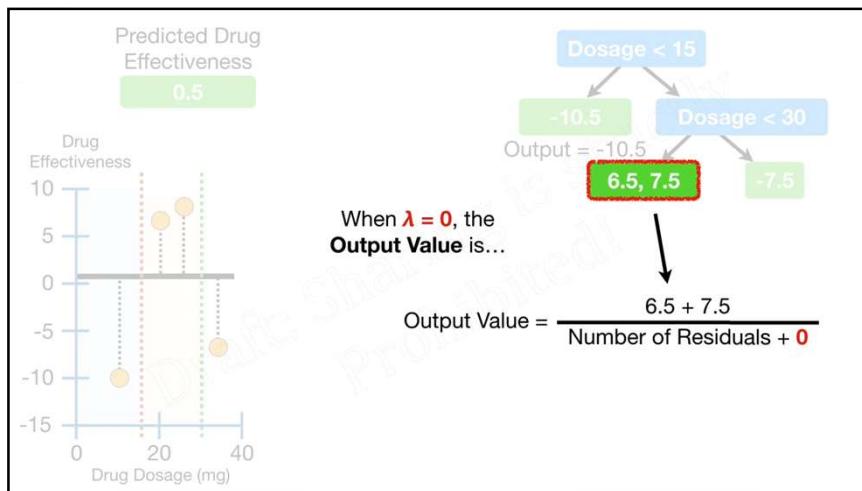


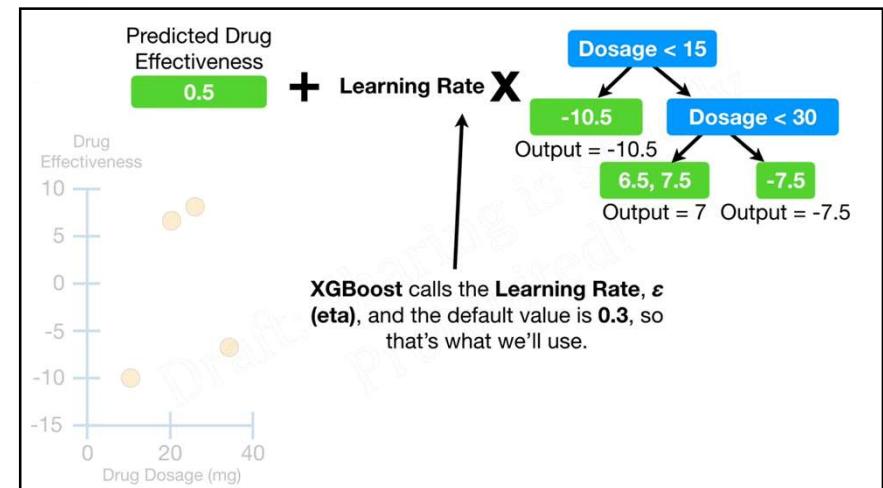
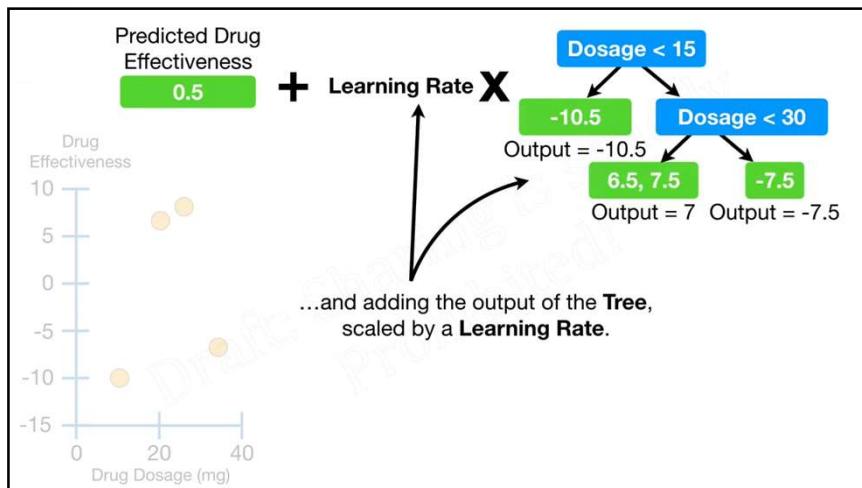
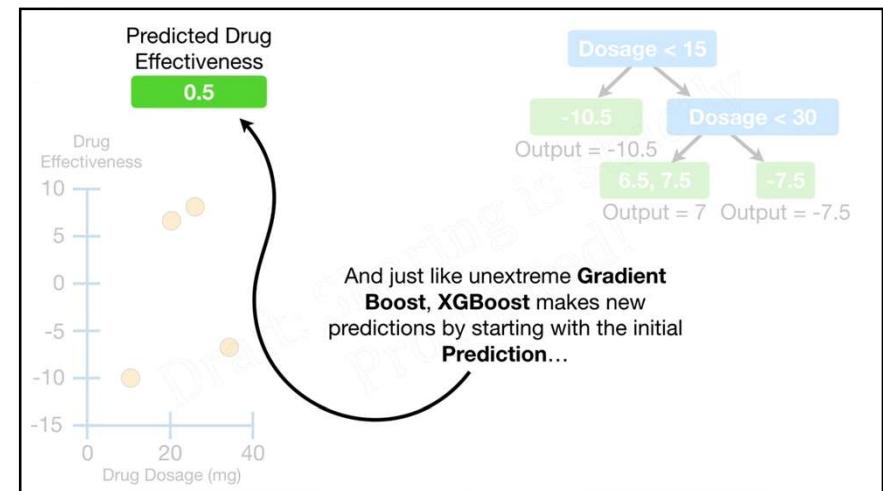
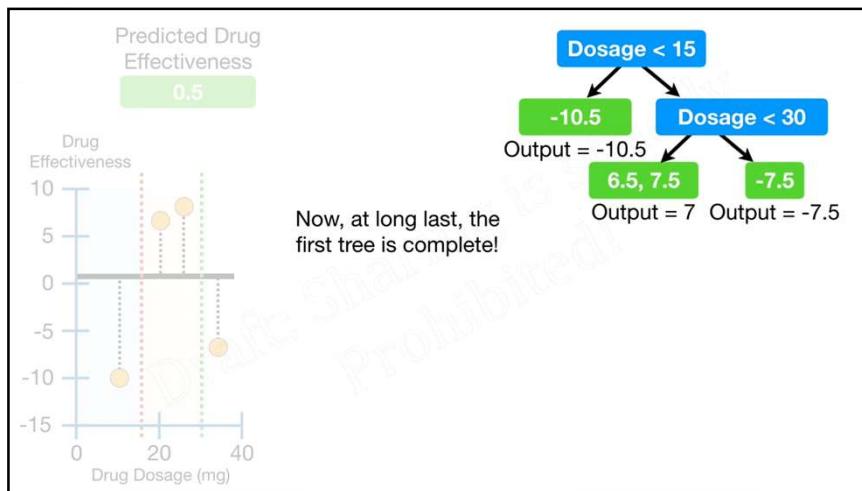


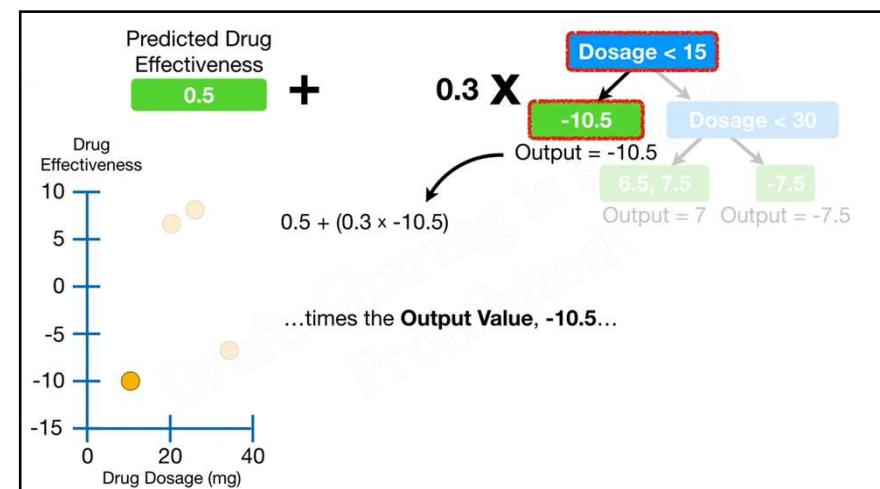
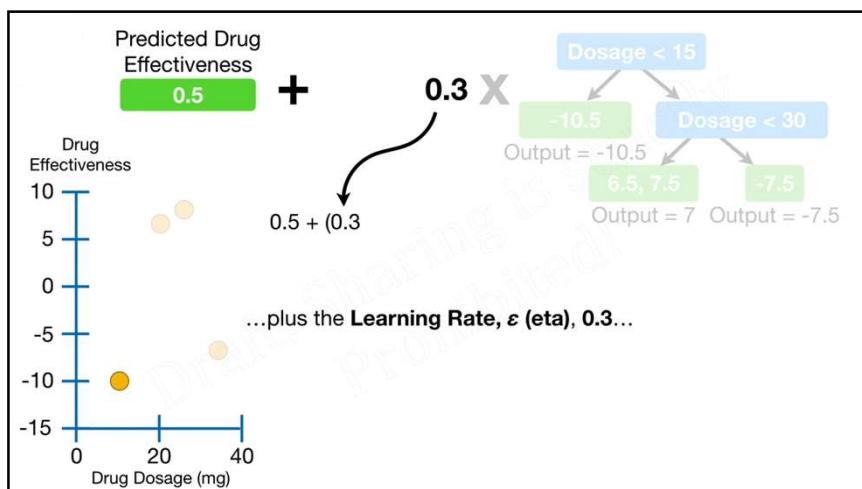
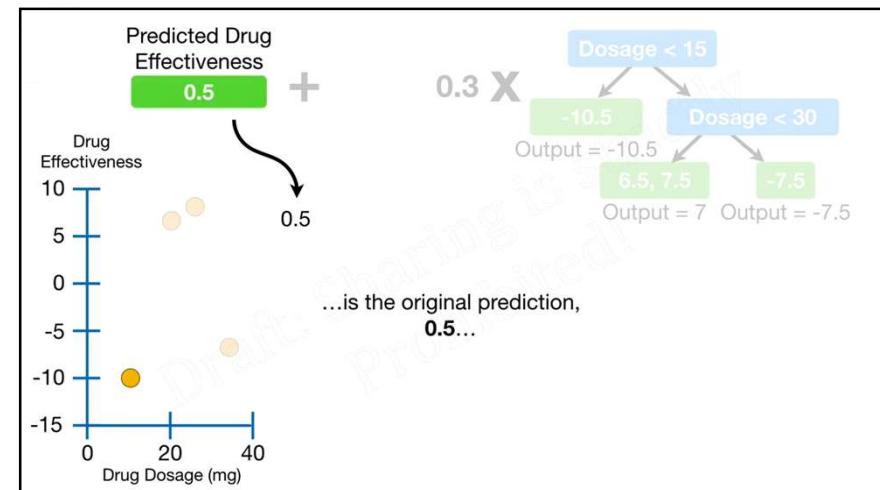
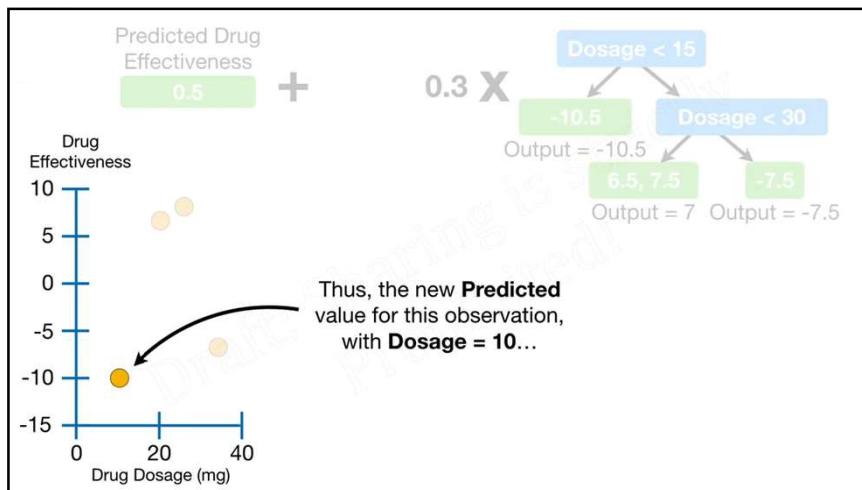


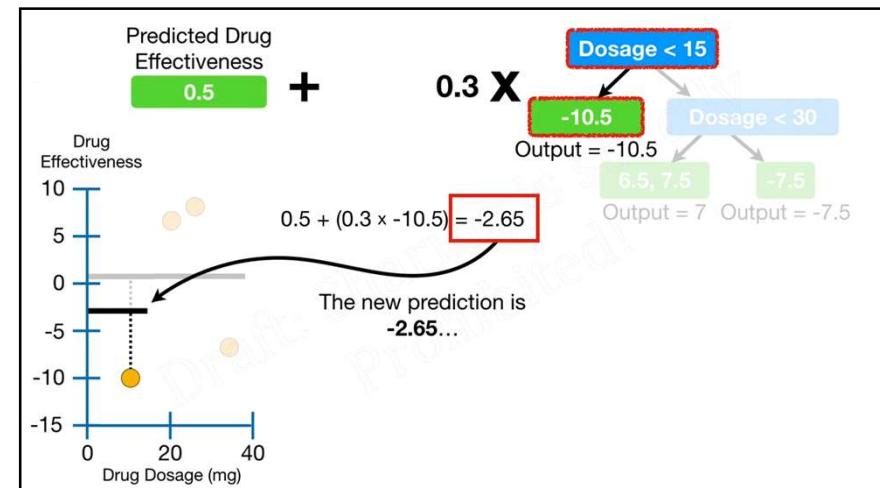
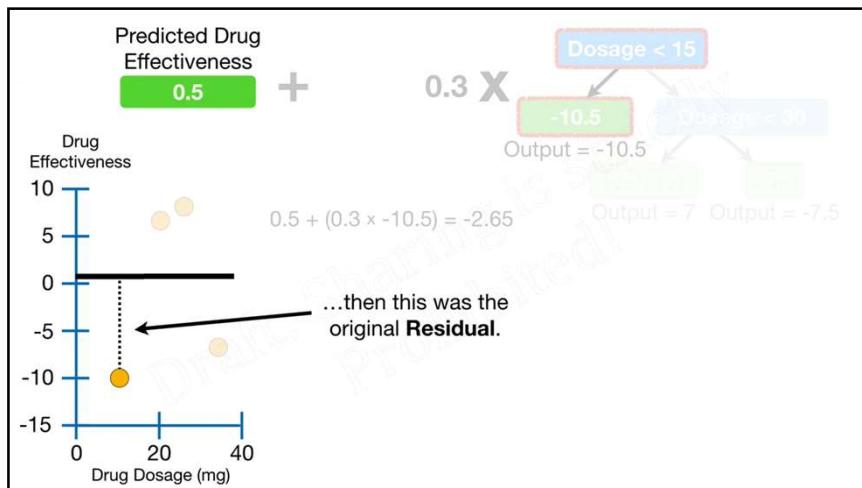
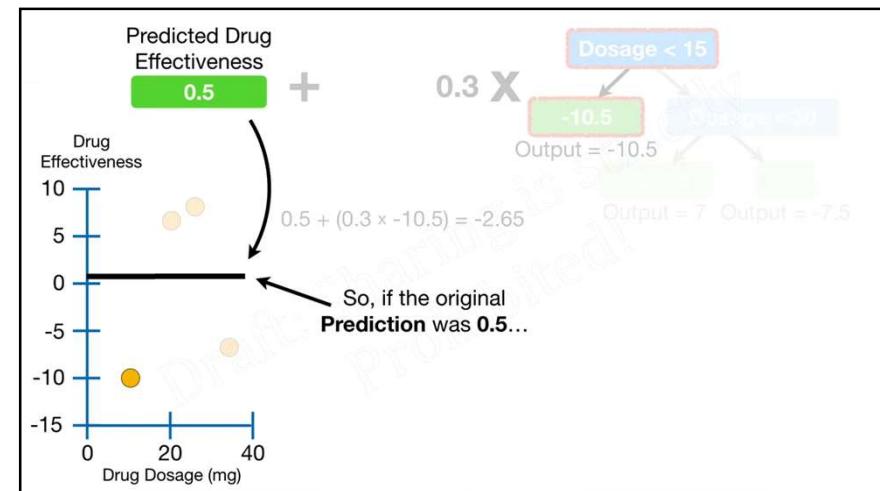
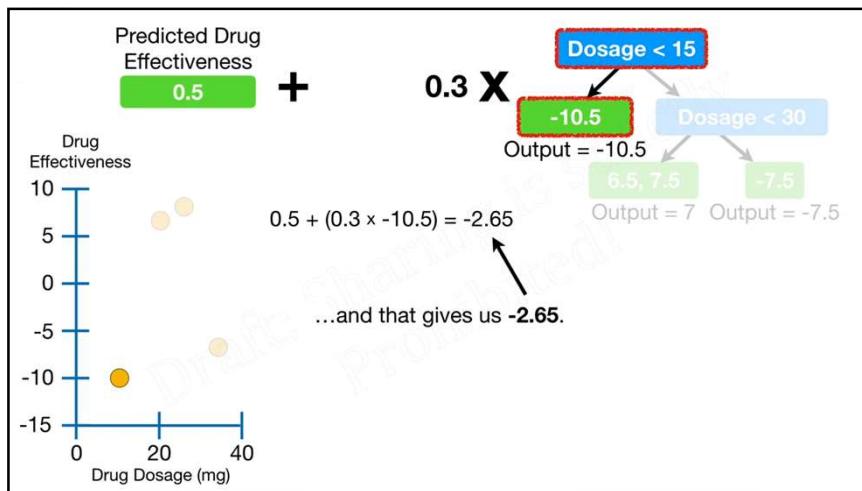


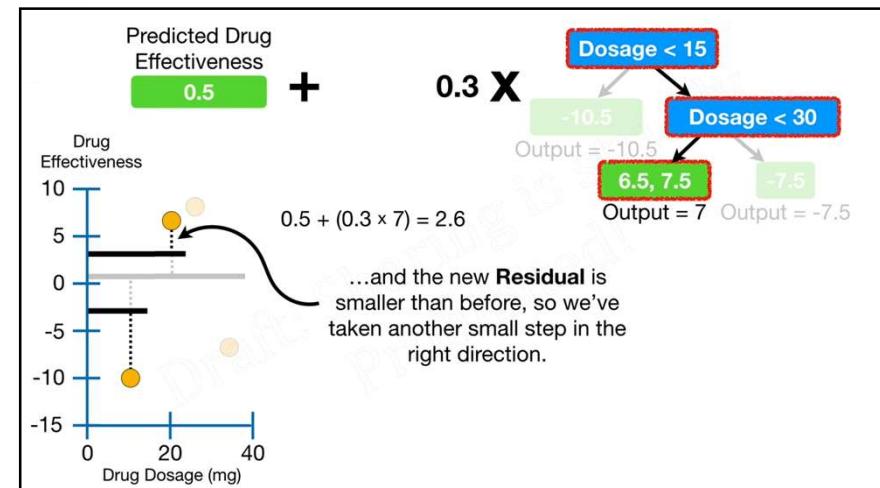
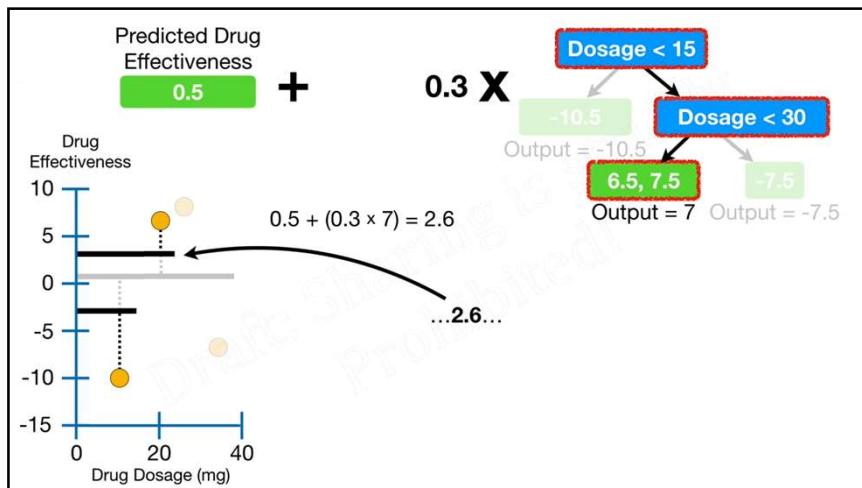
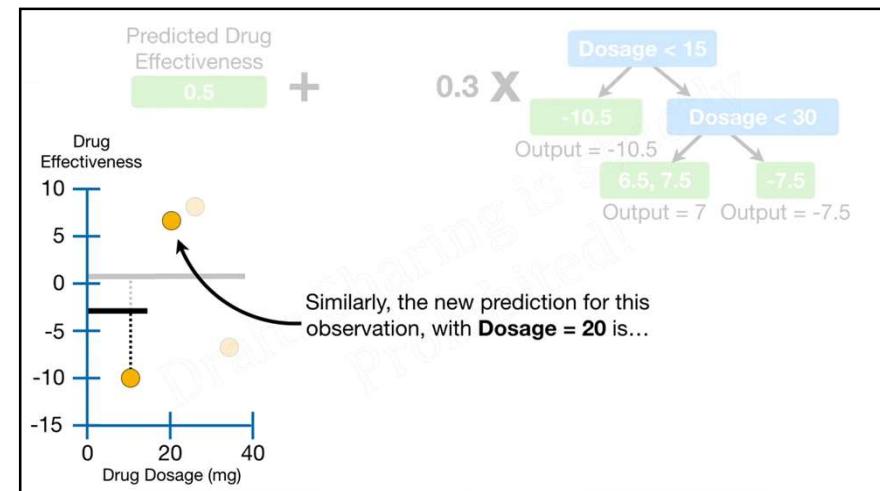
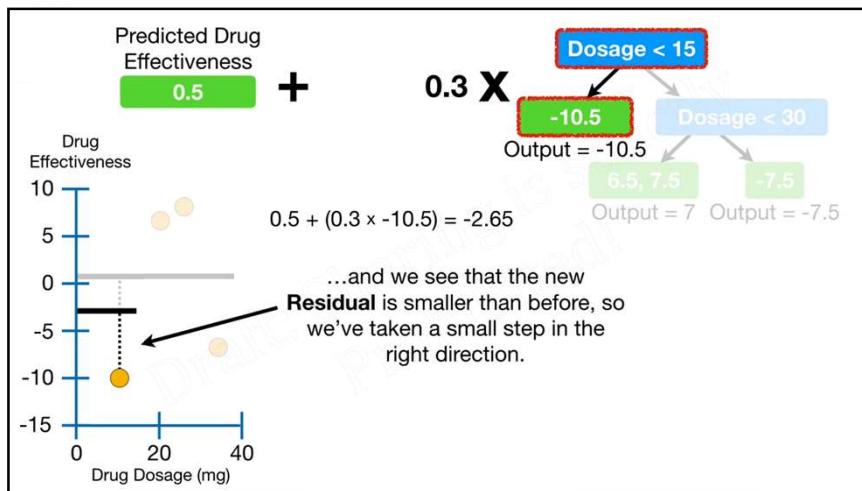


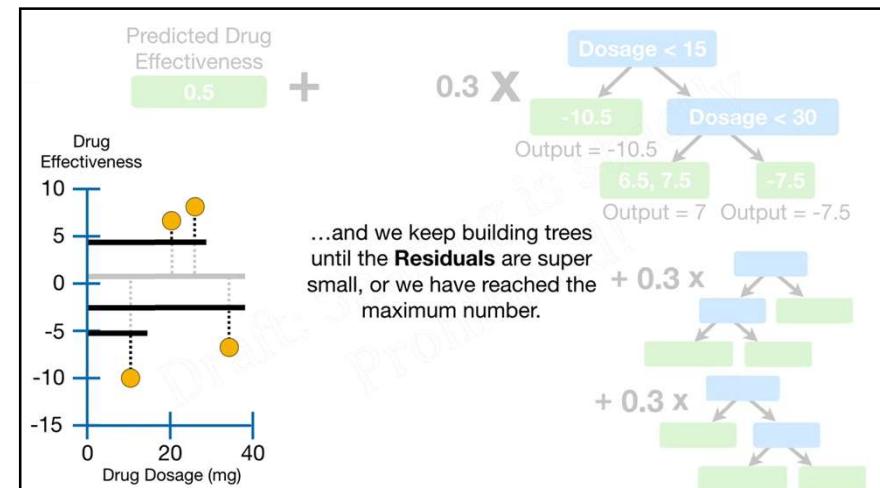
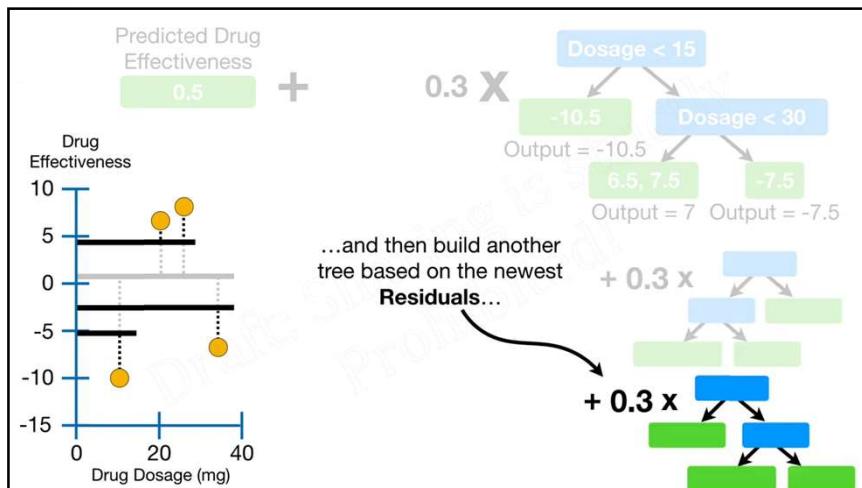
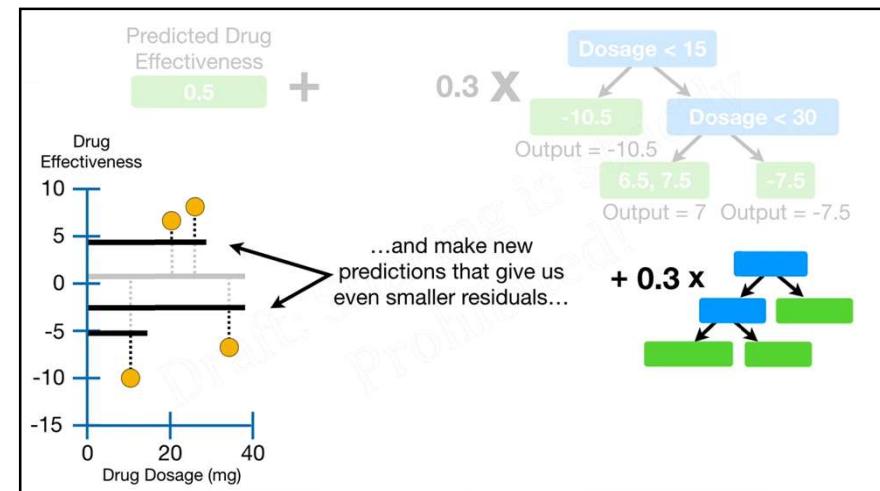
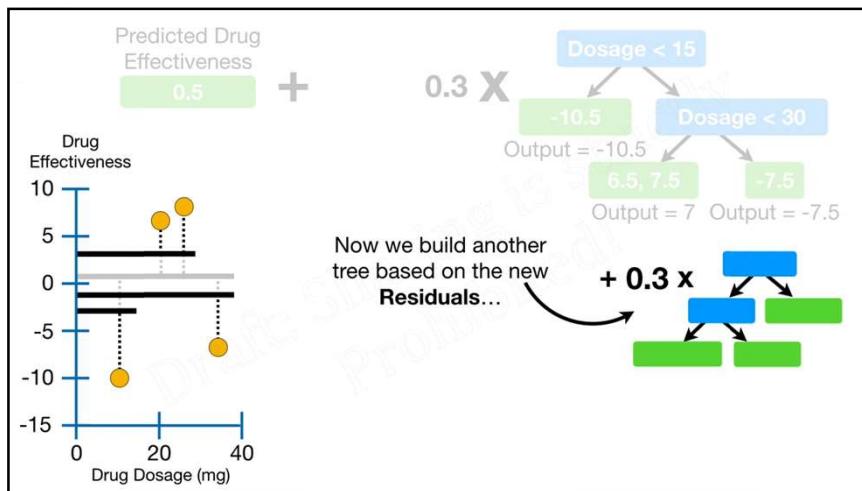












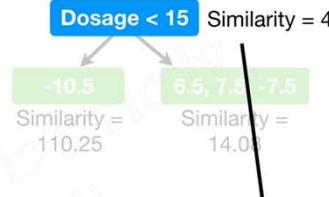
In summary, when building **XGBoost Trees for Regression...**

...we calculate **Similarity Scores...**

$$\text{Similarity Score} = \frac{\text{Sum of Residuals, Squared}}{\text{Number of Residuals} + \lambda}$$



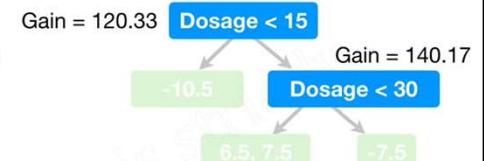
... and **Gain** to determine how to split the data...

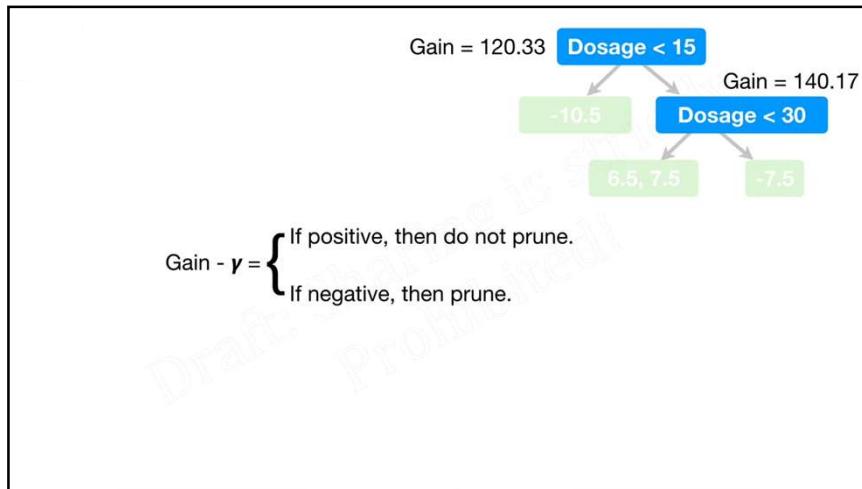


$$\text{Gain} = \text{LeftSimilarity} + \text{RightSimilarity} - \text{RootSimilarity}$$

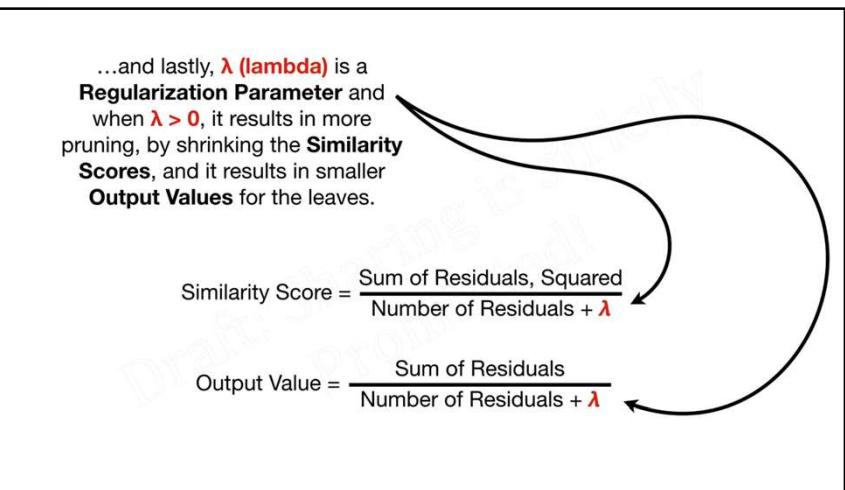
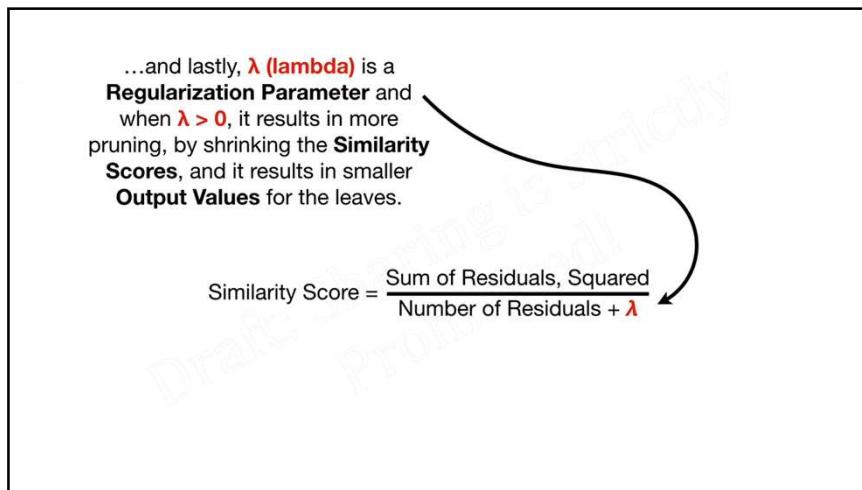
...and we prune the tree by calculating the differences between **Gain** values and a user defined **Tree Complexity Parameter, γ (gamma)**.

$$\text{Gain} - \gamma =$$





...and lastly, λ (lambda) is a **Regularization Parameter** and when $\lambda > 0$, it results in more pruning, by shrinking the **Similarity Scores**, and it results in smaller **Output Values** for the leaves.





THANK YOU

StatQuest with Josh Starmer
<http://www.youtube.com/watch?v=OJD8wVzPm6E>