

# Classification and Regression Trees... a bit more

Categorical and Numerical values  
in Classification trees

# Categorical and Numerical values in Classification trees (1/4)

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

Sort rows



Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

# Categorical and Numerical values in Classification trees (2/4)

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

Compute average for all adjacent rows



	Age	Loves Cool As Ice
9.5	7	No
15	12	No
26.5	18	Yes
36.5	35	Yes
44	38	Yes
66.5	50	No
	83	No

# Categorical and Numerical values in Classification trees (3/4)

	Age	Loves Cool As Ice
	7	No
9.5	12	No
15	18	Yes
26.5	35	Yes
36.5	38	Yes
44	50	No
66.5	83	No

Compute the Information Gain for all possible binary options.

Consider the value with the highest gain as representative of the feature

Age < 9.5

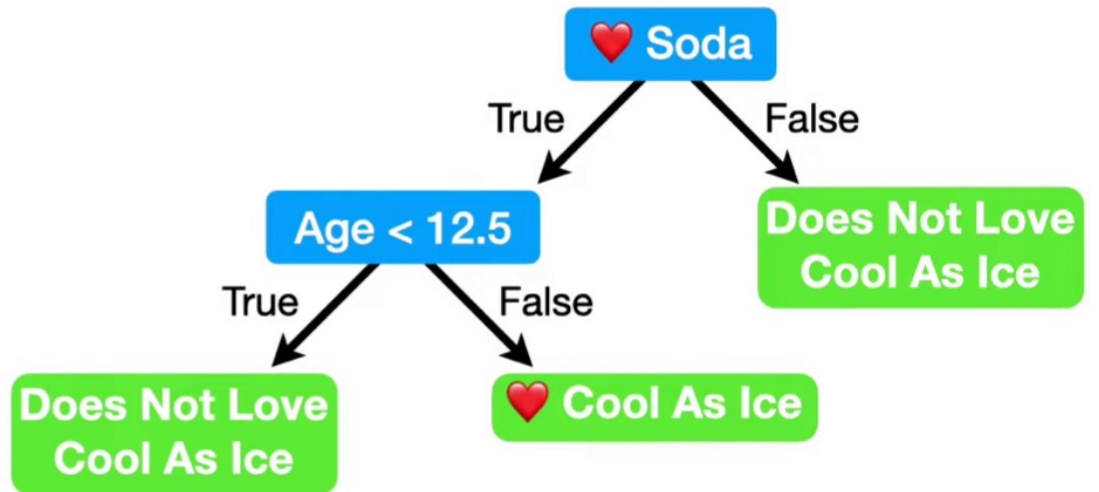
Age < 15

# Categorical and Numerical values in Classification trees (4/4)

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

Compare the Information Gain for all the features and select the one with the highest value.

Continue until you build the whole tree.



Missing values

# Missing categorical values

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...



# Missing categorical values

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	<b>YES</b>	Yes
etc...	etc...	etc...	etc...

Add the most frequent value

# Missing categorical values

Find a correlated feature and use it as guideline

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
No	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...

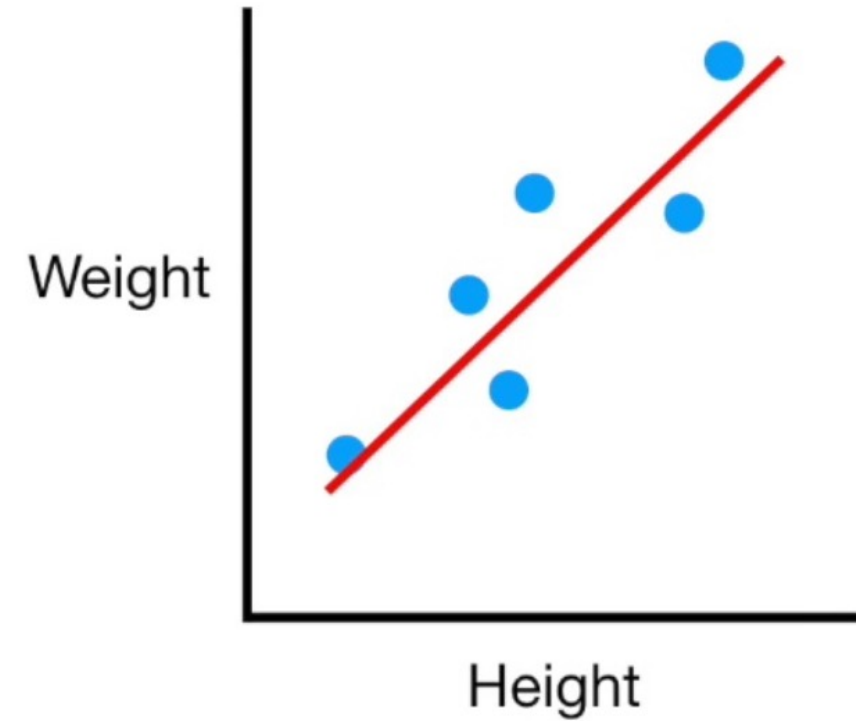
# Missing continuous values

Height	Good Blood Circulation	Weight	Heart Disease
5'7"	No	155	No
6'	Yes	180	Yes
5'4"	Yes	120	No
5'8"	No	???	Yes
etc...	etc...	etc...	etc...

could be not very useful

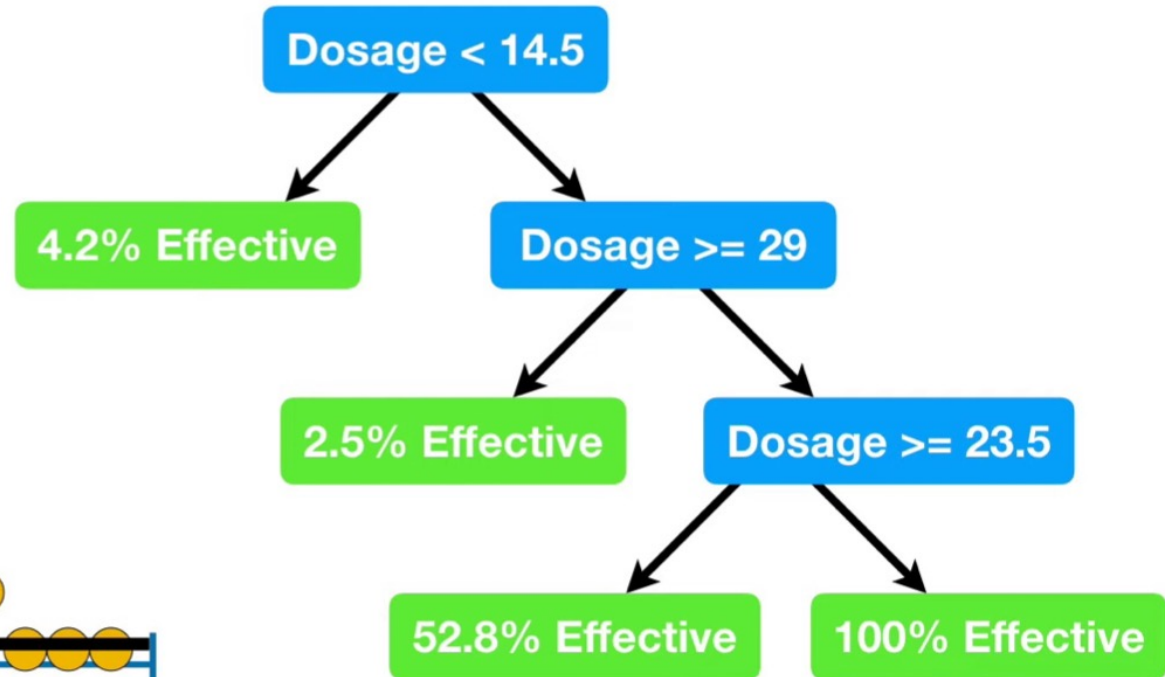
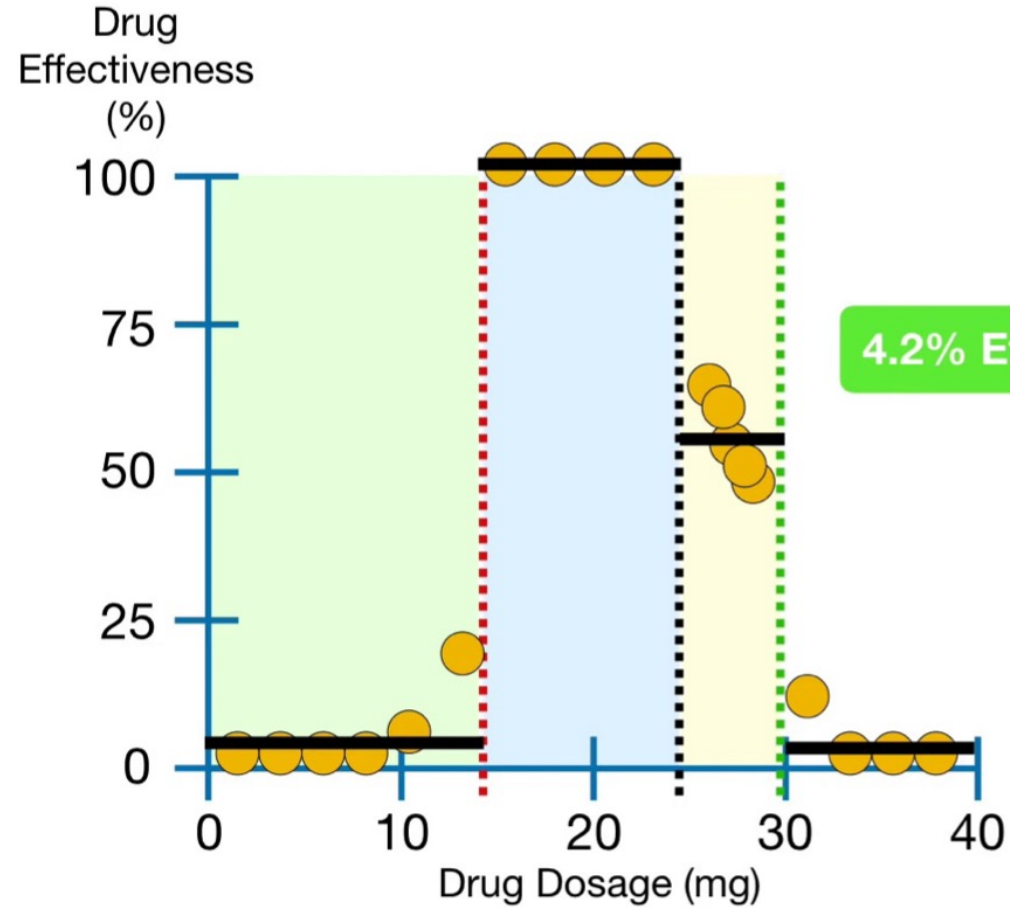
# Missing continuos values

Height	Good Blood Circulation	Weight	Heart Disease
5'7"	No	155	No
6'	Yes	180	Yes
5'4"	Yes	120	No
5'8"	No	???	Yes
etc...	etc...	etc...	etc...

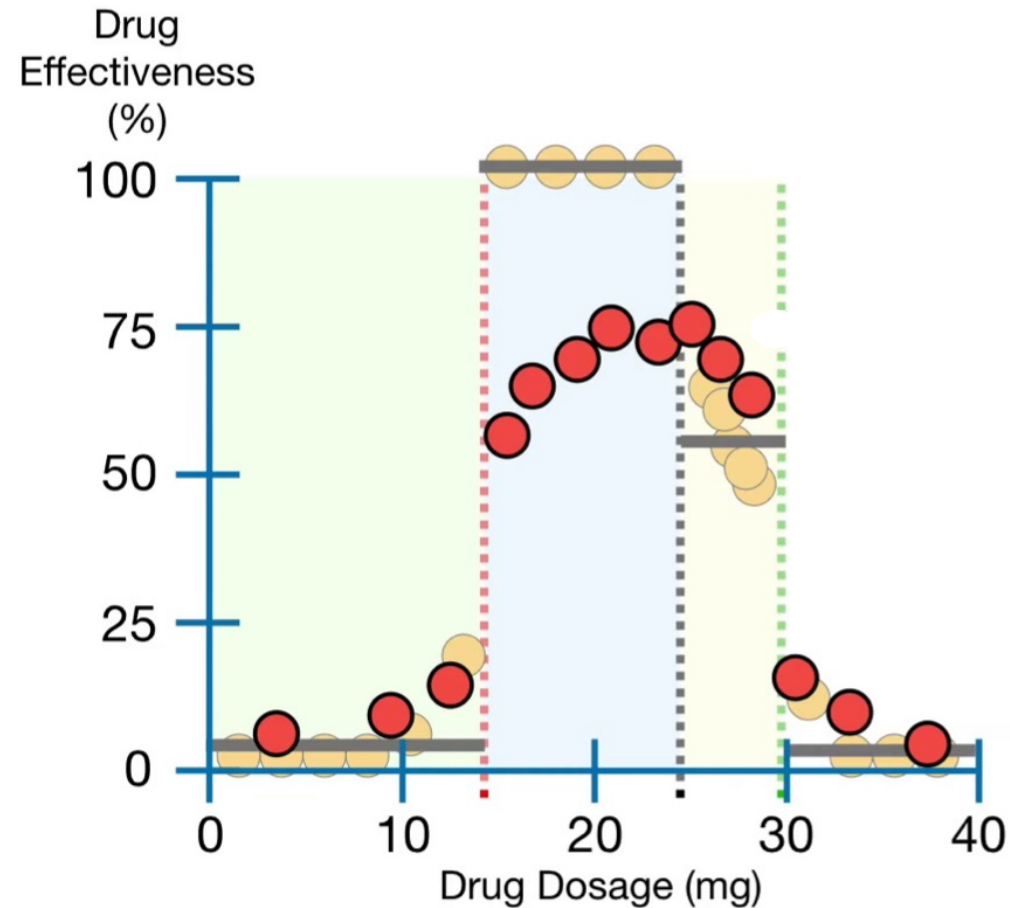


# Pruning Regression Trees

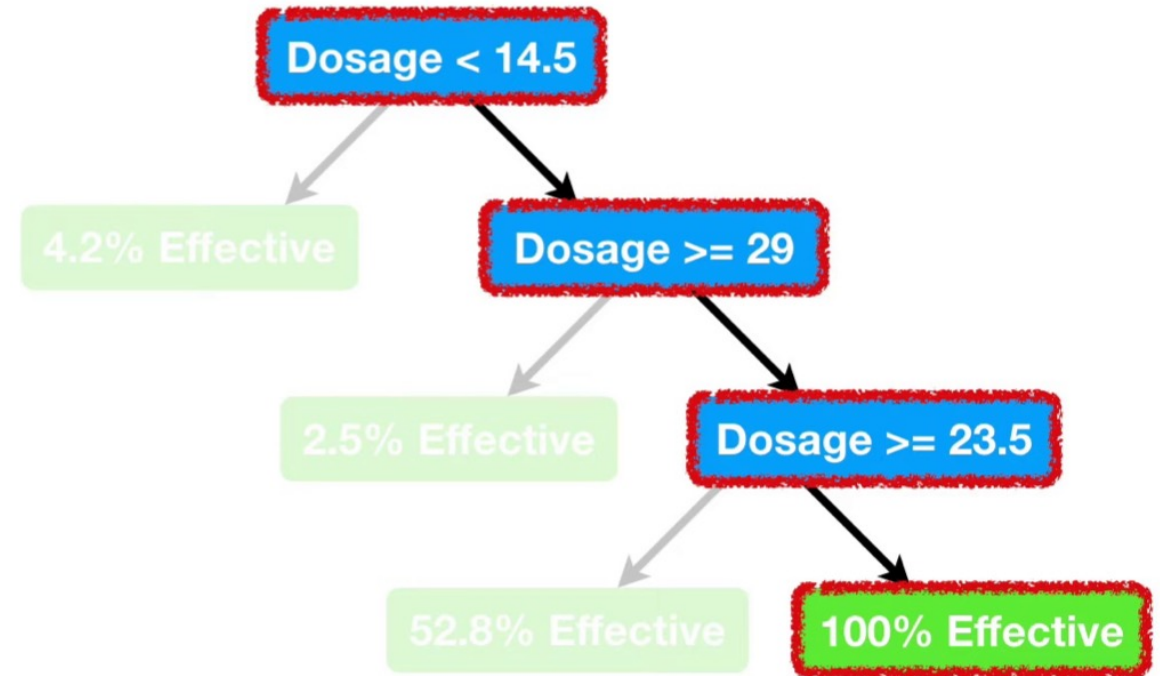
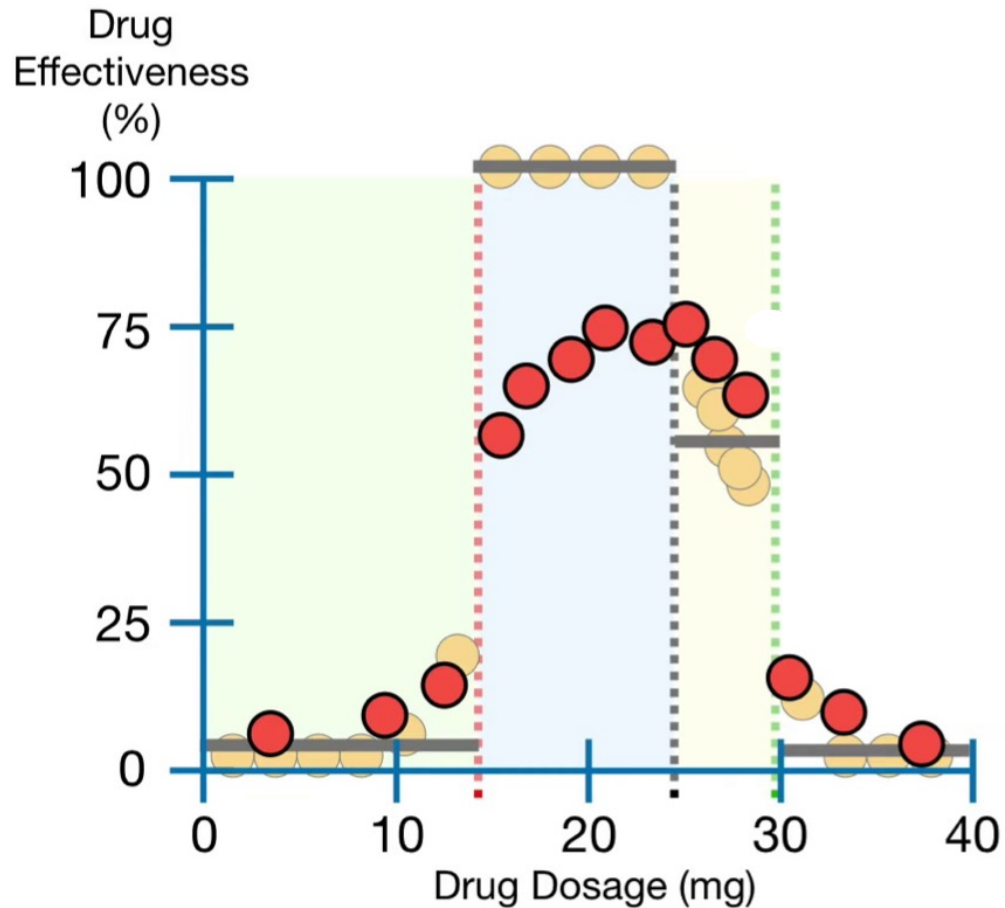
# Regression Tree



# Residual Sum of Squares on the Test Set

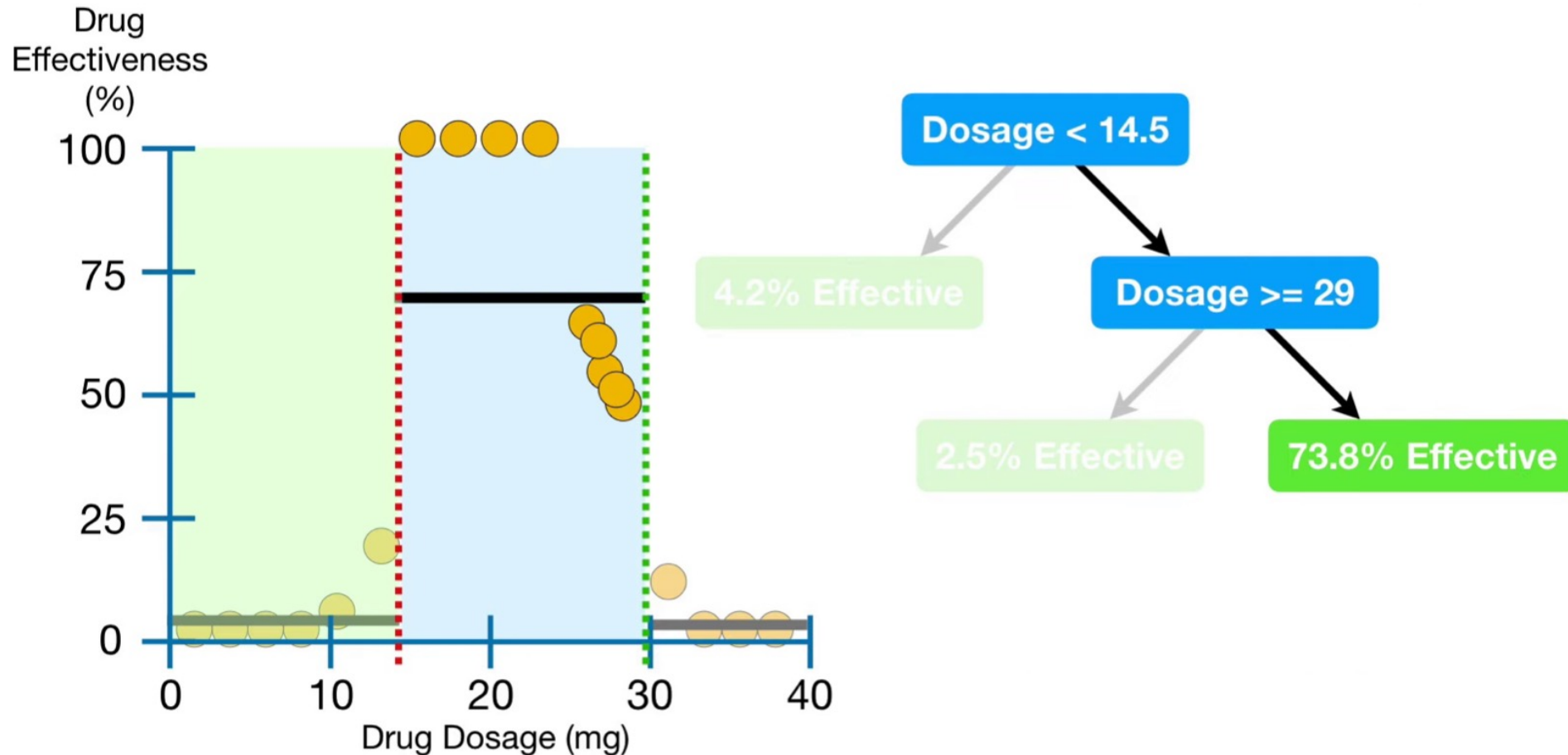


# Residual Sum of Squares on the Test Set

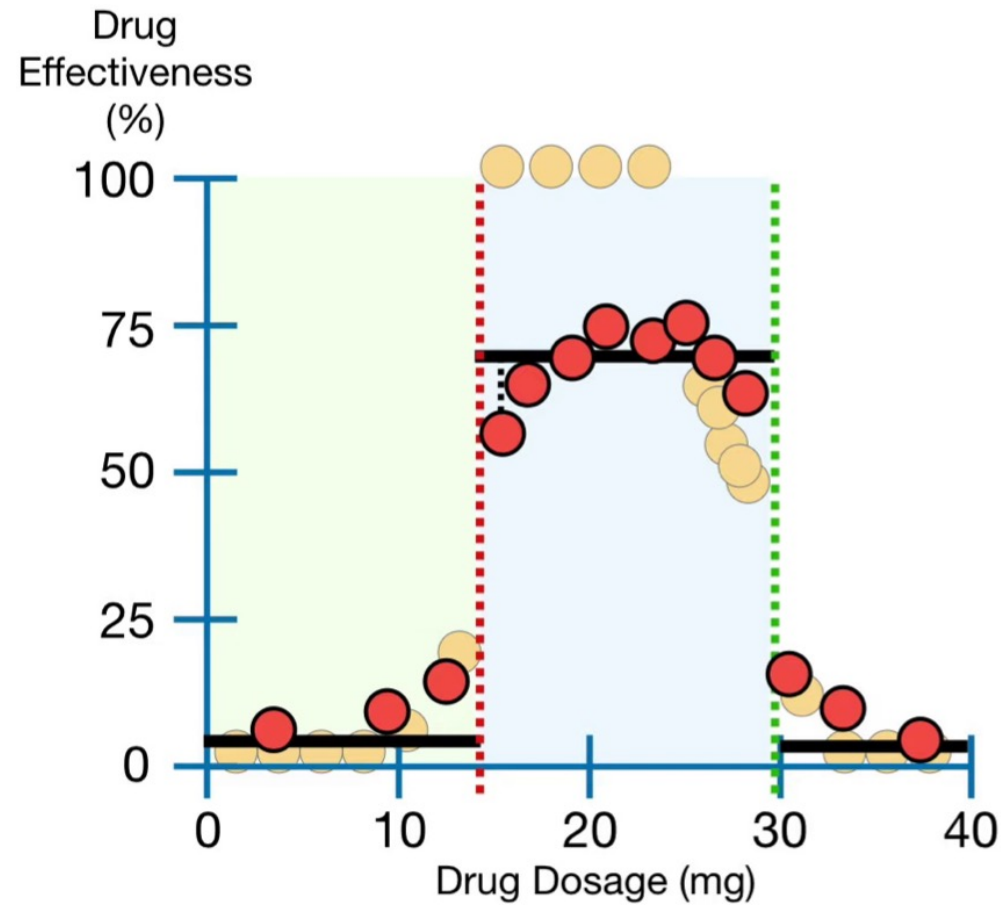




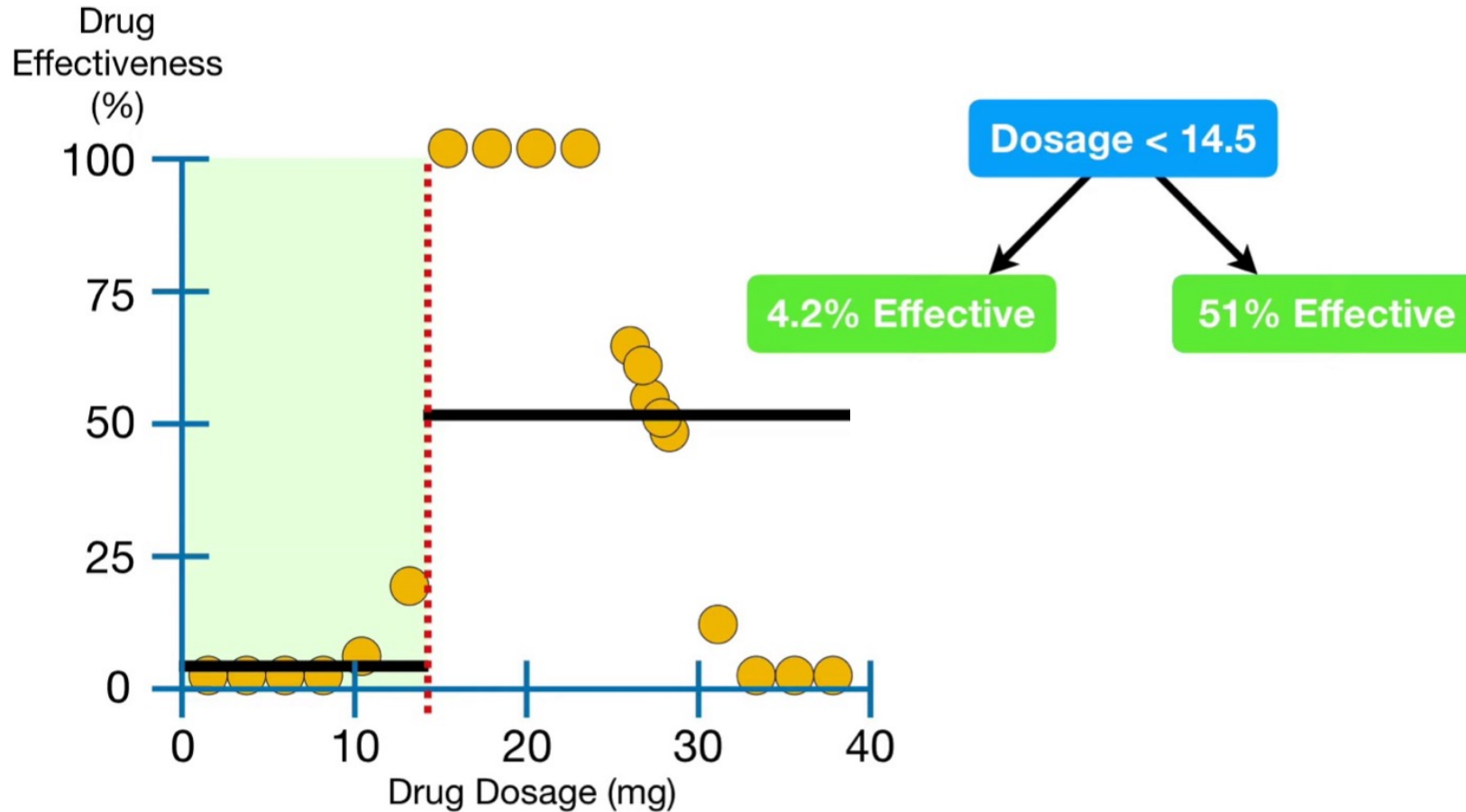
# Reduce the precision of the tree



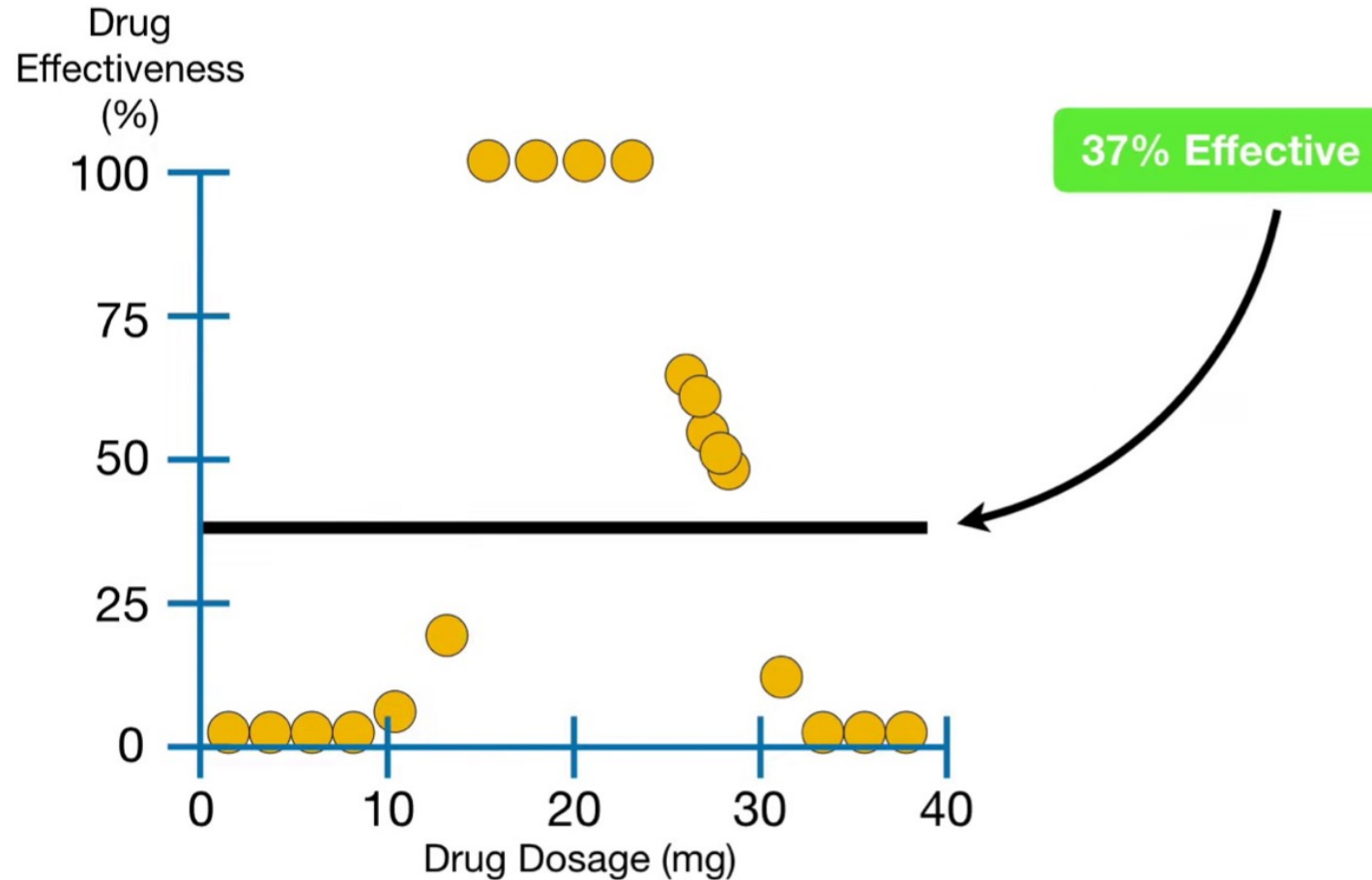
# Reduce overfitting



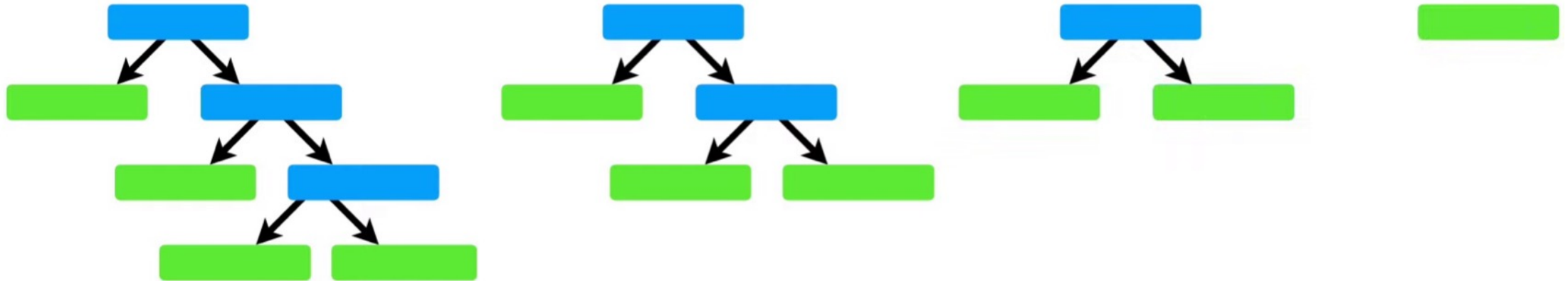
# Further pruning



# Further pruning

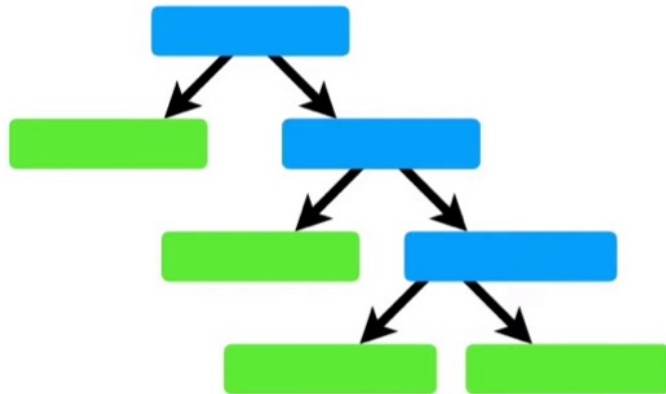


# What is the best tree?

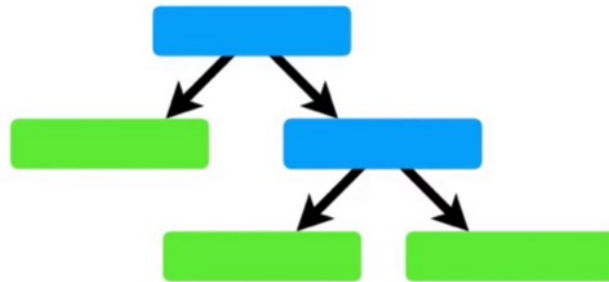


Compute the RSS for each tree

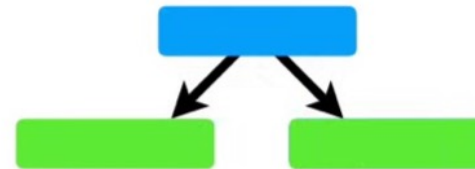
**RSS = 543.8**



**RSS = 5494.8**



**RSS = 19563.7**

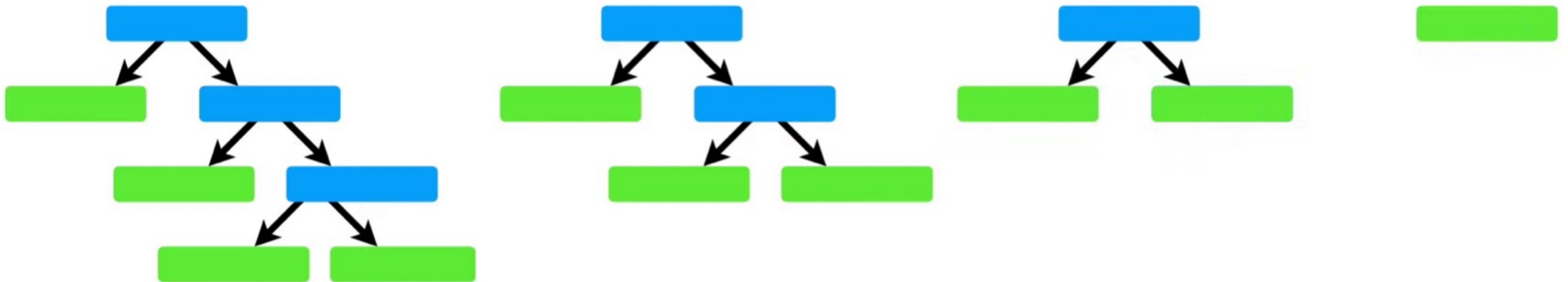


**RSS = 28897.2**



# Weakest Link Pruning

$$\begin{aligned}\text{Tree Score} &= \text{RSS} + \text{Tree Complexity Penalty} \\ &= \text{RSS} + \alpha \cdot T\end{aligned}$$

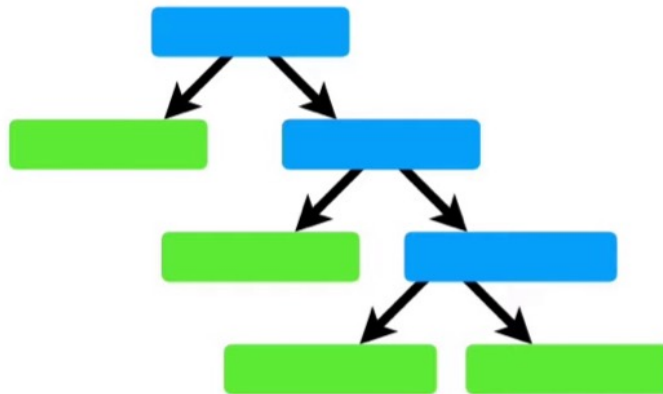


# Weakest Link Pruning

Tree Score =  $RSS + \alpha \cdot T$

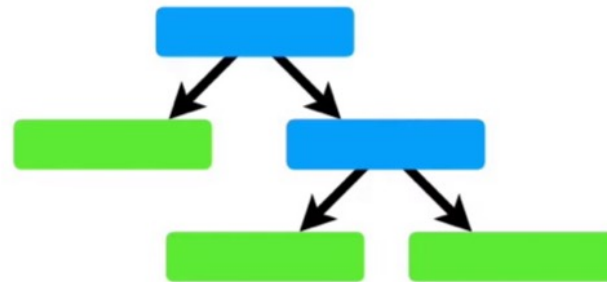
Example with  $\alpha = 10.000$

**RSS = 543.8**



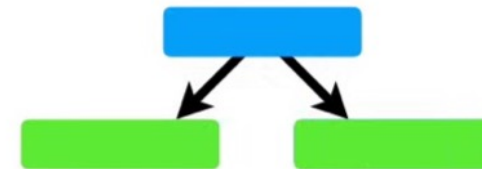
**Tree Score = 40,543.8**

**RSS = 5494.8**



**Tree Score = 35,494.8**

**RSS = 19563.7**



**Tree Score = 39,243.7**

**RSS = 28897.2**



**Tree Score = 39,243.7**

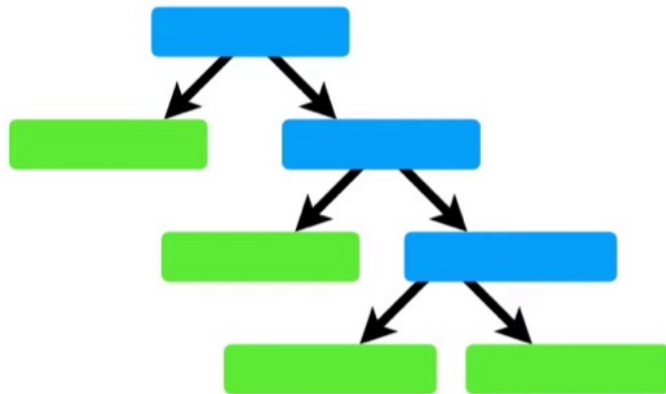


# Weakest Link Pruning

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

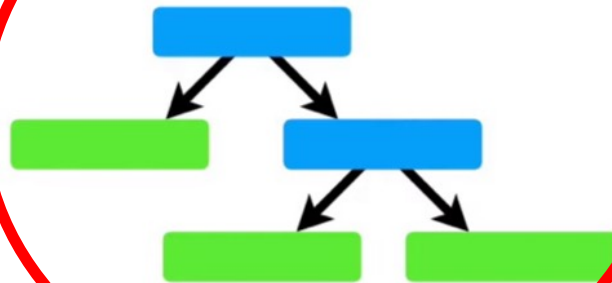
Example with  $\alpha = 10.000$

**RSS = 543.8**



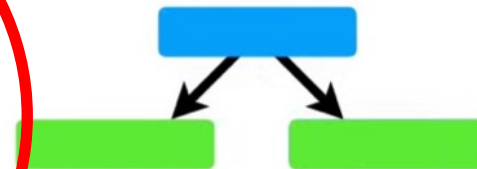
**Tree Score = 40,543.8**

**RSS = 5494.8**



**Tree Score = 35,494.8**

**RSS = 19563.7**



**Tree Score = 39,243.7**

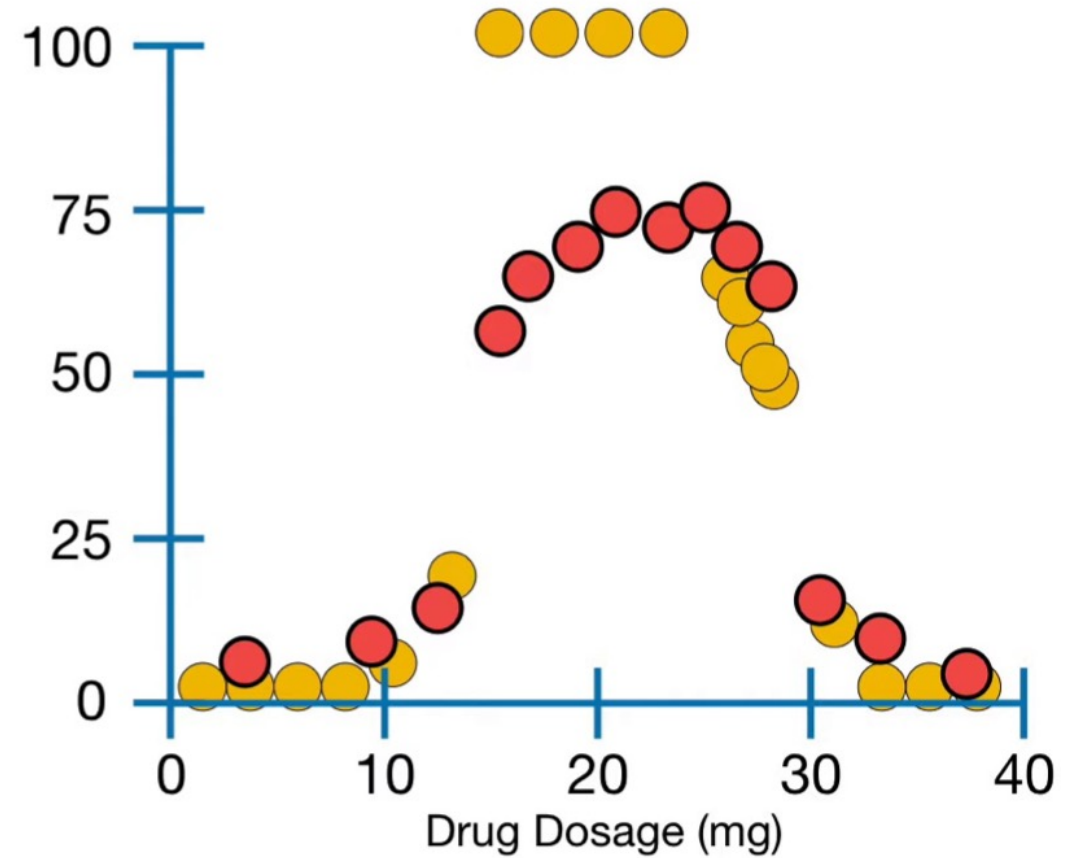
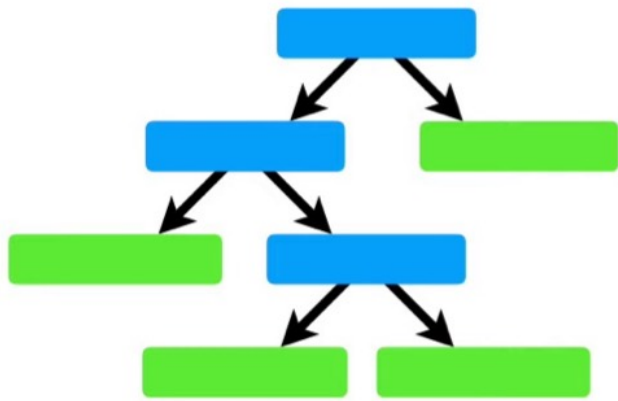
**RSS = 28897.2**



**Tree Score = 38,897.2**

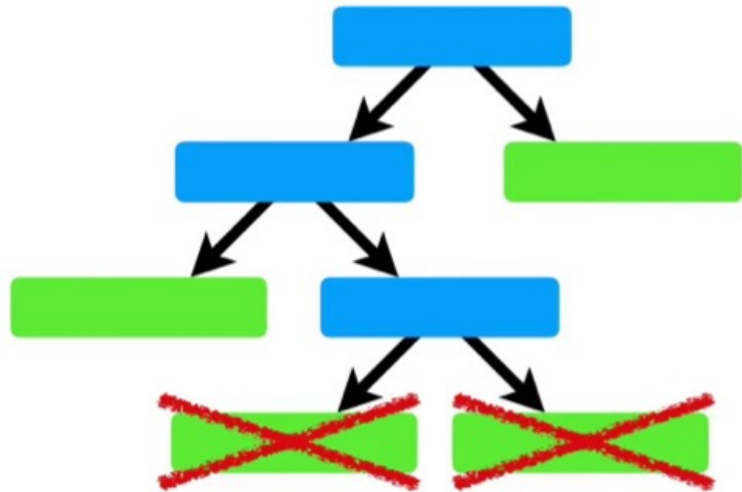
# How to evaluate the best $\alpha$

Build a tree considering all the data

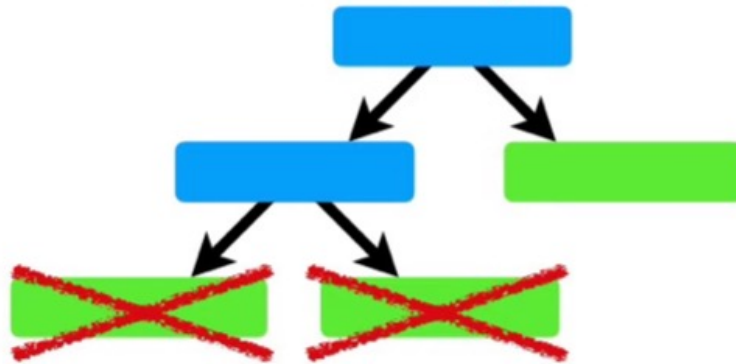


# How to evaluate the best $\alpha$

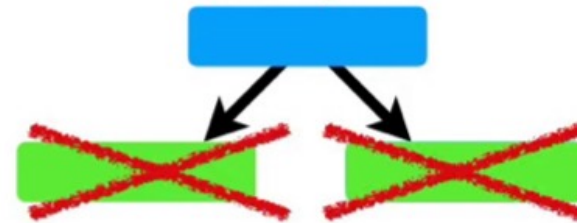
$\alpha = 0$



$\alpha = 10,000$



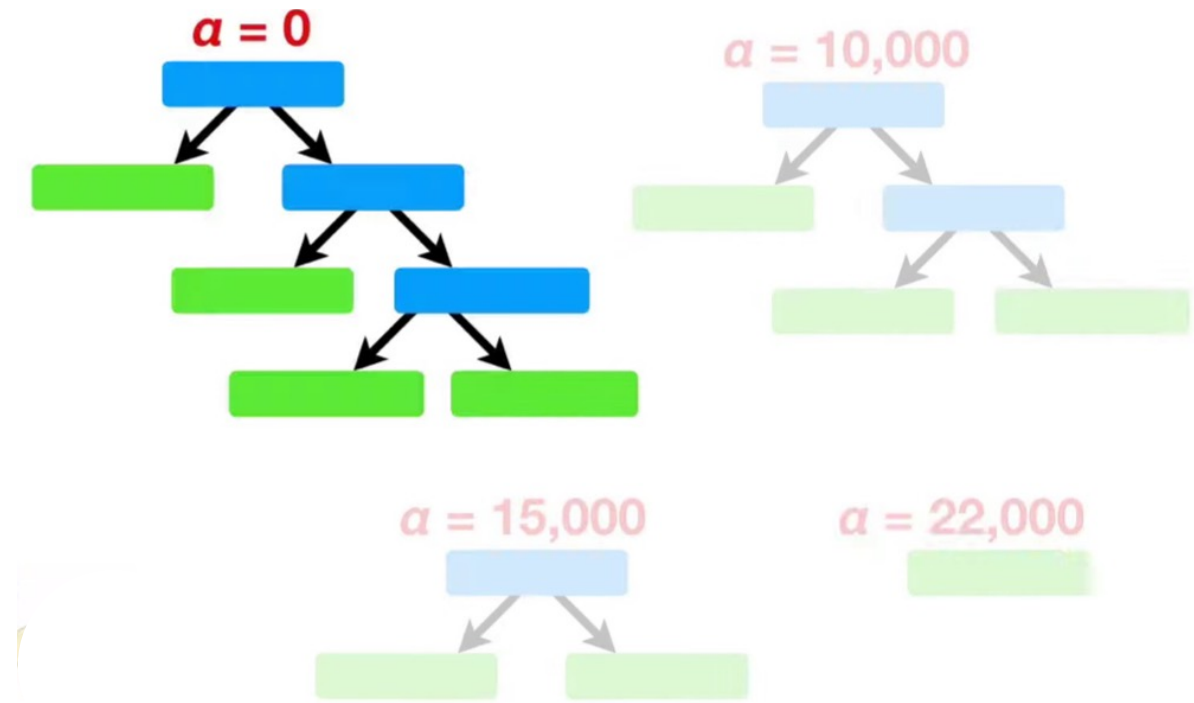
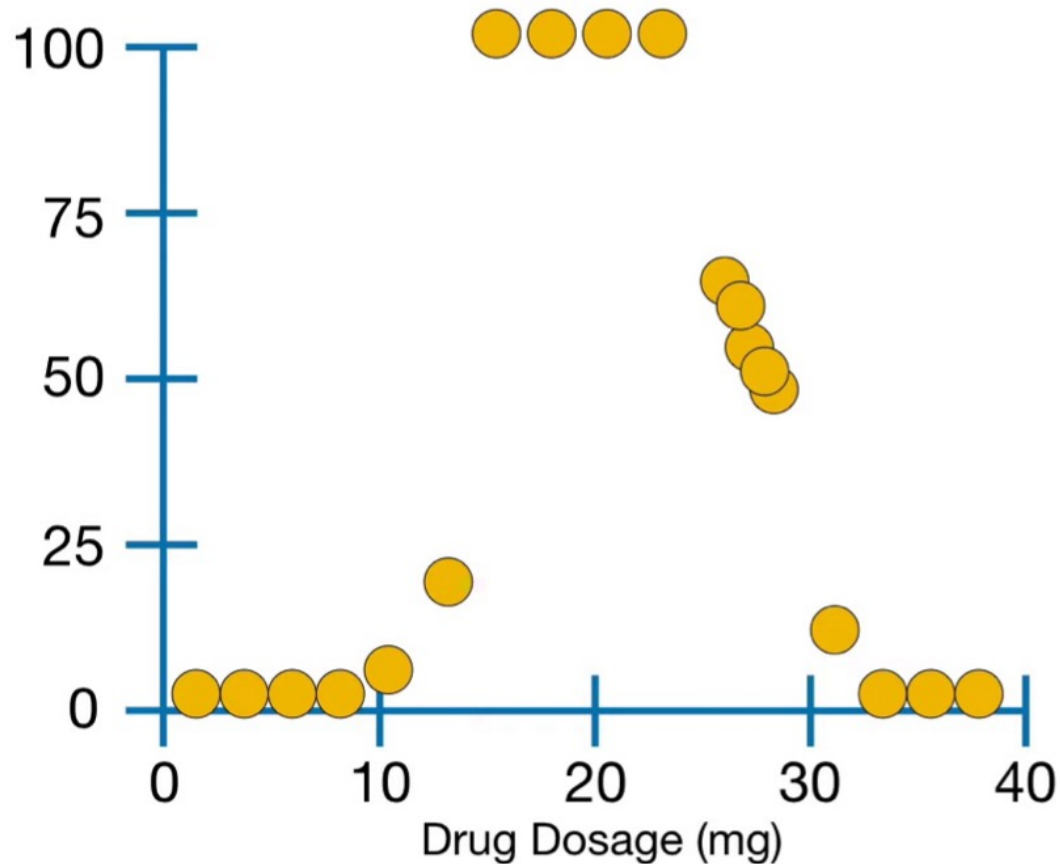
$\alpha = 15,000$



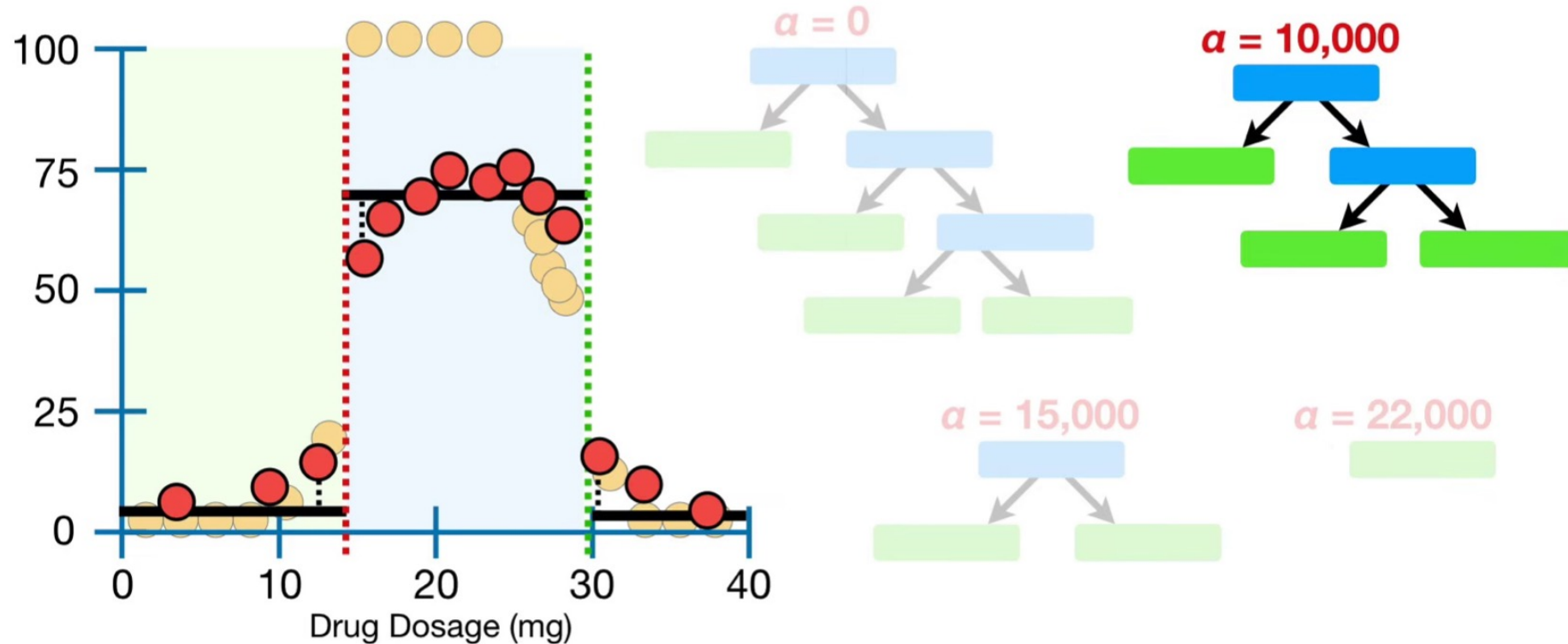
$\alpha = 22,000$



# Train the tree again using the Training set only

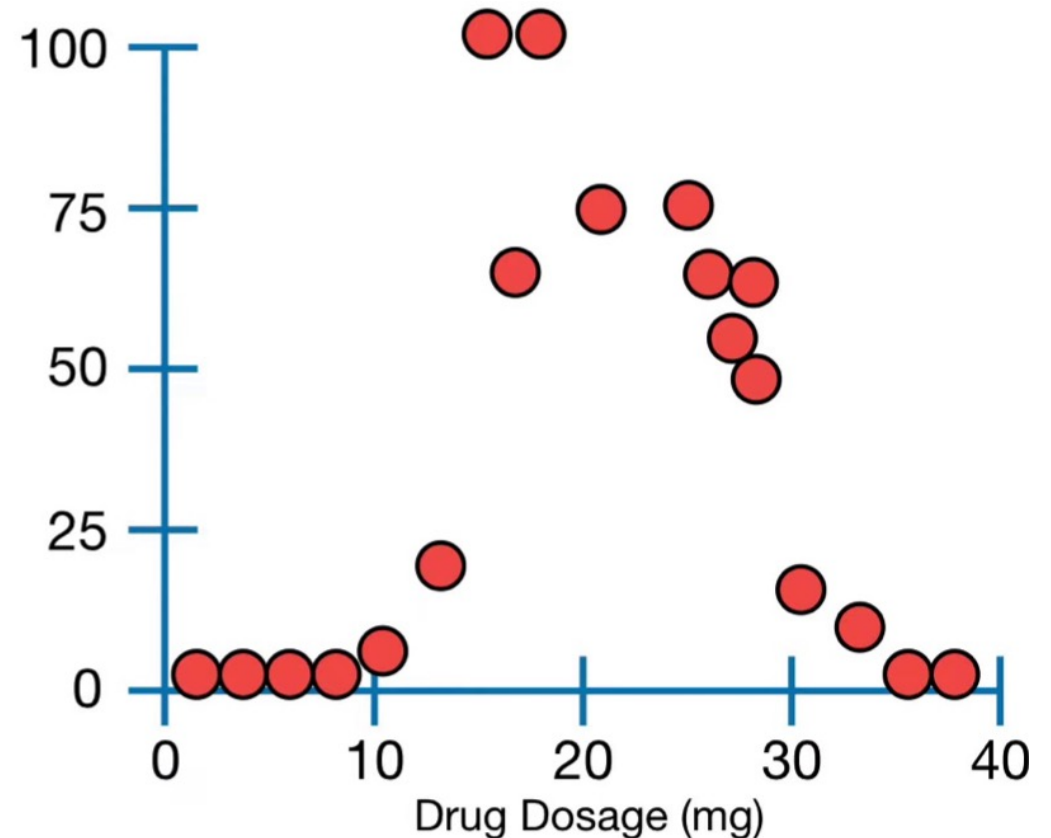
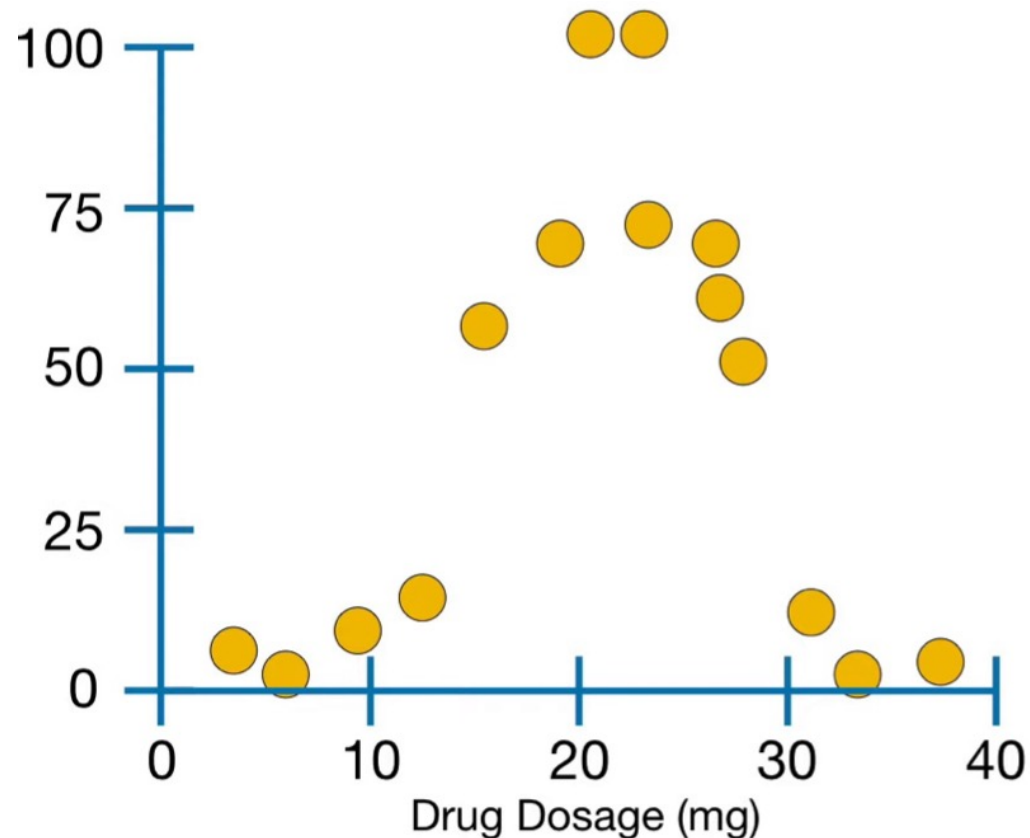


Compute RSS on the Test set for all the trees



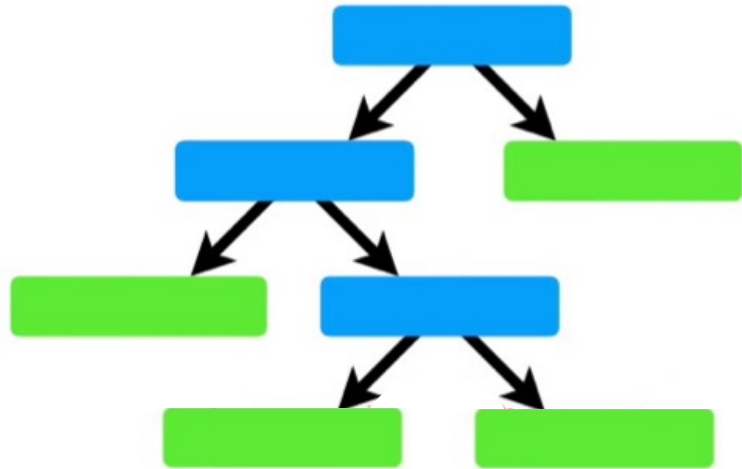
Vote for the one with the lowest RSS

Repeat the process with new Training and Test sets → k-fold Cross-Validation

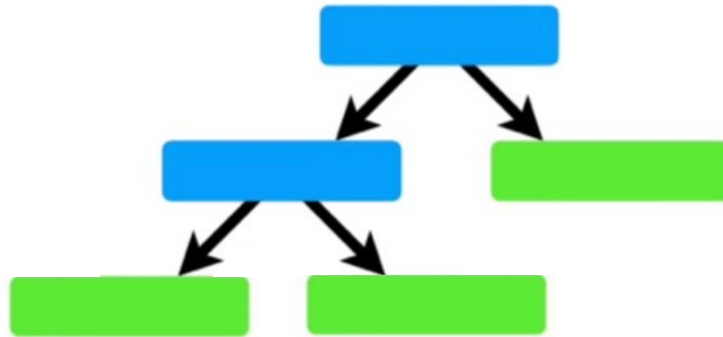


Go back to the trees built on the full dataset and select the one with the most voted  $\alpha$

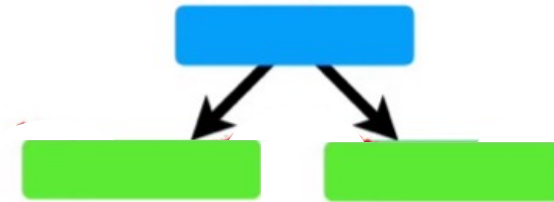
$\alpha = 0$



$\alpha = 10,000$



$\alpha = 15,000$



$\alpha = 22,000$

