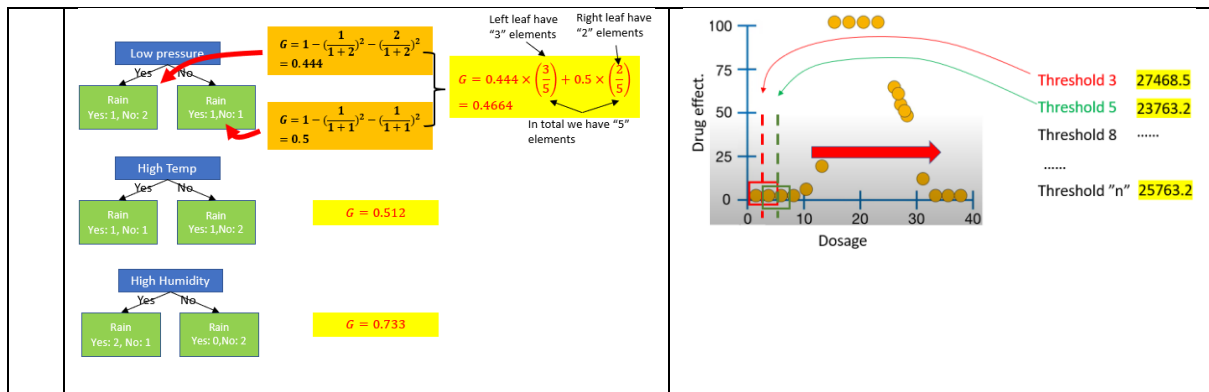
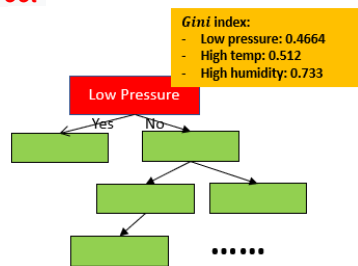


Summary: Decision Tree

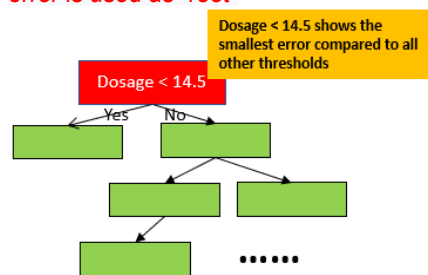
	Classification	Regression																																				
Sample data	<table><tr><th>Low pressure</th><th>High Temperature</th><th>High humidity</th><th>Rain</th></tr><tr><td>No</td><td>No</td><td>No</td><td>No</td></tr><tr><td>Yes</td><td>Yes</td><td>Yes</td><td>Yes</td></tr><tr><td>Yes</td><td>Yes</td><td>No</td><td>No</td></tr><tr><td>Yes</td><td>No</td><td>Yes</td><td>No</td></tr><tr><td>No</td><td>No</td><td>Yes</td><td>Yes</td></tr></table>	Low pressure	High Temperature	High humidity	Rain	No	No	No	No	Yes	Yes	Yes	Yes	Yes	Yes	No	No	Yes	No	Yes	No	No	No	Yes	Yes	<table><tr><th>Dosage</th><th>Drug effect.</th></tr><tr><td>10</td><td>58</td></tr><tr><td>20</td><td>60</td></tr><tr><td>35</td><td>57</td></tr><tr><td>5</td><td>44</td></tr><tr><td>...</td><td>...</td></tr></table>	Dosage	Drug effect.	10	58	20	60	35	57	5	44
Low pressure	High Temperature	High humidity	Rain																																			
No	No	No	No																																			
Yes	Yes	Yes	Yes																																			
Yes	Yes	No	No																																			
Yes	No	Yes	No																																			
No	No	Yes	Yes																																			
Dosage	Drug effect.																																					
10	58																																					
20	60																																					
35	57																																					
5	44																																					
...	...																																					
1		<p>Sort the predictors from small values to big ones, and plot it out:</p>																																				
2	<p>Grow trees for each predictor individually</p> <p>e.g., when we have "Low" pressure, there are:</p> <ul style="list-style-type: none">1 case has rain,2 cases do not have rain	<p>Grow trees for each threshold of every predictor: the threshold is the "mean" value for neighbouring points</p>																																				
3	<p>Obtain the "Gini" index for each tree</p> <p>The "Gini" index can be calculated as</p> $G = 1 - P_{yes}^2 - P_{no}^2$ <p>Where:</p> <ul style="list-style-type: none">P_{yes} is the probability of "yes" in a leafP_{no} is the probability of "no" in a leaf	<p>Obtain the error for each tree (or threshold)</p> <p>Using the tree (threshold based) to produce the prediction and calculate the error as</p> $(Pred - Actual)_{point1}^2 + (Pred - Actual)_{point2}^2 + \dots$																																				



4 **Based on the Gini index, determine the "root" to split the tree**
Usually, the predictor showing the smallest Gini is used as "root"

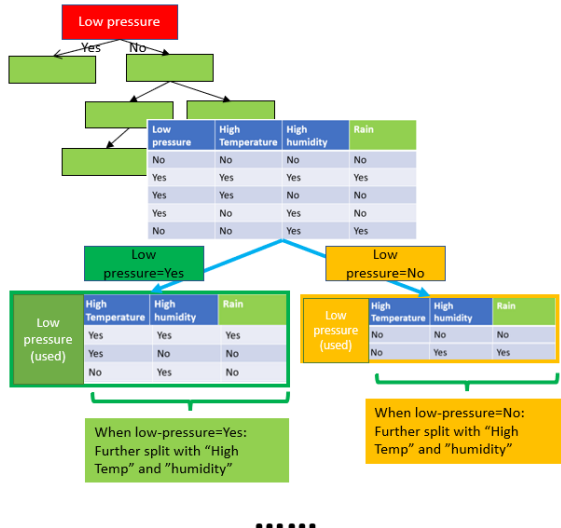


Based on the Error, determine the "root" to split the tree
Usually, the threshold showing the smallest error is used as "root"



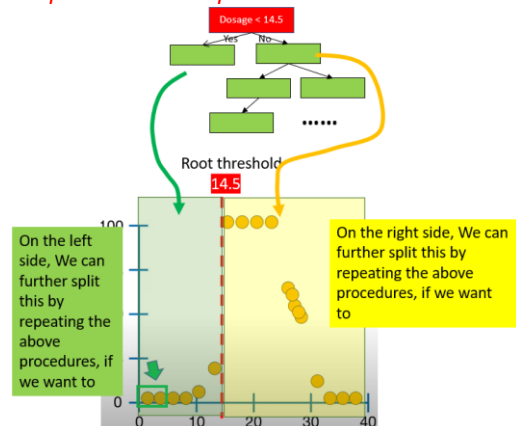
5 **Further split the tree using the similar method, until we are not able to split**

- Grow tree for every remained predictor
- Calculate "Gini" for each tree
- Get the tree with the smallest Gini
- Repeat the above process ...



Further split the tree using the similar method, until we are not able to split

- Calculate the threshold from neighbouring points, and grow trees based on thresholds
- Error for each tree
- Get the tree with the smallest error
- Repeat the above process ...



Note:

- In order to avoid overfitting, we may limit the number of levels of the tree to grow