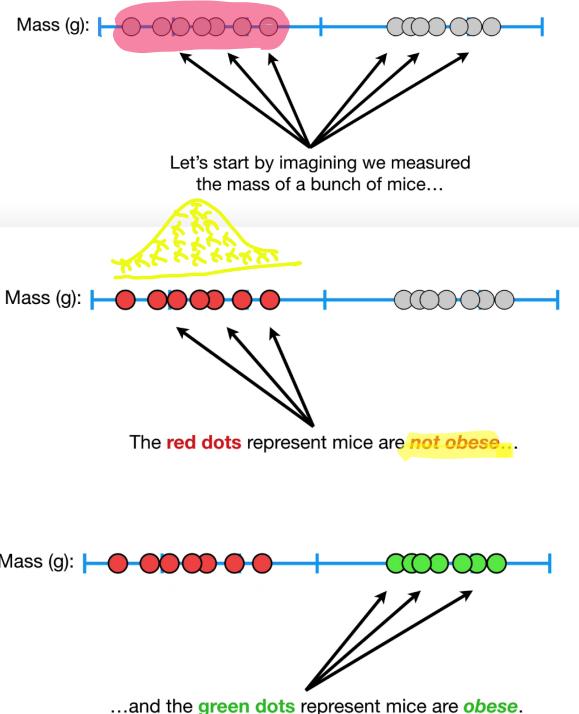
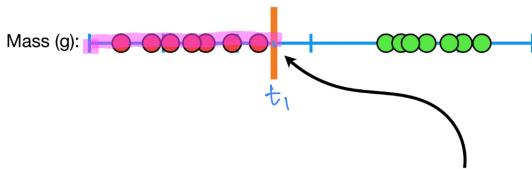


5 - Support Vector Machine - Overview

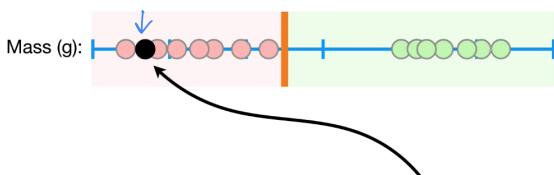
Sunday, 22 November BE 2563 19:41

External Resources: [Support Vector Machines, Clearly Explained!!!](#)

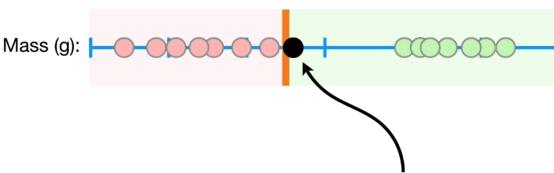




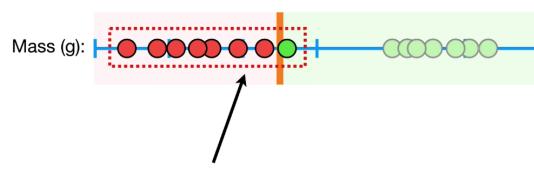
Based on these observations, we can pick a threshold...



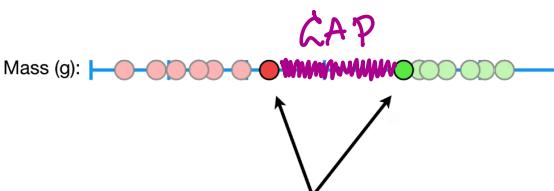
...and when we get a new observation that has less mass than the threshold...



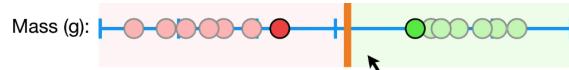
However, what if get a new observation here?



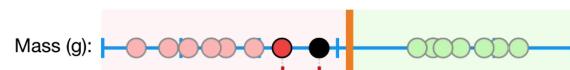
But that doesn't make sense, because it is much closer to the observations that are *not obese*.



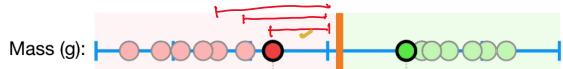
...we can focus on the observations on the edges of each cluster...



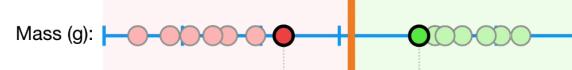
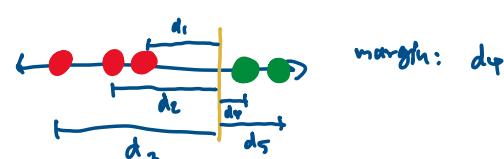
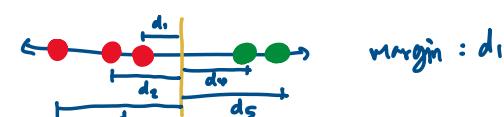
...and use the midpoint between them as the threshold.



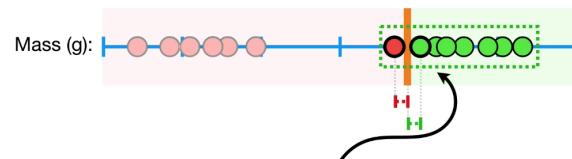
...it will be closer to the observations that are *not obese*...



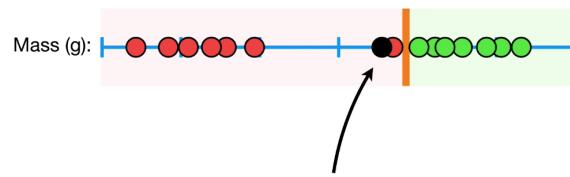
The shortest distance between the observations and the threshold is called the **margin**.



...we are using a **Maximal Margin Classifier**.



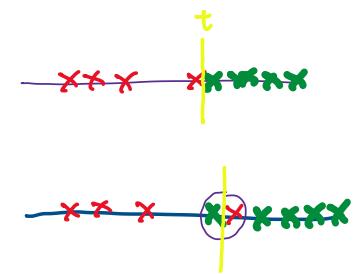
In this case, the **Maximum Margin Classifier** would be super close to the *obese* observations...



Now, if we got this new observation...

...we would classify it as *not obese*, even though most of the *not obese* observations are much further away than the *obese* observations.

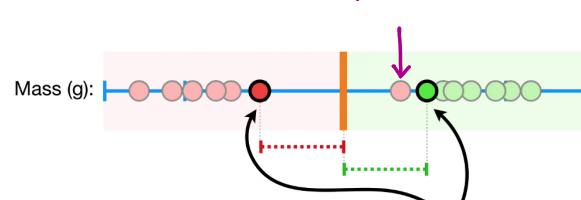
So **Maximal Margin Classifiers** are *super sensitive to outliers* in the training data and that makes them pretty lame.



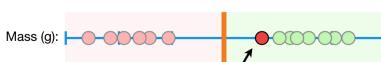
no possible t that can classify all data points correctly

To make a threshold that is not so sensitive to outliers we must **allow misclassifications**.

pretend that it is not here



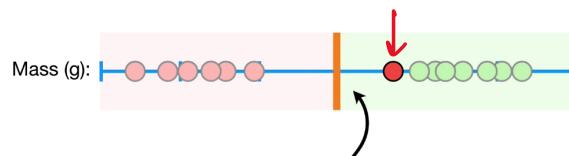
For example, if we put the threshold halfway between these two observations...



...then we will misclassify this observation.



However, now when we get a new observation here...



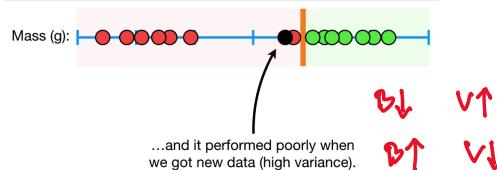
Choosing a threshold that allows misclassifications is an example of the Bias/Variance Tradeoff that plagues all of machine learning.

BIAS vs VARIANCE

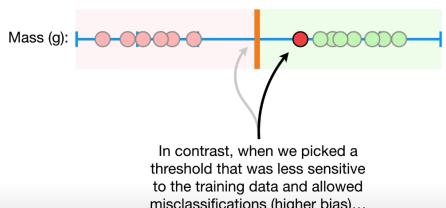
low bias = we just have to trust all the data points to do the job

high bias = we have some concerns/suggestions about the data points

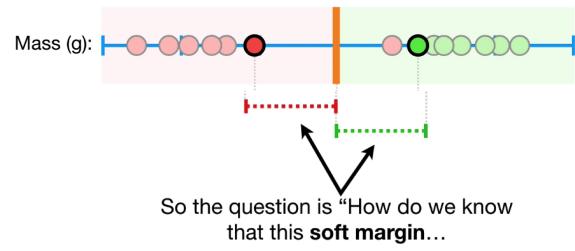
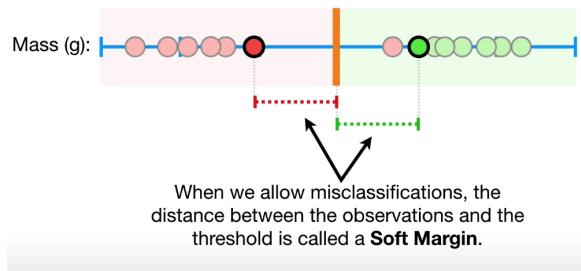
In other words, before we allowed misclassifications, we picked a threshold that was very sensitive to the training data (low bias)...



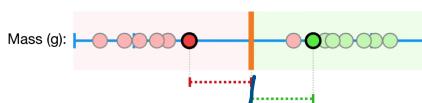
...and it performed poorly when we got new data (high variance).

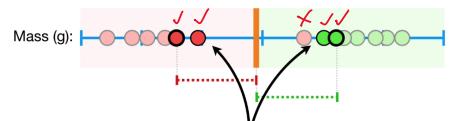
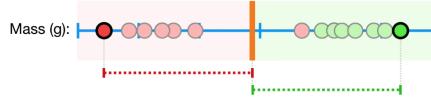
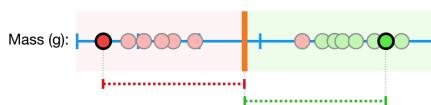
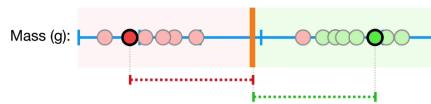
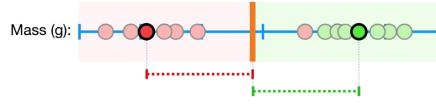
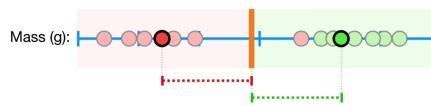
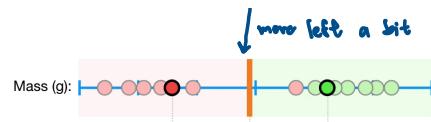


In contrast, when we picked a threshold that was less sensitive to the training data and allowed misclassifications (higher bias)...

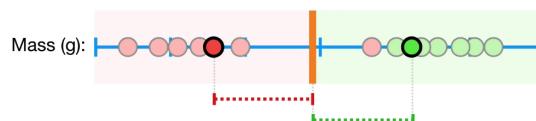


The answer is simple: We use **Cross Validation** to determine how many misclassifications and observations to allow inside of the **Soft Margin** to get the best classification.





...and two observations, that are correctly classified, to be within the **Soft Margin**.



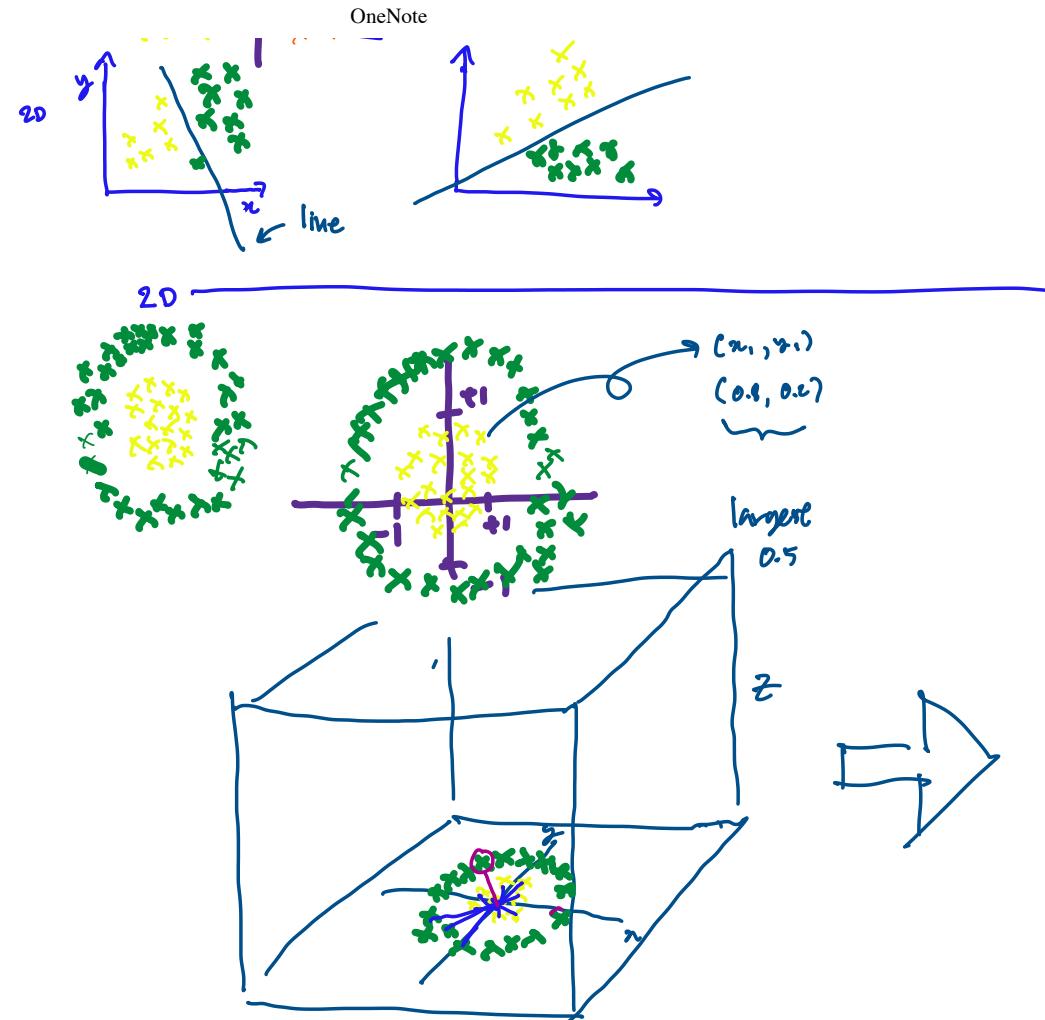
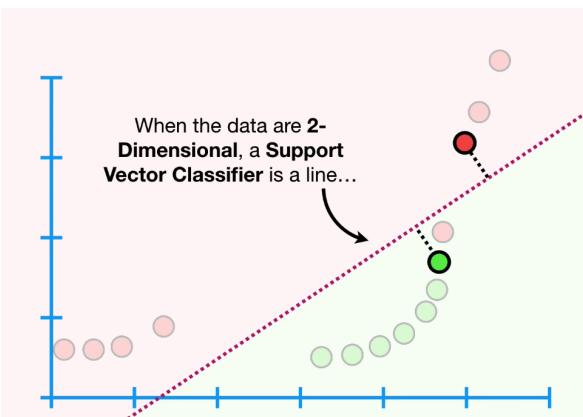
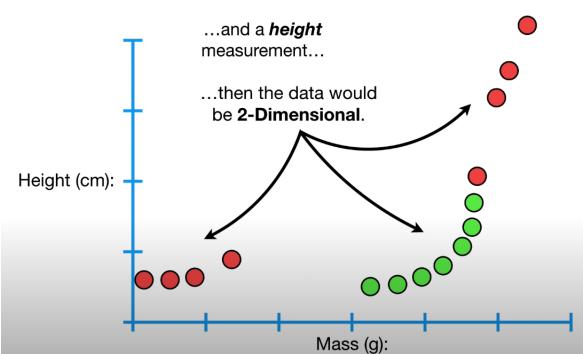
When we use a **Soft Margin** to determine the

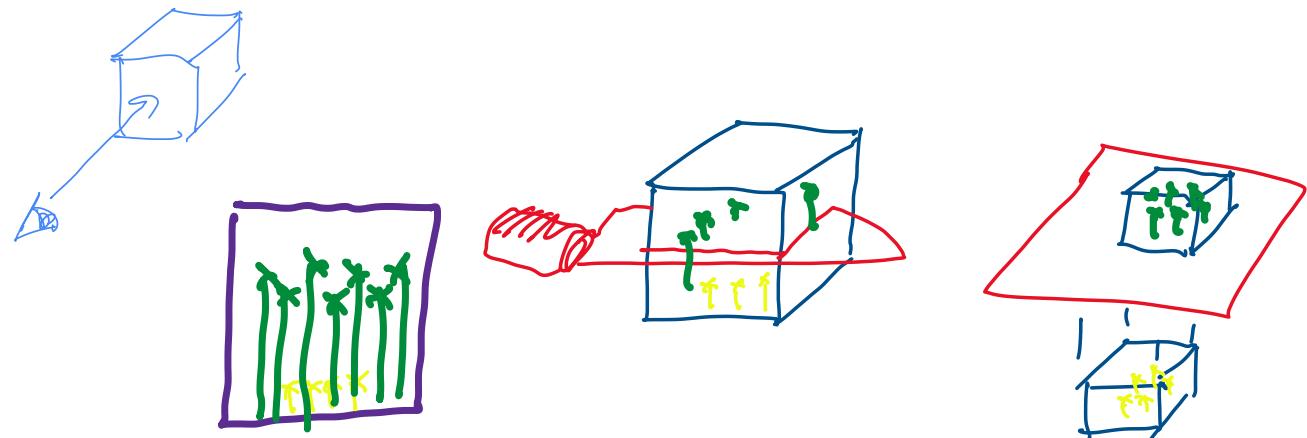
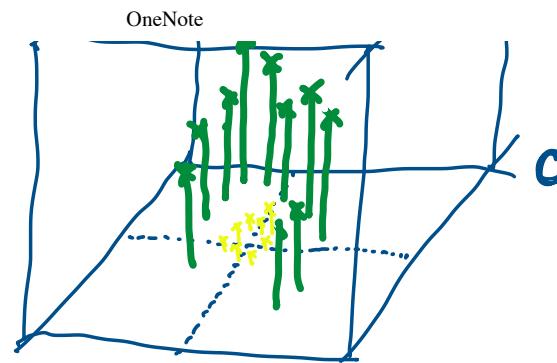
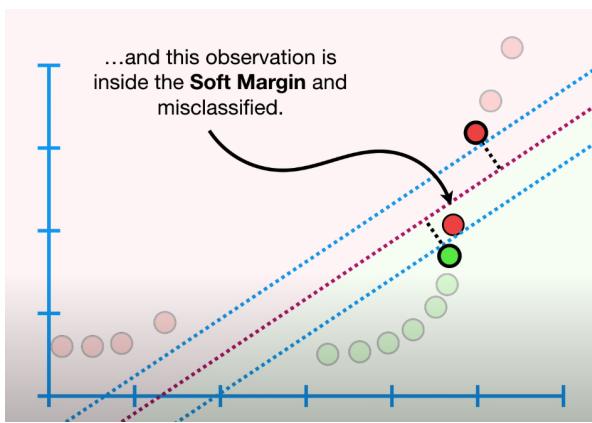
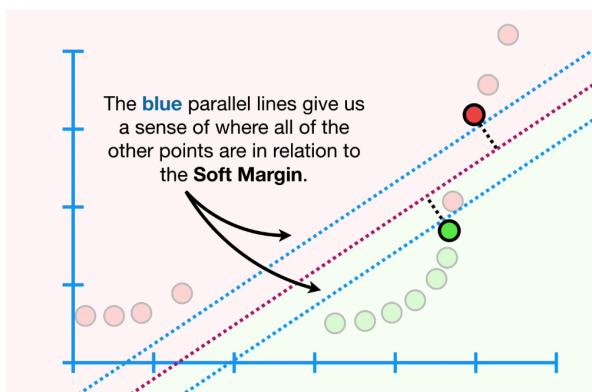
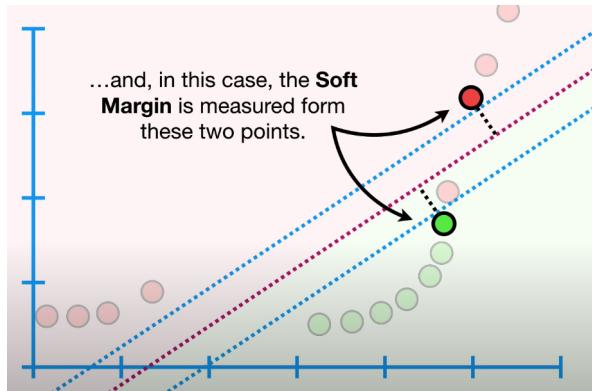
1D :

location of a threshold...

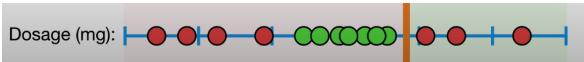
...then we are using a **Soft Margin Classifier** aka a **Support Vector Classifier** to classify observations.

The name **Support Vector Classifier** comes from the fact that the observations on the edge and *within* the **Soft Margin** are called **Support Vectors**.



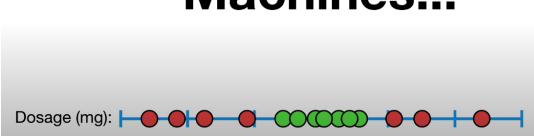


SVC.... Has serious limitation

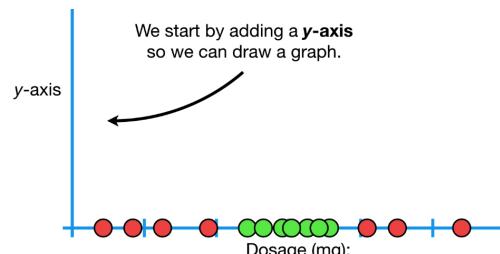


Since **Maximal Margin Classifiers** and
Support Vector Classifiers can't
handle this data, it's high time we talked
about...

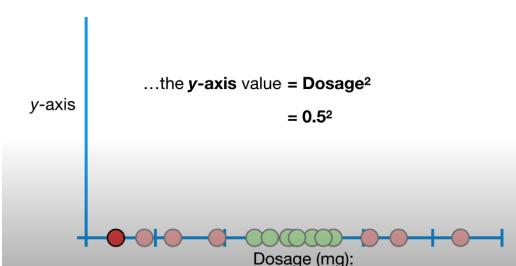
Support Vector Machines!!!



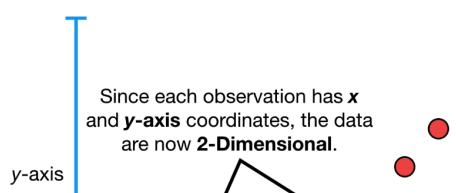
We start by adding a **y-axis**
so we can draw a graph.

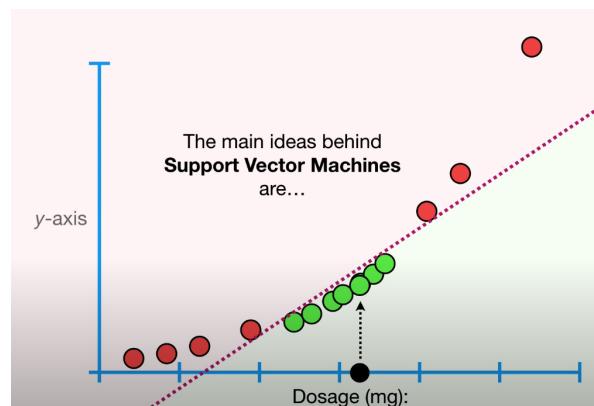
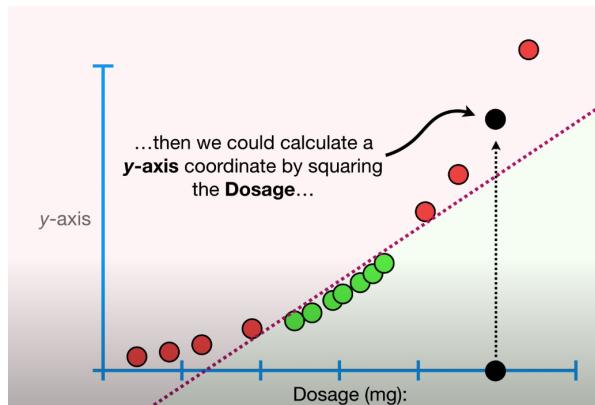
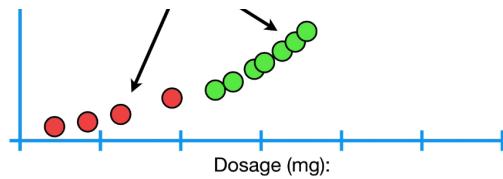


...the **y-axis value** = **Dosage²**
= **0.5²**

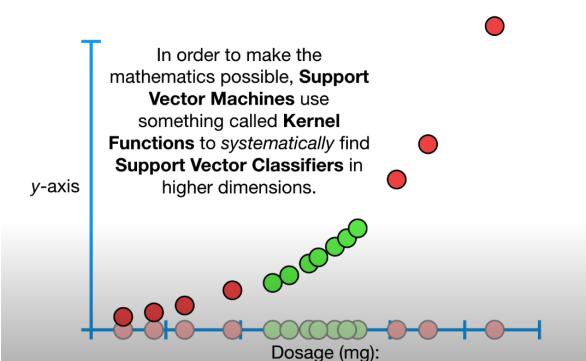
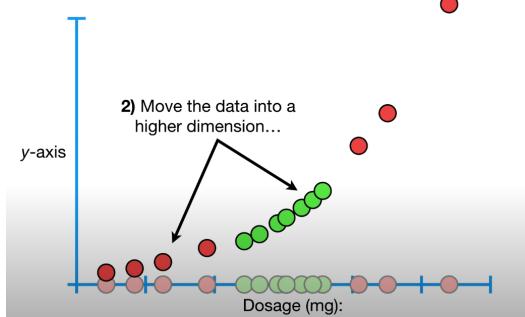
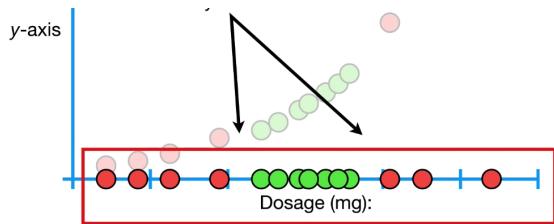


Since each observation has **x**
and **y-axis** coordinates, the data
are now **2-Dimensional**.



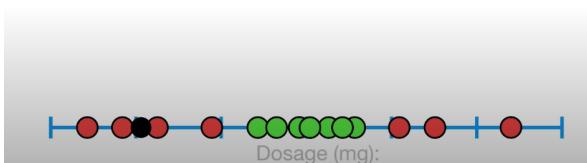


1) Start with data in a relatively low dimension...

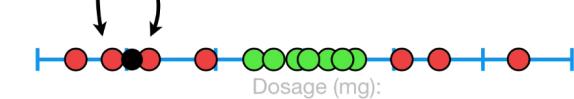


Another very commonly used **Kernel** is the **Radial Kernel**, also known as the **Radial Basis Function (RBF) Kernel**.

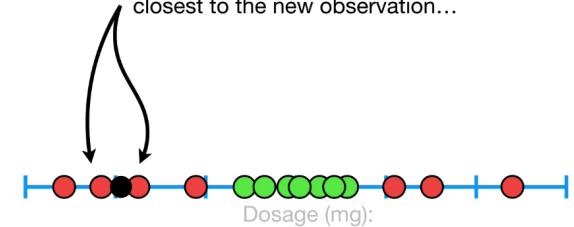
...the **Radial Kernel** behaves like a **Weighted Nearest Neighbor** model.



In other words, the closest observations (aka the nearest neighbors) have a lot of influence on how we classify the new observation...



So, since these observations are the closest to the new observation...



...the **Radial Kernel** uses their classification for the new observation.

