

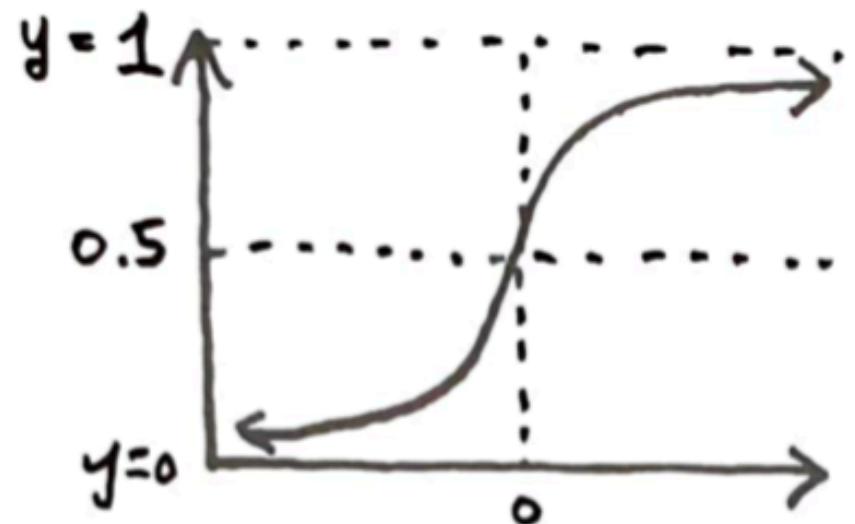
5.7 Supervised Learning Algorithms

박달님 03.19.21

5.7.1

Probabilistic Supervised Learning

Sigmoid function



$$f(x) = \frac{1}{1 + e^{-x}} \cdot e^x = \frac{e^x}{1 + e^x}$$

$$\begin{aligned} x = -\infty &: 0 \\ x = 0 &: \frac{1}{2} \\ x = \infty &: 1 \end{aligned}$$

$$\text{Loss} = \sum_i (y_i - f(x_i, w))^2$$

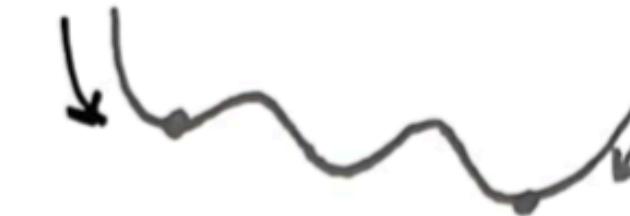
$$\text{Loss} = \sum_i (y_i - f(x_i, w))^2$$

Linear이면 제곱이 없음.
becomes very non-linear



gradient = 0로 minimize 가능.

제곱이 있으면 $f(x; -w)$ 가 weight에 대해 매우 non-linear.

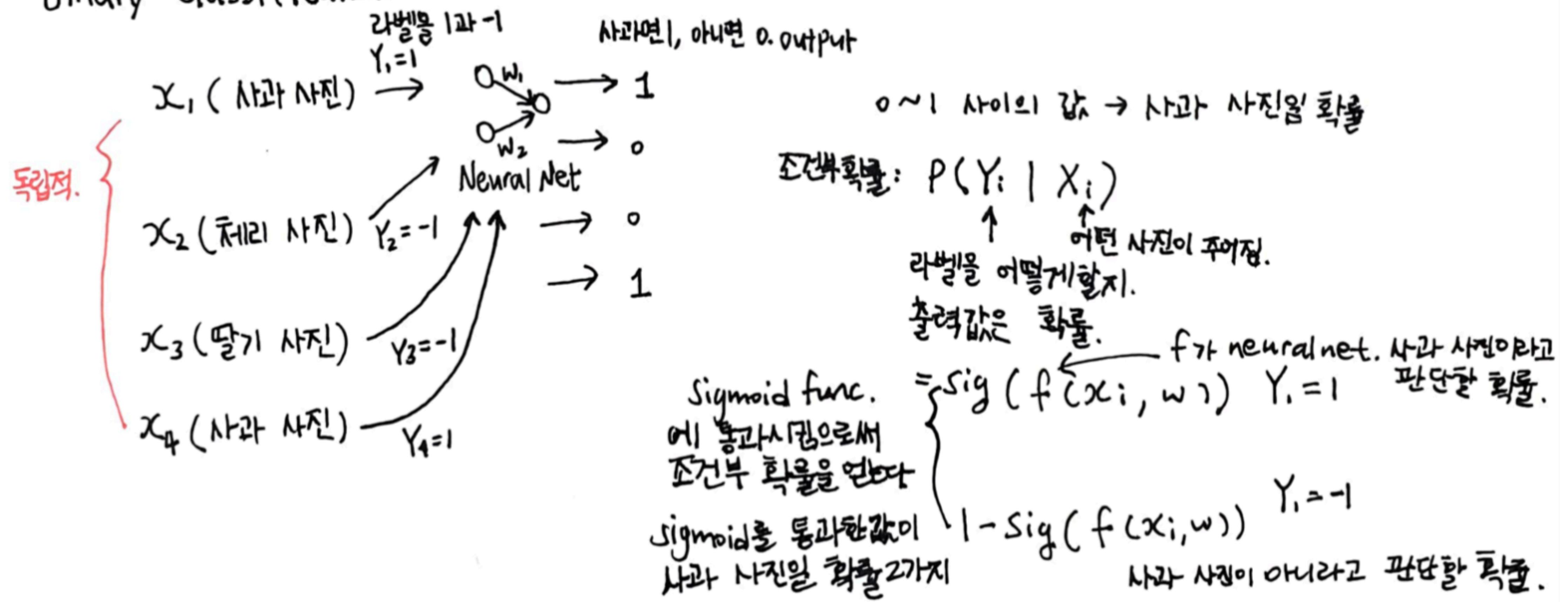


고여버린다. 이기 Gradient가 0으로 update 되어버려.

Sigmoid function의 문제점.

∴ Loss function은 $y = \frac{e^x}{1 + e^x}$ 에 맞는 likelihood를 가지고 loss function을 사용한다.

Binary Classification



사과 사진일 때는 사과가 출력될 확률을 높이고 아니면 아니라고 출력될 확률을 높여야 한다.
 확률이 꼽으로 표현되려면 독립이면 된다. \therefore 가능.

$$\text{Loss} = - \prod_{i=1}^N P(Y_i | X_i)$$

이 부분을 Maximize. 조건부 확률
Loss ↓

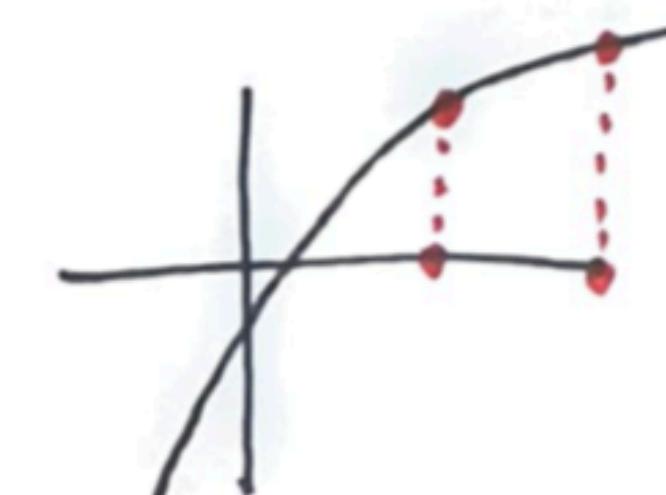
$P(Y_i | X_i)$ 부분이 알아서 1이면 $\text{sig}(f(x_i, w))$ 로 아니면 $1 - \text{sig}(f(x_i, w))$ 로 값을 계속 급해중.

급하기는 더하기로 바꿀 수 있다 → Log를 쓰우기.

Likelihood Loss = $-\sum_i \log P(Y_i | X_i)$

logistic regression이나
밀접한 cost function.

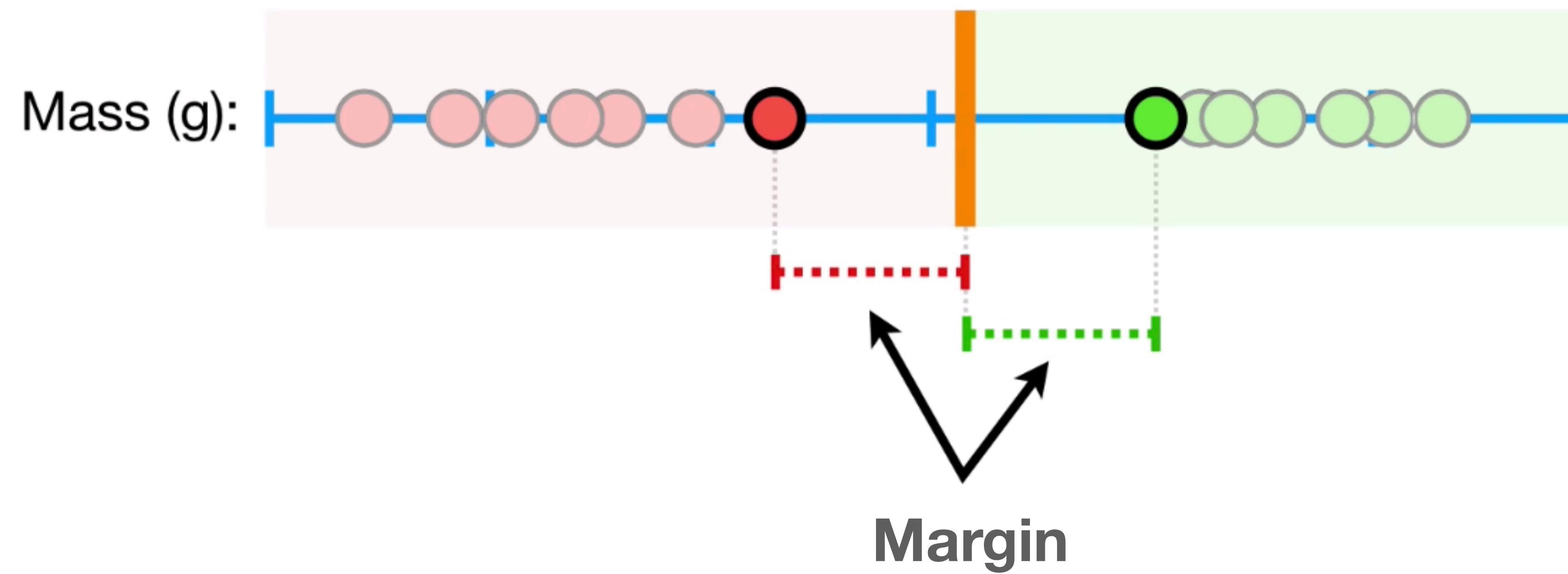
maximizing the log-likelihood.
= minimizing the negative log-likelihood (NLL)
using gradient descent.



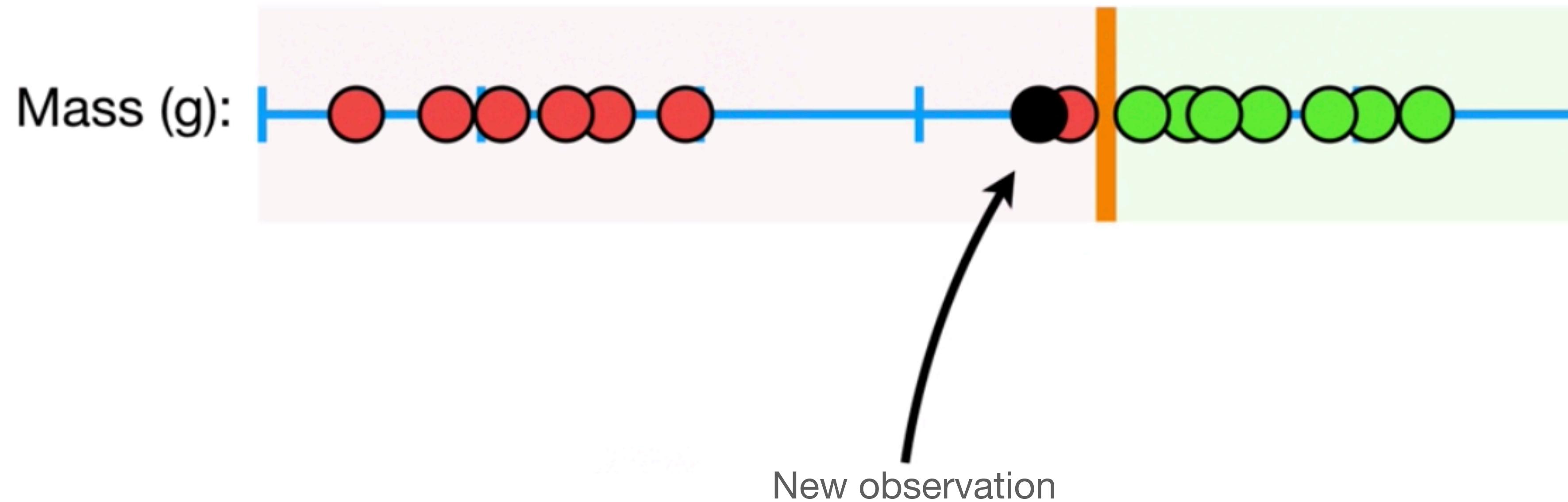
Log를 써우나마나
minimize하는 w를
찾는것은 같다.
양수에 대해서
단조증가함수이기 때문에
(Monotonic)
log를 쓰워도 된다.

5.7.2

Support Vector Machine

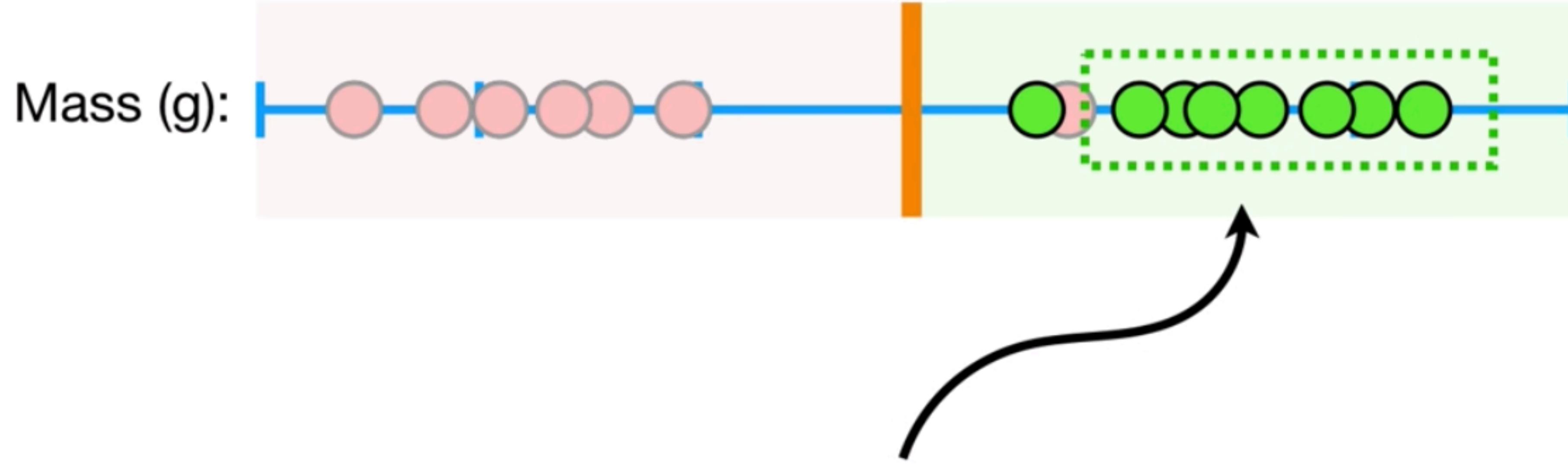


largest margin to make classifications:
Maximal Margin Classifiers

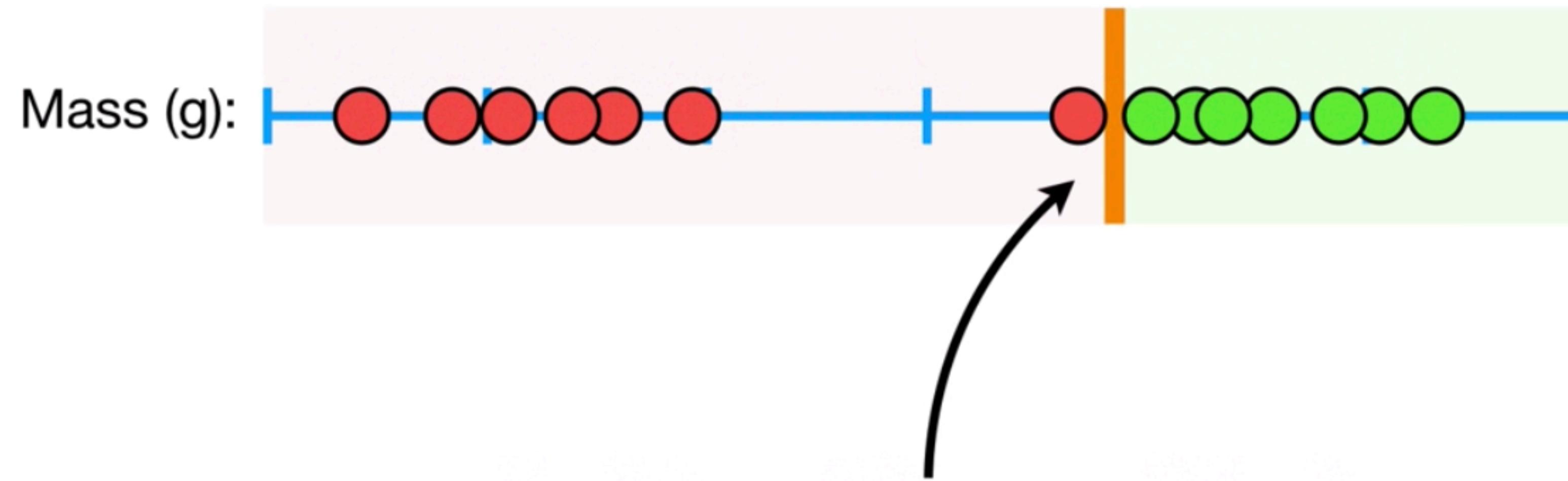


Maximal Margin Classifiers super **sensitive** to outliers in the training data

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

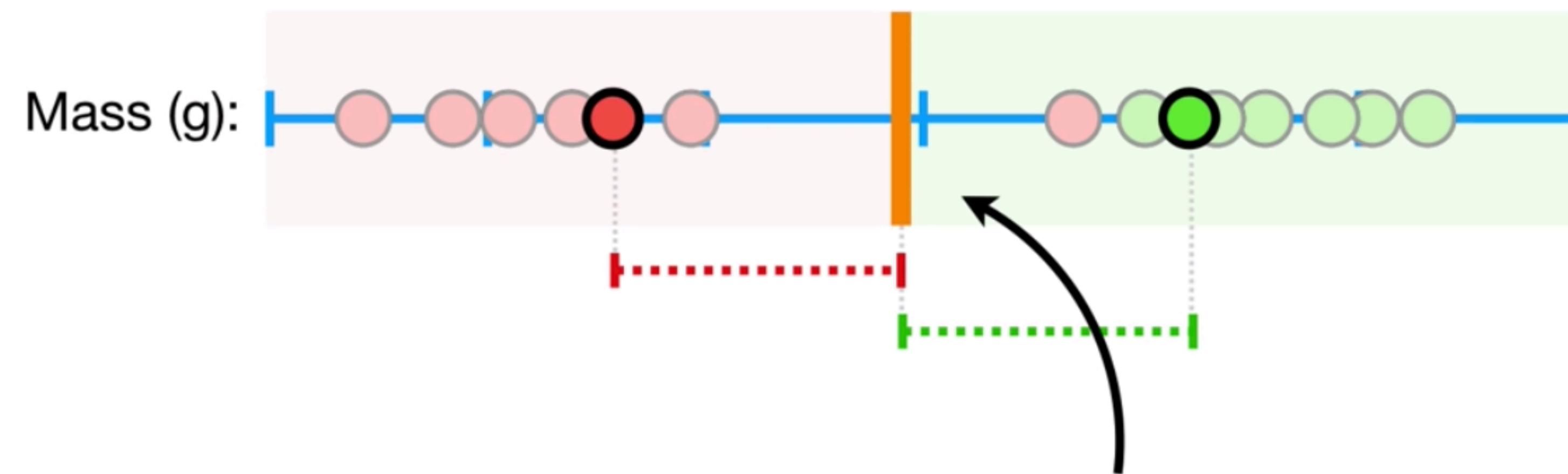


Choosing a threshold that allows misclassifications is an example of the **Bias/Variance Tradeoff** that plagues all of machine learning

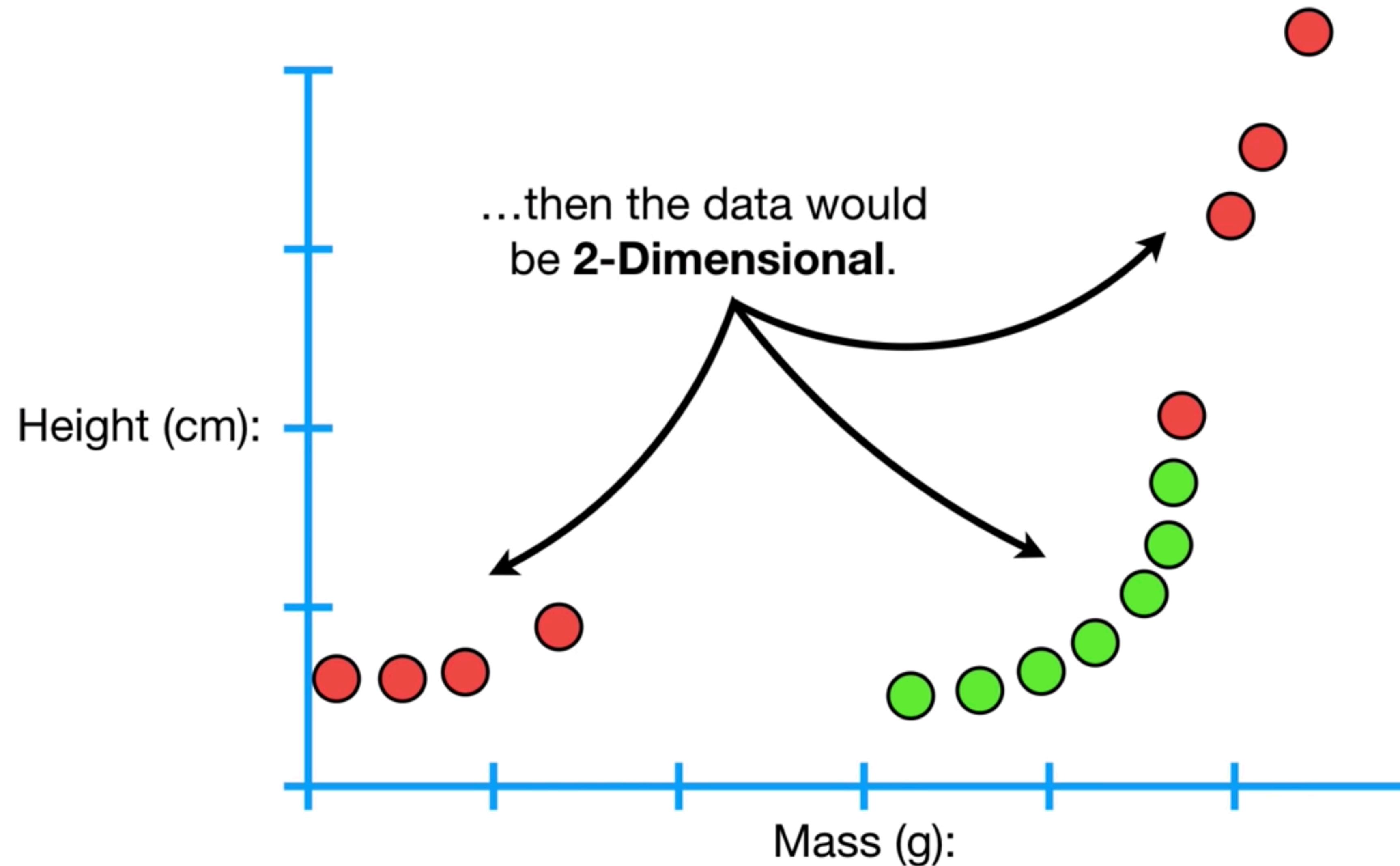


When we allow **misclassifications**, the distance between observations and the threshold is called a **Soft Margin**

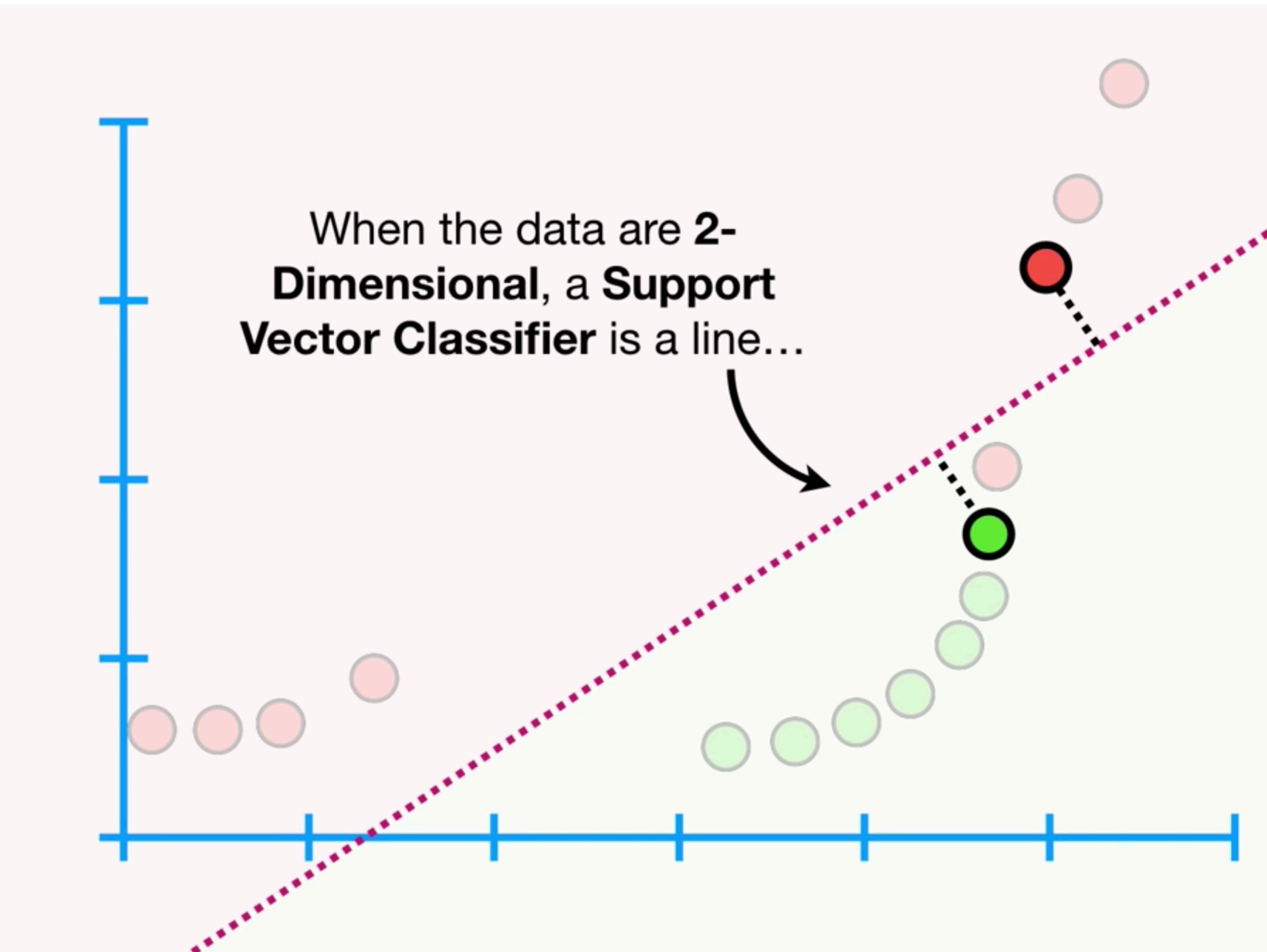
Decide which soft margin is better than the other by using **cross validation** to determine how many misclassifications and observations to allow inside of the soft margin to get the best classification



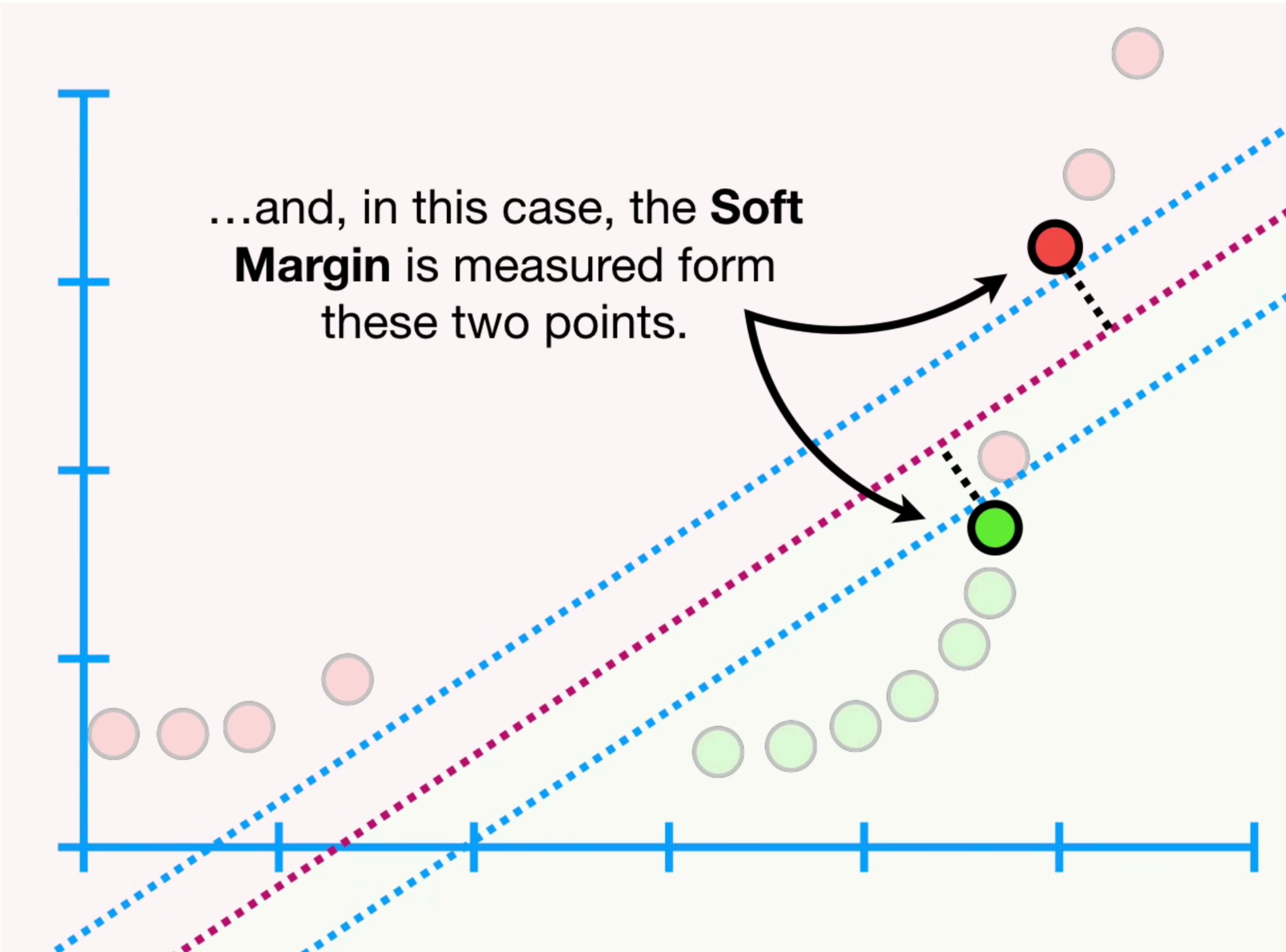
When we use a soft margin to determine the location of a threshold, then we are using a **Soft Margin classifier** aka a **Support Vector Classifier** to classify observations.



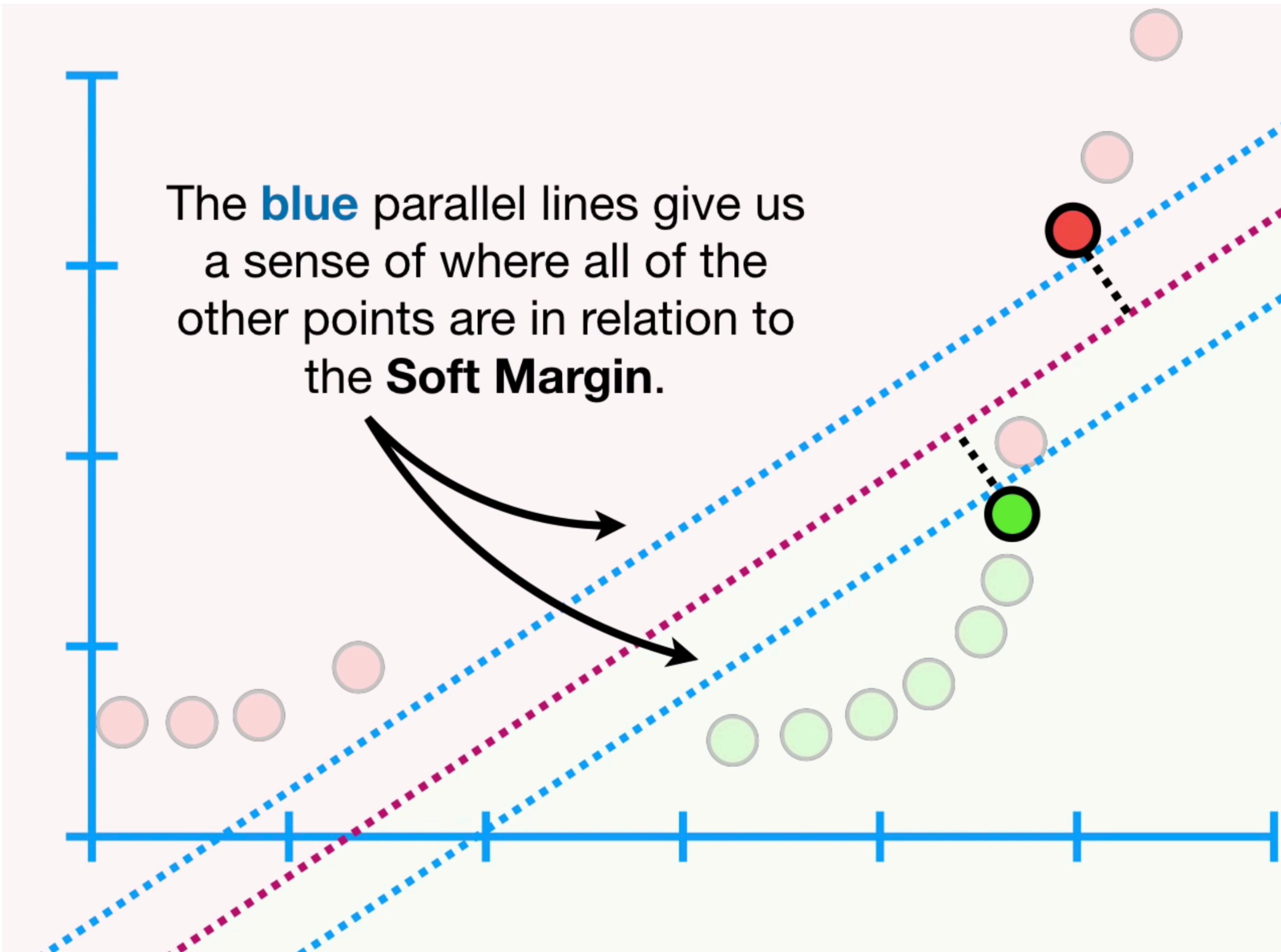
When the data are **2-Dimensional**, a **Support Vector Classifier** is a line...

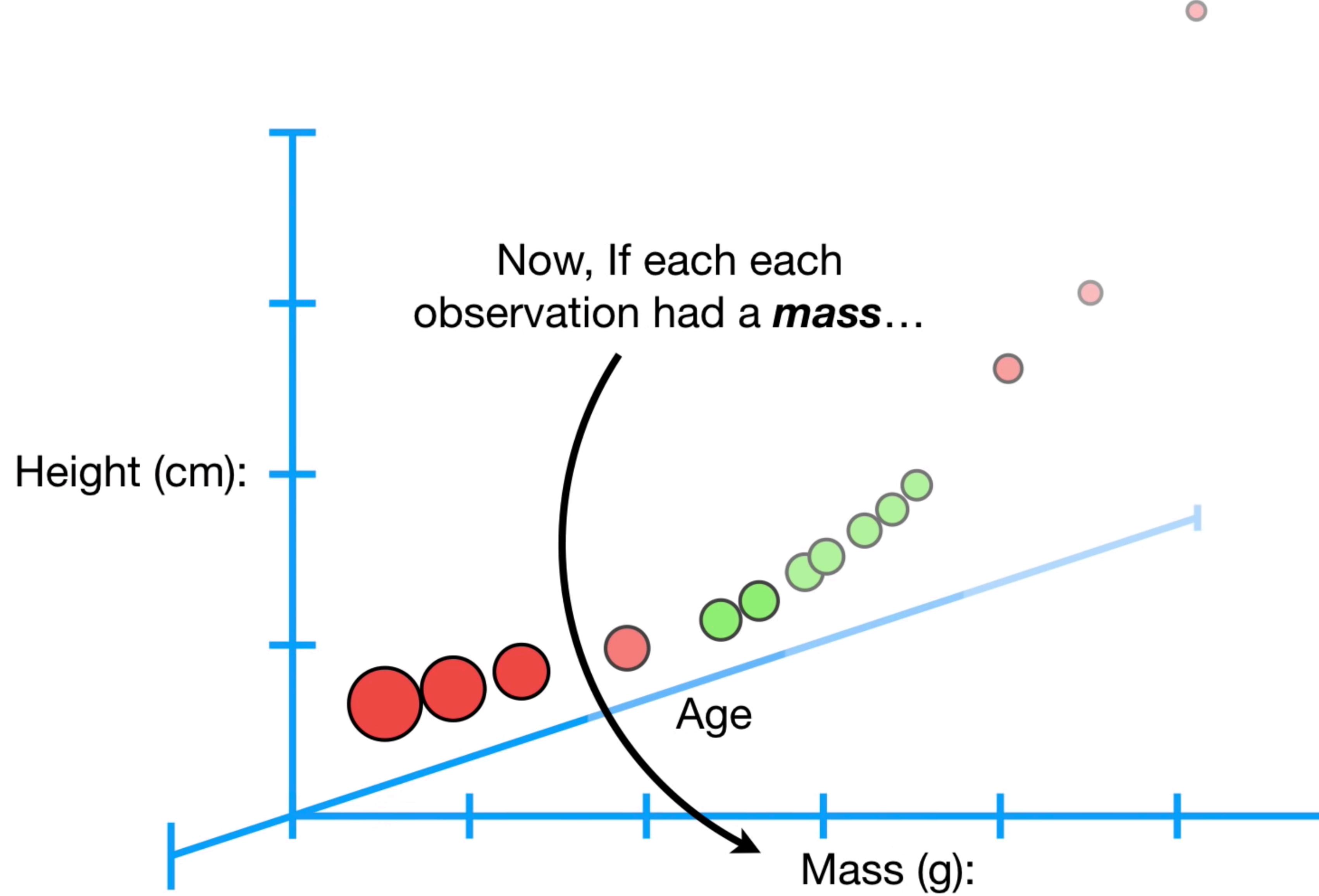


...and, in this case, the **Soft Margin** is measured from these two points.

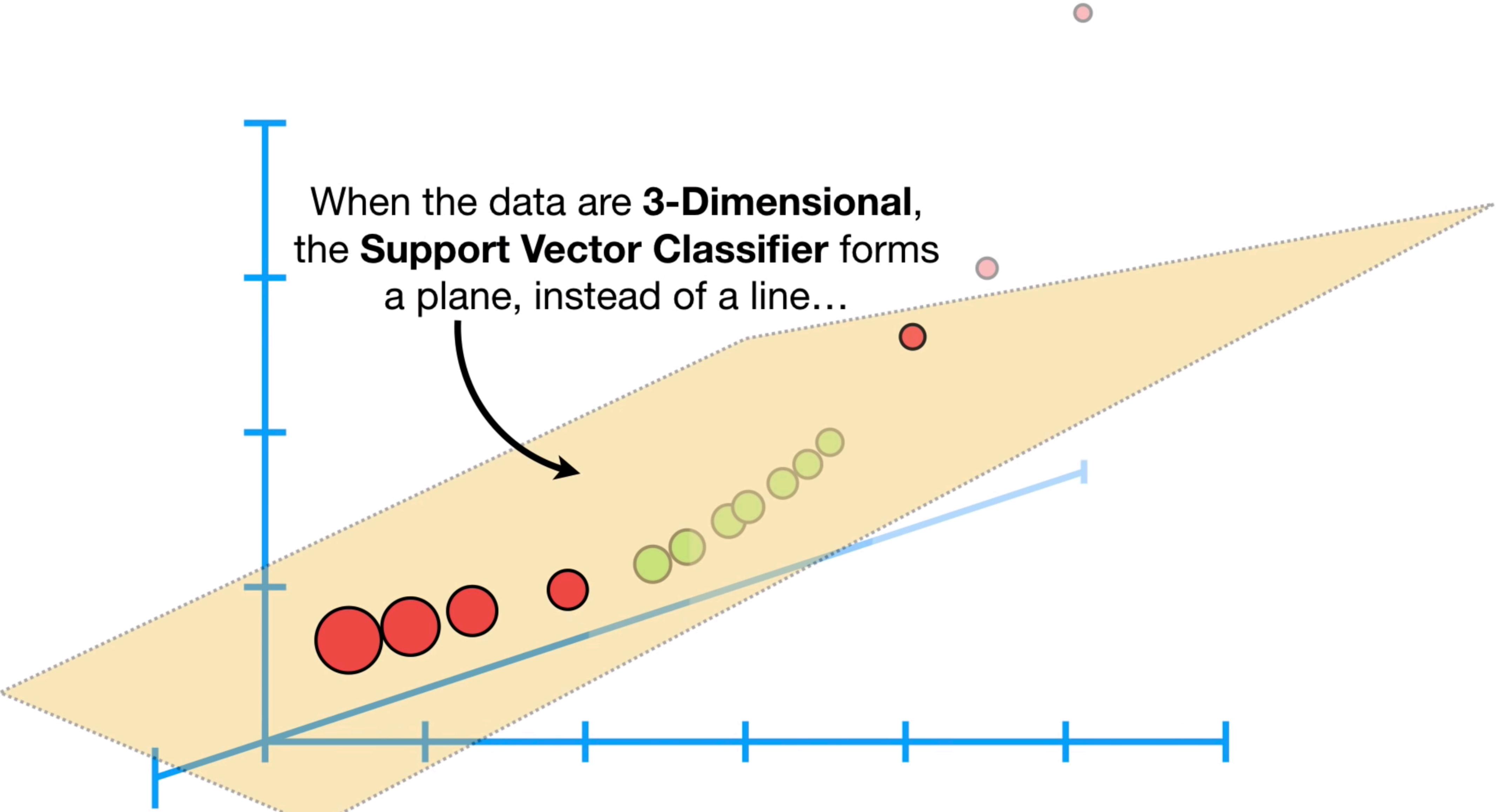


The **blue** parallel lines give us
a sense of where all of the
other points are in relation to
the **Soft Margin**.





When the data are **3-Dimensional**,
the **Support Vector Classifier** forms
a plane, instead of a line...



In this new example, with tons of overlap, we are now looking at
Drug Dosages...

