# MCSC201: Machine Learning

**Dr. Ankit Rajpal**

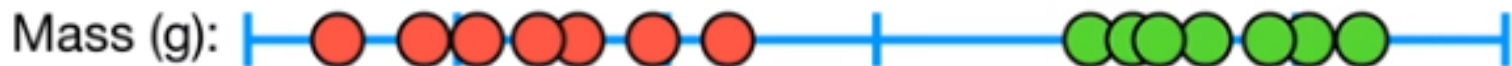**Department of Computer Science**

**University of Delhi**

# Support Vector Machine: Intuition

❧ Measurement of the **Mass of mice (g)**.
- ❧ Red dots → not obese
- ❧ Green dots → obese



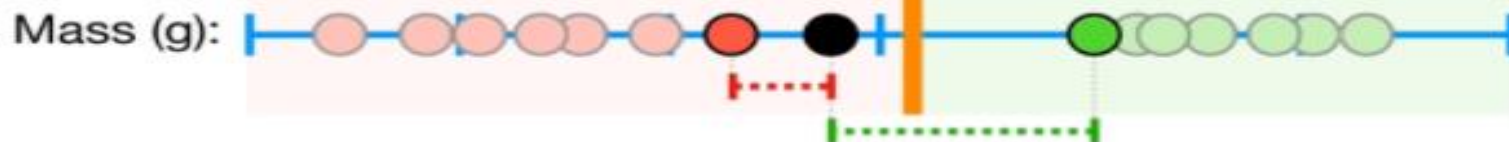In this case threshold is not a good estimator. The observation is above threshold, but closer to the non obese mice

# Support Vector Machine: Intuition

ℭℬ

ℭℬ Measurement of the **Mass of mice (g)**.

    ℭℬ Red dots → not obese

    ℭℬ Green dots → obese

# Support Vector Machine: Intuition

❧ Measurement of the **Mass of mice (g)**.

 ☙ Red dots → not obese

 ☙ Green dots → obese



The new observation is wrongly classified as not obese because the presence of an outlier

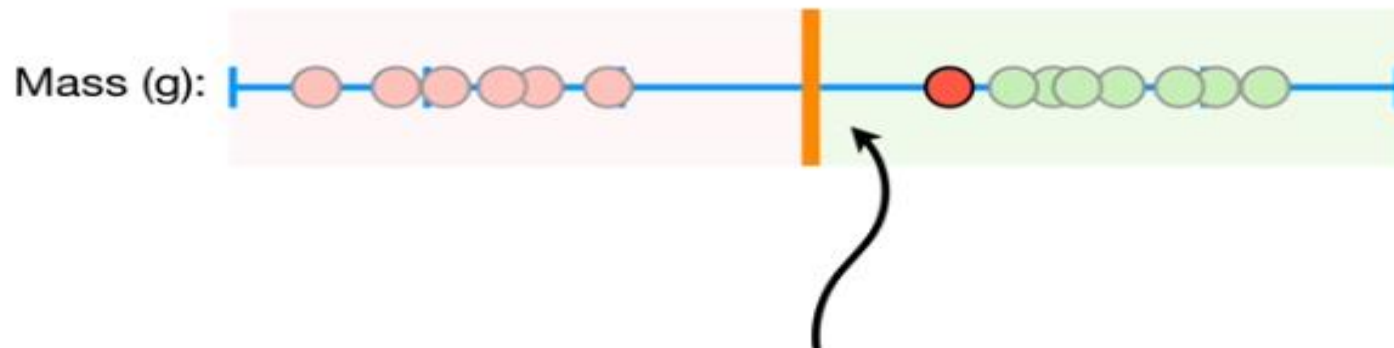Maximal Margin Classifier is super sensitive to **outliers**.

# Support Vector Machine: Intuition

℞ Measurement of the **Mass of mice (g)**.
  03 Red dots → not obese
  03 Green dots → obese

Mass (g):

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

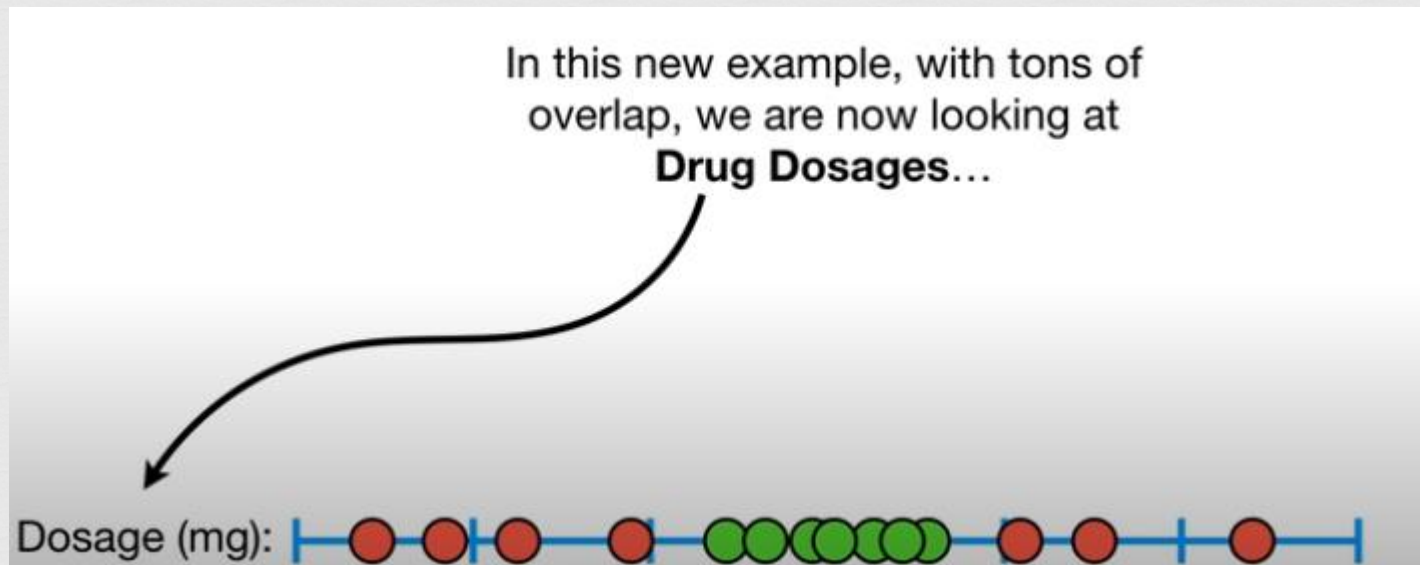Choosing a threshold that allows misclassifications.

# Support Vector Machine: Intuition

ℰ Measurement of the **Drug Dosage**.

   ℰ Red dots → patients were not cured

   ℰ Green dots → patients were cured

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

Dosage (mg):

# Support Vector Machine: Intuition

☙ Measurement of the **Drug Dosage**.

  ◦ Red dots → patients were not cured

  ◦ Green dots → patients were cured

# Support Vector Machine Intuition

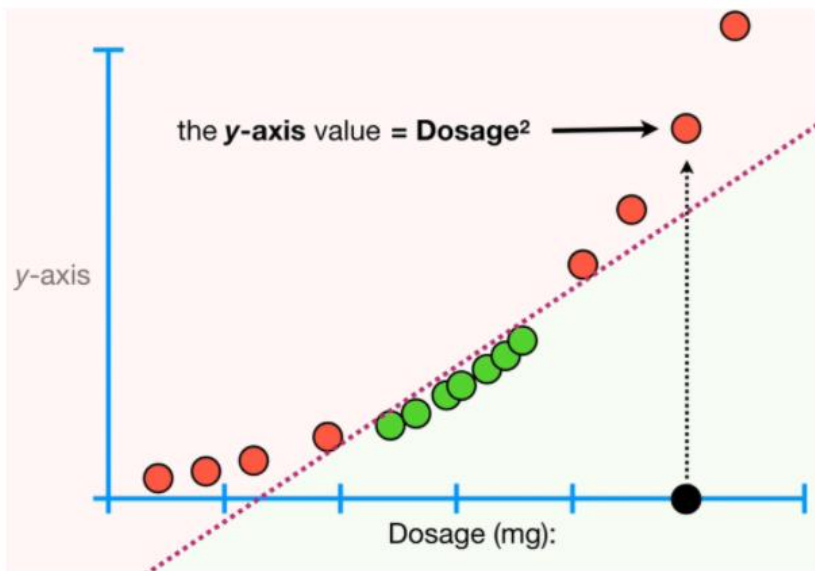ᘒ With the high overlapping depicted above, no matter where we put the classifier because will always make a lot of misclassifications.

ᘒ So, **Support Vector Classifiers** don't perform well with this type of data.

# Support Vector Machine Intuition

ﻌ Solution: We use the x-axis which represent the dosages we observed, but we also add an y-axis that will be the **square of the dosages**.



the **y-axis** value = Dosage²

y-axis

Dosage (mg):

**The main idea behind Support Vector Machines are:**
**1 – start with data in a relatively low dimension (in this example one dimension dosage in mg)**
**2 – move the data into a higher dimension (in this example from one to two dimensions)**
**3 – find a Support Vector Classifier that separates the higher dimensional data into two groups**
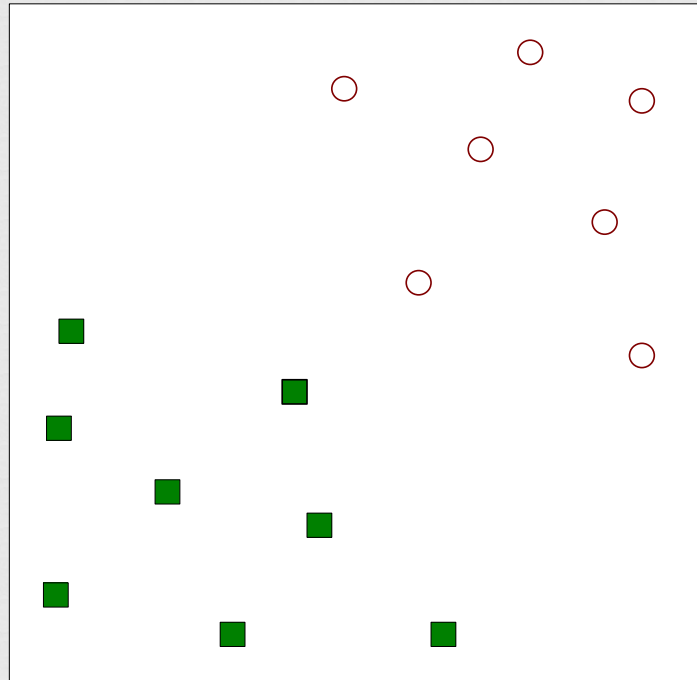
# Support Vector Machine: Intuition

 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.

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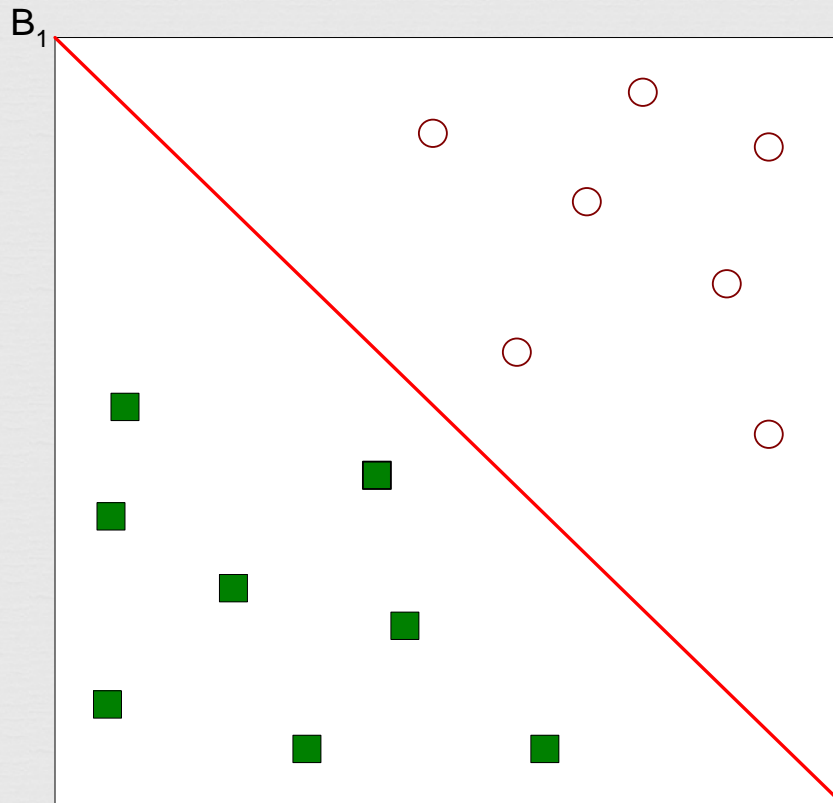
# Support Vector Machine

⍟ SVM is a method for the classification of both linear and nonlinear data.

⍟ SVM searches for the linear optimal separating hyperplane (i.e., a "decision boundary" separating the tuples of one class from another).

⍟ Extend to patterns that are not linearly separable by transformations of original data to map into new space – the Kernel function.

⍟ Support vectors are the data points that lie closest to the decision surface (or hyperplane).

  ⍟ They are the data points most difficult to classify

  ⍟ They have direct bearing on the optimum location of the decision surface

# Support Vector Machines
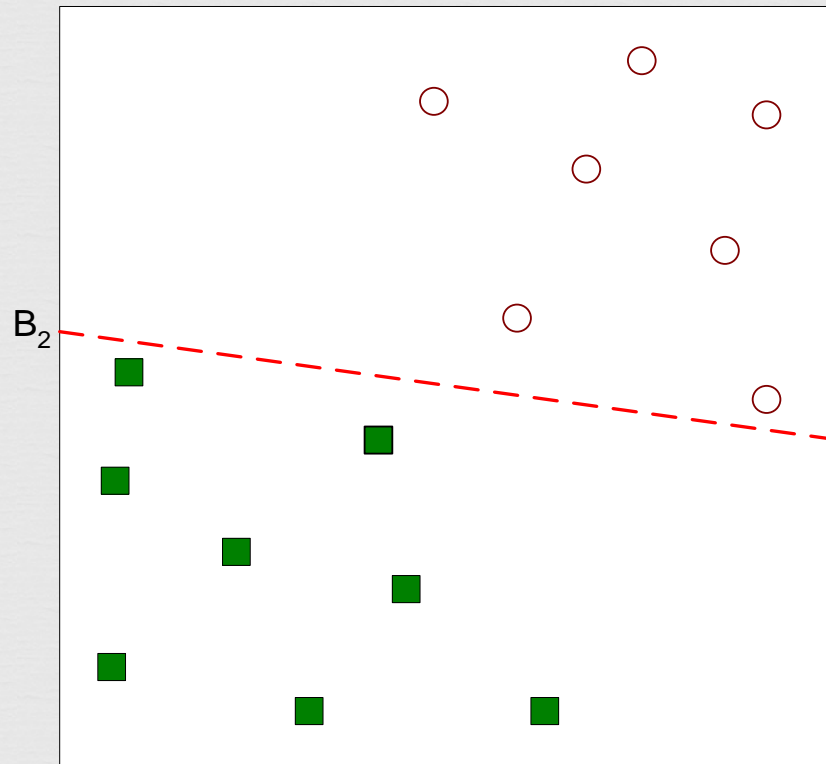


Find a linear hyperplane (decision boundary) that will separate the data

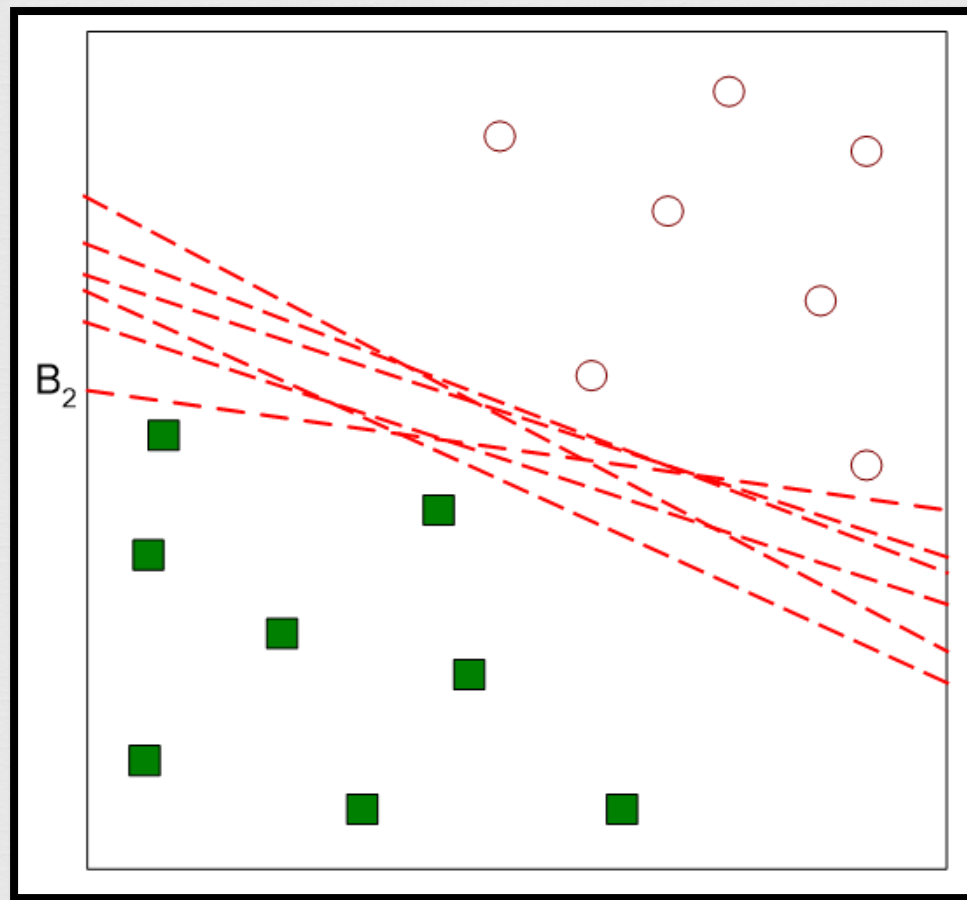# Support Vector Machines



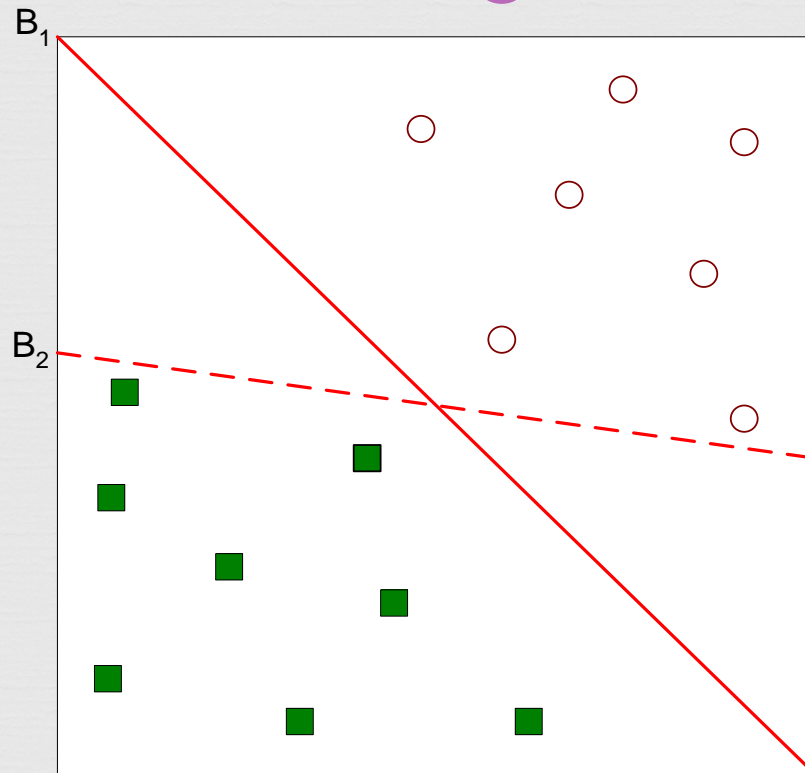One Possible Solution

# Support Vector Machines



B₂

Another possible solution

# Support Vector Machines



*Other possible solutions*

# Support Vector Machines



 Which one is better? B1 or B2?
 How do you define better?

# Support Vector Machines



Find hyperplane maximizes the margin => B1 is better than B2

# Support Vector Machines

$$\vec{w} \bullet \vec{x} + b = 0$$

$$\vec{w} \bullet \vec{x} + b = -1$$

$B_1$

$b_{11}$

$b_{12}$

$$\vec{w} \bullet \vec{x} + b = +1$$

$$f(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \bullet \vec{x} + b \geq 1 \\ -1 & \text{if } \vec{w} \bullet \vec{x} + b \leq -1 \end{cases}$$

# Rationale for Maximum Margin

ॐ Decision boundaries with large margins tend to have better generalization capability than those with smaller margins.

ॐ Intuitively, if the margin is small, then any slight perturbations to the decision boundary can have quite a significant impact on its classification.

ॐ Classifiers that produce decision boundaries with small margins are more susceptible to model overfitting and tend to generalize poorly on previously unseen examples.

# Maximum Margin Hyperplane

# Linear SVM: Separable Case

ରଷ Binary classification problem

ରଷ N training Examples

ରଷ Each example is denoted by a tuple $(x_i, y_i)(i = 1, 2, \ldots, N)$

ରଷ $x_i$ is d-dimensional and $y_i \in \{-1, 1\}$ such that $y_i$ is +1 for positive example and $y_i$ is -1 for negative example.

ରଷ Decision Boundary: $\vec{w}.\vec{x} + b = 0$, where $\vec{w}$ and $b$ are parameters of the model.

ରଷ All training instances from class $y = 1$ (i.e., the squares) must be located on or above the hyperplane $w.x + b = 1$, while those instances from class $y = -1$ (i.e., the circles) must be located on or below the hyperplane $w.x + b = -1$.

# Linear SVM: Separable Case

$\wp$ If $x_+$ and $x_-$ are any two points located above the positive and negative marginal boundaries, respectively, then

$$\vec{w}.x_+ + b \geq +1 \text{ and } \vec{w}.x_- + b \leq -1$$

$\wp$ Since, $y_i$ is +1 for positive example and $y_i$ is -1 for negative example. The compact form of the above two constraints is a follows:

$$y_i(\vec{w}.x + b) \geq +1$$

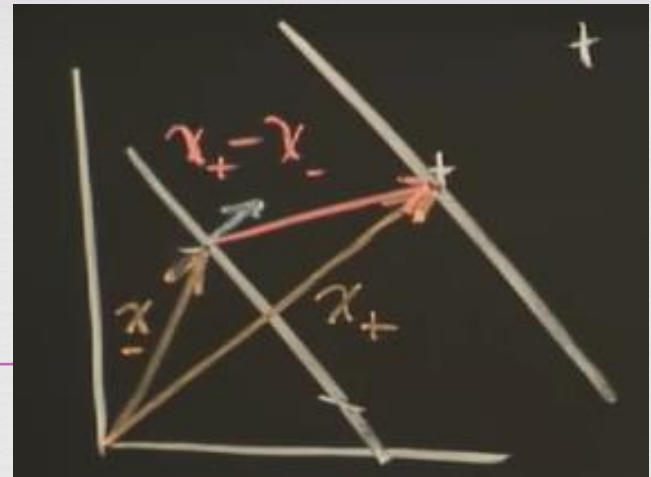# Learning a linear SVM Model

∞ Training Phase: Estimation of $w$ and $b$ of the decision boundary from training data.

∞ The Parameters are chosen such the following condition is met.

$$y_i(w.x_i + b) \geq 1, i = 1, 2, \ldots, N.$$

# Margin of a Linear Classifier



- Consider a square located on the positive marginal hyperplane => $b_{i1}: \vec{w}.x_+ + b = 1$

- Consider a circle located on the negative marginal hyperplane => $b_{i2}: \vec{w}.x_- + b = -1$

- The margin of decision boundary is the distance between these two hyperplanes.

- Since $x_+$ is a point located on $b_{i1}$ and $x_-$ is a point located on $b_{i2}$.

- $x_+ - x_-$ is a vector directed from $x_-$ to $x_+$.

- Direction of $\vec{w}$ is perpendicular to the decision boundary, therefore, normalizing $\vec{w}$ will yield a unit vector in the perpendicular direction.

- Therefore, the margin (distance between $b_{i1}$ and $b_{i2}$) can be computed as $\frac{\vec{w}}{\|w\|}.(x_+ - x_-)$   (Since $\vec{b}.\vec{a} = \|b\|\|a\|cos\theta$ => $\frac{\vec{b}}{\|b\|}.\vec{a} = \|a\|cos\theta$ )

$$= \frac{\vec{w}.x_+}{\|w\|} - \frac{\vec{w}.x_-}{\|w\|} = \frac{1-b-(-1-b)}{\|w\|} = \frac{2}{\|w\|}$$

# Support Vector Machines

ℰ We want to maximize: $Margin = \frac{2}{\|w\|}$

ℰ Which is equivalent to minimizing: $L(w) = \frac{\|w\|}{2}$

ℰ But subjected to the following constraints:

$$y_i(w.x_i + b) \geq 1, i = 1, 2, \ldots, N$$

ℰ This is a constrained optimization problem

ℰ Numerical approaches to solve it (e.g., quadratic programming)

# Hyperparameters of linear SVM
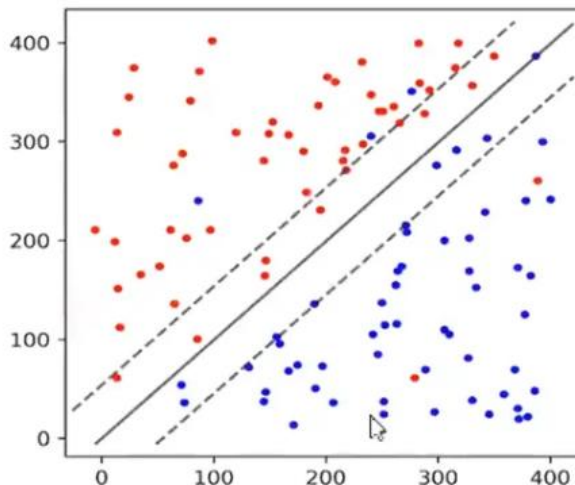
Objective (Cost) Function:

$$\min_{\vec{w},b} imize \; \frac{\|w\|}{2} + C \sum_{i=1}^{n} \xi_i$$

C is the hyperparameter that controls the number of misclassification points.

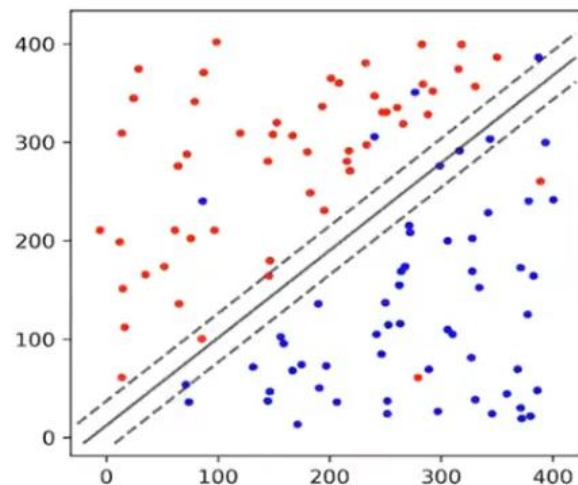Eta (greek symbol) is the summation of misclassified points from marginal plane.

# Hyperparameters of linear SVM



SVM Parameter C

C = 1

C = 100

**1.Smaller C:** When C is small => SVM places a higher priority on achieving a wide margin, even if that means allowing more misclassifications. In this case, the SVM is more tolerant of misclassified points and focuses on finding a larger margin. **2.Larger C:** When C is large, the SVM becomes more sensitive to misclassifications and tries to minimize them as much as possible. This can lead to a narrower margin to correctly classify more points.