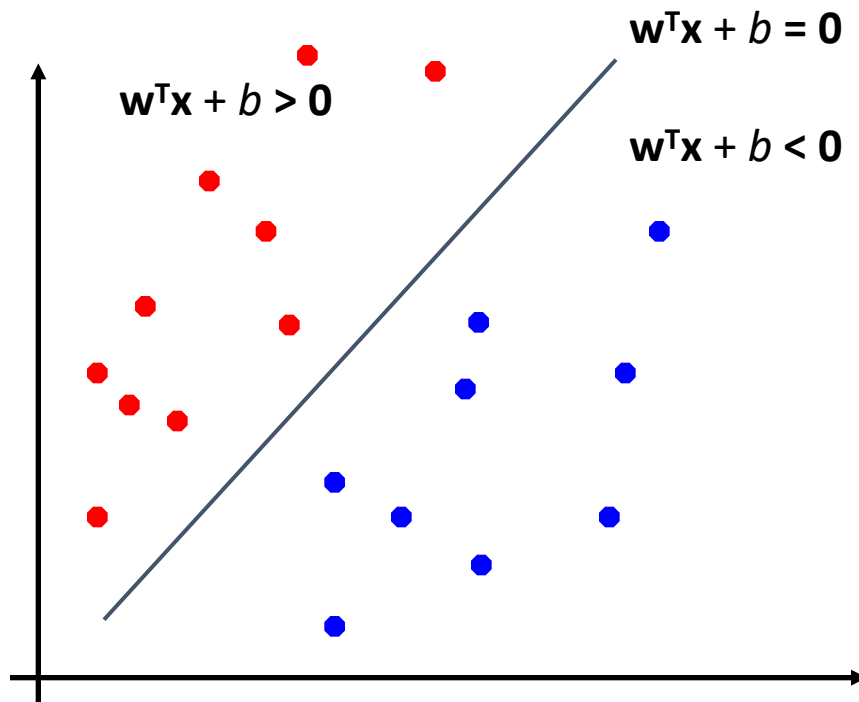


# Support Vector Machine

# Linear Separators

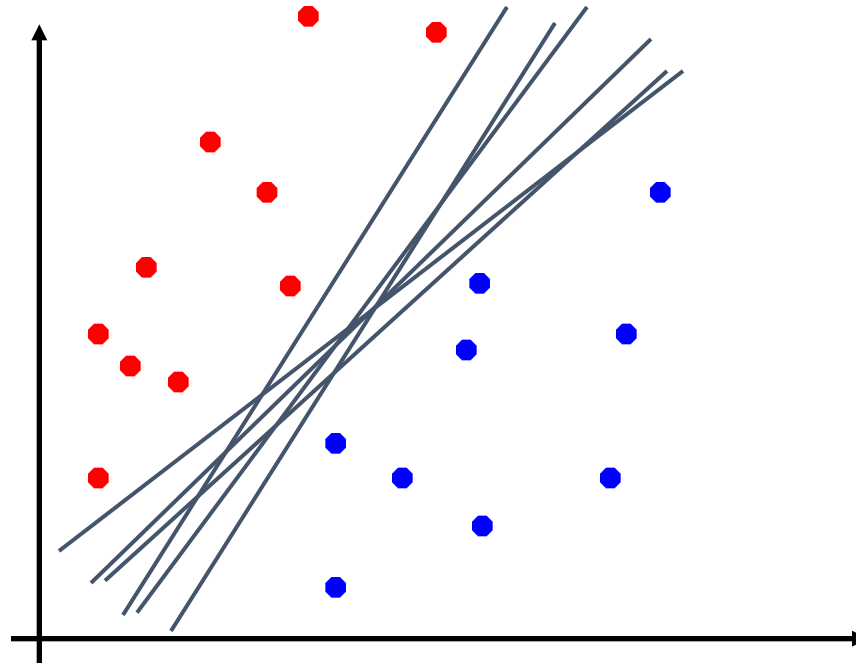
- Binary classification can be viewed as the task of separating classes in feature space:



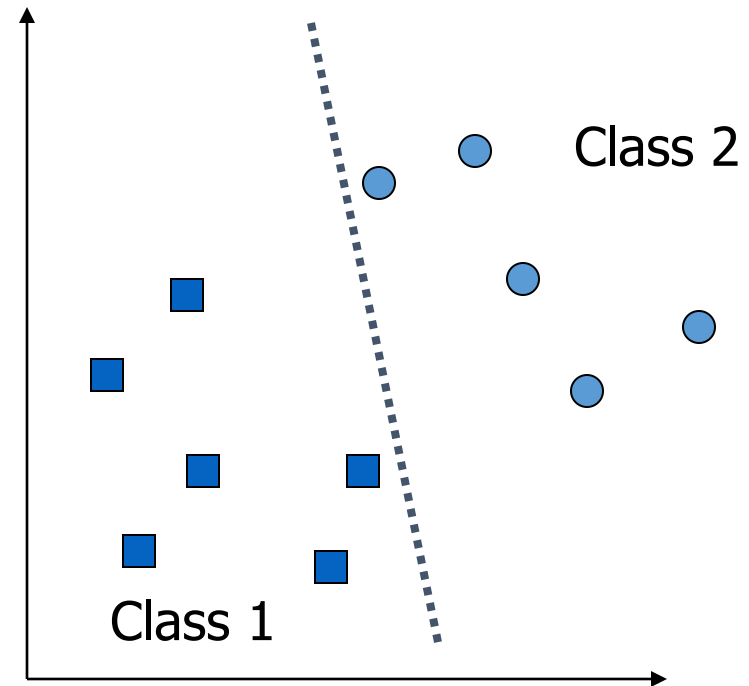
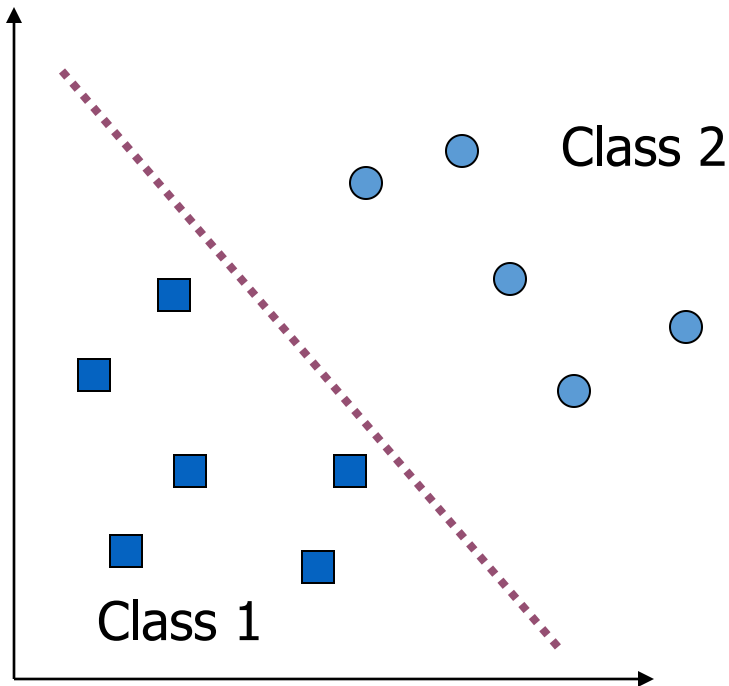
$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

# Linear Separators

- Which of the linear separators is optimal?

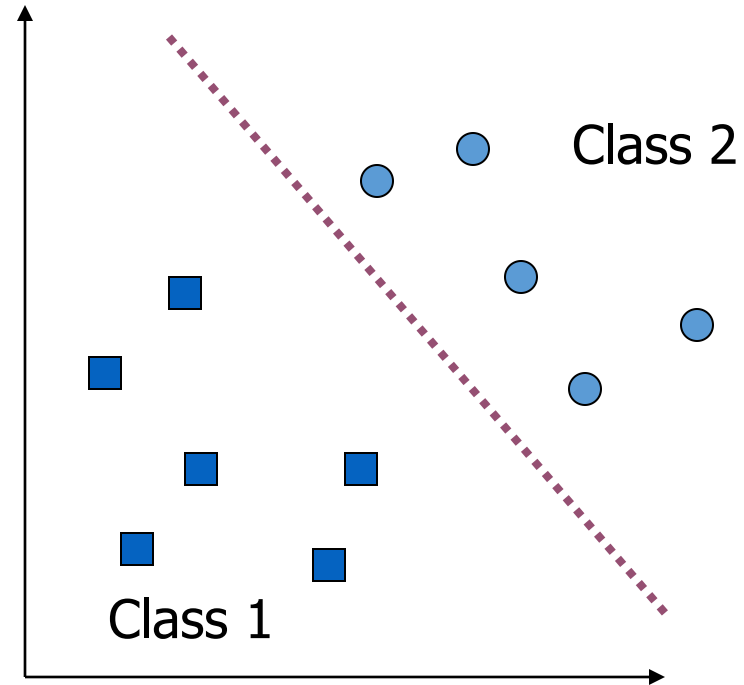


# Examples of Bad Decision Boundaries



# What is a good Decision Boundary?

- Many decision boundaries!
  - The Perceptron algorithm can be used to find such a boundary
- Are all decision boundaries equally good?



# Finding the Decision Boundary

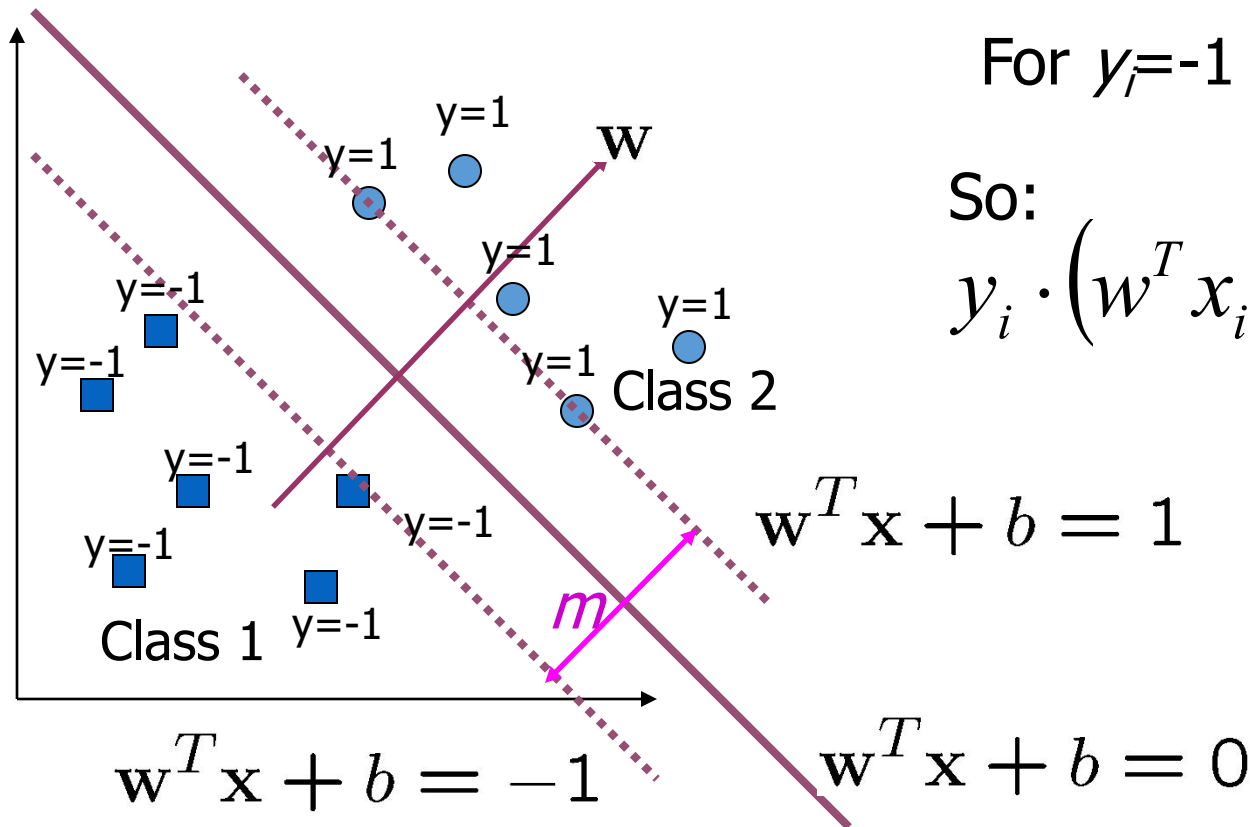
- Let  $\{x_1, \dots, x_n\}$  be our data set and let  $y_i \in \{1, -1\}$  be the class label of  $x_i$

$$\text{For } y_i = 1 \quad w^T x_i + b \geq 1$$

$$\text{For } y_i = -1 \quad w^T x_i + b \leq -1$$

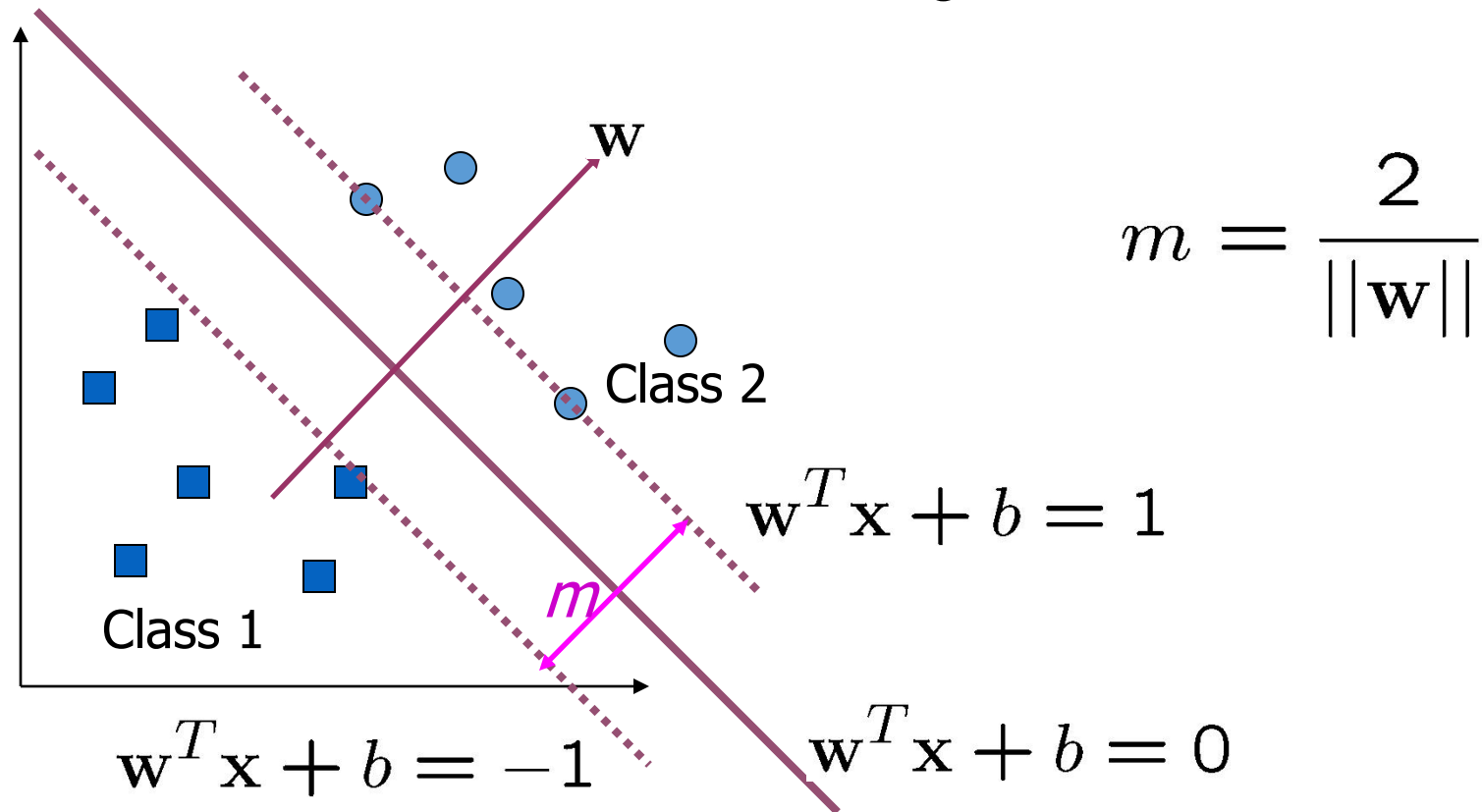
So:

$$y_i \cdot (w^T x_i + b) \geq 1, \forall (x_i, y_i)$$



# Large-margin Decision Boundary

- The decision boundary should be as far away from the data of both classes as possible
  - We should maximize the margin,  $m$



# Finding the Decision Boundary

- The decision boundary should classify all points correctly  $\vdash$

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad \forall i$$

- The decision boundary can be found by solving the following constrained optimization problem

$$\text{Minimize } \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{subject to } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$

- This is a constrained optimization problem. Solving it requires to use Lagrange multipliers



# What is a Support Vector Machine?

- It is a supervised machine learning problem where we try to find a hyperplane that best separates the two classes.
- SVM and logistic regression both the algorithms try to find the best hyperplane (decision boundary), but the main difference is logistic regression is a probabilistic approach whereas support vector machine is based on statistical approaches.
- SVM finds maximum margin between the hyperplanes that means maximum distances between the two classes.

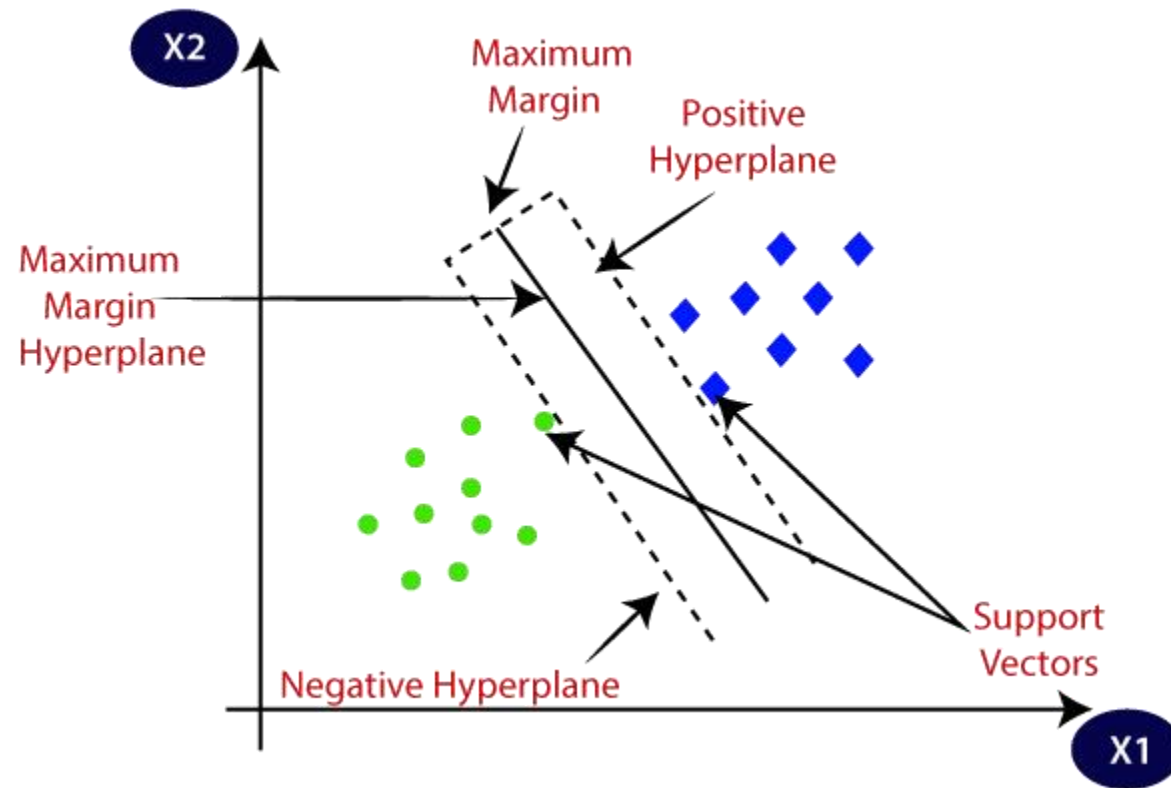
# Types of Support Vector Machine

- **Linear SVM**
- When the data is perfectly linearly separable only then we can use Linear SVM. Perfectly linearly separable means that the data points can be classified into 2 classes by using a single straight line(if 2D).
- **Non-Linear SVM**
- When the data is not linearly separable then we can use Non-Linear SVM, which means when the data points cannot be separated into 2 classes by using a straight line (if 2D) then we use some advanced techniques like kernel tricks to classify them. In most real-world applications we do not find linearly separable datapoints hence we use kernel trick to solve them.

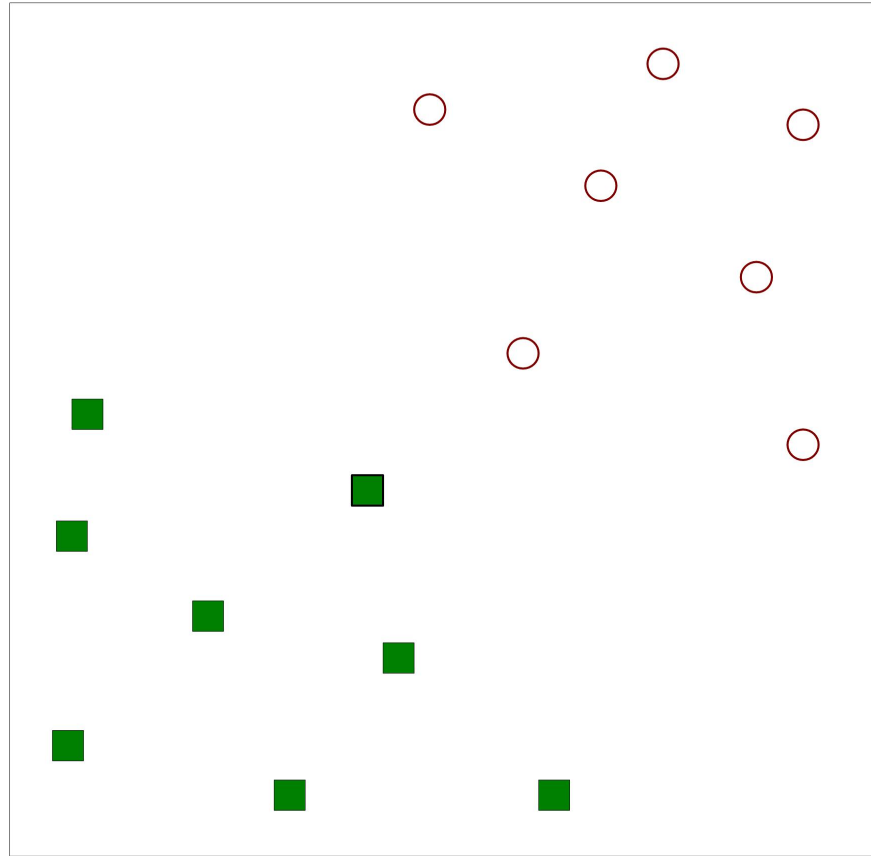
# Few mostly used terms

- **Support Vectors:** These are the points that are closest to the hyperplane. A separating line will be defined with the help of these data points.
- **Margin:** it is the distance between the hyperplane and the observations closest to the hyperplane (support vectors). In SVM large margin is considered a good margin. There are two types of margins **hard margin** and **soft margin**.

# SVM in figure

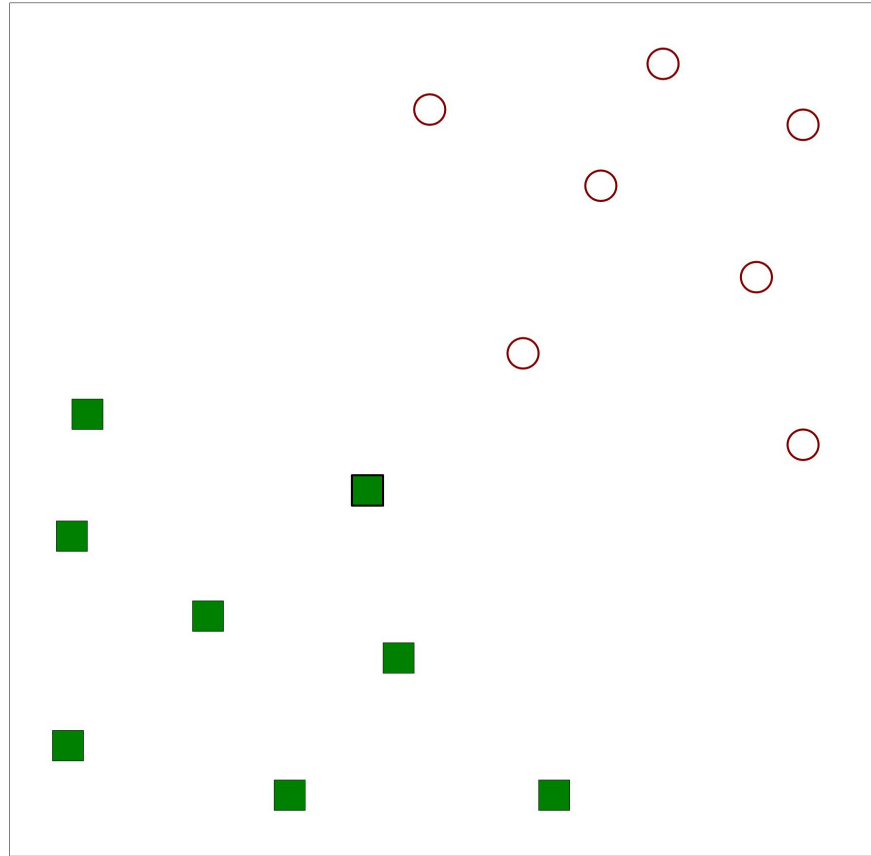


# Support Vector Machines



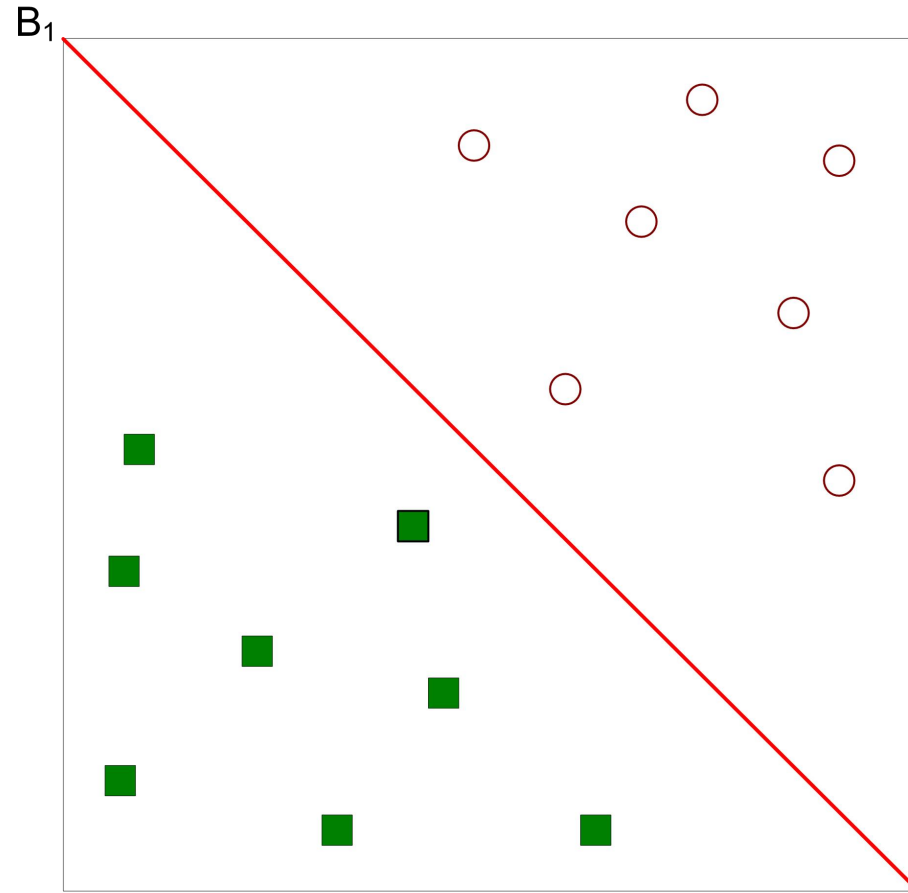
- Find a linear hyperplane (decision boundary) that will separate the data

# Support Vector Machines



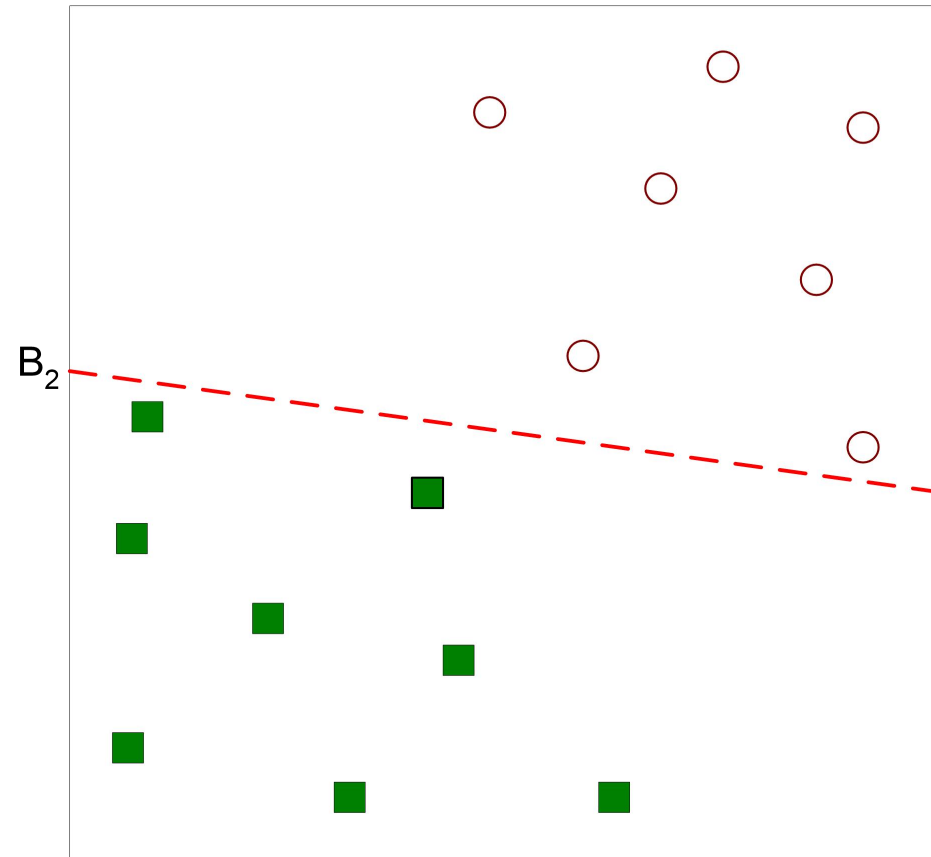
- Find a linear hyperplane (decision boundary) that will separate the data

# Support Vector Machines



- One Possible Solution

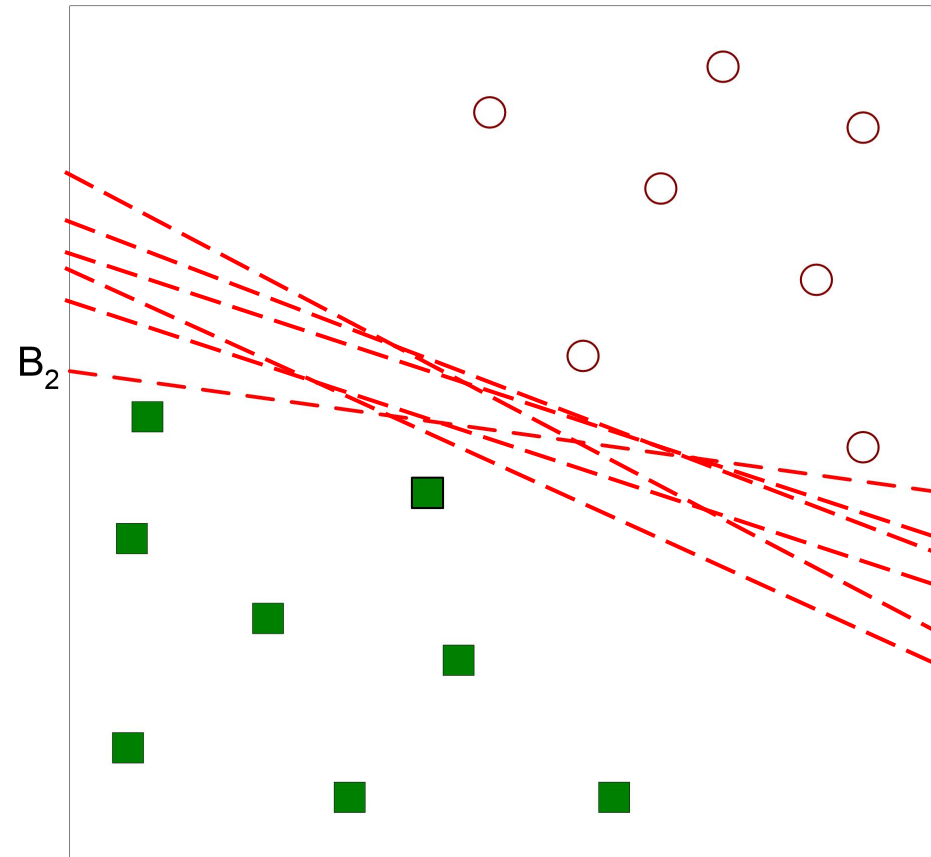
# Support Vector Machines



- Another possible solution

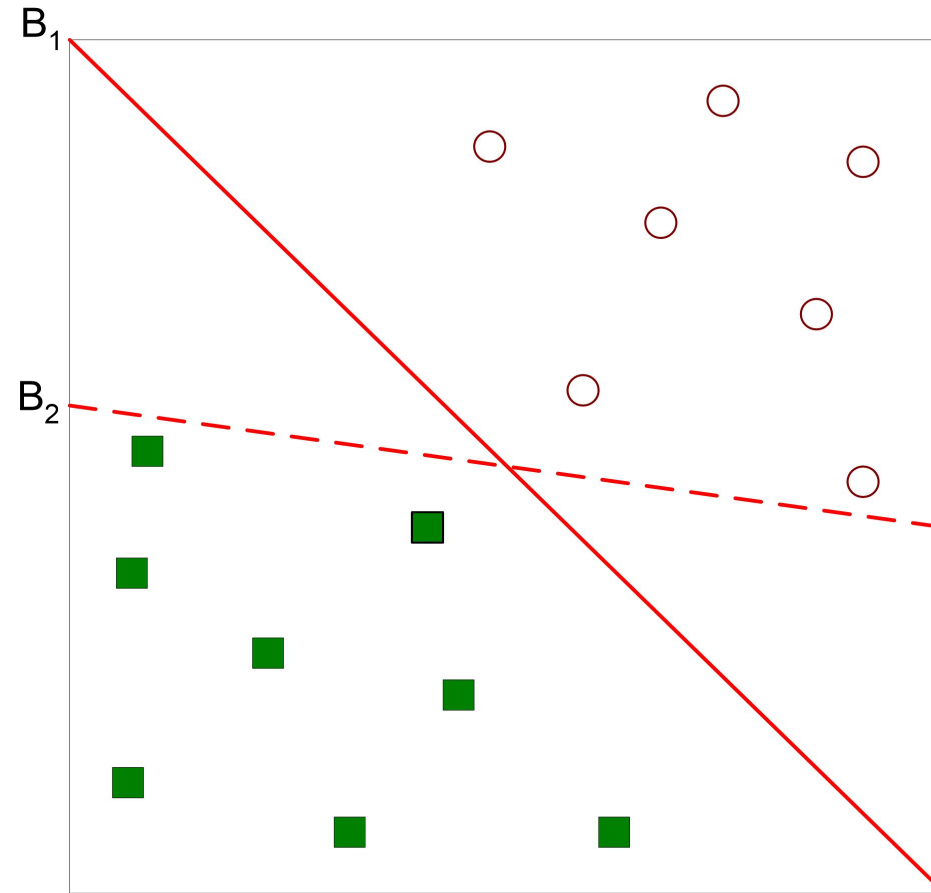


# Support Vector Machines



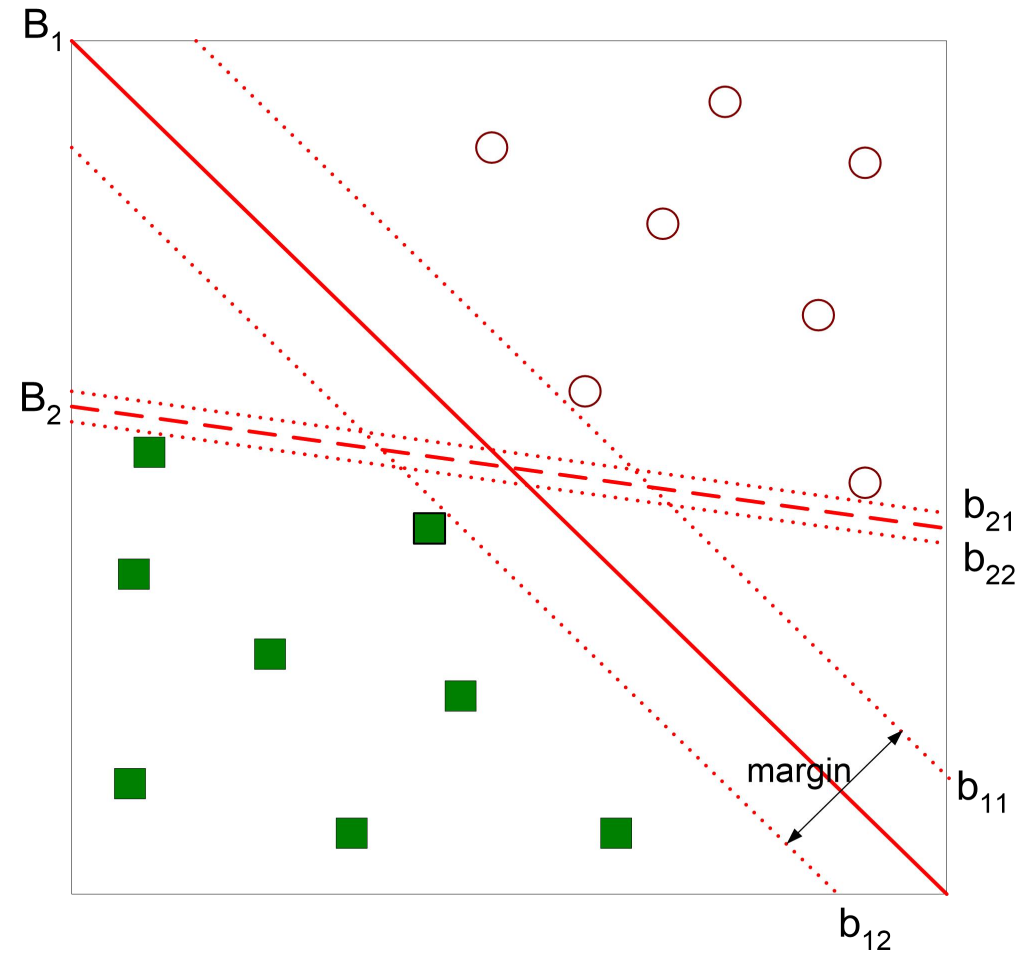
- Other possible solutions

# Support Vector Machines



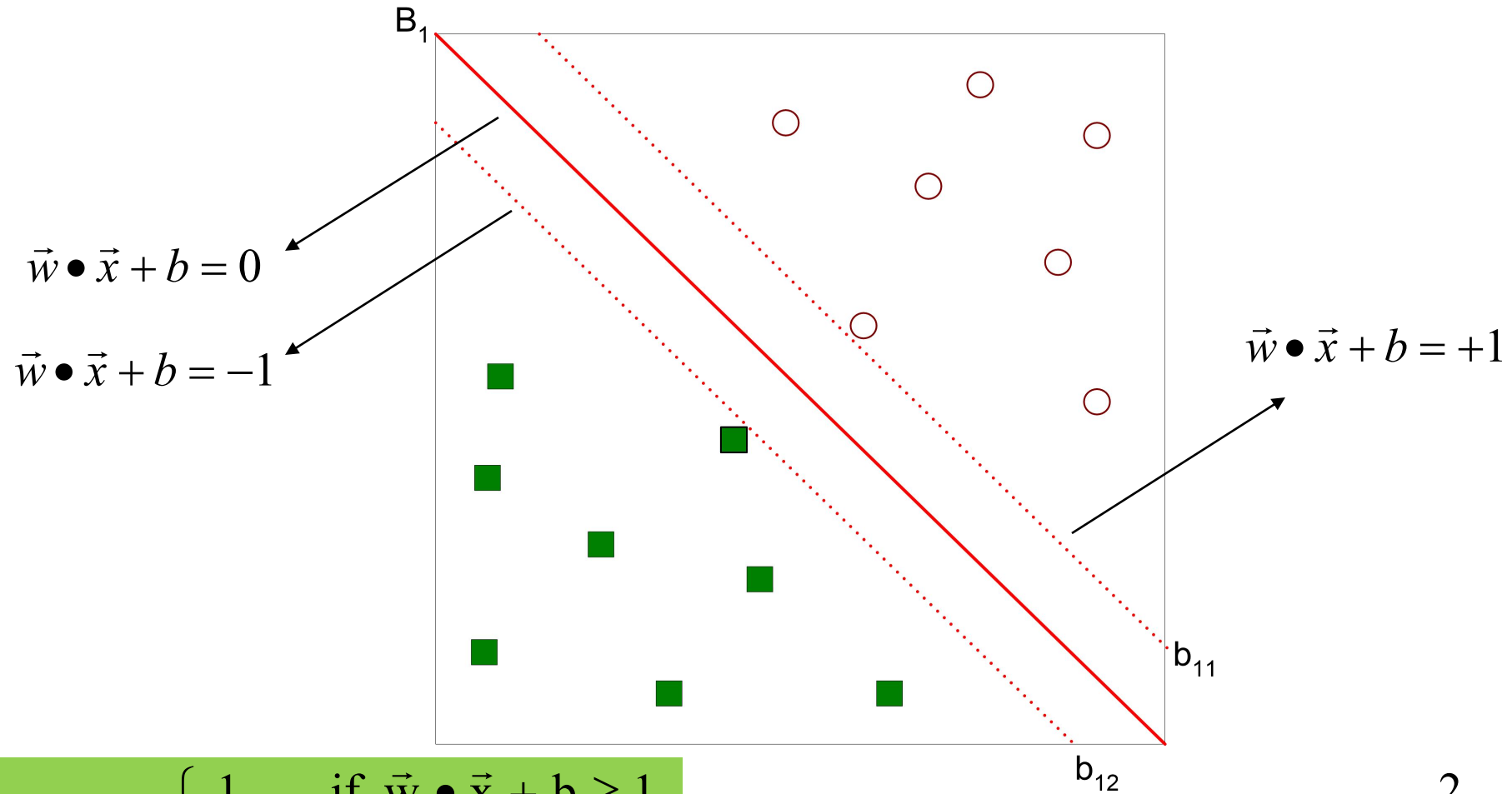
- Which one is better?  $B_1$  or  $B_2$ ?
- How do you define better?

# Support Vector Machines



- Find hyperplane **maximizes** the margin  $\Rightarrow$  B1 is better than B2

# Support Vector Machines



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

$$\text{Margin} = \frac{2}{\|\vec{w}\|}$$

# Support Vector Machines

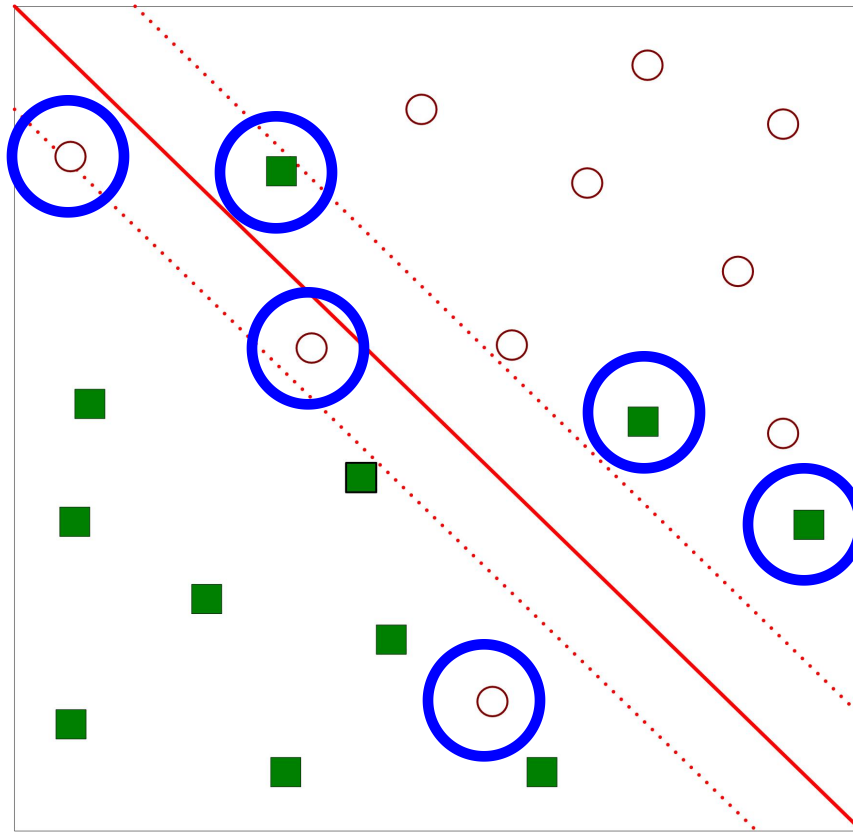
- We want to **maximize**:  $\text{Margin} = \frac{2}{\|\vec{w}\|^2}$
- Which is equivalent to **minimizing**:  $L(w) = \frac{\|\vec{w}\|^2}{2}$
- But subjected to the following **constraints**:

$$\begin{aligned}\vec{w} \cdot \vec{x}_i + b &\geq 1 \text{ if } y_i = 1 \\ \vec{w} \cdot \vec{x}_i + b &\leq -1 \text{ if } y_i = -1\end{aligned}$$

- This is a **constrained optimization problem**
  - Numerical approaches to solve it (e.g., **quadratic programming**)

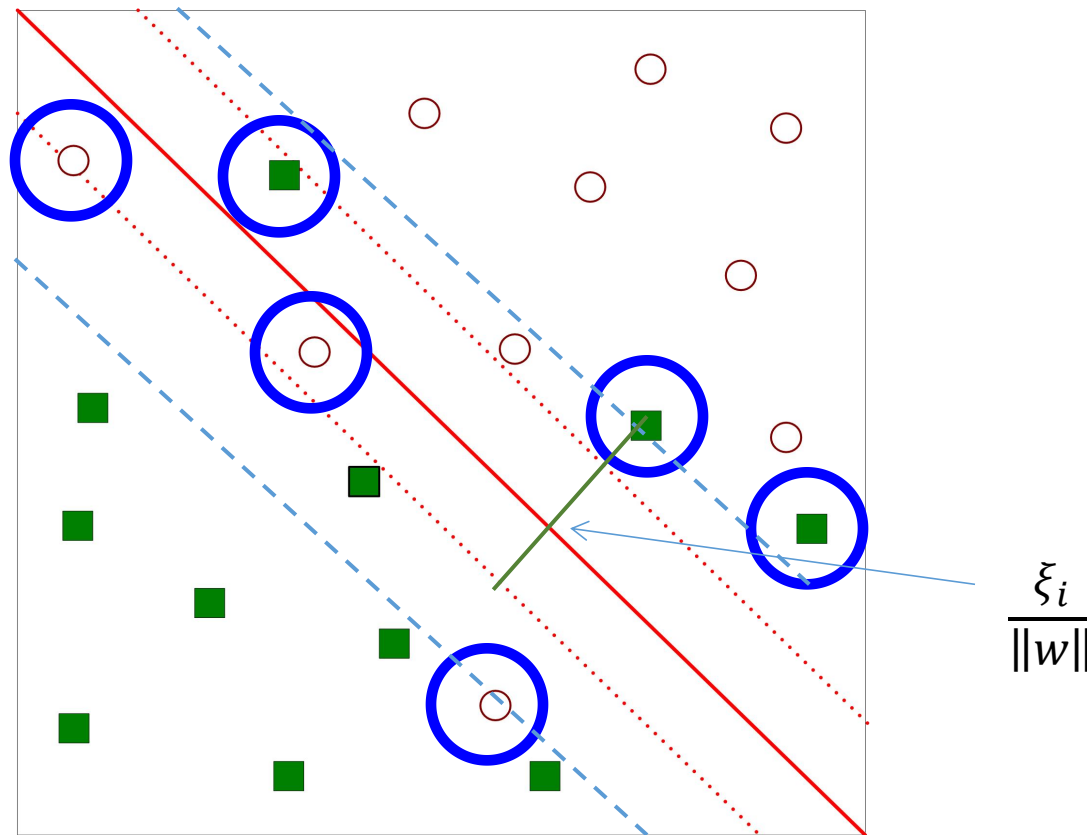
# Support Vector Machines

- What if the problem is **not linearly separable**?



# Support Vector Machines

- What if the problem is not linearly separable?



# Support Vector Machines

- What if the problem is not linearly separable?
  - Introduce slack variables
    - Need to minimize:

$$L(w) = \frac{\|\vec{w}\|^2}{2} + C \left( \sum_{i=1}^N \xi_i \right)$$

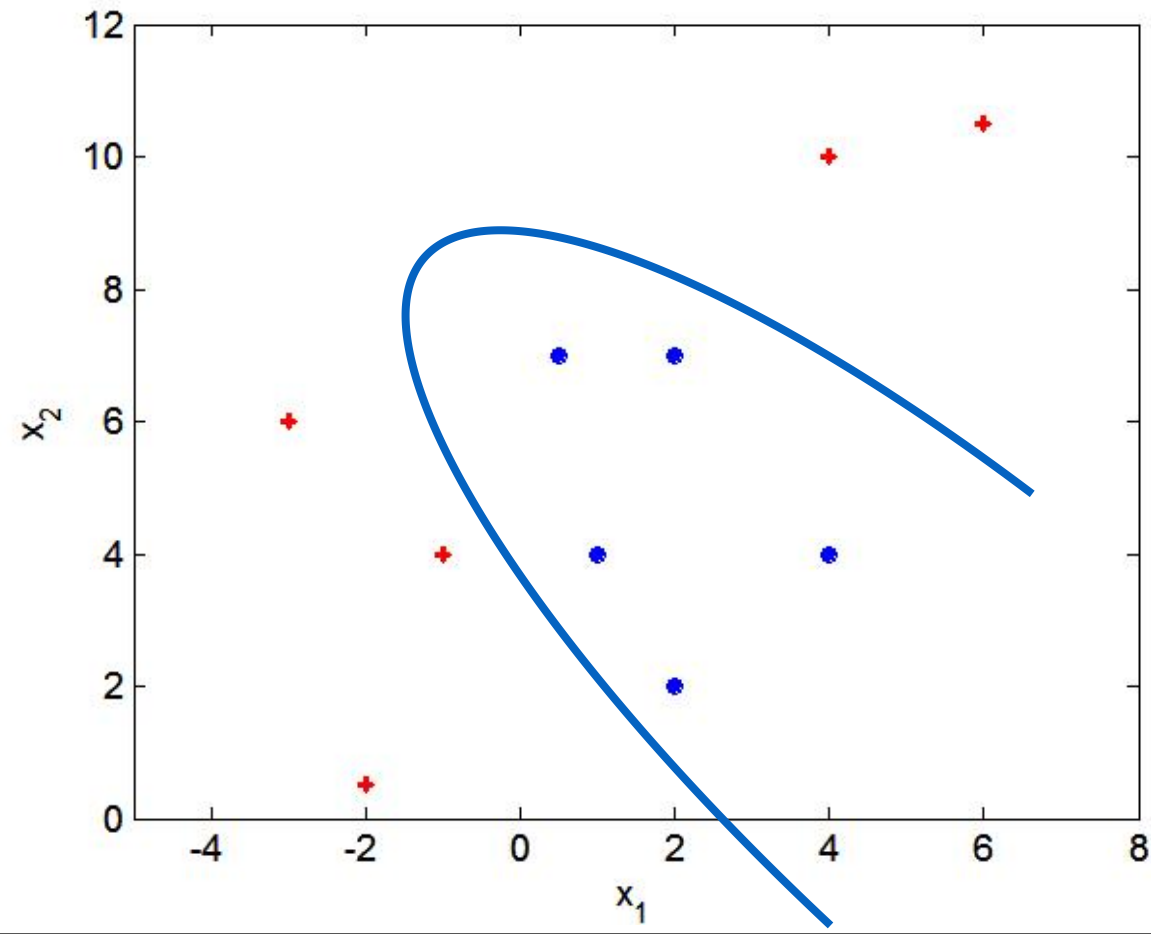
- Subject to:

$$\begin{aligned} \vec{w} \cdot \vec{x}_i + b &\geq 1 - \xi_i \text{ if } y_i = 1 \\ \vec{w} \cdot \vec{x}_i + b &\leq -1 + \xi_i \text{ if } y_i = -1 \end{aligned}$$



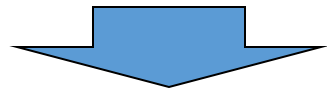
# Nonlinear Support Vector Machines

- What if decision boundary is not linear?

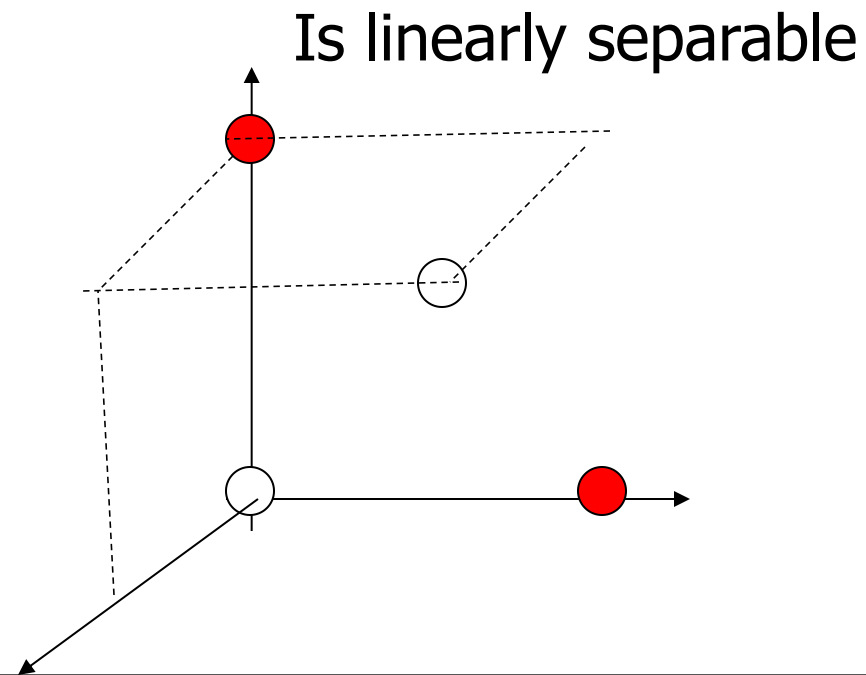
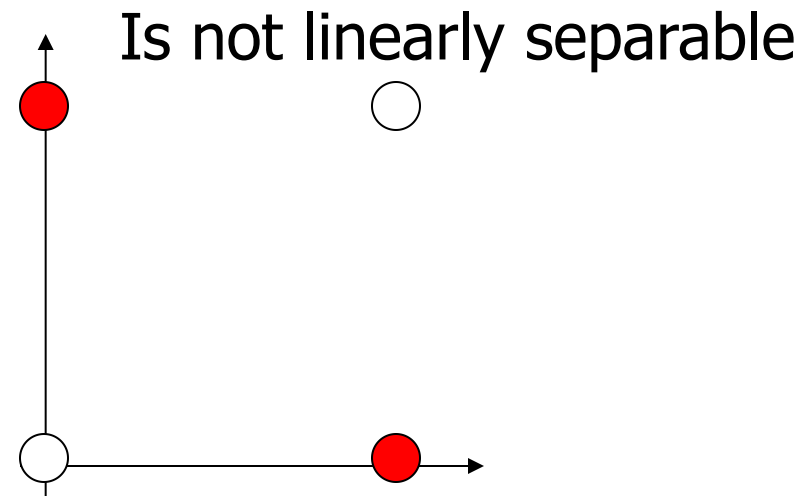


# XOR

X	Y	
0	0	0
0	1	1
1	0	1
1	1	0



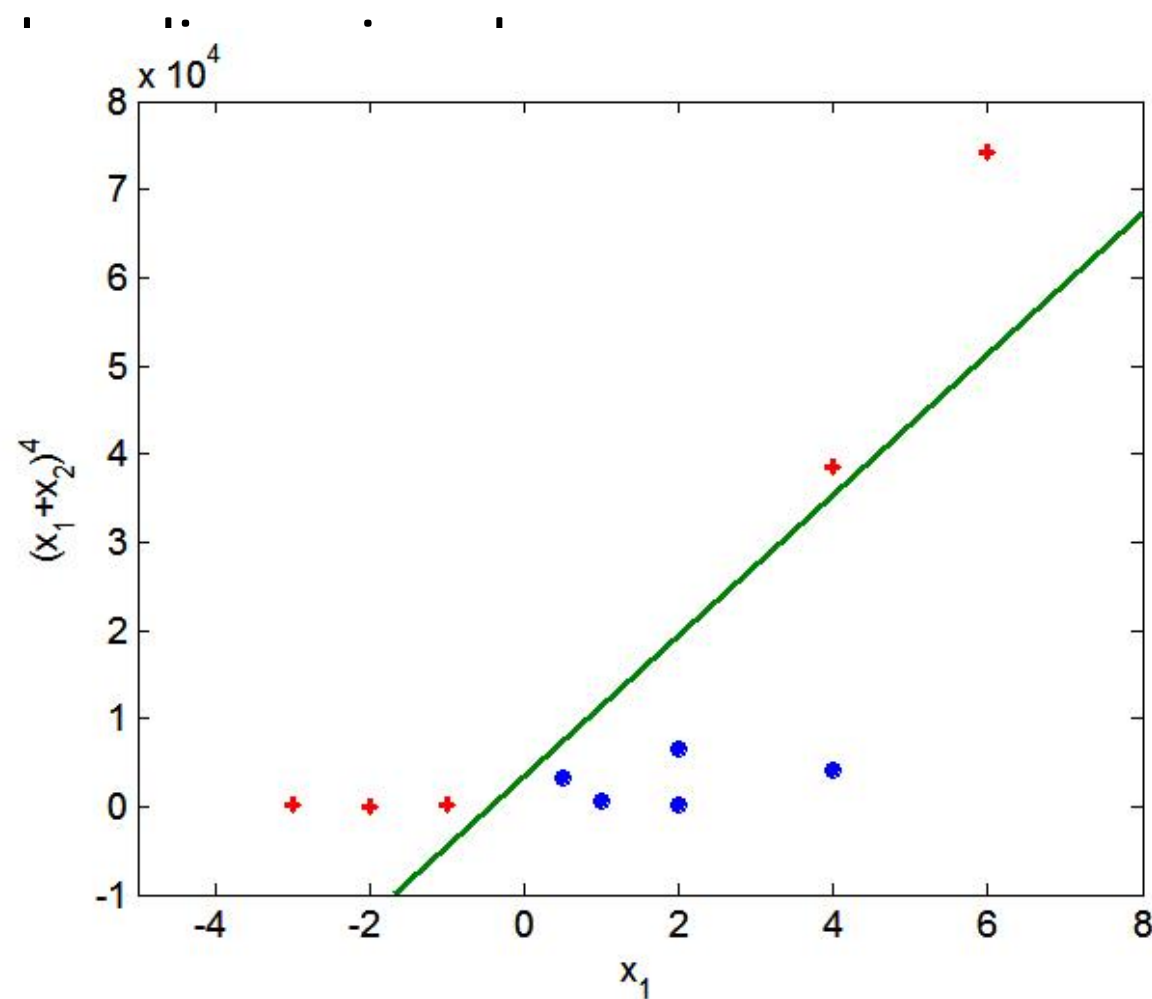
X	Y	XY	
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0



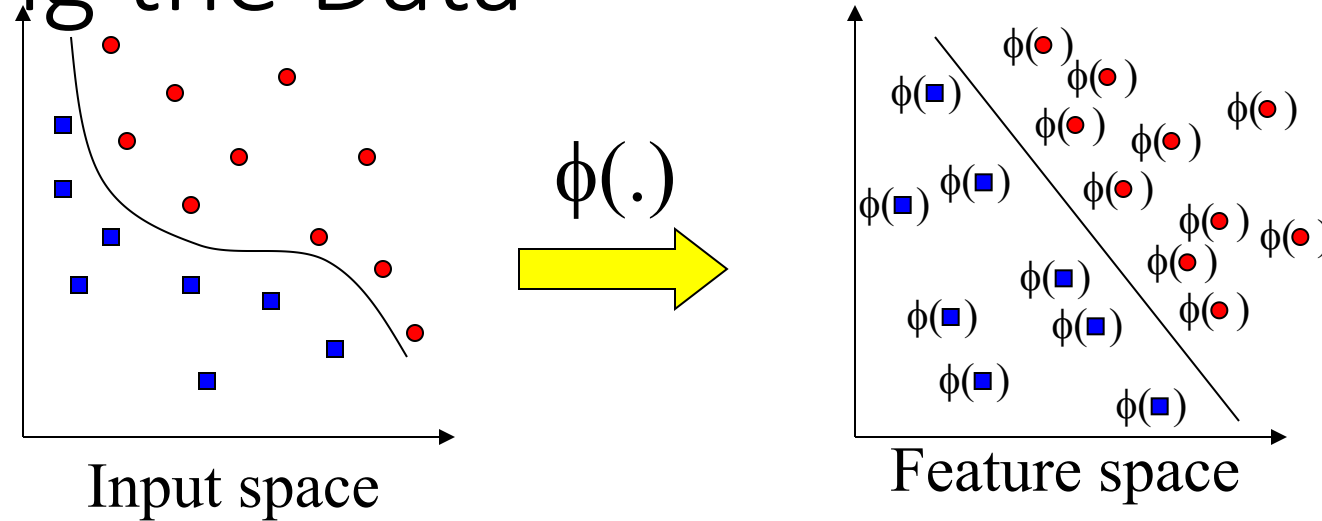
# Nonlinear Support Vector Machines

- Transform data into  $h$

Use the **Kernel Trick**



# Transforming the Data



Note: feature space is of higher dimension than the input space in practice

- Computation in the feature space can be costly because it is high dimensional
  - The feature space is typically infinite-dimensional!
- The kernel trick comes to rescue

# Kernel trick

- **Kernel:** A kernel is a method of placing a two dimensional plane into a higher dimensional space, so that it is curved in the higher dimensional space. (In simple terms, a kernel is a function from the low dimensional space into a higher dimensional space.)

# The Kernel Trick

- Recall the SVM optimization problem  $\max. W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$   
subject to  $C \geq \alpha_i \geq 0, \sum_{i=1}^n \alpha_i y_i = 0$

- The data points only appear as **inner product**
- As long as we can calculate the inner product in the feature space, we do not need the mapping explicitly
- Many common geometric operations (angles, distances) can be expressed by inner products
- Define the kernel function  $K$  by

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

# Strengths and Weaknesses of SVM

- Strengths

- Training is relatively easy
  - No local optimal, unlike in neural networks
- It scales relatively well to high dimensional data
- Tradeoff between classifier complexity and error can be controlled explicitly
- Non-traditional data like strings and trees can be used as input to SVM, instead of feature vectors

- Weaknesses

- Need to choose a “good” kernel function.