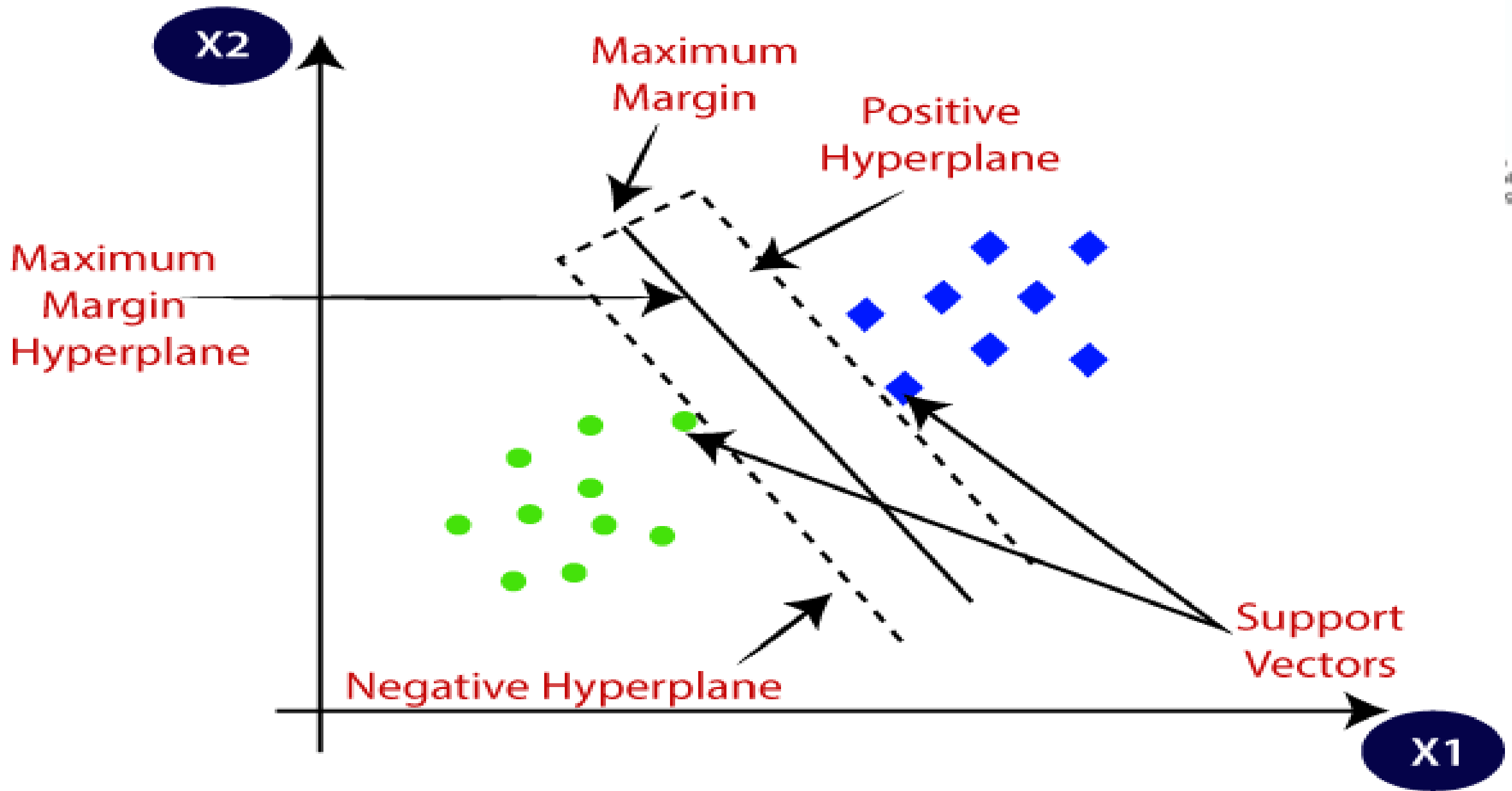# Practical Machine Learning

## Day 9: Mar22 DBDA

Kiran Waghmare

# Agenda

- SVM
- SVM-Kernel

# Support Vector Machine

- A new classification method for both linear and nonlinear data
- It uses a nonlinear mapping to transform the original training data into a higher dimension
- With the new dimension, it searches for the linear optimal separating hyperplane (i.e., "decision boundary")
- With an appropriate nonlinear mapping to a sufficiently high dimension, data from two classes can always be separated by a hyperplane
- SVM finds this hyperplane using support vectors ("essential" training tuples) and margins (defined by the support vectors)

# Support Vector Machine Algorithm

- **Goal :**

  - The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future.

  - This best decision boundary is called **a hyperplane**

  - SVM chooses the extreme points/vectors that help in creating the hyperplane

  - These extreme cases are called as **support vectors**

    and hence algorithm is termed as Support Vector Machine

- . Consider the below diagram in which there are two different categories that are classified using a decision boundary or hyperplane:

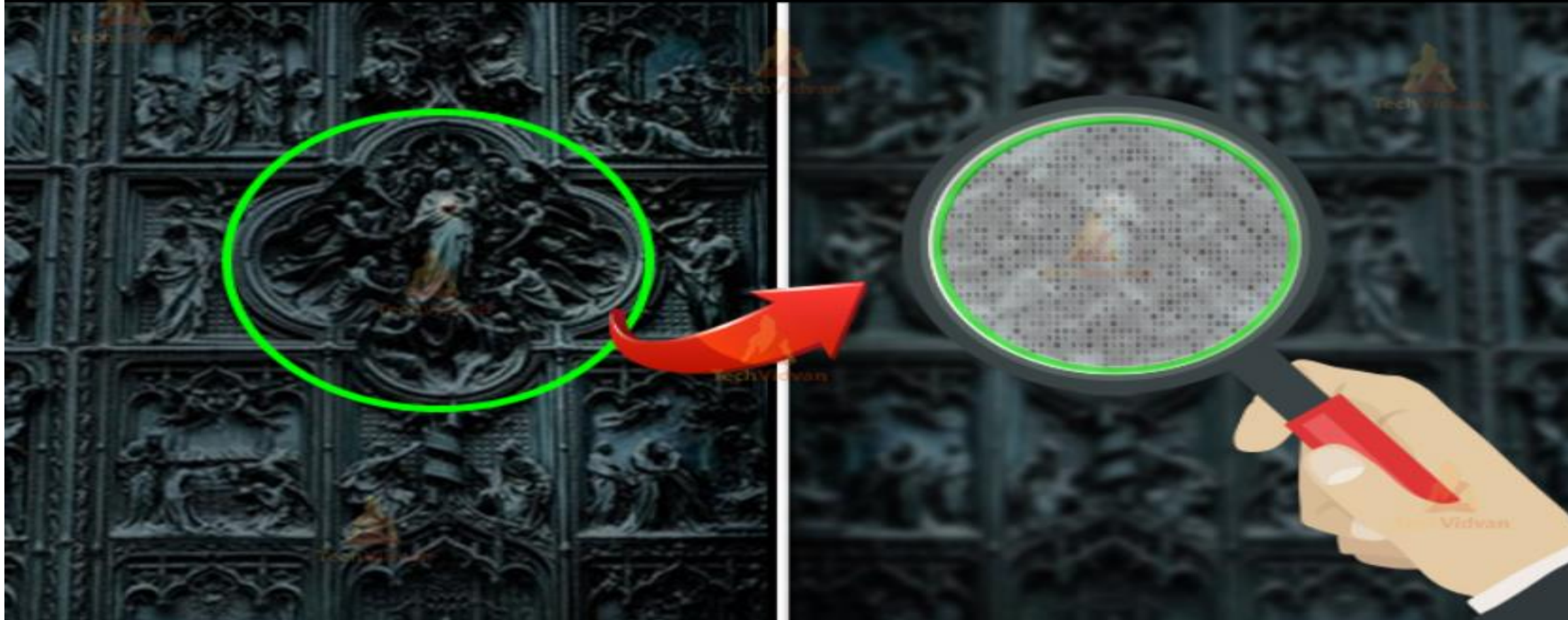# Text Classification using SVM



(a)

**Human Handwriting**

VS

(b)

**Computer Alphabets**

Text on tablets (both): It's supposed to be automatic, but actually you have to push this button.
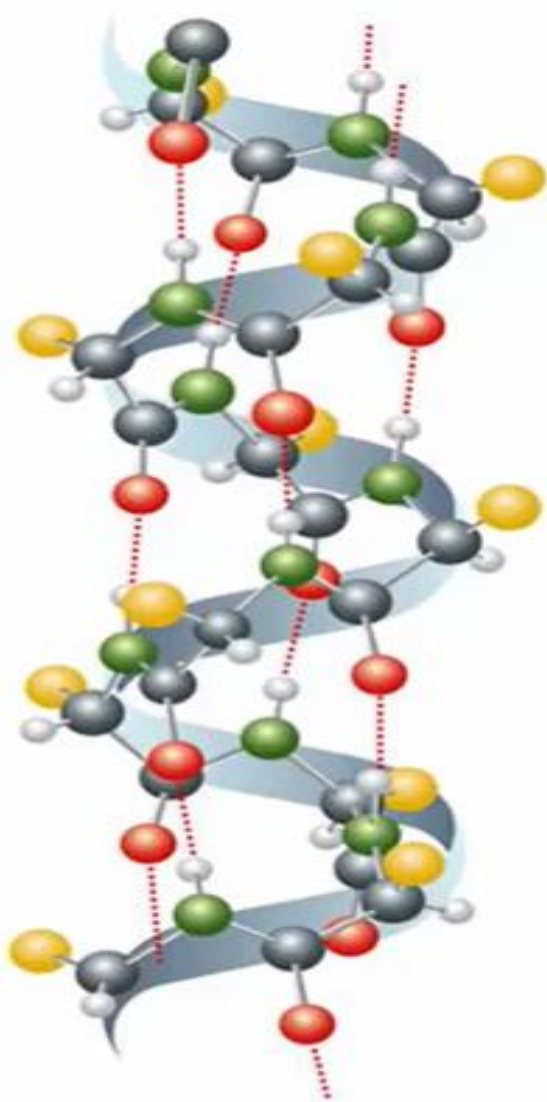
# Stenography Detection in Digital Images

La Proteina
nella sua struttura molecolare secondaria
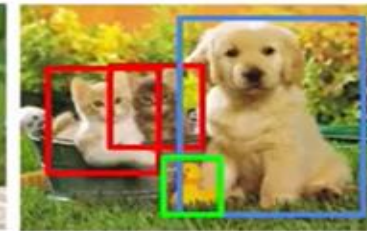*(secondary molecular structure of the protein)*

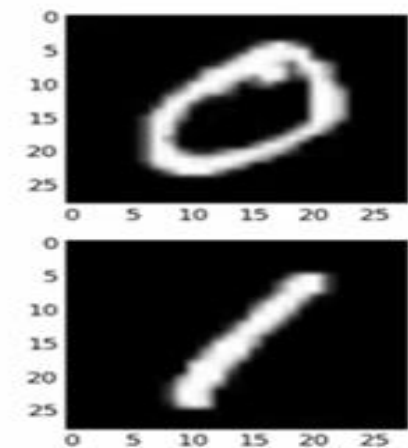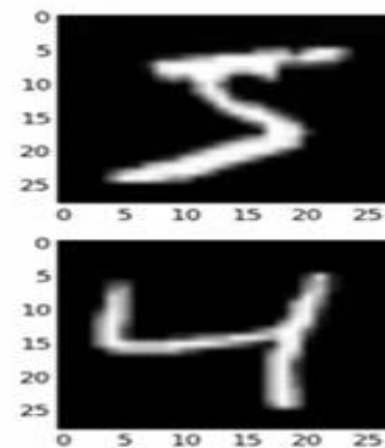**Classification** | **Classification + Localization** | **Object Detection** | **Instance Segmentation**

CAT | CAT | CAT, DOG, DUCK | CAT, DOG, DUCK

Single object | Multiple objects

Ossigeno *(oxygen)*
Carbonio *(carbon)*
Nitrogeno *(nitrogen)*
Aminoacidi *(aminoacid side chain)*
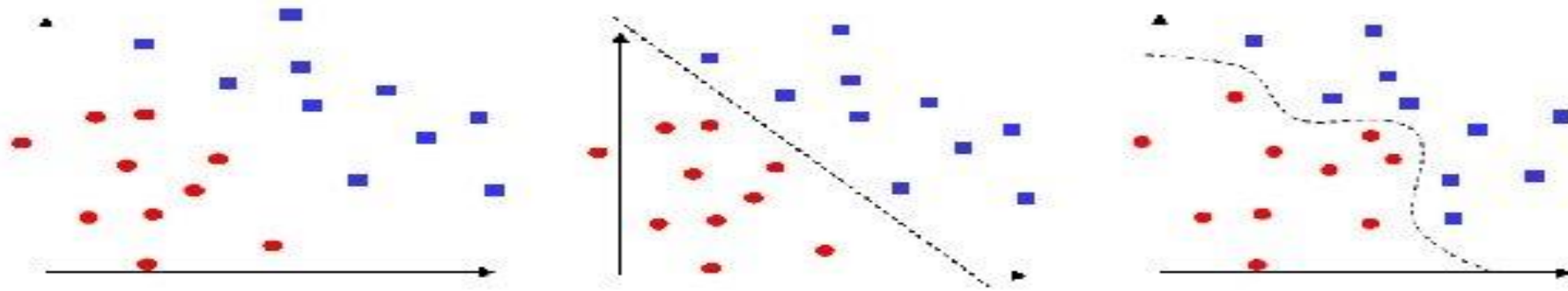Idrogeno *(hydrogen)*
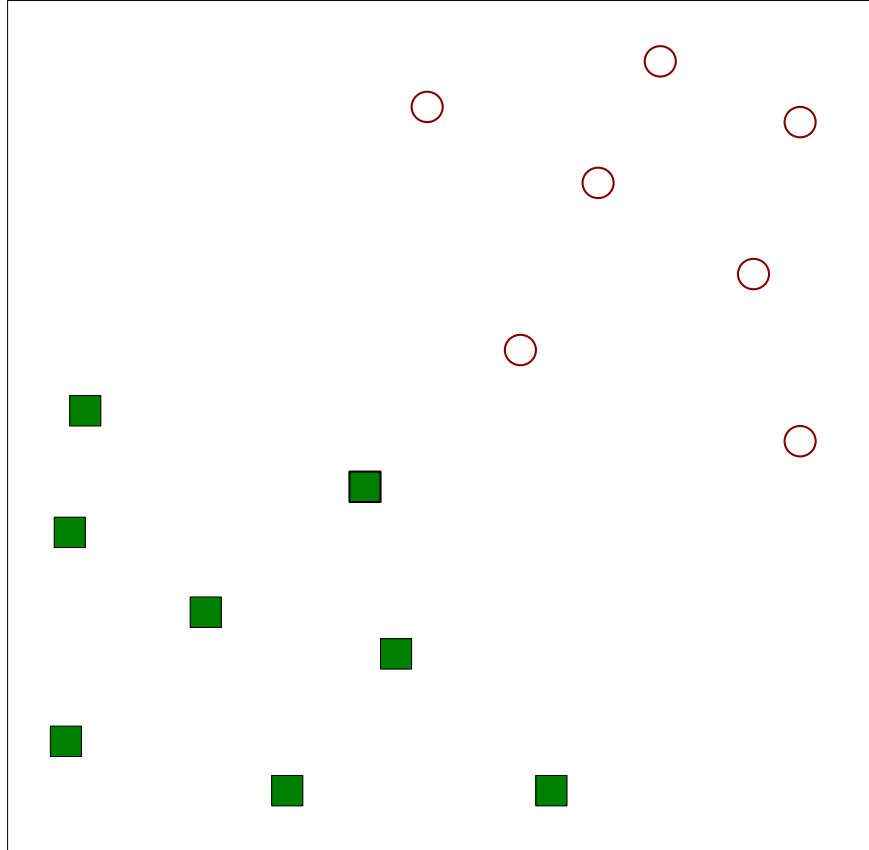
# Support Vector Machine (SVM)

-- classifier, forward neural network, supervised learning

**Difficulties** with SVM:

i) binary classifier,  ii) linearly separable patterns
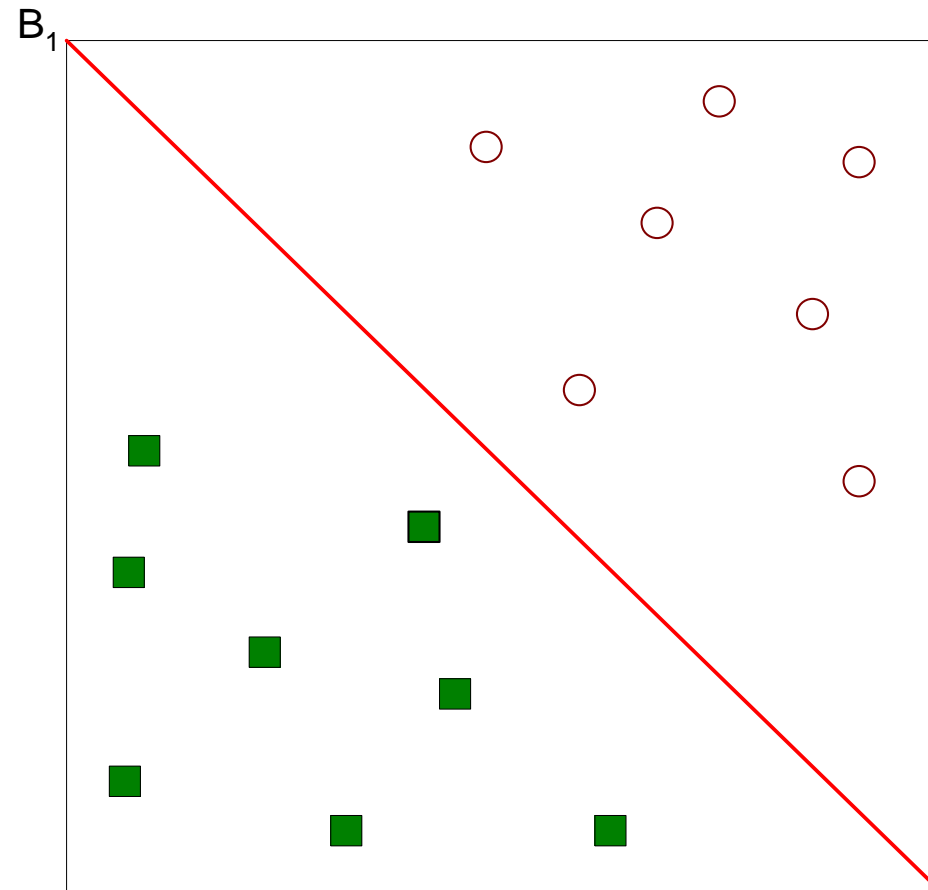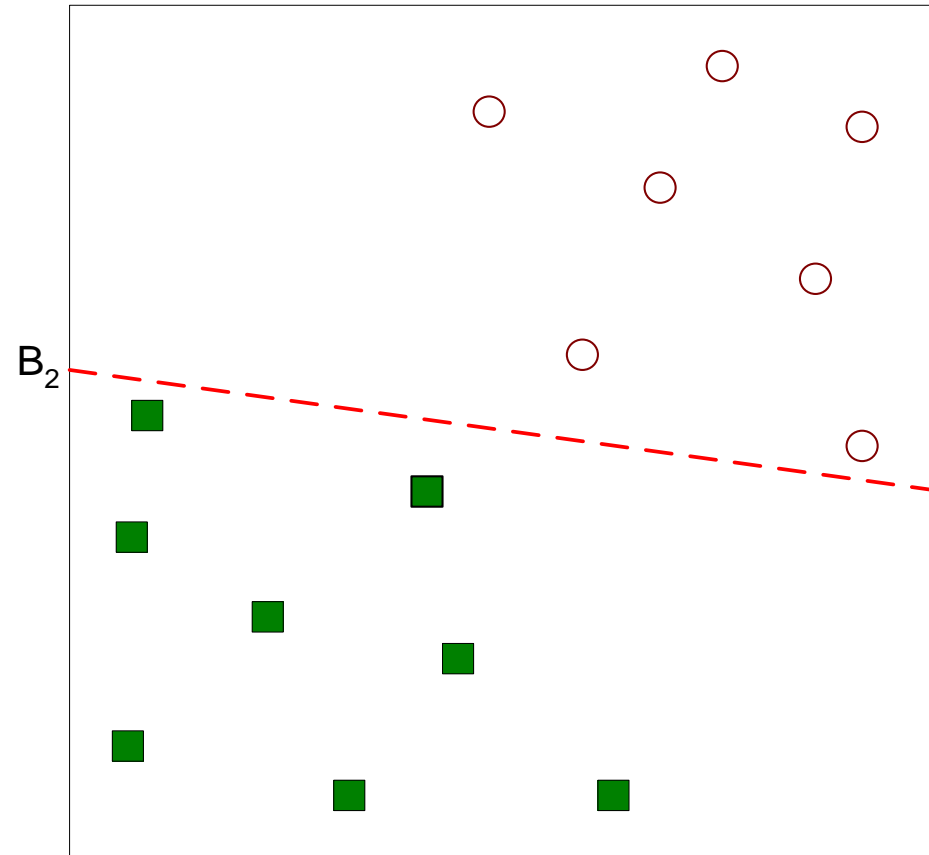


1

# Support Vector Machines



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

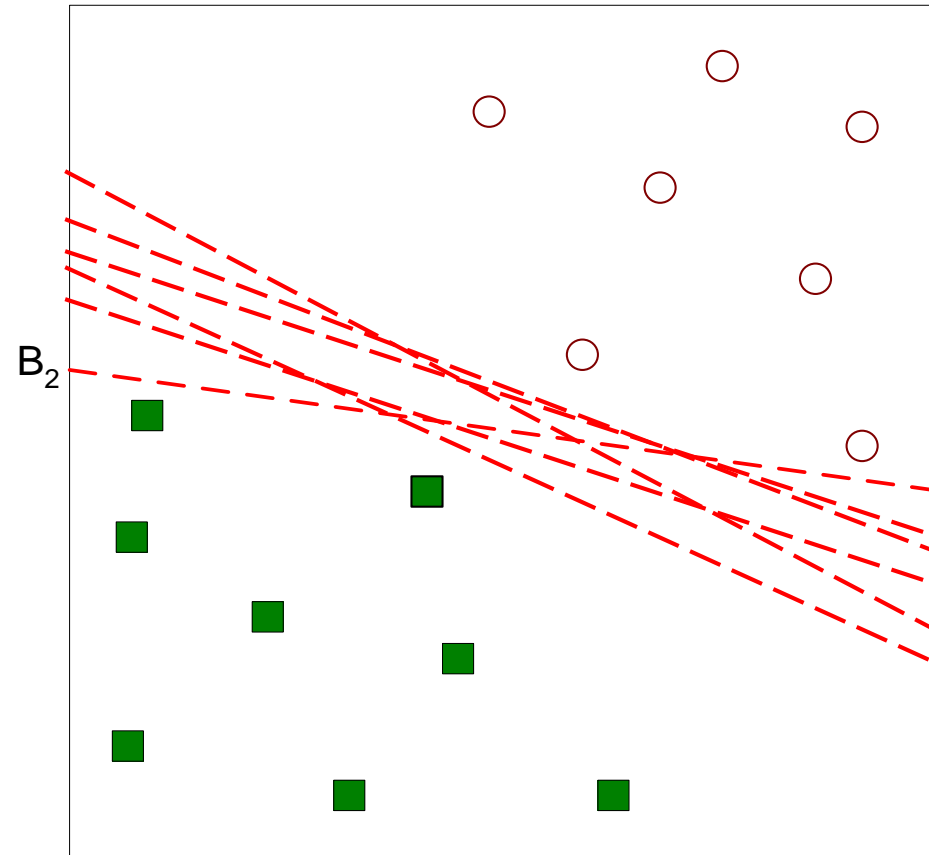# Support Vector Machines



- One Possible Solution
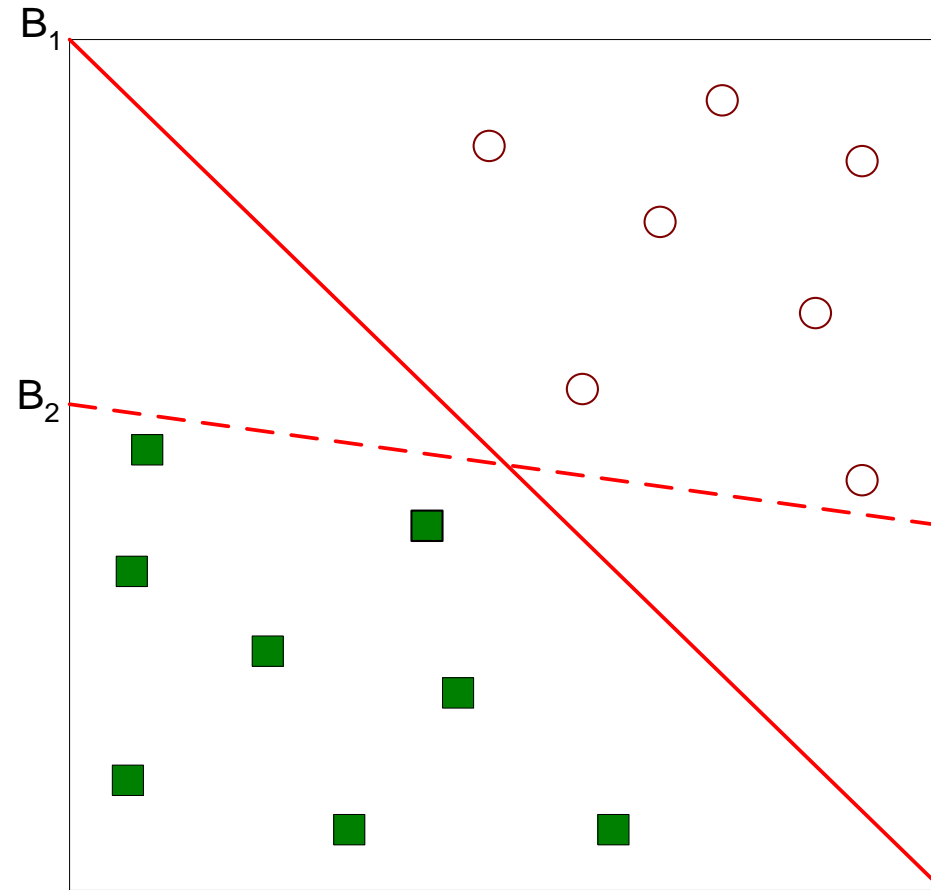
# Support Vector Machines



- Another possible solution
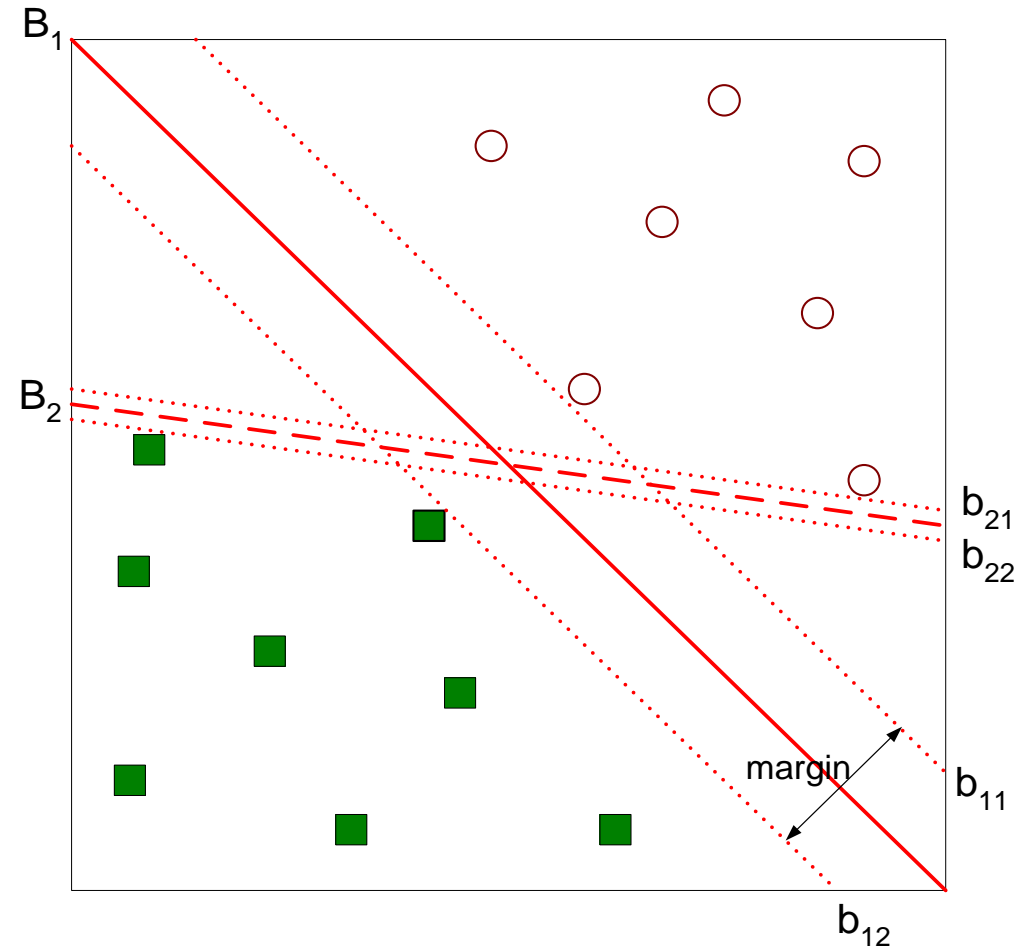
# Support Vector Machines
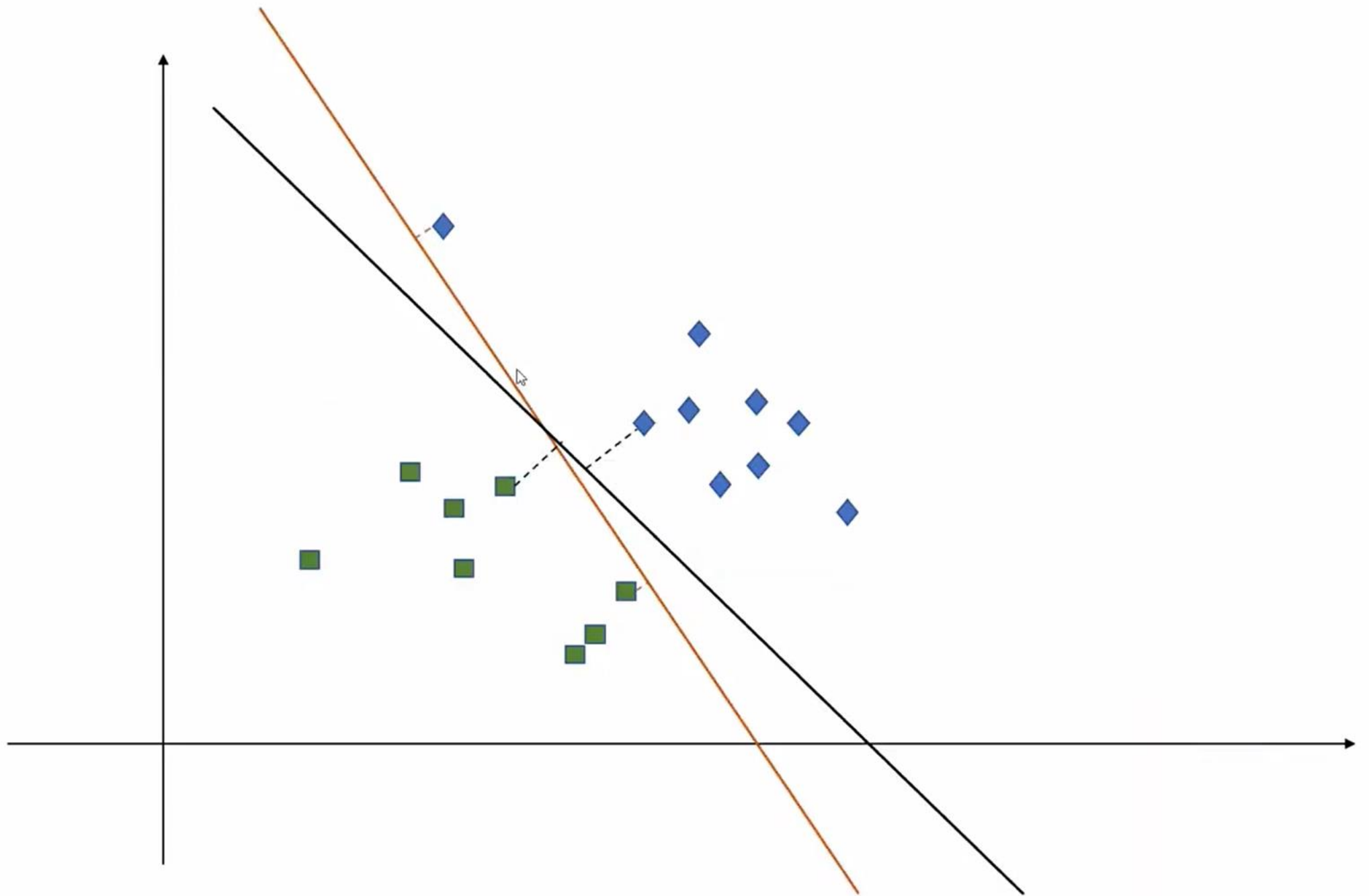


- Other possible solutions

# Support Vector Machines
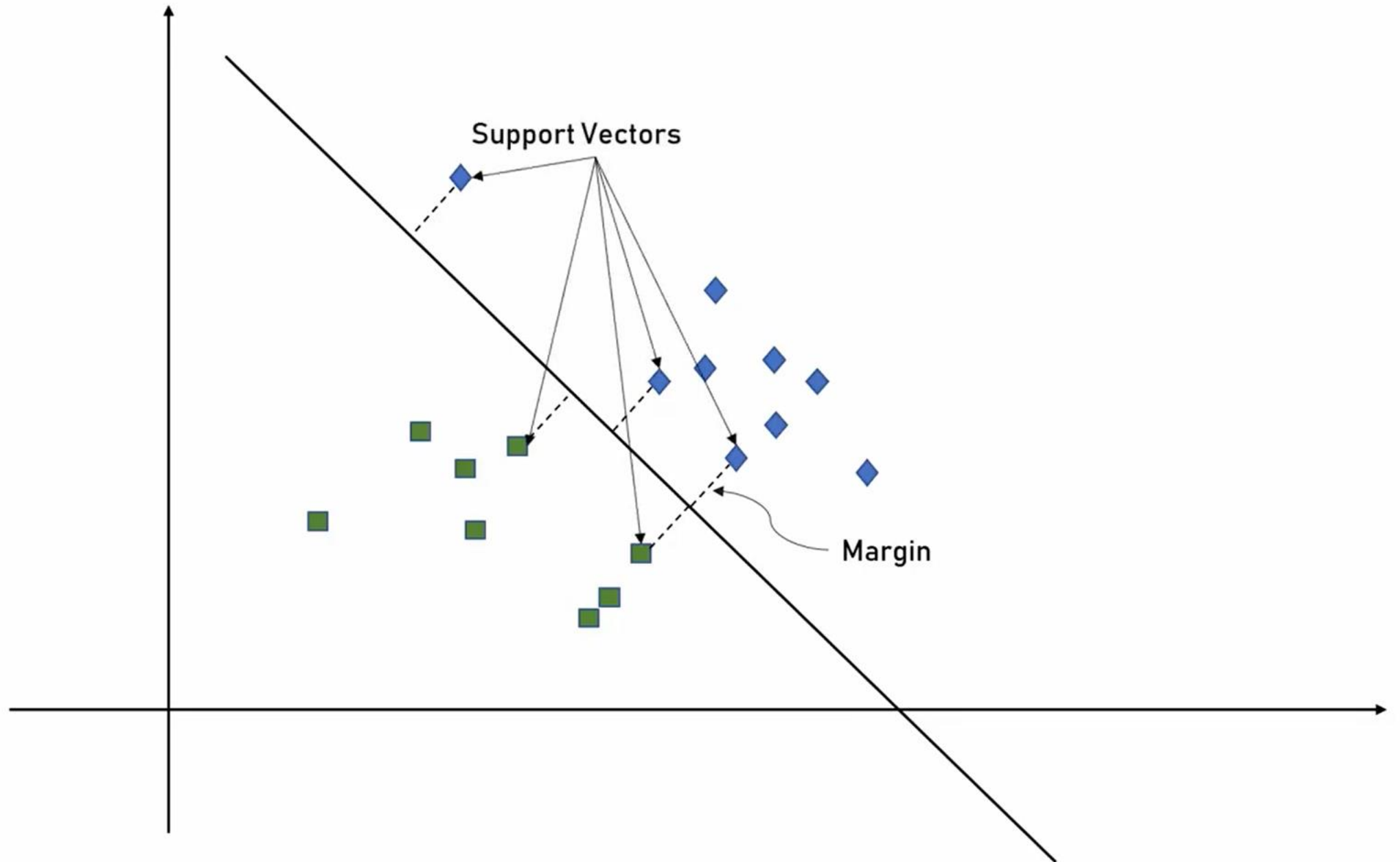


- Which one is better? B1 or B2?
- How do you define better?

# Support Vector Machines
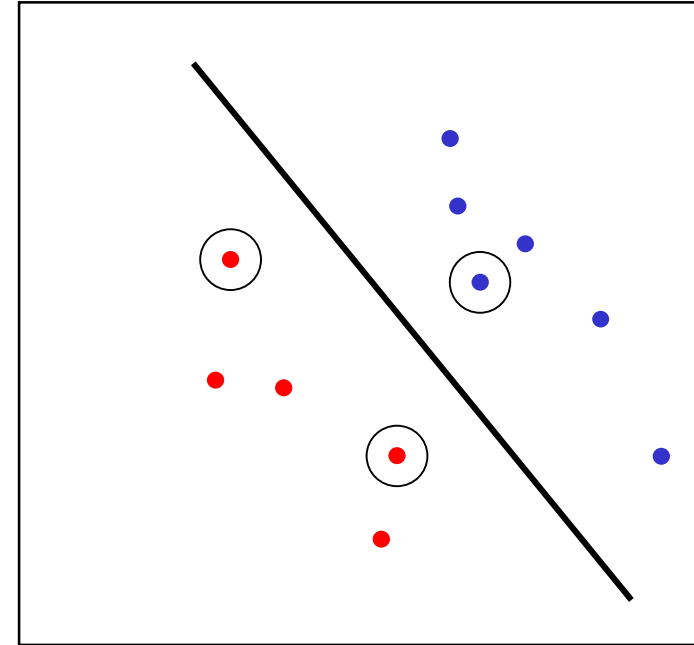


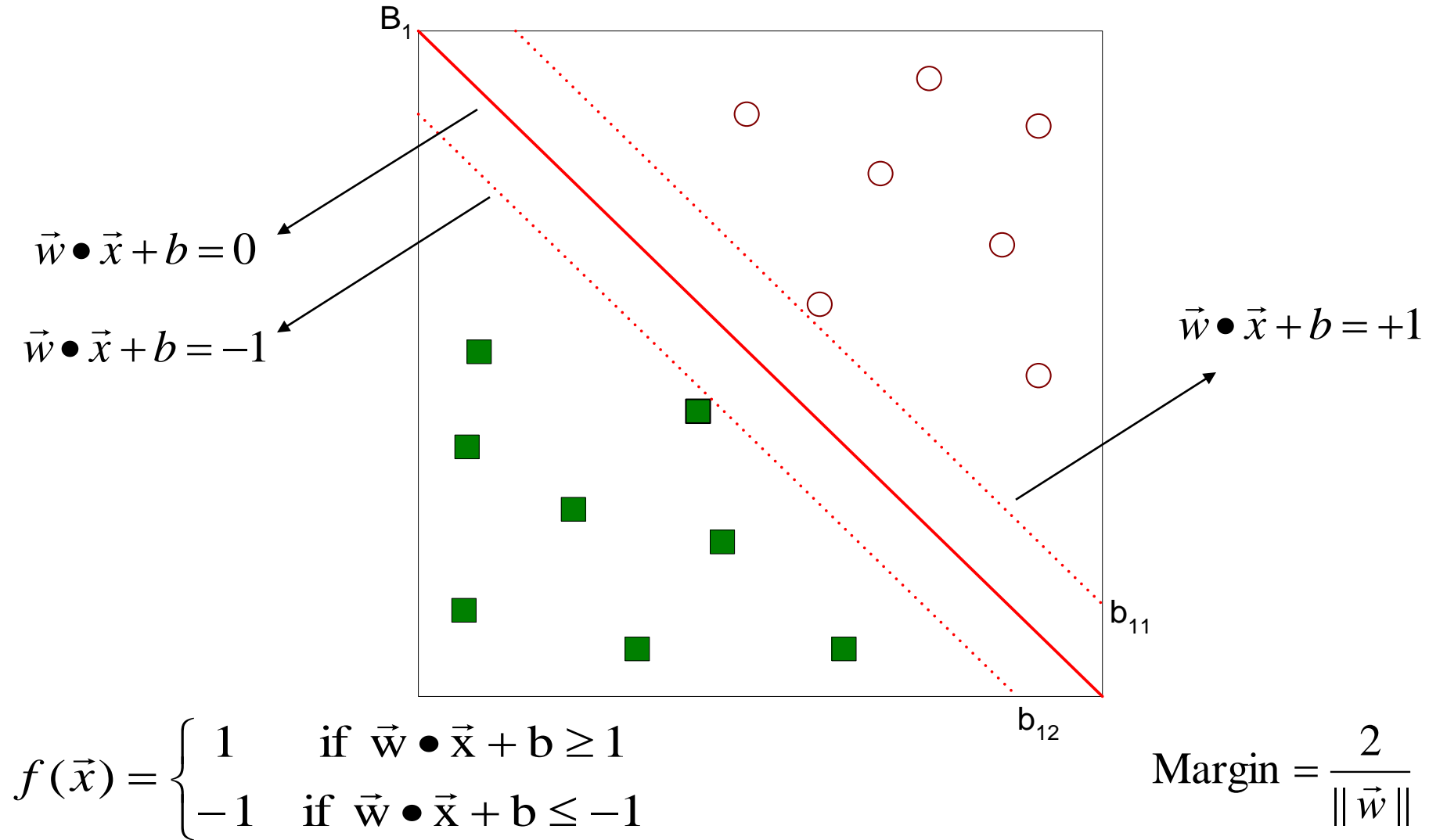- Find hyperplane maximizes the margin => B1 is better than B2

# Support Vector Machines

- The line that maximizes the minimum margin is a good bet.
  - The model class of "hyper-planes with a margin of m" has a low VC dimension if m is big.
- This maximum-margin separator is determined by a subset of the datapoints.
  - Datapoints in this subset are called "support vectors".
  - It will be useful computationally if only a small fraction of the datapoints are support vectors, because we use the support vectors to decide which side of the separator a test case is on.



The support vectors are indicated by the circles around them.

# Support Vector Machines



$B_1$

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

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

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

$b_{11}$

$b_{12}$

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

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

# Training a linear SVM

- To find the maximum margin separator, we have to solve the following optimization problem:

$$\mathbf{w}.\mathbf{x}^c + b > +1 \quad for \ \ positive \ cases$$

$$\mathbf{w}.\mathbf{x}^c + b < -1 \quad for \ negative \ cases$$

$$and \quad \| \mathbf{w} \|^2 \ \ is \ as \ small \ as \ possible$$

- This is tricky but it's a convex problem. There is only one optimum and we can find it without fiddling with learning rates or weight decay or early stopping.
  - Don't worry about the optimization problem. It has been solved. Its called quadratic programming.
  - It takes time proportional to N^2 which is really bad for very big datasets
    - so for big datasets we end up doing approximate optimization!
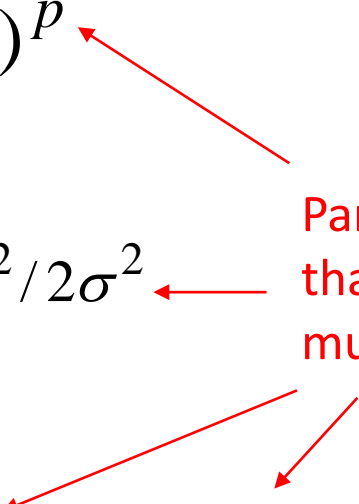
# Types of Kernel Functions

Polynomial Kernel Function — 01

Gaussian RBF Kernel Function — 02

Sigmoid Kernel Function — 03

Linear Kernel Function — 04

Hyperbolic Tangent Kernel Function — 05

Graph Kernel Function — 06

String Kernel Function — 07

Tree Kernel Function — 08

# Some commonly used kernels

Polynomial: $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}.\mathbf{y} + 1)^p$

Gaussian radial basis function $K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|^2 / 2\sigma^2}$

Parameters that the user must choose

Neural net: $K(\mathbf{x}, \mathbf{y}) = \tanh(k\,\mathbf{x}.\mathbf{y} - \delta)$

For the neural network kernel, there is one "hidden unit" per support vector, so the process of fitting the maximum margin hyperplane decides how many hidden units to use. Also, it may violate Mercer's condition.

# Introducing slack variables

- Slack variables are constrained to be non-negative. When they are greater than zero they allow us to cheat by putting the plane closer to the datapoint than the margin. So we need to minimize the amount of cheating. This means we have to pick a value for lamba (this sounds familiar!)

$$\mathbf{w}.\mathbf{x}^c + b \geq +1 - \xi^c \quad for \ positive \ cases$$

$$\mathbf{w}.\mathbf{x}^c + b \leq -1 + \xi^c \quad for \ negative \ cases$$

$$with \ \ \xi^c \geq 0 \quad for \ all \ c$$

$$and \ \ \frac{\|\mathbf{w}\|^2}{2} + \lambda \sum_c \xi^c \quad as \ small \ as \ possible$$

# A picture of the best plane with a slack variable

# 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:

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

    - Subject to:

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

    - If k is 1 or 2, this leads to similar objective function as linear SVM but with different constraints (see textbook)

# Support Vector Machines



- Find the hyperplane that optimizes both factors

# 0 dimensions:
## POINT

# 1 dimension:
## LINE SEGMENT

# 2 dimensions:
## SQUARE

# 3 dimensions:
## CUBE

# 4 dimensions:
## TESSERACT

# Kernal Trick (SVM)...



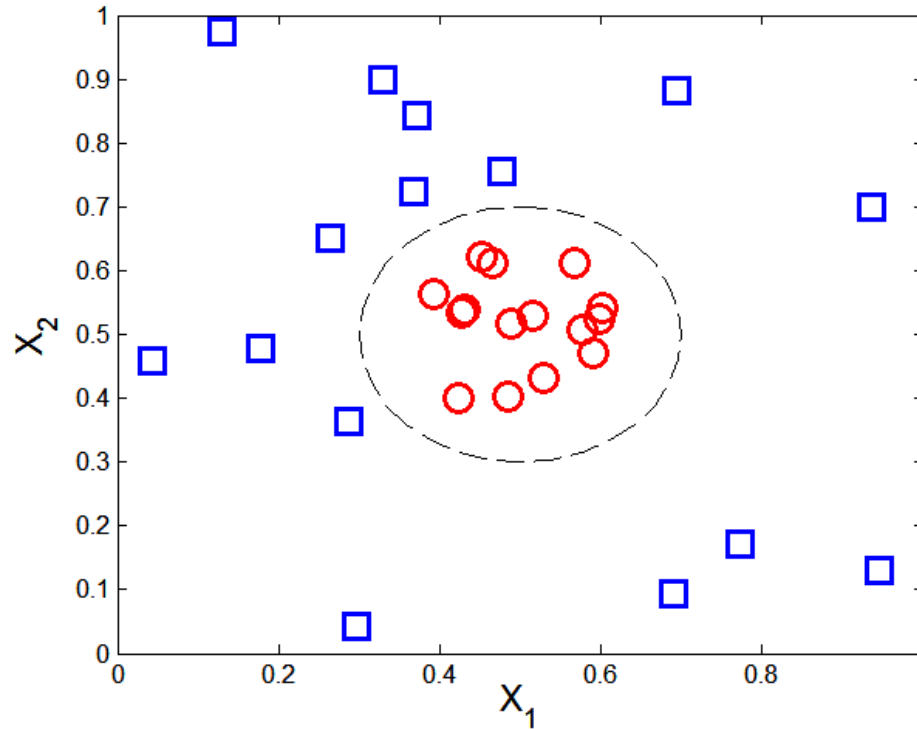(a) Input space ·········· map ·········· ▶ (b) Feature space

# Problems with linear SVM



=-1

=+1

What if the decision function is not a linear?

# Nonlinear Support Vector Machines

- What if decision boundary is not linear?



$$y(x_1, x_2) = \begin{cases} 1 & \text{if } \sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} > 0.2 \\ -1 & \text{otherwise} \end{cases}$$

# Nonlinear Support Vector Machines

- Transform data into higher dimensional space



$$x_1^2 - x_1 + x_2^2 - x_2 = -0.46.$$

$$\Phi : (x_1, x_2) \longrightarrow (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1).$$

$$w_4 x_1^2 + w_3 x_2^2 + w_2 \sqrt{2}x_1 + w_1 \sqrt{2}x_2 + w_0 = 0.$$

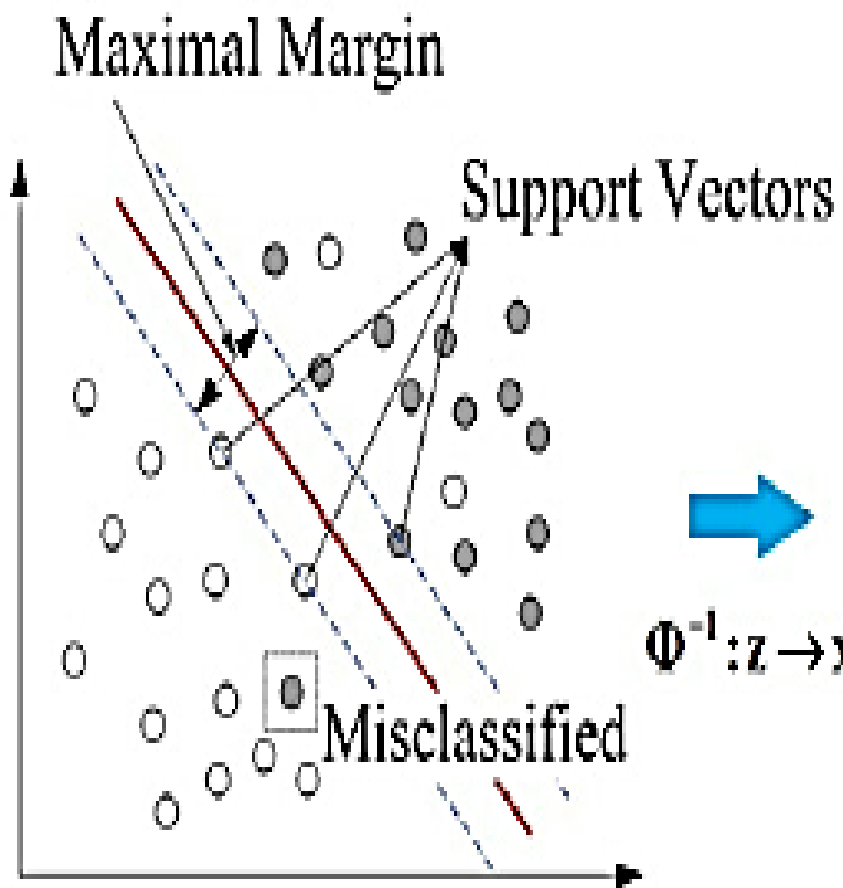Decision boundary:
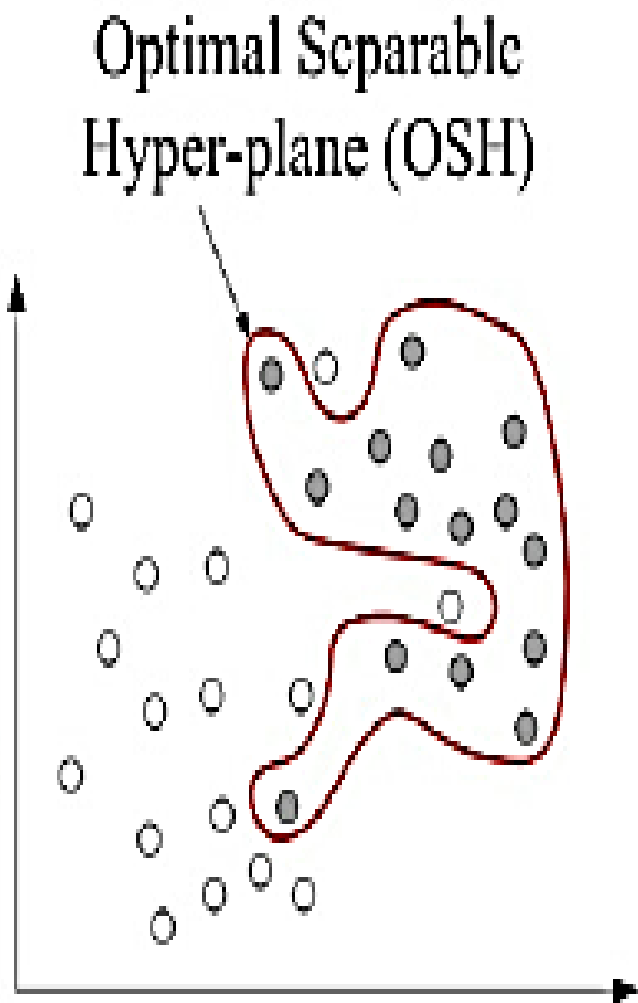
$$\vec{w} \bullet \Phi(\vec{x}) + b = 0$$

# 2D

# 3D



Image Credit: https://appliedmachinelearning.blog

Original feature space
$R^n : x$

$\Phi : x \rightarrow z$
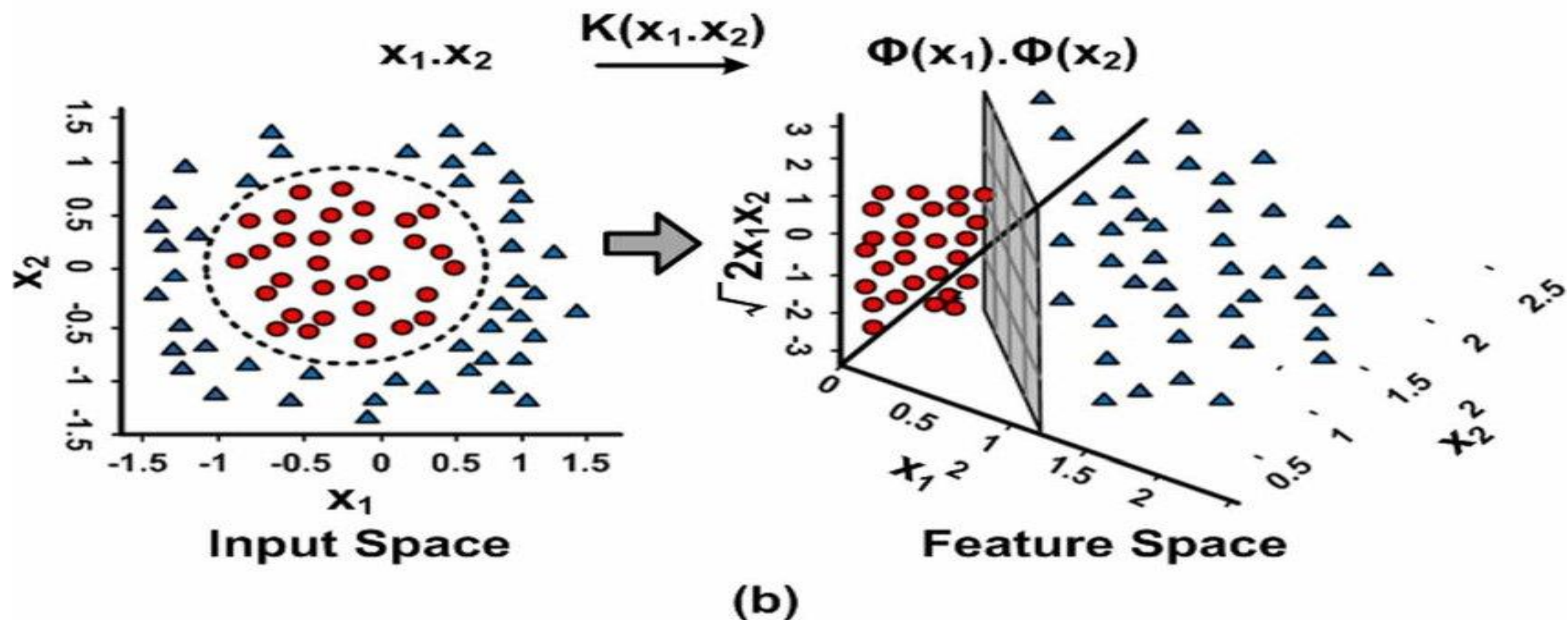
Maximal Margin

Support Vectors

Misclassified

High dimensional space
$R^n : z$ (Linear SVM)

$\Phi^{-1} : z \rightarrow x$

Optimal Separable
Hyper-plane (OSH)
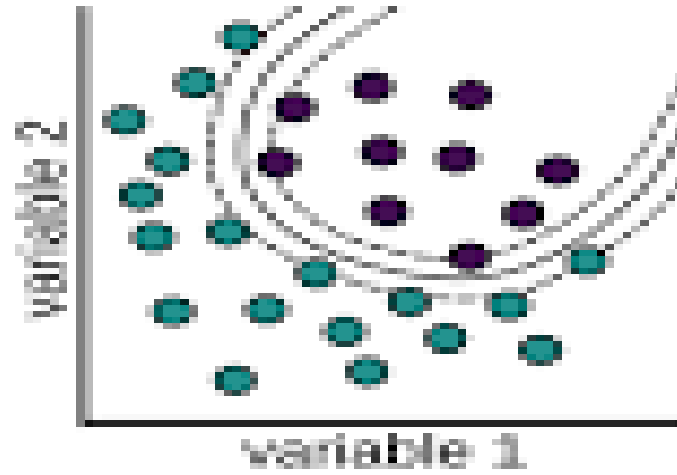
Original feature space
$R^n : x$ (Non-linear SVM)

**Input Space**

**(a)**

$x_1 . x_2$ $K(x_1 . x_2)$ $\Phi(x_1) . \Phi(x_2)$

**Input Space**

**Feature Space**

**(b)**