



# Modern Machine Learning — Support Vector Machines

Kenneth E. Barner

Department of Electrical and Computer Engineering  
University of Delaware

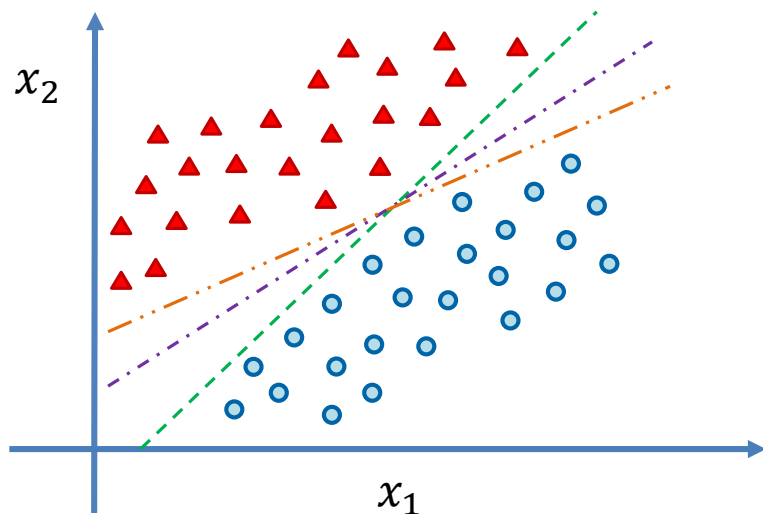


**Support Vector Machines** produce nonlinear boundaries by constructing linear boundaries in a large, transformed version of the feature space

- First derive the case for linear support vector machines with separable data
  - Constructing an **optimal separating hyperplane** between two perfectly separated classes
- Generalize result to the linear support vector machines with nonseparable data
  - Construct an **optimal separating hyperplane** focusing on samples near the decision boundary
- Finally, employ the kernel trick to transform the data, yielding nonlinear boundaries in the observation space [covered in a subsequent section]



**Assumption:** two linearly separable classes



**Hyperplane Equation**

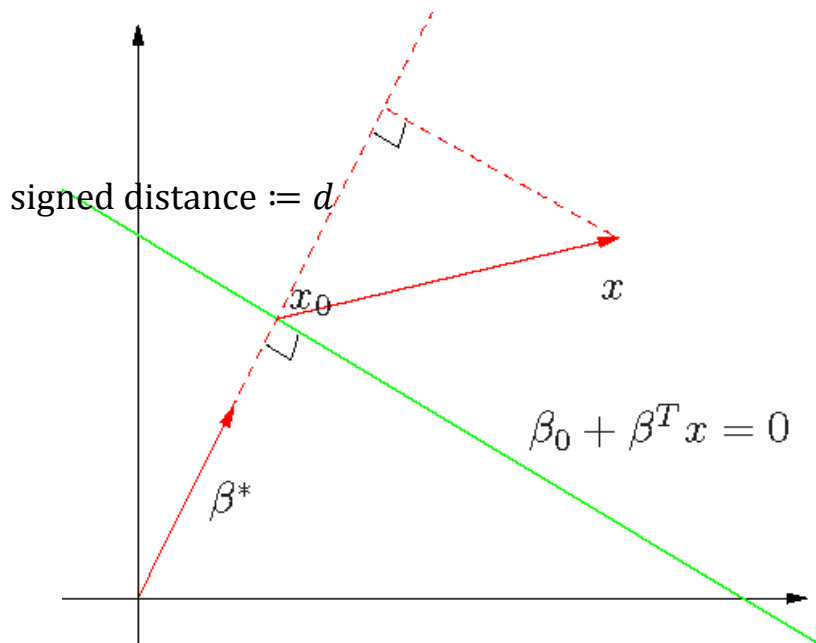
$\mathbf{x}^T \boldsymbol{\beta} + \beta_0 = 0$ , where  $\mathbf{x}^T = [x_1, \dots, x_p]$

and  $\boldsymbol{\beta} = [\beta_1, \dots, \beta_p]^T$

**Observation:** There are an infinite number of separating hyperplanes



To establish the optimal hyperplane, determine the **hyperplane vector normal** and the **signed distances** of samples to the hyperplane



The hyperplane, or *affine set*  $L$ , is defined by:

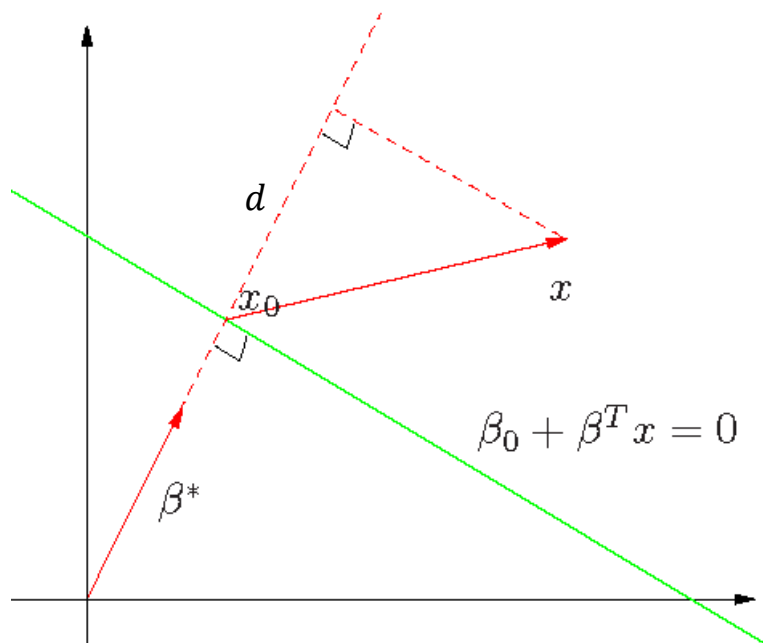
$$f(x) = \beta_0 + \mathbf{x}^T \boldsymbol{\beta} = 0$$

**Observations & Definitions:**

1. For any two points  $\mathbf{x}_1, \mathbf{x}_2 \in L$ , we have:  
 $f(\mathbf{x}_1) - f(\mathbf{x}_2) = \boldsymbol{\beta}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$   
 $\Rightarrow \boldsymbol{\beta}^* = \boldsymbol{\beta} / \|\boldsymbol{\beta}\|$  is the unit vector normal to  $L$
2. For any point  $\mathbf{x}_0 \in L$ ,  $\boldsymbol{\beta}^T \mathbf{x}_0 = -\beta_0$
3. This signed distance,  $d$ , of any point  $\mathbf{x}$  to  $L$  is

$$\begin{aligned} \boldsymbol{\beta}^{*T} (\mathbf{x} - \mathbf{x}_0) &= \frac{1}{\|\boldsymbol{\beta}\|} (\boldsymbol{\beta}^T \mathbf{x} + \beta_0) \\ &= \frac{1}{\|f'(\mathbf{x})\|} f(\mathbf{x}) \end{aligned}$$

$\Rightarrow f(\mathbf{x})$  is proportional to the signed distance from  $\mathbf{x}$  to the hyper plane defined by  $f(\mathbf{x})$ ; signed distance =  $f(\mathbf{x})$  when  $\|\boldsymbol{\beta}\| = 1$



**Observation:** Signed distances can be used as metrics for classification, i.e.

$$G(x) = \text{sgn}(d)$$

**Problem Setup:**

Binary class case:  $(y_i, \mathbf{x}_i) \in \{+1, -1\} \times \mathbb{R}^p$

Separating hyperplane:  $\beta_0 + \mathbf{x}^T \boldsymbol{\beta} = 0$ , where we enforce the normalization  $\|\boldsymbol{\beta}\| = 1$

Thus  $G(\mathbf{x}_i) = \text{sgn}(d) = \text{sgn}(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0)$

positive signed distance  $\Rightarrow G(\mathbf{x}_i) = 1$

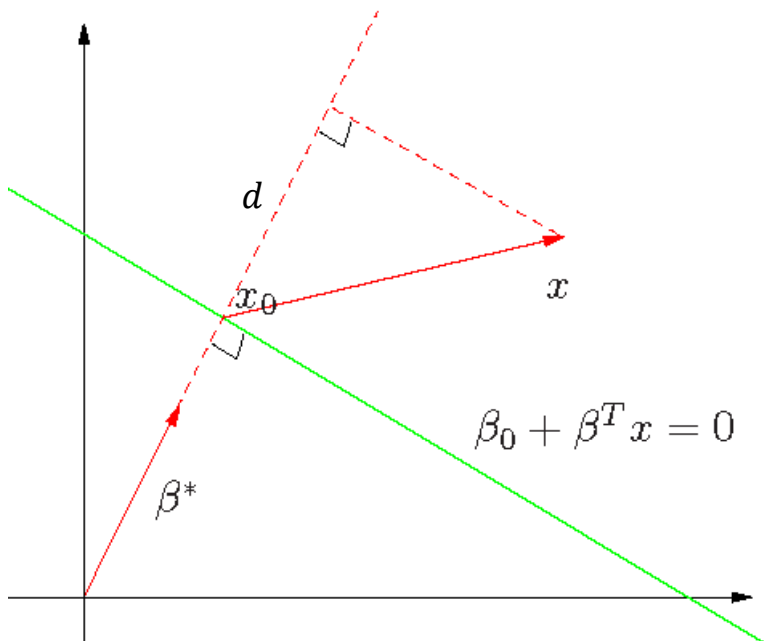
negative signed distance  $\Rightarrow G(\mathbf{x}_i) = -1$

Moreover

$G(\mathbf{x}_i)y_i = y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) > 0 \rightarrow$  correct classification

$G(\mathbf{x}_i)y_i = y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) < 0 \rightarrow$  incorrect classification

Finally, note that  $|y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0)|$  is the distance between  $\mathbf{x}_i$  and the separating hyperplane



**Perfectly Separable Case:** In this case,

$$G(x)y_i = y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) > 0 \forall i$$

**Question:** How large can we make the **margin**,  $M$ , separating points from the hyper plane?

Cast the problem as:

$$\max_{\beta_0, \boldsymbol{\beta}, \|\boldsymbol{\beta}\|=1} M$$

subject to

$$y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq M, i = 1, \dots, n$$



**Note:** The condition  $\|\boldsymbol{\beta}\| = 1$  can be included in the constraint as

$$\begin{aligned}\frac{1}{\|\boldsymbol{\beta}\|} y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) &\geq M \\ \Rightarrow y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) &\geq M \|\boldsymbol{\beta}\|\end{aligned}$$

We can always rescale  $(\beta_0, \boldsymbol{\beta})$  so that

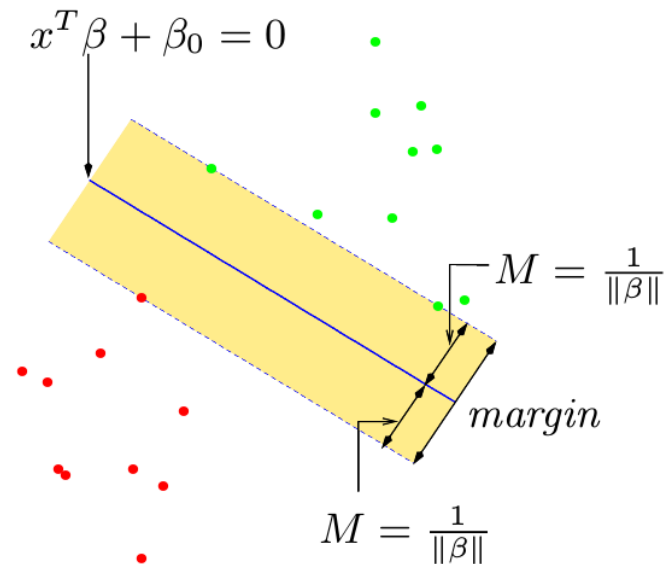
$$\|\boldsymbol{\beta}\| = \frac{1}{M}$$

Therefore, the optimization problem can be expressed as

$$\min_{\beta_0, \boldsymbol{\beta} \in \mathbb{R}^p} \|\boldsymbol{\beta}\|$$

subject to

$$y_i (\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq 1, i = 1, \dots, n$$





**Non-Separable Case:** Maximize  $M$ , but allow for some points to be on the wrong side of the margin.

Define the **slack variables**  $\xi = (\xi_1, \xi_2, \dots, \xi_n)$

Options to modify  $y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq M$

$$y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq M - \xi_i \quad (*)$$

or

$$y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq M(1 - \xi_i) \quad (**)$$

$\forall i, \xi_i \geq 0, \sum_{i=1}^n \xi_i \leq C$ , where  $C$  is a constant

**Observations:**

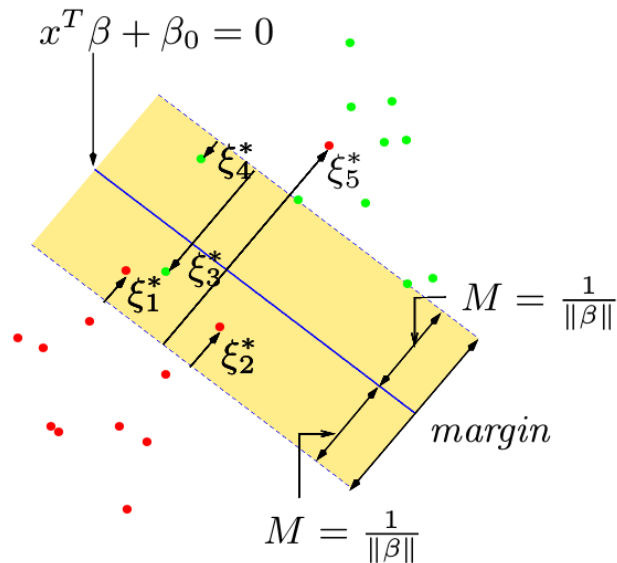
(\*) measures overlap in distance from the margin, but results in a nonconvex optimization

(\*\*) measures overlap in relative distance that scales with the margin  $M$ , but yields a convex optimization

Thus (\*\*) is referred to as the **standard support vector classifier**, which yields:

$$\min \|\boldsymbol{\beta}\| \quad \text{subject to} \quad \begin{cases} y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq 1 - \xi_i, \forall i \\ \xi_i \geq 0, \sum_{i=1}^n \xi_i \leq C \end{cases}$$

where as previously, we set  $M = 1/\|\boldsymbol{\beta}\|$





**Summary:**

- Linear **Support Vector Machines** determine the **optimal separating hyperplane** between two classes
- If the data is perfectly separable, the SVM problem statement is:

$$\min \|\boldsymbol{\beta}\| \text{ subject to } y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq 1, \forall i$$

- If the data is not separable, the SVM problem statement is:

$$\min \|\boldsymbol{\beta}\| \text{ subject to } \begin{cases} y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \geq 1 - \xi_i, \forall i \\ \xi_i \geq 0, \sum_{i=1}^n \xi_i \leq C \end{cases}$$

Where  $\xi = (\xi_1, \xi_2, \dots, \xi_n)$  are **slack variables** and  $C$  is a constant that controls the total amount of slack

- Classification is determined as  $G(x) = \text{sgn}(\mathbf{x}^T \boldsymbol{\beta} + \beta_0)$  for binary ( $\pm 1$ ) classes