

Support Vector Machine (SVM)

Machine Learning

Hamid R Rabiee – Zahra Dehghanian
Spring 2025



Sharif University
of Technology

Outline

- Margin concept
- Hard-Margin SVM
 - Dual Problem of Hard-Margin SVM
- Soft-Margin SVM
 - Dual Problem of Soft-Margin SVM

Hyperplanes

A *hyperplane* is a set of the form

$$\{x \mid a^T x = b\},$$

$a \in \mathbf{R}^n$, $a \neq 0$, and $b \in \mathbf{R}$.

a is the normal vector

Hyperplanes

A *hyperplane* is a set of the form

$$\{x \mid a^T x = b\},$$

$a \in \mathbf{R}^n$, $a \neq 0$, and $b \in \mathbf{R}$.

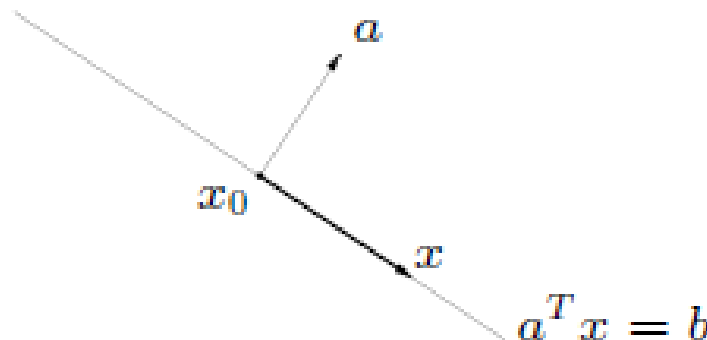
a is the normal vector

Geometrical interpretation

x_0 is any point in the hyperplane

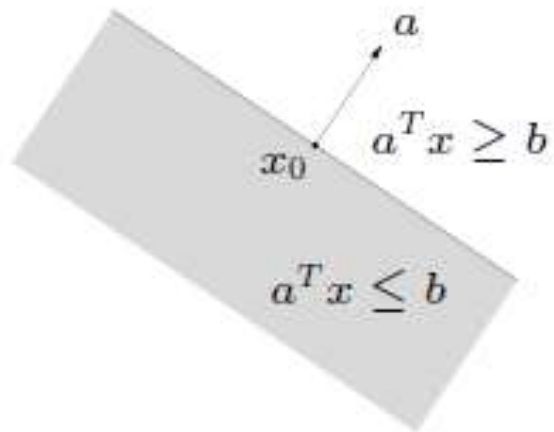
$$a^T x_0 = b$$

$$\{x \mid a^T (x - x_0) = 0\},$$



halfspaces

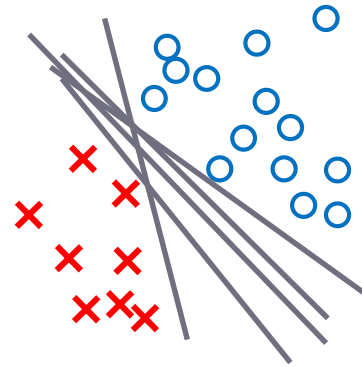
halfspace: set of the form $\{x \mid a^T x \leq b\}$ ($a \neq 0$)



Margin

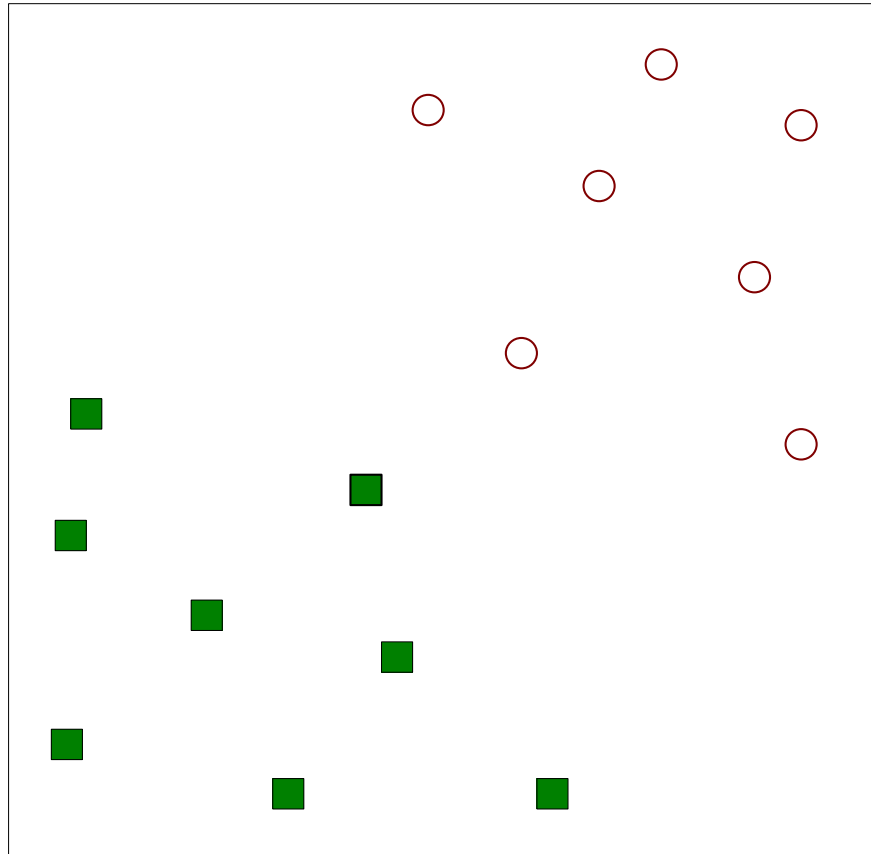
- Which line is better to select as the boundary to provide more generalization capability?

Larger margin provides better generalization to unseen data

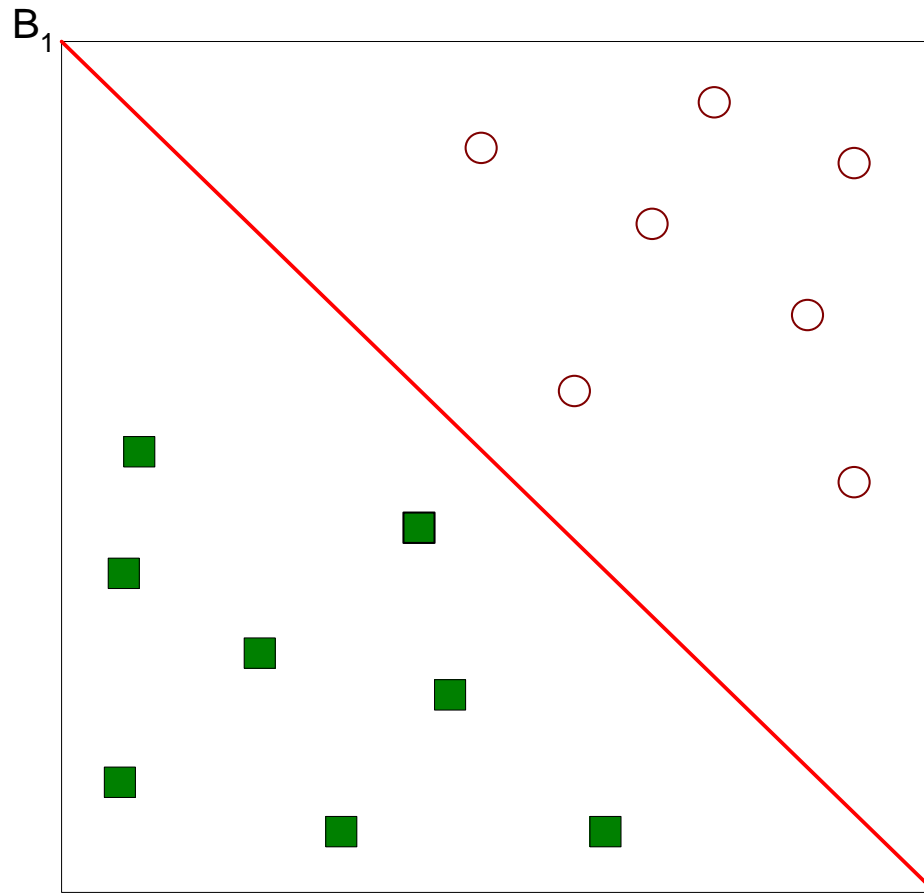


- **Margin** for a hyperplane that separates samples of two linearly separable classes is:
 - The smallest distance between the decision boundary and any of the training samples

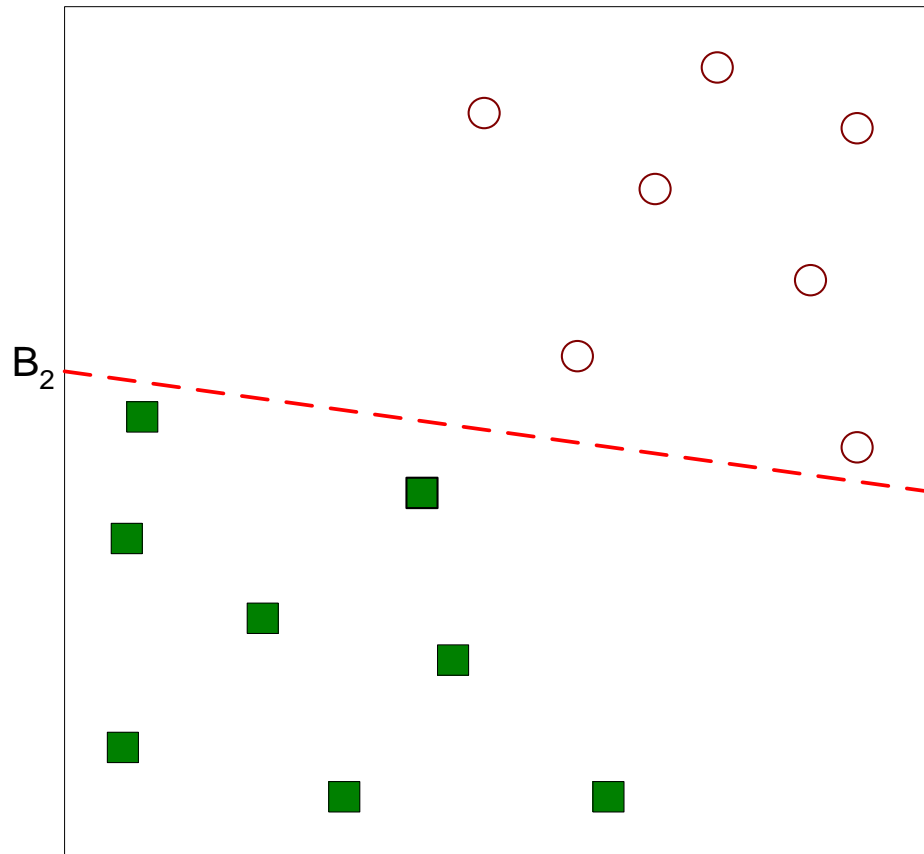
Maximum Margin Hyperplanes



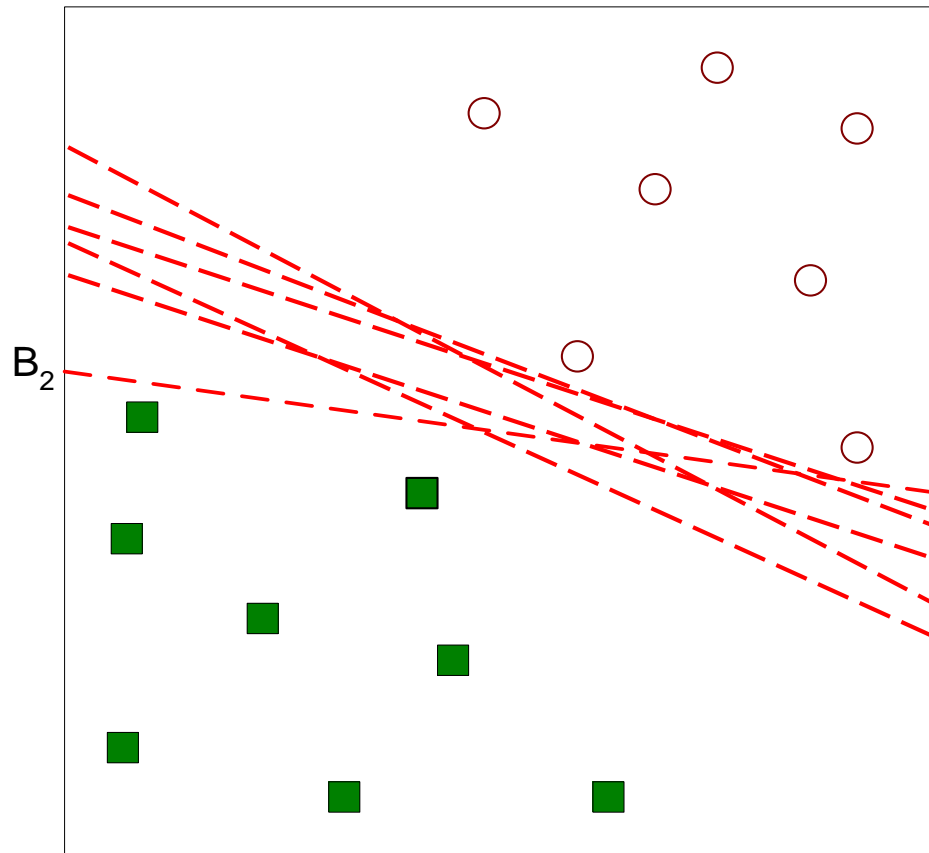
Maximum Margin Hyperplanes



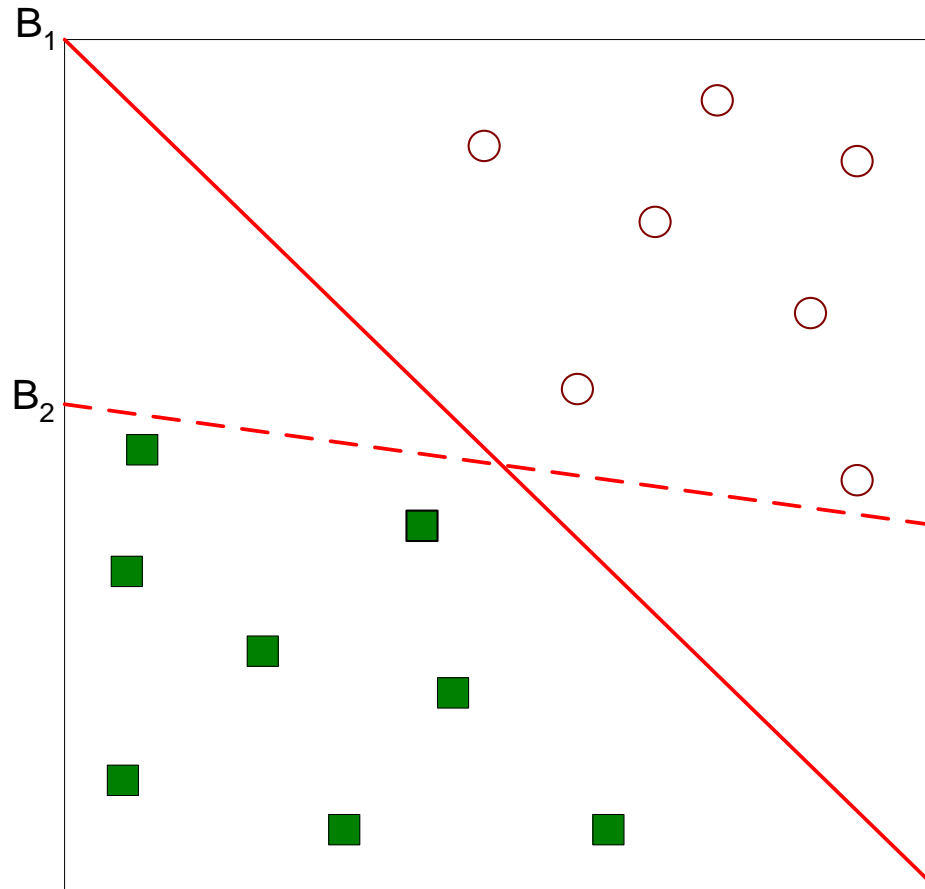
Maximum Margin Hyperplanes



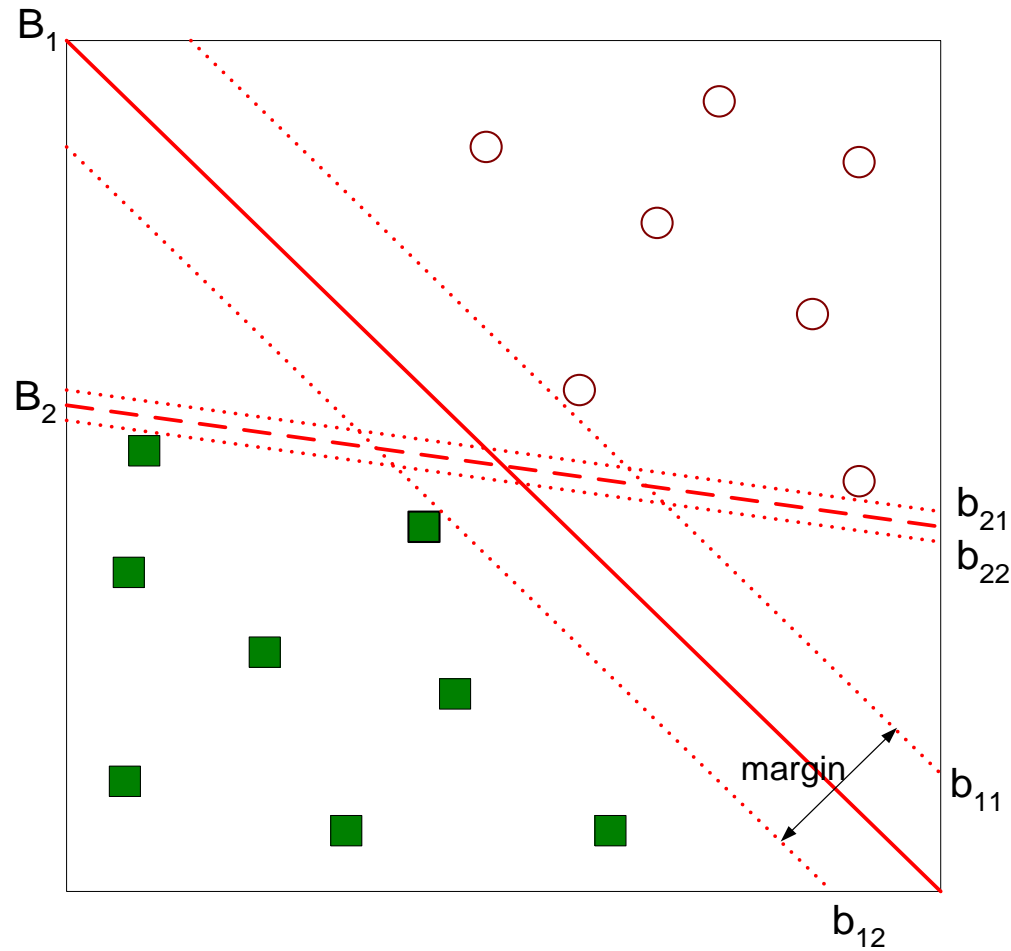
Maximum Margin Hyperplanes



Maximum Margin Hyperplanes



Maximum Margin Hyperplanes

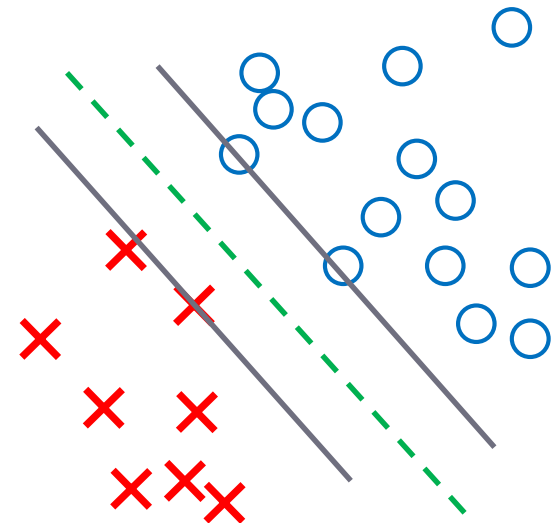
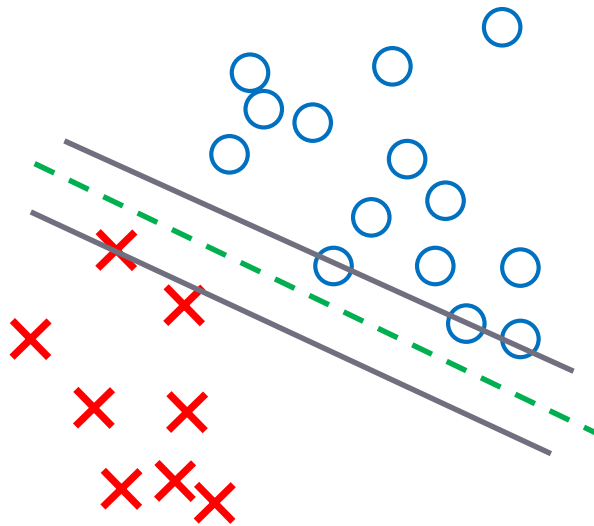


maximizes

SVM

Maximum margin

- SVM finds the solution with maximum margin
 - Solution: a hyperplane that is farthest from all training samples



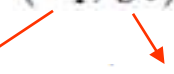
Larger margin

- The hyperplane with the largest margin has equal distances to the nearest sample of both classes

Linear SVM: Separable Case

A linear SVM is a classifier that searches for a hyperplane with the largest margin

$$(x_i, y_i) \quad (i = 1, 2, \dots, N)$$



$$(x_{i1}, x_{i2}, \dots, x_{id})^T \quad y_i \in \{-1, 1\}$$

decision boundary of a linear classifier

$$\mathbf{w} \cdot \mathbf{x} + b = 0,$$

\mathbf{w} and b are parameters of the model

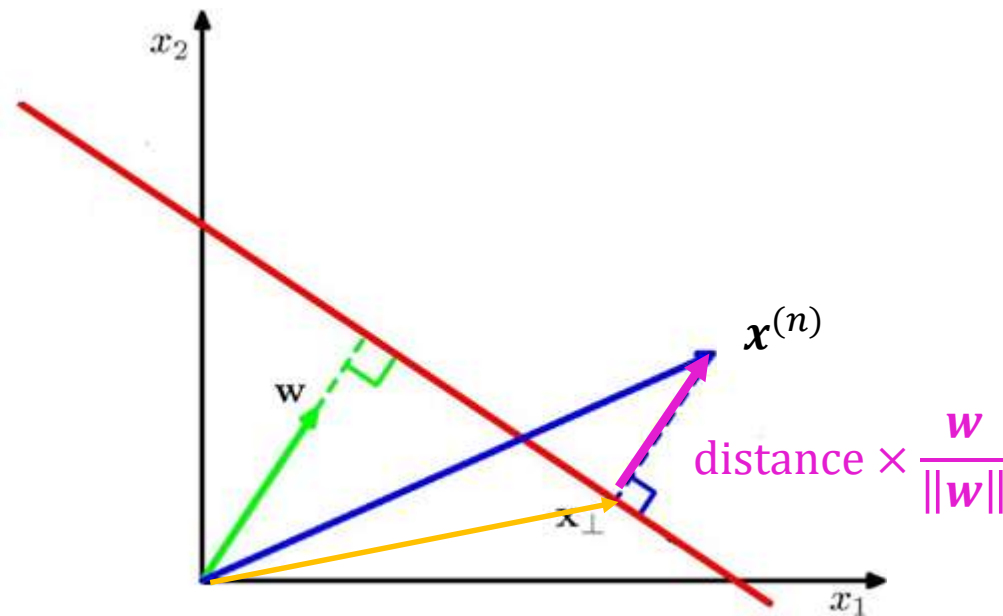
$$\mathbf{w} \cdot \mathbf{x}_s + b = k, \quad k > 0.$$

$$\mathbf{w} \cdot \mathbf{x}_c + b = k', \quad k' < 0.$$

$$y = \begin{cases} 1, & \text{if } \mathbf{w} \cdot \mathbf{z} + b > 0; \\ -1, & \text{if } \mathbf{w} \cdot \mathbf{z} + b < 0. \end{cases}$$

Distance between an $\mathbf{x}^{(n)}$ and the plane

$$\text{distance} = \frac{|\mathbf{w}^T \mathbf{x}^{(n)} + w_0|}{\|\mathbf{w}\|}$$



Hard-margin SVM: Optimization problem

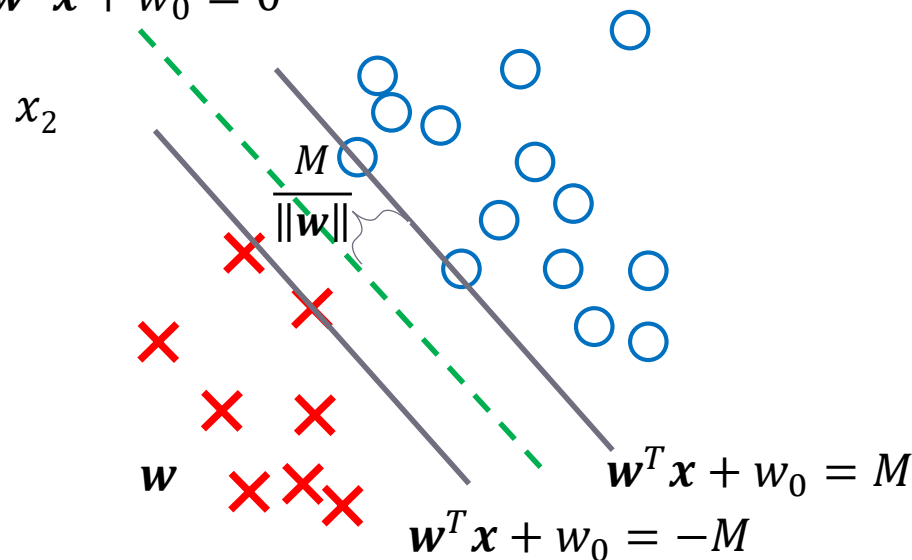
$$\max_{M, \mathbf{w}, w_0} \frac{2M}{\|\mathbf{w}\|}$$

$$\text{s. t. } (\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq M \quad \forall \mathbf{x}^{(i)} \in C_1 \longrightarrow y^{(i)} = 1$$

$$(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \leq -M \quad \forall \mathbf{x}^{(i)} \in C_2 \longrightarrow y^{(i)} = -1$$

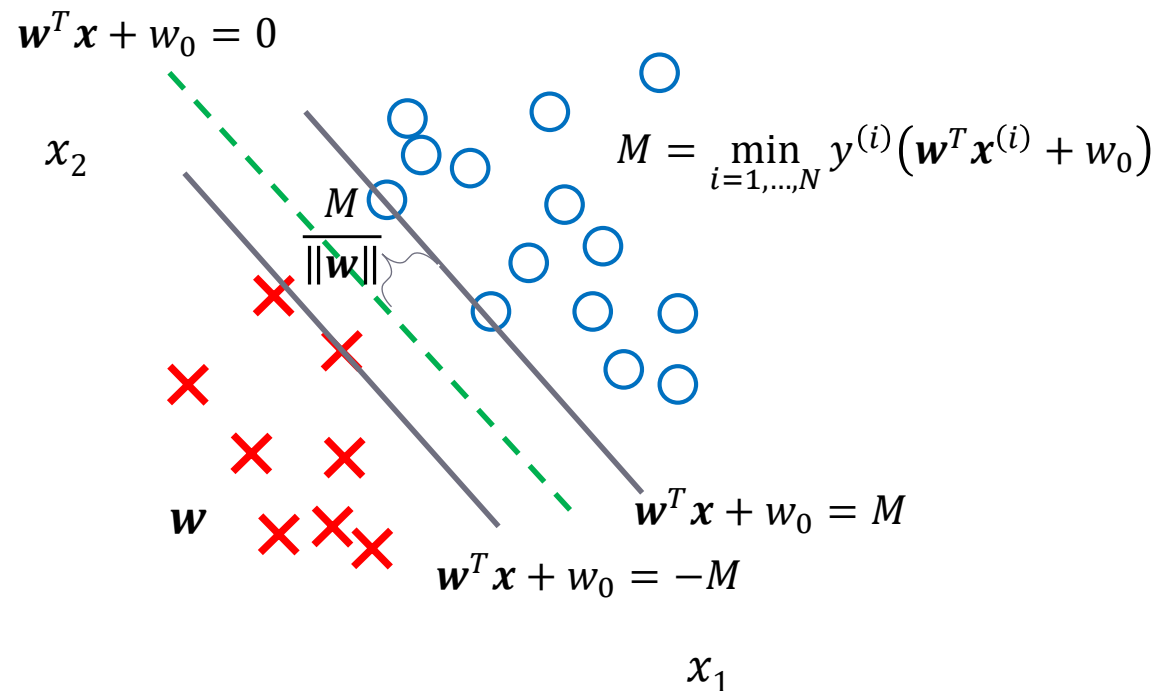
$$\mathbf{w}^T \mathbf{x} + w_0 = 0$$

Margin: $2 \frac{M}{\|\mathbf{w}\|}$

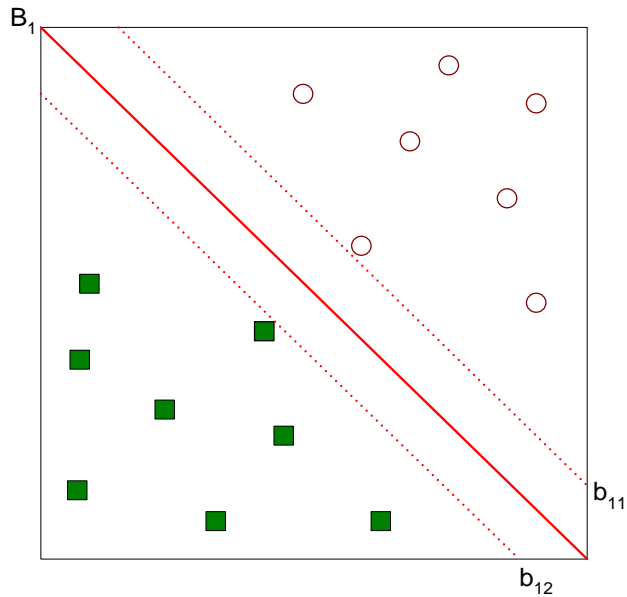


Hard-margin SVM: Optimization problem

$$\begin{aligned} & \max_{M, \mathbf{w}, w_0} \frac{2M}{\|\mathbf{w}\|} \\ \text{s. t. } & y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq M \quad i = 1, \dots, N \end{aligned}$$



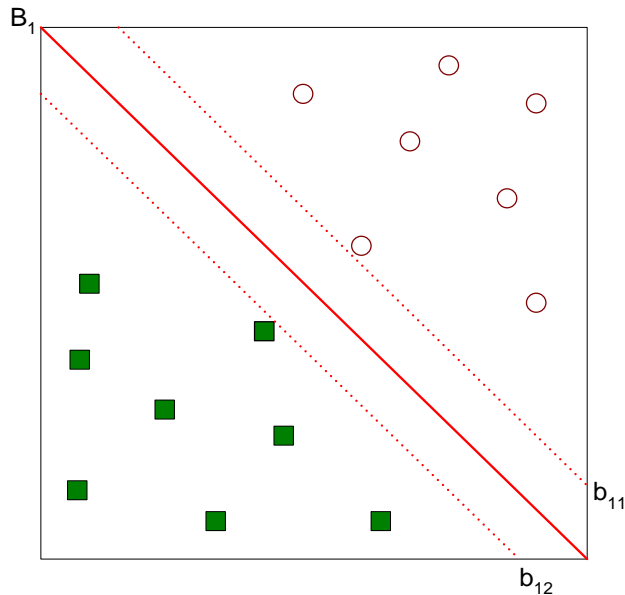
Linear SVM: Separable Case



$$b_{i1} : w \cdot x + b = 1,$$

$$b_{i2} : w \cdot x + b = -1.$$

Linear SVM: Separable Case



$$b_{i1} : \mathbf{w} \cdot \mathbf{x} + b = 1,$$

$$b_{i2} : \mathbf{w} \cdot \mathbf{x} + b = -1.$$

margin of the decision boundary
is given by the distance between
these two hyperplanes

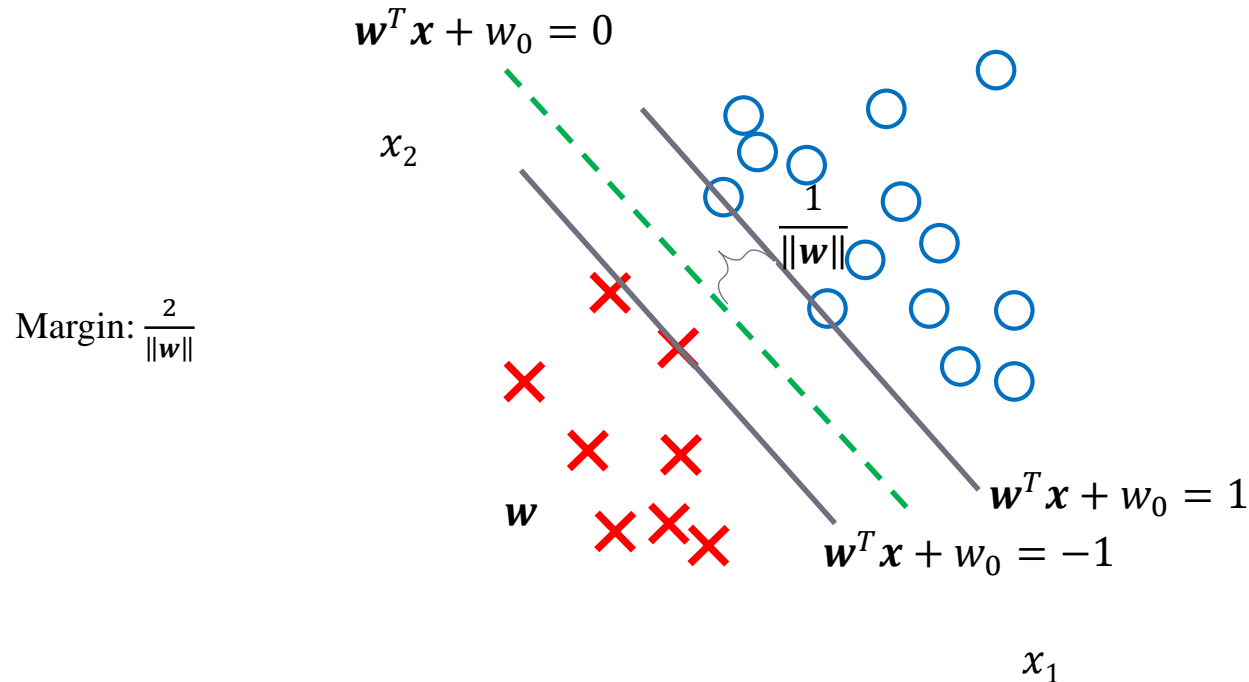
$$d = \frac{2}{\|\mathbf{w}\|}$$

Hard-margin SVM: Optimization problem

We can set $\mathbf{w}' = \frac{\mathbf{w}}{M}$, $w'_0 = \frac{w_0}{M}$:

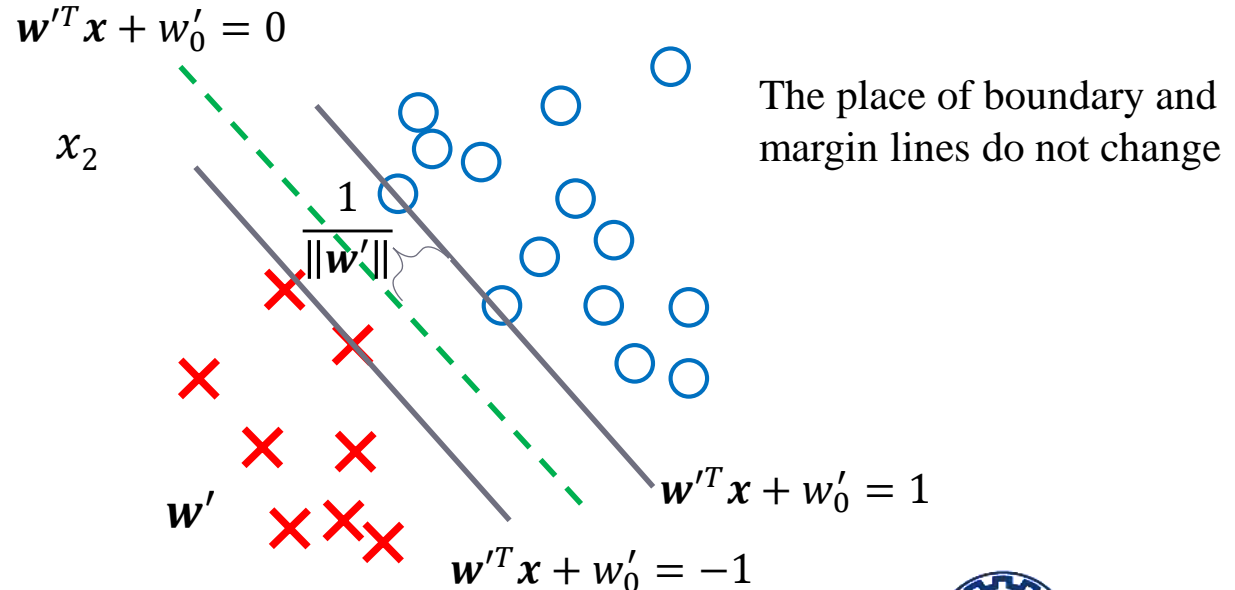
$$\max_{\mathbf{w}, w_0} \frac{2}{\|\mathbf{w}\|}$$

s. t. $(\mathbf{w}^T \mathbf{x}^{(n)} + w_0) \geq 1 \quad \forall y^{(n)} = 1$
 $(\mathbf{w}^T \mathbf{x}^{(n)} + w_0) \leq -1 \quad \forall y^{(n)} = -1$



Hard-margin SVM: Optimization problem

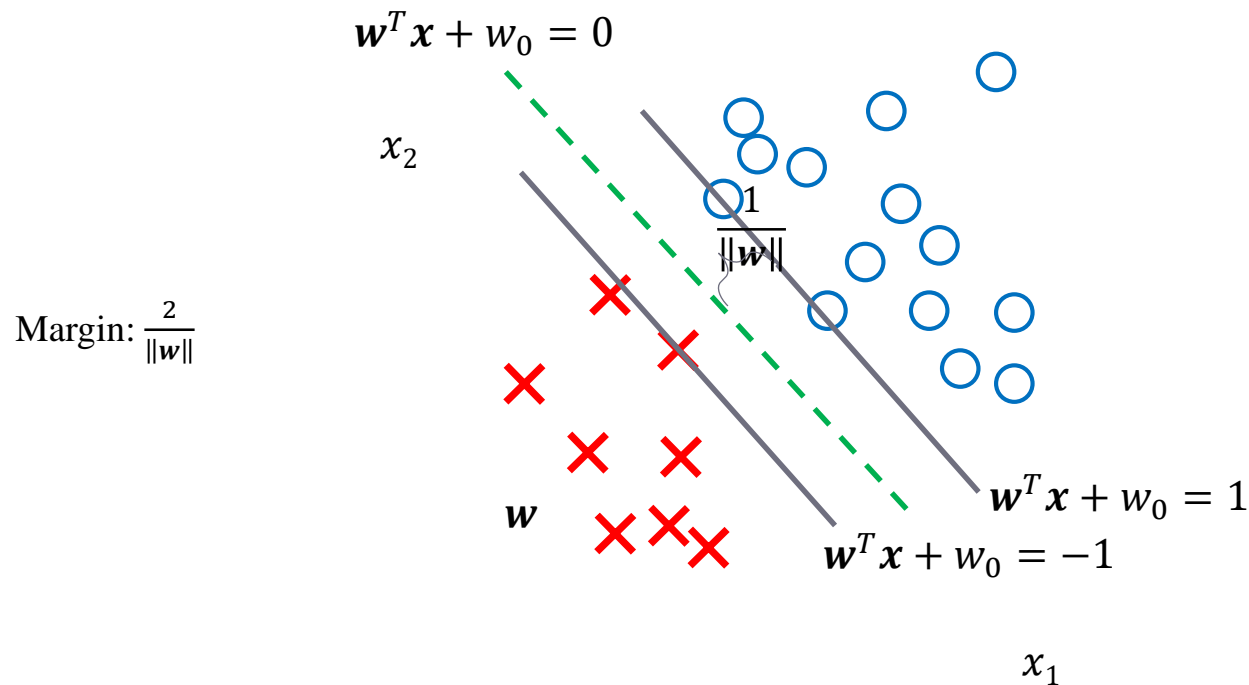
$$\begin{aligned} & \max_{w', w'_0} \frac{2}{\|w'\|} \\ \text{s.t. } & y^{(i)}(w'^T x^{(i)} + w'_0) \geq 1 \quad i = 1, \dots, N \end{aligned}$$



SVM

Hard-margin SVM: Optimization problem

$$\begin{aligned} & \max_{\mathbf{w}, w_0} \frac{2}{\|\mathbf{w}\|} \\ \text{s.t. } & y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1, n = 1, \dots, N \end{aligned}$$



SVM

$$\begin{aligned} \mathbf{w} \cdot \mathbf{x}_i + b &\geq 1 \text{ if } y_i = 1, & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) &\geq 1, \quad i = 1, 2, \dots, N. \\ \mathbf{w} \cdot \mathbf{x}_i + b &\leq -1 \text{ if } y_i = -1 \end{aligned}$$

Definition 5.1 (Linear SVM: Separable Case). The learning task in SVM can be formalized as the following constrained optimization problem:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} \quad & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, N. \end{aligned}$$

For Solving constrained optimization problem (like SVM Optimization) there exist Numerical approaches like **Quadratic Programming (QP)** !

Quadratic Programming (QP)

It is a convex Quadratic Programming (QP) problem

There are computationally efficient packages to solve it.

It has a global minimum (if any).

$$\begin{array}{ll}\min_x & \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x} \\ \text{s. t.} & \mathbf{A} \mathbf{x} \leq \mathbf{b} \\ & \mathbf{E} \mathbf{x} = \mathbf{d}\end{array}$$

Dual formulation of the SVM

- We are going to introduce the *dual* SVM problem which is equivalent to the original *primal* problem. The dual problem:
 - is often **easier**
 - It's **computationally more feasible** in high dimensional spaces where d is large
 - gives us further insights into the **optimal hyper-plane**
 - enable us to exploit the **kernel trick**

Optimization: Lagrangian multipliers

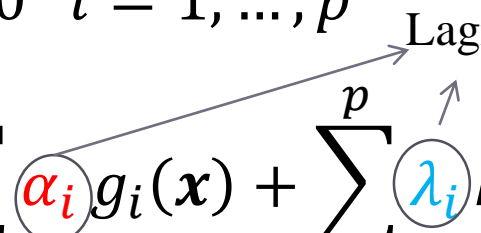
$$\begin{aligned} p^* &= \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{s.t. } g_i(\mathbf{x}) &\leq 0 \quad i = 1, \dots, m \\ h_i(\mathbf{x}) &= 0 \quad i = 1, \dots, p \end{aligned}$$

Optimization: Lagrangian multipliers

$$\begin{aligned} p^* &= \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{s. t. } g_i(\mathbf{x}) &\leq 0 \quad i = 1, \dots, m \\ h_i(\mathbf{x}) &= 0 \quad i = 1, \dots, p \end{aligned}$$

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{i=1}^m \alpha_i g_i(\mathbf{x}) + \sum_{i=1}^p \lambda_i h_i(\mathbf{x})$$

Lagrangian multipliers

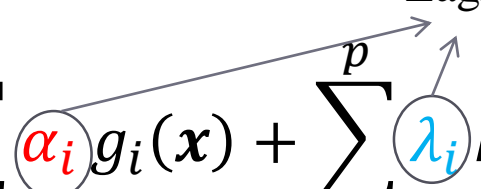


Optimization: Lagrangian multipliers

$$\begin{aligned} p^* &= \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{s. t. } g_i(\mathbf{x}) &\leq 0 \quad i = 1, \dots, m \\ h_i(\mathbf{x}) &= 0 \quad i = 1, \dots, p \end{aligned}$$

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{i=1}^m \alpha_i g_i(\mathbf{x}) + \sum_{i=1}^p \lambda_i h_i(\mathbf{x})$$

Lagrangian multipliers



$$\max_{\{\alpha_i \geq 0\}, \{\lambda_i\}} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\lambda}) = \begin{cases} \infty & \text{any } g_i(\mathbf{x}) > 0 \\ \infty & \text{any } h_i(\mathbf{x}) \neq 0 \\ f(\mathbf{x}) & \text{otherwise} \end{cases}$$

$$p^* = \min_{\mathbf{x}} \max_{\{\alpha_i \geq 0\}, \{\lambda_i\}} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\lambda})$$

Optimization: Dual problem

- In general, we have:

$$\max_x \min_y h(x, y) \leq \min_y \max_x h(x, y)$$

- **Primal problem:** $p^* = \min_x \max_{\{\alpha_i \geq 0\}, \{\lambda_i\}} \mathcal{L}(x, \alpha, \lambda)$

- **Dual problem:** $d^* = \max_{\{\alpha_i \geq 0\}, \{\lambda_i\}} \min_x \mathcal{L}(x, \alpha, \lambda)$

- Obtained by swapping the order of min and max
- $d^* \leq p^*$

- When the original problem is convex (f and g are convex functions and h is affine), we have strong duality $d^* = p^*$

Hard-margin SVM: Dual problem

$$\begin{aligned} & \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s. t. } & y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 \quad i = 1, \dots, N \end{aligned}$$

- By incorporating the constraints through Lagrangian multipliers, we will have:

$$\min_{\mathbf{w}, w_0} \max_{\{\alpha_n \geq 0\}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N \alpha_n (1 - y^{(n)} (\mathbf{w}^T \mathbf{x}^{(n)} + w_0)) \right\}$$

Hard-margin SVM: Dual problem

$$\begin{aligned} & \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s. t. } & y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 \quad i = 1, \dots, N \end{aligned}$$

- By incorporating the constraints through Lagrangian multipliers, we will have:

$$\min_{\mathbf{w}, w_0} \max_{\{\alpha_n \geq 0\}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N \alpha_n (1 - y^{(n)} (\mathbf{w}^T \mathbf{x}^{(n)} + w_0)) \right\}$$

- Dual problem (changing the order of min and max in the above problem):

$$\max_{\{\alpha_n \geq 0\}} \min_{\mathbf{w}, w_0} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N \alpha_n (1 - y^{(n)} (\mathbf{w}^T \mathbf{x}^{(n)} + w_0)) \right\}$$

Hard-margin SVM: Dual problem

$$\max_{\{\alpha_n \geq 0\}} \min_{\mathbf{w}, w_0} \mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha})$$

$$\mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N (1 - y^{(n)} (\mathbf{w}^T \mathbf{x}^{(n)} + w_0))$$

Hard-margin SVM: Dual problem

$$\max_{\{\alpha_n \geq 0\}} \min_{\mathbf{w}, w_0} \mathcal{L}(\mathbf{w}, w_0, \alpha)$$

$$\mathcal{L}(\mathbf{w}, w_0, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N (1 - y^{(n)} (\mathbf{w}^T \mathbf{x}^{(n)} + w_0))$$

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, w_0, \alpha) = 0 &\Rightarrow \mathbf{w} \sum_{n=1}^N \alpha_n y^{(n)} \mathbf{x}^{(n)} = \mathbf{0} \\ &\Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n y^{(n)} \mathbf{x}^{(n)} \end{aligned}$$

Hard-margin SVM: Dual problem

$$\max_{\{\alpha_n \geq 0\}} \min_{\mathbf{w}, w_0} \mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha})$$

$$\mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N (1 - y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + w_0))$$

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha}) = 0 &\Rightarrow \mathbf{w} \sum_{n=1}^N \alpha_n y^{(n)} \mathbf{x}^{(n)} = \mathbf{0} \\ &\Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n y^{(n)} \mathbf{x}^{(n)} \end{aligned}$$

$$\frac{\partial \mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha})}{\partial w_0} = 0 \Rightarrow -\underbrace{\sum_{n=1}^N \alpha_n y^{(n)}}_{\downarrow} = 0$$

w_0 do not appear, instead, a “global” constraint on $\boldsymbol{\alpha}$ is created.

Substituting

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y^{(n)} \mathbf{x}^{(n)} \quad \sum_{n=1}^N \alpha_n y^{(n)} = 0$$

In the Lagrangian

$$\mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n (1 - y^{(n)} (\mathbf{w}^T \mathbf{x}^{(n)} + w_0))$$

Substituting

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y^{(n)} \mathbf{x}^{(n)} \quad \sum_{n=1}^N \alpha_n y^{(n)} = 0$$

In the Lagrangian

$$\mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n (1 - y^{(n)} (\mathbf{w}^T \mathbf{x}^{(n)} + w_0))$$

We get

$$\mathcal{L}(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y^{(n)} y^{(m)} \mathbf{x}^{(n)T} \mathbf{x}^{(m)}$$

Maximize w.r.t. $\boldsymbol{\alpha}$ subject to $\alpha_n \geq 0$ for $n = 1, \dots, N$ and $\sum_{n=1}^N \alpha_n y^{(n)} = 0$

Hard-margin SVM: Dual problem

$$\max_{\alpha} \left\{ \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y^{(n)} y^{(m)} \mathbf{x}^{(n)T} \mathbf{x}^{(m)} \right\}$$

$$\text{Subject to } \sum_{n=1}^N \alpha_n y^{(n)} = 0$$

$$\alpha_n \geq 0 \quad n = 1, \dots, N$$

- It is a convex QP

Solution

- Quadratic programming:

$$\min_{\alpha} \frac{1}{2} \alpha^T \begin{bmatrix} y^{(1)}y^{(1)}\mathbf{x}^{(1)T}\mathbf{x}^{(1)} & \dots & y^{(1)}y^{(N)}\mathbf{x}^{(1)T}\mathbf{x}^{(N)} \\ \vdots & \ddots & \vdots \\ y^{(N)}y^{(1)}\mathbf{x}^{(N)T}\mathbf{x}^{(1)} & \dots & y^{(N)}y^{(N)}\mathbf{x}^{(N)T}\mathbf{x}^{(N)} \end{bmatrix} \alpha + (-\mathbf{1})^T \alpha$$

$$\text{s. t. } -\alpha \leq \mathbf{0} \\ \mathbf{y}^T \alpha = \mathbf{0}$$

Finding the hyperplane

- After finding α by QP, we find \mathbf{w} :

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y^{(n)} \mathbf{x}^{(n)}$$

- How to find \mathbf{w}_0 ?
 - we discuss it after introducing support vectors

Optimal Point

- Necessary conditions for the solution $[\mathbf{w}^*, w_0^*, \boldsymbol{\alpha}^*]$:
 - $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha})|_{\mathbf{w}^*, w_0^*, \boldsymbol{\alpha}^*} = 0$
 - $\frac{\partial \mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha})}{\partial w_0} |_{\mathbf{w}^*, w_0^*, \boldsymbol{\alpha}^*} = 0$
 - $\alpha_n^* \geq 0 \quad n = 1, \dots, N$

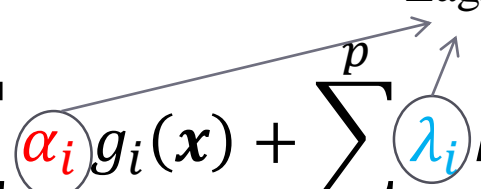
- $y^{(n)}(\mathbf{w}^{*T} \mathbf{x}^{(n)} + w_0^*) \geq 1 \quad n = 1, \dots, N$
- $\alpha_i^* \left(1 - y^{(n)}(\mathbf{w}^{*T} \mathbf{x}^{(n)} + w_0^*)\right) = 0 \quad n = 1, \dots, N$

Optimization: Lagrangian multipliers

$$\begin{aligned} p^* &= \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{s. t. } g_i(\mathbf{x}) &\leq 0 \quad i = 1, \dots, m \\ h_i(\mathbf{x}) &= 0 \quad i = 1, \dots, p \end{aligned}$$

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{i=1}^m \alpha_i g_i(\mathbf{x}) + \sum_{i=1}^p \lambda_i h_i(\mathbf{x})$$

Lagrangian multipliers



$$\max_{\{\alpha_i \geq 0\}, \{\lambda_i\}} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\lambda}) = \begin{cases} \infty & \text{any } g_i(\mathbf{x}) > 0 \\ \infty & \text{any } h_i(\mathbf{x}) \neq 0 \\ f(\mathbf{x}) & \text{otherwise} \end{cases}$$

$$p^* = \min_{\mathbf{x}} \max_{\{\alpha_i \geq 0\}, \{\lambda_i\}} \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\lambda})$$

Optimal Point

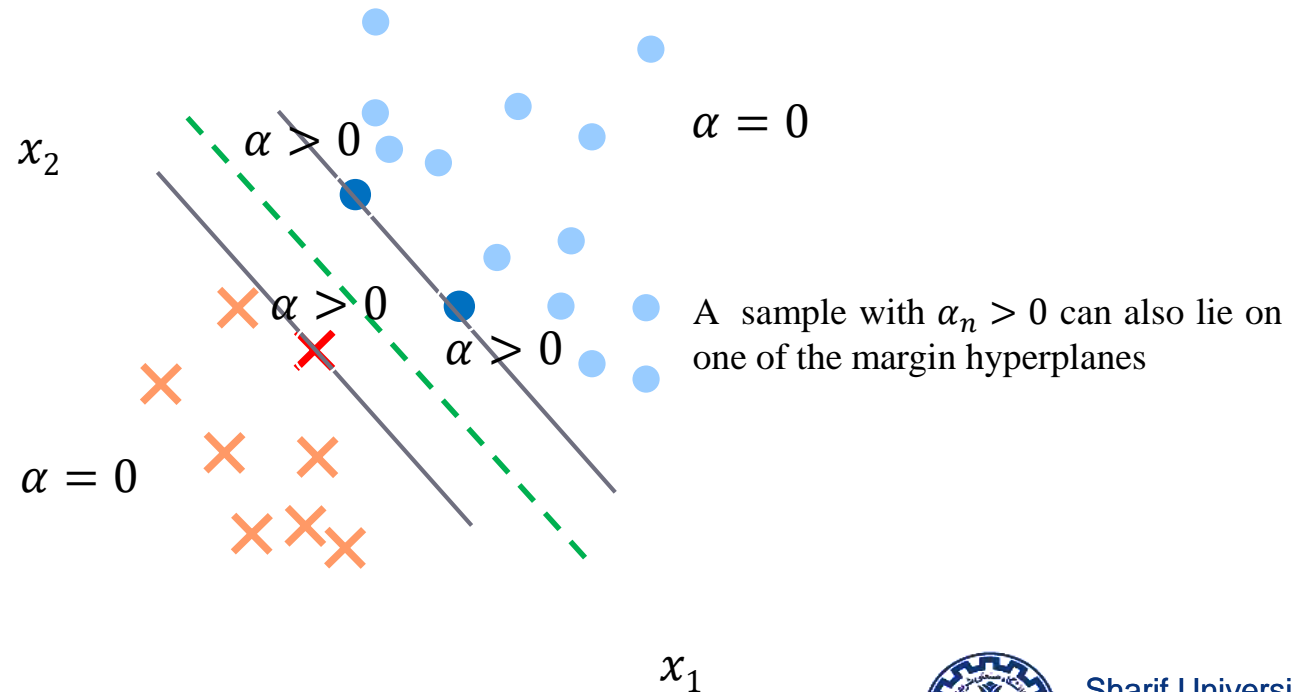
- Necessary conditions for the solution $[\mathbf{w}^*, w_0^*, \boldsymbol{\alpha}^*]$:
 - $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha})|_{\mathbf{w}^*, w_0^*, \boldsymbol{\alpha}^*} = 0$
 - $\frac{\partial \mathcal{L}(\mathbf{w}, w_0, \boldsymbol{\alpha})}{\partial w_0} |_{\mathbf{w}^*, w_0^*, \boldsymbol{\alpha}^*} = 0$
 - $\alpha_n^* \geq 0 \quad n = 1, \dots, N$

- $y^{(n)}(\mathbf{w}^{*T} \mathbf{x}^{(n)} + w_0^*) \geq 1 \quad n = 1, \dots, N$
- $\alpha_i^* \left(1 - y^{(n)}(\mathbf{w}^{*T} \mathbf{x}^{(n)} + w_0^*)\right) = 0 \quad n = 1, \dots, N$

Karush-Kuhn-Tucker
(KKT) conditions

Hard-margin SVM: Support vectors

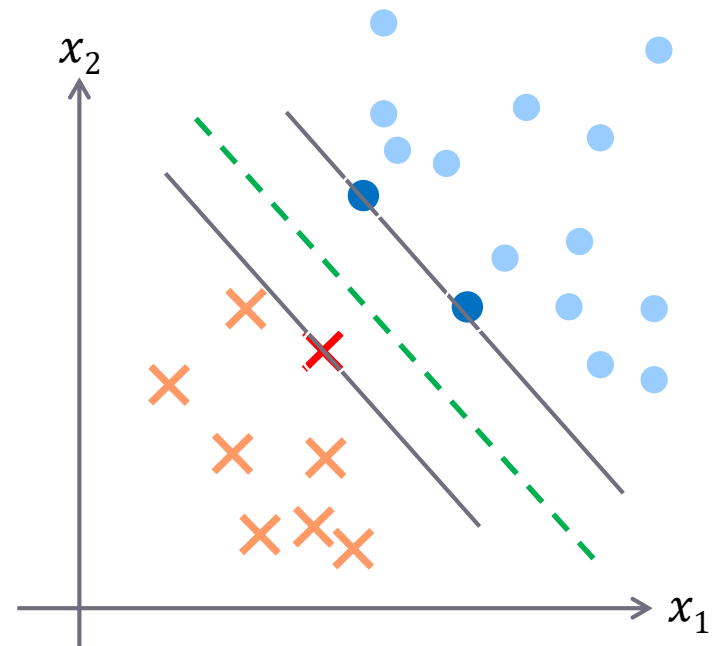
- **Inactive** constraint: $y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + w_0) > 1$
 - $\Rightarrow \alpha_n = 0$ and thus $\mathbf{x}^{(n)}$ is not a support vector.
- **Active** constraint: $y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + w_0) = 1$
 - $\Rightarrow \alpha_n$ can be greater than 0 and thus $\mathbf{x}^{(i)}$ can be a support vector.



Hard-margin SVM: Support vectors

- Support Vectors (SVs) = $\{\mathbf{x}^{(n)} \mid \alpha_n > 0\}$
- The **direction** of hyper-plane can be found only based on support vectors:

$$\mathbf{w} = \sum_{\alpha_n > 0} \alpha_n y^{(n)} \mathbf{x}^{(n)}$$



Finding the hyperplane

- After finding α by QP, we find \mathbf{w} :

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y^{(n)} \mathbf{x}^{(n)}$$

- How to find w_0 ?
 - Each of the samples that has α_s is on the margin, thus we solve for w_0 using any of SVs:

$$y^{(s)} (\mathbf{w}^T \mathbf{x}^{(s)} + w_0) = 1$$

$$\Rightarrow w_0 = y^{(s)} - \mathbf{w}^T \mathbf{x}^{(s)}$$

Hard-margin SVM: Dual problem

Classifying new samples using only SVs

- Classification of a new sample \mathbf{x} :

$$\hat{y} = \text{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$
$$\hat{y} = \text{sign} \left(w_0 + \left(\sum_{\alpha_n > 0} \alpha_n y^{(n)} \mathbf{x}^{(n)} \right)^T \mathbf{x} \right)$$

$$\hat{y} = \text{sign} \left(y^{(s)} - \underbrace{\sum_{\alpha_n > 0} \alpha_n y^{(n)} \mathbf{x}^{(n)T} \mathbf{x}^{(s)}}_{w_0} + \sum_{\alpha_n > 0} \alpha_n y^{(n)} \mathbf{x}^{(n)T} \mathbf{x} \right)$$

Support vectors are sufficient to predict labels of new samples

- The classifier is based on the expansion in terms of dot products of \mathbf{x} with support vectors.

Hard-margin SVM dual problem: An important property

$$\max_{\alpha} \left\{ \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y^{(n)} y^{(m)} \mathbf{x}^{(n)T} \mathbf{x}^{(m)} \right\}$$

$$\text{Subject to } \sum_{n=1}^N \alpha_n y^{(n)} = 0$$

$$\alpha_n \geq 0 \quad n = 1, \dots, N$$

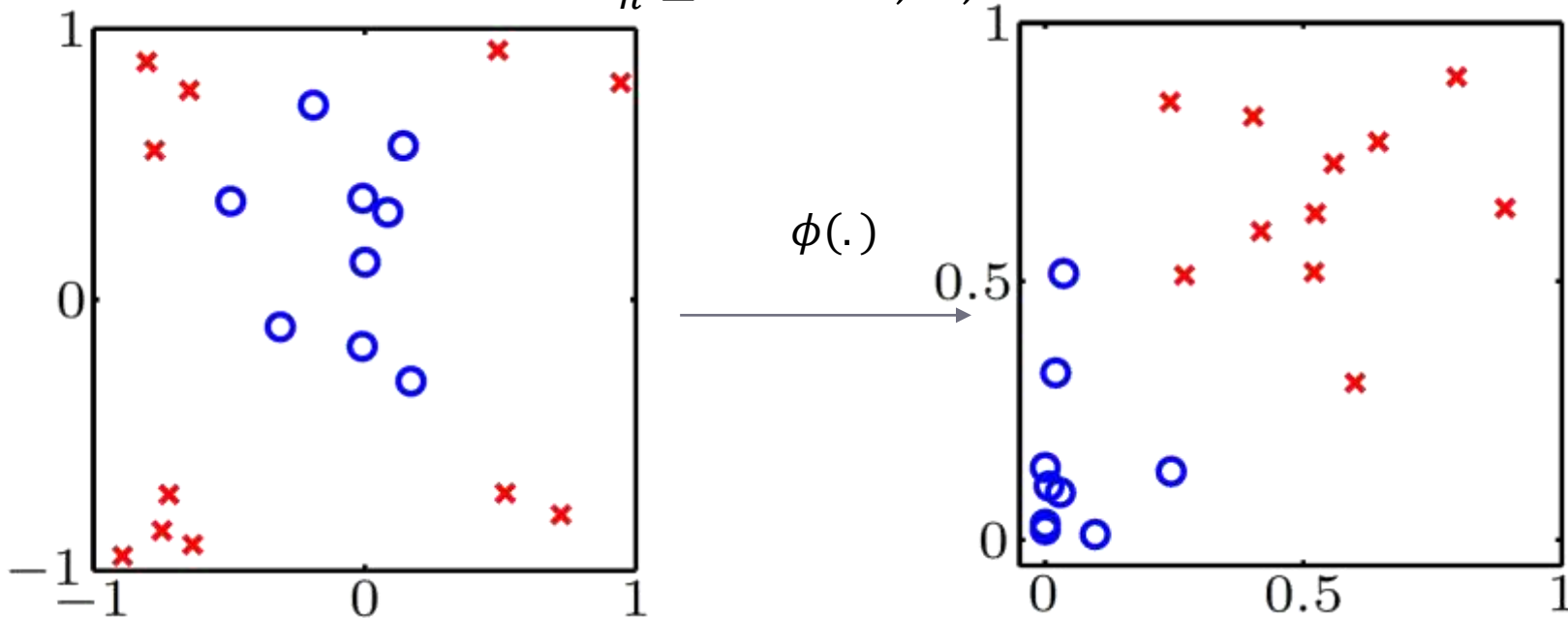
- Only the dot product of each pair of training data appears in the optimization problem
 - An important property that is helpful to extend to non-linear SVM
 - We will talk about it later (kernel-based methods)

In the transformed space

$$\max_{\alpha} \left\{ \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y^{(n)} y^{(m)} \phi(x^{(n)})^T \phi(x^{(m)}) \right\}$$

$$\text{Subject to } \sum_{n=1}^N \alpha_n y^{(n)} = 0$$

$$\alpha_n \geq 0 \quad n = 1, \dots, N$$

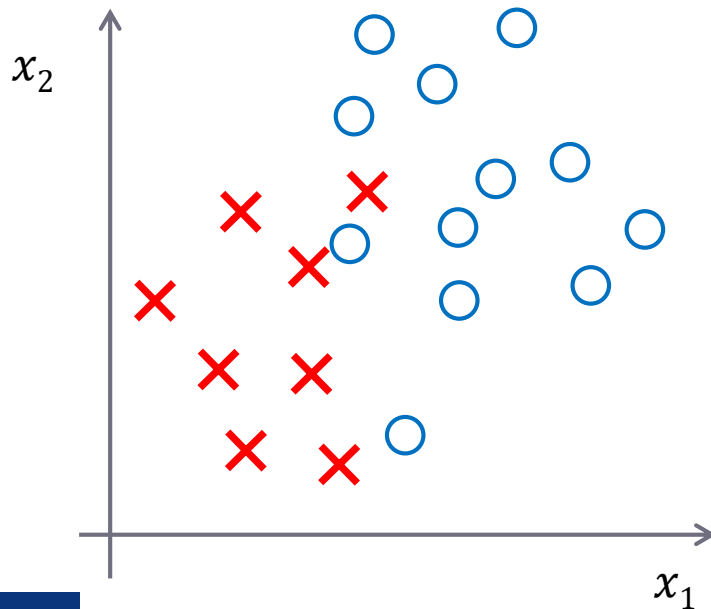


Beyond linear separability

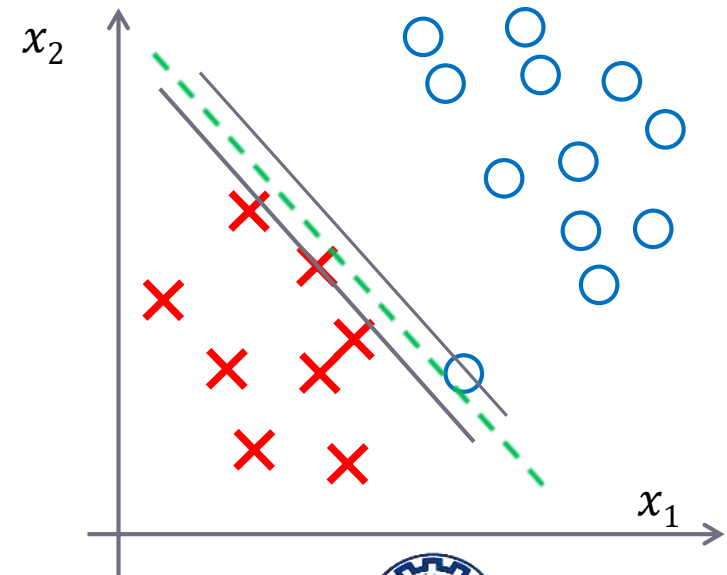
- When training samples are not linearly separable, it has no solution.
- How to extend it to find a solution even though the classes are not exactly linearly separable.

Near linear separability

- How to extend the hard-margin SVM to allow classification error
 - Overlapping classes that can be approximately separated by a linear boundary
 - Noise in the linearly separable classes



SVM



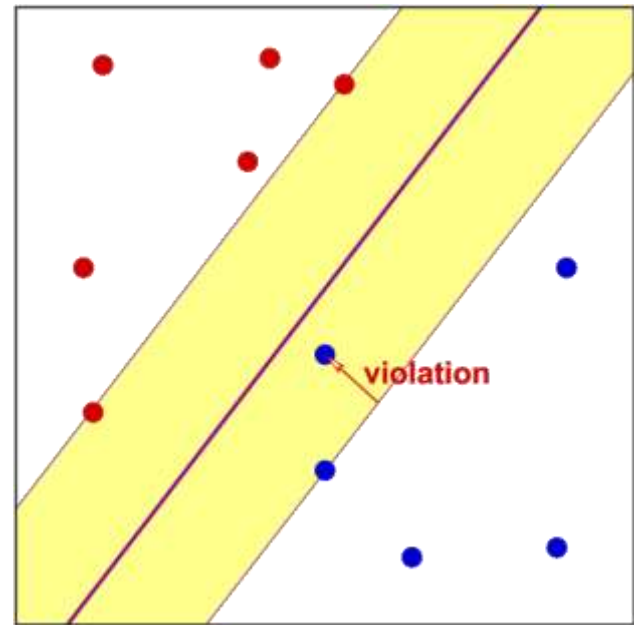
Sharif University
of Technology

Near linear separability: Soft-margin SVM

- Minimizing the number of misclassified points?!
 - NP-complete
- Soft margin:
 - Maximizing a margin while trying to minimize the *distance* between misclassified points and their correct margin plane

Error measure

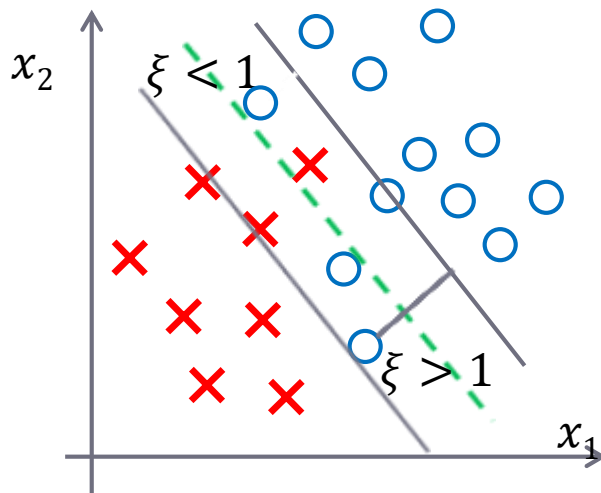
- Margin violation amount ξ_n ($\xi_n \geq 0$):
 - $y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + w_0) \geq 1 - \xi_n$
- Total violation: $\sum_{n=1}^N \xi_n$



Soft-margin SVM: Optimization problem

- SVM with slack variables: allows samples to fall within the margin, but penalizes them

$$\begin{aligned} \min_{\mathbf{w}, w_0, \{\xi_n\}_{n=1}^N} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \\ \text{s. t.} \quad & y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + w_0) \geq 1 - \xi_n \quad n = 1, \dots, N \\ & \xi_n \geq 0 \end{aligned}$$



ξ_n : **slack** variables

$0 < \xi_n < 1$: if $\mathbf{x}^{(n)}$ is correctly classified but inside margin

$\xi_n > 1$: if $\mathbf{x}^{(n)}$ is misclassified

Soft-margin SVM

- linear penalty (hinge loss) for a sample if it is misclassified or lied in the margin
 - tries to maintain ξ_n small while maximizing the margin.
 - always finds a solution (as opposed to hard-margin SVM)
 - more robust to the outliers
- Soft margin problem is still a convex QP

Soft-margin SVM: Parameter C

- C is a tradeoff parameter:
 - small C allows margin constraints to be easily ignored
 - large margin
 - large C makes constraints hard to ignore
 - narrow margin
- $C \rightarrow \infty$ enforces all constraints: hard margin
- C can be determined using a technique like cross-validation

Soft-margin SVM: Cost function

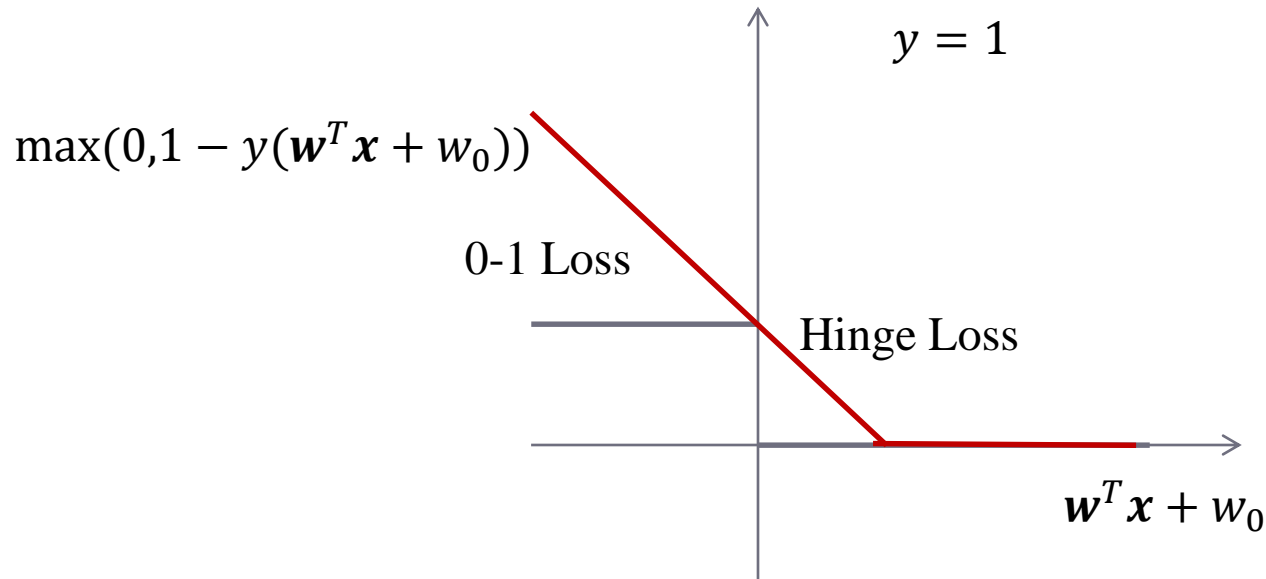
$$\begin{aligned} & \min_{\mathbf{w}, w_0, \{\xi_n\}_{n=1}^N} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \\ \text{s. t. } & y^{(n)} (\mathbf{w}^T \mathbf{x}^{(n)} + w_0) \geq 1 - \xi_n \quad n = 1, \dots, N \\ & \xi_n \geq 0 \end{aligned}$$

- It is equivalent to the unconstrained optimization problem:

$$\min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \max(0, 1 - y^{(n)} (\mathbf{w}^T \mathbf{x}^{(n)} + w_0))$$

SVM loss function

- Hinge loss vs. 0-1 loss



Lagrange formulation

$$\begin{aligned}\mathcal{L}(\mathbf{w}, w_0, \xi, \alpha, \beta) \\&= \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n (1 - \xi_n - y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + w_0)) \\&\quad - \sum_{n=1}^N \beta_n \xi_n\end{aligned}$$

- Minimize w.r.t. \mathbf{w}, w_0, ξ and maximize w.r.t. $\alpha_n \geq 0$ and $\beta_n \geq 0$

$$\begin{aligned}\min_{\mathbf{w}, w_0, \{\xi_n\}_{n=1}^N} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \\ \text{s. t.} \quad & y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + w_0) \geq 1 - \xi_n \quad n = 1, \dots, N \\ & \xi_n \geq 0\end{aligned}$$

Lagrange formulation

$$\mathcal{L}(\mathbf{w}, w_0, \xi, \alpha, \beta) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n (1 - \xi_n)$$

Soft-margin SVM: Dual problem

$$\max_{\alpha} \left\{ \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y^{(n)} y^{(m)} \mathbf{x}^{(n)T} \mathbf{x}^{(m)} \right\}$$

$$\text{Subject to } \sum_{n=1}^N \alpha_n y^{(n)} = 0$$

$$0 \leq \alpha_n \leq C \quad n = 1, \dots, N$$

- After solving the above quadratic problem, \mathbf{w} is found as:

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y^{(n)} \mathbf{x}^{(n)}$$

Soft-margin SVM: Support vectors

- Support Vectors: $\alpha_n > 0$

- If $0 < \alpha_n < C$ (**margin** support vector) SVs on the margin

$$y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + w_0) = 1 \quad (\xi_n = 0)$$

- If $\alpha = C$ (**non-margin** support vector) SVs on or over the margin

$$y^{(n)}(\mathbf{w}^T \mathbf{x}^{(n)} + w_0) < 1 \quad (\xi_n > 0)$$

$$C - \alpha_n - \beta_n = 0$$

SVM: Summary

- Hard margin: maximizing margin
- Soft margin: handling noisy data and overlapping classes
 - Slack variables in the problem
- Dual problems of hard-margin and soft-margin SVM
 - Classifier decision in terms of *support vectors*
- Dual problems lead us to non-linear SVM method easily by kernel substitution

Recourses

- C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 7.1.
- Yaser S.Abu-Mostafa, et al., “Learning from Data”, Chapter 8.
- Course CE-717, Dr. M.Soleymani
- Course cs231n, Fei Fei Li, Stanford 2017.