

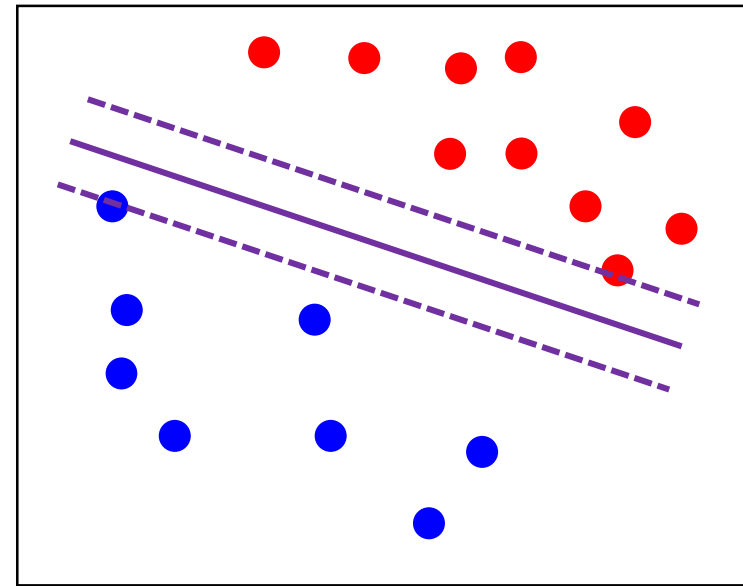
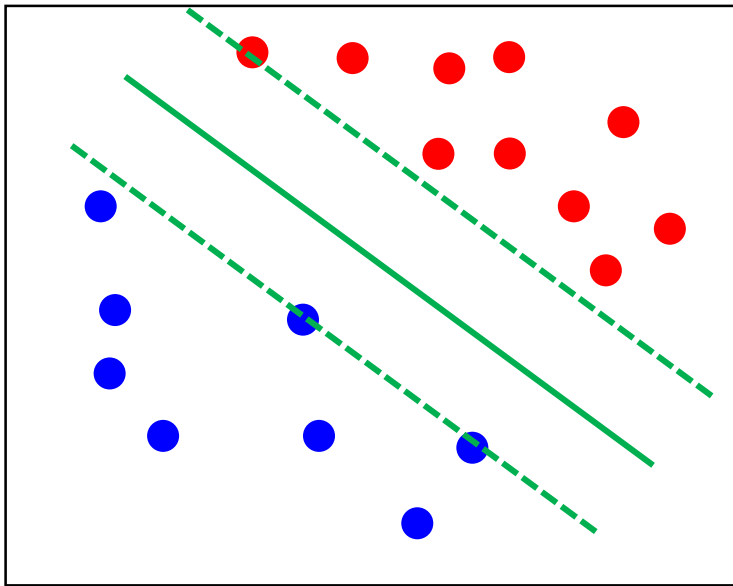
Support Vector Machine

Hanwool Jeong

hwjeong@kw.ac.kr

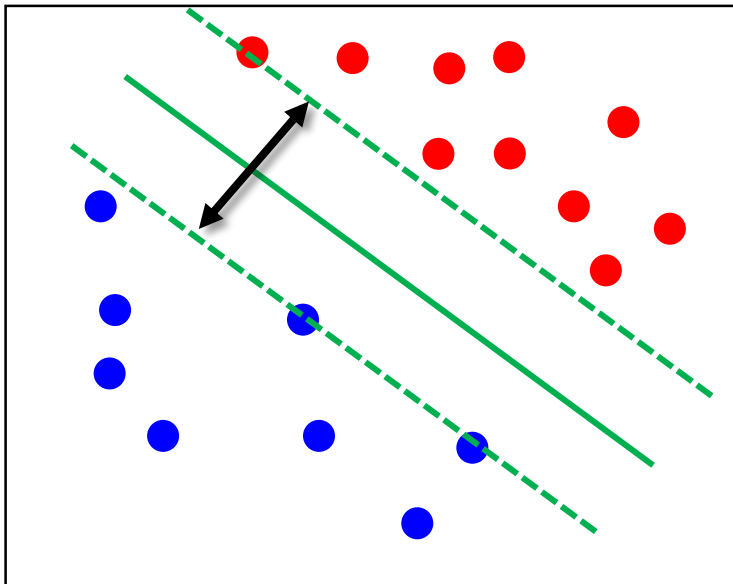
Large Margin vs. Small margin Classification

- Which is better? The optimal boundary for classification?
- The sample closest to the boundary in the dataset determines the boundary. Any other sample does not have effect.
→ Sample determining the boundary is called **Support vector**



Key Idea of Support Vector Machine (SVM)

- Key idea of SVM is to maximize the margin, which generally improves the accuracy of classification model.
- Then, we should formulate the margin then find the optimal parameter that maximize the margin.



- 1) Mathematically expressing the margin
- 2) Maximize the expressed term

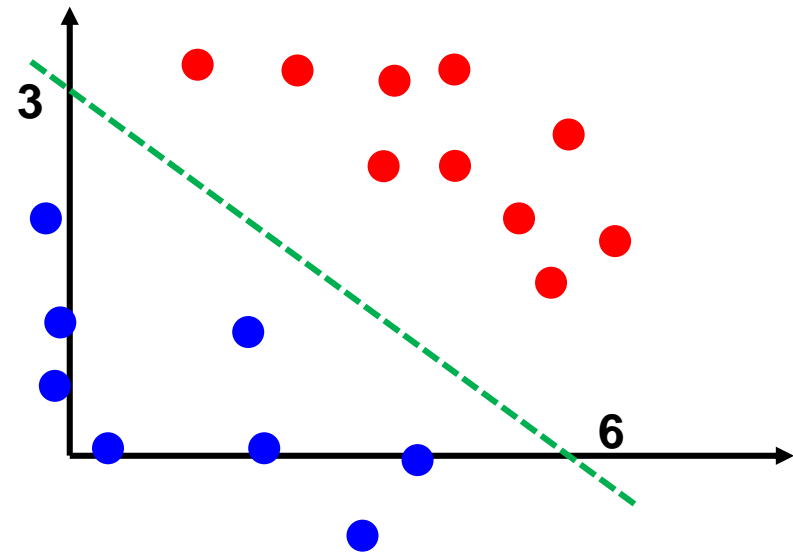
Boundary for Binary Classification

- First, let's assume linearly separable situation.
- Boundary hyperplane equation:

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$$

where \mathbf{x} is D dimensional point in the feature space.

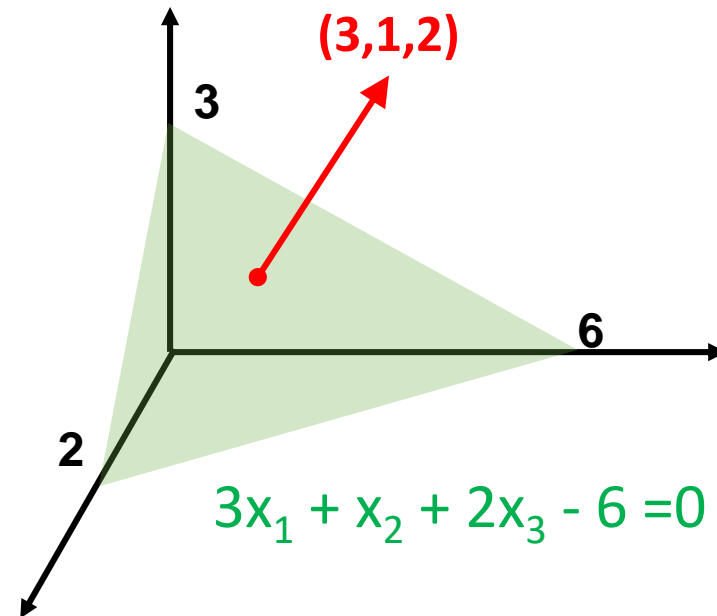
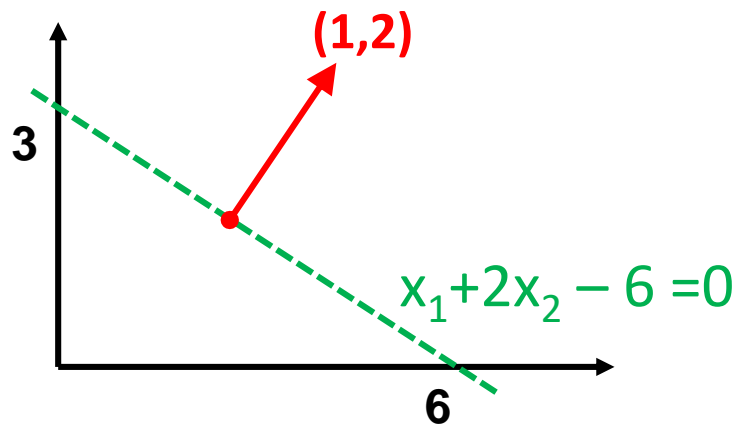
- We do classify as:



Looking into Hyperplane Expression

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

- 1) Graphical/geometrical meaning of \mathbf{w}
- 2) Is the mathematical expression unique?



General Form of Hyperplane

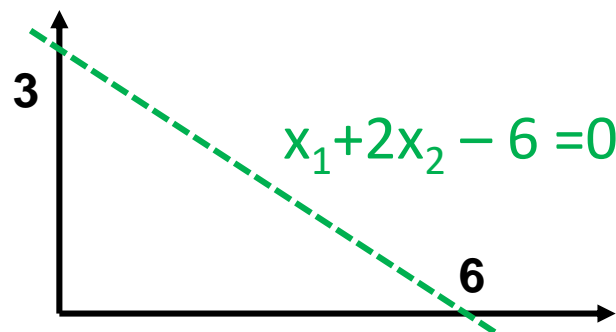
- General form : $h(\mathbf{x}) = 0$ where $h(\mathbf{x}) = ax_1 + bx_2 + c$
- $h(\mathbf{x})$ is not unique for the given hyperplane.

$$h(\mathbf{x}) = kx_1 + 2kx_2 - 6k$$

- We can freely choose $h(\mathbf{x})=0$ for mathematical convenience.

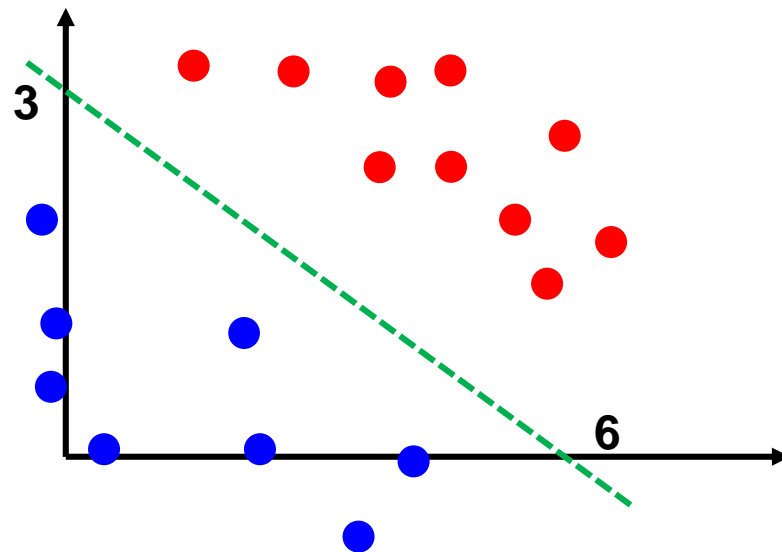
$$x_1 + 2x_2 - 6 = 0 \quad \text{vs.} \quad 2x_1 + 4x_2 - 12 = 0$$

- We can set any one condition of $h(\mathbf{x}_0)=h_0$, where \mathbf{x}_0 is not on the hyperplane, for the mathematical convenience.



Now, What do We have to do? Don't Forget Our Goal

- Mathematically expressing the margin.
- We should find out the distance from hyperplane to arbitrary point \mathbf{x} .



Distance Point to Hyperplane

- Boundary hyperplane equation:

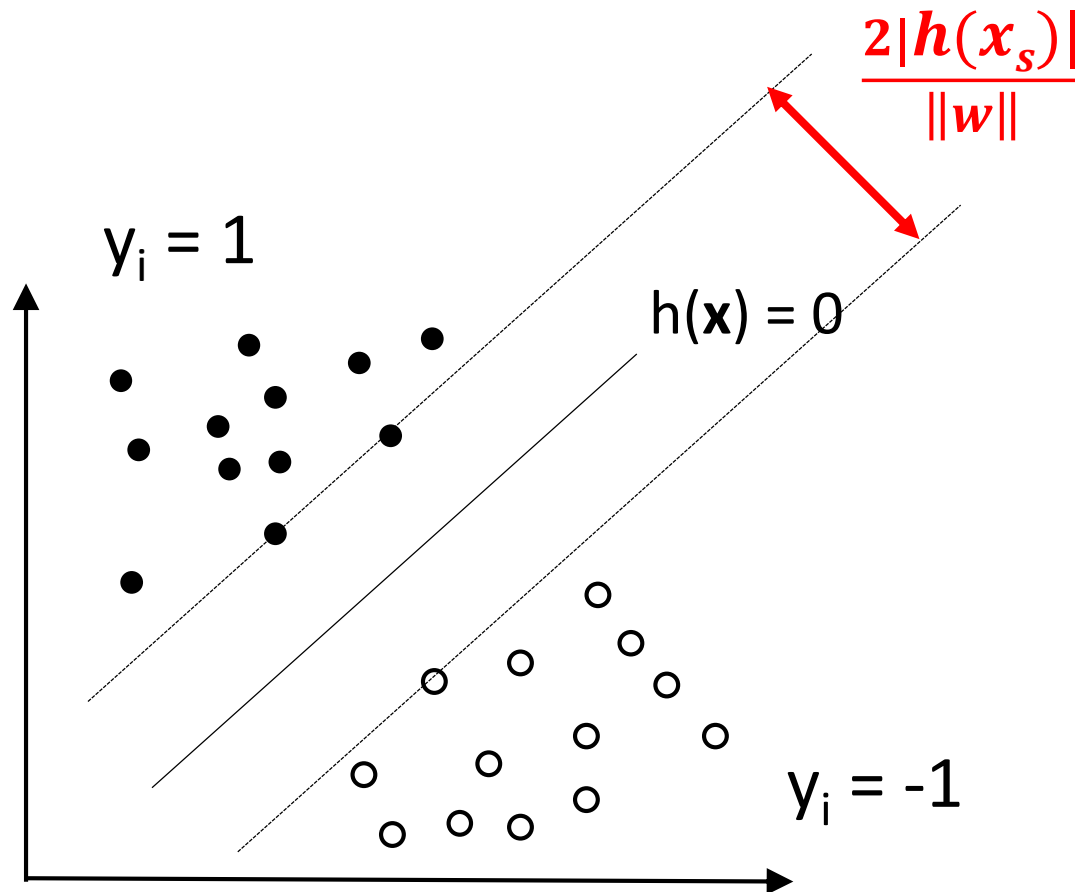
$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$$

- How can you derive the distance from the hyperplane to arbitrary point \mathbf{x} ?
- Considering \mathbf{w} means the normal vector of the plane, we can utilize it. Say the strategy for the derivation.

$$d = \frac{|h(\mathbf{x}')|}{\|\mathbf{w}\|}$$

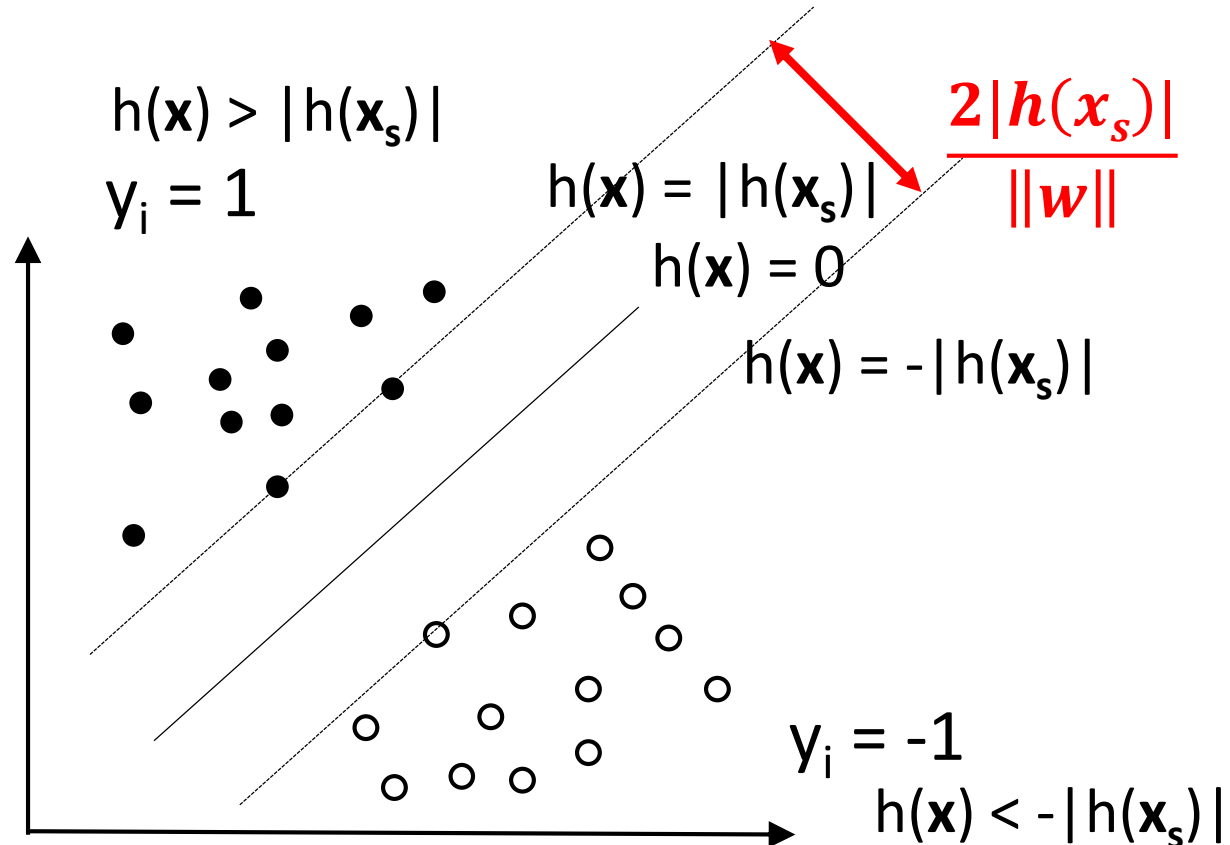
Let's See Training Set

- Maximizing the margin, is that all?
- Is there any additional **constraint** for the hyperplane?



Yes! Boundary Should Properly Classify the Dataset

- Can you express it mathematically? “Classify properly”
- Inserting all data into $h(\mathbf{x})$, then check...

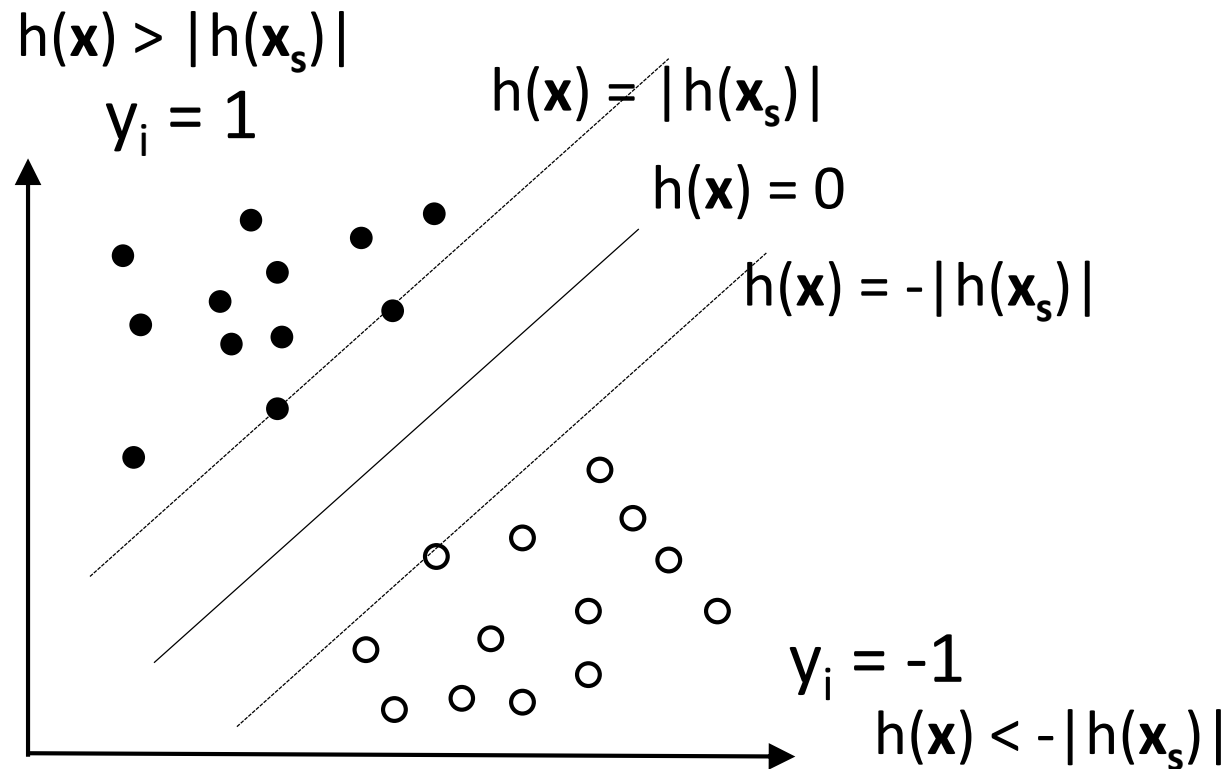


Mathematical Expression for Proper Classification Constraint

- Revisit the constraints

$$h(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b \geq |h(\mathbf{x}_s)| \text{ for } \mathbf{x}_i \text{ if } y_i = 1$$
$$h(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b \leq -|h(\mathbf{x}_s)| \text{ for } \mathbf{x}_i \text{ if } y_i = -1$$

➔ $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq |h(\mathbf{x}_s)|$



Now We Complete the Optimization Problem

- Maximize the margin
- With the constraint of the proper classification

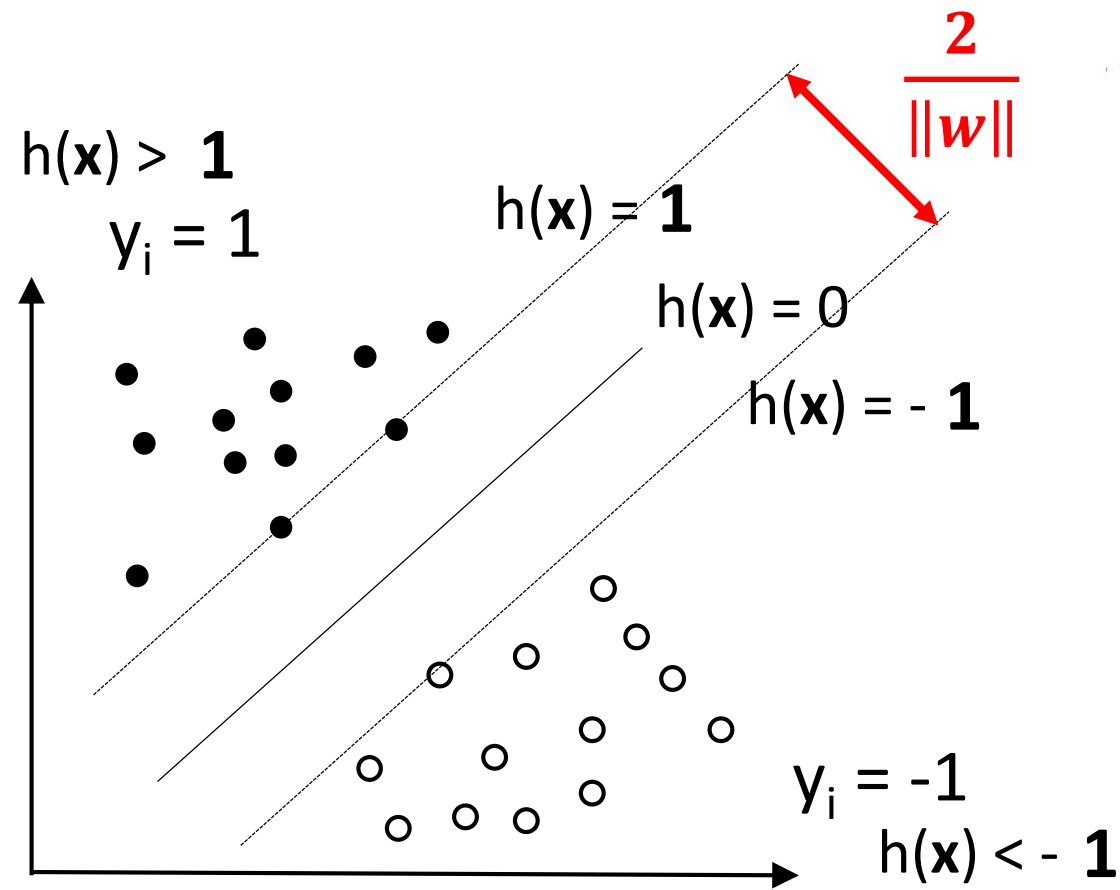
$$\begin{aligned} \text{maximize} \quad & d = \frac{2|h(x_s)|}{\|\mathbf{w}\|} \\ \text{subject to} \quad & (\mathbf{w}_i^T \mathbf{x}_i + b)y_i \geq |h(\mathbf{x}_s)| \end{aligned}$$

- Don't forget we can choose any $h(x_0) = h_0$, where x_0 is not on the hyperplane, for the mathematical convenience.
- How about setting

$$h(\mathbf{x}_s) = 1$$

Let $h(\mathbf{x}_s) = 1$

- To uniquely decide the form of $h(\mathbf{x})$



Finally Decided Optimized Problem

- Optimization for SVM boundary hyperplane.

$$\begin{array}{ll} \text{Maximize} & d = \frac{2}{\|\mathbf{w}\|} \\ \text{subject to} & (\mathbf{w}_i^T \mathbf{x}_i + b)y_i \geq \mathbf{1} \end{array}$$

- Oh, Lagrange multiplier! You remember?

Revisit KKT Condition w/ Lagrange Multiplier

- Suppose that

Minimize $f(\mathbf{x})$

subject to $\mathbf{g}_i(\mathbf{x}) \leq 0$ for $i = 1, \dots, M$

- Auxiliary function

$$L = f(\mathbf{x}) + \sum_i^M \mu_i g_i(\mathbf{x})$$

- Then solve not only $\nabla_{\mathbf{x}, \mu} L = 0$ but also

$$\mu_i g_i(\mathbf{x}^*) = 0$$

$$\mu_i \geq 0$$

Applying Lagrange Multiplier

- For the mathematical convenience , it can be reduced to

$$\begin{array}{l} \text{Maximize } \frac{2}{\|\mathbf{w}\|} \\ \text{subject to } (\mathbf{w}_i^T \mathbf{x}_i + b)y_i \geq 1 \end{array} \quad \rightarrow \quad \begin{array}{l} \text{Minimize } \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to } 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 0 \end{array}$$

- Lagrange auxiliary function is:

$$L(\mathbf{w}, b, \boldsymbol{\mu}) = \frac{\|\mathbf{w}\|^2}{2} + \sum_{i=1}^N \mu_i g(\mathbf{x}_i) = \frac{\|\mathbf{w}\|^2}{2} - \sum_{i=1}^N \mu_i \{y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1\}$$

- What is the parameter we should determine? \mathbf{w} and b
- Then, we need to find \mathbf{w} and b that satisfy $\nabla_{\mathbf{w}, b, \boldsymbol{\mu}} L = 0$ with $\mu_i \{y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1\} = 0$, and $\mu_i \geq 0$ (KKT condition)

With $L(\mathbf{w}, b, \boldsymbol{\mu}) = \frac{\|\mathbf{w}\|^2}{2} - \sum_{i=1}^N \mu_i \{y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1\}$,

- Differentiating L with \mathbf{w} then equalizing it to 0,

$$\mathbf{w} = \sum_{i=1}^N \mu_i y_i \mathbf{x}_i \quad (1)$$

- Differentiating L with b then equalizing it to 0,

$$\sum_{i=1}^N \mu_i y_i = 0 \quad (2)$$

- Before differentiating with μ , let's see KKT condition first,

$$\mu_i \geq 0 \quad (3)$$

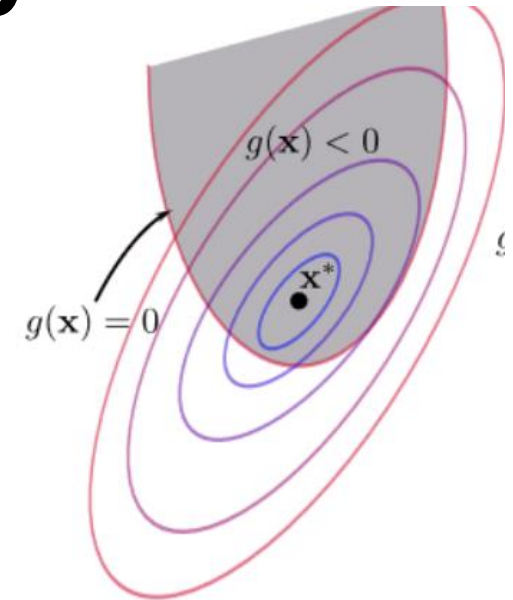
- Do you remember the meaning of $\mu_i = 0$ and $\mu_i > 0$? For i that correspond to support vectors, what would μ_i be?

$$\mu_i \{y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1\} = 0 \quad (4)$$

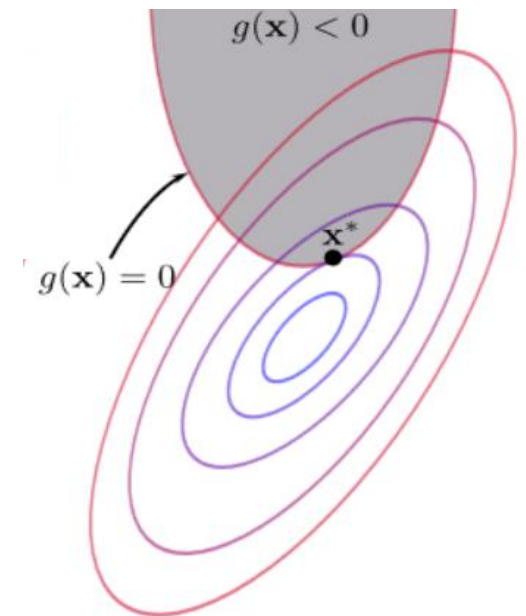
Revisit Meaning of $\mu \cdot g(\mathbf{x}) = 0$

1) $g(\mathbf{x}) \leq 0$ constraint has no meaning $\rightarrow \mu = 0$

$$g(\mathbf{x}) = 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)$$



2) Optimal point is $g(\mathbf{x}) = 0$



Let's Eliminate w & b in L

$$(1) \mathbf{w} = \sum_{i=1}^N \mu_i y_i \mathbf{x}_i \quad (2) \sum_{i=1}^N \mu_i y_i = 0 \quad (3) \mu_i \geq 0 \quad (4) \mu_i \{y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1\} = 0$$

$$\begin{aligned} L(\mathbf{w}, b, \boldsymbol{\mu}) &= \frac{\|\mathbf{w}\|^2}{2} - \sum_{i=1}^N \mu_i \{y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1\} \\ &= \frac{\|\mathbf{w}\|^2}{2} - \sum_{i=1}^N \mu_i (y_i \mathbf{w}^T \mathbf{x}_i - 1) = \frac{\|\mathbf{w}\|^2}{2} - \sum_{i=1}^N \mu_i (y_i \mathbf{w}^T \mathbf{x}_i) + \sum_{i=1}^N \mu_i \\ &= \sum_{i=1}^N \mu_i - \frac{\sum_{i=1}^N \sum_{j=1}^N \mu_i \mu_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j}{2} \end{aligned}$$

$$\begin{aligned} \frac{\|\mathbf{w}\|^2}{2} &= \frac{\mathbf{w}^T \mathbf{w}}{2} = \frac{\mathbf{w}^T \sum_{j=1}^N \mu_j y_j \mathbf{x}_j}{2} = \frac{\sum_{j=1}^N \mu_j y_j \mathbf{w}^T \mathbf{x}_j}{2} \\ &= \frac{\sum_{j=1}^N \mu_j y_j (\sum_{i=1}^N \mu_i y_i \mathbf{x}_i^T) \mathbf{x}_j}{2} = \frac{\sum_{i=1}^N \sum_{j=1}^N \mu_i \mu_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j}{2} \end{aligned}$$

Lagrange Auxiliary Function $L(\mu)$ & Determining w & b

$$L(\mu) = \sum_{i=1}^N \mu_i - \frac{\sum_{i=1}^N \sum_{j=1}^N \mu_i \mu_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j}{2}$$

- We should find μ making $\nabla_{\mu} L = 0$. From now, it is not our job.
- Does this maximize or minimize $L(\mu)$?
- Then, the problem is changed into finding μ

Maximize $L(\mu)$

Subject to $\sum_{i=1}^N \mu_i y_i = 0$ and $\mu_i \geq 0$

Determining w & b from Optimal μ

- Then we can find \mathbf{w}^*

$$\mathbf{w}^* = \sum_{i=1}^N \mu_i y_i \mathbf{x}_i$$

- Note that only for support vectors, $\mu_i > 0$ otherwise $\mu_i = 0$
- How about b ? You remember? $\mu_i \{y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1\} = 0$

$$\sum_{i=S.V} \{y_i(\mathbf{w}^{*T} \mathbf{x}_i + b) - 1\} = 0 \Rightarrow (1/N_{S.V}) \sum_{i=S.V} (y_i - \mathbf{w}^{*T} \mathbf{x}_i) = b$$

Prediction?

- With determined \mathbf{w} and b ,

$$\mathbf{w}^* = \sum_{i=1}^N \mu_i y_i \mathbf{x}_i$$

$$b = (1/N_{S.V.}) \sum_{i=S.V.} (y_i - \mathbf{w}^{*\top} \mathbf{x}_i)$$

Checkpoints

- Expressing the margin (to be maximized) mathematically
- Expressing the constraint mathematically
- Applying Lagrange multiplier to find out the optimal SVM boundary
- Coming up next : soft boundary & kernel trick