

Example

- $\mathbf{x}_1^T = (2,3)$, $\mathbf{x}_2^T = (4,1)$, $\mathbf{x}_3^T = (5,1)$
- $y_1 = +1$, $y_2 = -1$, $y_3 = -1$

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

Support Vector Machine w/ Soft Margin & Kernel Trick

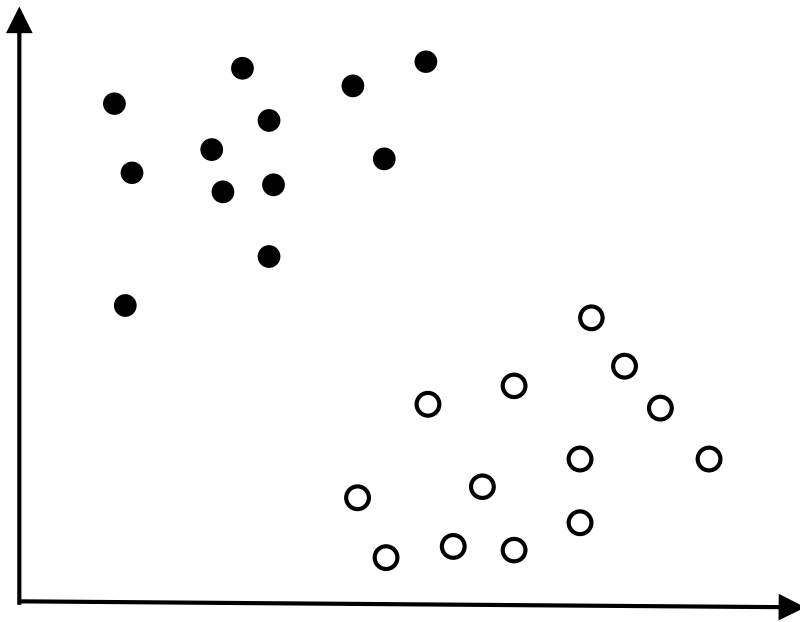
Hanwool Jeong

hwjeong@kw.ac.kr

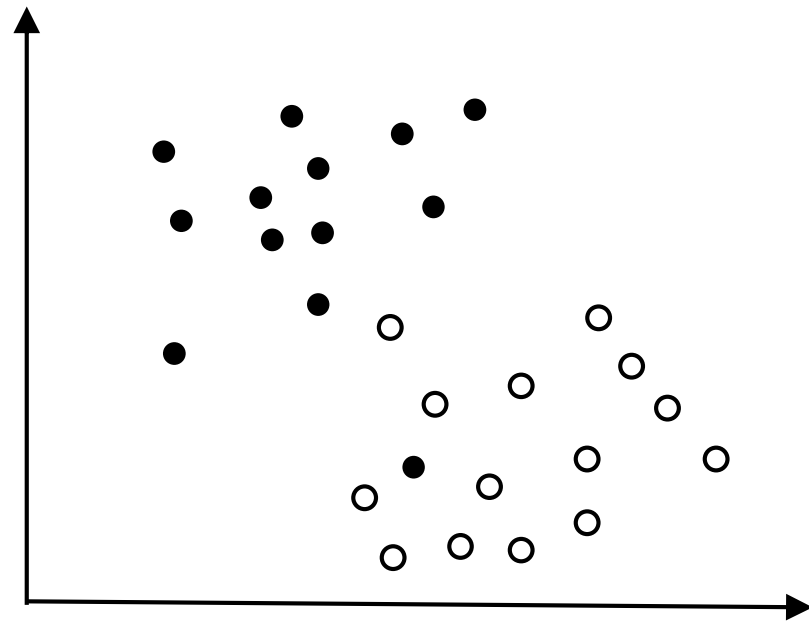
Soft Margin vs. Hard Margin

- There can be outliers

Hard margin

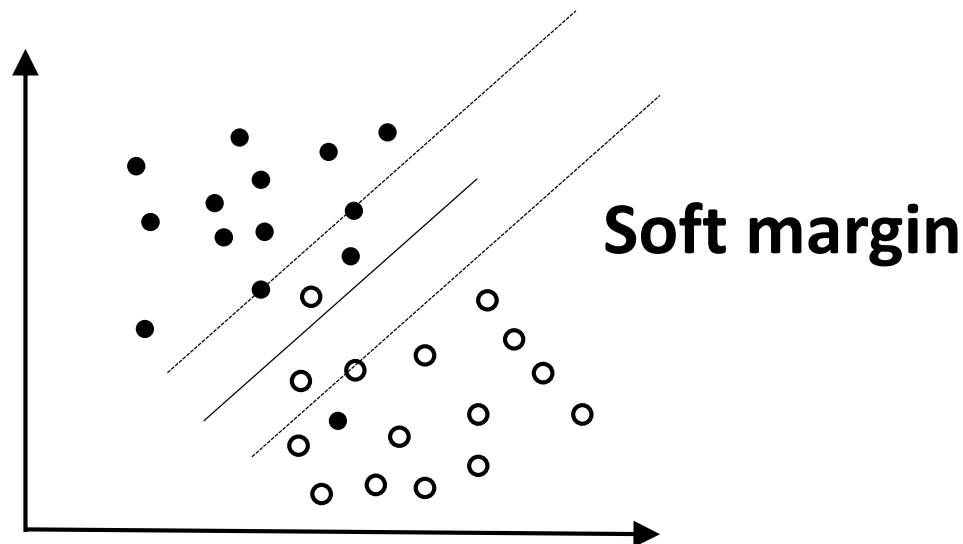


Soft margin



What is Optimal Boundary?

- We can categorize the samples into three types according to value of $y_i(\mathbf{w}^T \mathbf{x}_i + b)$
 - 1) $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$
 - 2) $0 < y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1$
 - 3) $y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$
- A slack variable ξ is defined using $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$.
- ➔ What is the meaning of $\sum_{i=1}^n \xi_i$ magnitude?



Problem Defining

- Maximize the number of samples 1) while the minimizing the number of samples of 2) & 3)
- We can revise the hard margin example into

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

$$g(\mathbf{x}_i) = 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 0$$



$$\text{Minimize } \frac{\|\mathbf{w}\|^2}{2} + C \sum_{i=1}^N \xi_i$$

$$\text{subject to } 0 \leq \xi_i$$

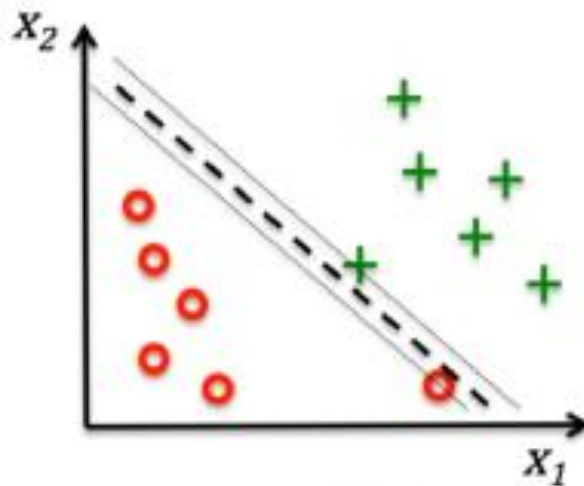
$$g(\mathbf{x}_i) = 1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 0$$

- What is the meaning of C?
- We can repeat the procedure by the Lagrange aux. function of

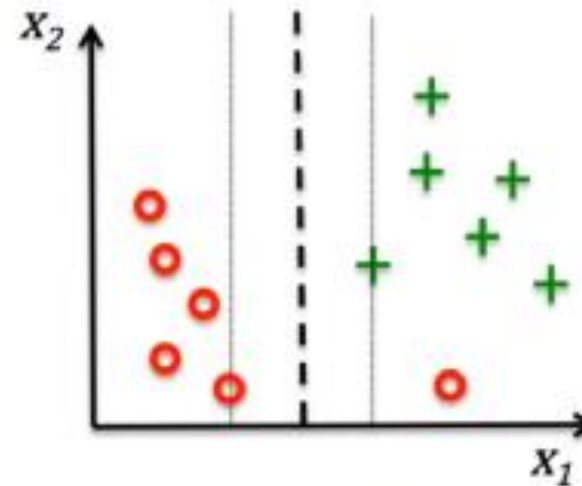
$$L = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \mu_i \{y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i\}$$

Impact of C

- $C = 0$ ignores the second term just maximizing margin and accuracy is not considered.
- With larger C accuracy becomes more taken into account.

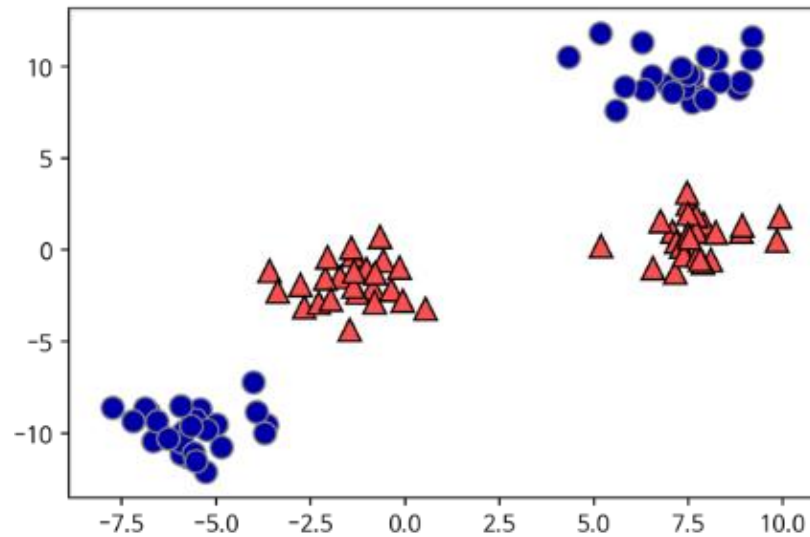


Large value for
parameter C



Small value for
parameter C

Linearly Separable vs. Inseparable



Kernel Trick!

- Revisit Lagrange auxiliary function:

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

- Using kernel trick, the following Lagrange auxiliary function should be maximized:

$$\begin{aligned} L(\mu) &= \sum_{i=1}^N \mu_i - \frac{\sum_{i=1}^N \sum_{j=1}^N \mu_i \mu_j y_i y_j \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)}{2} \\ &= \sum_{i=1}^N \mu_i - \frac{\sum_{i=1}^N \sum_{j=1}^N \mu_i \mu_j y_i y_j K(\mathbf{x}_i \bullet \mathbf{x}_j)}{2} \end{aligned}$$

➔ We can exploit the advantage of high dimension space without severe complexity increase.

Prediction?

- For prediction for \mathbf{x}_{new} , we can use

$$\mathbf{w}^* \text{ (in high dim)} = \sum_{i=1}^N \mu_i y_i \Phi(\mathbf{x}_i)$$

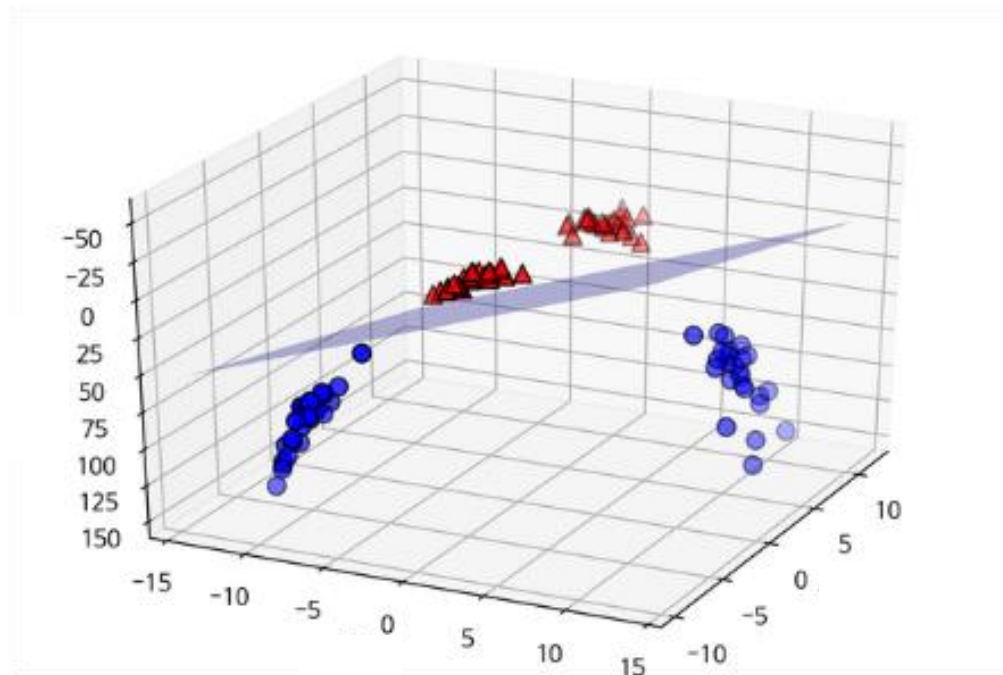
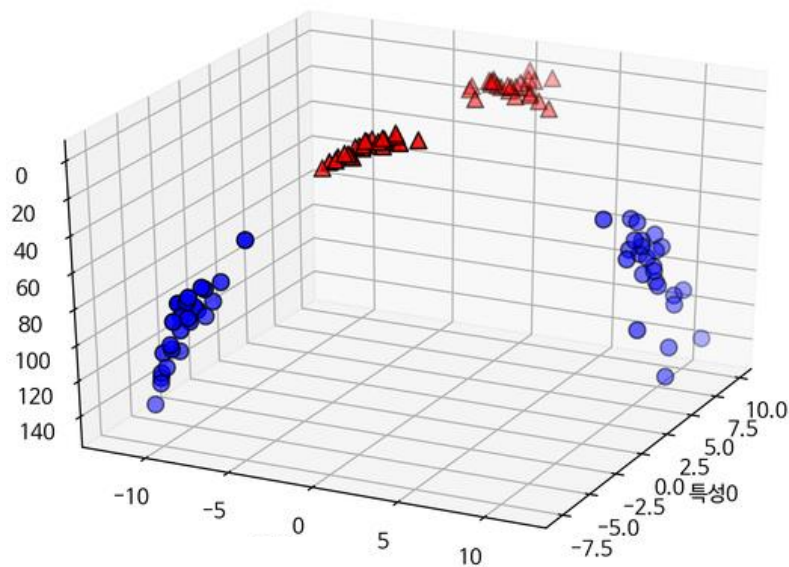
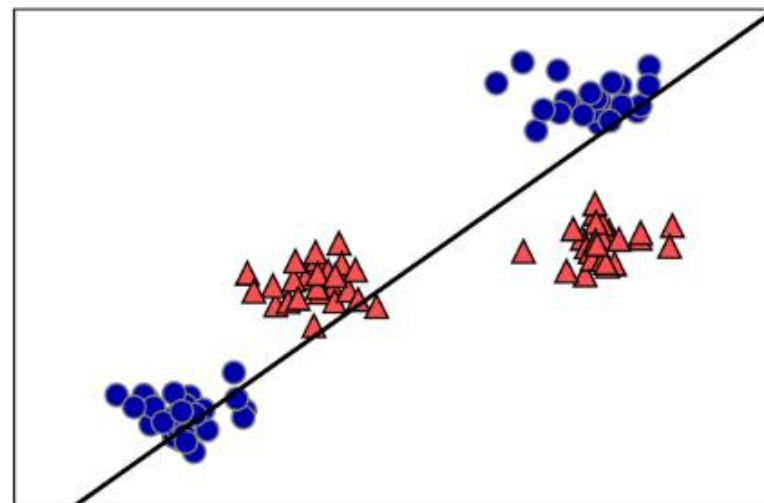
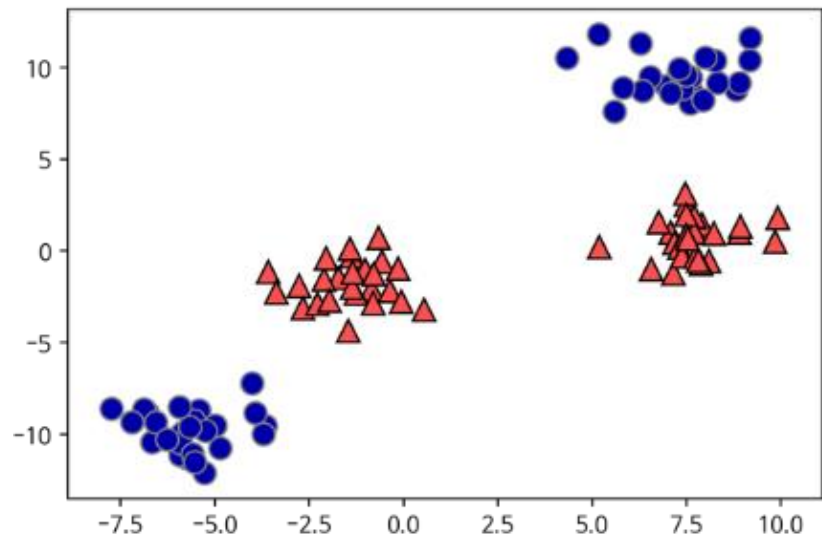
$$b = (1/N_{S.V.}) \sum_{j=S.V.} \{y_i - \mathbf{w}^{*\top} \Phi(\mathbf{x}_j)\}$$

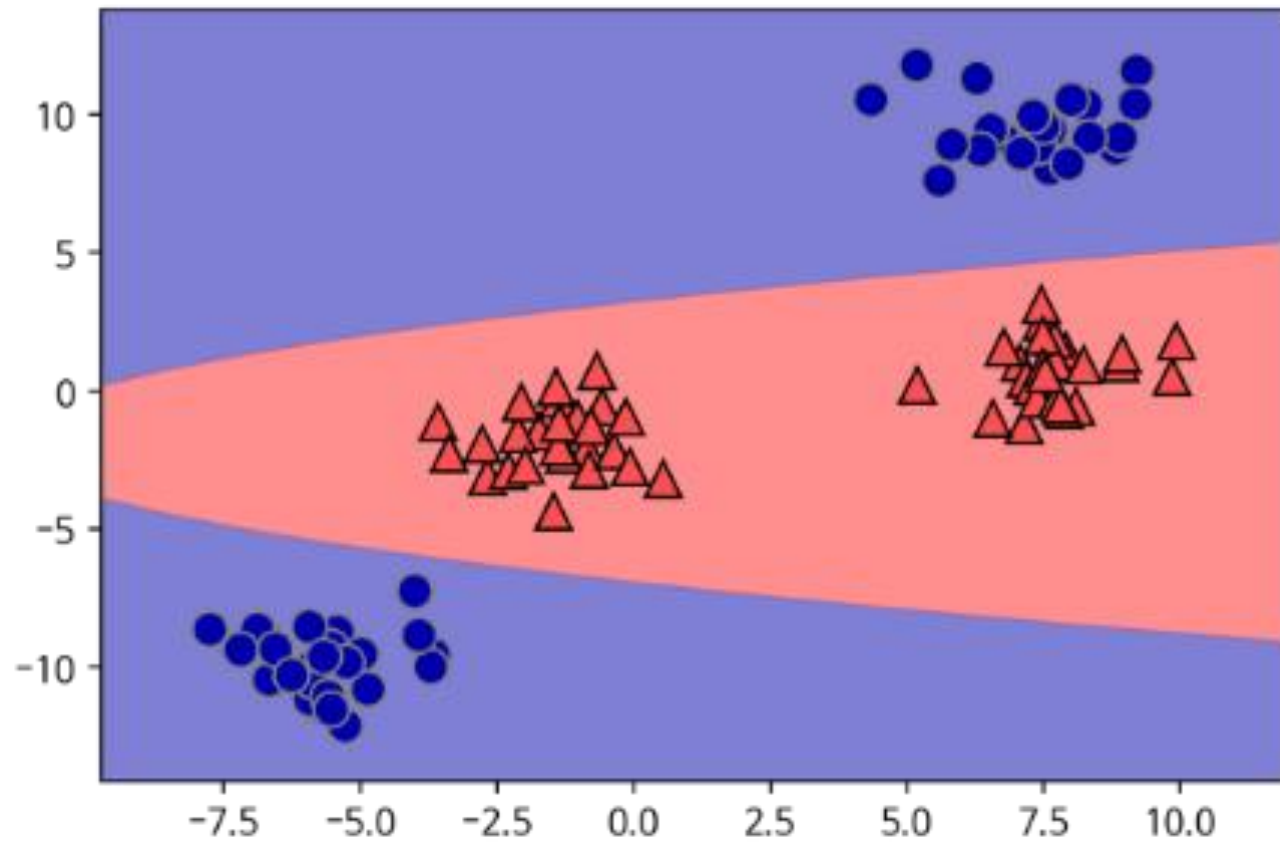
Then apply to below.

$$\text{sign}(\mathbf{w}^{*\top} \Phi(\mathbf{x}_{\text{new}}) + b^*)$$

- Do we need the mapping function and high dimension calculation?
- No!

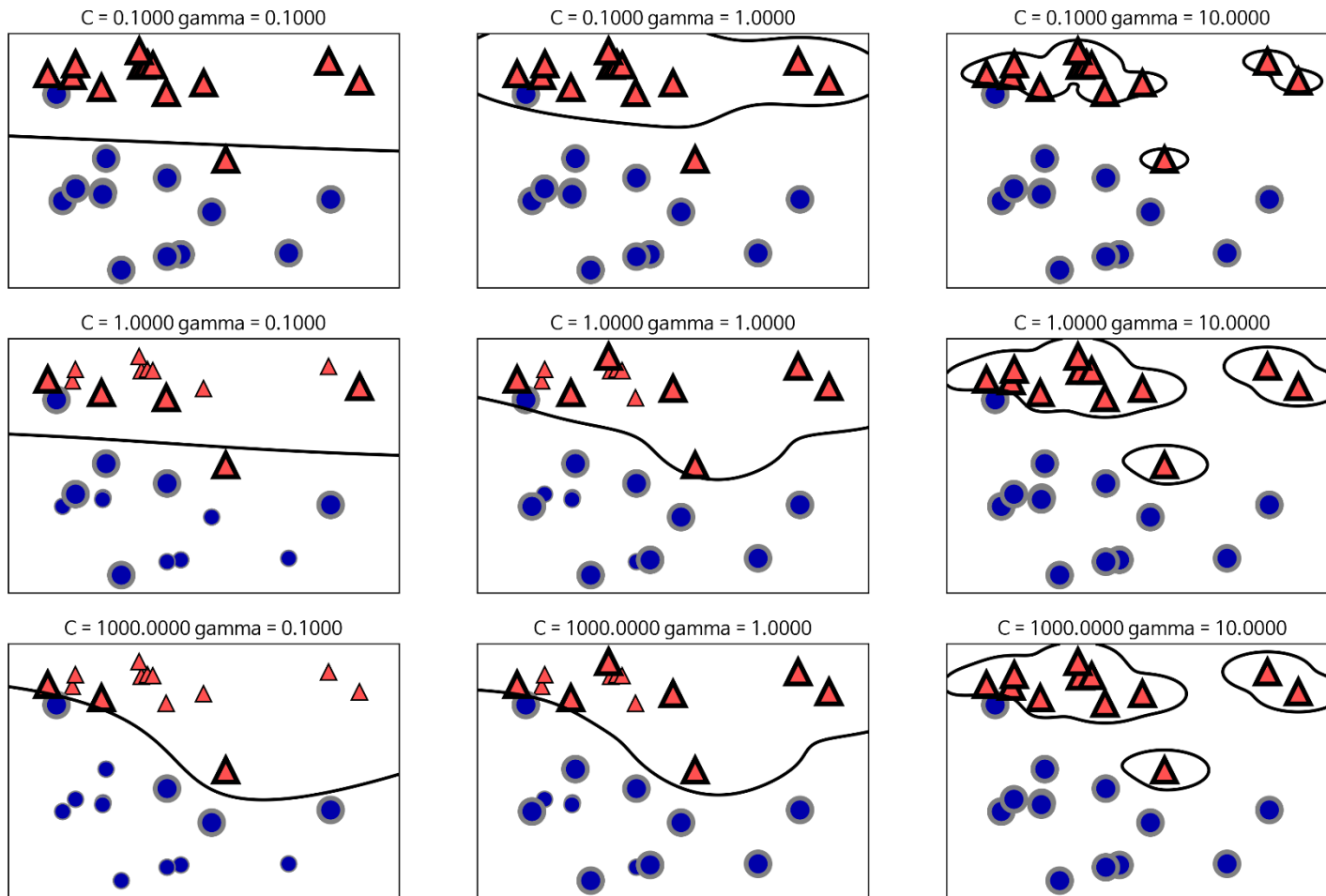
$$\begin{aligned} & \sum_{i=1}^N \mu_i y_i \Phi(\mathbf{x}_i) \bullet \Phi(\mathbf{x}_{\text{new}}) + (1/N_{S.V.}) \sum_{k=S.V.} \{y_i - \sum_{i=1}^N \mu_i y_i \Phi(\mathbf{x}_i) \Phi(\mathbf{x}_k)\} \\ & = \sum_{i=1}^N \mu_i y_i K(\mathbf{x}_i, \mathbf{x}_{\text{new}}) + (1/N_{S.V.}) \sum_{k=S.V.} \{y_i - \sum_{i=1}^N \mu_i y_i K(\mathbf{x}_i, \mathbf{x}_k)\} \end{aligned}$$





Applying Soft Margin & Kernel Trick

$$k_{rbf}(x_1, x_2) = \exp(-\gamma \|x_1 - x_2\|^2)$$



Checkpoints

- ✓ Implementing soft margins in SVM
- ✓ Applying kernel trick for SVM
- ✓ Coming up next : python coding for SVM