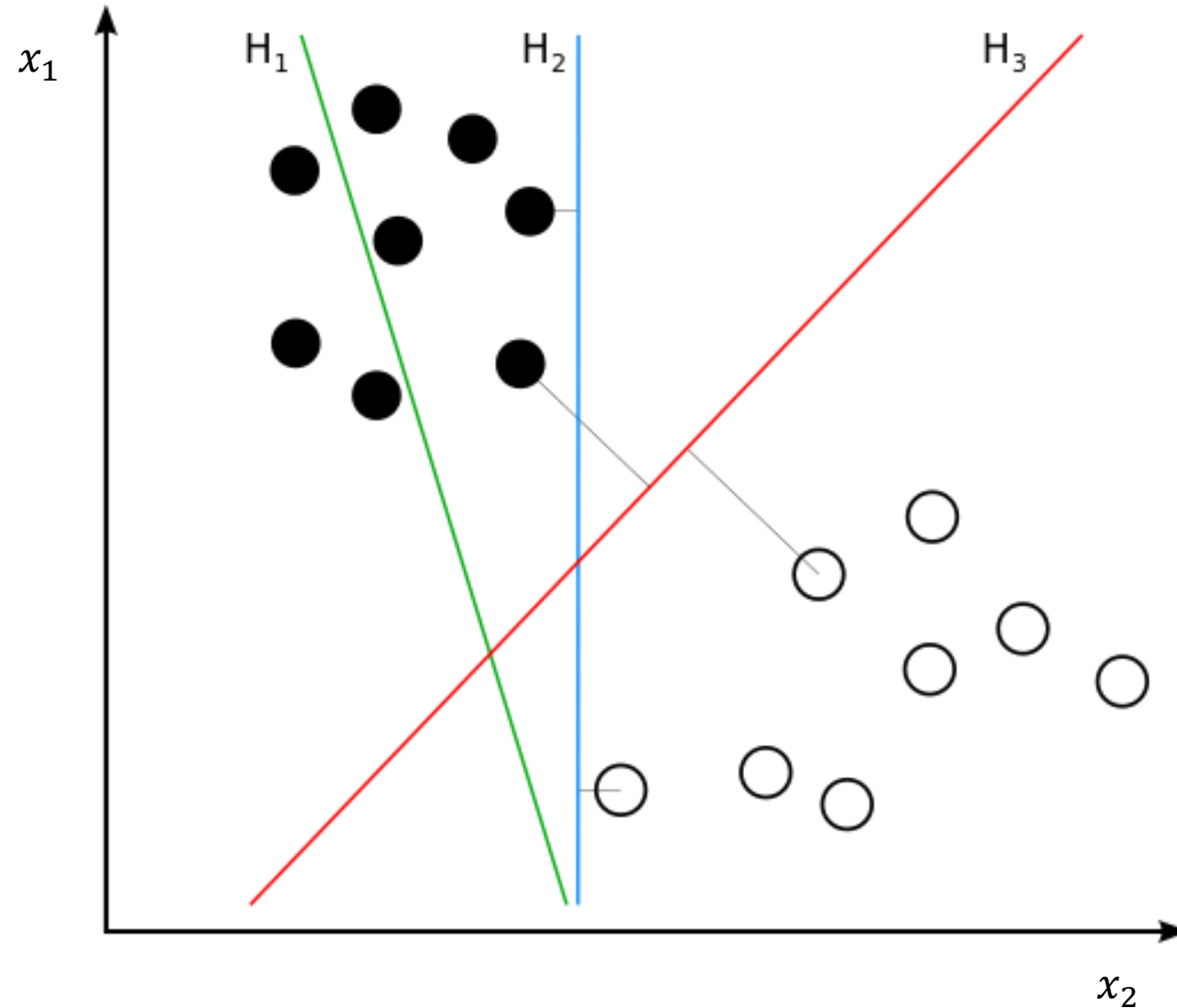


SVM

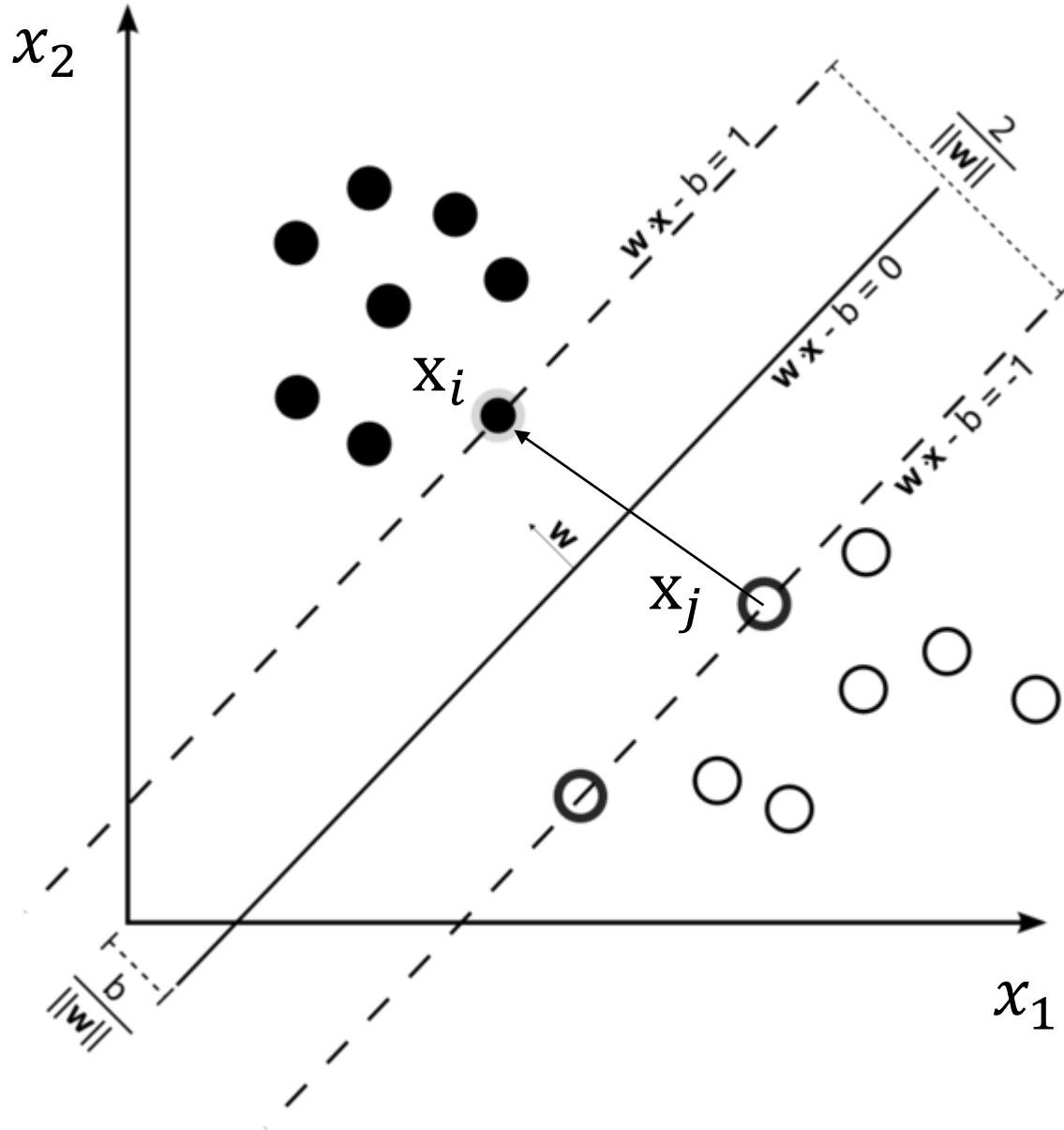
Support Vector Machines

# SVM (Support Vector Machines)



# Maximize the margin

## Linearly separable case



Scale  $\mathbf{w}$  and  $\mathbf{b}$  so that:

$$\min |\mathbf{w}^T \mathbf{x}_i - \mathbf{b}| = \min (y_i (\mathbf{w}^T \mathbf{x}_i - \mathbf{b})) = 1,$$

Then the width of the margin is:

$$\frac{\mathbf{w}^T (\mathbf{x}_i - \mathbf{x}_j)}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T \mathbf{x}_i - \mathbf{b} - (\mathbf{w}^T \mathbf{x}_j - \mathbf{b})}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$$

Our task become the maximizing of  $\frac{2}{\|\mathbf{w}\|}$ ,

or minimizing  $\|\mathbf{w}\|$ , or  $\mathbf{w}^T \mathbf{w}$

Under  $y_i (\mathbf{w}^T \mathbf{x}_i - \mathbf{b}) \geq 1$  constraints.

# Optimization task

$$\begin{cases} \frac{1}{2} \mathbf{w}^T \mathbf{w} \rightarrow \min \\ y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1 \end{cases}$$

# Karush-Kuhn-Tucker conditions

Our (primal) optimization under constraints problem:

$$\begin{cases} \min_z f(z) \\ g_i(z) \leq 0 \\ h_i(z) = 0 \end{cases}$$

If  $z^*$  – is a local minimum, then the exist Lagrangian multipliers  $\alpha_i^*$  and  $\beta_j^*$  for:

$$\mathcal{L}(z, \alpha, \beta) = f(z) + \sum_{i=1}^m \alpha_i g_i(z) + \sum_{j=1}^k \beta_j h_j(z),$$

such that:

$$\begin{cases} \frac{\partial}{\partial z_i} \mathcal{L}(z^*, \alpha^*, \beta^*) = 0, \\ \frac{\partial}{\partial \beta_i} \mathcal{L}(z^*, \alpha^*, \beta^*) = 0, \\ \alpha_i g_i(z^*) = 0, \\ \alpha_i^* \geq 0 \end{cases}$$

And the solution of the primal problem is the solution of the dual problem:  $\max_{\alpha, \beta} \mathcal{L}(z, \alpha, \beta)$

# Dual problem

$$\begin{cases} \frac{1}{2} \mathbf{w}^T \mathbf{w} \rightarrow \min \\ y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1 \end{cases} \rightarrow \begin{cases} \frac{1}{2} \mathbf{w}^T \mathbf{w} \rightarrow \min \\ -(y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1) \leq 0 \end{cases}$$



$$\mathcal{L}(\underbrace{\mathbf{w}, b}_{\mathbf{z}}, \alpha) = \frac{1}{2} (\mathbf{w}^T \mathbf{w}) - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1)$$

$$\alpha_i \geq 0; \quad \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1) = 0$$

# Dual problem solution

$$\mathcal{L}(w, \alpha, b) = \frac{1}{2}(w^T w) - \sum_{i=1}^N \alpha_i (y_i(w^T x_i - b) - 1); \quad \alpha_i \geq 0; \quad \alpha_i(y_i(w^T x_i - b) - 1) = 0$$

$$\nabla_w \mathcal{L}(w, \alpha, b) = w - \sum_{i=1}^N \alpha_i y_i x_i = 0 \rightarrow w = \sum_{i=1}^N \alpha_i y_i x_i$$

$$\frac{\partial}{\partial b} \mathcal{L}(w, \alpha, b) = \sum_{i=1}^N \alpha_i y_i = 0$$

$$\begin{aligned} \mathcal{L}(w, \alpha, b) &= \frac{1}{2}(w^T w) - \sum_{i=1}^N \alpha_i (y_i(w^T x_i - b) - 1) = \frac{1}{2} \sum \sum y_i y_j \alpha_i \alpha_j x_i^T x_j - \sum \sum y_i y_j \alpha_i \alpha_j x_i^T x_j + \sum \alpha_i = \\ &= \mathcal{L}(w, \alpha, b) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j x_i^T x_j \end{aligned}$$

Quadratic optimization problem under linear constraints is efficiently solved by quadratic programming.

# Support vectors

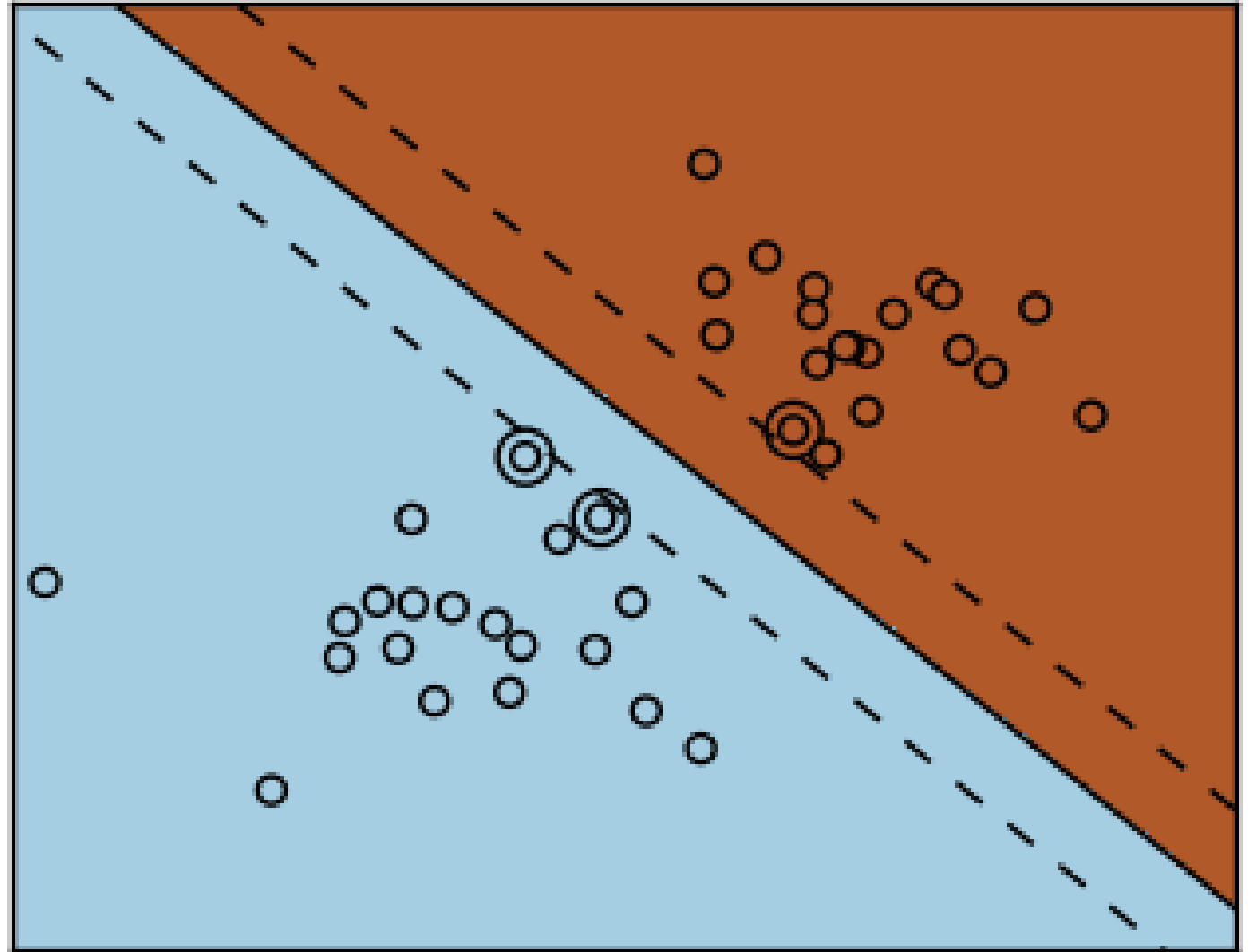
$$w = \sum_{i=1}^N \alpha_i y_i x_i$$

This is how we can find  $b$ .

$$\alpha_i (y_i (w^T x_i - b) - 1) = 0$$

$$x_i : \alpha_i > 0$$

$$y_i (w^T x_i - b) - 1 = 0$$





# Dual problem solution with CVXOPT Package

Task that QP solver solves:

$$\begin{cases} \frac{1}{2} \alpha^T P \alpha + q^T \alpha \rightarrow \min \\ G \alpha \leq h \\ A \alpha = b \end{cases}$$

$$\begin{cases} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \rightarrow \max \\ \alpha_i \geq 0 \\ \sum_{i=1}^N \alpha_i y_i = 0 \end{cases} \rightarrow \begin{cases} P = [y_i y_j \mathbf{x}_i^T \mathbf{x}_j] & q = [-1] \\ G = -I & h = [0] \\ A = y^T & b = 0 \end{cases}$$

# Kernel trick

Notice, that we only used dot products in all calculations (i.e.  $\mathbf{x}_i^T \mathbf{x}_j$ ). That means that we do not need to transition into higher dimensional space but rather only define a dot product operation there.

$K(\mathbf{x}, \mathbf{x}')$  is a kernel function if  $K(\mathbf{x}, \mathbf{x}') = \psi(\mathbf{x})^T \psi(\mathbf{x}')$ ,  
where  $\psi: X \rightarrow H$ , and  $H$  is some Hilbert space.

Example:

$$K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2 = 1 + x_1^2 x_1'^2 + x_2^2 x_2'^2 + 2x_1 x_1' + 2x_2 x_2' + 2x_1 x_1' x_2 x_2'$$

$$\psi(\mathbf{x}) = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2)$$

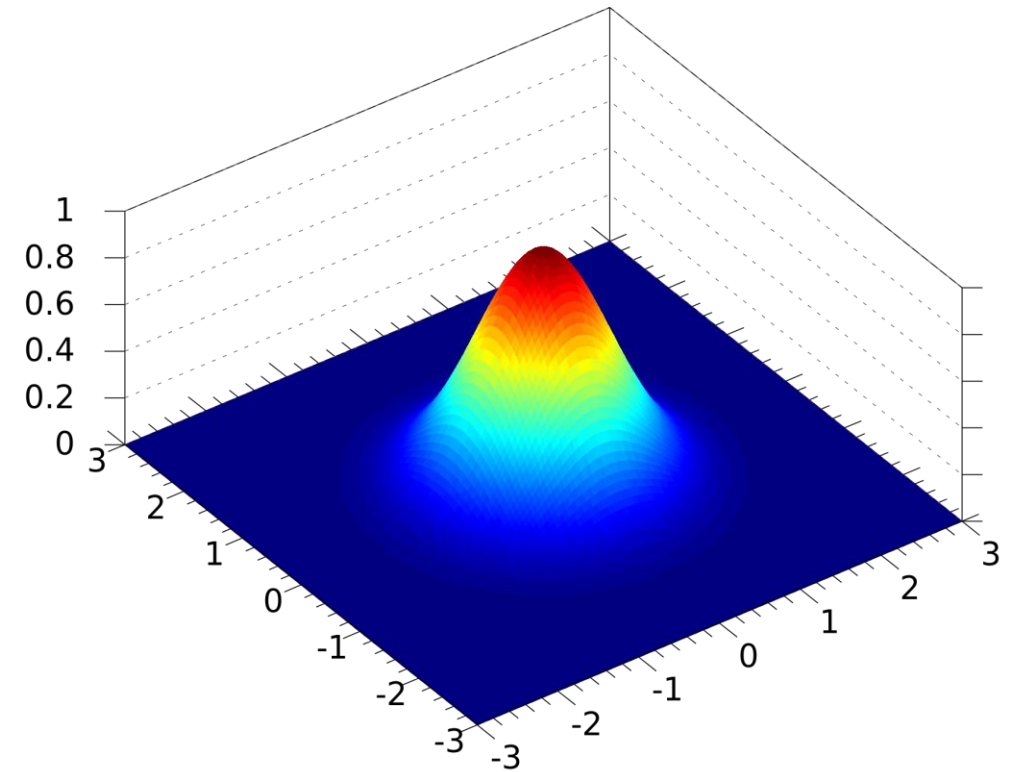
# Kernels

Linear kernel:  $\langle \mathbf{x}, \mathbf{x}' \rangle$

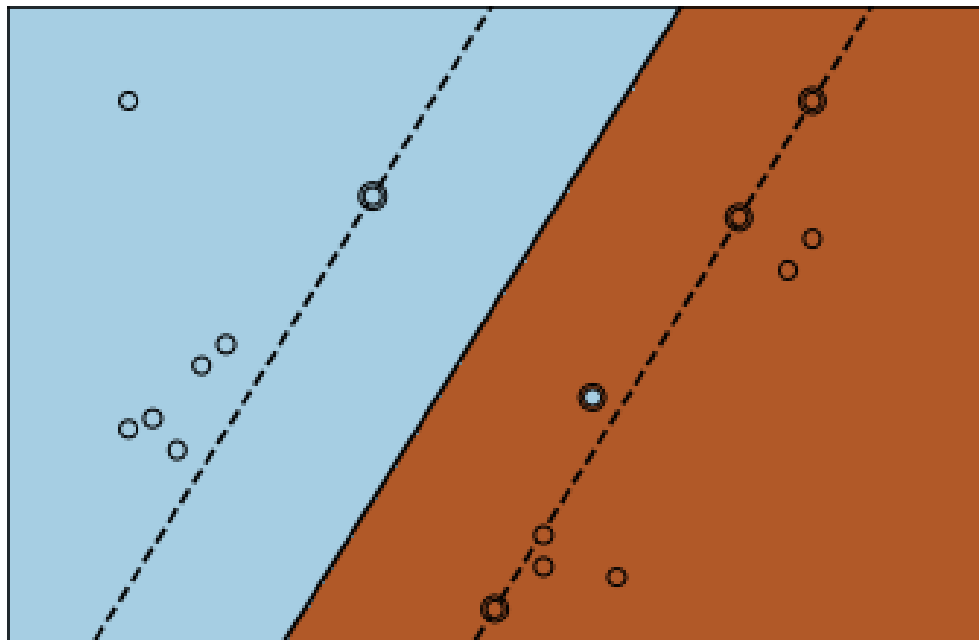
Polynomial kernel:  $(r + \gamma \langle \mathbf{x}, \mathbf{x}' \rangle)^d$

Radial basis function:  $e^{-\gamma \|\mathbf{x} - \mathbf{x}'\|^2}$

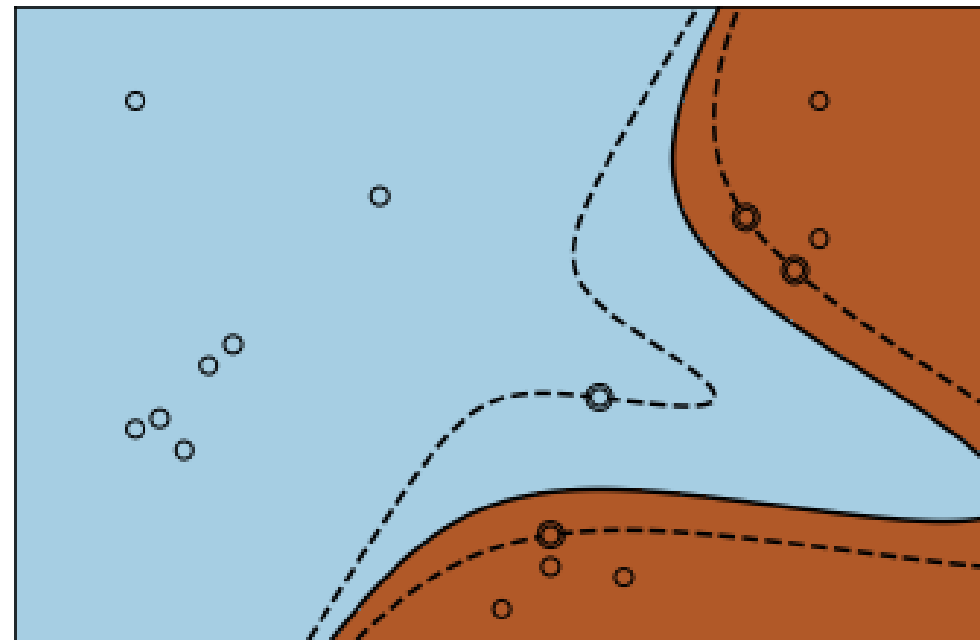
Sigmoid kernel:  $\text{th}(\gamma \langle \mathbf{x}, \mathbf{x}' \rangle + r)$



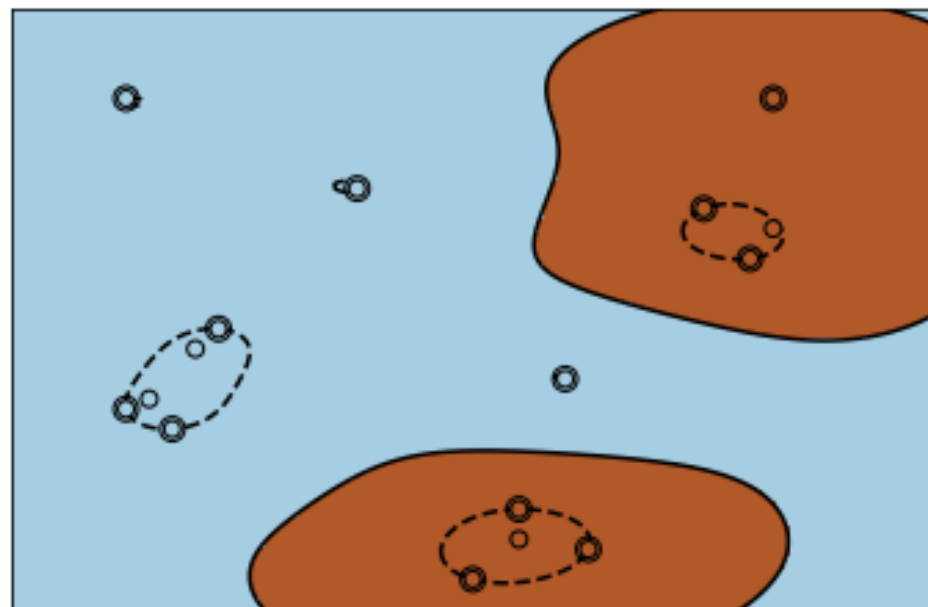
Radial Basis Function (RBF)



Linear kernel



Polynomial kernel



Radial Basis Function

# Soft margin

Linearly separable case:

$$\begin{cases} \frac{1}{2} \mathbf{w}^T \mathbf{w} \rightarrow \min \\ y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1 \end{cases}$$

Linearly inseparable case:

$$\begin{cases} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \rightarrow \min \\ y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{cases}$$

# Dual problem for soft margin

$$\begin{aligned}\mathcal{L}(\overbrace{w, b, \xi}^z, \underbrace{\alpha, r}_{\alpha}) &= \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i (y_i (w^T x_i - b) - 1 + \xi_i) - \sum_{i=1}^N r_i \xi_i \\ &= \frac{1}{2} w^T w - \sum \alpha_i (y_i (w^T x_i - b) - 1) - \sum \xi_i (r_i + \alpha_i - C)\end{aligned}$$

$$\alpha_i, r_i \geq 0; \quad \alpha_i (y_i (w^T x_i - b) - 1 + \xi_i) = 0; \quad r_i \xi_i = 0$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial \xi} = 0 \Rightarrow \begin{cases} w = \sum \alpha_i y_i x_i \\ \sum \alpha_i y_i = 0 \\ \alpha_i = C - r_i \Rightarrow 0 \leq \alpha_i \leq C \end{cases}$$

# Dual problem for soft margin

$$\left\{ \begin{array}{l} \mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i \\ \sum \alpha_i y_i = 0 \\ \alpha_i = C - r_i \end{array} \right. \Rightarrow \begin{aligned} \mathcal{L}(\alpha) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1) - \sum \xi_i (r_i + \alpha_i - C) = \\ &= \frac{1}{2} \sum \sum y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j - \sum \sum y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j + \sum \alpha_i \end{aligned}$$

$$\left\{ \begin{array}{l} \max_{\alpha} \mathcal{L}(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \\ 0 \leq \alpha_i \leq C \\ \sum \alpha_i y_i = 0 \end{array} \right.$$

# Vector types

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i; \quad y_i(\mathbf{w}^T \mathbf{x}_i - b) \geq 1 - \xi_i$$

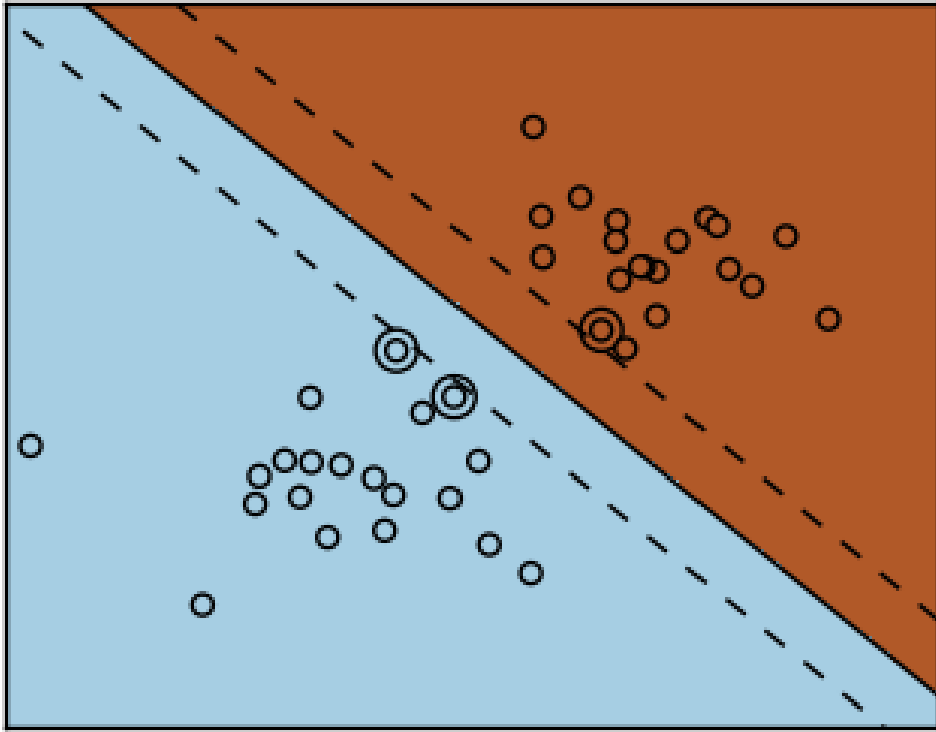
$$\alpha_i = C - r_i; \quad \alpha_i (y_i(\mathbf{w}^T \mathbf{x}_i - b) - 1 + \xi_i) = 0; \quad r_i \xi_i = 0$$

1.  $\alpha_i = 0; \xi_i = 0; y_i(\mathbf{w}^T \mathbf{x}_i - b) \geq 1$  – Inside Vectors
2.  $0 < \alpha_i < C; \xi_i = 0; y_i(\mathbf{w}^T \mathbf{x}_i - b) = 1$  – **"Good" support vectors**
3.  $\alpha_i = C; \xi_i > 0; y_i(\mathbf{w}^T \mathbf{x}_i - b) \leq 1$  – **"Bad" support vectors**

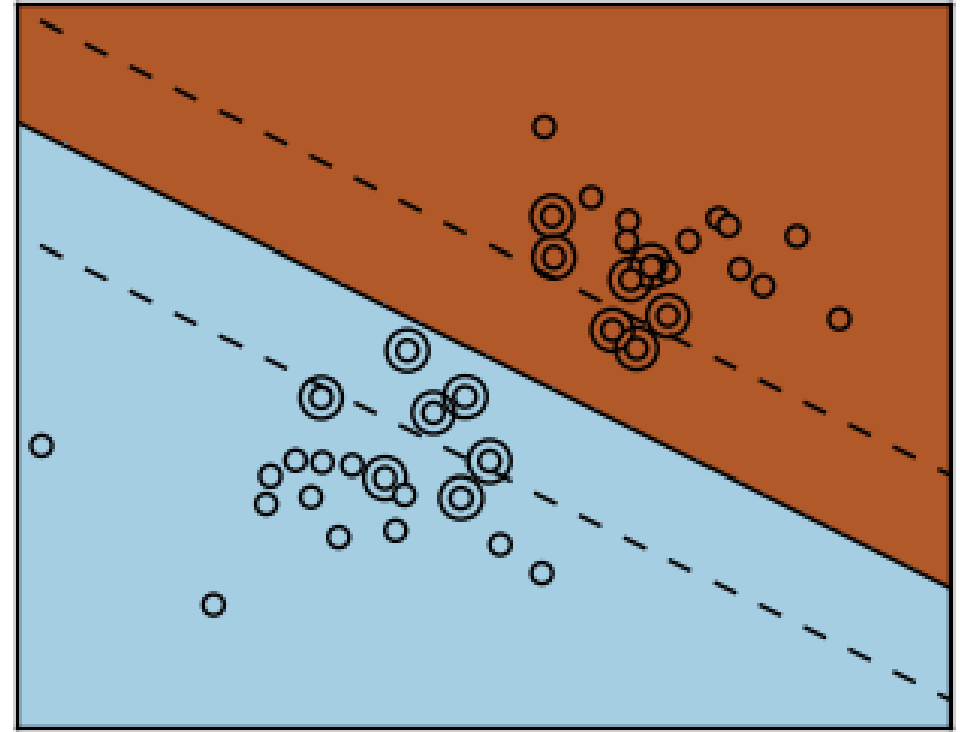


# Soft margin coefficient

$$\min \left( \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \right)$$

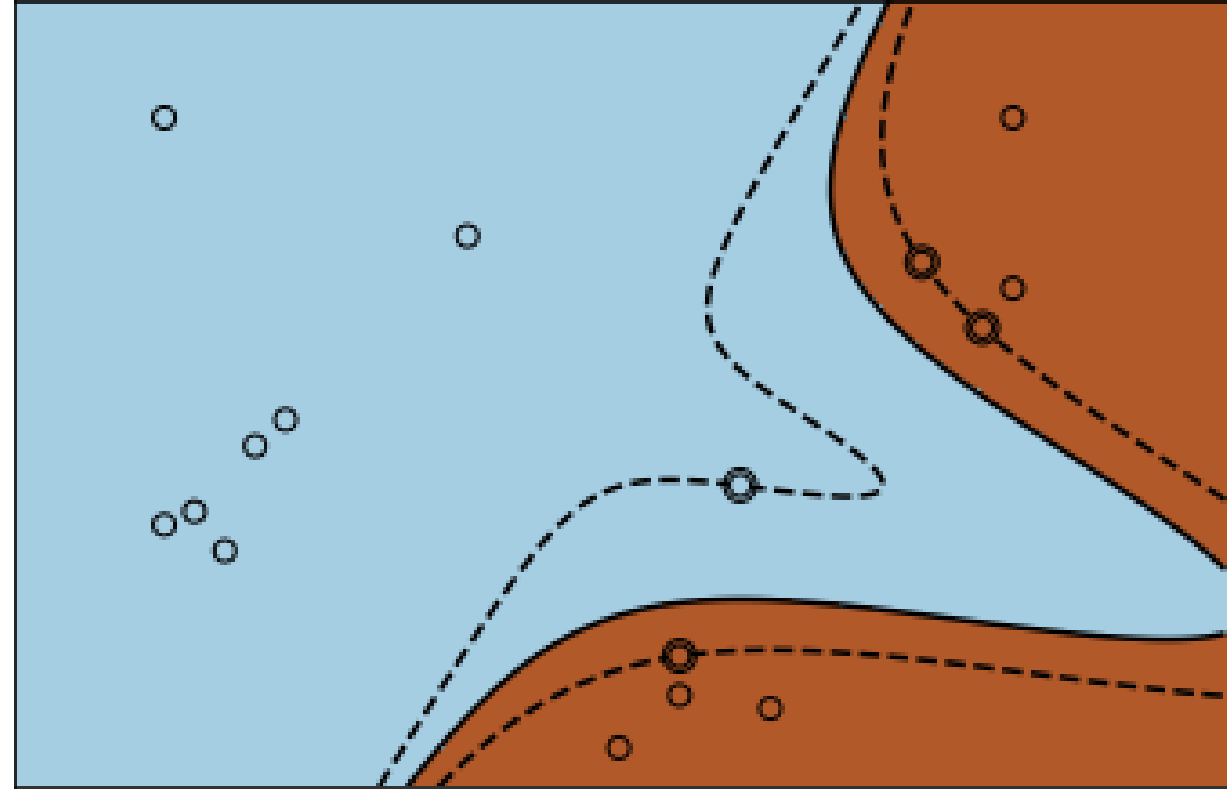
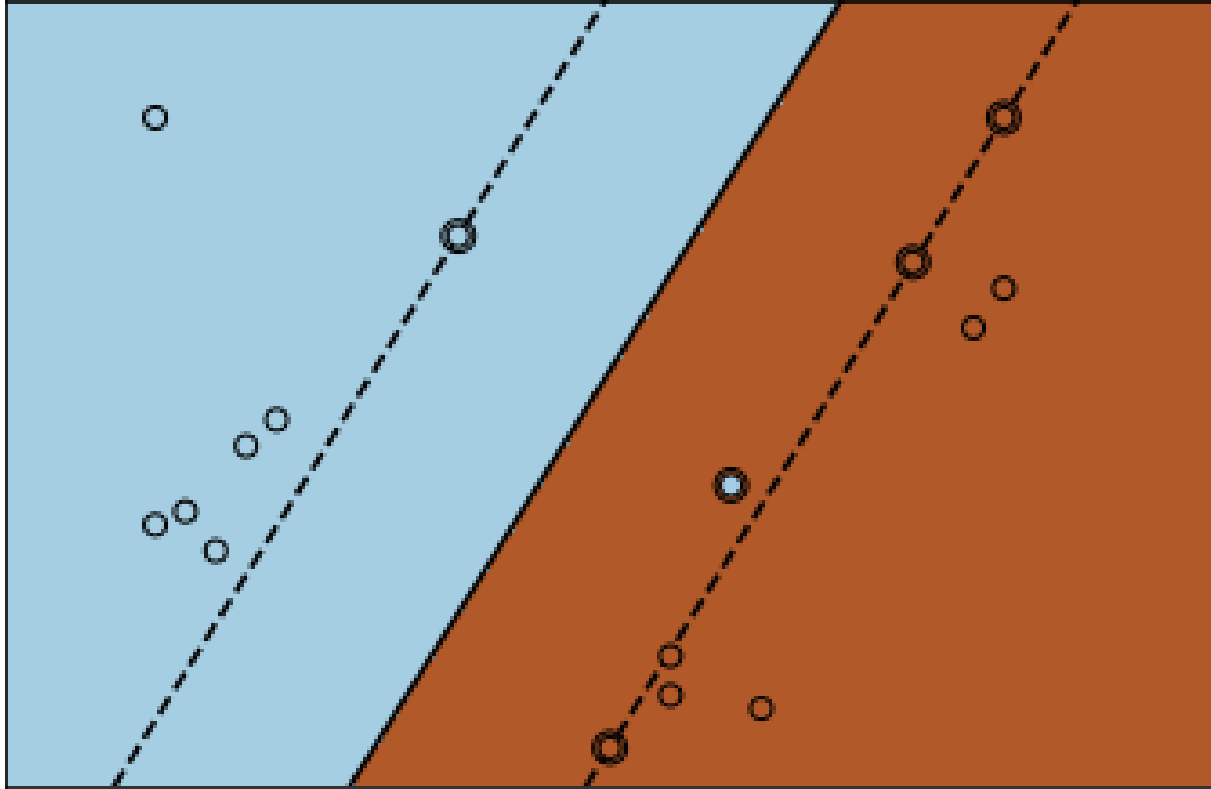


$C \uparrow$



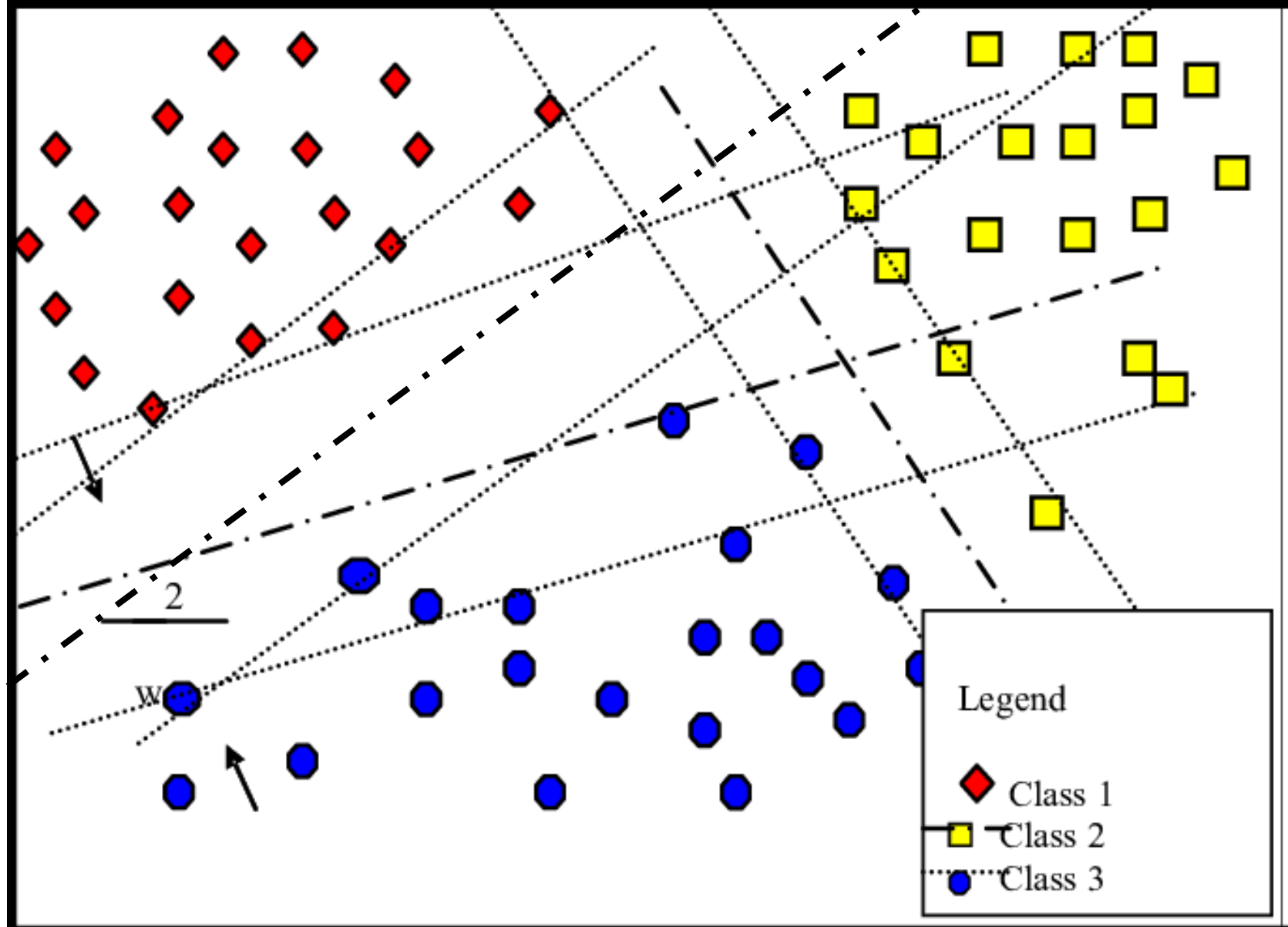
$C \downarrow$

# Higher dimensionality and generalization



$$E_{gen} \leq \frac{\# \text{ of support vectors}}{\# \text{ of training vectors}}$$

# Multi-class

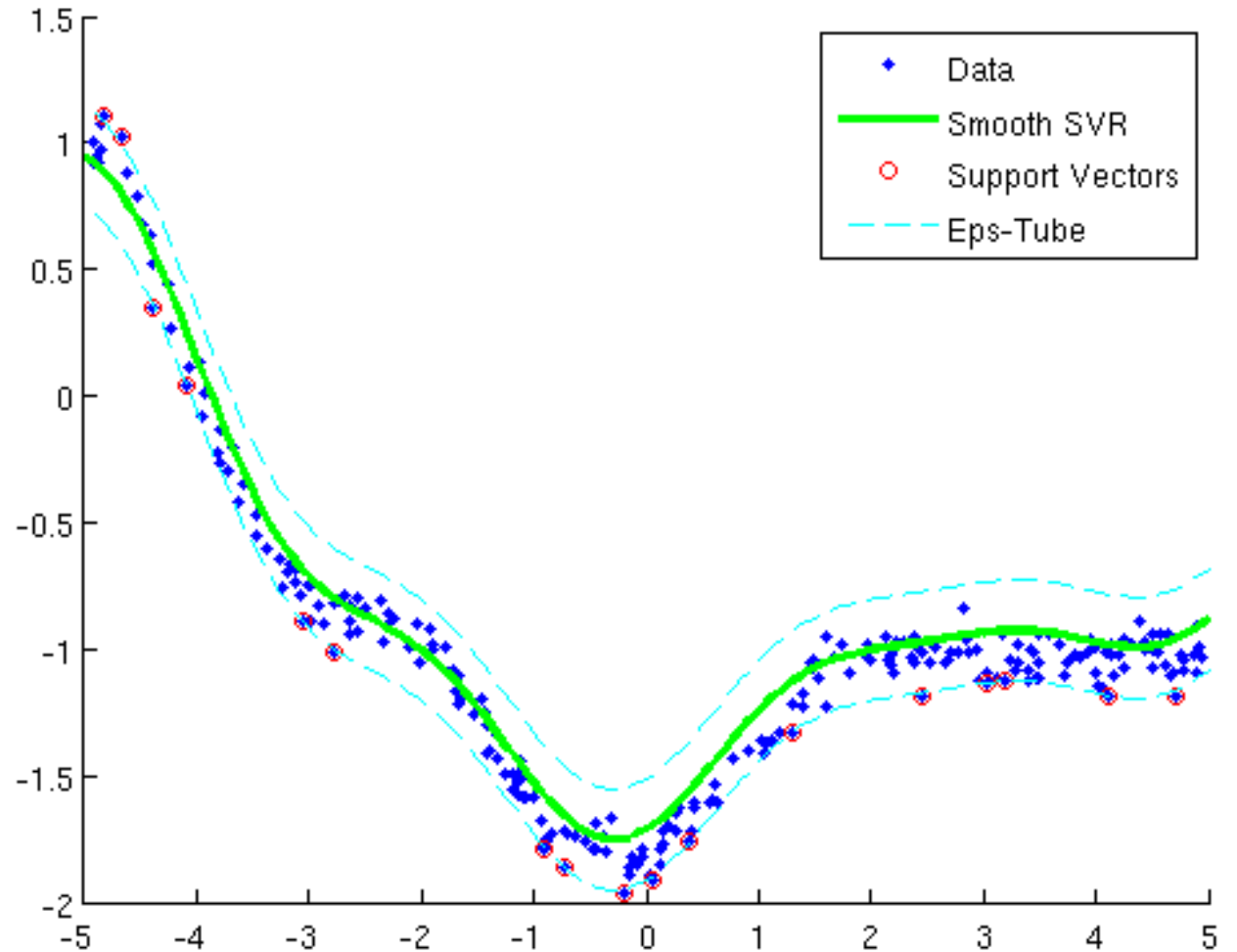


1) One-vs-all classifiers

2)  $\max_y (w_y^T x - b_y)$

# Support Vector Regression Machine

$$\begin{cases} \frac{1}{2} \|w\|^2 + C \sum (\xi_i + \xi_i^*) \rightarrow \min \\ y_i - w^T x_i - b \leq \epsilon + \xi_i \\ w^T x_i + b - y_i \leq \epsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0 \end{cases}$$



# Support Vector Networks

Cortes and Vapnik, 1995

