# Support Vector Machines
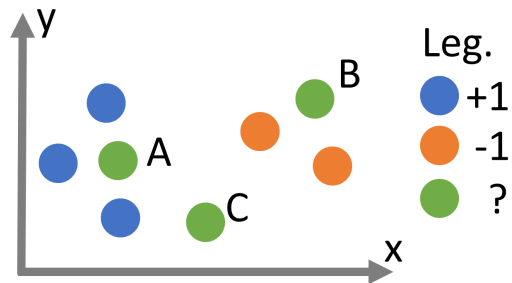## Machine Learning - ENS 2022

Thomas Romary
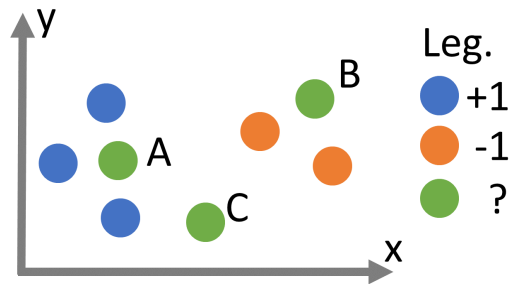
Mines Paris
(from a course of Tim Schlottmann, Hendrik Sieck, Jonas Kru)

10.26.2022

- Object classification
- Simple and efficient
- As accurate as possible

- Support Vector Machine
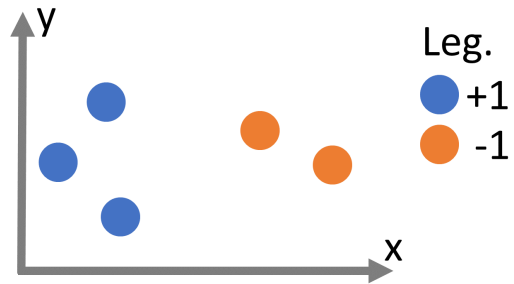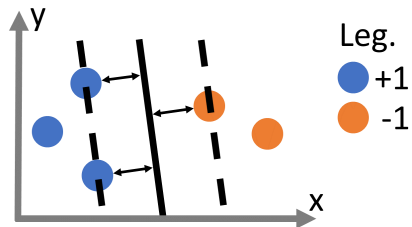- Binary classification

# Contents

▶ How do I separate the two classes?

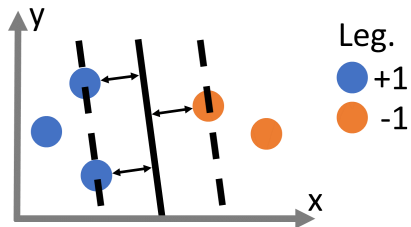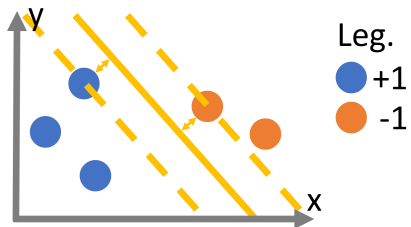- What is the best way to set the hyperplane?
- Other name of Support Vector Machines: Large Margin Classifier

# Problem



▶ What is the best way to set the hyperplane?

▶ Other name of Support Vector Machines: Large Margin Classifier

# Support vectors (s.v.)

## Definition

- $m \in \mathbb{R}$ data points
- Input $\boldsymbol{x} \in \mathbb{R}^N$
- Output $y \in \{-1, +1\}$
- Training set $S \in (\mathbb{R}^N \times \{-1, +1\})^m$

- Hypothesis

$$h : \mathbb{R}^N \to \{+1, -1\}$$
$$\boldsymbol{x} \mapsto y$$



- $m = 5$
- Training set $S$:

$$S = \begin{pmatrix} 1 & 2.5 & +1 \\ 5 & 2 & -1 \\ & \vdots & \end{pmatrix}$$

# Hyperplane



- Hyperplane: $H = \boldsymbol{w}^T \boldsymbol{x} + b = 0$
- Hyperplanes $H_+$ and $H_-$:

$$H_+ := \boldsymbol{w}^T \boldsymbol{x}_p + b = +1, \quad \forall \text{ s.v., lying on } H_+$$

$$H_- := \boldsymbol{w}^T \boldsymbol{x}_n + b = -1, \quad \forall \text{ s.v., lying on } H_-$$

$$y_i(\boldsymbol{w}^T \boldsymbol{x}_i + b) = 1 \quad \forall \text{ s.v.}$$

- Hyperplane: $H = \boldsymbol{w}^T \boldsymbol{x} + b = 0$
- Classification:

$$h(x_i) = \begin{cases} +1 & \text{when } \boldsymbol{w}^T \boldsymbol{x}_i + b \geq 0 \\ -1 & \text{when } \boldsymbol{w}^T \boldsymbol{x}_i + b \leq 0 \end{cases}$$

## Example I



Leg.
- ● +1
- ● -1

▶ Hyperplane: $H = \boldsymbol{w}^T \boldsymbol{x} + b = 0$
▶ Constraints: $\boldsymbol{w}^T \boldsymbol{x}_i + b = y_i, \quad \forall$ s. v.

▶ Graphical determination of the hyperplane parameters:

$$x = 3$$

$$\boldsymbol{w} = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} 2 & 0 \end{pmatrix}^T \begin{pmatrix} 3 \\ 0 \end{pmatrix} + b = 0 \Rightarrow b = -6$$

# Example I



- Hyperplane: $H = \boldsymbol{w}^T \boldsymbol{x} + b = 0$
- Constraints: $\boldsymbol{w}^T \boldsymbol{x}_i + b = y_i, \quad \forall \text{ s. v.}$

▶ Graphical determination of the hyperplane parameters:

$$x = 3$$

$$\boldsymbol{w} = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} 2 & 0 \end{pmatrix}^T \begin{pmatrix} 3 \\ 0 \end{pmatrix} + b = 0 \Rightarrow b = -6$$

# Example II



- Hyperplane: $H = \boldsymbol{w}^T \boldsymbol{x} + b = 0$
- Constraints: $\boldsymbol{w}^T \boldsymbol{x}_i + b = y_i, \quad \forall$ s. v.

- Consider the constraint using the example of point $A$:

$$c \left( \begin{pmatrix} 2 & 0 \end{pmatrix}^T \begin{pmatrix} 2 \\ 1 \end{pmatrix} - 6 \right) \stackrel{!}{=} +1 \Rightarrow c = -0.5$$

- Thus for the canonical hyperplane:

$$\boldsymbol{w} = \begin{pmatrix} -1 \\ 0 \end{pmatrix} \quad b = 3$$

- Control with point $D$:

## Example II



▶ Hyperplane: $H = \boldsymbol{w}^T \boldsymbol{x} + b = 0$
▶ Constraints: $\boldsymbol{w}^T \boldsymbol{x}_i + b = y_i, \quad \forall$ s. v.

▶ Consider the constraint using the example of point $A$:

$$c \left( \begin{pmatrix} 2 & 0 \end{pmatrix}^T \begin{pmatrix} 2 \\ 1 \end{pmatrix} - 6 \right) \overset{!}{=} +1 \Rightarrow c = -0.5$$

▶ Thus for the canonical hyperplane:

$$\boldsymbol{w} = \begin{pmatrix} -1 \\ 0 \end{pmatrix} \quad b = 3$$

▶ Control with point $D$:

# Example II



- Hyperplane: $H = \boldsymbol{w}^T \boldsymbol{x} + b = 0$
- Constraints: $\boldsymbol{w}^T \boldsymbol{x}_i + b = y_i, \quad \forall$ s. v.

- Consider the constraint using the example of point $A$:

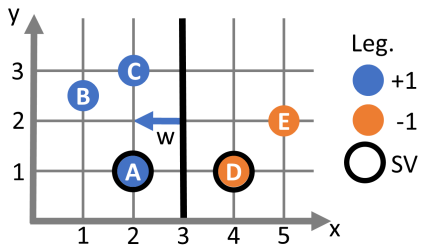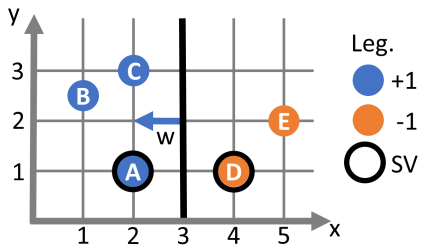$$c \left( \begin{pmatrix} 2 & 0 \end{pmatrix}^T \begin{pmatrix} 2 \\ 1 \end{pmatrix} - 6 \right) \overset{!}{=} +1 \Rightarrow c = -0.5$$

- Thus for the canonical hyperplane:

$$\boldsymbol{w} = \begin{pmatrix} -1 \\ 0 \end{pmatrix} \quad b = 3$$

- Control with point $D$:

# Example III



- Classification:

$$h(x_i) = \begin{cases} +1 & \text{if } \boldsymbol{w}^T \boldsymbol{x}_i + b \geq 0 \\ -1 & \text{if } \boldsymbol{w}^T \boldsymbol{x}_i + b \leq 0 \end{cases}$$

- Parameters of the canonical hyperplane

$$\boldsymbol{w} = \begin{pmatrix} -1 \\ 0 \end{pmatrix} \quad b = 3$$

- Classification of point $F$:

$$\begin{pmatrix} -1 & 0 \end{pmatrix}^T \begin{pmatrix} 5 \\ 3 \end{pmatrix} + 3 = -2$$

- So point $F$ belongs to the class $-1$.

# Minimization problem



- Projection property of the scalar product
- Width of the margin $\rho$:

$$\rho = (\mathbf{x}_p - \mathbf{x}_n)^T \frac{\mathbf{w}}{\|\mathbf{w}\|} \Leftrightarrow \rho = (\mathbf{x}_p^T \mathbf{w} - \mathbf{x}_n^T \mathbf{w}) \frac{1}{\|\mathbf{w}\|}$$

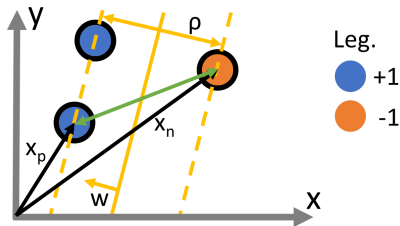- Constraint: $y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1 \quad \forall \text{ s.v.}$
- $\rho = \frac{2}{\|\mathbf{w}\|}$
- Goal of an SVM: Maximize the margin

$$\max_{\mathbf{w},b} \frac{2}{\|\mathbf{w}\|} \Leftrightarrow \min_{\mathbf{w},b} \|\mathbf{w}\| \Leftrightarrow \min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|^2$$

with the constraint $y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0$

## Dual problem: Lagrange Multipliers

▶ Lagrange Multipliers:

$$L = \frac{1}{2}\|\mathbf{w}\|^2 - \sum \alpha_i \left[ y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \right]$$

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum \alpha_i y_i \mathbf{x}_i = 0 \qquad\qquad \Rightarrow \mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L}{\partial b} = -\sum \alpha_i y_i = 0 \qquad\qquad \Rightarrow \sum_i \alpha_i y_i = 0$$

▶ Replace $\mathbf{w}$ in $L$:

$$L = \frac{1}{2}(\sum_i \alpha_i y_i \mathbf{x}_i)^T (\sum_j \alpha_j y_j \mathbf{x}_j) - (\sum_i \alpha_i y_i \mathbf{x}_i)^T (\sum_j \alpha_j y_j \mathbf{x}_j) - \sum_i \alpha_i y_i b + \sum_i \alpha_i$$

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

▶ Target: $\max_{\boldsymbol{\alpha}} L$
▶ Decision function:

$$h(\mathbf{x}) = \begin{cases} +1 & \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \geq 0 \\ -1 & \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \leq 0 \end{cases}$$

▶ Characteristics:

$$\boldsymbol{\alpha} \geq 0$$

$$\sum_{\substack{\text{pos. s.v.} \\ p}} \alpha_p = \sum_{\substack{\text{neg. s.v.} \\ n}} \alpha_n$$

$$\alpha_i \begin{cases} > 0 & \text{when } x_i \text{ is a support vector} \\ = 0 & \text{otherwise} \end{cases}$$

# Support variables $\alpha$

- Characteristics:

$$\boldsymbol{\alpha} \geq 0$$

$$\sum_{\substack{\text{pos. s.v.} \\ p}} \alpha_p = \sum_{\substack{\text{neg. s.v.} \\ n}} \alpha_n$$

$$\alpha_i \begin{cases} > 0 & \text{when } x_i \text{ is a support vector} \\ = 0 & \text{otherwise} \end{cases}$$
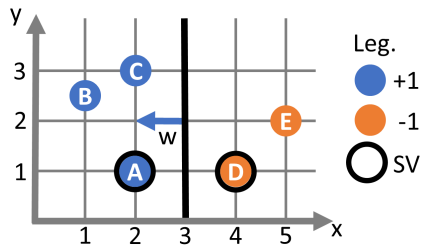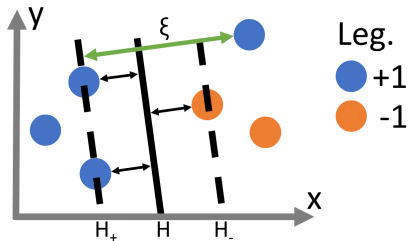


- Computation of $\alpha_A$ and $\alpha_D$ results:

$$\alpha_A = 0.5$$
$$\alpha_D = 0.5$$

- Computation of $\alpha_C$ result:

$$\alpha_C = 0$$

# Compensation for a faulty training set



- ▶ Relaxation variable $\xi$ to compensate for false classification
- ▶ It tolerates some outliers in the classification
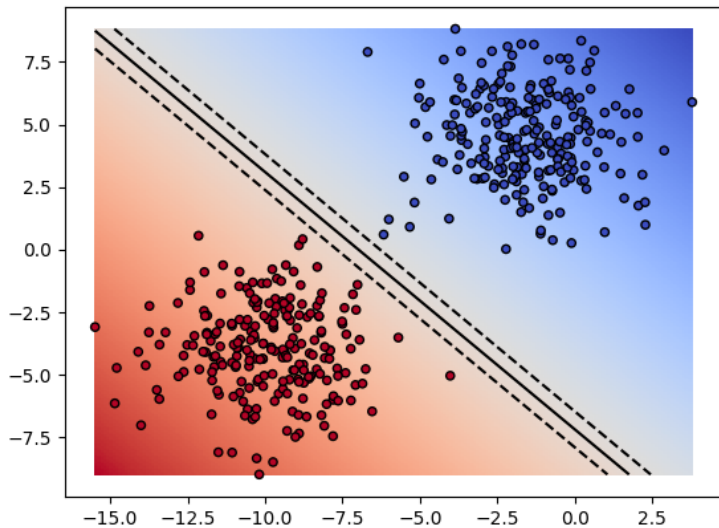
▶ Constraints:

$$y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall \mathbf{x}_i \in S \qquad\qquad \xi_i \geq 0$$
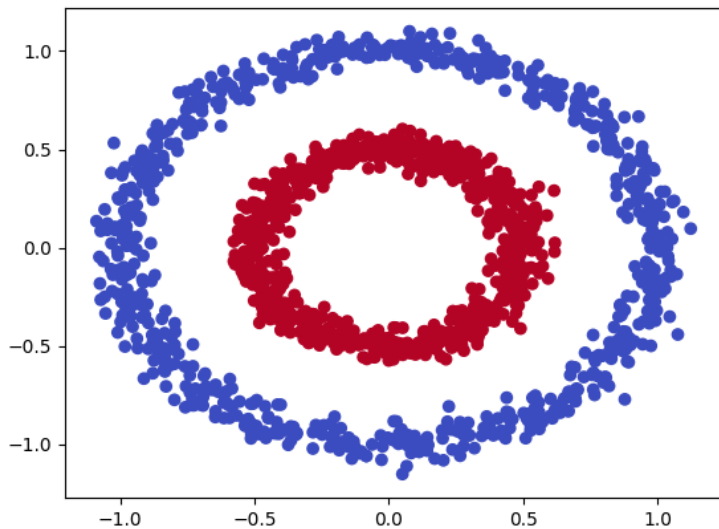
▶ Minimization problem:

$$\min_{\mathbf{w},b,\boldsymbol{\xi}} \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_i \xi_i$$

with the constraint $y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i$ and $\xi_i \geq 0$

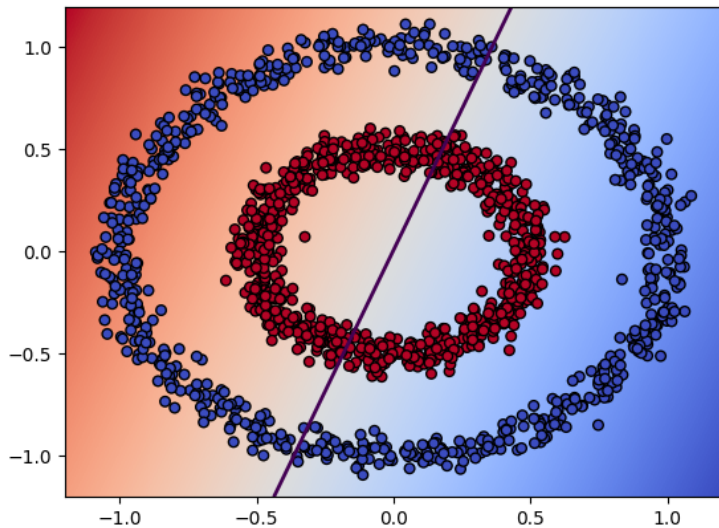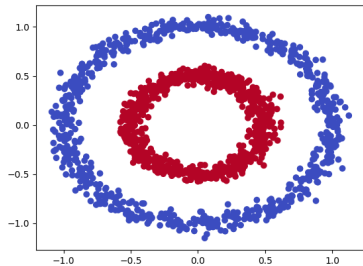# Introduction: Linear separable data

$$\phi(\boldsymbol{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{bmatrix}$$

$$\phi(\boldsymbol{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{bmatrix}$$

# Approach: image function $\phi(x)$

# Change in dual problem

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

$$h(\mathbf{x}) = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

$$h(\mathbf{x}) = \sum_i \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + b$$

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

$$h(\mathbf{x}) = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

$$h(\mathbf{x}) = \sum_i \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + b$$

Image function

$$\phi : \mathbb{R}^n \to \mathbb{R}^m$$
$$\boldsymbol{x} \mapsto \boldsymbol{f}$$

After transformation, the data are linearly separable
Problem:

- $m > n$ high computational effort
  Above a certain size, it's difficult to use

- Moreover, note that $\phi$ is only needed for scalar product

# Definition: $\phi$

Image function

$$\phi : \mathbb{R}^n \to \mathbb{R}^m$$
$$\boldsymbol{x} \mapsto \boldsymbol{f}$$

After transformation, the data are linearly separable
Problem:

- $m > n$ high computational effort
  Above a certain size, it's difficult to use
- Moreover, note that $\phi$ is only needed for scalar product

# Kernel example

$$\phi(\boldsymbol{x}) = (1 \ \sqrt{2}x_1 \ \sqrt{2}x_2 \ \ldots \ x_1^2 \ x_2^2 \ \ldots \ \sqrt{2}x_1x_2 \ \sqrt{2}x_1x_3 \ \ldots)$$

$$\phi(\boldsymbol{v})^T \phi(\boldsymbol{w}) = \sum_j 2v_j w_j + \sum_j v_j^2 w_j^2 + \sum_j \sum_{k>j} 2v_j v_k w_j w_k + \ldots$$

$$= (1 + \sum_j v_j w_j)^2$$

$$= (1 + \boldsymbol{v}^T \boldsymbol{w})^2$$

$$= K(\boldsymbol{v}, \boldsymbol{w})$$

$$K(\mathbf{v}, \mathbf{w}) = \phi(\mathbf{v})^T \phi(\mathbf{w})$$
$$K : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$$

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$h(\mathbf{x}) = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

$$K(\mathbf{v}, \mathbf{w}) = \phi(\mathbf{v})^T \phi(\mathbf{w})$$

$$K : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$$

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$h(\mathbf{x}) = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

$\{\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(m)}\}$
$\mathcal{K}_{i,j} = K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$

$$
\begin{aligned}
\mathcal{K}_{i,j} &= K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \\
&= \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)}) \\
&= \phi(\mathbf{x}^{(j)})^T \phi(\mathbf{x}^{(i)}) \\
&= K(\mathbf{x}^{(j)}, \mathbf{x}^{(i)}) \\
&= \mathcal{K}_{j,i}
\end{aligned}
$$

$\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$
$\mathcal{K}_{i,j} = K(\boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)})$

$$\begin{aligned}
\mathcal{K}_{i,j} &= K(\boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)}) \\
&= \phi(\boldsymbol{x}^{(i)})^T \phi(\boldsymbol{x}^{(j)}) \\
&= \phi(\boldsymbol{x}^{(j)})^T \phi(\boldsymbol{x}^{(i)}) \\
&= K(\boldsymbol{x}^{(j)}, \boldsymbol{x}^{(i)}) \\
&= \mathcal{K}_{j,i}
\end{aligned}$$

# Mercer's Theorem: 2

$\{\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(m)}\}$

$\mathcal{K}_{i,j} = K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$
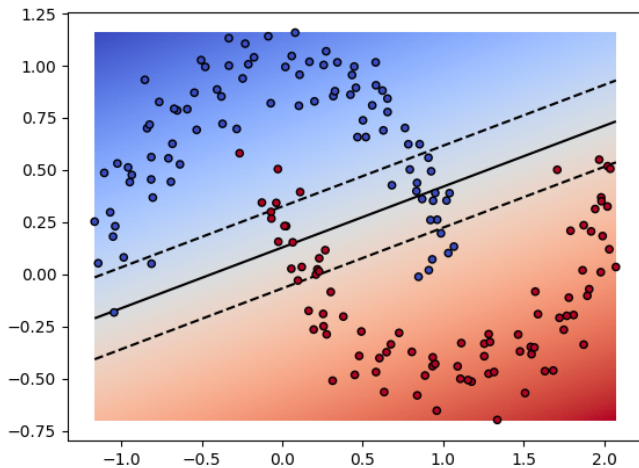
Choose any $\mathbf{z}$:

$$
\begin{aligned}
\mathbf{z}^T \mathcal{K} \mathbf{z} &= \sum_i \sum_j \mathbf{z}_i \mathcal{K}_{i,j} \mathbf{z}_j \\
&= \sum_i \sum_j \mathbf{z}_i \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)}) \mathbf{z}_j \\
&= \sum_i \sum_j \mathbf{z}_i \sum_k \phi_k(\mathbf{x}^{(i)}) \phi_k(\mathbf{x}^{(j)}) \mathbf{z}_j \\
&= \sum_k \sum_i \sum_j \mathbf{z}_i \phi_k(\mathbf{x}^{(i)}) \phi_k(\mathbf{x}^{(j)}) \mathbf{z}_j \\
&= \sum_k \left( \sum_i \mathbf{z}_i \phi_k(\mathbf{x}^{(i)}) \right)^2 \\
&\geq 0
\end{aligned}
$$

# Different kernels in practice

- ▶ Linear Kernel
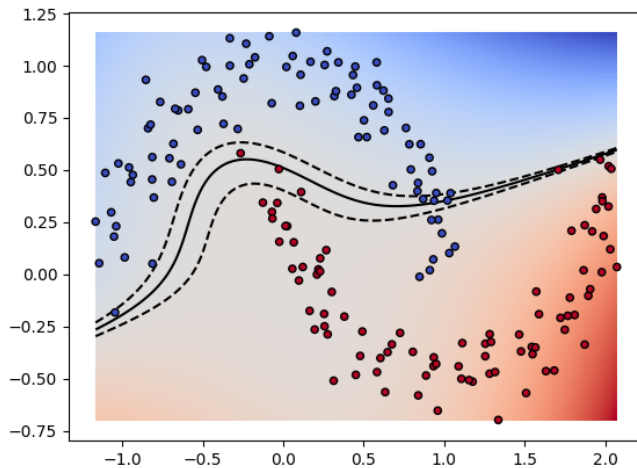- ▶ Polynomial Kernel
- ▶ Gaussian Kernel

# Different kernels in practice: Linear

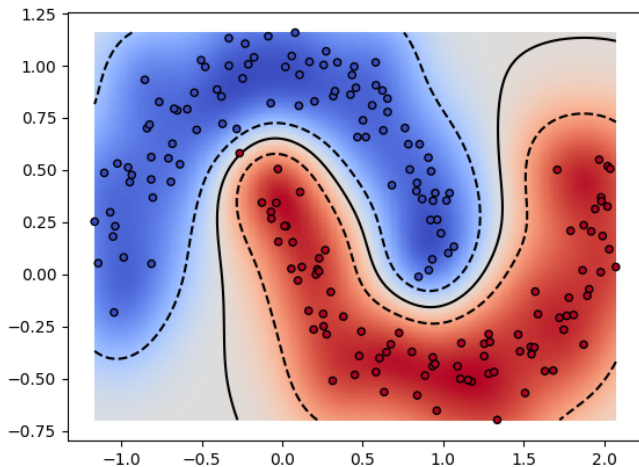$$K(\mathbf{v}, \mathbf{w}) = \mathbf{v}^T \mathbf{w}$$

# Different kernels in practice: Polynomial

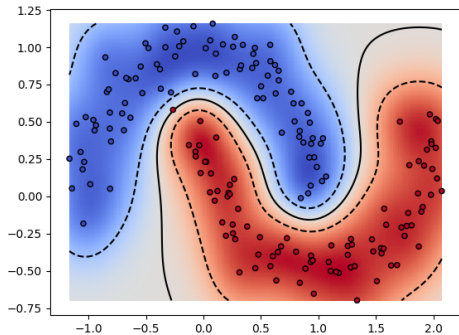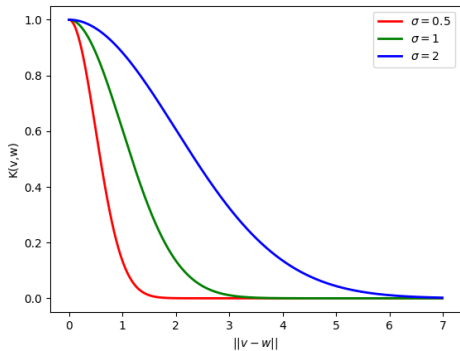$$K(\mathbf{v}, \mathbf{w}) = (\mathbf{v}^T \mathbf{w} + c)^d$$

# Different kernels in practice: Gaussian

$$K(\mathbf{v}, \mathbf{w}) = \exp\left(-\frac{||\mathbf{v}-\mathbf{w}||^2}{2\sigma^2}\right)$$

# Different kernels in practice: Gaussian

$$K(\mathbf{v}, \mathbf{w}) = \exp\left(-\frac{||\mathbf{v} - \mathbf{w}||^2}{2\sigma^2}\right)$$

# Summary

- SVMs in their standard form have problems classifying non-linearly separable datasets
- Use $\phi(x)$ to map data into a space where this is possible
- Use the kernel to simplify the resulting calculation