



Ciência de Dados Quântica  
2021/22

# Kernel Based Methods and Feature Maps

LUÍS PAULO SANTOS



# Material de Consulta

2

- ▶ [Schuld2021] – Secs. 2.5.4, 3.6.1; Chap. 6
- ▶ “Support Vector Machines: All you need to know!”  
<https://youtu.be/ny1iZ5A8ilA>
- ▶ “The Kernel Trick in Support Vector Machine (SVM)”  
<https://youtu.be/Q7vT0--5VII>
- ▶ “Supervised learning with quantum enhanced feature spaces”  
Vojtech Havlicek et al.  
<https://arxiv.org/pdf/1804.11326.pdf>
- ▶ Qiskit Global Summer School 2021: Lectures 6.1, 6.2  
[https://www.youtube.com/watch?v=xgA4Dx\\_7q34&list=PLOFEBzvs-VvqJwybFxxkTiDzhf5E11p8BI](https://www.youtube.com/watch?v=xgA4Dx_7q34&list=PLOFEBzvs-VvqJwybFxxkTiDzhf5E11p8BI)



# Kernel Methods: concept

3

- ▶ Kernel methods are based on a **similarity measure** between data points
- ▶ **Definition:** for a data domain  $\mathcal{X}$  a kernel is a positive semi-definite bivariate function  $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ 
  - ▶ positive semi-definite means:
    - ▶  $\kappa(x, x') \geq 0$
    - ▶  $\kappa(x, x') = \kappa(x', x)^*$



# Support Vector Machines (SVM)

4

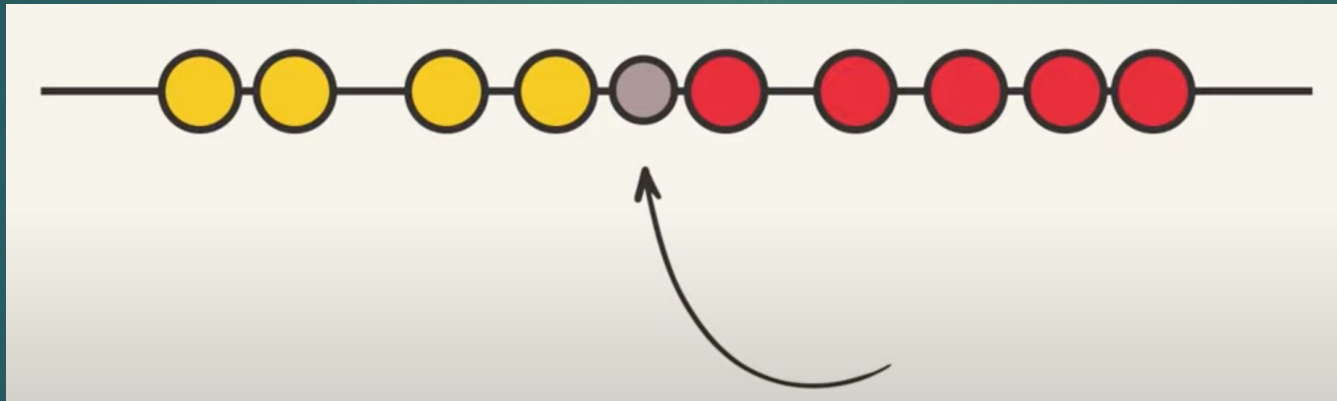
- ▶ Family of kernel based methods for classification
  - ▶ we focus on binary classification, but multi-class (SVC) and regression (SVR) also possible
- ▶ Supervised method which finds an hyperplane separating the classes
  - ▶ training input ( $M$  labelled data points):  $\mathcal{D} = \{(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^M, y^M)\}$
  - ▶ each  $\mathbf{x}^i \in \mathbb{R}^N$ ,  $\mathbf{x}^i = (x_1^i, \dots, x_N^i)^T$  for  $N$  features
  - ▶ each  $y^i \in \{-1, 1\}$  for binary classification
- ▶ The training stage finds the class separating hyperplane
- ▶ The classification stage classifies a previously unseen unlabelled data point  $\mathbf{u} \in \mathbb{R}^N$  as  $y^u \in \{-1, 1\}$



# Support Vector Machines (SVM)

5

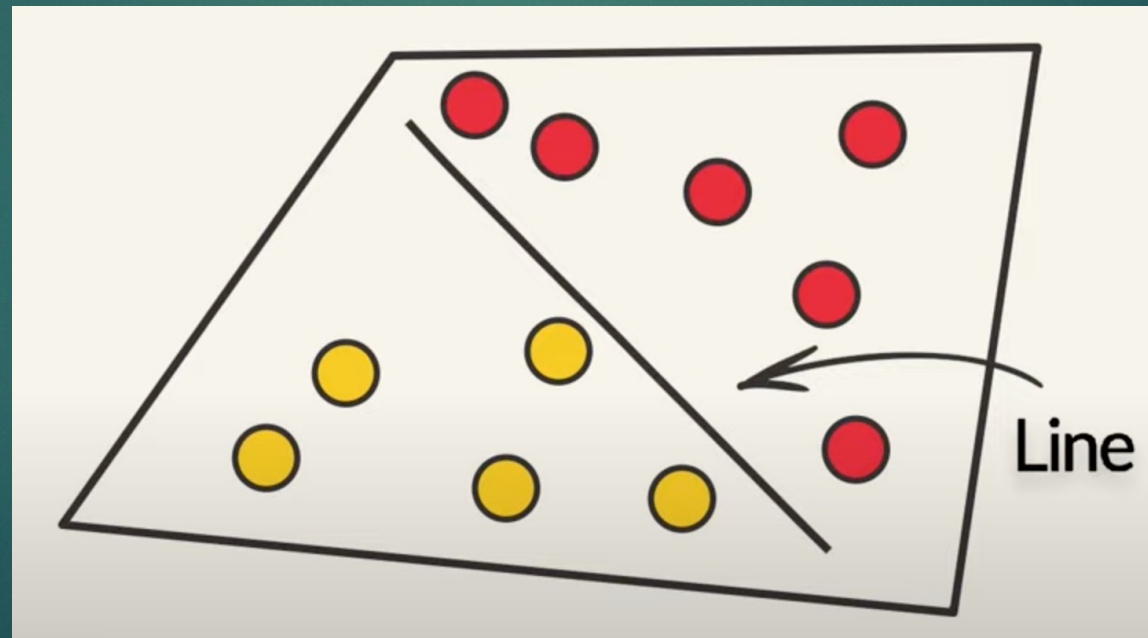
- ▶ The hyperplane is defined by a vector  $\mathbf{w} \in \mathbb{R}^N$  with the same dimensionality  $N$  as the data points
- ▶ 1D – The hyperplane is a point



# Support Vector Machines (SVM)

6

- ▶ The hyperplane is defined by a vector  $\mathbf{w} \in \mathbb{R}^N$  with the same dimensionality  $N$  as the data points
- ▶ 2D – The hyperplane is a line

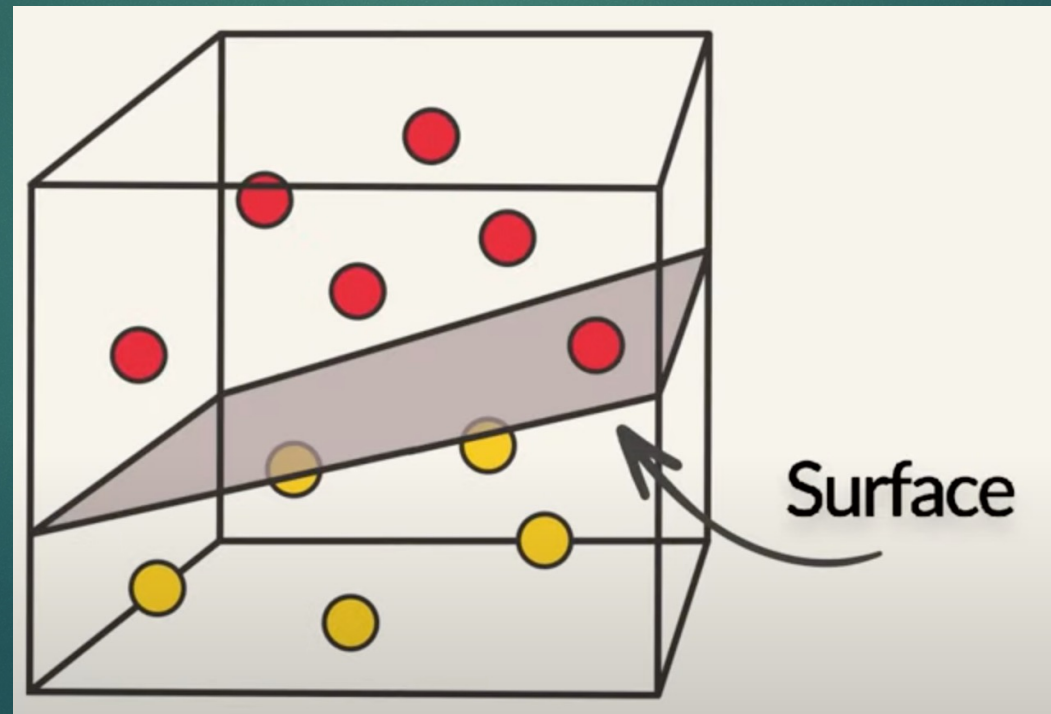




# Support Vector Machines (SVM)

7

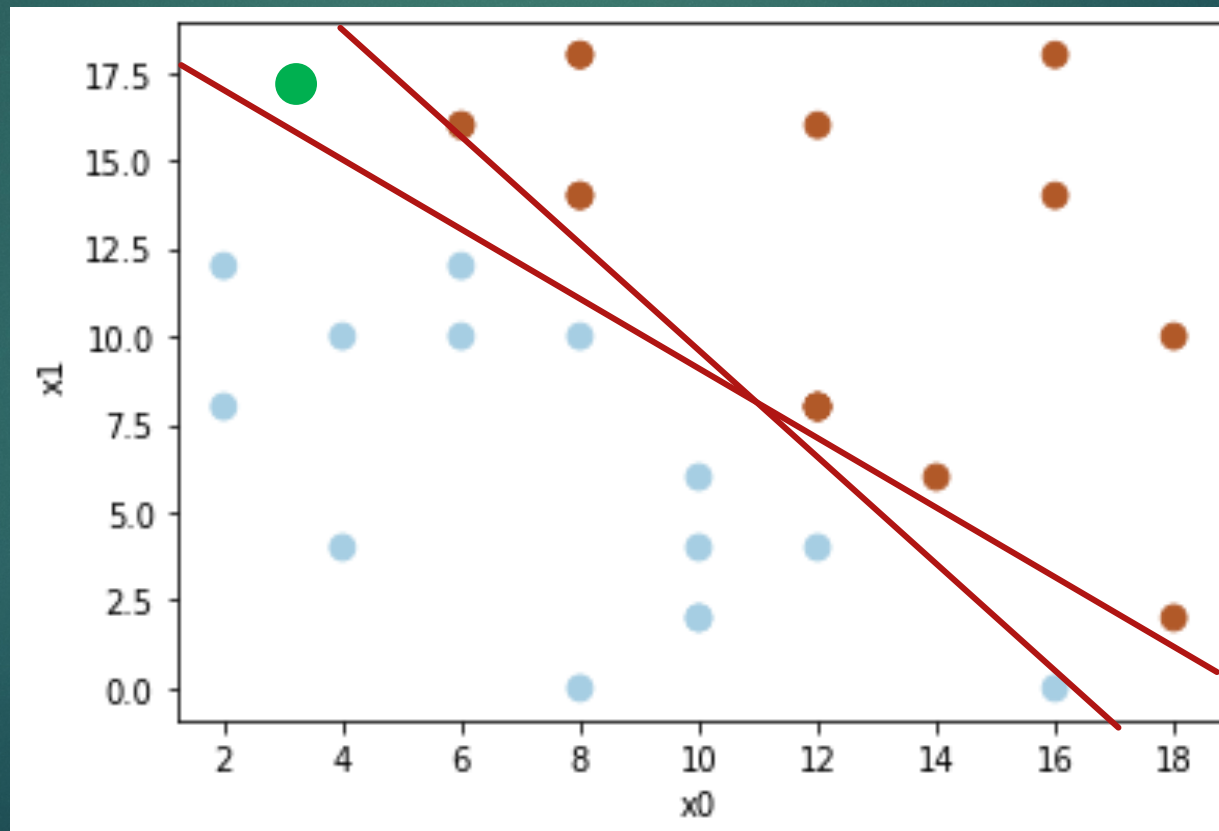
- ▶ The hyperplane is defined by a vector  $\mathbf{w} \in \mathbb{R}^N$  with the same dimensionality  $N$  as the data points
- ▶ 3D – The hyperplane is a plane



# Support Vector Machines (SVM)

8

- There are infinite hyperplanes separating the two classes. Which hyperplane to select?

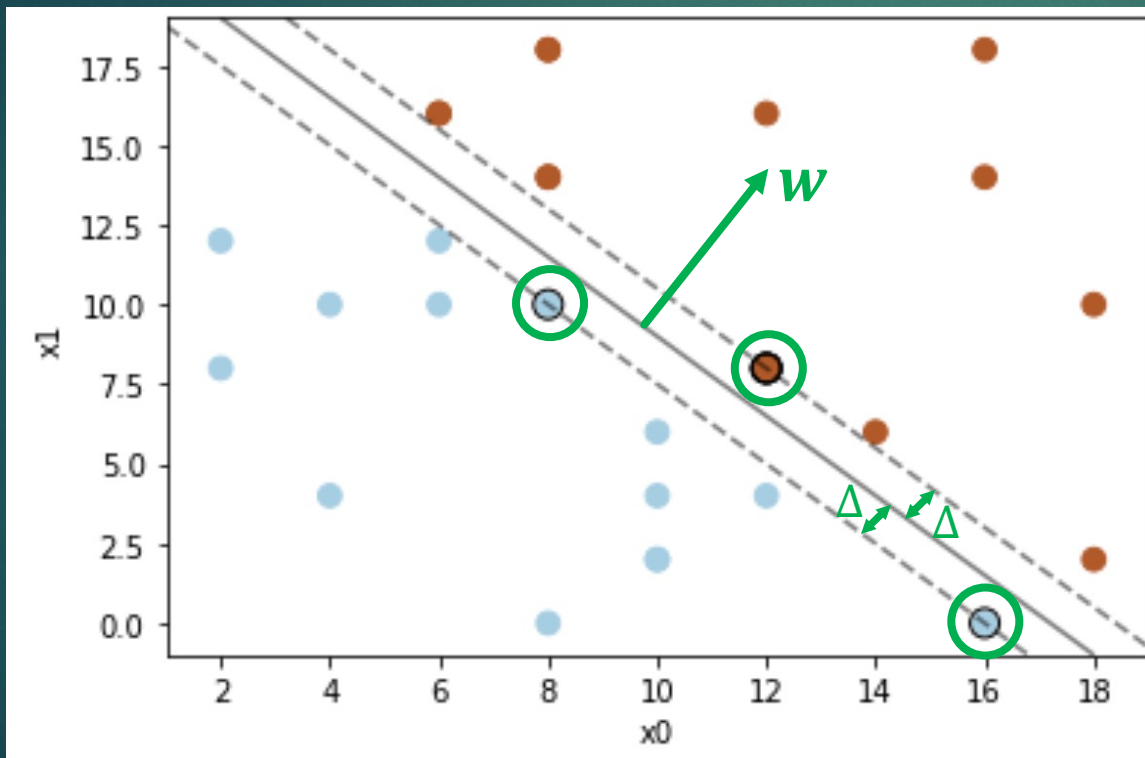




# Support Vector Machines (SVM)

9

- ▶ Select the hyperplane which maximizes the “margin”, i.e, the distance to the nearest points on each class.
- ▶ These are the support vectors

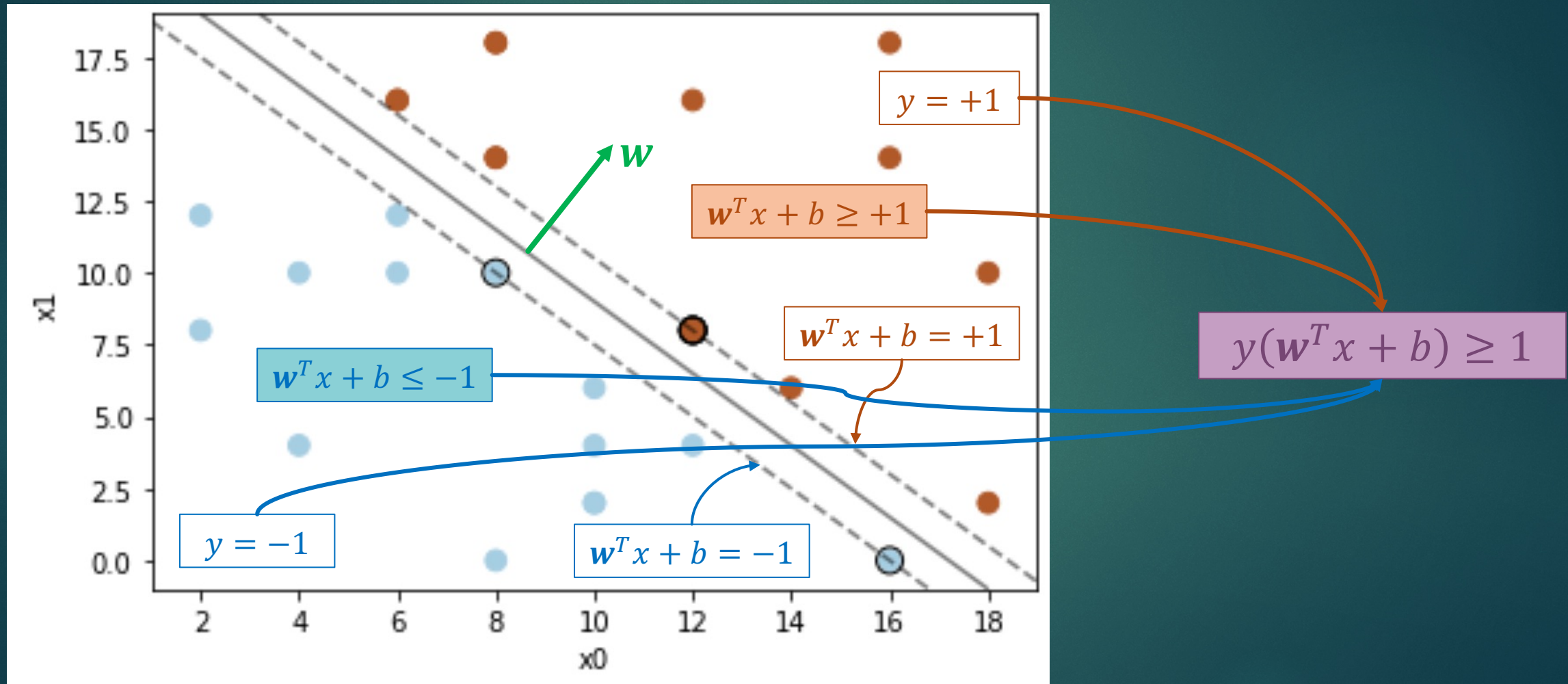


- ▶ The vector  $\mathbf{w} \in \mathbb{R}^N$  is perpendicular to the hyperplane
- ▶ The points in the hyperplane obey
$$f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$$
- ▶ The training stage optimizes over  $\mathbf{w}$  and  $b$ , maximizing the margin  $\Delta$



# Support Vector Machines (SVM)

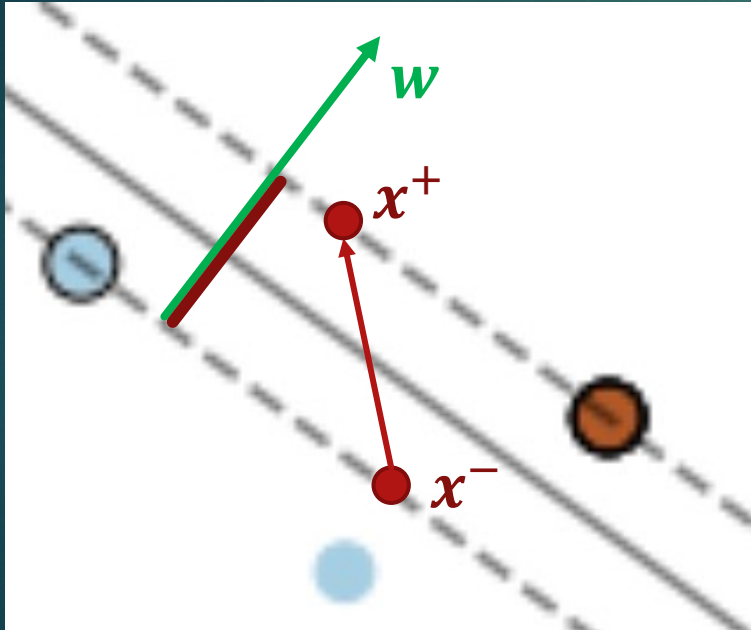
10





# Support Vector Machines (SVM)

11



- ▶  $\max_{w,b} \left( (x^+ - x^-) \frac{w}{\|w\|} \right) \text{ s.t. } y^m (w^T x^m + b) \geq 1$
- ▶  $\max_{w,b} \left( \frac{x^+ w - x^- w}{\|w\|} \right) \text{ s.t. } y^m (w^T x^m + b) \geq 1$
- ▶  $\max_{w,b} \left( \frac{(1-b) - (-1-b)}{\|w\|} = \frac{2}{\|w\|} \right) \text{ s.t. } y^m (w^T x^m + b) \geq 1$
- ▶  $\min_{w,b} \left( \frac{1}{2} \|w\|^2 \right) \text{ s.t. } y^m (w^T x^m + b) \geq 1$



# Support Vector Machines (SVM)

12

►  $\min_{\mathbf{w}, b} \left( \frac{1}{2} \|\mathbf{w}\|^2 \right) \text{ s.t. } y^m (\mathbf{w}^T \mathbf{x}^m + b) \geq 1$

Depends on  $\mathbf{w}$ ,  
thus  $\mathcal{O}(NM)$

►  $\min_{\mathbf{w}, b, \alpha} \left[ \mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^2 - \sum_{m=1}^M (\alpha^m y^m (\mathbf{w}^T \mathbf{x}^m + b) - 1) \right]$

►  $\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \sum_{m=1}^M (\alpha^m y^m \mathbf{x}^m) = \mathbf{0} \Leftrightarrow \mathbf{w} = \sum_{m=1}^M (\alpha^m y^m \mathbf{x}^m)$

►  $\frac{\partial \mathcal{L}}{\partial b} = \sum_{m=1}^M (\alpha^m y^m) = 0$

Doesn't depend on  $\mathbf{w}$ ,  
thus  $\mathcal{O}(M)$

►  $\min_{\alpha} \left( \mathcal{L}(\alpha) = \sum_{m=1}^M \alpha^m - \frac{1}{2} \sum_{m, m'=1}^M \alpha^m \alpha^{m'} y^m y^{m'} \mathbf{x}^{mT} \mathbf{x}^{m'} \right)$

s.t.  $\sum_{m=1}^M \alpha^m y^m = 0$

$\mathbf{x}^{mT} \mathbf{x}^{m'}$   
 $\langle \mathbf{x}^m, \mathbf{x}^{m'} \rangle$



# Support Vector Machines (SVM)

13

- ▶  $\min_{\alpha} \left( \mathcal{L}(\alpha) = \sum_{m=1}^M \alpha^m - \frac{1}{2} \sum_{m,m'=1}^M \alpha^m \alpha^{m'} y^m y^{m'} \langle \mathbf{x}^m, \mathbf{x}^{m'} \rangle \right) \text{ s.t. } \sum_{m=1}^M \alpha^m y^m = 0$
- ▶ Let  $\kappa(\mathbf{x}^m, \mathbf{x}^{m'}) = \langle \mathbf{x}^m, \mathbf{x}^{m'} \rangle$  then:
  - ▶  $\min_{\alpha} \left( \mathcal{L}(\alpha) = \sum_{m=1}^M \alpha^m - \frac{1}{2} \sum_{m,m'=1}^M \alpha^m \alpha^{m'} y^m y^{m'} \kappa(\mathbf{x}^m, \mathbf{x}^{m'}) \right) \text{ s.t. } \sum_{m=1}^M \alpha^m y^m = 0$
- ▶  $\mathbf{w} = \sum_{m=1}^M (\alpha^m y^m \mathbf{x}^m)$ 
  - ▶  $\alpha^m \neq 0$  only for the support vectors, Let  $S$  be such set  $S = \{s = \alpha^m \mid \alpha^m \neq 0, m = 1 \cdots M\}$
- ▶ Classification of an unseen point  $\mathbf{u} \in \mathbb{R}^N$ :

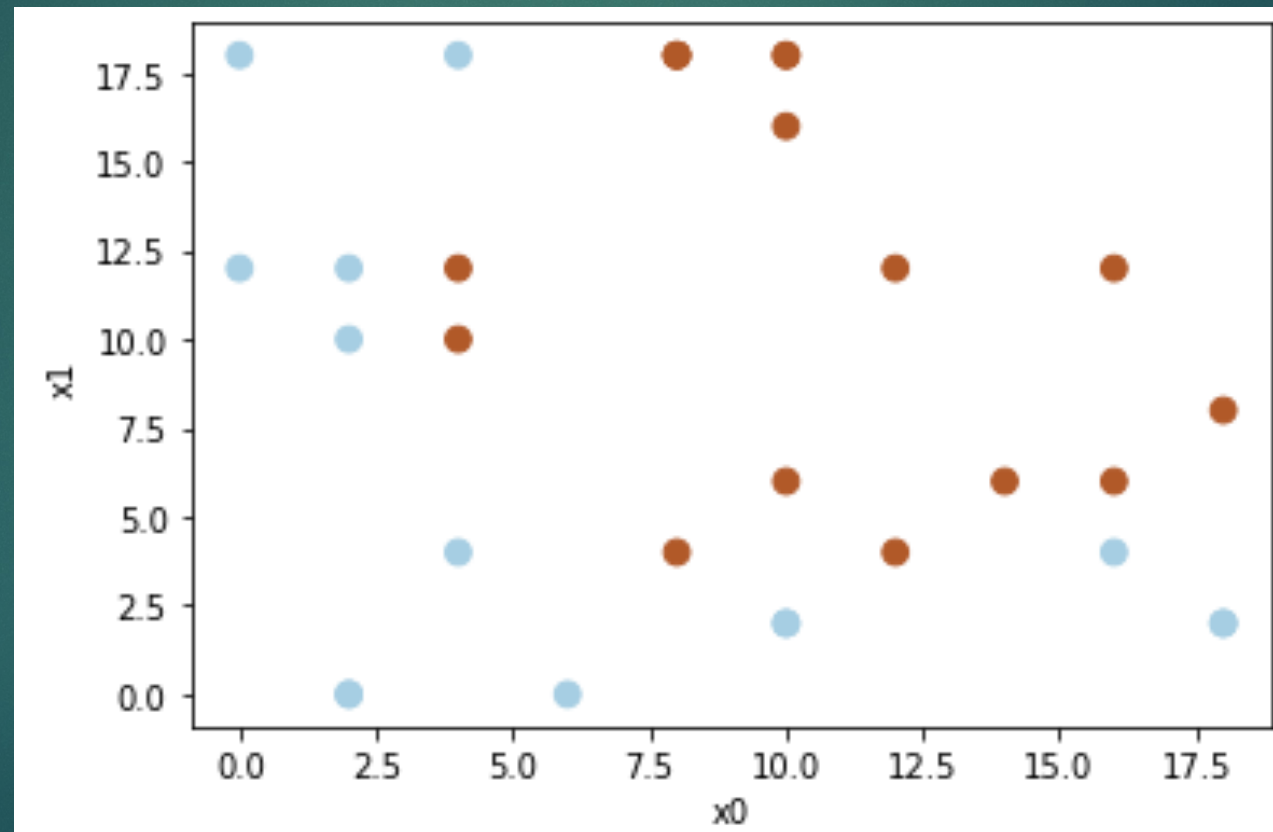
$$f_{\mathbf{w}}(\mathbf{u}) = \text{sign} \left[ \left( \sum_{s \in S} (\alpha^s y^s \mathbf{x}^s) \right)^T \mathbf{u} + b \right] = \text{sign}[\kappa(\mathbf{w}, \mathbf{u}) + b]$$



# Feature Maps

14

- What if the data is not linearly separable on its original domain?



# Feature Maps

15

- ▶ A feature map  $\phi(\mathbf{x}): \mathbb{R}^N \rightarrow \mathbb{R}^F$  is a non-linear transformation of the data, which increases its dimensionality ( $F > N$ )
- ▶ The goal is to reach linearly separability in the higher dimensional feature space
- ▶ The SVM training equations become:

- ▶  $\min_{\alpha} \left( \mathcal{L}(\alpha) = \sum_{m=1}^M \alpha^m - \frac{1}{2} \sum_{m,m'=1}^M \alpha^m \alpha^{m'} y^m y^{m'} \langle \phi(\mathbf{x}^m), \phi(\mathbf{x}^{m'}) \rangle \right) \text{ s.t. } \sum_{m=1}^M \alpha^m y^m = 0$

- ▶  $\mathbf{w} = \sum_{s \in S} (\alpha^s y^s \phi(\mathbf{x}^s))$

- ▶ Classification of an unseen point  $\mathbf{u} \in \mathbb{R}^N$ :

$$f_{\mathbf{w}}(\mathbf{u}) = \text{sign} \left[ \left( \sum_{s \in S} (\alpha^s y^s \phi(\mathbf{x}^s)) \right)^T \phi(\mathbf{u}) + b \right] = \text{sign}[\langle \phi(\mathbf{w}), \phi(\mathbf{u}) \rangle + b]$$



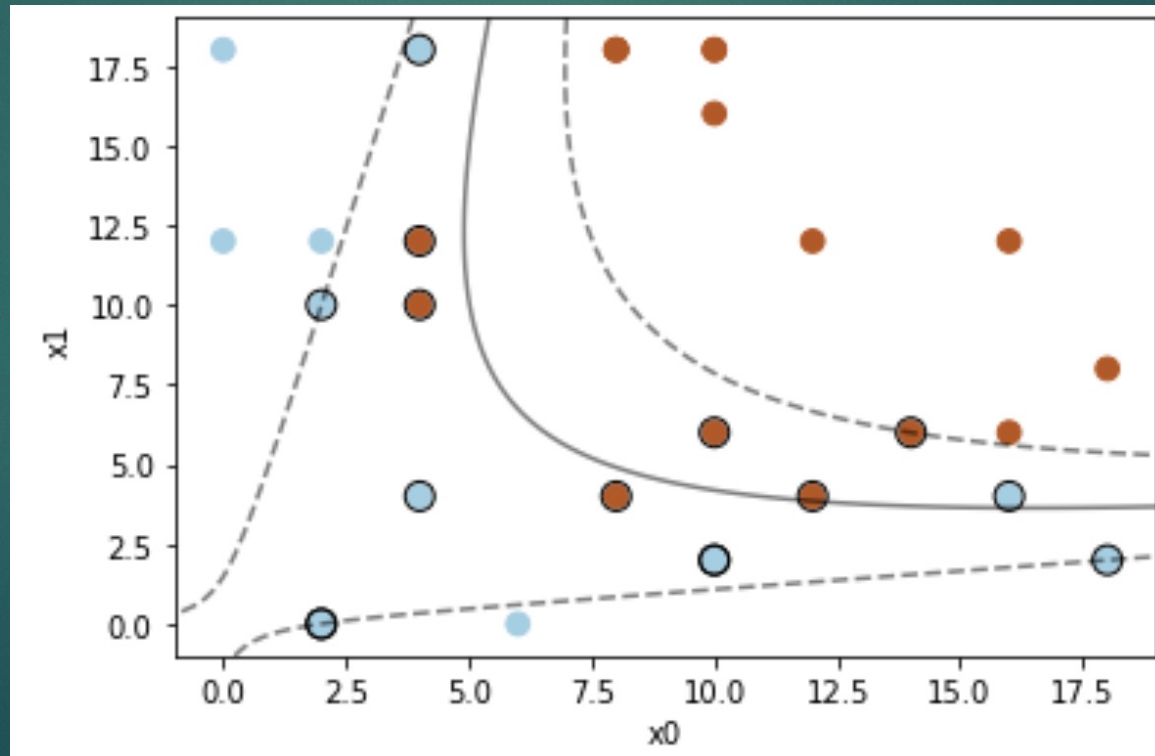
# Feature Maps

16

- Polynomial feature map (N=2, degree=2, F=6)

$$\phi(\mathbf{x}) : (x_0, x_1) \rightarrow (1, \sqrt{2}x_0, \sqrt{2}x_1, \sqrt{2}x_0x_1, x_0^2, x_1^2)$$

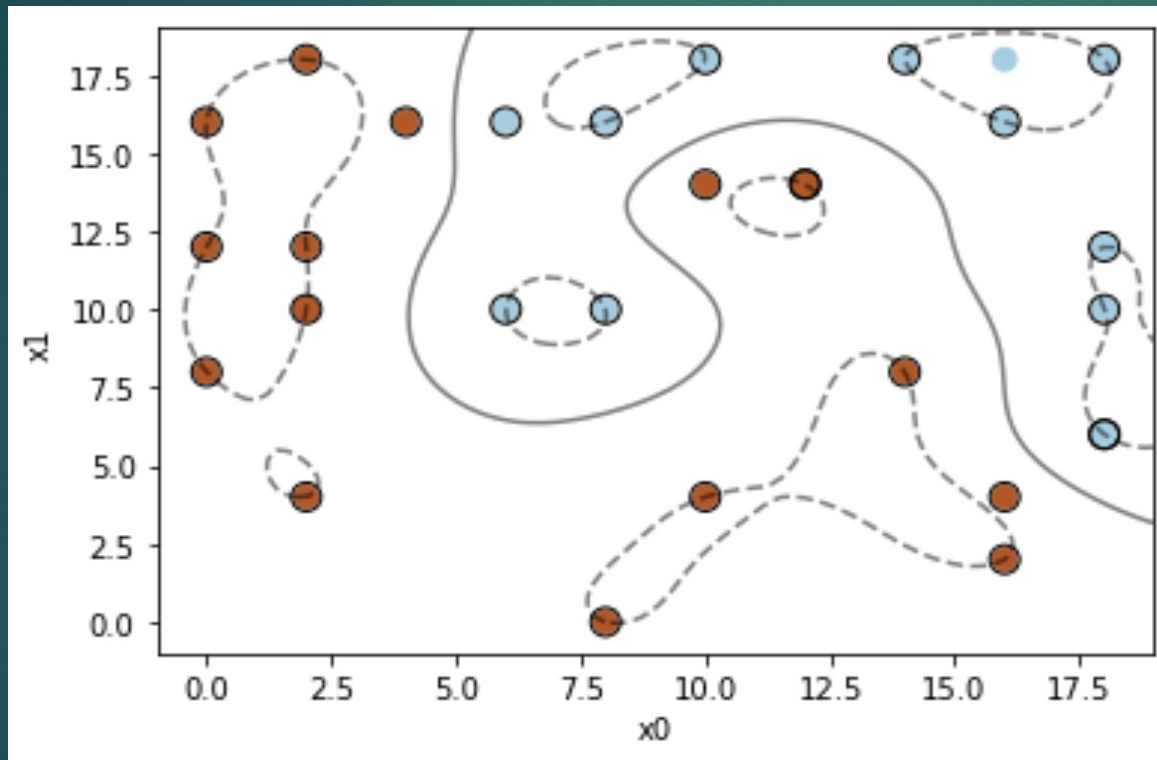
$$\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = 1 + 2x_0z_0 + 2x_1z_1 + 2x_0z_0x_1z_1 + x_0^2z_0^2 + x_1^2z_1^2$$



# Feature Maps

17

- ▶ The dimensionality of the feature space can be very large (even infinite) compromising computational efficiency



polynomial feature map  
degree=7



# Feature Maps and “Kernel Trick”

18

- ▶ Certain feature maps can be captured analytically by a kernel, without ever computing or representing  $\phi(\mathbf{x})$

- ▶ Polynomial kernel of degree  $d$ :

$$\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = \kappa(\mathbf{x}, \mathbf{z}) = (1 + \langle \mathbf{x}, \mathbf{z} \rangle)^d$$

- ▶ Radial Basis Function kernel (infinitely dimensional):

$$\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = \kappa(\mathbf{x}, \mathbf{z}) = e^{-\gamma \|\mathbf{x} - \mathbf{z}\|^2}$$

- ▶ The “kernel trick” allows for memory and computation efficient implementations of high dimensional feature maps

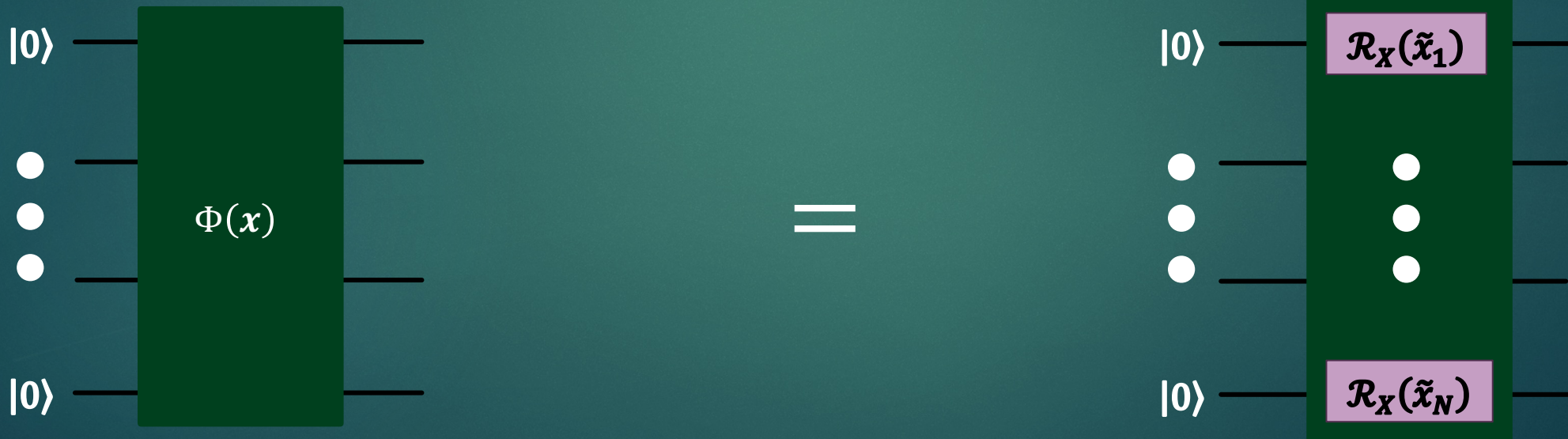


# Quantum Feature Map

19

- ▶ The quantum feature map, operator  $\Phi(x)$  might be seen as the encoding of data into a quantum state.
- ▶ Example: angle (or rotation) encoding

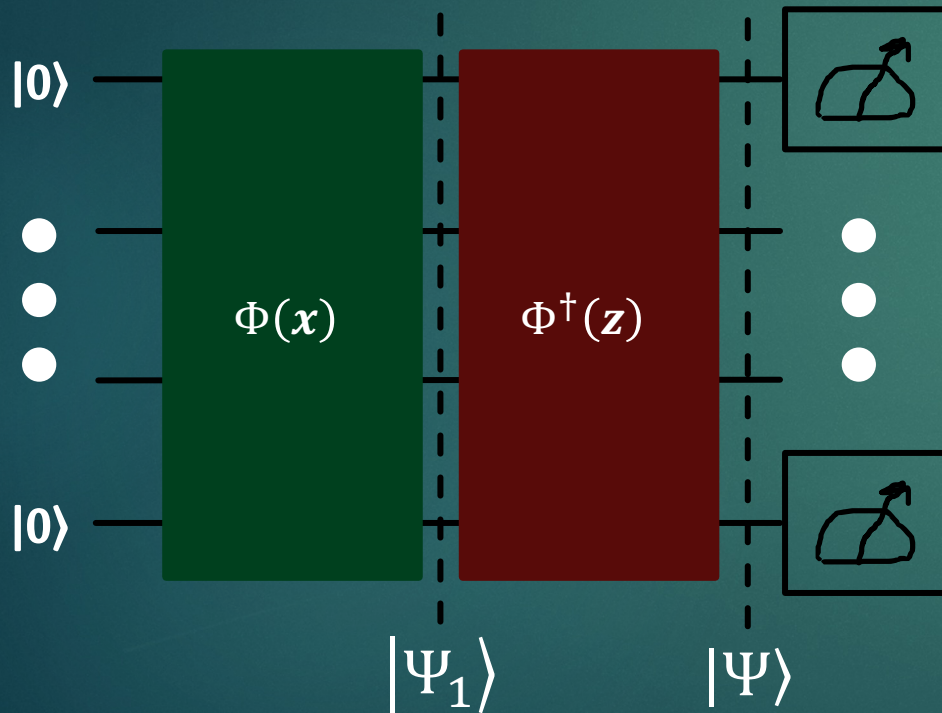
$$\tilde{x}_i = \frac{x_i}{\max_j (abs(x_j))} * \pi$$





# Quantum $\kappa(x, z)$

20



- ▶  $|\Psi_1\rangle = \Phi(x)|0\rangle$
- ▶  $|\Psi\rangle = \Phi^\dagger(z)\Phi(x)|0\rangle$
- ▶  $P(|0\rangle) = \langle 0|\Phi^\dagger(z)\Phi(x)|0\rangle = |\langle\Phi(x)|\Phi(z)\rangle|^2$
- ▶  $|\langle\Phi(x)|\Phi(z)\rangle|^2 = P(|0\rangle)$

## Quantum advantage:

$\Phi(x)$  must be computationally hard to evaluate classically

Necessary but not sufficient condition