



# Statistical Learning II

Lecture 12 - Kernel methods

---

**Bruno Loureiro**  
@ CSD, DI-ENS & CNRS

[brloureiro@gmail.com](mailto:brloureiro@gmail.com)

# Feature maps

---



Idea: Introduce a **feature map**:

$$\begin{aligned}\boldsymbol{\varphi} &: \mathbb{R}^d \rightarrow \mathbb{R}^p \\ \mathbf{x} &\mapsto \boldsymbol{\varphi}(\mathbf{x})\end{aligned}$$

And consider a linear predictor in **feature space**:

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(\mathbf{x}) \rangle$$

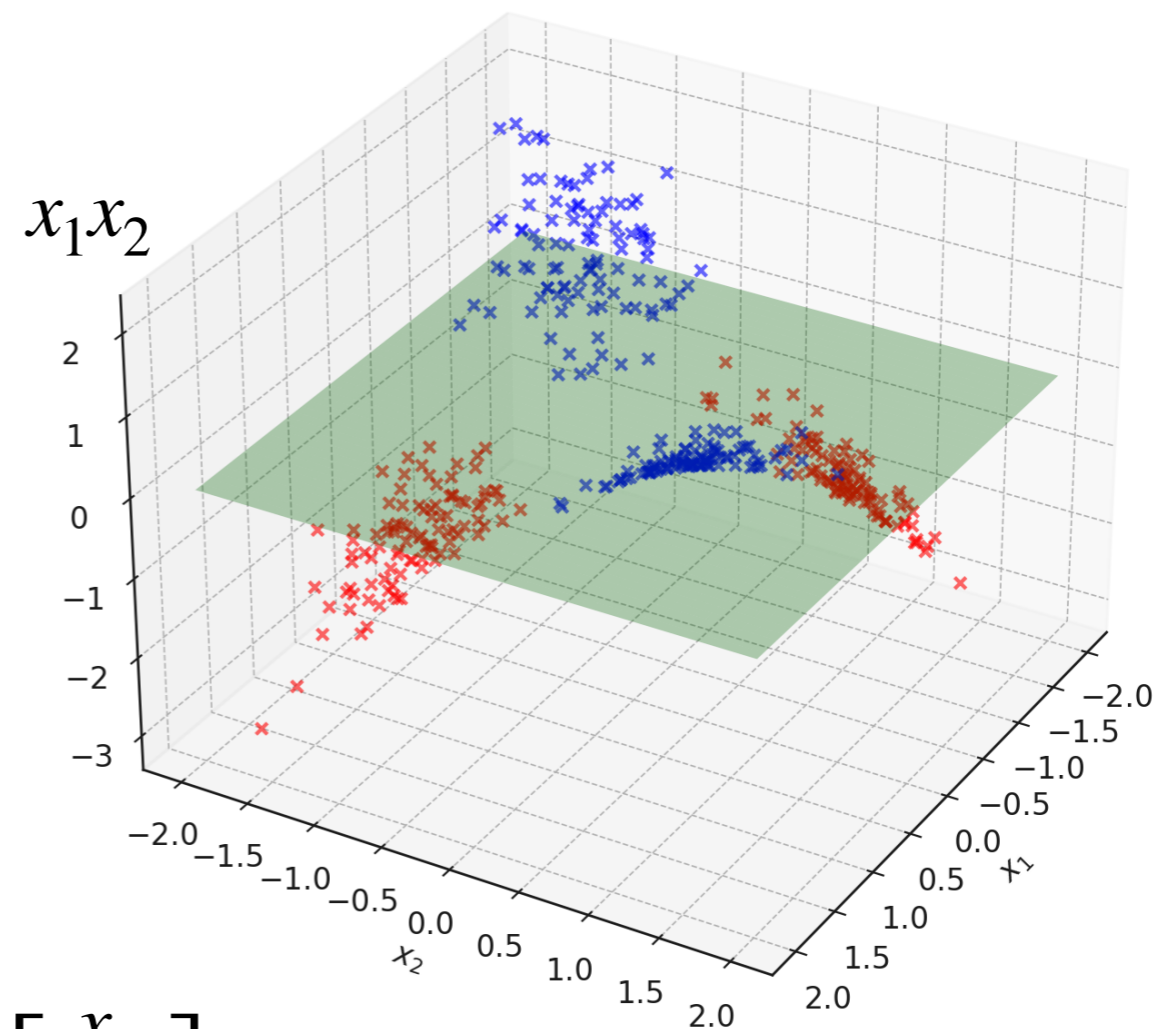
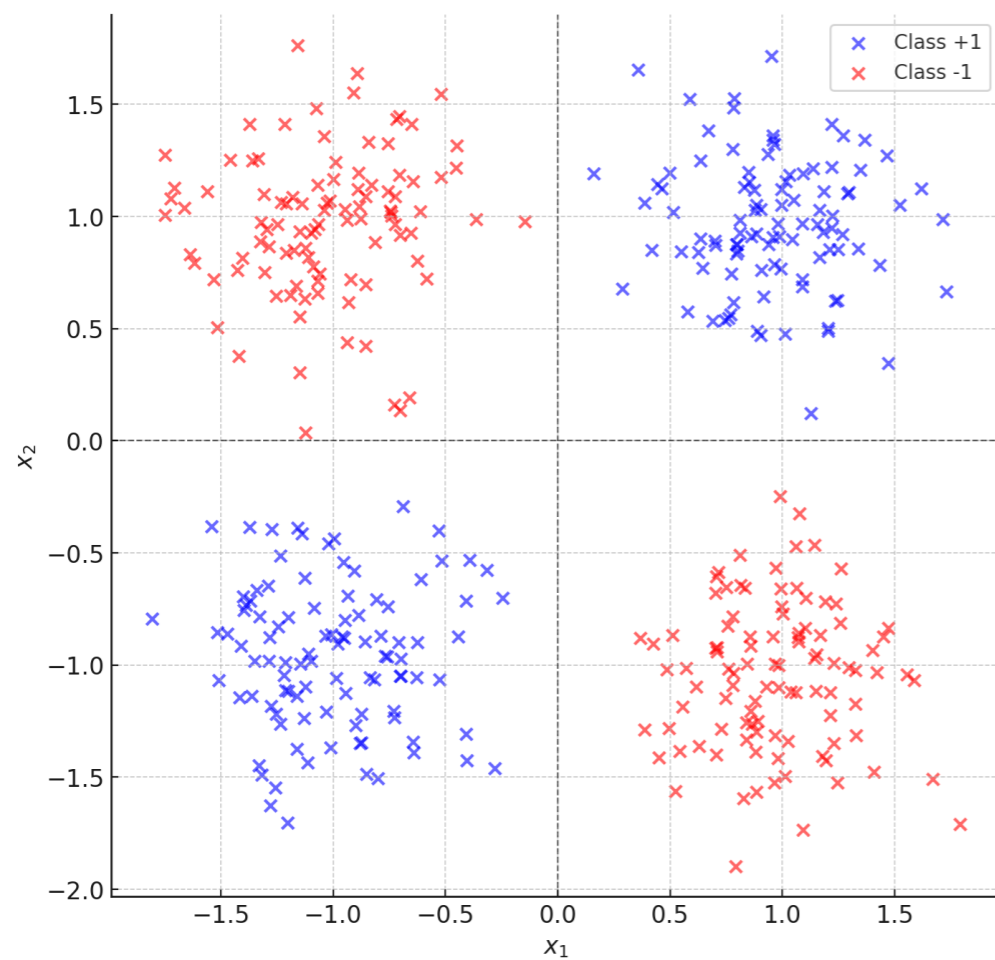


- Now we have  $\boldsymbol{\theta} \in \mathbb{R}^p$ .
- $f_{\boldsymbol{\theta}}$  still a linear function of  $\boldsymbol{\theta}$ .
- Typically  $p > d$ .
- More generally, we can consider  $\boldsymbol{\varphi} : \mathcal{X} \rightarrow \mathbb{R}^p$

Example:  $\mathcal{X}$  a collection of books.

# Examples: XOR Gaussian mixture

$$x \in \mathbb{R}^2 \quad (d = 2) \quad p(x) = \frac{1}{4} \sum_{k=1}^4 \mathcal{N}(\mu_k, I_2)$$



$$\varphi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \end{bmatrix}$$

# Ridge regression on feature space

Let  $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i \in [n]\}$  denote training data and  $\varphi : \mathcal{X} \rightarrow \mathbb{R}^p$  a feature map.

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(x_i) \rangle)^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

# Ridge regression on feature space

Let  $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i \in [n]\}$  denote training data and  $\varphi : \mathcal{X} \rightarrow \mathbb{R}^p$  a feature map.

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(x_i) \rangle)^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

Defining the feature matrix and label vector:

$$\mathbf{\Phi} = \begin{bmatrix} \boldsymbol{\varphi}(x_1) \\ \vdots \\ \boldsymbol{\varphi}(x_n) \end{bmatrix} \in \mathbb{R}^{n \times p} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

# Ridge regression on feature space

Let  $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i \in [n]\}$  denote training data and  $\varphi : \mathcal{X} \rightarrow \mathbb{R}^p$  a feature map.

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(x_i) \rangle)^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

Defining the feature matrix and label vector:

$$\mathbf{\Phi} = \begin{bmatrix} \boldsymbol{\varphi}(x_1) \\ \vdots \\ \boldsymbol{\varphi}(x_n) \end{bmatrix} \in \mathbb{R}^{n \times p} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

The above admits an explicit solution:

$$\hat{\boldsymbol{\theta}}_{\lambda}(\mathbf{\Phi}, \mathbf{y}) = (\mathbf{\Phi}^{\top} \mathbf{\Phi} + n\lambda \mathbf{I}_p)^{-1} \mathbf{\Phi}^{\top} \mathbf{y}$$

# Ridge regression on feature space

Let  $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i \in [n]\}$  denote training data and  $\varphi : \mathcal{X} \rightarrow \mathbb{R}^p$  a feature map.

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \boldsymbol{\theta}, \boldsymbol{\varphi}(x_i) \rangle)^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

Note we can equivalently write:

$$\hat{\boldsymbol{\theta}}_{\lambda}(\boldsymbol{\Phi}, \mathbf{y}) = \begin{cases} (\boldsymbol{\Phi}^{\top} \boldsymbol{\Phi} + n\lambda \mathbf{I}_p)^{-1} \boldsymbol{\Phi}^{\top} \mathbf{y} \\ \boldsymbol{\Phi}^{\top} (\boldsymbol{\Phi} \boldsymbol{\Phi}^{\top} + n\lambda \mathbf{I}_n)^{-1} \mathbf{y} \end{cases}$$



Same result, but one might be cheaper than the other.

# Kernels

---

Note that the solution:

$$\hat{\theta}_\lambda(\Phi, \mathbf{y}) = \Phi^\top (\Phi \Phi^\top + n\lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

Actually lives in the  $\text{span}(\boldsymbol{\varphi}(x_1), \dots, \boldsymbol{\varphi}(x_n))$ .



# Kernels

---

Note that the solution:

$$\hat{\boldsymbol{\theta}}_{\lambda}(\boldsymbol{\Phi}, \mathbf{y}) = \boldsymbol{\Phi}^{\top} (\boldsymbol{\Phi} \boldsymbol{\Phi}^{\top} + n\lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

Actually lives in the  $\text{span}(\boldsymbol{\varphi}(x_1), \dots, \boldsymbol{\varphi}(x_n))$ . This means we can also write:

$$\hat{\boldsymbol{\theta}}_{\lambda} = \boldsymbol{\Phi}^{\top} \hat{\boldsymbol{\alpha}}_{\lambda} \quad \hat{\boldsymbol{\alpha}}_{\lambda}(\boldsymbol{\Phi}, \mathbf{y}) = (\boldsymbol{\Phi} \boldsymbol{\Phi}^{\top} + n\lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

# Kernels

---

Note that the solution:

$$\hat{\boldsymbol{\theta}}_{\lambda}(\boldsymbol{\Phi}, \mathbf{y}) = \boldsymbol{\Phi}^{\top} (\boldsymbol{\Phi} \boldsymbol{\Phi}^{\top} + n\lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

Actually lives in the  $\text{span}(\boldsymbol{\varphi}(x_1), \dots, \boldsymbol{\varphi}(x_n))$ . This means we can also write:

$$\hat{\boldsymbol{\theta}}_{\lambda} = \boldsymbol{\Phi}^{\top} \hat{\boldsymbol{\alpha}}_{\lambda} \quad \hat{\boldsymbol{\alpha}}_{\lambda}(\boldsymbol{\Phi}, \mathbf{y}) = (\boldsymbol{\Phi} \boldsymbol{\Phi}^{\top} + n\lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

And the predictor:

$$f_{\theta}(x) = \langle \hat{\boldsymbol{\theta}}_{\lambda}, \boldsymbol{\varphi}(x) \rangle = \langle \hat{\boldsymbol{\alpha}}_{\lambda}, \boldsymbol{\Phi} \boldsymbol{\varphi}(x) \rangle$$

# Kernels

---

$$f_{\theta}(x) = \langle \hat{\theta}_{\lambda}, \varphi(x) \rangle = \langle \hat{\alpha}_{\lambda}, \Phi \varphi(x) \rangle$$

$$\hat{\alpha}_{\lambda}(\Phi, \mathbf{y}) = (\Phi \Phi^{\top} + n\lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

Note everything only depends on the scalar product of features

$$K(x, x') = \langle \varphi(x), \varphi(x') \rangle$$

This is also known as a *kernel*.

# Kernels

---

$$f_{\theta}(x) = \langle \hat{\theta}_{\lambda}, \varphi(x) \rangle = \langle \hat{\alpha}_{\lambda}, \Phi \varphi(x) \rangle$$

$$\hat{\alpha}_{\lambda}(\Phi, \mathbf{y}) = (\Phi \Phi^{\top} + n\lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

Note everything only depends on the scalar product of features

$$K(x, x') = \langle \varphi(x), \varphi(x') \rangle$$

This is also known as a *kernel*.



This is true for any linear predictor, and goes under the name of “representer theorem”

# Kernel methods

---

# Hilbert space

---

As we have shown in the previous examples, it is easier to linearly separate a function in higher dimensions.



Key idea: Take the number of features to infinity ( $p \rightarrow \infty$ )

# Hilbert space

---

As we have shown in the previous examples, it is easier to linearly separate a function in higher dimensions.



Key idea: Take the number of features to infinity ( $p \rightarrow \infty$ )



First, we need to make sense of  $\mathbb{R}^\infty$ ...

# Hilbert space

---

As we have shown in the previous examples, it is easier to linearly separate a function in higher dimensions.



Key idea: Take the number of features to infinity ( $p \rightarrow \infty$ )



First, we need to make sense of  $\mathbb{R}^\infty$ ...

## Definition (Hilbert space)

A Hilbert space  $\mathcal{H}$  is a **vector space** (over  $\mathbb{R}$  or  $\mathbb{C}$ ) with an **inner product**  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  which is **complete**.



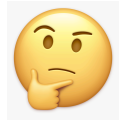
# Hilbert space

---

As we have shown in the previous examples, it is easier to linearly separate a function in higher dimensions.



Key idea: Take the number of features to infinity ( $p \rightarrow \infty$ )



First, we need to make sense of  $\mathbb{R}^\infty$ ...

## Definition (Hilbert space)

A Hilbert space  $\mathcal{H}$  is a **vector space** (over  $\mathbb{R}$  or  $\mathbb{C}$ ) with an **inner product**  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  which is **complete**.

Informally, an inner product is the minimum we need to do linear algebra in infinite dimensions

# Hilbert space

---

## Definition (Hilbert space)

A Hilbert space  $\mathcal{H}$  is a **vector space** (over  $\mathbb{R}$  or  $\mathbb{C}$ ) with an **inner product**  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  which is **complete**.

- Vector space (over  $\mathbb{R}$ ): Let  $a, b \in \mathbb{R}$  and  $f, g \in \mathcal{H}$

$$af + bg \in \mathcal{H} \quad + \text{ usual properties of the sum}$$

# Hilbert space

---

## Definition (Hilbert space)

A Hilbert space  $\mathcal{H}$  is a **vector space** (over  $\mathbb{R}$  or  $\mathbb{C}$ ) with an **inner product**  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  which is **complete**.

- Vector space (over  $\mathbb{R}$ ): Let  $a, b \in \mathbb{R}$  and  $f, g \in \mathcal{H}$

$$af + bg \in \mathcal{H} \quad + \text{ usual properties of the sum}$$

- Inner product: a function  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  such that:
  - $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$
  - $\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} \geq 0$  with equality iff  $f = 0$
  - $\langle af + bg, h \rangle_{\mathcal{H}} = a\langle f, h \rangle_{\mathcal{H}} + b\langle g, h \rangle_{\mathcal{H}}$

# Hilbert space

---

## Definition (Hilbert space)

A Hilbert space  $\mathcal{H}$  is a **vector space** (over  $\mathbb{R}$  or  $\mathbb{C}$ ) with an **inner product**  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  which is **complete**.

- Vector space (over  $\mathbb{R}$ ): Let  $a, b \in \mathbb{R}$  and  $f, g \in \mathcal{H}$

$$af + bg \in \mathcal{H} \quad + \text{ usual properties of the sum}$$

- Inner product: a function  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  such that:
  - $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$
  - $\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} \geq 0$  with equality iff  $f = 0$
  - $\langle af + bg, h \rangle_{\mathcal{H}} = a\langle f, h \rangle_{\mathcal{H}} + b\langle g, h \rangle_{\mathcal{H}}$



Inner product induces norm, but converse not always true.

# Hilbert space

---

## Definition (Hilbert space)

A Hilbert space  $\mathcal{H}$  is a **vector space** (over  $\mathbb{R}$  or  $\mathbb{C}$ ) with an **inner product**  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  which is **complete**.

- Vector space (over  $\mathbb{R}$ ): Let  $a, b \in \mathbb{R}$  and  $f, g \in \mathcal{H}$

$$af + bg \in \mathcal{H} \quad + \text{ usual properties of the sum}$$

- Inner product: a function  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  such that:

- $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$

- $\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} \geq 0$  with equality iff  $f = 0$

- $\langle af + bg, h \rangle_{\mathcal{H}} = a\langle f, h \rangle_{\mathcal{H}} + b\langle g, h \rangle_{\mathcal{H}}$

- Complete: Cauchy sequences  $f_n \in \mathcal{H}$  converge  $f_{\infty} \in \mathcal{H}$

# Examples of Hilbert spaces

- $\mathcal{H} = \mathbb{R}^d$  with the usual Euclidean inner product:

$$\langle \mathbf{u}, \mathbf{v} \rangle_2 = \sum_{i=1}^d v_i u_i$$

# Examples of Hilbert spaces

- $\mathcal{H} = \mathbb{R}^d$  with the usual Euclidean inner product:

$$\langle \mathbf{u}, \mathbf{v} \rangle_2 = \sum_{i=1}^d v_i u_i$$

- $\ell^2(\mathbb{R})$ : sequences  $\mathbf{u} = (u_1, u_2, \dots)$  with

Such that: 
$$\|\mathbf{u}\|_{\ell^2}^2 = \langle \mathbf{u}, \mathbf{u} \rangle_{\ell^2} = \sum_{i=1}^{\infty} |u_i|^2 < \infty$$

# Examples of Hilbert spaces

- $\mathcal{H} = \mathbb{R}^d$  with the usual Euclidean inner product:

$$\langle \mathbf{u}, \mathbf{v} \rangle_2 = \sum_{i=1}^d v_i u_i$$

- $\ell^2(\mathbb{R})$ : sequences  $\mathbf{u} = (u_1, u_2, \dots)$  with

Such that: 
$$\|\mathbf{u}\|_{\ell^2}^2 = \langle \mathbf{u}, \mathbf{u} \rangle_{\ell^2} = \sum_{i=1}^{\infty} |u_i|^2 < \infty$$

- $L^2(\mathbb{R})$ : functions  $f: \mathbb{R} \rightarrow \mathbb{R}$  with 
$$\langle f, g \rangle_{L^2(\mathbb{R})} = \int_{-\infty}^{\infty} f(x)g(x)dx$$

Such that: 
$$\|f\|_{L^2(\mathbb{R})}^2 = \langle f, f \rangle_{L^2(\mathbb{R})} = \int_{-\infty}^{\infty} |f(x)|^2 dx < \infty$$



# Infinite dimensional features

---

This provides the right structure to define infinite dimensions features.

# Infinite dimensional features

This provides the right structure to define infinite dimensional features.

Let  $\mathcal{H}$  denote a Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ .



Idea: Given data  $x \in \mathcal{X}$ , define features:

$$\begin{aligned}\varphi : \mathcal{X} &\rightarrow \mathcal{H} \\ x &\mapsto \varphi(x)\end{aligned}$$

and predictors:  $f_{\theta}(x) = \langle \theta, \varphi(x) \rangle_{\mathcal{H}}$  with  $\theta \in \mathcal{H}$ .

# Infinite dimensional features

This provides the right structure to define infinite dimensional features.

Let  $\mathcal{H}$  denote a Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ .



Idea: Given data  $x \in \mathcal{X}$ , define features:

$$\begin{aligned}\varphi : \mathcal{X} &\rightarrow \mathcal{H} \\ x &\mapsto \varphi(x)\end{aligned}$$

and predictors:  $f_{\theta}(x) = \langle \theta, \varphi(x) \rangle_{\mathcal{H}}$  with  $\theta \in \mathcal{H}$ .



Problems:

- In general  $f \notin \mathcal{H}$ .
- Class of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  defined this way can be small.

# Example

---

Let  $\mathcal{H} \subset \mathbb{R}^2$  with standard Euclidean inner product.

Let  $\mathcal{X} = \{x_1, x_2, x_3\}$  be a discrete data space. Define  $\varphi : \mathcal{X} \rightarrow \mathcal{H}$

$$\varphi(x_1) = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \varphi(x_2) = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \varphi(x_3) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

# Example

---

Let  $\mathcal{H} \subset \mathbb{R}^2$  with standard Euclidean inner product.

Let  $\mathcal{X} = \{x_1, x_2, x_3\}$  be a discrete data space. Define  $\varphi : \mathcal{X} \rightarrow \mathcal{H}$

$$\varphi(x_1) = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \varphi(x_2) = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \varphi(x_3) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

For any  $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \in \mathcal{H}$ , define the function:  $f(x) = \langle \theta, \varphi(x) \rangle$

We have:  $f(x_1) = \theta_1$     $f(x_2) = \theta_2$     $f(x_3) = \theta_1 + \theta_2$

# Example

---

Let  $\mathcal{H} \subset \mathbb{R}^2$  with standard Euclidean inner product.

Let  $\mathcal{X} = \{x_1, x_2, x_3\}$  be a discrete data space. Define  $\varphi : \mathcal{X} \rightarrow \mathcal{H}$

$$\varphi(x_1) = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \varphi(x_2) = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \varphi(x_3) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

For any  $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \in \mathcal{H}$ , define the function:  $f(x) = \langle \theta, \varphi(x) \rangle$

We have:  $f(x_1) = \theta_1$     $f(x_2) = \theta_2$     $f(x_3) = \theta_1 + \theta_2$

Only few functions on  $\mathcal{X}$  can be expressed this way.

e.g. can't express  $f(x_1) = 1$     $f(x_2) = 0$     $f(x_3) = 2$

# Reproducing property

To make the Hilbert space compatible with  $\mathcal{X}$ , we need the following **reproducing property**:

## Definition (RKHS)

A Hilbert space  $\mathcal{H}$  of functions over  $\mathcal{X}$  is said to be a “**Reproducing Kernel Hilbert Space**” (RKHS) if there exists  $\varphi \in \mathcal{H}$  such that:

$$\forall x \in \mathcal{X} \quad \forall f \in \mathcal{H} \quad f(x) = \langle f, \varphi(x) \rangle_{\mathcal{H}}$$

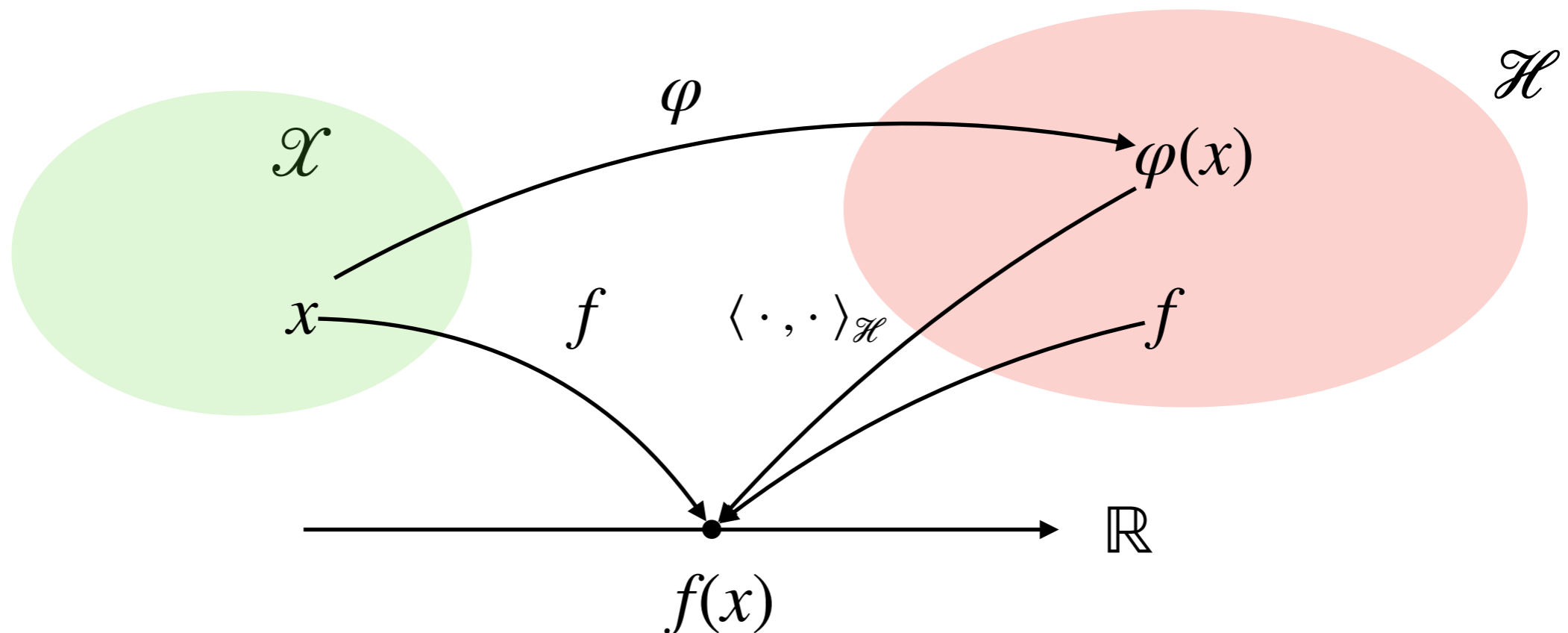
# Reproducing property

To make the Hilbert space compatible with  $\mathcal{X}$ , we need the following **reproducing property**:

## Definition (RKHS)

A Hilbert space  $\mathcal{H}$  of functions over  $\mathcal{X}$  is said to be a “**Reproducing Kernel Hilbert Space**” (RKHS) if there exists  $\varphi \in \mathcal{H}$  such that:

$$\forall x \in \mathcal{X} \quad \forall f \in \mathcal{H} \quad f(x) = \langle f, \varphi(x) \rangle_{\mathcal{H}}$$





# Kernel ridge regression

Let  $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i \in [n]\}$  denote training data. We now have everything we need to define ERM on a RKHS.

$$\min_{f \in \mathcal{H}} \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

# Kernel ridge regression

Let  $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i \in [n]\}$  denote training data. We now have everything we need to define ERM on a RKHS.

$$\min_{f \in \mathcal{H}} \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

By using the feature map  $\varphi : \mathcal{X} \rightarrow \mathcal{H}$ , this can be equivalently written:

$$\min_{\theta \in \mathcal{H}} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \theta, \varphi(x_i) \rangle_{\mathcal{H}})^2 + \frac{\lambda}{2} \|\theta\|_{\mathcal{H}}^2$$

# Kernel ridge regression

Let  $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i \in [n]\}$  denote training data. We now have everything we need to define ERM on a RKHS.

$$\min_{f \in \mathcal{H}} \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

By using the feature map  $\varphi : \mathcal{X} \rightarrow \mathcal{H}$ , this can be equivalently written:

$$\min_{\theta \in \mathcal{H}} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \theta, \varphi(x_i) \rangle_{\mathcal{H}})^2 + \frac{\lambda}{2} \|\theta\|_{\mathcal{H}}^2$$



Closed-form in terms of “infinite dimensional” matrices “ $\Phi \in \mathbb{R}^{n \times \infty}$ ”?

# Kernel ridge regression

Let  $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i \in [n]\}$  denote training data. We now have everything we need to define ERM on a RKHS.

$$\min_{\theta \in \mathcal{H}} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \theta, \varphi(x_i) \rangle_{\mathcal{H}})^2 + \frac{\lambda}{2} \|\theta\|_{\mathcal{H}}^2$$

As before, defining the kernel function and matrix

$$K(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}} \quad \mathbf{K}_{ij} = \langle \varphi(x_i), \varphi(x_j) \rangle_{\mathcal{H}}$$

# Kernel ridge regression

Let  $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i \in [n]\}$  denote training data. We now have everything we need to define ERM on a RKHS.

$$\min_{\theta \in \mathcal{H}} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \theta, \varphi(x_i) \rangle_{\mathcal{H}})^2 + \frac{\lambda}{2} \|\theta\|_{\mathcal{H}}^2$$

As before, defining the kernel function and matrix

$$K(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}} \quad \mathbf{K}_{ij} = \langle \varphi(x_i), \varphi(x_j) \rangle_{\mathcal{H}}$$

The solution can be written as:

$$\hat{f}(x) = \sum_{i=1}^n \hat{\alpha}_{\lambda, i} K(x, x_i) \quad \hat{\alpha}_{\lambda}(\Phi, \mathbf{y}) = (\mathbf{K} + n\lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

# Kernels

---

Note that in practice, to do ridge regression on  $\mathcal{H}$  we don't even need to know what  $\varphi$  is. It suffices to have  $K$ .

# Kernels

---

Note that in practice, to do ridge regression on  $\mathcal{H}$  we don't even need to know what  $\varphi$  is. It suffices to have  $K$ .

## Theorem (Aronszajn, 1950)

A function  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  defines a positive definite Kernel if and only if there exists a Hilbert space  $\mathcal{H}$  and a map  $\varphi : \mathcal{X} \rightarrow \mathcal{H}$  such that:

$$K(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}} \quad \forall x, x' \in \mathcal{X}$$

# Kernels

---

Note that in practice, to do ridge regression on  $\mathcal{H}$  we don't even need to know what  $\varphi$  is. It suffices to have  $K$ .

## Theorem (Aronszajn, 1950)

A function  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  defines a positive definite Kernel if and only if there exists a Hilbert space  $\mathcal{H}$  and a map  $\varphi : \mathcal{X} \rightarrow \mathcal{H}$  such that:

$$K(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}} \quad \forall x, x' \in \mathcal{X}$$

In words: specifying  $\mathcal{H}$  and  $\varphi$  is completely equivalent to specifying  $K$ ,



# Kernels

---

Note that in practice, to do ridge regression on  $\mathcal{H}$  we don't even need to know what  $\varphi$  is. It suffices to have  $K$ .

## Theorem (Aronszajn, 1950)

A function  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  defines a positive definite Kernel if and only if there exists a Hilbert space  $\mathcal{H}$  and a map  $\varphi : \mathcal{X} \rightarrow \mathcal{H}$  such that:

$$K(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}} \quad \forall x, x' \in \mathcal{X}$$

In words: specifying  $\mathcal{H}$  and  $\varphi$  is completely equivalent to specifying  $K$ ,

 A kernel can correspond to several feature maps. e.g.  $\mathcal{X} = \mathbb{R}$

$$\varphi(x) = x \quad \varphi(x) = \frac{1}{\sqrt{2}} \begin{bmatrix} x \\ x \end{bmatrix} \quad K(x, x') = xx'$$

# Examples of Kernels

- Gaussian kernel:  $K(\mathbf{x}, \mathbf{x}') = e^{-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|_2^2}$  (a.k.a. RBF)
- Laplace kernel:  $K(\mathbf{x}, \mathbf{x}') = e^{-\lambda \|\mathbf{x} - \mathbf{x}'\|_2}$
- Polynomial kernel:  $K(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle + b)^k$

# Examples of Kernels

- Gaussian kernel:  $K(\mathbf{x}, \mathbf{x}') = e^{-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|_2^2}$  (a.k.a. RBF)
- Laplace kernel:  $K(\mathbf{x}, \mathbf{x}') = e^{-\lambda \|\mathbf{x} - \mathbf{x}'\|_2}$
- Polynomial kernel:  $K(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle + b)^k$
- Translational invariant kernels  
 $K(\mathbf{x}, \mathbf{x}') = \kappa(\mathbf{x} - \mathbf{x}')$
- Rotationally invariant kernels  
 $K(\mathbf{x}, \mathbf{x}') = \kappa(\langle \mathbf{x}, \mathbf{x}' \rangle)$

# Examples of Kernels

- Gaussian kernel:  $K(\mathbf{x}, \mathbf{x}') = e^{-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|_2^2}$  (a.k.a. RBF)
- Laplace kernel:  $K(\mathbf{x}, \mathbf{x}') = e^{-\lambda \|\mathbf{x} - \mathbf{x}'\|_2}$
- Polynomial kernel:  $K(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle + b)^k$
- Translational invariant kernels  
 $K(\mathbf{x}, \mathbf{x}') = \kappa(\mathbf{x} - \mathbf{x}')$
- Rotationally invariant kernels  
 $K(\mathbf{x}, \mathbf{x}') = \kappa(\langle \mathbf{x}, \mathbf{x}' \rangle)$

Or any other positive-definite function...

# Examples of Kernels

- Gaussian kernel:  $K(\mathbf{x}, \mathbf{x}') = e^{-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|_2^2}$  (a.k.a. RBF)
- Laplace kernel:  $K(\mathbf{x}, \mathbf{x}') = e^{-\lambda \|\mathbf{x} - \mathbf{x}'\|_2}$
- Polynomial kernel:  $K(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle + b)^k$
- Translational invariant kernels:  $K(\mathbf{x}, \mathbf{x}') = \kappa(\mathbf{x} - \mathbf{x}')$
- Rotationally invariant kernels:  $K(\mathbf{x}, \mathbf{x}') = \kappa(\langle \mathbf{x}, \mathbf{x}' \rangle)$

Or any other positive-definite function...



In general, finding  $\varphi$  associated to these is not obvious.

# Examples of Kernels

$$y_i = \sin(x) + \varepsilon$$

$$n = 100$$

$$\varepsilon \sim \mathcal{N}(0, 0.2^2)$$

$$\lambda = 0.1$$

