# Lecture 14: Kernels

Prof. Dr. Mario Fritz

2022 06 08

https://fritz.cispa.saarland
Trustworthy AI
CISPA Helmholtz Center for Information Security, Saarbrücken, Germany

# Outline

## Bibliography

Motivation

Kernels

Properties

Vector space structure

Appendix

# Main references

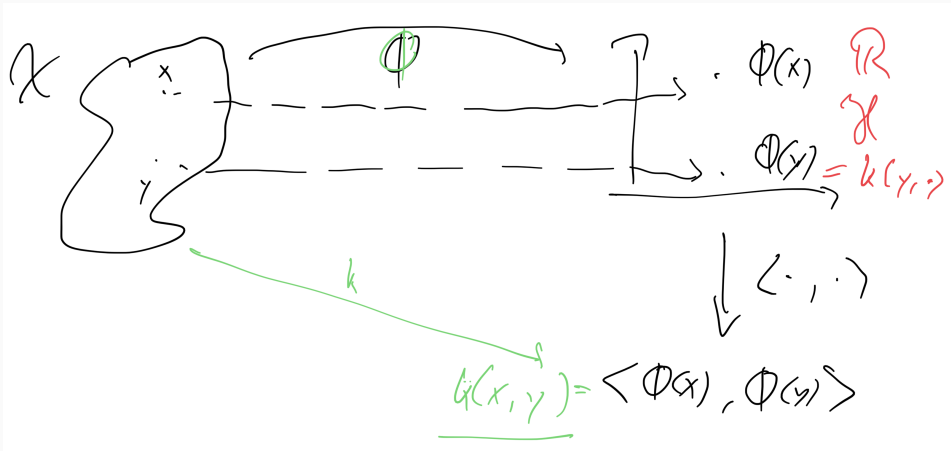- Learning with Kernels - Chapter 2
- Bishop - Chapter 6

# Outline

Insights:

(1) SVM dual objective

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{ij} \alpha_i \alpha_j \langle x_i, x_j \rangle \longleftarrow \langle \phi(x), \phi(y) \rangle$$

$$= k(x, y)$$

(2) SVM solution : $w = \sum_i \alpha_i y_i x_i$

SVM prediction: $y = \langle w, x \rangle = \sum_i \alpha_i y_i \langle x_i, x \rangle$

... refer to video lecture ...

... refer to video lecture ...

## Recap (I): Classification Problems

- An input space $\mathcal{X} \subset \mathbb{R}^d$ and an output space $\mathcal{Y} = \{-1, +1\}$.

- A training data $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)\}$

- So far, we considered

$$\mathcal{F} = \left\{ f(x) = w^\top x \mid w \in \mathbb{R}^d \right\}.$$

- Given any $f \in \mathcal{F}$, our decision function is

$$g(x) = \begin{cases} -1 & \text{if } f(x) < 0 \\ +1 & \text{if } f(x) \geq 0 \end{cases}.$$

- We also write $g(x) = \text{sgn}(f(x)) = \text{sgn}(w^\top x)$.

- **Goal:** Learn a predictor $f^* \in \mathcal{F}$ from the training data such that
    1. $\text{sgn}(f^*(x)) = y_i$ for $i = 1, \ldots, n$ and
    2. the predictor $f$ **generalizes** well to previously unseen data.

## Perceptron Algorithm

**Perceptron learning rule:**

- Initialize $w_0 = 0$. For $t = 1, \ldots, T$:
    1. $\hat{y}_t = \text{sgn}(w_t^\top x_t)$.
    2. If $\hat{y}_t = y_t$, do nothing.
    3. If $\hat{y}_t \neq y_t$, update $w_{t+1} \leftarrow w_t + y_t x_t$.

## Perceptron Algorithm

**Perceptron learning rule:**

- Initialize $w_0 = 0$. For $t = 1, \ldots, T$:
    1. $\hat{y}_t = \text{sgn}(w_t^\top x_t)$.
    2. If $\hat{y}_t = y_t$, do nothing.
    3. If $\hat{y}_t \neq y_t$, update $w_{t+1} \leftarrow w_t + y_t x_t$.
- Any solution can be expressed as

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

## Perceptron Algorithm

**Perceptron learning rule:**

- Initialize $w_0 = 0$. For $t = 1, \ldots, T$:
    1. $\hat{y}_t = \text{sgn}(w_t^\top x_t)$.
    2. If $\hat{y}_t = y_t$, do nothing.
    3. If $\hat{y}_t \neq y_t$, update $w_{t+1} \leftarrow w_t + y_t x_t$.

- Any solution can be expressed as

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

- Perceptron **prediction** rule:

$$\hat{y} = \text{sgn}(w^\top x) = \text{sgn}\left(\left(\sum_{i=1}^{n} \alpha_i y_i x_i\right)^\top x\right) = \text{sgn} \sum_{i=1}^{n} \alpha_i y_i \langle x_i, x \rangle.$$

# Perceptron Algorithm

**Perceptron learning rule:**

- Initialize $w_0 = 0$. For $t = 1, \ldots, T$:
    1. $\hat{y}_t = \text{sgn}(w_t^\top x_t)$.
    2. If $\hat{y}_t = y_t$, do nothing.
    3. If $\hat{y}_t \neq y_t$, update $w_{t+1} \leftarrow w_t + y_t x_t$.
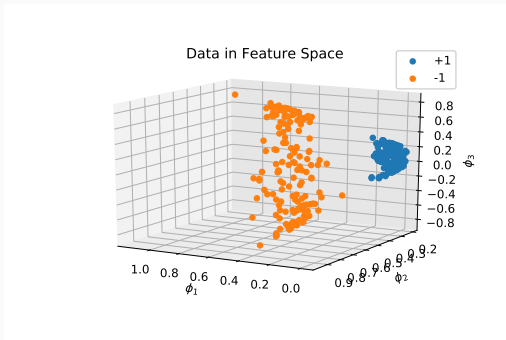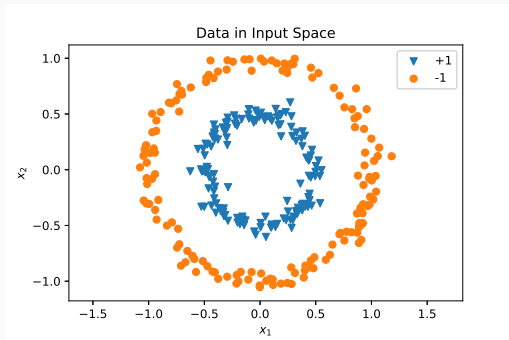- Any solution can be expressed as

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

- Perceptron **prediction** rule:

$$\hat{y} = \text{sgn}(w^\top x) = \text{sgn}\left(\left(\sum_{i=1}^n \alpha_i y_i x_i\right)^\top x\right) = \text{sgn}\sum_{i=1}^n \alpha_i y_i \langle x_i, x\rangle.$$

- **Step 3:** (Dual form) If $\hat{y}_t \neq y_i$, update $\alpha_i \leftarrow \alpha_i + 1$.

# Non-linear Classification Problem

**Question:** How would you solve this classification problem?



**Solution:** Feature Map (Embedding).

$$\phi \; : \; \Re^2 \longrightarrow \Re^3$$

$$(x_1, x_2) \longmapsto (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$$

## Solving Non-linear Classification Problems

1. Construct a non-linear feature map $\phi : \Re^d \to \Re^m$.
2. Evaluate $D_\phi = \{\phi(x_1), \phi(x_2), \ldots, \phi(x_n)\}$.
3. Learn classifier using $D_\phi$.

**Example:** In the dual perceptron:

$$w = \sum_{i=1}^{n} \alpha_i y_i \phi(x_i).$$

$$\hat{y} = \text{sgn} \sum_{i=1}^{n} \alpha_i y_i \langle \phi(x_i), \phi(x) \rangle.$$

*Observation:* In general, $m >> d$ or even $m = \infty$!

## Dual Perceptron Revisited

- Recall our feature map $\phi \; : \; (x_1, x_2) \longmapsto (x_1^2, x_2^2, \sqrt{2}x_1x_2)$
- The prediction rule of **dual perceptron algorithm** becomes

$$\hat{y} = \mathsf{sgn}(\mathsf{w}^\top \phi(\mathsf{x})) = \mathsf{sgn}\left(\sum_{i=1}^{n} \alpha_i y_i \langle \phi(\mathsf{x}_i), \phi(\mathsf{x}) \rangle\right).$$

## Dual Perceptron Revisited

- Recall our feature map $\phi : (x_1, x_2) \longmapsto (x_1^2, x_2^2, \sqrt{2}x_1x_2)$
- The prediction rule of **dual perceptron algorithm** becomes

$$\hat{y} = \text{sgn}(\mathsf{w}^\top \phi(\mathsf{x})) = \text{sgn}\left(\sum_{i=1}^{n} \alpha_i y_i \langle \phi(\mathsf{x}_i), \phi(\mathsf{x}) \rangle\right).$$

- An inner product between $\phi(\mathsf{x})$ and $\phi(\mathsf{z})$ in $\mathbb{R}^3$

$$\begin{aligned}
\langle \phi(\mathsf{x}), \phi(\mathsf{z}) \rangle_{\mathbb{R}^3} &= (x_1^2, x_2^2, \sqrt{2}x_1x_2)^\top (z_1^2, z_2^2, \sqrt{2}z_1z_2) \\
&= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1x_2z_1z_2
\end{aligned}$$

## Dual Perceptron Revisited

- Recall our feature map $\phi : (x_1, x_2) \longmapsto (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$
- The prediction rule of **dual perceptron algorithm** becomes

$$\hat{y} = \text{sgn}(\mathsf{w}^\top \phi(\mathsf{x})) = \text{sgn}\left(\sum_{i=1}^{n} \alpha_i y_i \langle \phi(\mathsf{x}_i), \phi(\mathsf{x}) \rangle\right).$$

- An inner product between $\phi(\mathsf{x})$ and $\phi(\mathsf{z})$ in $\mathbb{R}^3$

$$\begin{aligned}
\langle \phi(\mathsf{x}), \phi(\mathsf{z}) \rangle_{\mathbb{R}^3} &= (x_1^2, x_2^2, \sqrt{2}x_1 x_2)^\top (z_1^2, z_2^2, \sqrt{2}z_1 z_2) \\
&= x_1^2 z_1^2 + x_2^2 z_2^2 + 2 x_1 x_2 z_1 z_2
\end{aligned}$$

## Dual Perceptron Revisited

- Recall our feature map $\phi : (x_1, x_2) \longmapsto (x_1^2, x_2^2, \sqrt{2}x_1x_2)$

- The prediction rule of **dual perceptron algorithm** becomes

$$\hat{y} = \text{sgn}(w^\top \phi(x)) = \text{sgn}\left(\sum_{i=1}^n \alpha_i y_i \langle \phi(x_i), \phi(x) \rangle\right).$$

- An inner product between $\phi(x)$ and $\phi(z)$ in $\mathbb{R}^3$

$$\begin{aligned}
\langle \phi(x), \phi(z) \rangle_{\mathbb{R}^3} &= (x_1^2, x_2^2, \sqrt{2}x_1x_2)^\top (z_1^2, z_2^2, \sqrt{2}z_1z_2) \\
&= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1x_2z_1z_2 \\
&= (x_1z_1 + x_2z_2)^2
\end{aligned}$$

## Dual Perceptron Revisited

- Recall our feature map $\phi \, : \, (x_1, x_2) \longmapsto (x_1^2, x_2^2, \sqrt{2}x_1x_2)$

- The prediction rule of **dual perceptron algorithm** becomes

$$\hat{y} = \text{sgn}(w^\top \phi(x)) = \text{sgn}\left(\sum_{i=1}^{n} \alpha_i y_i \langle \phi(x_i), \phi(x) \rangle\right).$$

- An inner product between $\phi(x)$ and $\phi(z)$ in $\mathbb{R}^3$

$$
\begin{aligned}
\langle \phi(x), \phi(z) \rangle_{\mathbb{R}^3} &= (x_1^2, x_2^2, \sqrt{2}x_1x_2)^\top (z_1^2, z_2^2, \sqrt{2}z_1z_2) \\
&= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1x_2z_1z_2 \\
&= (x_1z_1 + x_2z_2)^2 \\
&= (x \cdot z)^2.
\end{aligned}
$$

## Dual Perceptron Revisited

- Recall our feature map $\phi : (x_1, x_2) \longmapsto (x_1^2, x_2^2, \sqrt{2}x_1x_2)$

- The prediction rule of **dual perceptron algorithm** becomes

$$\hat{y} = \text{sgn}(w^\top \phi(x)) = \text{sgn}\left(\sum_{i=1}^{n} \alpha_i y_i \langle \phi(x_i), \phi(x) \rangle\right).$$

- An inner product between $\phi(x)$ and $\phi(z)$ in $\mathbb{R}^3$

$$\begin{aligned}
\langle \phi(x), \phi(z) \rangle_{\mathbb{R}^3} &= (x_1^2, x_2^2, \sqrt{2}x_1x_2)^\top (z_1^2, z_2^2, \sqrt{2}z_1z_2) \\
&= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1x_2z_1z_2 \\
&= (x_1z_1 + x_2z_2)^2 \\
&= (x \cdot z)^2.
\end{aligned}$$

- For $x, x' \in \mathcal{X}$, define a **kernel** function:

$$k(x, x') = \langle \phi(x_i), \phi(x) \rangle = (x \cdot x')^2$$

## Dual Perceptron Revisited

- Recall our feature map $\phi : (x_1, x_2) \longmapsto (x_1^2, x_2^2, \sqrt{2}x_1x_2)$

- The prediction rule of **dual perceptron algorithm** becomes

$$\hat{y} = \text{sgn}(w^\top \phi(x)) = \text{sgn}\left(\sum_{i=1}^n \alpha_i y_i \langle \phi(x_i), \phi(x)\rangle\right).$$

- An inner product between $\phi(x)$ and $\phi(z)$ in $\mathbb{R}^3$

$$\begin{aligned}
\langle \phi(x), \phi(z)\rangle_{\mathbb{R}^3} &= (x_1^2, x_2^2, \sqrt{2}x_1x_2)^\top (z_1^2, z_2^2, \sqrt{2}z_1z_2) \\
&= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1x_2z_1z_2 \\
&= (x_1z_1 + x_2z_2)^2 \\
&= (x \cdot z)^2.
\end{aligned}$$

- For $x, x' \in \mathcal{X}$, define a **kernel** function:

$$k(x, x') = \langle \phi(x_i), \phi(x)\rangle = (x \cdot x')^2$$

- Hence, $\hat{y} = \text{sgn}\left(\sum_{i=1}^n \alpha_i y_i k(x_i, x)\right)$. **We do not need the embedding $\phi$ but only a kernel!**

## Kernel Trick

**Theorem (Kernel implies embedding)**

*A function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a* kernel *on $\mathcal{X}$ if and only if there exists a Hilbert space $\mathcal{H}$ and a (feature) map $\phi : \mathcal{X} \to \mathcal{H}$ such that for all $x, x' \in \mathcal{X}$ we have*

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}. \tag{1}$$

## Kernel Trick

**Theorem (Kernel implies embedding)**

*A function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a* kernel *on $\mathcal{X}$ if and only if there exists a Hilbert space $\mathcal{H}$ and a (feature) map $\phi : \mathcal{X} \to \mathcal{H}$ such that for all $x, x' \in \mathcal{X}$ we have*

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}. \tag{1}$$

Main properties of a **Hilbert space** $\mathcal{H}$ to keep in mind:

- $\mathcal{H}$ is a vector space with scarlar product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$.
- Space is complete (all Cauchy sequences converge).
- Scalar product leads to a norm: $\|x\|_{\mathcal{H}} = \langle x, x \rangle_{\mathcal{H}}$

*Observation:* If you are not familiar with Hilbert spaces, think of $\mathbb{R}^d$.

**Definition**

A **reproducing kernel Hilbert space (RKHS)** $\mathcal{H}$ on $\mathcal{X}$ is a Hilbert space of linear functions from $\mathcal{X}$ to $\mathbb{R}$ with a reproducing kernel $k(x, x')$ on $\mathcal{X} \times \mathcal{X}$ such that

$$\forall x \in \mathcal{X}, \quad k(x, \cdot) \in \mathcal{H}$$
$$\forall f \in \mathcal{H}, \quad \langle f, k(x, \cdot) \rangle_{\mathcal{H}} = f(x) \quad \text{(reproducing property)}$$

*Important:* RKHS defines a space of pointwise defined functions! In ML, we do have to do predictions for each point!

**Theorem (Moore)**

*If $k$ is a positive definite kernel then there exists a* unique *reproducing kernel Hilbert space $\mathcal{H}$ whose kernel is $k$.*

**CISPA**
HELMHOLTZ CENTER FOR
INFORMATION SECURITY

**Theorem (Moore)**

*If k is a positive definite kernel then there exists a* unique *reproducing kernel Hilbert space $\mathcal{H}$ whose kernel is k.*

**Important observations:**

- there is a **one-to-one** relation between reproducing kernel Hilbert spaces and positive definite kernels.
- there exist an embedding $\Phi : \mathcal{X} \to \mathcal{H}$, with $x \to \Phi(x) \in \mathcal{H}$ of the input space into a Hilbert space (a.k.a. **feature space**),
- the embedding is not unique given a kernel but there exists one and they are all isometric isomorphic, the easiest is

$$\Phi : x \to k(x, \cdot) \in \mathcal{H}_k.$$

- kernel allows us to compute the inner product between the embedded vectors $k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}},$

## Positive Definite Kernels

---

**Definition (Positive definite (PD) kernel)**

A kernel $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is **positive definite (PD)** if for all $m \geq 1$, $x_1, \ldots, x_m \in \mathcal{X}$, $c_1, \ldots, c_m \in \mathbb{R}$

$$\sum_{i,j=1}^{m} c_i c_j k(x_i, x_j) \geq 0$$

The set of all real-valued positive definite kernels on $\mathcal{X}$ is denoted $\mathbb{R}_+^{\mathcal{X} \times \mathcal{X}}$.

---

**Remark:**

- In this lecture a kernel is always positive definite if not stated otherwise.
- Note that $\mathcal{X}$ is a general set $\implies$ later on we will define kernels on structured domains (graphs, histograms, etc.).

## Gram Matrix

**Definition (Gram Matrix)**

Given a kernel $k$ and a set of $n$ points $x_1, \ldots, x_n \in \mathcal{X}$ the $n \times n$ matrix

$$K = (k(x_i, x_j))_{ij},$$

is called the **kernel matrix** (or **Gram Matrix**) $K$ of the kernel $k$ with respect to $x_1, \ldots, x_n$.

$$\phi^\top \phi = K = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \ldots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \ldots & k(x_2, x_n) \\ \vdots & & \vdots & \\ k(x_n, x_1) & k(x_n, x_2) & \ldots & k(x_n, x_n) \end{bmatrix}$$

*Observation:* The Gram matrix is therefore positive definite!

- **Linear kernel:** $k(x, x') = x \cdot x'$
  - $\phi(x) = x$ and $\mathcal{H} = \mathbb{R}$
  - $\phi(x) = (x/\sqrt{2}, x/\sqrt{2})$ and $\mathcal{H} = \mathbb{R}^2$

- **Linear kernel:** $k(\mathsf{x}, \mathsf{x}') = \mathsf{x} \cdot \mathsf{x}'$
  - $\phi(\mathsf{x}) = \mathsf{x}$ and $\mathcal{H} = \mathbb{R}$
  - $\phi(\mathsf{x}) = (\mathsf{x}/\sqrt{2}, \mathsf{x}/\sqrt{2})$ and $\mathcal{H} = \mathbb{R}^2$
- **Polynomial kernel:** $k(\mathsf{x}, \mathsf{x}') = (\mathsf{x} \cdot \mathsf{x}' + c)^m$ for $c \geq 0$
  - $\phi(\mathsf{x}) = \left( \sqrt{\binom{m}{n_1, \ldots, n_{d+1}}} x_1^{n_1} \cdots x_d^{n_d} \cdot c^{n_{d+1}/2} \right)_{n_i \geq 0, \sum_{i=1}^{d+1} n_i = m}$
  - $\dim(\mathcal{H}) = \binom{d+m}{m}$

- **Linear kernel:** $k(\mathsf{x}, \mathsf{x}') = \mathsf{x} \cdot \mathsf{x}'$
  - $\phi(\mathsf{x}) = \mathsf{x}$ and $\mathcal{H} = \mathbb{R}$
  - $\phi(\mathsf{x}) = (\mathsf{x}/\sqrt{2}, \mathsf{x}/\sqrt{2})$ and $\mathcal{H} = \mathbb{R}^2$
- **Polynomial kernel:** $k(\mathsf{x}, \mathsf{x}') = (\mathsf{x} \cdot \mathsf{x}' + c)^m$ for $c \geq 0$
  - $\phi(x) = \left( \sqrt{\binom{m}{n_1, \ldots, n_{d+1}}} x_1^{n_1} \cdots x_d^{n_d} \cdot c^{n_{d+1}/2} \right)_{n_i \geq 0, \sum_{i=1}^{d+1} n_i = m}$
  - $\dim(\mathcal{H}) = \binom{d+m}{m}$
- **Exponential kernel:** $k(\mathsf{x}, \mathsf{x}') = \exp(\langle \mathsf{x}, \mathsf{x}' \rangle)$
  - Assume $\mathsf{x} \in \mathbb{R}$ and use Taylor series expansion of $e^{\mathsf{x}}$,

$$\phi(x) = \left[ 1, x, \sqrt{\frac{1}{2!}} x^2, \sqrt{\frac{1}{3!}} x^3, \ldots \right], \quad \mathcal{H} = \mathbb{R}^\infty$$

# Examples of Kernels

- **Linear kernel:** $k(\mathsf{x}, \mathsf{x}') = \mathsf{x} \cdot \mathsf{x}'$
  - $\phi(\mathsf{x}) = \mathsf{x}$ and $\mathcal{H} = \mathbb{R}$
  - $\phi(\mathsf{x}) = (\mathsf{x}/\sqrt{2}, \mathsf{x}/\sqrt{2})$ and $\mathcal{H} = \mathbb{R}^2$
- **Polynomial kernel:** $k(\mathsf{x}, \mathsf{x}') = (\mathsf{x} \cdot \mathsf{x}' + c)^m$ for $c \geq 0$
  - $\phi(x) = \left( \sqrt{\binom{m}{n_1, \ldots, n_{d+1}}} x_1^{n_1} \cdots x_d^{n_d} \cdot c^{n_{d+1}/2} \right)_{n_i \geq 0, \sum_{i=1}^{d+1} n_i = m}$
  - $\dim(\mathcal{H}) = \binom{d+m}{m}$
- **Exponential kernel:** $k(\mathsf{x}, \mathsf{x}') = \exp(\langle \mathsf{x}, \mathsf{x}' \rangle)$
  - Assume $\mathsf{x} \in \mathbb{R}$ and use Taylor series expansion of $e^x$,
  $$\phi(x) = \left[ 1, x, \sqrt{\frac{1}{2!}} x^2, \sqrt{\frac{1}{3!}} x^3, \ldots \right], \quad \mathcal{H} = \mathbb{R}^\infty$$
- **Gaussian RBF kernel:** $k(\mathsf{x}, \mathsf{x}') = \exp(-\gamma \|\mathsf{x} - \mathsf{x}'\|_2^2)$
  - Assume $\mathsf{x} \in \mathbb{R}$ and use Taylor series expansion of $e^x$,
  $$\phi(x) = \exp(-\gamma x^2) \left[ 1, \sqrt{\frac{2\gamma}{1!}} x, \sqrt{\frac{(2\gamma)^2}{2!}} x^2, \sqrt{\frac{(2\gamma)^3}{3!}} x^3, \ldots \right], \quad \mathcal{H} = \mathbb{R}^\infty$$

# Outline

**Properties of kernels:** Transformations of kernels which preserve the property of positive definiteness are important for

1. the construction of new kernels,
2. the verification that a given function $k(x, x')$ is positive definite, as we often have a similarity measure, we would like to use.

## Properties

Let $k_1, k_2$ be kernels on $\mathcal{X}$. Then, the following properties allow us to build complicated kernels from simpler ones $(k_1, k_2)$.

1. $k(x, x') = \alpha k_1(x, x')$ for $\alpha > 0$.

## Properties

Let $k_1, k_2$ be kernels on $\mathcal{X}$. Then, the following properties allow us to build complicated kernels from simpler ones $(k_1, k_2)$.

1. $k(x, x') = \alpha k_1(x, x')$ for $\alpha > 0$.
2. $k(x, x') = f(x) k_1(x, x') f(x')$, where $f(\cdot)$ is any function.
3. $k(x, x') = f(k_1(x, x'))$, where $f(\cdot)$ is a polynomial with non-negative coefficients.
4. $k(x, x') = \exp(k_1(x, x'))$

## Properties

Let $k_1, k_2$ be kernels on $\mathcal{X}$. Then, the following properties allow us to build complicated kernels from simpler ones $(k_1, k_2)$.

1. $k(x, x') = \alpha k_1(x, x')$ for $\alpha > 0$.
2. $k(x, x') = f(x) k_1(x, x') f(x')$, where $f(\cdot)$ is any function.
3. $k(x, x') = f(k_1(x, x'))$, where $f(\cdot)$ is a polynomial with non-negative coefficients.
4. $k(x, x') = \exp(k_1(x, x'))$
5. $k(x, x') = k_1(x, x') + k_2(x, x')$.
6. $k(x, x') = k_1(x, x') k_2(x, x')$.

Let $k_1, k_2$ be kernels on $\mathcal{X}$. Then, the following properties allow us to build complicated kernels from simpler ones $(k_1, k_2)$.

1. $k(x, x') = \alpha k_1(x, x')$ for $\alpha > 0$.
2. $k(x, x') = f(x) k_1(x, x') f(x')$, where $f(\cdot)$ is any function.
3. $k(x, x') = f(k_1(x, x'))$, where $f(\cdot)$ is a polynomial with non-negative coefficients.
4. $k(x, x') = \exp(k_1(x, x'))$
5. $k(x, x') = k_1(x, x') + k_2(x, x')$.
6. $k(x, x') = k_1(x, x') k_2(x, x')$.
7. $k(x, x') = k_1(\phi(x), \phi(x'))$, where $\phi : \mathcal{X} \mapsto \mathbb{R}^m$.
8. $k(x, x') = x^T A x'$ for a symmetric positive definite matrix $A$.

Let $k_1, k_2$ be kernels on $\mathcal{X}$. Then, the following properties allow us to build complicated kernels from simpler ones $(k_1, k_2)$.

1. $k(x, x') = \alpha k_1(x, x')$ for $\alpha > 0$.
2. $k(x, x') = f(x) k_1(x, x') f(x')$, where $f(\cdot)$ is any function.
3. $k(x, x') = f(k_1(x, x'))$, where $f(\cdot)$ is a polynomial with non-negative coefficients.
4. $k(x, x') = \exp(k_1(x, x'))$
5. $k(x, x') = k_1(x, x') + k_2(x, x')$.
6. $k(x, x') = k_1(x, x') k_2(x, x')$.
7. $k(x, x') = k_1(\phi(x), \phi(x'))$, where $\phi : \mathcal{X} \mapsto \mathbb{R}^m$.
8. $k(x, x') = x^T A x'$ for a symmetric positive definite matrix $A$.
9. $k(x, x') = k_1(x_a, x'_a) + k_2(x_b, x'_b)$ or $k(x, x') = k_1(x_a, x'_a) k_2(x_b, x'_b)$, with $x = (x_a, x_b)$ and variables $x_a$ and $x_b$ are not necessarily disjoint.

## Properties

Let $k_1, k_2$ be kernels on $\mathcal{X}$. Then, the following properties allow us to build complicated kernels from simpler ones $(k_1, k_2)$.

1. $k(x, x') = \alpha k_1(x, x')$ for $\alpha > 0$.
2. $k(x, x') = f(x)k_1(x, x')f(x')$, where $f(\cdot)$ is any function.
3. $k(x, x') = f(k_1(x, x'))$, where $f(\cdot)$ is a polynomial with non-negative coefficients.
4. $k(x, x') = \exp(k_1(x, x'))$
5. $k(x, x') = k_1(x, x') + k_2(x, x')$.
6. $k(x, x') = k_1(x, x')k_2(x, x')$.
7. $k(x, x') = k_1(\phi(x), \phi(x'))$, where $\phi : \mathcal{X} \mapsto \mathbb{R}^m$.
8. $k(x, x') = x^T A x'$ for a symmetric positive definite matrix $A$.
9. $k(x, x') = k_1(x_a, x'_a) + k_2(x_b, x'_b)$ or $k(x, x') = k_1(x_a, x'_a)k_2(x_b, x'_b)$, with $x = (x_a, x_b)$ and variables $x_a$ and $x_b$ are not necessarily disjoint.

**Norm of a feature vector $\Phi(x)$**

$$\|\Phi(x)\|_{\mathcal{H}} = \sqrt{\|\Phi(x)\|_{\mathcal{H}}^2} = \sqrt{k(x, x)}.$$

**Distance/Metric:**

$$d(x, x') = \sqrt{\|\Phi(x) - \Phi(x')\|_{\mathcal{H}}^2} = \sqrt{\|\Phi(x)\|_{\mathcal{H}}^2 + \|\Phi(x')\|_{\mathcal{H}}^2 - 2\langle \Phi(x), \Phi(x')\rangle_{\mathcal{H}}}$$
$$= \sqrt{k(x, x) + k(x', x') - 2k(x, x')}.$$

**Angle:**

$$\cos\left(\angle(\Phi(x), \Phi(x'))\right) = \frac{\langle \Phi(x), \Phi(x')\rangle}{\|\Phi(x)\| \|\Phi(x')\|} = \frac{k(x, x')}{\sqrt{k(x, x)k(x', x')}}.$$

**The kernel function as similarity measure between** $x$ **and** $x'$:
This interpretation is motivated by the fact that the kernel function $k(x, x')$ can be seen as an inner product

$$k(x, x') = \langle \phi(x), \phi(x') \rangle,$$

where $\phi : \mathcal{X} \to \mathcal{H}$ and $\mathcal{H}$ is a **Hilbert space**.

A new kernel $\tilde{k}$ defined as

$$\tilde{k}(x, x') = \frac{k(x, x')}{\sqrt{k(x, x)k(x', x')}} = \frac{\langle \phi(x), \phi(x') \rangle}{\|\phi(x)\| \, \|\phi(x')\|} = \cos(\angle(\phi(x), \phi(x'))),$$

is again a positive definite kernel.

- The cosine is a common similarity measure (text classification).
- $|\tilde{k}(x, x')| \leq 1$ and $\tilde{k}(x, x') = 1$ if and only if $x = x'$.

**Why are (dis)similarity measures useful for learning?**

- one needs only to define the similarity $k(x, x')$ of two points $x$ and $x'$ instead of a set of functions $\phi(x)$ in the feature maps approach.

- construction of a similarity measure between structured objects e.g. graphs is conceptually much easier than defining a certain set of feature maps on these structured objects.

**Center of mass/Centroid/Mean vector:**

$$\Phi_m = \frac{1}{n} \sum_{i=1}^{n} \Phi(x_i).$$

**Distance of $\Phi(x)$ to the center of mass:**

$$\|\Phi(x) - \Phi_m\|^2 = \langle \Phi(x), \Phi(x) \rangle + \langle \Phi_m, \Phi_m \rangle - 2 \langle \Phi(x), \Phi_m \rangle$$
$$= k(x,x) + \frac{1}{n^2} \sum_{i,j=1}^{n} k(x_i, x_j) - \frac{2}{n} \sum_{i=1}^{n} k(x, x_i).$$

**Centering of datapoints in the feature space:**

$$\tilde{\Phi}(x) = \Phi(x) - \Phi_m = \Phi(x) - \frac{1}{n} \sum_{i=1}^{n} \Phi(x_i).$$
$$\left\langle \tilde{\Phi}(x), \tilde{\Phi}(x') \right\rangle = k(x,x') - \frac{1}{n} \sum_{i=1}^{n} k(x, x_i) - \frac{1}{n} \sum_{i=1}^{n} k(x', x_i) + \frac{1}{n^2} \sum_{i,j=1}^{n} k(x_i, x_j).$$

## Example: Parzen window classifier

**Center of mass/Centroid/Mean vector:**

$$\Phi_m^+ = \frac{1}{n_+} \sum_{i=1}^{n_+} \Phi(x_i^+), \quad \Phi_m^- = \frac{1}{n_-} \sum_{j=1}^{n_-} \Phi(x_j^-).$$

**Classify by assigning point to class of closest centroid:**

$$f(x) = \mathrm{sign} \left( \left\| \Phi(x) - \Phi_m^- \right\|^2 - \left\| \Phi(x) - \Phi_m^+ \right\|^2 \right)$$

$$= \mathrm{sign} \left( \frac{2}{n_+} \sum_{i=1}^{n_+} k(x, x_i^+) - \frac{2}{n_-} \sum_{j=1}^{n_-} k(x, x_j^-) + b \right)$$

where

$$b = \frac{1}{n_-^2} \sum_{i,j=1}^{n_-} k(x_i^-, x_i^-) - \frac{1}{n_+^2} \sum_{i,j=1}^{n_+} k(x_i^+, x_j^+).$$

This is a so called **Parzen window classifier** (with offset).

*Disclosure:* The following appendix aims to provide further mathematical details on how to construct RKHS for those that are interested. However, the following material is not necessary or required to prepare for the exam.

**Definition**

A **metric space** is a set $\mathcal{X}$ with a distance function $d : \mathcal{X} \times X \to \mathbb{R}$ such that:

- $d(x, x') \geq 0$,
- $d(x, x') = 0$ if and only if $x = x'$,
- $d(x, x') = d(\mathrm{x}, x')$, (symmetry)
- $d(x, x') \leq d(x, z) + d(z, x')$. (triangle inequality)

It is denoted as $(\mathcal{X}, d)$.
For a **semi-metric** $d(x, x') = 0$ does not imply $x = x'$.

**Remark:** any semi-metric space $(\mathcal{X}, d)$ can be turned into a metric space by identifying points which have zero distance.

## Convergence and Cauchy sequences

**Definition**

A sequence of elements $\{x_n\}_{n\in\mathbb{N}}$ of a metric space $(\mathcal{X}, d)$ is said to **converge** to an element $x \in \mathcal{X}$ if $\lim_{n\to\infty} d(x, x_n) = 0$. We will denote this either as $x_n \xrightarrow{d} x$ or $\lim_{n\to\infty} x_n = x$.

**Definition**

A sequence of elements $\{x_n\}$ of a metric space $(\mathcal{X}, d)$ is called a **Cauchy sequence** if $\forall\, \epsilon > 0, \quad \exists N$ such that $d(x_n, x_m) < \epsilon, \ \forall\, n, m > N$.

**Proposition:** Every convergent sequence is a Cauchy sequence.

**Definition**

A metric space in which all Cauchy sequences converge is called **complete**.

Example: $\mathbb{R}$ is complete, but not $\mathbb{Q}$.

Sets of functions $\mathcal{F} : \mathcal{X} \to \mathbb{R}$ as vector spaces - apply vector axioms pointwise.

Three functions $f, g, h \in \mathcal{F}$, $\alpha, \beta \in \mathbb{R}$,

$$(f + g)(x) := f(x) + g(x), \quad \forall x \in \mathcal{X},$$
$$(\alpha f)(x) := \alpha f(x), \quad \forall x \in \mathcal{X}.$$

| | |
|---|---|
| Associativity | $\big(f(x) + g(x)\big) + h(x) = f(x) + \big(g(x) + h(x)\big),$ |
| Commutativity | $f(x) + g(x) = g(x) + f(x),$ |
| Identity (addition) | $f(x) + 0 = f(x), \Rightarrow$ zero function $h(x) = 0, \forall x \in \mathcal{X},$ |
| Distributivity I | $(\alpha + \beta)f(x) = \alpha f(x) + \beta f(x),$ |
| Distributivity II | $\alpha(f(x) + g(x)) = \alpha f(x) + \alpha g(x),$ |
| Compatibility | $(\alpha\beta)f(x) = \alpha(\beta f(x)),$ |
| Identity (multiplication) | $(1 f(x)) = f(x).$ |

## Examples of function spaces

**Sets of functions as vector spaces:**

- all linear functions (finite dimensional),
- all polynomials (infinite dimensional),
- given a set of functions $\{\phi_1, \ldots, \phi_D\}$, they generate an $D$-dimensional vector space by taking all linear combinations:

$$\mathcal{F} = \mathrm{span}\{\phi_1, \ldots, \phi_D\} := \Big\{ \sum_{i=1}^{D} \alpha_i \phi_i \,\Big|\, \alpha_i \in \mathbb{R}, \quad i = 1, \ldots, D \Big\}.$$

$\implies$ given that the functions are linearly independent,

$$\sum_{i=1}^{D} c_i \phi_i(x) = 0, \quad \forall x \in \mathcal{X}, \quad \implies \quad c_i = 0, \quad i = 1, \ldots, D,$$

$\{\phi_1, \ldots, \phi_D\}$ is then also a basis of $\mathcal{F}$ (by definition).

### Definition

A real vector space $V$ is called an **inner product space** if there is a function
$\langle \cdot, \cdot \rangle : V \times V \to \mathbb{R}$ that satisfies the following four conditions $\forall x, y, z \in V$ and $\forall \alpha \in \mathbb{R}$:

1. $\langle x, x \rangle \geq 0$ and $\langle x, x \rangle = 0$ if and only if $x = 0$,

2. $\langle x, y + z \rangle = \langle x, y \rangle + \langle x, z \rangle$,

3. $\langle x, \alpha y \rangle = \alpha \langle x, y \rangle$,

4. $\langle x, y \rangle = \langle y, x \rangle$.

The function $\langle \cdot, \cdot \rangle$ is called **inner product**.

Every inner product defines a norm, $\|x\| := \sqrt{\langle x, x \rangle}$, and a metric,
$d(x, y) = \|x - y\| = \sqrt{\langle x - y, x - y \rangle}$.
On inner product spaces we have the **Cauchy-Schwarz inequality:**

$$|\langle x, y \rangle| \leq \|x\| \, \|y\|.$$

## Example - Inner product spaces

**An inner product on functions:** Let $f, g : \mathcal{X} \to \mathbb{R}$,

$$\langle f, g \rangle := \int_{\mathcal{X}} f(x)g(x)dx.$$

We obtain:

- $\langle f, f \rangle = \int_{\mathcal{X}} \big(f(x)\big)^2 dx \geq 0$,
- $\langle f, f \rangle = 0$ if and only if $f = 0$ (almost everywhere),
- $\langle f, g + h \rangle = \int_{\mathcal{X}} f(x)\big(g(x) + h(x)\big)dx = \int_{\mathcal{X}} f(x)g(x)dx + \int_{\mathcal{X}} f(x)h(x)dx$,
- $\langle f, \alpha g \rangle = \int_{\mathcal{X}} f(x)(\alpha\, g)(x)dx = \alpha \int_{\mathcal{X}} f(x)g(x)dx = \alpha \langle f, g \rangle$,
- $\langle f, g \rangle = \int_{\mathcal{X}} f(x)g(x)dx = \int_{\mathcal{X}} g(x)f(x)dx = \langle g, f \rangle$.

Induced norm

$$\|f\| = \sqrt{\langle f, f \rangle} = \Big( \int_{\mathcal{X}} \big(f(x)\big)^2 dx \Big)^{\frac{1}{2}}.$$

## Hilbert Space I

### Definition

A complete inner product space is called a **Hilbert space**.

**The space $L_2(\mathcal{X})$ of square-integrable functions:**

$$L_2(\mathcal{X}) := \Big\{ f : \mathcal{X} \to \mathbb{R} \, \Big| \, \int_{\mathcal{X}} \big( f(x) \big)^2 dx < \infty \Big\},$$

is a Hilbert space together with the inner product,

$$\langle f, g \rangle := \int_{\mathcal{X}} f(x) g(x) dx.$$

**but !**

- $\|f\| = 0$ if and only if $f$ is zero almost everywhere !
- Functions which agree almost everywhere are identified (the relation equal almost everywhere defines equivalence classes of functions),
- $L_2(\mathcal{X})$ is not a space of pointwise defined functions !

**Definition**

If $S$ is an orthonormal set in a Hilbert space $\mathcal{H}$ and no other orthonormal set contains $S$ as a proper subset, then $S$ is called an **orthonormal basis** (or a **complete orthonormal system**) for $\mathcal{H}$.

**Theorem**

Let $\mathcal{H}$ be a Hilbert space and $S = \{x_\alpha\}_{\alpha \in A}$ an orthonormal basis. Then for each $y \in \mathcal{H}$,

$$y = \sum_{\alpha \in A} \langle x_\alpha, y \rangle \, x_\alpha, \qquad \|y\|^2 = \sum_{\alpha \in A} |\langle x_\alpha, y \rangle|^2.$$

The coefficients $\langle x_\alpha, y \rangle$ are called the **coefficients** of $y$ with respect to the basis $\{x_\alpha\}$.

## Construction of RKHSs I

**Steps for the construction of a RKHS:**

- consider the set of all finite linear combinations of the kernel:

$$\mathcal{G} = \mathrm{Span}\{k(x, .) : x \in \mathcal{X}\}$$

- Let $f(x) = \sum_i a_i k(x_i, x)$ and $g(x) = \sum_j b_j k(z_j, x)$. Then

$$\left\langle \sum_i a_i k(x_i, .), \sum_j b_j k(z_j, .) \right\rangle_{\mathcal{G}} := \sum_{i,j} a_i b_j k(x_i, z_j).$$

- check that $\langle \cdot, \cdot \rangle$ is well-defined.

$$\sum_{i,j} a_i b_j k(x_i, z_j) = \sum_i a_i g(x_i) = \sum_j b_j f(z_j)$$

The value of the inner product does not depend on the expansion of $f$ or $g \Rightarrow$ (semi)-inner product with the reproducing property on $\mathcal{G}$.

**Steps for the construction of a RKHS:**

- construct the **semi-norm** associated to this inner product,

$$\|f\|_{\mathcal{G}}^2 = \sum_{i,j=1}^{n} a_i a_j k(x_i, x_j).$$

The Cauchy-Schwarz inequality holds also on semi-inner product spaces,

$$|f(x)| = |\langle f, k(x,.)\rangle_{\mathcal{G}}| \leq \|f\|_{\mathcal{G}} \|k(x,.)\|_{\mathcal{G}} = \|f\|_{\mathcal{G}} \sqrt{k(x,x)}.$$

$\|f\|_{\mathcal{G}} = 0$ implies $f \equiv 0 \Rightarrow$ inner product on $\mathcal{G}$ and $\mathcal{G}$ is an **inner product space**.

- Standard completion by adding all limits of Cauchy sequences in $\mathcal{G}$
  - one has to check that the inner product as well as the reproducing property carries over to the limit elements.

$$\langle \underbrace{k(x,\cdot)}_{\phi(x)}, \underbrace{k(y,\cdot)}_{\phi(y)} \rangle = k(x,y)$$

$$\langle k(x,\cdot), f \rangle = \langle \underbrace{k(x,\cdot)}_{\phi(x)}, \sum_i a_i \underbrace{k(x_i,\cdot)}_{\phi(x_i)} \rangle$$

$$= \underbrace{\sum_i a_i\, k(x_i,\cdot)}_{f}$$

$$= f(x)$$