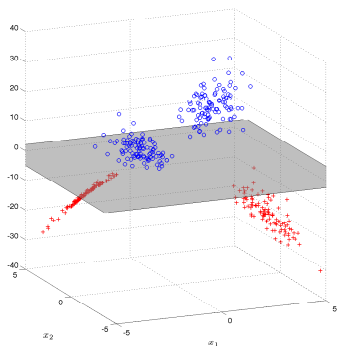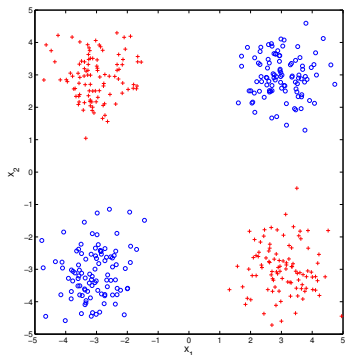# Lecture 1: Introduction to RKHS
## MLSS Tübingen, 2015

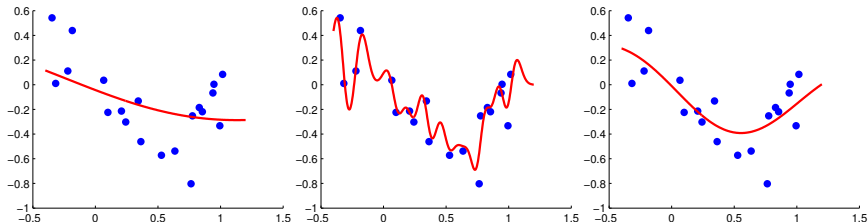Gatsby Unit, CSML, UCL

July 22, 2015

# Kernels and feature space (1): XOR example



- No linear classifier separates red from blue
- Map points to **higher dimensional feature space**:
  $$\phi(x) = \begin{bmatrix} x_1 & x_2 & x_1 x_2 \end{bmatrix} \in \mathbb{R}^3$$

# Kernels and feature space (2): smoothing



Kernel methods can control **smoothness** and **avoid overfitting/underfitting**.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Outline: reproducing kernel Hilbert space

We will describe in order:

1. Hilbert space
2. Kernel (lots of examples: e.g. you can build kernels from simpler kernels)
3. Reproducing property

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Hilbert space

## Definition (Inner product)

Let $\mathcal{H}$ be a vector space over $\mathbb{R}$. A function $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \to \mathbb{R}$ is an inner product on $\mathcal{H}$ if

1. Linear: $\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle_{\mathcal{H}} = \alpha_1 \langle f_1, g \rangle_{\mathcal{H}} + \alpha_2 \langle f_2, g \rangle_{\mathcal{H}}$
2. Symmetric: $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$
3. $\langle f, f \rangle_{\mathcal{H}} \geq 0$ and $\langle f, f \rangle_{\mathcal{H}} = 0$ if and only if $f = 0$.

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Hilbert space

## Definition (Inner product)

Let $\mathcal{H}$ be a vector space over $\mathbb{R}$. A function $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \to \mathbb{R}$ is an inner product on $\mathcal{H}$ if

1. Linear: $\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle_{\mathcal{H}} = \alpha_1 \langle f_1, g \rangle_{\mathcal{H}} + \alpha_2 \langle f_2, g \rangle_{\mathcal{H}}$

2. Symmetric: $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$

3. $\langle f, f \rangle_{\mathcal{H}} \geq 0$ and $\langle f, f \rangle_{\mathcal{H}} = 0$ if and only if $f = 0$.

Norm induced by the inner product: $\|f\|_{\mathcal{H}} := \sqrt{\langle f, f \rangle_{\mathcal{H}}}$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Hilbert space

## Definition (Inner product)

Let $\mathcal{H}$ be a vector space over $\mathbb{R}$. A function $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \to \mathbb{R}$ is an inner product on $\mathcal{H}$ if

1. Linear: $\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle_{\mathcal{H}} = \alpha_1 \langle f_1, g \rangle_{\mathcal{H}} + \alpha_2 \langle f_2, g \rangle_{\mathcal{H}}$
2. Symmetric: $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$
3. $\langle f, f \rangle_{\mathcal{H}} \geq 0$ and $\langle f, f \rangle_{\mathcal{H}} = 0$ if and only if $f = 0$.

Norm induced by the inner product: $\|f\|_{\mathcal{H}} := \sqrt{\langle f, f \rangle_{\mathcal{H}}}$

## Definition (Hilbert space)

Inner product space containing Cauchy sequence limits.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Kernel

### Definition

Let $\mathcal{X}$ be a non-empty set. A function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a **kernel** if there exists an $\mathbb{R}$-Hilbert space and a map $\phi : \mathcal{X} \to \mathcal{H}$ such that $\forall x, x' \in \mathcal{X}$,

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

- Almost no conditions on $\mathcal{X}$ (eg, $\mathcal{X}$ itself doesn't need an inner product, eg. documents).
- A single kernel can correspond to several possible features. A trivial example for $\mathcal{X} := \mathbb{R}$:

$$\phi_1(x) = x \qquad \text{and} \qquad \phi_2(x) = \left[ \begin{array}{c} x/\sqrt{2} \\ x/\sqrt{2} \end{array} \right]$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
**Constructing new kernels**
Positive definite functions
Reproducing kernel Hilbert space

# New kernels from old: sums, transformations

### Theorem (Sums of kernels are kernels)

*Given $\alpha > 0$ and $k$, $k_1$ and $k_2$ all kernels on $\mathcal{X}$, then $\alpha k$ and $k_1 + k_2$ are kernels on $\mathcal{X}$.*

(Proof via positive definiteness: later!) A difference of kernels may not be a kernel (**why?**)

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
**Constructing new kernels**
Positive definite functions
Reproducing kernel Hilbert space

# New kernels from old: sums, transformations

## Theorem (Sums of kernels are kernels)

*Given $\alpha > 0$ and $k$, $k_1$ and $k_2$ all kernels on $\mathcal{X}$, then $\alpha k$ and $k_1 + k_2$ are kernels on $\mathcal{X}$.*

(Proof via positive definiteness: later!) A difference of kernels may not be a kernel (**why?**)

## Theorem (Mappings between spaces)

*Let $\mathcal{X}$ and $\widetilde{\mathcal{X}}$ be sets, and define a map $A : \mathcal{X} \to \widetilde{\mathcal{X}}$. Define the kernel $k$ on $\widetilde{\mathcal{X}}$. Then the kernel $k(A(x), A(x'))$ is a kernel on $\mathcal{X}$.*

Example: $k(x, x') = x^2 (x')^2$.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
**Constructing new kernels**
Positive definite functions
Reproducing kernel Hilbert space

# New kernels from old: products

### Theorem (Products of kernels are kernels)

*Given $k_1$ on $\mathcal{X}_1$ and $k_2$ on $\mathcal{X}_2$, then $k_1 \times k_2$ is a kernel on $\mathcal{X}_1 \times \mathcal{X}_2$. If $\mathcal{X}_1 = \mathcal{X}_2 = \mathcal{X}$, then $k := k_1 \times k_2$ is a kernel on $\mathcal{X}$.*

**Proof:** Main idea only!

$\mathcal{H}_1$ space of kernels between **shapes**,

$$\phi_1(x) = \begin{bmatrix} \mathbb{I}_\square \\ \mathbb{I}_\triangle \end{bmatrix} \qquad \phi_1(\square) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \qquad k_1(\square, \triangle) = 0.$$

$\mathcal{H}_2$ space of kernels between **colors**,

$$\phi_2(x) = \begin{bmatrix} \mathbb{I}_\bullet \\ \mathbb{I}_\bullet \end{bmatrix} \qquad \phi_2(\bullet) = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \qquad k_2(\bullet, \bullet) = 1.$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
**Constructing new kernels**
Positive definite functions
Reproducing kernel Hilbert space

# New kernels from old: products

"Natural" feature space for **colored shapes**:

$$\Phi(x) = \left[ \begin{array}{cc} \mathbb{I}_\square & \mathbb{I}_\triangle \\ \mathbb{I}_\square & \mathbb{I}_\triangle \end{array} \right] = \left[ \begin{array}{c} \mathbb{I}_\bullet \\ \mathbb{I}_\bullet \end{array} \right] \left[ \begin{array}{cc} \mathbb{I}_\square & \mathbb{I}_\triangle \end{array} \right] = \phi_2(x)\phi_1^\top(x)$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
**Constructing new kernels**
Positive definite functions
Reproducing kernel Hilbert space

# New kernels from old: products

"Natural" feature space for **colored shapes**:

$$\Phi(x) = \left[ \begin{array}{cc} \mathbb{I}_{\square} & \mathbb{I}_{\triangle} \\ \mathbb{I}_{\square} & \mathbb{I}_{\triangle} \end{array} \right] = \left[ \begin{array}{c} \mathbb{I}_{\bullet} \\ \mathbb{I}_{\bullet} \end{array} \right] \left[ \begin{array}{cc} \mathbb{I}_{\square} & \mathbb{I}_{\triangle} \end{array} \right] = \phi_2(x)\phi_1^{\top}(x)$$

Kernel is:

$$k(x, x') = \sum_{i \in \{\bullet, \bullet\}} \sum_{j \in \{\square, \triangle\}} \Phi_{ij}(x)\Phi_{ij}(x') = \mathrm{tr}\left( \phi_1(x) \underbrace{\phi_2^{\top}(x)\phi_2(x')}_{k_2(x,x')} \phi_1^{\top}(x') \right)$$

$$= \mathrm{tr}\left( \underbrace{\phi_1^{\top}(x')\phi_1(x)}_{k_1(x,x')} \right) k_2(x, x') = k_1(x, x')k_2(x, x')$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Sums and products $\implies$ polynomials

### Theorem (Polynomial kernels)

*Let $x, x' \in \mathbb{R}^d$ for $d \geq 1$, and let $m \geq 1$ be an integer and $c \geq 0$ be a positive real. Then*

$$k(x, x') := (\langle x, x' \rangle + c)^m$$

*is a valid kernel.*

**To prove**: expand into a sum (with non-negative scalars) of kernels $\langle x, x' \rangle$ raised to integer powers. These individual terms are valid kernels by the product rule.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Infinite sequences

The kernels we've seen so far are dot products between finitely many features. E.g.

$$k(x, y) = \begin{bmatrix} \sin(x) & x^3 & \log x \end{bmatrix}^\top \begin{bmatrix} \sin(y) & y^3 & \log y \end{bmatrix}$$

where $\phi(x) = \begin{bmatrix} \sin(x) & x^3 & \log x \end{bmatrix}$

Can a kernel be a dot product between infinitely many features?

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Infinite sequences

### Definition

The space $\ell_2$ (**square** summable sequences) comprises all sequences $a := (a_i)_{i \geq 1}$ for which

$$\|a\|_{\ell_2}^2 = \sum_{i=1}^{\infty} a_i^2 < \infty.$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Infinite sequences

### Definition

The space $\ell_2$ (**square** summable sequences) comprises all sequences $a := (a_i)_{i \geq 1}$ for which

$$\|a\|_{\ell_2}^2 = \sum_{i=1}^{\infty} a_i^2 < \infty.$$

### Definition

Given sequence of functions $(\phi_i(x))_{i \geq 1}$ in $\ell_2$ where $\phi_i : \mathcal{X} \to \mathbb{R}$ is the $i$th coordinate of $\phi(x)$. Then

$$k(x, x') := \sum_{i=1}^{\infty} \phi_i(x)\phi_i(x') \tag{1}$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
**Constructing new kernels**
Positive definite functions
Reproducing kernel Hilbert space

# Infinite sequences (proof)

Why square summable? By Cauchy-Schwarz,

$$\left| \sum_{i=1}^{\infty} \phi_i(x)\phi_i(x') \right| \leq \|\phi(x)\|_{\ell_2} \|\phi(x')\|_{\ell_2},$$

so the sequence defining the inner product converges for all $x, x' \in \mathcal{X}$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
**Constructing new kernels**
Positive definite functions
Reproducing kernel Hilbert space

## Taylor series kernels

### Definition (Taylor series kernel)

For $r \in (0, \infty]$, with $a_n \geq 0$ for all $n \geq 0$

$$f(z) = \sum_{n=0}^{\infty} a_n z^n \qquad |z| < r, \ z \in \mathbb{R},$$

Define $\mathcal{X}$ to be the $\sqrt{r}$-ball in $\mathbb{R}^d$, so $\|x\| < \sqrt{r}$,

$$k(x, x') = f\left(\langle x, x' \rangle\right) = \sum_{n=0}^{\infty} a_n \langle x, x' \rangle^n.$$

### Example (Exponential kernel)

$$k(x, x') := \exp\left(\langle x, x' \rangle\right).$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Taylor series kernel (proof)

Proof: Non-negative weighted sums of kernels are kernels, and products of kernels are kernels, so the following is a kernel **if it converges**:

$$k(x, x') = \sum_{n=0}^{\infty} a_n \left( \langle x, x' \rangle \right)^n$$

By Cauchy-Schwarz,

$$\left| \langle x, x' \rangle \right| \leq \|x\| \|x'\| < r,$$

so the sum converges.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Gaussian kernel

### Example (Gaussian kernel)

The Gaussian kernel on $\mathbb{R}^d$ is defined as

$$k(x, x') := \exp\left(-\gamma^{-2} \left\| x - x' \right\|^2\right).$$

**Proof**: an exercise! Use product rule, mapping rule, exponential kernel.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Positive definite functions

If we are given a function of two arguments, $k(x, x')$, how can we determine if it is a valid kernel?

1. Find a feature map?
   1. Sometimes this is not obvious (eg if the feature vector is infinite dimensional, e.g. the Gaussian kernel in the last slide)
   2. The feature map is not unique.

2. A direct property of the function: positive definiteness.

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
**Positive definite functions**
Reproducing kernel Hilbert space

# Positive definite functions

### Definition (Positive definite functions)

A symmetric function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is positive definite if
$\forall n \geq 1$, $\forall (a_1, \ldots a_n) \in \mathbb{R}^n$, $\forall (x_1, \ldots, x_n) \in \mathcal{X}^n$,

$$\sum_{i=1}^{n} \sum_{j=1}^{n} a_i a_j k(x_i, x_j) \geq 0.$$

The function $k(\cdot, \cdot)$ is strictly positive definite if for mutually distinct $x_i$, the equality holds only when all the $a_i$ are zero.

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
**Positive definite functions**
Reproducing kernel Hilbert space

# Kernels are positive definite

## Theorem

*Let $\mathcal{H}$ be a Hilbert space, $\mathcal{X}$ a non-empty set and $\phi : \mathcal{X} \to \mathcal{H}$. Then $\langle \phi(x), \phi(y) \rangle_{\mathcal{H}} =: k(x, y)$ is positive definite.*

## Proof.

$$
\begin{aligned}
\sum_{i=1}^{n} \sum_{j=1}^{n} a_i a_j k(x_i, x_j) &= \sum_{i=1}^{n} \sum_{j=1}^{n} \langle a_i \phi(x_i), a_j \phi(x_j) \rangle_{\mathcal{H}} \\
&= \left\| \sum_{i=1}^{n} a_i \phi(x_i) \right\|_{\mathcal{H}}^2 \geq 0.
\end{aligned}
$$

Reverse also holds: positive definite $k(x, x')$ is inner product in a unique $\mathcal{H}$ (Moore-Aronsajn: coming later!).

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Sum of kernels is a kernel

Consider two kernels $k_1(x, x')$ and $k_2(x, x')$. Then

$$\sum_{i=1}^{n} \sum_{j=1}^{n} a_i a_j \left[ k_1(x_i, x_j) + k_2(x_i, x_j) \right]$$
$$= \sum_{i=1}^{n} \sum_{j=1}^{n} a_i a_j k_1(x_i, x_j) + \sum_{i=1}^{n} \sum_{j=1}^{n} a_i a_j k_2(x_i, x_j)$$
$$\geq 0$$

# The reproducing kernel Hilbert space

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# First example: finite space, polynomial features

Reminder: XOR example:

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# First example: finite space, polynomial features

Reminder: Feature space from XOR motivating example:

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$x = \left[ \begin{array}{c} x_1 \\ x_2 \end{array} \right] \mapsto \phi(x) = \left[ \begin{array}{c} x_1 \\ x_2 \\ x_1 x_2 \end{array} \right],$$

with kernel

$$k(x, y) = \left[ \begin{array}{c} x_1 \\ x_2 \\ x_1 x_2 \end{array} \right]^\top \left[ \begin{array}{c} y_1 \\ y_2 \\ y_1 y_2 \end{array} \right]$$

(the standard inner product in $\mathbb{R}^3$ between features). Denote this feature space by $\mathcal{H}$.

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# First example: finite space, polynomial features

Define a linear function of the inputs $x_1, x_2$, and their product $x_1 x_2$,

$$f(x) = f_1 x_1 + f_2 x_2 + f_3 x_1 x_2.$$

$f$ in a space of functions mapping from $\mathcal{X} = \mathbb{R}^2$ to $\mathbb{R}$. Equivalent representation for $f$,

$$f(\cdot) = \begin{bmatrix} f_1 & f_2 & f_3 \end{bmatrix}^\top.$$

$f(\cdot)$ refers to the function as an object (here as a vector in $\mathbb{R}^3$)
$f(x) \in \mathbb{R}$ is function evaluated at a point (a real number).

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# First example: finite space, polynomial features

Define a linear function of the inputs $x_1, x_2$, and their product $x_1 x_2$,

$$f(x) = f_1 x_1 + f_2 x_2 + f_3 x_1 x_2.$$

$f$ in a space of functions mapping from $\mathcal{X} = \mathbb{R}^2$ to $\mathbb{R}$. Equivalent representation for $f$,

$$f(\cdot) = \begin{bmatrix} f_1 & f_2 & f_3 \end{bmatrix}^\top.$$

$f(\cdot)$ refers to the function as an object (here as a vector in $\mathbb{R}^3$)
$f(x) \in \mathbb{R}$ is function evaluated at a point (a real number).

$$f(x) = f(\cdot)^\top \phi(x) = \langle f(\cdot), \phi(x) \rangle_{\mathcal{H}}$$

Evaluation of $f$ at $x$ is an **inner product in feature space** (here standard inner product in $\mathbb{R}^3$)
$\mathcal{H}$ is a space of functions mapping $\mathbb{R}^2$ to $\mathbb{R}$.

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# First example: finite space, polynomial features

I give you a vector:

$$g(\cdot) = [\ 1 \quad -1 \quad -1\ ]$$

Is this a function? Or is it a feature map $\phi(y) = [\ y_1 \quad y_2 \quad y_1 y_2\ ]$?

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# First example: finite space, polynomial features

I give you a vector:

$$g(\cdot) = [ \begin{array}{ccc} 1 & -1 & -1 \end{array} ]$$

Is this a function? Or is it a feature map $\phi(y) = [ \begin{array}{ccc} y_1 & y_2 & y_1 y_2 \end{array} ]$?
Both! All feature maps are also functions.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## First example: finite space, polynomial features

I give you a vector:

$$g(\cdot) = \begin{bmatrix} 1 & -1 & -1 \end{bmatrix}$$

Is this a function? Or is it a feature map $\phi(y) = \begin{bmatrix} y_1 & y_2 & y_1 y_2 \end{bmatrix}$?
Both! All feature maps are also functions.
I give you a vector:

$$h(\cdot) = \begin{bmatrix} 1 & -1 & 2 \end{bmatrix}$$

Is this a function or a feature map?

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# First example: finite space, polynomial features

I give you a vector:

$$g(\cdot) = \left[\begin{array}{ccc} 1 & -1 & -1 \end{array}\right]$$

Is this a function? Or is it a feature map $\phi(y) = \left[\begin{array}{ccc} y_1 & y_2 & y_1 y_2 \end{array}\right]$?
Both! All feature maps are also functions.
I give you a vector:

$$h(\cdot) = \left[\begin{array}{ccc} 1 & -1 & 2 \end{array}\right]$$

Is this a function or a feature map?
It is a function but not a feature map.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# First example: finite space, polynomial features

I give you a vector:

$$g(\cdot) = [\ 1 \quad -1 \quad -1\ ]$$

Is this a function? Or is it a feature map $\phi(y) = [\ y_1 \quad y_2 \quad y_1 y_2\ ]$?
Both! All feature maps are also functions.
I give you a vector:

$$h(\cdot) = [\ 1 \quad -1 \quad 2\ ]$$

Is this a function or a feature map?
It is a function but not a feature map.
All feature maps are also functions. But the space of functions is larger: some functions are not feature maps.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# First example: finite space, polynomial features

$\phi(y)$ is a mapping from $\mathbb{R}^2$ to $\mathbb{R}^3\ldots$

...which also parametrizes a *function* mapping $\mathbb{R}^2$ to $\mathbb{R}$.

$$k(\cdot, y) := \begin{bmatrix} y_1 & y_2 & y_1 y_2 \end{bmatrix}^\top = \phi(y),$$

We can *evaluate* this function at $x$

$$\langle k(\cdot, y), \phi(x) \rangle_{\mathcal{H}} = ax_1 + bx_2 + cx_1 x_2,$$

where $a = y_1$, $b = y_2$, and $c = y_1 y_2$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# First example: finite space, polynomial features

$\phi(y)$ is a mapping from $\mathbb{R}^2$ to $\mathbb{R}^3$...
...which also parametrizes a function mapping $\mathbb{R}^2$ to $\mathbb{R}$.

$$k(\cdot, y) := \begin{bmatrix} y_1 & y_2 & y_1 y_2 \end{bmatrix}^\top = \phi(y),$$

We can *evaluate* this function at $x$

$$\langle k(\cdot, y), \phi(x) \rangle_{\mathcal{H}} = a x_1 + b x_2 + c x_1 x_2,$$

where $a = y_1$, $b = y_2$, and $c = y_1 y_2$
...but due to symmetry,

$$\begin{aligned} \langle k(\cdot, x), \phi(y) \rangle &= u y_1 + v y_2 + w y_1 y_2 \\ &= k(x, y). \end{aligned}$$

We can write $\phi(x) = k(\cdot, x)$ and $\phi(y) = k(\cdot, y)$ without ambiguity:
canonical feature map

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

## The kernel trick



# Statistics Professors <u>HATE</u> Him!

*Doctor's discovery revealed the secret to learning any problem with just 10 training samples. Watch this shocking video and learn how rapidly you can find a solution to your learning problems using this one sneaky kernel trick! Free from overfitting!*
   *http://www.oneweirdkerneltrick.com*

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# The kernel trick

This example illustrates the two defining features of an RKHS:

- **The reproducing property:** (kernel trick)
  $\forall x \in \mathcal{X}, \forall f(\cdot) \in \mathcal{H}, \ \langle f(\cdot), k(\cdot, x) \rangle_{\mathcal{H}} = f(x)$
  . . .or use shorter notation $\langle f, \phi(x) \rangle_{\mathcal{H}}$.

- In particular, for any $x, y \in \mathcal{X}$,

$$k(x, y) = \langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}}.$$

Note: the feature map of every point is in the feature space:
$\forall x \in \mathcal{X}, \ k(\cdot, x) = \phi(x) \in \mathcal{H},$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# First example: finite space, polynomial features

Another, more subtle point: $\mathcal{H}$ can be larger than all $\phi(x)$.

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# First example: finite space, polynomial features

Another, more subtle point: $\mathcal{H}$ can be larger than all $\phi(x)$.



E.g. $f = [1\,1\,-1] \in \mathcal{H}$ cannot be obtained by $\phi(x) = [x_1\,x_2\,(x_1 x_2)]$.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Second example: infinite feature space

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Second example: infinite feature space

Reproducing property for function with Gaussian kernel:
$$f(x) := \sum_{i=1}^{m} \alpha_i k(x_i, x) = \langle \sum_{i=1}^{m} \alpha_i \phi(x_i), \phi(x) \rangle_{\mathcal{H}}.$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Second example: infinite feature space

Reproducing property for function with Gaussian kernel:
$$f(x) := \sum_{i=1}^{m} \alpha_i k(x_i, x) = \langle \sum_{i=1}^{m} \alpha_i \phi(x_i), \phi(x) \rangle_{\mathcal{H}}.$$



- What do the features $\phi(x)$ look like (warning: there are infinitely many of them!)
- What do these features have to do with smoothness?

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Second example: infinite feature space

Under certain conditions (Mercer's theorem and extensions), we can write

$$k(x, x') = \sum_{i=1}^{\infty} \lambda_i e_i(x) e_i(x'), \qquad \int_{\mathcal{X}} e_i(x) e_j(x) d\mu(x) = \begin{cases} 1 & i = j \\ 0 & i \neq j. \end{cases}$$

where this sum is guaranteed to converge whatever the $x$ and $x'$.

Infinite dimensional feature map: $\qquad \phi(x) = \begin{bmatrix} \vdots \\ \sqrt{\lambda_i} e_i(x) \\ \vdots \end{bmatrix} \in \ell_2.$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Second example: infinite feature space

Under certain conditions (Mercer's theorem and extensions), we can write

$$k(x, x') = \sum_{i=1}^{\infty} \lambda_i e_i(x) e_i(x'), \qquad \int_{\mathcal{X}} e_i(x) e_j(x) d\mu(x) = \begin{cases} 1 & i = j \\ 0 & i \neq j. \end{cases}$$

where this sum is guaranteed to converge whatever the $x$ and $x'$.

Infinite dimensional feature map: $\quad \phi(x) = \begin{bmatrix} \vdots \\ \sqrt{\lambda_i} e_i(x) \\ \vdots \end{bmatrix} \in \ell_2.$

Define $\mathcal{H}$ to be the space of functions: for $\{f_i\}_{i=1}^{\infty} \in \ell_2$,

$$f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}} = \sum_{i=1}^{\infty} f_i \sqrt{\lambda_i} e_i(x).$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Second example: infinite feature space

Gaussian kernel, $k(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$,

$$\begin{aligned}
\lambda_k &\propto b^k \qquad b < 1 \\
e_k(x) &\propto \exp(-(c - a)x^2)H_k(x\sqrt{2c}),
\end{aligned}$$

$a, b, c$ are functions of $\sigma$, and $H_k$ is $k$th order Hermite polynomial.



$$k(x, x') = \sum_{i=1}^{\infty} \lambda_i e_i(x) e_i(x')$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Second example: infinite feature space

Example RKHS function, Gaussian kernel:

$$f(x) := \sum_{i=1}^{m} \alpha_i k(x_i, x) = \sum_{i=1}^{m} \alpha_i \left[ \sum_{j=1}^{\infty} \lambda_j e_j(x_i) e_j(x) \right] = \sum_{j=1}^{\infty} f_j \underbrace{\left[ \sqrt{\lambda_j} e_j(x) \right]}_{\phi_j(x)}$$

where $f_j = \sum_{i=1}^{m} \alpha_i \sqrt{\lambda_j} e_j(x_i)$.



NOTE that this enforces smoothing: $\lambda_j$ decay as $e_j$ become rougher, $f_j$ decay since $\sum_j f_j^2 < \infty$.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Third (infinite) example: fourier series

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Third (infinite) example: fourier series

Function on the torus $\mathbb{T} := [-\pi, \pi]$ with periodic boundary. Fourier series:

$$f(x) = \sum_{\ell=-\infty}^{\infty} \hat{f}_\ell \exp(\imath \ell x) = \sum_{l=-\infty}^{\infty} \hat{f}_\ell \left( \cos(\ell x) + \imath \sin(\ell x) \right).$$

Example: "top hat" function,

$$f(x) = \begin{cases} 1 & |x| < T, \\ 0 & T \leq |x| < \pi. \end{cases}$$

Fourier series:

$$\hat{f}_\ell := \frac{\sin(\ell T)}{\ell \pi} \qquad f(x) = \sum_{\ell=0}^{\infty} 2\hat{f}_\ell \cos(\ell x).$$

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# Fourier series for top hat function

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Fourier series for top hat function

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Fourier series for top hat function

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# Fourier series for top hat function

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# Fourier series for top hat function

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Fourier series for top hat function

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# Fourier series for top hat function

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Fourier series for kernel function

Kernel takes a single argument,

$$k(x, y) = k(x - y),$$

Define the Fourier series representation of $k$

$$k(x) = \sum_{\ell=-\infty}^{\infty} \hat{k}_\ell \exp\left(\imath \ell x\right),$$

$k$ and its Fourier transform are real and symmetric. E.g. ,

$$k(x) = \frac{1}{2\pi} \vartheta\left(\frac{x}{2\pi}, \frac{\imath \sigma^2}{2\pi}\right), \qquad \hat{k}_\ell = \frac{1}{2\pi} \exp\left(\frac{-\sigma^2 \ell^2}{2}\right).$$

$\vartheta$ is the Jacobi theta function, close to Gaussian when $\sigma^2$ sufficiently narrower than $[-\pi, \pi]$.

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
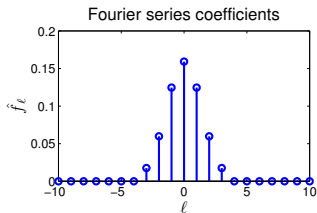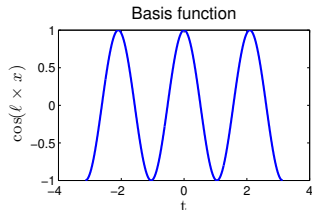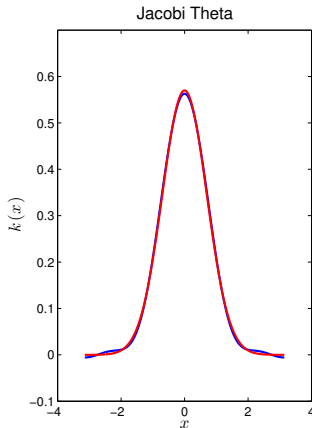Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# Fourier series for Gaussian-spectrum kernel

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression
What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# Fourier series for Gaussian-spectrum kernel

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Fourier series for Gaussian-spectrum kernel

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# Fourier series for Gaussian-spectrum kernel

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Feature space via fourier series

Define $\mathcal{H}$ to be the space of functions with (infinite) feature space representation

$$f(\cdot) = \begin{bmatrix} \ldots & \hat{f}_\ell / \sqrt{\hat{k}_\ell} & \ldots \end{bmatrix}^\top.$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Feature space via fourier series

Define $\mathcal{H}$ to be the space of functions with (infinite) feature space representation

$$f(\cdot) = \left[ \ \ldots \ \ \hat{f}_\ell / \sqrt{\hat{k}_\ell} \ \ \ldots \ \right]^\top.$$

Define the feature map

$$k(\cdot, x) = \phi(x) = \left[ \ \ldots \ \ \sqrt{\hat{k}_\ell} \exp(-\imath \ell x) \ \ \ldots \ \right]^\top$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Feature space via fourier series

The reproducing theorem holds,

$$
\begin{aligned}
\langle f(\cdot), k(\cdot, x) \rangle_{\mathcal{H}} &= \sum_{\ell=-\infty}^{\infty} \left( \frac{\hat{f}_\ell}{\sqrt{\hat{k}_\ell}} \right) \overline{\sqrt{\hat{k}_\ell} \exp(-\imath \ell x)} \\
&= \sum_{\ell=-\infty}^{\infty} \hat{f}_\ell \exp(\imath \ell x) = f(x),
\end{aligned}
$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Feature space via fourier series

The reproducing theorem holds,

$$
\begin{aligned}
\langle f(\cdot), k(\cdot, x) \rangle_{\mathcal{H}} &= \sum_{\ell=-\infty}^{\infty} \left( \frac{\hat{f}_{\ell}}{\sqrt{\hat{k}_{\ell}}} \right) \overline{\sqrt{\hat{k}_{\ell}} \exp(-i\ell x)} \\
&= \sum_{\ell=-\infty}^{\infty} \hat{f}_{\ell} \exp(i\ell x) = f(x),
\end{aligned}
$$

. . .including for the kernel itself,

$$
\begin{aligned}
\langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}} &= \sum_{\ell=-\infty}^{\infty} \left( \sqrt{\hat{k}_{\ell}} \exp(-i\ell x) \right) \left( \overline{\sqrt{\hat{k}_{\ell}} \exp(-i\ell y)} \right) \\
&= \sum_{\ell=-\infty}^{\infty} \hat{k}_{\ell} \exp(i\ell(y - x)) = k(x - y).
\end{aligned}
$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Fourier series and smoothness

The squared norm of a function $f$ in $\mathcal{H}$ is:

$$\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \sum_{l=-\infty}^{\infty} \frac{\hat{f}_\ell \overline{\hat{f}_\ell}}{\hat{k}_\ell}.$$

If $\hat{k}_\ell$ decays fast, then so must $\hat{f}_\ell$ if we want $\|f\|_{\mathcal{H}}^2 < \infty$.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Fourier series and smoothness

The squared norm of a function $f$ in $\mathcal{H}$ is:

$$\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \sum_{l=-\infty}^{\infty} \frac{\hat{f}_\ell \overline{\hat{f}_\ell}}{\hat{k}_\ell}.$$

If $\hat{k}_\ell$ decays fast, then so must $\hat{f}_\ell$ if we want $\|f\|_{\mathcal{H}}^2 < \infty$. Recall

$$f(x) = \sum_{\ell=-\infty}^{\infty} \hat{f}_\ell \left( \cos(\ell x) + \imath \sin(\ell x) \right).$$

Enforces smoothness.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Fourier series and smoothness

The squared norm of a function $f$ in $\mathcal{H}$ is:

$$\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \sum_{l=-\infty}^{\infty} \frac{\hat{f}_\ell \overline{\hat{f}_\ell}}{\hat{k}_\ell}.$$

If $\hat{k}_\ell$ decays fast, then so must $\hat{f}_\ell$ if we want $\|f\|_{\mathcal{H}}^2 < \infty$.
Recall

$$f(x) = \sum_{\ell=-\infty}^{\infty} \hat{f}_\ell \left( \cos(\ell x) + \imath \sin(\ell x) \right).$$

Enforces smoothness.

Question: is the **top hat** function in the Gaussian-spectrum RKHS?

# Some reproducing kernel Hilbert space theory

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Reproducing kernel Hilbert space (1)

## Definition

$\mathcal{H}$ a Hilbert space of $\mathbb{R}$-valued functions on non-empty set $\mathcal{X}$. A function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a reproducing kernel of $\mathcal{H}$, and $\mathcal{H}$ is a reproducing kernel Hilbert space, if

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

In particular, for any $x, y \in \mathcal{X}$,

$$k(x, y) = \langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}}. \qquad (2)$$

Original definition: kernel an inner product between feature maps. Then $\phi(x) = k(\cdot, x)$ a valid feature map.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# Reproducing kernel Hilbert space (2)

Another RKHS definition:

Define $\delta_x$ to be the operator of evaluation at $x$, i.e.

$$\delta_x f = f(x) \quad \forall f \in \mathcal{H},\, x \in \mathcal{X}.$$

## Definition (Reproducing kernel Hilbert space)

$\mathcal{H}$ is an RKHS if the evaluation operator $\delta_x$ is bounded: $\forall x \in \mathcal{X}$ there exists $\lambda_x \geq 0$ such that for all $f \in \mathcal{H}$,

$$|f(x)| = |\delta_x f| \leq \lambda_x \|f\|_{\mathcal{H}}$$

$\implies$ two functions identical in RHKS norm agree at every point:

$$|f(x) - g(x)| = |\delta_x (f - g)| \leq \lambda_x \|f - g\|_{\mathcal{H}} \quad \forall f, g \in \mathcal{H}.$$

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

# RKHS definitions equivalent

### Theorem (Reproducing kernel equivalent to bounded $\delta_x$ )

$\mathcal{H}$ is a reproducing kernel Hilbert space (i.e., its evaluation operators $\delta_x$ are bounded linear operators), if and only if $\mathcal{H}$ has a reproducing kernel.

**Proof:** If $\mathcal{H}$ has a reproducing kernel $\implies \delta_x$ bounded

$$
\begin{aligned}
|\delta_x[f]| &= |f(x)| \\
&= |\langle f, k(\cdot, x)\rangle_{\mathcal{H}}| \\
&\leq \|k(\cdot, x)\|_{\mathcal{H}} \|f\|_{\mathcal{H}} \\
&= \langle k(\cdot, x), k(\cdot, x)\rangle_{\mathcal{H}}^{1/2} \|f\|_{\mathcal{H}} \\
&= k(x, x)^{1/2} \|f\|_{\mathcal{H}}
\end{aligned}
$$

Cauchy-Schwarz in 3rd line . Consequently, $\delta_x : \mathcal{F} \to \mathbb{R}$ bounded with $\lambda_x = k(x, x)^{1/2}$.

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## RKHS definitions equivalent

Proof: $\delta_x$ bounded $\implies \mathcal{H}$ has a reproducing kernel
We use. . .

### Theorem

*(Riesz representation) In a Hilbert space $\mathcal{H}$, all bounded linear functionals are of the form $\langle \cdot, g \rangle_{\mathcal{H}}$, for some $g \in \mathcal{H}$.*

If $\delta_x : \mathcal{F} \to \mathbb{R}$ is a bounded linear functional, by Riesz $\exists f_{\delta_x} \in \mathcal{H}$ such that

$$\delta_x f = \langle f, f_{\delta_x} \rangle_{\mathcal{H}}, \ \forall f \in \mathcal{H}.$$

*Define* $k(x', x) = f_{\delta_x}(x'), \ \forall x, x' \in \mathcal{X}$. By its definition, both $k(\cdot, x) = f_{\delta_x} \in \mathcal{H}$ and $\langle f, k(\cdot, x) \rangle_{\mathcal{H}} = \delta_x f = f(x)$. Thus, $k$ is the reproducing kernel.

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# Moore-Aronszajn Theorem

> **Theorem (Moore-Aronszajn)**
>
> *Let $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ be positive definite. There is a **unique RKHS** $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$ with reproducing kernel $k$.*

Recall feature map is *not* unique (as we saw earlier): only kernel is.

Feature space
**Basics of reproducing kernel Hilbert spaces**
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
**Reproducing kernel Hilbert space**

# Main message #1

Feature space
Basics of reproducing kernel Hilbert spaces
Kernel Ridge Regression

What is a kernel?
Constructing new kernels
Positive definite functions
Reproducing kernel Hilbert space

## Main message #2

Small RKHS norm results in smooth functions.
E.g. kernel ridge regression with Gaussian kernel:

$$f^* \;=\; \arg\min_{f \in \mathcal{H}} \left( \sum_{i=1}^{n} (y_i - \langle f, \phi(x_i) \rangle_{\mathcal{H}})^2 + \lambda \|f\|_{\mathcal{H}}^2 \right).$$



$\lambda=0.1, \sigma=0.6$      $\lambda=10, \sigma=0.6$      $\lambda=1e{-}07, \sigma=0.6$

# Kernel Ridge Regression

# Kernel ridge regression



Very simple to implement, works well when no outliers.

## Kernel ridge regression

Use features of $\phi(x_i)$ in the place of $x_i$:

$$f^* = \arg\min_{f \in \mathcal{H}} \left( \sum_{i=1}^{n} (y_i - \langle f, \phi(x_i) \rangle_{\mathcal{H}})^2 + \lambda \|f\|_{\mathcal{H}}^2 \right).$$

E.g. for finite dimensional feature spaces,

$$\phi_p(x) = \begin{bmatrix} x \\ x^2 \\ \vdots \\ x^\ell \end{bmatrix} \qquad \phi_s(x) = \begin{bmatrix} \sin x \\ \cos x \\ \sin 2x \\ \vdots \\ \cos \ell x \end{bmatrix}$$

$a$ is a vector of length $\ell$ giving weight to each of these features so as to find the mapping between $x$ and $y$. Feature vectors can also have *infinite* length (more soon).

# Kernel ridge regression

Solution easy if we already know $f$ is a linear combination of feature space mappings of points: representer theorem.

$$f = \sum_{i=1}^{n} \alpha_i \phi(x_i) = \sum_{i=1}^{n} \alpha_i k(x_i, \cdot).$$

## Representer theorem

Given a set of paired observations $(x_1, y_1), \ldots (x_n, y_n)$ (regression or classification).
Find the function $f^*$ in the RKHS $\mathcal{H}$ which satisfies

$$J(f^*) = \min_{f \in \mathcal{H}} J(f), \qquad (3)$$

where

$$J(f) = L_y(f(x_1), \ldots, f(x_n)) + \Omega\left(\|f\|_{\mathcal{H}}^2\right),$$

$\Omega$ is non-decreasing, and $y$ is the vector of $y_i$.

- Classification: $L_y(f(x_1), \ldots, f(x_n)) = \sum_{i=1}^n \mathbb{I}_{y_i f(x_i) \leq 0}$
- Regression: $L_y(f(x_1), \ldots, f(x_n)) = \sum_{i=1}^n (y_i - f(x_i))^2$

## Representer theorem

**The representer theorem:**(simple version) solution to

$$\min_{f \in \mathcal{H}} \left[ L_y(f(x_1), \ldots, f(x_n)) + \Omega \left( \|f\|_{\mathcal{H}}^2 \right) \right]$$

takes the form

$$f^* = \sum_{i=1}^{n} \alpha_i k(x_i, \cdot).$$

If $\Omega$ is strictly increasing, all solutions have this form.

## Representer theorem: proof

**Proof:** Denote $f_s$ projection of $f$ onto the subspace

$$\operatorname{span}\left\{k(x_i, \cdot) : \ 1 \leq i \leq n\right\}, \tag{4}$$

such that

$$f = f_s + f_\perp,$$

where $f_s = \sum_{i=1}^{n} \alpha_i k(x_i, \cdot)$.

**Regularizer**:

$$\|f\|_{\mathcal{H}}^2 = \|f_s\|_{\mathcal{H}}^2 + \|f_\perp\|_{\mathcal{H}}^2 \geq \|f_s\|_{\mathcal{H}}^2,$$

then

$$\Omega\left(\|f\|_{\mathcal{H}}^2\right) \geq \Omega\left(\|f_s\|_{\mathcal{H}}^2\right),$$

so this term is minimized for $f = f_s$.

## Representer theorem: proof

**Proof (cont.):** Individual terms $f(x_i)$ in the loss:

$$f(x_i) = \langle f, k(x_i, \cdot) \rangle_{\mathcal{H}} = \langle f_s + f_\perp, k(x_i, \cdot) \rangle_{\mathcal{H}} = \langle f_s, k(x_i, \cdot) \rangle_{\mathcal{H}},$$

so

$$L_y(f(x_1), \ldots, f(x_n)) = L_y(f_s(x_1), \ldots, f_s(x_n)).$$

Hence

- Loss $L(\ldots)$ only depends on the component of $f$ in the data subspace,

- Regularizer $\Omega(\ldots)$ minimized when $f = f_s$.

- If $\Omega$ is strictly non-decreasing, then $\|f_\perp\|_{\mathcal{H}} = 0$ is required at the minimum.

# Kernel ridge regression: proof

We *begin* knowing $f$ is a linear combination of feature space mappings of points (<span style="color:red">representer theorem</span>)

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

Then

$$\sum_{i=1}^{n} (y_i - \langle f, \phi(x_i) \rangle_{\mathcal{H}})^2 + \lambda \|f\|_{\mathcal{H}}^2 \;=\; \|y - K\alpha\|^2 + \lambda \alpha^\top K \alpha$$

Differentiating wrt $\alpha$ and setting this to zero, we get

$$\alpha^* = (K + \lambda I_n)^{-1} y.$$

## Reminder: smoothness

What does $\|a\|_{\mathcal{H}}$ have to do with smoothing?
Example 1: The Fourier series representation on torus $\mathbb{T}$:

$$f(x) = \sum_{l=-\infty}^{\infty} \hat{f}_l \exp(\imath l x),$$

and

$$\langle f, g \rangle_{\mathcal{H}} = \sum_{l=-\infty}^{\infty} \frac{\hat{f}_l \overline{\hat{g}_l}}{\hat{k}_l}.$$
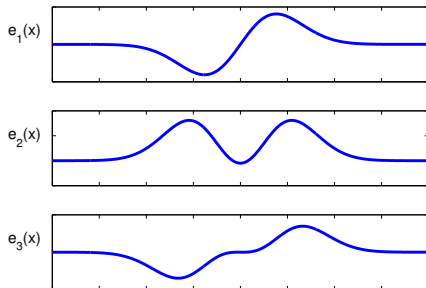
Thus,

$$\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \sum_{l=-\infty}^{\infty} \frac{\left| \hat{f}_l \right|^2}{\hat{k}_l}.$$

# Reminder: smoothness

What does $\|a\|_{\mathcal{H}}$ have to do with smoothing?
Example 2: The Gaussian kernel on $\mathbb{R}$. Recall

$$f(x) = \sum_{i=1}^{\infty} a_i \sqrt{\lambda_i} e_i(x), \qquad \|f\|_{\mathcal{H}}^2 = \sum_{i=1}^{\infty} a_i^2.$$
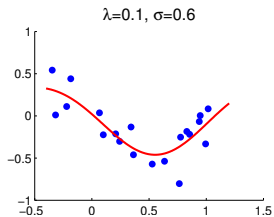
## Parameter selection for KRR

Given the objective

$$f^* \;=\; \arg\min_{f \in \mathcal{H}} \left( \sum_{i=1}^{n} \left( y_i - \langle f, \phi(x_i) \rangle_{\mathcal{H}} \right)^2 + \lambda \|f\|_{\mathcal{H}}^2 \right).$$
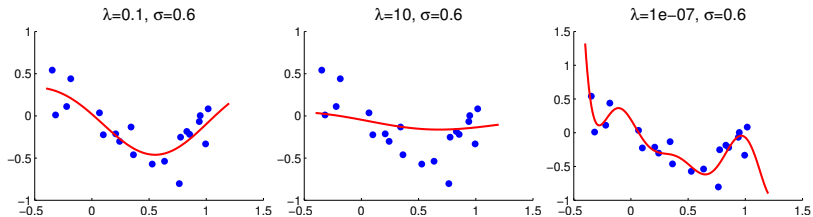
How do we choose

- The regularization parameter $\lambda$?
- The kernel parameter: for Gaussian kernel, $\sigma$ in

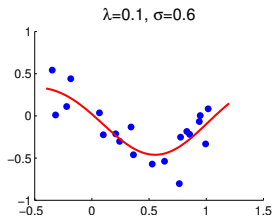$$k(x, y) = \exp\left( \frac{-\|x - y\|^2}{\sigma} \right).$$

# Choice of $\lambda$



$\lambda$=0.1, $\sigma$=0.6

# Choice of $\lambda$

# Choice of $\sigma$



λ=0.1, σ=0.6

# Choice of $\sigma$