A thesis submitted to attain the degree of

# DOCTOR OF COMPUTER SCIENCE

Entitled

# DATA ARE NOT REAL!

Presented by

## ROMAIN BRAULT *

About

## LARGE-SCALE LEARNING ON STRUCTURED INPUT-OUTPUT DATA WITH OPERATOR-VALUED KERNELS .

Under supervision of
Professor (Prof.) FLORENCE D'ALCHÉ-BUC †



| | | |
|---|---|---|
| M. John DOE | UEVE | examinator, |
| M. John DOE | UEVE | director, |
| M. John DOE | UEVE | examinator, |
| M. John DOE | UEVE | examinator, |
| M. John DOE | UEVE | reporter, |

accepted on the recommendation of M. John DOE (UEVE) and M. John DOE (UEVE).

Computer Science
IBISC
Université d'Évry val d'Essonne

Septembre 2016 – version 0.1

* Email: romain.brault@ibisc.fr
† Email: florence.dalche@telecom-paristech.fr

# ABSTRACT

Short summary of the contents...a great guide by Kent Beck how to write good abstracts can be found here:

`https://plg.uwaterloo.ca/~migod/research/beckOOPSLA.html`

## PUBLICATIONS

Some ideas and figures have appeared previously in the following publications:

Put your publications from the thesis here. The packages `multibib` or `bibtopic` etc. can be used to handle multiple different bibliographies in your document.

*We have seen that computer programming is an art,*
*because it applies accumulated knowledge to the world,*
*because it requires skill and ingenuity, and especially*
*because it produces objects of beauty.*

## ACKNOWLEDGEMENTS

---

[1] Members of GuIT (Gruppo Italiano Utilizzatori di TeX e LaTeX)

# CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LISTINGS

## ACRONYMS

OVK  Operator-Valued Kernel.

ORFF  Operator-valued Random Fourier Feature.

POVM  Positive Operator-Valued Measure

RKHS  Reproducing Kernel Hilbert Space.

vv-RKHS  vector-valued Reproducing Kernel Hilbert Space.

LCA  Locally Compact Abelian.

FT  Fourier transform.

IFT  inverse Fourier transform.

Part I

# INTRODUCTION

You can put some informational part preamble text here. Illo principalmente su nos. Non message *occidental* angloromanic da. Debitas effortio simplificate sia se, auxiliar summarios da que, se avantiate publicationes via. Pan in terra summarios, capital interlingua se que. Al via multo esser specimen, campo responder que da. Le usate medical addresses pro, europa origine sanctificate nos se.

# 1

## MOTIVATIONS

# 2

## BACKGROUND

## 2.1 NOTATIONS

The euclidean inner product in $\mathbb{R}^d$ is denoted $\langle \cdot, \cdot \rangle$ and the euclidean norm is denoted $\|\cdot\|$. The unit pure imaginary number $\sqrt{-1}$ is denoted i. $\mathcal{B}(\mathbb{R}^d)$ is the Borel $\sigma$-algebra on $\mathbb{R}^d$. If $\mathcal{X}$ and $\mathcal{Y}$ are two vector spaces, we denote by $\mathcal{F}(\mathcal{X}; \mathcal{Y})$ the vector space of functions $f : \mathcal{X} \to \mathcal{Y}$ and $\mathcal{C}(\mathcal{X}; \mathcal{Y}) \subset \mathcal{F}(\mathcal{X}; \mathcal{Y})$ the subspace of continuous functions. If $\mathcal{H}$ is an Hilbert space we denote its scalar product by $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ and its norm by $\|\cdot\|_{\mathcal{H}}$. We set $\mathcal{L}(\mathcal{H}) = \mathcal{L}(\mathcal{H}; \mathcal{H})$ to be the space of linear operators from $\mathcal{H}$ to itself. If $W \in \mathcal{L}(\mathcal{H})$, Ker $W$ denotes the nullspace, Im $W$ the image and $W^* \in \mathcal{L}(\mathcal{H})$ the adjoint operator (transpose when $W$ is a real matrix). All these notations are summarized in table 1.

## 2.2 ABOUT STATISTICAL LEARNING

## 2.3 ON LARGE-SCALE LEARNING

## 2.4 ELEMENTS OF ABSTRACT HARMONIC ANALYSIS

### 2.4.1 *Locally compact Abelian groups*

**Definition 1.** *Locally Compact Abelian group. A group $(\mathcal{X}, \star)$ is said to be Locally Compact Abelian if it is a topological commutative group $\mathcal{X}$ for which every point has a compact neighborhood and is Hausdorff.*

Locally Compact Abelian (LCA) groups are central to the general definition of Fourier Transform which is related to the concept of Pontryagin duality [9]. Let $(\mathcal{X}, \star)$ be a LCA group with $e$ its neutral element and the notation, $x^{-1}$, for the inverse of $x \in \mathcal{X}$. A *character* is a complex continuous homomorphism $\omega : \mathcal{X} \to \mathbb{U}$ from $\mathcal{X}$ to the set of complex numbers of unit module $\mathbb{U}$. The set of all characters of $\mathcal{X}$ forms the Pontryagin *dual group* $\hat{\mathcal{X}}$. The dual group of an LCA group is an LCA group and the dual group operation is defined by

$$(\omega_1 \star \omega_2)(x) = \omega_1(x)\omega_2(x) \in \mathbb{U}.$$

The Pontryagin duality theorem states that $\hat{\hat{\mathcal{X}}} \cong \mathcal{X}$. I.e. there is a canonical isomorphism between any LCA group and its double dual. To emphasize this duality the following notation is usually adopted: $\omega(x) = (x, \omega) = (\omega, x)$, where $x \in \mathcal{X}$, $\omega \in \hat{\mathcal{X}}$. Another important property involves the complex conjugate of the pairing which is defined as $\overline{(x, \omega)} = (x^{-1}, \omega)$.

Table 1: Mathematical symbols used throughout the parper and their signi-
fication.

| Symbol | Meaning |
| --- | --- |
| $i$ | Unit pure imaginary number $\sqrt{-1}$. |
| $e$ | Euler constant. |
| $\langle \cdot, \cdot \rangle$ | Euclidean inner product. |
| $\|\cdot\|$ | Euclidean norm. |
| $\mathcal{X}$ | Input space (). |
| $\hat{\mathcal{X}}$ | The Pontryagin dual of $\mathcal{X}$. |
| $\mathcal{Y}$ | Output space (Hilbert space). |
| $\mathcal{H}$ | Feature space (Hilbert space). |
| $\langle \cdot, \cdot \rangle_{\mathcal{Y}}$ | The canonical inner product of the Hilbert space $\mathcal{Y}$. |
| $\|\cdot\|_{\mathcal{Y}}$ | The canonical norm induced by the inner product of the Hilbert space $\mathcal{Y}$. |
| $\mathcal{F}(\mathcal{X}; \mathcal{Y})$ | Vector space of function from $\mathcal{X}$ to $\mathcal{Y}$. |
| $\mathcal{C}(\mathcal{X}; \mathcal{Y})$ | The vector subspace of $\mathcal{F}$ of continuous function from $\mathcal{X}$ to $\mathcal{Y}$. |
| $\mathcal{L}(\mathcal{H}; \mathcal{Y})$ | The set of bounded linear operator from a Hilbert space $\mathcal{H}$ to a Hilbert space $\mathcal{Y}$. |
| $\mathcal{L}(\mathcal{Y})$ | The set of bounded linear operator from a Hilbert space $\mathcal{H}$ to itself. |
| $\mathcal{L}_+(\mathcal{Y})$ | The set of non-negative bounded linear operator from a Hilbert space $\mathcal{H}$ to itself. |
| $\mathcal{B}(\mathcal{X})$ | Borel $\sigma$-algebra on $\mathcal{X}$. |
| $\mu(\mathcal{X})$ | A scalar positive measure of $\mathcal{X}$. |
| $p_\mu(x)$ | The Radon-Nikodym derivative of $\mu$ w.r.t. the Lebesgue measure. |
| $dx, d\omega$ | The canonical Haar measure of the LCA group $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$. (resp. $(\hat{\mathcal{X}}, \mathcal{B}(\hat{\mathcal{X}}))$. |
| $L^p(\mathcal{X}, dx)$ | The Banach space of $|\cdot|^p$-integrable function from $(\mathcal{X}, \mathcal{B}(\mathcal{X}), dx)$ to $\mathbb{C}$. |
| $L^p(\mathcal{X}, dx; \mathcal{Y})$ | The Banach space of $\|\cdot\|_{\mathcal{Y}^p}$ (Bochner)-integrable function from $(\mathcal{X}, \mathcal{B}(\mathcal{X}), dx)$ to $\mathcal{Y}$. |

We notice that for any pairing depending of $\omega$, there exists a function $h_\omega : \mathcal{X} \to \mathbb{R}$ such that: $(x, \omega) = \exp(-ih_\omega(x))$ since any pairing maps into $\mathbb{U}$. Moreover,

$$(x \star z^{-1}, \omega) = \omega(x)\omega(z^{-1}) = \exp(-ih_\omega(x))\exp(-ih_\omega(z^{-1}))$$
$$= \exp(-ih_\omega(x))\exp(+ih_\omega(z)).$$

Table 2: Classification of Fourier transforms in terms of their domain and transform domain.

| $\mathcal{X}$ | $\hat{\mathcal{X}}$ | Operation | Pairing |
|---|---|---|---|
| $\mathbb{R}^d$ | $\mathbb{R}^d$ | $+$ | $(x, \omega) = \exp\left(i\langle x, \omega \rangle\right)$ |
| $\mathbb{R}^d_{*,+}$ | $\mathbb{R}^d$ | $\cdot$ | $(x, \omega) = \exp\left(i\langle \log(x), \omega \rangle\right)$ |
| $(-c; +\infty)^d$ | $\mathbb{R}^d$ | $\odot$ | $(x, \omega) = \exp\left(i\langle \log(x+c), \omega \rangle\right)$ |

Table 2 provide an explicit list of pairings for various groups based on $\mathbb{R}^d$ or its subsets. We especially mention the duality pairing associated to the skewed multiplicative LCA group $\mathcal{X} = ((-c; +\infty)^d, \odot)$ $(x\_k+c)(z\_k+c)$ - c, Hence $h_\omega(x) = \sum_{k=1}^{d} \omega_k \log(x_k + c)$. This group together with the operation $\odot$ has been proposed by [11] to handle histograms features especially useful in image recognition applications.

### 2.4.2  *The Fourier transform*

For a function with values in a separable Hilbert space $f \in L^1(\mathcal{X}, dx; \mathcal{Y})$, where $dx$ is the Haar measure on $\mathcal{X}$, we denote $\mathcal{F}[f]$ its Fourier transform (FT) which is defined by

$$\forall \omega \in \hat{\mathcal{X}}, \quad \mathcal{F}[f](\omega) = \int_{\hat{\mathcal{X}}} \overline{(x, \omega)} f(x) dx.$$

For a measure defined on $\mathcal{X}$, there exists a unique suitably normalized measure $d\omega$ on $\hat{\mathcal{X}}$ such that $\forall f \in L^1(\mathcal{X}, dx; \mathcal{Y})$ and if $\mathcal{F}[f] \in L^1(\hat{\mathcal{X}}, d\omega, \mathcal{Y})$ we have

$$\forall x \in \mathcal{X}, \quad f(x) = \int_{\hat{\mathcal{X}}} \mathcal{F}[f](\omega)(x, \omega) d\omega. \tag{1}$$

Moreover if $d\omega$ is normalized, $\mathcal{F}$ extends to a unitary operator from $L^2(\mathcal{X}, dx, \mathcal{Y})$ onto $L^2(\hat{\mathcal{X}}, d\omega, \mathcal{Y})$ Then the inverse Fourier transform (IFT) of a function $g \in L^1(\hat{\mathcal{X}}, d\omega, \mathcal{Y})$ (where $d\omega$ is a Haar measure on $\hat{\mathcal{X}}$ suitably normalize w. r. t. the Haar measure $dx$) is noted $\mathcal{F}^{-1}[g]$ defined by

$$\forall x \in \mathcal{X}, \quad \mathcal{F}^{-1}[g](x) = \int_{\hat{\mathcal{X}}} (x, \omega) g(\omega) d\omega,$$

Equation (0) gives some examples of real Abelian groups with their associated dual and pairing. The interested reader can refer to Folland [9] for a more detailed construction of LCA, Pontryagin duality and Fourier transforms on LCA. For the familiar case of a scalar-valued function $f$ on the LCA group $(\mathbb{R}^d, +)$, we have:

$$\forall \omega \in \hat{\mathcal{X}}, \quad \mathcal{F}[f](\omega) = \int_{\mathbb{R}^d} e^{-i\langle \omega, x-z \rangle} f(x) dx, \tag{2}$$

the Haar measure being here the Lebesgue measure.

## 2.5 ON OPERATOR-VALUED KERNELS

We now introduce the theory of vector-valued Reproducing Kernel Hilbert Space (vv-RKHS) that provides a flexible framework to study and learn vector-valued functions.

### 2.5.1 *Definitions and properties*

An operator-valued kernel is defined here as a $\mathcal{Y}$-reproducing kernel Carmeli et al. [6].

**Definition 2.** *Given $\mathcal{X}$ a Polish space and $\mathcal{Y}$ a Hilbert Space, a map $K : \mathcal{X} \times \mathcal{X} \to \mathcal{L}(\mathcal{Y})$ is called a $\mathcal{Y}$-reproducing kernel if*

$$\sum_{i,j=1}^{N} \langle K(x_i, x_j) y_j, y_i \rangle_{\mathcal{Y}} \geqslant 0,$$

*for all $x_1, \ldots, x_N$ in $\mathcal{X}$, all $y_1, \ldots, y_N$ in $\mathcal{Y}$ and $N \geqslant 1$.*

Given $x \in \mathcal{X}$, $K_x : \mathcal{Y} \to \mathcal{F}(\mathcal{X}; \mathcal{Y})$ denotes the linear operator whose action on a vector $y$ is the function $K_x y \in \mathcal{F}(\mathcal{X}; \mathcal{Y})$ defined for all $z \in \mathcal{X}$ by

$$(K_x y)(z) = K(z, x) y. \tag{3}$$

Additionally, given a $\mathcal{Y}$-reproducing kernel $K$, there is a unique Hilbert space $\mathcal{H}_K \subset \mathcal{F}(\mathcal{X}; \mathcal{Y})$ satisfying $K_x \in \mathcal{L}(\mathcal{Y}; \mathcal{H}_K)$, for all $x \in \mathcal{X}$ and $\forall x \in \mathcal{X}, \forall f \in \mathcal{H}_K, \ f(x) = K_x^* f$, where $K_x^* : \mathcal{H}_K \to \mathcal{Y}$ is the adjoint of $K_x$. The space $\mathcal{H}_K$ is called the *vector-valued Reproducing Kernel Hilbert Space* associated with $K$. The corresponding product and norm are denoted by $\langle ., . \rangle_K$ and $\|.\|_K$, respectively. As a consequence [6] we have

$$K(x, z) = K_x^* K_z \ \ \forall x, z \in \mathcal{X}, \tag{4}$$
$$\mathcal{H}_K = \overline{\text{span}} \{ K_x y \mid \forall x \in \mathcal{X}, \ \forall y \in \mathcal{Y} \}.$$

From eq. (4) we deduce that

$$K(x, z)^* = (K_x^* K_z)^* = K_z^* K_x = K(z, x). \tag{5}$$

Another way to describe functions of $\mathcal{H}_K$ consists in using a suitable feature map.

**Proposition 3** (Feature Operator Carmeli et al. [6]). *Let $\mathcal{H}$ be a Hilbert space and $\phi : \mathcal{X} \to \mathcal{L}(\mathcal{Y}; \mathcal{H})$, with $\phi_x := \phi(x)$. Then the operator $W : \mathcal{H} \to \mathcal{F}(\mathcal{X}; \mathcal{Y})$ defined for all $g \in \mathcal{H}$, and for all $x \in \mathcal{X}$ by $(Wg)(x) = \phi_x^* g$ is a partial isometry from $\mathcal{H}$ onto the vv-RKHS $\mathcal{H}_K$ with reproducing kernel*

$$K(x, z) = \phi_x^* \phi_z, \ \ \forall x, z \in \mathcal{X}.$$

*$W^* W$ is the orthogonal projection onto*

$$Ker \ W^\perp = \overline{\text{span}} \{ \phi_x y \mid \forall x \in \mathcal{X}, \ \forall y \in \mathcal{Y} \}.$$

*Then $\|f\|_K = \inf \{ \|g\|_{\mathcal{H}} \mid \forall g \in \mathcal{H}, \ Wg = f \}$.*

We call $\phi$ a *feature map*, $W$ a *feature operator* and $\mathcal{H}$ a *feature space*. Since $W$ is a partial isometry from $(\ker W)^\perp$ onto $\mathcal{H}_K$, the map $W$ allows us to identify $\mathcal{H}_K$ with the closed subspace $(\ker W)^\perp$ of $\mathcal{H}$. With mild abuse of language, we say that $\mathcal{H}_K$ is embed into $\mathcal{H}$ by mean of the feature operator $W$. If we choose the trivial feature map $\phi_x = K_x$ then the feature operator is the identity.

### 2.5.2 *Shift-Invariant operator-valued kernels*

The main subjects of interest of chapter 3 are shift-invariant Operator-Valued Kernel. When referring to a shift-invariant OVK $K : \mathcal{X} \times \mathcal{X} \to \mathcal{L}(\mathcal{Y})$ we assume that $\mathcal{X}$ is a locally compact second countable topological group with identity $e$.

**Definition 4** (Shift-invariant OVK). *A reproducing Operator-Valued Kernel $K : \mathcal{X} \times \mathcal{X} \to \mathcal{L}(\mathcal{Y})$ is called shift-invariant[1] if for all $x$, $z$, $t \in \mathcal{X}$,*

$$K(x \star t, z \star t) = K(x, z). \tag{6}$$

A shift-invariant kernel can be characterize by a function of one variable $K_e$ called the signature of $K$. Here $e$ denotes the neutral element of the LCA group $\mathcal{X}$ endowed with the operator $\star$.

To study shift-invariant kernels on LCA groups we introduce the left regular representation of $\mathcal{X}$ acting on $\mathcal{H}_K$. For all $x, z \in \mathcal{X}$ and for all $f \in \mathcal{H}_K$,

$$(\lambda_z f)(x) := f(x \star z^{-1}).$$

A group representation $\lambda_z$ describe the group by making it act on a vector space (here $\mathcal{H}_K$) in a linear manner. In other word, the group representation let us see a group as a linear operator which are well studied mathematical object.

**Proposition 5** (Kernel signature). *Let $K : \mathcal{X} \times \mathcal{X} \to \mathcal{L}(\mathcal{Y})$ be a reproducing kernel. The following conditions are equivalents.*

1. *$K$ is a shift-invariant reproducing kernel.*

2. *There is unique a function $K_e : \mathcal{X} \to \mathcal{L}(\mathcal{Y})$ of completely positive type such that $K(x, z) = K_e(x \star z^{-1})$.*

*If one of the above conditions is satisfied, then the representation $\lambda$ leaves invariant $\mathcal{H}_K$, its action on $\mathcal{H}_K$ is unitary and*

$$K(x, z) = K_e^* \lambda_{x^{-1} \star z} K_e \qquad \forall (x, z) \in \mathcal{X}^2. \tag{7a}$$
$$\|K(x, x)\| = \|K_e(e)\| \qquad \forall x \in \mathcal{X} \tag{7b}$$

*Proof.* Assume eq. (7) holds true. Given $x, z \in \mathcal{X}$, eq. (3) and eq. (6) yields

$$K_e(x \star z^{-1}) = K(x \star z^{-1}, e) = K(x, z).$$

Since $K$ is a reproducing kernel, $K_e$ is of completely positive type, so that proposition 5 item 2 holds true. Besides if proposition 5 item 2 holds true obviously the definition of a reproducing kernel (definition 2) is fulfilled so that holds true. Finally let $K_e^1(x \star z^{-1}) = K(x \star z^{-1}, e)$ and $K_e^2(x \star z^{-1}) = K(x \star z^{-1}, e)$. Obviously $K_e^1(x \star z^{-1}) = K_e^2(x \star z^{-1})$ by transitivity of the relation equal.

Suppose that $K$ is a shift-invariant reproducing kernel. Given $t \in \mathcal{X}$ and $y \in \mathcal{Y}$, for all $x, z \in \mathcal{X}$,

$$
\begin{aligned}
(\lambda_x K_t y)(z) = (K_t y)(x^{-1} \star z) &= K(x^{-1} \star z, t) \\
&= K(z, x \star t) \\
&= (K_{x \star t} y)z,
\end{aligned}
$$

that is $\lambda_x K_t = K_{x \star t}$. Besides $\qquad\square$

**Lemma 6** (Shift-invariant $\mathcal{Y}$-Mercer kernels). *Let $K_e : \mathcal{X} \to \mathcal{Y}$ be a function of completely positive type and let $K$ be the corresponding translation invariant reproducing kernel. The following conditions are equivalent.*

1. *The map $K$ is a $\mathcal{Y}$-Mercer kernel.*

2. *For all $y \in \mathcal{Y}$, $K_e(\cdot)y \in \mathcal{C}(\mathcal{X}; \mathcal{Y})$.*

3. *The representation $\lambda$ is continuous on $\mathcal{H}_K$.*

*Proof.* $\qquad\square$

### 2.5.3  *Examples of operator-valued kernels*

Operator-valued kernels have been first introduced in Machine Learning to solve multi-task regression problems. Multi-task regression is encountered in many fields such as structured classification when classes belong to a hierarchy for instance. Instead of solving independently $p$ single output regression task, one would like to take advantage of the relationships between output variables when learning and making a decision.

**Definition 7** (Decomposable kernel). *Let $A$ be a positive semi-definite operator of $\mathcal{L}(\mathcal{Y})$. $K$ is said to be a $\mathcal{Y}$-Mercer decomposable kernel[2] if for all $(x, z) \in \mathcal{X}^2$,*

$$
K(x, z) = k(x, z)A,
$$

*where $k$ is a scalar Mercer kernel.*

[2] *Some authors also refer to as separable kernels.*

When $\mathcal{Y} = \mathbb{R}^p$, the matrix $A$ is interpreted as encoding the relationships between the outputs coordinates. If a graph coding for the proximity between tasks is known, then it is shown in Álvarez,

Rosasco, and Lawrence [1], Baldassarre et al. [2], and Evgeniou, Micchelli, and Pontil [8] that $A$ can be chosen equal to the pseudo inverse $L^\dagger$ of the graph Laplacian such that the norm in $\mathcal{H}_K$ is a graph-regularizing penalty for the outputs (tasks). When no prior knowledge is available, $A$ can be set to the empirical covariance of the output training data or learned with one of the algorithms proposed in the literature [7, 13, 16]. Another interesting property of the decomposable kernel is its universality (a kernel which may approximate an arbitrary continuous target function uniformly on any compact subset of the input space). A reproducing kernel $K$ is said *universal* if the associated vv-RKHS $\mathcal{H}_K$ is dense in the space $\mathcal{C}(\mathcal{X}, \mathcal{Y})$. The conditions for a kernel to be universal have been discussed in Caponnetto et al. [4] and Carmeli et al. [6]. In particular they show that a decomposable kernel is universal provided that the scalar kernel $k$ is universal and the operator $A$ is injective.

Curl-free and divergence-free kernels provide an interesting application of operator-valued kernels [3, 14, 15] to *vector field* learning, for which input and output spaces have the same dimensions ($d = p$). Applications cover shape deformation analysis [15] and magnetic fields approximations [17]. These kernels discussed in [10] allow encoding input-dependent similarities between vector-fields.

**Definition 8** (Curl-free and Div-free kernel). *Assume* $\mathcal{X} = (\mathbb{R}^d, +)$ *and* $\mathcal{Y} = \mathbb{R}^p$ *with* $d = p$. *The divergence-free kernel is defined as*

$$K^{div}(x, z) = K_0^{div}(\delta) = (\nabla \nabla^T - \Delta I) k_0(\delta)$$

*and the curl-free kernel as*

$$K^{curl}(x, z) = K_0^{curl}(\delta) = -\nabla \nabla^T k_0(\delta),$$

*where* $\nabla \nabla^T$ *is the Hessian operator and* $\Delta$ *is the Laplacian operator.*

Although taken separately these kernels are not universal, a convex combination of the curl-free and divergence-free kernels allows to learn any vector field that satisfies the Helmholtz decomposition theorem [3, 14].

❧

Part II

You can put some informational part preamble text here. Illo principalmente su nos. Non message *occidental* angloromanic da. Debitas effortio simplificate sia se, auxiliar summarios da que, se avantiate publicationes via. Pan in terra summarios, capital interlingua se que. Al via multo esser specimen, campo responder que da. Le usate medical addresses pro, europa origine sanctificate nos se.

# OPERATOR-VALUED RANDOM FOURIER FEATURES

### 3.1  MOTIVATIONS

Random Fourier Features have been proved useful to implement efficiently kernel methods in the scalar case, allowing to learn a linear model based on an approximated feature map. In this work, we are interested to construct approximated operator-valued feature maps to learn vector-valued functions. With an explicit (approximated) feature map, one converts the problem of learning a function $f$ in the vector-valued Reproducing Kernel Hilbert Space $\mathcal{H}_K$ into the learning of a linear model $\tilde{f}$ defined by:

$$\tilde{f}(x) = \tilde{\phi}_{1:D}(x)^*\theta,$$

where $\phi : \mathcal{X} \to \mathcal{L}(\mathcal{H}, \mathcal{Y})$ and $\theta \in \mathcal{H}$. The methodology we propose works for operator-valued kernels defined on any Locally Compact Abelian (LCA) group, noted $(\mathcal{X}, \star)$, for some operation noted $\star$. This allows us to use the general context of Pontryagin duality for Fourier transform of functions on LCA groups. Building upon a generalization of Bochner's theorem for operator-valued measures, an operator-valued kernel is seen as the *Fourier transform* of an operator-valued positive measure. From that result, we extend the principle of Random Fourier Feature for scalar-valued kernels and derive a general methodology to build Operator Random Fourier Feature when operator-valued kernels are shift-invariant according to the chosen group operation.

### 3.2  CONSTRUCTION

We present a construction of Operator-valued Random Fourier Feature (ORFF) such that $f : x \mapsto \tilde{\phi}_{1:D}(x)^*\theta$ is a continuous function that maps an arbitrary LCA group $\mathcal{X}$ as input space to an arbitrary output Hilbert space $\mathcal{Y}$. First we define a functional *Fourier feature map*, and then propose a Monte-Carlo sampling from this feature map to construct an approximation of a shift-invariant $\mathcal{Y}$-Mercer kernel. Then, we prove the convergence of the kernel approximation $\tilde{K}(x, z) = \tilde{\phi}_{1:D}(x)^*\tilde{\phi}_{1:D}(z)$ with high probability on *compact* subsets of the LCA $\mathcal{X}$, when $\mathcal{Y}$ is *finite dimensional*. Eventually we conclude with some numerical experiments.

### 3.2.1  *Theoretical study*

The following proposition of Carmeli et al. [6] and Zhang, Xu, and Zhang [18] extends Bochner's theorem to any shift-invariant $\mathcal{Y}$-Mercer kernel.

**Proposition 9** (Operator-valued Bochner's theorem [18])**.** *If a continuous function K from $\mathcal{X} \times \mathcal{X}$ to $\mathcal{Y}$ is a shift-invariant $\mathcal{Y}$-Mercer kernel*

*on $\mathcal{X}$, then there exists a unique positive operator-valued measure* $M : \mathcal{B}(\mathcal{X}) \to \mathcal{L}_+(\mathcal{Y})$ *such that for all* $x, z \in \mathcal{X}$,

$$K(x, z) = \int_{\hat{\mathcal{X}}} \overline{(x \star z^{-1}, \omega)} dM(\omega), \tag{8}$$

*where* $M$ *belongs to the set of all the* $\mathcal{L}_+(\mathcal{Y})$*-valued measures of bounded variation on the* $\sigma$*-algebra of Borel subsets of* $\hat{\mathcal{X}}$*. Conversely, from any positive operator-valued measure* $M$*, a shift-invariant kernel* $K$ *can be defined by proposition* 9.

Although this theorem is central to the spectral decomposition of shift-invariant $\mathcal{Y}$-Mercer OVK, the following results proved by Carmeli et al. [6] provides insights about this decomposition that are more relevant in practice. It first gives the necessary conditions to build shift-invariant $\mathcal{Y}$-Mercer kernel with a pair $(A, \mu)$ where $A$ is an operator-valued function on $\hat{\mathcal{X}}$ and $\mu$ is a real-valued positive measure on $\hat{\mathcal{X}}$. Note that obviously such a pair is not unique ad the choice of this paper may have an impact on theoretical properties as well as practical computations. Secondly it also states that any OVK have such a spectral decomposition when $\mathcal{Y}$ is finite dimensional or $\mathcal{X}$.

**Proposition 10** (Carmeli et al. [6]). *Let* $\mu$ *be a positive measure on* $\mathcal{B}(\hat{\mathcal{X}})$ *and* $A : \hat{\mathcal{X}} \to \mathcal{L}(\mathcal{Y})$ *such that* $\langle A(.)y, y' \rangle \in L^1(\mathcal{X}, d\mu)$ *for all* $y, y' \in \mathcal{Y}$ *and* $A(\omega) \succcurlyeq 0$ *for* $\mu$*-almost all* $\omega$*. Then, for all* $\delta \in \mathcal{X}$ *and for all* $y, y' \in \mathcal{Y}$,

$$K_e(\delta) = \int_{\hat{\mathcal{X}}} \overline{(\delta, \omega)} A(\omega) d\mu(\omega) \tag{9}$$

*is the kernel signature of a shift-invariant* $\mathcal{Y}$*-Mercer kernel* $K$ *such that* $K(x, z) = K_e(x \star z^{-1})$*. The* vv-RKHS $\mathcal{H}_K$ *is embed in* $L^2(\hat{\mathcal{X}}, d\mu; \mathcal{Y}')$ *by mean of the feature operator*

$$(Wg)(x) = \int_{\hat{\mathcal{X}}} \overline{(x, \omega)} B(\omega) g(\omega) d\mu(\omega), \tag{10}$$

*Where* $B(\omega)B(\omega)^* = A(\omega)$ *and both integral converges in the weak sense. If* $\mathcal{Y}$ *is finite dimensional or* $\mathcal{X}$ *is compact, any shift-invariant kernel is of the above form for some pair* $(A, d\mu)$.

When $p = 1$ one can always assume $A$ is reduced to the scalar 1, $\mu$ is still a bounded positive measure and we retrieve the Bochner theorem applied to the scalar case (**??**).

Proposition 10 shows that a given pair $(A, d\mu)$ characterize an OVK. Namely given a measure $d\mu$ and a function $A$ such that $\langle A(.)y, y' \rangle \in L^1(\mathcal{X}, d\mu)$ for all $y, y' \in \mathcal{Y}$ and $A(\omega) \succcurlyeq 0$ for $\mu$-almost all $\omega$, it gives rise to an OVK. Since $(A, d\mu)$ determine a unique kernel we can write $\mathcal{H}_{(A, d\mu)} \Longrightarrow \mathcal{H}_K$ where $K$ is defined as in eq. (9). However the converse is to true: Given a $\mathcal{Y}$-Mercer shift invariant Operator-Valued Kernel, there exist infinitely many pairs $(A, d\mu)$ that characterize an OVK.

The main difference between proposition 9 and proposition 10 is that the first one characterize an ovk by a unique Positive Operator-Valued Measure (povm), while the second one shows that the povm that uniquely characterize a $\mathcal{Y}$-Mercer ovk has an operator-valued density with respect to a *scalar* measure $\mu$; and that this operator-valued density is not unique.

Finally proposition 10 does not provide any *constructive* way to obtain the pair $(A, d\mu)$ that characterize an ovk. The following section 3.2.2 is based on an other proposition of Carmeli, De Vito, and Toigo and show that if the kernel signature $K_e(\delta)$ of an ovk is in $L^1$ then it is possible to construct *explicitly* a pair $(C, d\omega)$ from it. We show that we can always extract a scalar-valued *probability* density function from $C$ such that we obtain a pair $(A, d\mu)$ where $\mu$ is a probability measure.

### 3.2.2 *Sufficient conditions of existence*

While proposition 10 gives some insights on how to build an approximation of a $\mathcal{Y}$-Mercer kernel, we need a theorem that provides an explicit construction of the pair $A(\omega), \mu(\omega)$ from the kernel signature. Proposition 14 in Carmeli et al. [6] gives the solution, and also provide a sufficient condition for proposition 10 to apply.

**Proposition 11** (Carmeli et al. [6]). *Let* K *be a shift-invariant $\mathcal{Y}$-Mercer kernel. Suppose that $\forall z \in \mathcal{X}$ and $\forall y, y' \in \mathcal{Y}$, $\langle K_e(.)y, y' \rangle \in L^1(\mathcal{X}, dx)$ where $dx$ denotes the Haar measure on $(\mathcal{X}, \star)$. Define C such that for all $\omega \in \hat{\mathcal{X}}$ and for all $y$, $y'$ in $\mathcal{Y}$,*

$$\langle y, C(\omega)y' \rangle = \int_{\mathcal{X}} (\delta, \omega)\langle y, K_e(\delta)y' \rangle d\delta = \mathcal{F}^{-1}\left[\langle y, K_e(\cdot)y' \rangle\right](\omega) \text{ (11)}$$

*Then*

1. *$C(\omega)$ is a bounded non-negative operator for all $\omega \in \hat{\mathcal{X}}$,*

2. *$\langle y, C(.)y' \rangle \in L^1(\hat{\mathcal{X}}, d\omega)$ for all $y, y' \in \mathcal{X}$,*

3. *for all $\delta \in \mathcal{X}$ and for all $y$, $y'$ in $\mathcal{Y}$,*

$$\langle y, K_e(\delta)y' \rangle = \int_{\hat{\mathcal{X}}} \overline{(\delta, \omega)}\langle y, C(\omega)y' \rangle d\omega.$$

The following proposition allows to build a spectral decomposition of a shift-invariant $\mathcal{Y}$-Mercer kernel on a lca group $(\mathcal{X}, \star)$ with respect to a scalar probability measure, by extracting a scalar probability density from $C$.

**Proposition 12** (Sufficient condition for shift-invariant $\mathcal{Y}$-Mercer kernel spectral decomposition). *Let* $K_e$ *be the signature of a shift-invariant $\mathcal{Y}$-Mercer kernel on $(\mathcal{X}, \star)$.*

*If for all* $y, y' \in \mathcal{Y}$, $\langle K_e(.)y, y' \rangle \in L^1(\mathcal{X}, dx)$, *then there exists a positive measure* $\mu$ *with density* $p_\mu$ *on* $\mathcal{B}(\hat{\mathcal{X}})$ *and* $A : \hat{\mathcal{X}} \to \mathcal{L}_+(\mathcal{Y})$ *an operator-valued function such that for all* $y, y' \in \mathcal{Y}$ $\langle A(.)y, y' \rangle \in L^1(\mathcal{X}, d\mu)$ *and for all* $y, y'$ *in* $\mathcal{Y}$,

$$\langle y, K_e(\delta)y' \rangle = \int_{\hat{\mathcal{X}}} \overline{(\delta, \omega)} \langle y, A(\omega)y' \rangle p_\mu(\omega) d\omega.$$

*where* $\langle y, A(\omega)y' \rangle p_\mu(\omega) = \mathcal{F}^{-1}\left[\langle y, K_e(\cdot)y' \rangle\right](\omega)$.

*Proof.* From eq. (11) we can write the following equality concerning the OVK signature $K_e$.

$$\int_{\hat{\mathcal{X}}} \overline{(\delta, \omega)} \langle y, C(\omega)y' \rangle d\omega. = \int_{\hat{\mathcal{X}}} \overline{(\delta, \omega)} \left\langle y, \frac{C(\omega)}{p(\omega)}y' \right\rangle p(\omega) d\omega$$

Where $p$ is a function mapping the measured space $(\hat{\mathcal{X}}, \mathcal{B}(\hat{\mathcal{X}}), d\omega)$ to $\mathbb{R}$. It is always possible to choose $p(\omega)$ such that $\int_{\hat{\mathcal{X}}} p(\omega) d\omega = 1$ in this case $p(\omega)$ is the density of a probability measure $\mu$. In this case we note $p(\omega) = p_\mu(\omega)$. Conclude by taking $A(\omega) = C(\omega)/p_\mu(\omega)$ and $d\mu(\omega) = p_\mu(\omega) d\omega$. □

In the case where $\mathcal{Y} = \mathbb{R}^p$, we rewrite proposition 12 coefficient-wise by choosing an orthonormal basis $(e_1, \ldots, e_p)$ of $\mathcal{Y}$, such that for all $i, j \in \{1, \ldots, p\}$,

$$\langle e_i, C(\omega)e_j \rangle = C(\omega)_{ij} = A(\omega)_{ij} p_\mu(\omega) = \mathcal{F}^{-1}\left[K_e(\delta)_{ij}\right]. \tag{12}$$

It follows that for all $i, j \in \{1, \ldots, p\}$,

$$K_e(x \star z^{-1})_{ij} = \mathcal{F}\left[A(\cdot)_{ij}\right] \tag{13}$$

**Remark 13.** *Note that although the inverse Fourier transform of* $K_e$ *yields a unique operator-valued function* $C(\cdot)$, *the decomposition of* $C(\omega)$ *into* $A(\omega)p_\mu(\omega)$ *is again not unique. The choice of the decomposition may be justified by the computational cost or by the nature of the constants involved in the uniform convergence of the estimator.*

### 3.2.3 Regularization property

We have shown so far that it is always possible to construct a feature map that allows to approximate a shift-invariant $\mathcal{Y}$-Mercer kernel. However we could also propose a construction of such map by studying the regularization induced with respect to the Fourier transform of a target function $f \in \mathcal{H}_K$. In other words, what is the norm in $L^2(\hat{\mathcal{X}}, d\omega, \mathcal{Y}')$ induced by $\|\cdot\|_K$?

**Proposition 14.** *Let* $K$ *be a shift-invariant* $\mathcal{Y}$-*Mercer Kernel such that for all* $y, y'$ *in* $\mathcal{Y}$, $\langle y, K_e(\cdot)y' \rangle \in L^1(\mathcal{X}, dx)$ *and*

*Let $\langle y, A(\omega)y' \rangle p_\mu(\omega) := \mathcal{F}^{-1}\left[\langle y, K_e(\cdot)y'\rangle\right](\omega)$ and let $f \in \mathcal{H}_K$. Then*

$$\|f\|_K^2 = \int_{\hat{\mathcal{X}}} \frac{\langle \mathcal{F}[f](\omega), A(\omega)^\dagger \mathcal{F}[f](\omega)\rangle_{\mathcal{Y}}}{p_\mu(\omega)} d\omega. \tag{14}$$

*Proof.* We first show how the Fourier transform relates to the feature operator. Since $\mathcal{H}_K$ is embed into $\mathcal{H} = L^2(\hat{\mathcal{X}}, \mu, \mathcal{Y})$ by mean of the feature operator $W$, we have:

$$\mathcal{F}\left[\mathcal{F}^{-1}[f]\right](x) = \int_{\hat{\mathcal{X}}} \overline{(x,\omega)} \mathcal{F}^{-1}[f](\omega) d\omega = f(x)$$

$$(Wg)(x) = \int_{\hat{\mathcal{X}}} \overline{(x,\omega)} p_\mu(\omega) B(\omega) g(\omega) d\omega = f(x).$$

By injectivity of the Fourier transform, $\mathcal{F}^{-1}[f](\omega) = p_\mu(\omega) B(\omega) g(\omega)$ $\mu$-almost everywhere. From proposition 3 we have

$$\|f\|_K^2 = \inf\left\{\|g\|_{\mathcal{H}}^2 \mid \forall g \in \mathcal{H}, \ Wg = f\right\}$$

$$= \inf\left\{\int_{\hat{\mathcal{X}}} \|g\|_{\mathcal{Y}}^2 d\mu \mid \forall g \in \mathcal{H}, \ \mathcal{F}^{-1}[f] = p_\mu(\cdot)B(\cdot)g(\cdot)\right\}.$$

The pseudo inverse of the operator $B(\omega)$ (noted $B(\omega)^\dagger$) is the unique solution of the system $\mathcal{F}^{-1}[f](\omega) = p_\mu(\omega)B(\omega)g(\omega)$ w.r.t. $g(\omega)$ with minimal norm. Eventually,

$$\|f\|_K^2 = \int_{\hat{\mathcal{X}}} \frac{\left\|B(\omega)^\dagger \mathcal{F}^{-1}[f](\omega)\right\|_{\mathcal{Y}}^2}{p_\mu(\omega)^2} d\mu(\omega)$$

$$= \int_{\hat{\mathcal{X}}} \frac{\left\|B(\omega)^\dagger \mathcal{F}[f](\omega)\right\|_{\mathcal{Y}}^2}{p_\mu(\omega)^2} d\mu(\omega) \tag{15}$$

Conclude the proof by taking $d\mu(\omega) = p_\mu(\omega)d\omega$ and rewriting the integral as an expectation. $\qquad\square$

Note that if $K(x,z) = k(x,z)$ is a scalar kernel then for all $\omega$ in $\hat{\mathcal{X}}$, $A(\omega) = 1$. Therefore we recover a well known results for kernels that is for any $f \in \mathcal{H}_k$ we have $\|f\|_k = \int_{\hat{\mathcal{X}}} \mathcal{F}[k_e](\omega)^{-1}\mathcal{F}[f](\omega)^2 d\omega$. We also note that the regularization property in $\mathcal{H}_K$ does not depends (as expected) on the decomposition of $A(\omega)$ into $B(\omega)B(\omega)^*$. Therefore the decomposition should be chosen such that it optimizes the computation cost. For instance if $A(\omega) \in \mathcal{L}(\mathbb{R}^p)$ has rank $r$, one could find an operator $B(\omega) \in \mathcal{L}(\mathbb{R}^p, \mathbb{R}^r)$ such that $A(\omega) = B(\omega)B(\omega)^*$.

### 3.2.4    *Functional Fourier feature map*

Let us introduce a functional feature map, we call here *Fourier Feature map*, defined by the following proposition as a direct consequence of proposition 10.

**Proposition 15** (Fourier feature map). *If there exist an operator-valued function* $B : \hat{\mathcal{X}} \to \mathcal{L}(\mathcal{Y}, \mathcal{Y}')$ *such that for all* $y, y' \in \mathcal{Y}$, $\langle y, B(\omega)B(\omega)^* y' \rangle = \langle y, A(\omega)y' \rangle$ *$\mu$-almost everywhere and* $\langle y, A(\omega)y' \rangle \in L^1(\hat{\mathcal{X}}, d\mu)$ *then the operator* $\phi_x$ *defined for all* $y$ *in* $\mathcal{Y}$ *by*

$$(\phi_x y)(\omega) = (x, \omega)B(\omega)^* y, \tag{16}$$

*is a feature map[3] of some shift-invariant kernel* $K$.

> [3] *I. e. it satisfies for all* $x, z \in \mathcal{X}$, $\phi_x^* \phi_z = K(x, z)$ *where* $K$ *is a* $\mathcal{Y}$-*Mercer* OVK.

*Proof.* For all $y, y' \in \mathcal{Y}$ and $x, z \in \mathcal{X}$,

$$
\begin{aligned}
\langle y, \phi_x^* \phi_z y' \rangle &= \langle \phi_x y, \phi_z y' \rangle_{L^2(\hat{\mathcal{X}}, \mu, \mathcal{Y}')} \\
&= \int_{\hat{\mathcal{X}}} \overline{(x, \omega)} \langle y, B(\omega)(z, \omega)B(\omega)^* y' \rangle d\mu(\omega) \\
&= \int_{\hat{\mathcal{X}}} \overline{(x \star z^{-1}, \omega)} \langle y B(\omega)B(\omega)^* y' \rangle d\mu(\omega) \\
&= \int_{\hat{\mathcal{X}}} \overline{(x \star z^{-1}, \omega)} \langle y, A(\omega)y' \rangle d\mu(\omega),
\end{aligned}
$$

which defines a $\mathcal{Y}$-Mercer according to proposition 10 of Carmeli et al. [6]. □

With this notation notice that $\phi : \mathcal{X} \to \mathcal{L}(\mathcal{Y}; L^2(\hat{\mathcal{X}}, \mu; \mathcal{Y}'))$ such that $\phi_x \in \mathcal{L}(\mathcal{Y}; L^2(\hat{\mathcal{X}}, \mu; \mathcal{Y}'))$ where $\phi_x := \phi(x)$.

### 3.2.5 *Building Operator-valued Random Fourier Features*

Throughout the document, without loss of generality, we assume that $\int_{\mathcal{X}} d\mu(\omega) = 1$ and thus $d\mu$ is a probability measure with density $p_\mu$. As shown in proposition 15 it is always possible to find a pair $(A(\omega), d\mu)$ such that $d\mu$ is a probability measure and $\mathbf{E}_\mu \overline{(\delta, \omega)} A(\omega)$.

Given a $\mathcal{Y}$-Mercer shift-invariant kernel $K$ on $\mathcal{X}$, an approximation of $K$ can be obtained using a decomposition $(A, \mu)$ and a plug-in Monte-Carlo estimator instead of the expectation. However, for efficient computations, as motivated in the introduction, we are interested in finding an approximated feature map more than a kernel approximation. Indeed, an approximated feature map will allow to build linear models in regression tasks. The following proposition provides the general form of an Operator-valued Random Fourier Feature.

**Proposition 16.** *If one can find* $B : \hat{\mathcal{X}} \to \mathcal{L}(\mathcal{Y}', \mathcal{Y})$ *and a probability measure* $\mu$ *on* $\mathcal{B}(\hat{\mathcal{X}})$, *such that for all* $y \in \mathcal{Y}$ *and all* $y' \in \mathcal{Y}'$, $\langle y, B(\cdot)y' \rangle \in L^2(\hat{\mathcal{X}}, d\mu)$, *then the operator-valued function*

$$\tilde{\phi}_{1:D}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^{D} (x, \omega_j)B(\omega_j)^*, \qquad \omega_j \sim \mu \tag{17}$$

*is an approximated feature map of an Operator-Valued Kernel[4] .*

> [4] *I. e. it satisfies* $\tilde{\phi}_{1:D}(x)^* \tilde{\phi}_{1:D}(z) \xrightarrow[D \to \infty]{a.\,s.} K(x, z)$ *where* $K$ *is a* $\mathcal{Y}$-*Mercer* OVK.

*Proof.* Let $\omega_1, \ldots, \omega_D$ be $D$ i. i. d. random vectors following the law $\mu$. For all $x, z \in \mathcal{X}$ and all $y, y' \in \mathcal{Y}$,

$$
\begin{aligned}
& \left\langle \tilde{\phi}_{1:D}(x)y, \tilde{\phi}_{1:D}(z)y' \right\rangle \\
& = \left\langle y, \left( \frac{1}{\sqrt{D}} \bigoplus_{j=1}^{D} (z, \omega_j) B(\omega_j)^* \right)^* \left( \frac{1}{\sqrt{D}} \bigoplus_{j=1}^{D} (x, \omega_j) B(\omega_j)^* \right) y' \right\rangle \\
& = \frac{1}{D} \sum_{j=1}^{D} \overline{(x \star z^{-1}, \omega_j)} A(\omega_j),
\end{aligned}
$$

where $A(\omega) = B(\omega)B(\omega)^*$. By assumption $\langle y, A(\cdot)y' \rangle \in L^1(\hat{\mathcal{X}}, \mu)$ and $\omega_j$ are i. i. d. . Hence from the strong law of large numbers and proposition 10,

$$
\frac{1}{D} \sum_{j=1}^{D} \overline{(x \star z^{-1}, \omega_j)} A(\omega_j) \xrightarrow[D \to \infty]{\text{a. s.}} \mathbf{E}_\mu [\overline{(x \star z^{-1}, \omega_j)} A(\omega)] = K_e(x \star z^{-1})
$$

in the weak operator topology. $\qquad \square$

The approximate feature map proposed in proposition 16 has direct link with the functional feature map defined in proposition 15 since we have for all $y \in \mathcal{Y}$

$$
\begin{aligned}
\tilde{\phi}_{1:D}(x)y & = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^{D} (\phi_x y)(\omega_j) \\
& = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^{D} (x, \omega_j) B(\omega_j)^* y, \qquad \omega_j \sim \mu.
\end{aligned}
\tag{18}
$$

Therefore $\tilde{\phi}_{1:D}(x)$ can be seen as an "operator-valued vector" corresponding the "stacking" of $D$ i. i. d. operator-valued realization of $\phi_x$, the functional feature map. In the same way we can define an approximate feature operator $\tilde{W}_D$.

**Definition 17** (Random Fourier feature operator)**.** *Let* $\theta = \bigoplus_{j=1}^{D} \theta_j \in (\mathcal{Y}')^D$, *where* $\theta \in \mathcal{Y}'$. *We call random Fourier feature operator the linear application* $W : (\mathcal{Y}')^D \to \mathcal{F}(\mathcal{X}; \mathcal{Y})$ *defined as follow.*

$$
(\tilde{W}_D \theta)(x) := \tilde{\phi}_{1:D}(x)^* \theta = \frac{1}{D} \sum_{j=1}^{D} \overline{(x, \omega_j)} B(\omega_j) \theta_j, \qquad \omega_j \sim \mu. \tag{19}
$$

The approximate feature operator is useful to show the relations between the approximate feature map with the functional feature map defined in proposition 15.

**Proposition 18.** *Let* $g \in \mathcal{H} = L^2(\hat{\mathcal{X}}, d\mu; \mathcal{Y})$ *and let* $\theta := \bigoplus_{j=1}^{D} g(\omega_j)$ *where* $\omega_j \sim \mu$. *Then for all* $g \in \mathcal{H}$,

1. $\tilde{\phi}_{1:D}(x)^*\theta \xrightarrow[D\to\infty]{a.\,s.} \phi_x^* g,$

2. $\|\theta\|_{\mathcal{Y}}^2 \xrightarrow[D\to\infty]{a.\,s.} \|g\|_{L^2(\hat{\mathcal{X}}, d\mu; \mathcal{Y}')}^2.$

*Proof.* Proof of proposition 18 item 1: since $\omega_1, \ldots \omega_D$ are i. i. d. random vectors, for all $y \in \mathcal{Y}$ and for all $y' \in \mathcal{Y}'$, $\langle y, B(\cdot)y'\rangle \in L^2(\hat{\mathcal{X}}, d\mu)$ and $g \in L^2(\hat{\mathcal{X}}, d\mu; \mathcal{Y})$, from the strong law of large numbers

$$\tilde{\phi}_{1:D}(x)^*\theta = \frac{1}{D}\sum_{j=1}^{D} \overline{(x, \omega_j)} B(\omega_j) g(\omega_j), \qquad \omega_j \sim \mu$$

$$\xrightarrow[D\to\infty]{a.\,s.} \int_{\hat{\mathcal{X}}} \overline{(x, \omega)} B(\omega) g(\omega) d\mu(\omega) = (Wg)(x) := \phi_x^* g.$$

Proof of proposition 18 item 2: again, since $\omega_1, \ldots \omega_D$ are i. i. d. random vectors and $g \in L^2(\hat{\mathcal{X}}, d\mu; \mathcal{Y})$, from the strong law of large numbers

$$\|\theta\|_{\mathcal{Y}}^2 = \sum_{j=1}^{D} \|g(\omega_j)\|_{\mathcal{Y}}^2, \qquad \omega_j \sim \mu$$

$$\xrightarrow[D\to\infty]{a.\,s.} \int_{\hat{\mathcal{X}}} \|g(\omega)\|_{\mathcal{Y}}^2 d\mu(\omega) := \|g\|_{L^2(\hat{\mathcal{X}}, d\mu; \mathcal{Y}')}^2$$

$\square$

Hence the sequence of function $\tilde{f}_D := \tilde{\phi}_{1:D}(\cdot)^*\theta$ converges almost surely to a function $f \in \mathcal{H}_{(A, d\mu)} \implies \mathcal{H}_K$. Therefore, in light of eq. (14), it is possible to define an approximate feature map of an Operator-Valued Kernel from its regularization properties in the vv-RKHS. Otherwise corollary 19 exhibit a construction of an ORFF directly from an OVK.

**Corollary 19.** *If* $K(x, z)$ *is a shift-invariant* $\mathcal{Y}$*-Mercer kernel such that for all* $y, y' \in \mathcal{Y}$, $\langle y, K_e(\delta)y'\rangle \in L^1(\mathcal{X}, dx)$. *Then*

$$\tilde{\phi}_{1:D}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^{D}(x, \omega_j) B(\omega_j)^*, \qquad \omega_j \sim \mu, \tag{20}$$

*where* $\langle y, B(\omega)B(\omega)^*y'\rangle p_\mu(\omega) = \mathcal{F}^{-1}\left[\langle y, K_e(\cdot)y'\rangle\right](\omega)$, *is an approximated feature map of* K.

*Proof.* Find $(A(\omega), d\mu)$ from proposition 12 and apply proposition 16.

$\square$

We write $\tilde{\phi}_{1:D}(x)^*\tilde{\phi}_{1:D}(x) \approx K(x, z)$ when $\tilde{\phi}_{1:D}(x)^*\tilde{\phi}_{1:D}(x) \xrightarrow{a.\,s.} K(x, z)$ in the weak operator topology when $D$ tends to infinity. With mild abuse of notation we say that $\tilde{\phi}_{1:D}(x)$ is an approximate feature map of $\phi_x$ i. e. $\tilde{\phi}_{1:D}(x) \approx \phi_x$, when for all $y \in \mathcal{Y}$, $\langle y, K(x, z)y'\rangle = \langle \phi_x y, \phi_z y'\rangle \approx \langle \tilde{\phi}_{1:D}(x)y, \tilde{\phi}_{1:D}(x)y'\rangle := \tilde{K}(x, z)$ where $\phi_x$ is defined in the sense of proposition 15.

The kernel approximation $\tilde{K}$ can be seen as the sample mean of the product of functional feature map. Indeed $\tilde{K} = 1/D \sum_{j=1}^{D}$.

**Remark 20.** *We find a decomposition such that for all* $j = 1, \ldots, D$, $A(\omega_j) = B(\omega_j)B(\omega_j)^*$ *either by exhibiting an analytic closed-form or using a numerical decomposition.*

Corollary 19 allows us to define algorighm 1 for constructing ORFF from an operator valued kernel.

---

**Algorithm 1:** Construction of ORFF from OVK

**Input** : $K(x, z) = K_e(\delta)$ a $\mathcal{Y}$-shift-invariant Mercer kernel such that $\forall y, y' \in \mathcal{Y}, \langle y, K_e(\delta)y' \rangle \in L^1(\mathbb{R}^d, dx)$ and $D$ the number of features.

**Output:** A random feature $\tilde{\phi}_{1:D}(x)$ such that
$\tilde{\phi}_{1:D}(x)^* \tilde{\phi}_{1:D}(z) \approx K(x, z)$

1 Define the pairing $(x, \omega)$ from the LCA group $(\mathcal{X}, \star)$;
2 Find a decomposition $(B(\omega), p_\mu(\omega))$ such that
  $B(\omega)B(\omega)^* p_\mu(\omega) = \mathcal{F}^{-1}[K_e](\omega)$;
3 Draw $D$ random vectors $\omega_j, j = 1, \ldots, D$ from the probability distribution $\mu$;
4 **return** $\tilde{\phi}_{1:D}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^{D}(x, \omega_j)B(\omega_j)^*$;

---

### 3.2.6  *Examples of Operator Random Fourier Feature maps*

We now give two examples of operator-valued random Fourier feature map when $\mathcal{Y} \subset \mathbb{R}^p$. First we introduce the general form of an approximated feature map for a matrix-valued kernel on the additive group $(\mathbb{R}^d, +)$.

**Example 1** (Matrix-valued kernel on the additive group). *In the following,* $K(x, z) = K_0(x - z)$ *is a* $\mathbb{R}^p$-*Mercer matrix-valued kernel invariant w.r.t. the group operation* +. *An Operator Random Fourier feature map of an* $\mathbb{R}^p$-*Mercer shift-invariant matrix-valued kernel takes the general form:*

$$\tilde{\phi}_{1:D}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^{D} \begin{pmatrix} \cos \langle x, \omega_j \rangle B(\omega_j)^* \\ \sin \langle x, \omega_j \rangle B(\omega_j)^* \end{pmatrix}, \quad \omega_j \sim \mu.$$

*Proof.* The (Pontryagin) dual of $\mathcal{X} = \mathbb{R}^d$ is $\hat{\mathcal{X}} = \mathbb{R}^d$, and the duality pairing is $(x - z, \omega) = \exp(i\langle x - z, \omega \rangle)$. A $\mathbb{R}^p$-operator-valued function has a real operator-valued Fourier transform if and only if $A(\omega)$ is

even with respect to $\omega$. Taking this point into account, the kernel approximation yields:

$$\tilde{K}(x, z) = \tilde{\phi}_{1:D}(x)^* \tilde{\phi}_{1:D}(z)$$

$$= \frac{1}{D} \sum_{j=1}^{D} \cos \langle x, \omega_j \rangle \cos \langle z, \omega_j \rangle A(\omega_j) + \sin \langle x, \omega_j \rangle \sin \langle z, \omega_j \rangle A(\omega_j)$$

$$= \frac{1}{D} \sum_{j=1}^{D} \cos \langle x - z, \omega_j \rangle A(\omega_j)$$

$$= \frac{1}{D} \sum_{j=1}^{D} \exp(-i\langle x - z, \omega_j \rangle) A(\omega_j).$$

which tends to $\mathbf{E}_\mu[\exp(-i\langle x - z, \omega \rangle) A(\omega)] = \mathbf{E}_\mu[\overline{(x - z, \omega)} A(\omega)] = K(x, z)$ when $D$ tends to infinity. $\qquad \square$

The second example extends scalar-valued Random Fourier Features on the skewed multiplicative group described in **??** [11]) to the operator-valued case.

**Example 2** (Matrix-valued kernel on the skewed multiplicative group). *In the following, suppose that $\mathcal{X} = (-c; +\infty)^d$, $\mathcal{Y} = \mathbb{R}^p$ and $K(x, z) = K_{1-c}(x \odot z^{-1})$ is a $\mathbb{R}^p$-Mercer matrix-valued kernel invariant w.r.t. the group operation $\odot$ defined in **??**. The group operation is defined coefficient-wise for all $k \in \{1, \ldots, d\}$ as $x_k \odot z_k := (x_k + c)(z_k + c) - c$. As a consequence $z_k^{-1} = 1/(z_k + c) - c$. The following function $\tilde{\phi}_{1:D}$ is an operator-valued Random Fourier Feature map built following the construction principle:*

$$\tilde{\phi}_{(x)} = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^{D} \begin{pmatrix} \cos \langle \log(x + c), \omega_j \rangle B(\omega_j)^* \\ \sin \langle \log(x + c), \omega_j \rangle B(\omega_j)^* \end{pmatrix}, \quad \omega_j \sim \mu.$$

*Proof.* The dual of $\mathcal{X} = (-c; +\infty)^d$ is $\hat{\mathcal{X}} = \mathbb{R}^d$, and the duality pairing is $(x \odot z^{-1}, \omega) = \exp(i\langle \log(x \odot z^{-1} + c), \omega \rangle)$ (see Li, Ionescu, and Sminchisescu [12]). Following the proof of example 1, we have

$$\tilde{K}(x, z) = \frac{1}{D} \sum_{j=1}^{D} e^{i\langle \log(\frac{x+c}{z+c}), \omega_j \rangle} A(\omega_j).$$

which tends to $\mathbf{E}_\mu[\exp(-i\langle \log(x \odot z^{-1} + c) \rangle) A(\omega)] = \mathbf{E}_\mu[\overline{(x \odot z^{-1}, \omega)} A(\omega)] = K(x, z)$ when $D$ tends to infinity. $\qquad \square$

## 3.3 LEARNING WITH OPERATOR-VALUED RANDOM-FOURIER FEATURES

❦

## 3.4 UNIFORM BOUND ON THE APPROXIMATION

## 3.5 CONSISTENCY AND GENERALIZATION BOUNDS

## 3.6 CONCLUSIONS

# CONCURRENT METHODS

Part III

You can put some informational part preamble text here. Illo principalmente su nos. Non message *occidental* angloromanic da. Debitas effortio simplificate sia se, auxiliar summarios da que, se avantiate publicationes via. Pan in terra summarios, capital interlingua se que. Al via multo esser specimen, campo responder que da. Le usate medical addresses pro, europa origine sanctificate nos se.

CONLUSIONS

Part IV

APPENDIX

PROOFS OF THEOREMS

# BIBLIOGRAPHY

[1]  M. A. Álvarez, L. Rosasco, and N. D. Lawrence. "Kernels for vector-valued functions: a review." In: *Foundations and Trends in Machine Learning* 4.3 (2012), pp. 195–266.

[2]  L. Baldassarre, L. Rosasco, A. Barla, and A. Verri. "Vector Field Learning via Spectral Filtering." In: *ECML/PKDD*. Ed. by J. Balcazar, F. Bonchi, A. Gionis, and M. Sebag. Vol. 6321. LNCS. Springer Berlin / Heidelberg, 2010, pp. 56–71.

[3]  L. Baldassarre, L. Rosasco, A. Barla, and A. Verri. "Multi-output learning via spectral filtering." In: *Machine Learning* 87.3 (2012), pp. 259–301.

[4]  A. Caponnetto, C. A. Micchelli, M., and Y. Ying. "Universal Multi-Task Kernels." In: *Journal of Machine Learning Research* 9 (2008), pp. 1615–1646.

[5]  C. Carmeli, E. De Vito, and A. Toigo. "Vector valued reproducing kernel Hilbert spaces of integrable functions and Mercer theorem." In: *Analysis and Applications* 4.04 (2006), pp. 377–408.

[6]  C. Carmeli, E. De Vito, A. Toigo, and V. Umanità. "Vector valued reproducing kernel Hilbert spaces and universality." In: *Analysis and Applications* 8 (2010), pp. 19–61.

[7]  F. Dinuzzo, C.S. Ong, P. Gehler, and G. Pillonetto. "Learning Output Kernels with Block Coordinate Descent." In: *Proc. of the 28th Int. Conf. on Machine Learning*. 2011.

[8]  T. Evgeniou, C. A. Micchelli, and M. Pontil. "Learning Multiple Tasks with kernel methods." In: *JMLR* 6 (2005), pp. 615–637.

[9]  Gerald B Folland. *A course in abstract harmonic analysis.* CRC press, 1994.

[10]  E. Fuselier. "Refined Error Estimates for Matrix-Valued Radial Basis Functions." PhD thesis. Texas A&M University, 2006.

[11]  F. Li, C. Ionescu, and C. Sminchisescu. "Pattern Recognition: 32nd DAGM Symposium, Darmstadt, Germany, September 22-24, 2010. Proc." In: ed. by M. Goesele, S. Roth, A. Kuijper, B. Schiele, and K. Schindler. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010. Chap. Random Fourier Approximations for Skewed Multiplicative Histogram Kernels, pp. 262–271. ISBN: 978-3-642-15986-2. DOI: 10.1007/978-3-642-15986-2_27. URL: http://dx.doi.org/10.1007/978-3-642-15986-2_27.

[12]   F. Li, C. Ionescu, and C. Sminchisescu. "Pattern Recognition: 32nd DAGM Symposium, Darmstadt, Germany, September 22-24, 2010. Proc." In: ed. by M/ Goesele, S. Roth, A. Kuijper, B. Schiele, and K. Schindler. 2010. Chap. Random Fourier Approximations for Skewed Multiplicative Histogram Kernels.

[13]   N. Lim, F. d'Alché-Buc, C. Auliac, and G. Michailidis. "Operator-valued kernel-based vector autoregressive models for network inference." In: *Machine Learning* 99.3 (2015), pp. 489–513.

[14]   Y. Macedo and R. Castro. *Learning Div-Free and Curl-Free Vector Fields by Matrix-Valued Kernels.* Tech. rep. Preprint A 679/2010 IMPA, 2008.

[15]   M. Micheli and J. Glaunes. *Matrix-valued kernels for shape deformation analysis.* Tech. rep. Arxiv report, 2013.

[16]   V. Sindhwani, H. Q. Minh, and A.C. Lozano. "Scalable Matrix-valued Kernel Learning for High-dimensional Nonlinear Multivariate Regression and Granger Causality." In: *Proc. of UAI'13, Bellevue, WA, USA, August 11-15, 2013.* AUAI Press, Corvallis, Oregon, 2013.

[17]   N. Wahlström, M. Kok, T.B. Schön, and Fredrik Gustafsson. "Modeling magnetic fields using Gaussian processes." In: *in Proc. of the 38th ICASSP.* 2013.

[18]   Haizhang Zhang, Yuesheng Xu, and Qinghui Zhang. "Refinement of Operator-valued Reproducing Kernels." In: *Journal of Machine Learning Research* 13 (2012), pp. 91–136.

## DECLARATION

Put your declaration here.

*15, Rue Plumet, 75015 - Paris, France, Septembre 2016*

Romain Brault