

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

Romain Brault: *Data are not real!* Large-scale learning on structured input-output data with operator-valued kernels, © Septembre 2016

SUPERVISOR:

Professor (Prof.) Florence d'Alché-Buc

LOCATION:

15, Rue Plumet, 75015 - Paris, France

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

Put your acknowledgements here.

Many thanks to everybody who already sent me a postcard!

Regarding the typography and other help, many thanks go to Marco Kuhlmann, Philipp Lehman, Lothar Schlesier, Jim Young, Lorenzo Pantieri and Enrico Gregorio¹, Jörg Sommer, Joachim Köstler, Daniel Gottschlag, Denis Aydin, Paride Legovini, Steffen Prochnow, Nicolas Repp, Heinrich Harms, Roland Winkler, and the whole \LaTeX -community for support, ideas and some great software.

Regarding \LyX : The \LyX port was initially done by *Nicholas Mariette* in March 2009 and continued by *Ivo Pletikosić* in 2011. Thank you very much for your work and the contributions to the original style.

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

CONTENTS

I	INTRODUCTION	1
1	MOTIVATIONS	3
2	BACKGROUND	5
2.1	Notations	6
2.2	About statistical learning	6
2.3	On large-scale learning	6
2.4	Elements of abstract harmonic analysis	6
2.4.1	Locally compact Abelian groups	6
2.4.2	The Haar measure	6
2.4.3	Representation of Groups	8
2.4.4	Characters	9
2.4.5	The Fourier transform	10
2.5	On operator-valued kernels	11
2.5.1	Definitions and properties	11
2.5.2	Shift-Invariant operator-valued kernels	16
2.5.3	Examples of operator-valued kernels . .	18
II	CONTRIBUTIONS	21
3	OPERATOR-VALUED RANDOM FOURIER FEATURES	23
3.1	Motivations	24
3.2	Construction	24
3.2.1	Theoretical study	24
3.2.2	Sufficient conditions of existence	26
3.2.3	Regularization property	29
3.2.4	Functional Fourier feature map	30
3.2.5	Building Operator-valued Random Fourier Features	30
3.2.6	Examples of Operator Random Fourier Feature maps	33
3.3	Learning with operator-valued random-Fourier features	35
3.4	Uniform bound on the approximation	35
3.5	Consistency and generalization bounds	35
3.6	Conclusions	35
4	CONCURRENT METHODS	37
4.1	Background	38
4.2	The Nyström method	38
4.3	Sub-sampling the data	38
4.4	Conclusions	38

III	FINAL WORDS	39
5	CONCLUSIONS	41
IV	APPENDIX	43
A	OPERATOR-VALUED FUNCTIONS AND INTEGRATION	45
B	PROOFS OF THEOREMS	47
	BIBLIOGRAPHY	49

LIST OF FIGURES

LIST OF TABLES

Table 1	Mathematical symbols used throughout the paper and their signification.	7
Table 2	Classification of Fourier transforms in terms of their domain and transform domain. .	10

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.

MOTIVATIONS



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 σ -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 $f : \mathcal{X} \rightarrow \mathcal{Y}$ 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 $W \in \mathcal{L}(\mathcal{H})$, $\text{Ker } W$ denotes the nullspace, $\text{Im } W$ 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} endowed with a binary operation \star is said to be *Locally Compact Abelian* if \mathcal{X} is a topological commutative group w. r. t. \star for which every point has a compact neighborhood and is Hausdorff.

2.4.2 The Haar measure

Definition 2 (The Haar measure). *There is, up to a positive multiplicative constant, a unique countably additive, nontrivial measure dx on the Borel subsets of \mathcal{X} satisfying the following properties¹:*

1. *The measure dx is translation-invariant: $dx(z \star E) = dx(E)$ for every x in \mathcal{X} and Borel set $E \in \mathcal{B}(\mathcal{X})$.*
2. *The measure dx is finite on every compact set: $dx(K) < \infty$ for all compact $K \in \mathcal{B}(\mathcal{X})$.*
3. *The measure dx is outer regular on Borel sets E :*

$$dx(E) = \inf \{ dx(U) \mid E \subseteq U, U \text{ open} \}$$

4. *The measure dx is inner regular on open sets E :*

$$dx(E) = \sup \{ dx(K) \mid K \subseteq E, K \text{ compact} \}.$$

Such a measure on G is called a Haar measure².

¹ For more details and constructive proofs see [1, 8, 11].

² If \mathcal{X} was not supposed to be Abelian, we should have defined a left Haar measure and a right Haar measure. In our case both measure are the same, so we refer to both of them as Haar measure

Table 1: Mathematical symbols used throughout the paper and their signification.

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 (\cdot) .
$\hat{\mathcal{X}}$	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 σ -algebra on \mathcal{X} .
$\mu(\mathcal{X})$	A scalar positive measure of \mathcal{X} .
$p_{\mu}(x)$	The Radon-Nikodym derivative of μ 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} .

It can be shown as a consequence of the above properties that $dx(U) > 0$ for every non-empty open subset U . In particular, if \mathcal{X} is compact then $dx(\mathcal{X})$ is finite and positive, so we can uniquely specify a left Haar measure on \mathcal{X} by adding the normalization condition

$dx(\mathcal{X}) = 1$. Another useful property is that given a continuous function $f \in \mathcal{C}(\mathcal{X})$,

$$\int_{\mathcal{X}} f(x) dx = \int_{\mathcal{X}} f(x^{-1}) dx.$$

We call measured space the space $(\mathcal{X}, \mathcal{B}(\mathcal{X}), dx)$ the space \mathcal{X} endowed with its Borel sigma algebra and some measure dx . If $dx(\mathcal{X}) = 1$ then the space $(\mathcal{X}, \mathcal{B}(\mathcal{X}), dx)$ is called a probability space.

We say that a function $f : \mathcal{X} \rightarrow \mathcal{C}$ is even if for all $x \in \mathcal{X}$, $f(x) = f(x^{-1})$ and odd if $f(x) = -f(x^{-1})$. The definition can be extended to operator-valued functions.

Definition 3 (Even and odd function on a LCA group). *Let \mathcal{X} be a measured LCA group. A function $f : \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ is (weakly) even if for all $x \in \mathcal{X}$ and all $y, y' \in \mathcal{Y}$*

$$\langle y, f(x^{-1})y' \rangle = \langle y, f(x)y' \rangle \quad (1)$$

and (weakly) odd if

$$\langle y, f(x^{-1})y' \rangle = -\langle y, f(x)y' \rangle \quad (2)$$

It is easy to check that if f is odd then $\int_{\mathcal{X}} \langle y, f(x)y' \rangle dx = 0$.

$$\begin{aligned} \int_{\mathcal{X}} \langle y, f(x)y' \rangle dx &= \frac{1}{2} \left(\int_{\mathcal{X}} \langle y, f(x^{-1})y' \rangle dx + \int_{\mathcal{X}} \langle y, f(x)y' \rangle dx \right) \\ &= \frac{1}{2} \left(-\int_{\mathcal{X}} \langle y, f(x)y' \rangle dx + \int_{\mathcal{X}} \langle y, f(x)y' \rangle dx \right) = 0. \end{aligned}$$

Besides the product of an even and an odd function is odd. Indeed for all $f, g \in \mathcal{F}(\mathcal{X}; \mathcal{L}(\mathcal{Y}))$ where f is even and g odd, define $h(x) = \langle y, f(x)y' \rangle \langle y, g(x)y' \rangle$. Then we have

$$\begin{aligned} h(x^{-1}) &= \langle y, f(x^{-1})y' \rangle \langle y, g(x^{-1})y' \rangle \\ &= -\langle y, f(x)y' \rangle \langle y, g(x)y' \rangle = -h(x). \end{aligned} \quad (3)$$

2.4.3 Representation of Groups

Representations of groups are convenient tools that allows group-theoretic problems to be replaced by linear algebra problems. Let $GL(\mathcal{H})$ be the group of continuous isomorphism of \mathcal{H} , a Hilbert space, onto itself. A representation π of a LCA group \mathcal{X} in \mathcal{H} is an homomorphism π :

$$\pi : \mathcal{X} \rightarrow GL(\mathcal{H})$$

for which all the maps $\mathcal{X} \rightarrow \mathcal{H}$ defined for all $v \in \mathcal{H}$ as $x \mapsto \pi(x)v$, are continuous. The space \mathcal{H} in which the representation takes place is

called the representation space of π . A representation π of a group \mathcal{X} in a vector space \mathcal{H} defines an action defined for all $x \in \mathcal{X}$ by

$$\pi_x : \begin{cases} \mathcal{H} & \rightarrow \mathcal{H} \\ v & \mapsto \pi(x)v. \end{cases}$$

If for all $x \in \mathcal{X}$, $\pi(x)$ is a unitary operator, then the group representation π is said to be unitary (i.e. $\forall x \in \mathcal{X}$, $\pi(x)$ is isometric and surjective). Thus π is unitary when for all $x \in \mathcal{X}$

$$\pi(x)^* = \pi(x)^{-1} = \pi(x^{-1}).$$

The representation π of \mathcal{X} in \mathcal{H} is said to be irreducible when $\mathcal{H} \neq \{0\}$ and are the only two stable invariant subspaces under all operators $\pi(x)$ for all $x \in \mathcal{X}$. I.e. $\forall \mathcal{U} \subset \mathcal{H}$, $\mathcal{U} \neq \{0\}$, $\{\pi(x)v \mid \forall x \in \mathcal{X}, \forall v \in \mathcal{U}\} \neq \mathcal{U}$.

To study LCA groups we also introduce the left regular representation of \mathcal{X} acting on a Hilbert space of function $\mathcal{H} \subset \mathcal{F}(\mathcal{X}; \mathbb{C})$. For all $x, z \in \mathcal{X}$ and for all $f \in \mathcal{H}$,

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

The representation λ of \mathcal{X} defines an action λ_x on \mathcal{H} which is the translation of $f(x)$ by z^{-1} . With this definition one has for all $x, z \in \mathcal{X}$, $\lambda_x \lambda_z = \lambda_{x^{-1} \star z}$. Such representations are faithful, that is $\lambda_x = 1 \iff x = e$.

2.4.4 Characters

Locally Compact Abelian (LCA) groups are central to the general definition of Fourier Transform which is related to the concept of Pontryagin duality [11]. 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 $\omega : \mathcal{X} \rightarrow \mathbb{T}$ from \mathcal{X} to the set of complex numbers of unit module \mathbb{T} . 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

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

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 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	$+$	$(x, \omega) = \exp(i\langle x, \omega \rangle)$
$\mathbb{R}_{*,+}^d$	\mathbb{R}^d	\cdot	$(x, \omega) = \exp(i\langle \log(x), \omega \rangle)$
$(-c; +\infty)^d$	\mathbb{R}^d	\odot	$(x, \omega) = \exp(i\langle \log(x+c), \omega \rangle)$

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

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

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 [13] to handle histograms features especially useful in image recognition applications.

2.4.5 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

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

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

$$\forall x \in \mathcal{X}, \quad f(x) = \int_{\hat{\mathcal{X}}} \mathcal{F}[f](\omega)(x, \omega) d\omega. \quad (4)$$

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 (3) gives some examples of real Abelian groups with their associated dual and pairing. The interested reader can refer to Folland [11] 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:

$$\forall \omega \in \hat{\mathcal{X}}, \quad \mathcal{F}[f](\omega) = \int_{\mathbb{R}^d} e^{-i\langle \omega, x-z \rangle} f(x) dx, \quad (5)$$

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

A vector-valued Reproducing Kernel Hilbert Space (vv-RKHS) is a functional vector space defined as follow.

Definition 4 (vector-valued Reproducing Kernel Hilbert Space [6, 17]). *Let \mathcal{Y} be a (real or complex) Hilbert space. A vector-valued Reproducing Kernel Hilbert Space on \mathcal{X} , a locally compact second countable topological space is a Hilbert space \mathcal{H} such that*

1. *the elements of \mathcal{H} are functions from \mathcal{X} to \mathcal{Y} (i. e. $\mathcal{H} \subset \mathcal{F}(\mathcal{X}, \mathcal{Y})$);*
2. *for all $x \in \mathcal{X}$, there exists a positive constant C_x such that for all $f \in \mathcal{H}$*

$$\|f(x)\|_{\mathcal{Y}} \leq C_x \|f\|_{\mathcal{H}}. \quad (6)$$

An operator-valued kernel is defined here as a Operator-Valued Kernel of positive type Carmeli et al. [7].

Definition 5 (Operator-Valued Kernel of positive type acting on a complex space). *Given \mathcal{X} a locally compact second countable topological space and \mathcal{Y} a complex Hilbert Space, a map $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ is called an Operator-Valued Kernel of positive type kernel if*

$$\sum_{i,j=1}^N \langle K(x_i, x_j) y_j, y_i \rangle_{\mathcal{Y}} \geq 0, \quad (7)$$

for all x_1, \dots, x_N in \mathcal{X} , all y_1, \dots, y_N in \mathcal{Y} and $N \geq 1$.

if \mathcal{Y} is a complex Hilbert space, an Operator-Valued Kernel of positive type is always Hermitian, i. e. $K(x, z) = K(z, x)^*$. This gives rise to the following definition of Operator-Valued Kernel of positive type acting on a real Hilbert space.

Definition 6 (Operator-Valued Kernel of positive type acting on a real Hilbert space). *Given \mathcal{X} a locally compact second countable topological space and \mathcal{Y} a real Hilbert Space, a map $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ is called an Operator-Valued Kernel of positive type kernel if*

$$K(x, z) = K(z, x) \quad (8)$$

and

$$\sum_{i,j=1}^N \langle K(x_i, x_j) y_j, y_i \rangle_{\mathcal{Y}} \geq 0, \quad (9)$$

for all x_1, \dots, x_N in \mathcal{X} , all y_1, \dots, y_N in \mathcal{Y} and $N \geq 1$.

As in the scalar case any vector-valued Reproducing Kernel Hilbert Space defines a unique Operator-Valued Kernel of positive type and conversely an Operator-Valued Kernel of positive type defines a unique vector-valued Reproducing Kernel Hilbert Space.

Proposition 7. *Given a vector-valued Reproducing Kernel Hilbert Space there is a unique Operator-Valued Kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ of positive type.*

Proof. Given $x \in \mathcal{X}$, eq. (6) ensure that the evaluation map at x defined as

$$\text{ev}_x : \begin{cases} \mathcal{H} \rightarrow \mathcal{Y} \\ f \mapsto f(x) \end{cases}$$

is a bounded operator and the Operator-Valued Kernel K associated to \mathcal{H} is defined as

$$K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}) \quad K(x, z) = \text{ev}_x \text{ev}_z^*.$$

Since for all $y_i, y_j \in \mathcal{Y}$,

$$\sum_{i,j=1}^N \langle K(x_i, x_j) y_j, y_i \rangle_{\mathcal{Y}} = \left\langle \sum_{i=1}^N \text{ev}_{x_i}^* y_i, \sum_{j=1}^N \text{ev}_{x_j}^* y_j \right\rangle_{\mathcal{Y}} \geq 0,$$

the map K is of positive type. □

Given $x \in \mathcal{X}$, $K_x : \mathcal{Y} \rightarrow \mathcal{F}(\mathcal{X}; \mathcal{Y})$ 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 $K_x = \text{ev}_x^*$. As a consequence we have that

$$K(x, z) y = \text{ev}_x \text{ev}_z^* y = K_x^* K_z y = (K_z y)(x). \quad (10)$$

Some direct consequences follows from the definition.

1. The kernel reproduces the value of a function $f \in \mathcal{H}$ at a point $x \in \mathcal{X}$ since for all $y \in \mathcal{Y}$ and $x \in \mathcal{X}$, $\text{ev}_x^* y = K_x y = K(\cdot, x)y$ so that $\langle f(x), y \rangle_{\mathcal{Y}} = \langle f, K(\cdot, x)y \rangle_{\mathcal{H}}$.
2. The set $\{K_x y \mid \forall x \in \mathcal{X}, \forall y \in \mathcal{Y}\}$ is total in \mathcal{H} . Namely,

$$\left(\bigcup_{x \in \mathcal{X}} \text{Im } K_x \right)^{\perp} = \{0\}.$$

If $f \in (\bigcup_{x \in \mathcal{X}} \text{Im } K_x)^{\perp}$, then for all $x \in \mathcal{X}$, $f \in (\text{Im } K_x)^{\perp} = \text{Ker } K_x^*$, hence $f(x) = 0$ for all $x \in \mathcal{X}$ that is $f = 0$.

3. Finally for all $x \in \mathcal{X}$ and all $f \in \mathcal{H}$, $\|f(x)\|_{\mathcal{Y}} \leq \sqrt{\|K(x, x)\|_{\mathcal{Y}, \mathcal{Y}}} \|f\|_{\mathcal{H}}$. This come from the fact that $\|K_x\| = \|K_x^*\| = \sqrt{\|K(x, x)\|_{\mathcal{Y}, \mathcal{Y}}}$ and the operator norm is sub-multiplicative.

Additionally given an Operator-Valued Kernel of positive type, it defines a unique [vv-RKHS](#).

Proposition 8. *Given an Operator-Valued Kernel of positive type there $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$, there is a unique vector-valued Reproducing Kernel Hilbert Space \mathcal{H} on \mathcal{X} with reproducing kernel K .*

Proof. Let $K_{x,y} = K(\cdot, x)y \in \mathcal{F}(\mathcal{X}; \mathcal{Y})$. Let $\mathcal{H}_0 = \text{span}\{K_{x,y} \mid \forall x \in \mathcal{X}, \forall y \in \mathcal{Y}\} \subset \mathcal{F}(\mathcal{X}; \mathcal{Y})$. If $f = \sum_{i=1}^N c_i K_{x_i, y_i}$ and $g = \sum_{i=1}^N d_i K_{z_i, y'_i}$ are elements of \mathcal{H}_0 we have that

$$\sum_{j=1}^N \overline{d_j} \langle f(z_j), y'_j \rangle_{\mathcal{Y}} = \sum_{i,j=1}^N c_i \overline{d_j} \langle K(z_j, x_i) y_i, y'_j \rangle_{\mathcal{Y}} = \sum_{i=1}^N c_i \langle y_i, g(x_i) \rangle_{\mathcal{Y}},$$

so that the sesquilinear form on $\mathcal{H}_0 \times \mathcal{H}_0$

$$\langle f, g \rangle_{\mathcal{H}_0} = \sum_{i,j=1}^N c_i \overline{d_j} \langle K(z_j, x_i) y_i, y'_j \rangle_{\mathcal{Y}}$$

is well defined. Since K is an Operator-Valued Kernel of positive type, we have for all $f \in \mathcal{H}_0$ that $\langle f, f \rangle_{\mathcal{H}_0} \geq 0$. Since the sesquilinear form is positive, it is also Hermitian. Choosing $g = K_{x,y}$ in the above definition yields for all $x \in \mathcal{X}$, all $f \in \mathcal{H}_0$ and all $y \in \mathcal{Y}$

$$\langle f, K_{x,y} \rangle_{\mathcal{H}_0} = \langle f(x), y \rangle_{\mathcal{Y}}.$$

Besides if $f \in \mathcal{H}_0$ for all unitary vector $y \in \mathcal{Y}$, by the Cauchy-Schwartz inequality we have

$$\begin{aligned} |\langle f(x), y \rangle_{\mathcal{Y}}| &= |\langle f, K_{x,y} \rangle_{\mathcal{H}_0}| \leq \sqrt{\langle f, f \rangle_{\mathcal{H}_0}} \sqrt{\langle K_{x,y}, K_{x,y} \rangle_{\mathcal{H}_0}} \\ &= \sqrt{\langle f, f \rangle_{\mathcal{H}_0}} \sqrt{\langle K(x, x)y, y \rangle_{\mathcal{Y}}} \leq \sqrt{\langle f, f \rangle_{\mathcal{H}_0}} \sqrt{\|K(x, x)\|_{\mathcal{Y}, \mathcal{Y}}}, \end{aligned}$$

which implies that

$$\|f(x)\|_{\mathcal{Y}} \leq \|f\|_{\mathcal{H}_0} \sqrt{\|K(x, x)\|_{\mathcal{Y}, \mathcal{Y}}}$$

Hence if $\langle f, f \rangle_{\mathcal{H}_0}$ then $f = 0$. Eventually we deduce that $\langle \cdot, \cdot \rangle_{\mathcal{H}_0}$ is a scalar product on \mathcal{H}_0 . Hence \mathcal{H}_0 is a pre-Hilbert space. To make it a (complete) Hilbert space we need to take the completion of this space. Let \mathcal{H} be the completion of \mathcal{H}_0 . Moreover let $K_x : \mathcal{Y} \rightarrow \mathcal{H}$ where $K_x y = K_{x,y}$. By construction K_x is bounded and let $W : \mathcal{H} \rightarrow \mathcal{F}(\mathcal{X}; \mathcal{Y})$ where $(Wf)(x) = K_x^* f$. The operator W is injective. Indeed if $Wf = 0$ then for all $x \in \mathcal{X}$, $f \in \text{Ker } K_x^* = (\text{Im } f)^\perp$. Since the set $\cup_{x \in \mathcal{X}} \text{Im } K_x = \{K_x y \mid \forall x \in \mathcal{X}, \forall y \in \mathcal{Y}\}$ generates by definition \mathcal{H}_0 , we have $f = 0$. Since W is injective, we have for all $f_1, f_2 \in \mathcal{H}_0$ $(Wf_1)(x) = (Wf_2)(x) \implies f_1(x) = f_2(x)$ pointwise in \mathcal{H} so that we can identify \mathcal{H} with a subspace of $\mathcal{F}(\mathcal{X}; \mathcal{Y})$. Hence $K_x^* f = (Wf)(x) = f(x) = \text{ev}_x f$, showing that \mathcal{H} is a vector-valued Reproducing Kernel Hilbert Space with reproducing kernel

$$K_{\mathcal{H}}(x, z)y = (\text{ev}_z^* y)(x) = K(x, z)y.$$

The uniqueness of \mathcal{H} comes from the uniqueness of the completion of \mathcal{H}_0 up to an isometry. \square

The above theorem holds also if \mathcal{Y} is a real vector space provided we add the assumption that K is symmetric, i.e. $K(x, z) = K(z, x)$.

Since an Operator-Valued Kernel of positive type defines a unique [vv-RKHS](#) and conversely a [vv-RKHS](#) defines a unique Operator-Valued Kernel, we denote the Hilbert space \mathcal{H} endowed with the scalar product $\langle \cdot, \cdot \rangle$ respectively \mathcal{H}_K and $\langle \cdot, \cdot \rangle_K$. From now we refer to Operator-Valued Kernel of positive type as Operator-Valued Kernel whether they act on complex or real Hilbert spaces. As a consequence, given K an Operator-Valued Kernel of positive type, define $K_x = K(\cdot, x)$ we have

$$K(x, z) = K_x^* K_z \quad \forall x, z \in \mathcal{X}, \quad (11a)$$

$$\mathcal{H}_K = \overline{\text{span}} \{K_x y \mid \forall x \in \mathcal{X}, \forall y \in \mathcal{Y}\}. \quad (11b)$$

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

Proposition 9 (Feature Operator [7]). *Let \mathcal{H} be any Hilbert space and $\phi : \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}; \mathcal{H})$, with $\phi_x := \phi(x)$. Then the operator $W : \mathcal{H} \rightarrow \mathcal{F}(\mathcal{X}; \mathcal{Y})$ 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*

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

W^*W is the orthogonal projection onto

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

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

Proof. The operator $(Wg)(x) = \phi(x)^*g$ ensure that the nullspace of W is $\mathcal{N} = \text{Ker } W = \bigcap_{x \in \mathcal{X}} \text{Ker } \phi(x)^*$. Since $\phi(x)$ is bounded, $\phi(x)$ is a continuous operator, thus for all $x \in \mathcal{X}$, $\text{Ker } \phi(x)^*$ is closed so that \mathcal{N} is closed. Moreover,

$$\mathcal{N} = \text{Ker } W = \bigcap_{x \in \mathcal{X}} \text{Ker } \phi(x)^* = \bigcap_{x \in \mathcal{X}} (\text{Im } \phi(x))^\perp = \left(\bigcup_{x \in \mathcal{X}} \text{Im } \phi(x) \right)^\perp$$

So that $\mathcal{N}^\perp = \bigcup_{x \in \mathcal{X}} \text{Im } \phi(x)$ and the restriction of W to \mathcal{N}^\perp is injective.

Let $\mathcal{H}_K = \text{Im } W$ be a vector space. Define the unique inner product on \mathcal{H}_K such that W becomes a partial isometry from \mathcal{H} onto \mathcal{H}_K . This is possible by taking $\langle \cdot, \cdot \rangle_K := \langle (W|_{\mathcal{N}^\perp})^{-1} \cdot, (W|_{\mathcal{N}^\perp})^{-1} \cdot \rangle_{\mathcal{H}} = \langle W \cdot, W \cdot \rangle_{\mathcal{H}}$ which exist and is unique since the restriction of W to \mathcal{N}^\perp is injective. We call again this new partial isometry W . We show that \mathcal{H}_K is a vector-valued Reproducing Kernel Hilbert Space. Since W^*W is a projection on \mathcal{N}^\perp , given $f \in \mathcal{H}_K$, where $f = Wg$ and $g \in \mathcal{N}^\perp$ we have for all $x \in \mathcal{X}$

$$f(x) = (Wg)(x) = \phi(x)^*g = \phi(x)^*W^*Wg = (W\phi(x))^*f.$$

Since $\text{Ker } W$ is closed, W is bounded, and $\phi(x)$ is bounded by definition so that the evaluation map

$$\text{ev}_x = (W\phi(x))^*$$

is bounded so continuous. Then the reproducing kernel is given for all $x, z \in \mathcal{X}$ by

$$K(x, z) = \text{ev}_x \text{ev}_z^* = (W\phi(x))^* (W\phi(z)) = \phi(x)^* W^* W \phi(z) = \phi(x)^* \phi(z),$$

Since W^*W is the identity on $\text{Im } \phi(z)$. Hence \mathcal{H}_K is a [vv-RKHS](#) (see proof of proposition 7). \square

We call ϕ a *feature map*, W a *feature operator* and \mathcal{H} a *feature space*. Since W is a partial isometry from $(\text{Ker } W)^\perp$ onto \mathcal{H}_K , the map W allows us to identify \mathcal{H}_K with the closed subspace $(\text{Ker } W)^\perp$ of \mathcal{H} . In this work we mainly focus on the class kernels inducing a [vv-RKHS](#) of continuous functions. Such kernels are named \mathcal{Y} -Mercer kernels.

Definition 10 (\mathcal{Y} -Mercer kernel). *A reproducing kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ is called \mathcal{Y} -Mercer kernel if \mathcal{H}_K is a subspace of $\mathcal{C}(\mathcal{X}; \mathcal{Y})$.*

The following proposition characterize \mathcal{Y} -Mercer kernel in terms of the properties of a kernels rather than properties of the [vv-RKHS](#).

Proposition 11 (Characterization of \mathcal{Y} -Mercer kernel [7]). *Let K be a reproducing kernel. The kernel K is Mercer if and only if the function $x \mapsto \|K(x, x)\|$ is locally bounded and for all $x \in \mathcal{X}$ and all $y \in \mathcal{Y}$, $K_x y \in \mathcal{C}(\mathcal{X}; \mathcal{Y})$.*

Proof. If $\mathcal{H}_K \subset \mathcal{C}(\mathcal{X}; \mathcal{Y})$, then for all $x \in \mathcal{X}$ and all $y \in \mathcal{Y}$, $K_x y$ is an element of $\mathcal{C}(\mathcal{X}; \mathcal{Y})$ (see eq. (11b)). In addition for all $f \in \mathcal{H}_K$, $\|K_x^* f\| = \|f(x)\| \leq \|f\|_\infty$. Hence there exist a constant $M \in \mathbb{R}_+$ such that for all $x \in \mathcal{X}$, $\|K_x\| \leq M$. Therefore from eq. (11a), for all $x \in \mathcal{X}$, $\|K(x, x)\| = \|K_x^*\|^2 \leq M^2$. Conversely assume that the function $x \mapsto \|K(x, x)\|$ is locally bounded and $K_x y \in \mathcal{C}(\mathcal{X}; \mathcal{Y})$. For all $f \in \mathcal{H}_K$ and all $x \in \mathcal{X}$,

$$\|f(x)\| = \|f\|_K \sqrt{\|K(x, x)\|} \leq M \|f\|_K.$$

Thus convergence in \mathcal{H}_K implies uniform convergence. Since by assumption $\{K_x t \mid \forall x \in \mathcal{X}, \forall y \in \mathcal{Y}\} \subset \mathcal{C}(\mathcal{X}; \mathcal{Y})$, then the vector-valued Reproducing Kernel Hilbert Space $\mathcal{H}_K = \overline{\text{span}} \{K_x y \mid \forall x \in \mathcal{X}, \forall y \in \mathcal{Y}\} \subset \mathcal{C}$ is also a subset of $\mathcal{C}(\mathcal{X}; \mathcal{Y})$ by the uniform convergence theorem. \square

Lemma 12. *Let \mathcal{H}_K be a vector-valued Reproducing Kernel Hilbert Space of continuous function $f : \mathcal{X} \rightarrow \mathcal{Y}$. If \mathcal{X} and \mathcal{Y} are separable then \mathcal{H}_K is separable.*

Proof. The separability of \mathcal{X} assure that there exist a countable dense subset $\mathcal{X}_0 \subseteq \mathcal{X}$. Since \mathcal{Y} is separable,

$$\mathcal{S} = \bigcup_{x \in \mathcal{X}_0} \text{Im } K_x = \{K_x y \mid \forall x \in \mathcal{X}_0, \forall y \in \mathcal{Y}\} \subset \mathcal{H}_K$$

is separable too. We show that \mathcal{S} is total in \mathcal{H}_K so that \mathcal{H}_K is separable. If for all $x \in \mathcal{X}_0$, $f \in \mathcal{S}^\perp$, then $f \in \text{Ker } K_x^*$. Namely $f(x) = \text{ev}_x f = 0$. Since f is continuous and \mathcal{X}_0 is dense in \mathcal{X} , for all $x \in \mathcal{X}$, $f(x) = 0$ so $f = 0$. \square

Proposition 13 (Countable orthonormal basis for \mathcal{Y} -Mercer kernel [6]). *Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{Y}$ be a reproducing kernel where \mathcal{X} and \mathcal{Y} are separable spaces. If K is a \mathcal{Y} -Mercer kernel then \mathcal{H}_K has a countable orthonormal basis.*

Proof. From proposition 11 K is a \mathcal{Y} -Mercer kernel if and only if $\mathcal{H}_K \subset \mathcal{C}(\mathcal{X}; \mathcal{Y})$. Applying lemma 12 of [6], we have that \mathcal{H}_K is separable. Thus since \mathcal{H}_K is also a Hilbert space it has a countable orthonormal basis. \square

An important consequence is that if K is a \mathcal{Y} -Mercer and \mathcal{X} and \mathcal{Y} are separable then \mathcal{H}_K is isometrically isomorphic to ℓ^2 .

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** $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ we assume that \mathcal{X} is a locally compact second countable topological group with identity e .

Definition 14 (Shift-invariant **ovk**). A reproducing Operator-Valued Kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ is called *shift-invariant*³ if for all $x, z, t \in \mathcal{X}$,

$$K(x \star t, z \star t) = K(x, z). \quad (12)$$

³ Also referred to as translation-invariant **ovk**.

A shift-invariant kernel can be characterized 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 .

We recall the definition of left regular representation of \mathcal{X} acting on \mathcal{H}_K which is useful to study **LCA** groups. For all $x, z \in \mathcal{X}$ and for all $f \in \mathcal{H}_K$,

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

A group representation λ_z describes the group by making it act on a vector space (here \mathcal{H}_K) in a linear manner. In other words, the group representation lets us see a group as a linear operator which are well studied mathematical objects.

Proposition 15 (Kernel signature [7]). Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ be a reproducing kernel. The following conditions are equivalent.

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

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

$$K(x, z) = K_e^* \lambda_{x^{-1} \star z} K_e \quad \forall (x, z) \in \mathcal{X}^2. \quad (13a)$$

$$\|K(x, x)\| = \|K_e(e)\| \quad \forall x \in \mathcal{X} \quad (13b)$$

Proof. Assume eq. (13) holds true. Given $x, z \in \mathcal{X}$, eq. (10) and eq. (12) yields

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

Since K is a reproducing kernel, K_e is of completely positive type, so that proposition 15 item 2 holds true. Besides if proposition 15 item 2 holds true obviously the definition of a reproducing kernel (definition 5) is fulfilled so that holds true.

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

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

that is $\lambda_x K_t = K_{x \star t}$. Besides for all $y, y' \in \mathcal{Y}$ and all $x, z, t, t' \in \mathcal{X}$,

$$\begin{aligned} \langle \lambda_x K_t y, \lambda_x K_{t'} y' \rangle_K &= \langle K_{x \star t} y, K_{x \star t'} y' \rangle_K = \langle K(x \star t', x \star t) y, y' \rangle \\ &= \langle K(t', t) y, y' \rangle = \langle K_t y, K_{t'} y' \rangle_K \end{aligned}$$

This means that λ leaves the set $\{K_x y \mid \forall x \in \mathcal{X}, \forall y \in \mathcal{Y}\}$ invariant. Since $\{K_x y \mid \forall x \in \mathcal{X}, \forall y \in \mathcal{Y}\}$ is total in \mathcal{H}_K (see eq. (11b)), λ is surjective and because it also leaves the inner product invariant, the first two claims follow. \square

The notation K_e for the function of completely positive type associated with the reproducing kernel K is consistent with the definition given by eq. (10) since for all $x \in \mathcal{X}$ and all $y \in \mathcal{Y}$

$$(K_e y)(x) = K_e(x)y.$$

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 16 (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 kernel⁴ if for all $(x, z) \in \mathcal{X}^2$,*

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

where k is a scalar Mercer kernel.

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 [2], Baldassarre et al. [3], and Evgeniou, Michelli, and Pontil [10] 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 [9, 15, 19]. 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. [5] and Carmeli et al. [7]. In particular they show that a decomposable kernel is universal provided that the scalar kernel k is universal and the operator A is injective.

⁴ Some authors also refer to as separable kernels.

Proposition 17 (Kernels and Regularizers [2]). *Let $K(x, z) := k(x, z)A$ for all $x, z \in \mathcal{X}$ be a decomposable kernel where A is a matrix of size $p \times p$. Then for all $f \in \mathcal{H}_K$,*

$$\|f\|_K = \sum_{i,j=1}^p A_{ij}^\dagger \langle f_i, f_j \rangle_K \quad (14)$$

where $f_i = \langle f, e_i \rangle$ (resp $f_j = \langle f, e_j \rangle$), denotes the i -th (resp j -th) component of f .

Curl-free and divergence-free kernels provide an interesting application of operator-valued kernels [4, 16, 18] to *vector field* learning, for which input and output spaces have the same dimensions ($d = p$). Applications cover shape deformation analysis [18] and magnetic fields approximations [20]. These kernels discussed in [12] allow encoding input-dependent similarities between vector-fields.

Definition 18 (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*

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

and the curl-free kernel as

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

where $\nabla \nabla^T$ is the Hessian operator and I 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 [4, 16].



Part II

CONTRIBUTIONS

You can put some informational part preamble text here. Illo principalmente su nos. Non message *occidental* an-gloromanic 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.

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:

$$\tilde{f}(x) = \tilde{\Phi}_{\tau;D}(x) * \theta,$$

where $\phi : \mathcal{X} \rightarrow \mathcal{L}(\mathcal{H}, \mathcal{Y})$ 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 $f : x \mapsto \tilde{\Phi}_{\tau;D}(x) * \theta$ 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 $\tilde{K}(x, z) = \tilde{\Phi}_{\tau;D}(x) * \tilde{\Phi}_{\tau;D}(z)$ 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. [7] and Zhang, Xu, and Zhang [21] extends Bochner's theorem to any shift-invariant \mathcal{Y} -Mercer kernel.

Proposition 19 (Operator-valued Bochner's theorem [21]). *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 $M : \mathcal{B}(\mathcal{X}) \rightarrow \mathcal{L}_+(\mathcal{Y})$ such that for all $x, z \in \mathcal{X}$,

$$K(x, z) = \int_{\hat{\mathcal{X}}} \overline{(x \star z^{-1}, \omega)} dM(\omega), \quad (15)$$

where M belongs to the set of all the $\mathcal{L}_+(\mathcal{Y})$ -valued measures of bounded variation on the σ -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 19.

Although this theorem is central to the spectral decomposition of shift-invariant \mathcal{Y} -Mercer [ovk](#), the following results proved by Carmeli et al. [7] 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, μ) where A is an operator-valued function on $\hat{\mathcal{X}}$ and μ is a real-valued positive measure on $\hat{\mathcal{X}}$. Note that obviously such a pair is not unique and 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 20 (Carmeli et al. [7]). *Let μ be a positive measure on $\mathcal{B}(\hat{\mathcal{X}})$ and $A : \hat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$ such that $\langle A(\cdot)y, y' \rangle \in L^1(\mathcal{X}, d\mu)$ for all $y, y' \in \mathcal{Y}$ and $A(\omega) \succcurlyeq 0$ for μ -almost all ω . Then, for all $\delta \in \mathcal{X}$ and for all $y, y' \in \mathcal{Y}$,*

$$K_e(\delta) = \int_{\hat{\mathcal{X}}} \overline{(\delta, \omega)} A(\omega) d\mu(\omega) \quad (16)$$

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 [vw-RKHS](#) \mathcal{H}_K is embed in $L^2(\hat{\mathcal{X}}, d\mu; \mathcal{Y}')$ by mean of the feature operator

$$(Wg)(x) = \int_{\hat{\mathcal{X}}} \overline{(x, \omega)} B(\omega) g(\omega) d\mu(\omega), \quad (17)$$

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 , μ is still a bounded positive measure and we retrieve the Bochner theorem applied to the scalar case (??).

Proposition 20 shows that a given pair $(A, d\mu)$ characterize an [ovk](#). Namely given a measure $d\mu$ and a function A such that $\langle A(\cdot)y, y' \rangle \in L^1(\mathcal{X}, d\mu)$ for all $y, y' \in \mathcal{Y}$ and $A(\omega) \succcurlyeq 0$ for μ -almost all ω , it gives rise to an [ovk](#). Since $(A, d\mu)$ determine a unique kernel we can write $\mathcal{H}_{(A, d\mu)} \Rightarrow \mathcal{H}_K$ where K is defined as in eq. (16). 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 19 and proposition 20 is that the first one characterizes an **OVK** by a unique Positive Operator-Valued Measure (**POVM**), while the second one shows that the **POVM** that uniquely characterizes a \mathcal{Y} -Mercer **OVK** has an operator-valued density with respect to a scalar measure μ ; and that this operator-valued density is not unique.

Finally proposition 20 does not provide any *constructive* way to obtain the pair $(A, d\mu)$ that characterizes an **OVK**. The following section 3.2.2 is based on another proposition of Carmeli, De Vito, and Toigo and shows 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 μ is a probability measure.

3.2.2 Sufficient conditions of existence

While proposition 20 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. [7] gives the solution, and also provides a sufficient condition for proposition 20 to apply.

Proposition 21 (Carmeli et al. [7]). *Let K be a shift-invariant \mathcal{Y} -Mercer kernel. Suppose that $\forall z \in \mathcal{X}$ and $\forall y, y' \in \mathcal{Y}$, $\langle K_e(\cdot)y, y' \rangle \in L^1(\mathcal{X}, dx)$ where dx denotes the Haar measure on \mathcal{X} endowed with the group law \star . Define C such that for all $\omega \in \hat{\mathcal{X}}$ and for all y, y' in \mathcal{Y} ,*

$$\begin{aligned} \langle y, C(\omega)y' \rangle &= \int_{\mathcal{X}} (\delta, \omega) \langle y, K_e(\delta)y' \rangle d\delta \\ &= \mathcal{F}^{-1} [\langle y, K_e(\cdot)y' \rangle] (\omega) \end{aligned} \tag{18}$$

Then

1. $C(\omega)$ is a bounded non-negative operator for all $\omega \in \hat{\mathcal{X}}$,
2. $\langle y, C(\cdot)y' \rangle \in L^1(\hat{\mathcal{X}}, d\omega)$ for all $y, y' \in \mathcal{Y}$,
3. for all $\delta \in \mathcal{X}$ and for all y, y' in \mathcal{Y} ,

$$\begin{aligned} \langle y, K_e(\delta)y' \rangle &= \int_{\hat{\mathcal{X}}} \overline{(\delta, \omega)} \langle y, C(\omega)y' \rangle d\omega \\ &= \mathcal{F} [\langle y, C(\cdot)y' \rangle] (\delta). \end{aligned}$$

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

There have been a lot of confusion in the literature whether a kernel is the Fourier transform or inverse Fourier transform of a measure. However ?? clarify the relation between the Fourier transform and inverse Fourier transform for a translation invariant Operator-Valued Kernel. Notice that in the real scalar case the Fourier transform and inverse Fourier transform of a shift-invariant kernel are the same, while the difference is significant for [OVK](#).

The following lemma is a direct consequence of the definition of $C(\omega)$ as the Fourier transform of the adjoint of K_e and also helps simplifying the definition of [ORFF](#).

Lemma 22. *Let K_e be the signature of a shift-invariant \mathcal{Y} -Mercer kernel and let $\langle y, C(\cdot)y' \rangle = \mathcal{F}^{-1}[\langle y, K_e(\cdot)y' \rangle]$ for all $y, y' \in \mathcal{Y}$. Then*

1. $C(\omega)$ is self-adjoint and C is even.
2. $\mathcal{F}^{-1}[\langle y, K_e(\cdot)y' \rangle] = \mathcal{F}[\langle y, K_e(\cdot)y' \rangle]$.
3. $K_e(\delta)$ is self-adjoint and K_e is even.

Proof. For any function $f \in L^1(\mathcal{X}, dx)$ define the flip operator \mathcal{R} by

$$\mathcal{R}f(x) := f(x^{-1}).$$

for any shift invariant \mathcal{Y} -Mercer kernel, we have for all $\delta \in \mathcal{X}$, $K_e(\delta) = K_e(\delta^{-1})^*$. Indeed,

$$\begin{aligned} \mathcal{R}\langle y, K_e(x \star z^{-1})y' \rangle &= \langle y, K_e((x \star z^{-1})^{-1})y' \rangle = \langle y, K_e(x^{-1} \star z)y' \rangle \\ &= \langle y, K_e(x \star z^{-1})^*y' \rangle. \end{aligned}$$

Proposition 22 item 1: taking the Fourier transform yields,

$$\mathcal{F}^{-1}[\langle y, K_e(\cdot)y' \rangle] = \mathcal{F}^{-1}\mathcal{R}[\langle K_e(\cdot)y, y' \rangle] = \mathcal{R}\langle C(\cdot)y, y' \rangle.$$

Hence $C(\omega) = C(\omega^{-1})^*$.

Suppose that \mathcal{Y} is a complex Hilbert space. Since for all $\omega \in \hat{\mathcal{X}}$ $C(\omega)$ is bounded and non-negative so $C(\omega)$ is self-adjoint. Besides we have $C(\omega) = C(\omega^{-1})^*$ to C must be pair.

Suppose that \mathcal{Y} is a real Hilbert space. Then we have the additional hypothesis that $K_e(\delta) = K_e(\delta)^*$. Taking the Fourier transform yields that $C(\omega) = C(\omega)^*$. Since for any shift invariant \mathcal{Y} -Mercer kernel $C(\omega) = C(\omega^{-1})^*$ we also conclude that $C(\omega^{-1}) = C(\omega)$.

Proposition 22 item 2: simply, for all $y, y' \in \mathcal{Y}$, $\langle y, C(\omega^{-1})y' \rangle = \langle y, C(\omega)y' \rangle$ thus $\mathcal{F}^{-1}[\langle y, C(\cdot)y' \rangle] = \mathcal{F}\mathcal{R}[\langle y, C(\cdot)y' \rangle] = \mathcal{F}[\langle y, C(\cdot)y' \rangle]$.

Proposition 22 item 3: from proposition 22 item 2, $\mathcal{F}^{-1}[\langle y, K_e(\cdot)y' \rangle] = \mathcal{F}^{-1}\mathcal{R}\langle y, K_e(\cdot)y' \rangle$. By injectivity of the Fourier transform, K_e is even. Since $K_e(\delta) = K_e(\delta^{-1})^*$, we must have $K_e(\delta) = K_e(\delta)^*$. \square

Proposition 23 (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} endowed with the group law \star .*

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

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

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

Proof. From proposition 21 and lemma 22, by taking $\langle y, C(\omega)y' \rangle = \mathcal{F}^{-1}[\langle y, K_e(\cdot)y' \rangle](\omega) = \mathcal{F}[\langle y, K_e(\cdot)y' \rangle](\omega)$, we can write the following equality concerning the OVK signature K_e .

$$\begin{aligned} \langle y, K_e(\cdot)y' \rangle(\omega) &= \int_{\hat{\mathcal{X}}} \overline{(\delta, \omega)} \langle y, C(\omega)y' \rangle d\omega \\ &= \int_{\hat{\mathcal{X}}} \overline{(\delta, \omega)} \left\langle y, \frac{C(\omega)}{p(\omega)} y' \right\rangle p(\omega) d\omega \end{aligned}$$

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 μ . In this case we note $p(\omega) = p_\mu(\omega)$. Conclude by taking $A(\omega) = C(\omega)/p_\mu(\omega)$ and $d\mu(\omega) = p_\mu(\omega) d\omega$. \square

In the case where $\mathcal{Y} = \mathbb{R}^p$, we rewrite proposition 23 coefficient-wise by choosing an orthonormal basis (e_1, \dots, e_p) of \mathcal{Y} , such that for all $i, j \in \{1, \dots, p\}$,

$$\langle e_i, C(\omega)e_j \rangle = C(\omega)_{ij} = A(\omega)_{ij} p_\mu(\omega) = \mathcal{F}[K_e(\delta)_{ij}]. \quad (19)$$

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

$$K_e(x \star z^{-1})_{ij} = \mathcal{F}[A(\cdot)_{ij} p_\mu(\cdot)] \quad (20)$$

Remark 24. 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 25. *Let K be a shift-invariant \mathcal{Y} -Mercer Kernel such that for all y, y' in \mathcal{Y} , $\langle y, K_e(\cdot)y' \rangle \in L^1(\mathcal{X}, dx)$.*

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

$$\|f\|_K^2 = \int_{\hat{\mathcal{X}}} \frac{\left\langle \mathcal{F}[f](\omega), A(\omega)^\dagger \mathcal{F}[f](\omega) \right\rangle_{\mathcal{Y}}}{p_\mu(\omega)} d\omega. \quad (21)$$

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:

$$\begin{aligned} \mathcal{F}[\mathcal{F}^{-1}[f]](x) &= \int_{\hat{\mathcal{X}}} \overline{(x, \omega)} \mathcal{F}^{-1}[f](\omega) d\omega = f(x) \\ (Wg)(x) &= \int_{\hat{\mathcal{X}}} \overline{(x, \omega)} p_\mu(\omega) B(\omega) g(\omega) d\omega = f(x). \end{aligned}$$

By injectivity of the Fourier transform, $\mathcal{F}^{-1}[f](\omega) = p_\mu(\omega) B(\omega) g(\omega)$ μ -almost everywhere. From proposition 9 we have

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

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,

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

Using the fact that $\mathcal{F}^{-1}[\cdot] = \mathcal{FR}[\cdot]$ and $\mathcal{F}^2[\cdot] = \mathcal{R}[\cdot]$,

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

Conclude the proof by taking $d\mu(\omega) = p_\mu(\omega)d\omega$ and rewriting the integral as an expectation and using the fact that $A(\omega)p_\mu(\omega) = C(\omega) = C(\omega^{-1})$. \square

Note that if $K(x, z) = k(x, z)$ is a scalar kernel then for all ω in $\hat{\mathcal{X}}$, $A(\omega) = \mathbf{1}$. 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 20.

Proposition 26 (Fourier feature map). *If there exist an operator-valued function $B : \hat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y}, \mathcal{Y}')$ such that for all $y, y' \in \mathcal{Y}$, $\langle y, B(\omega)B(\omega)^*y' \rangle = \langle y, A(\omega)y' \rangle$ μ -almost everywhere and $\langle y, A(\omega)y' \rangle \in L^1(\hat{\mathcal{X}}, d\mu)$ then the operator ϕ_x defined for all y in \mathcal{Y} by*

$$(\phi_x y)(\omega) = (x, \omega)B(\omega)^*y, \quad (22)$$

is a feature map⁵ of some shift-invariant kernel K .

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 20 of Carmeli et al. [7]. \square

With this notation notice that $\phi : \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}; L^2(\hat{\mathcal{X}}, \mu; \mathcal{Y}'))$ 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_μ . As shown in proposition 26 it is always possible to find a pair $(A(\omega), d\mu)$ such that $d\mu$ is a probability measure and $\mathbf{E}_\mu(\overline{(\delta, \omega)}A(\omega))$.

⁵ 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](#).

Given a \mathcal{Y} -Mercer shift-invariant kernel K on \mathcal{X} , an approximation of K can be obtained using a decomposition (A, μ) 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 27. *If one can find $B : \hat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y}', \mathcal{Y})$ and a probability measure μ 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*

$$\tilde{\Phi}_{1:D}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D (x, \omega_j) B(\omega_j)^*, \quad \omega_j \sim \mu \quad (23)$$

is an approximated feature map of an Operator-Valued Kernel⁶.

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

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

where $A(\omega) = B(\omega)B(\omega)^*$. By assumption $\langle y, A(\cdot)y' \rangle \in L^1(\hat{\mathcal{X}}, \mu)$ and ω_j are i.i.d.. Hence from the strong law of large numbers and proposition 20,

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

in the weak operator topology. \square

The approximate feature map proposed in proposition 27 has direct link with the functional feature map defined in proposition 26 since we have for all $y \in \mathcal{Y}$

$$\begin{aligned} \tilde{\Phi}_{1:D}(x)y &= \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D (\phi_x y)(\omega_j) \\ &= \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D (x, \omega_j) B(\omega_j)^* y, \quad \omega_j \sim \mu. \end{aligned} \quad (24)$$

⁶ I.e. it satisfies $\tilde{\Phi}_{1:D}(x)^* \tilde{\Phi}_{1:D}(z) \xrightarrow[D \rightarrow \infty]{\text{a.s.}} K(x, z)$ where K is a \mathcal{Y} -Mercer [OVK](#).

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

Definition 28 (Random Fourier feature operator). *Let $\theta = \bigoplus_{j=1}^D \theta_j \in (\mathcal{Y}')^D$, where $\theta_j \in \mathcal{Y}'$. We call random Fourier feature operator the linear application $W : (\mathcal{Y}')^D \rightarrow \mathcal{F}(\mathcal{X}; \mathcal{Y})$ defined as follow.*

$$(\tilde{W}_D \theta)(\mathbf{x}) := \tilde{\Phi}_{\tau:D}(\mathbf{x})^* \theta = \frac{1}{D} \sum_{j=1}^D \overline{(\mathbf{x}, \omega_j)} B(\omega_j) \theta_j, \quad \omega_j \sim \mu. \quad (25)$$

The approximate feature operator is useful to show the relations between the approximate feature map with the functional feature map defined in proposition 26.

Proposition 29. *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. $\tilde{\Phi}_{\tau:D}(\mathbf{x})^* \theta \xrightarrow[D \rightarrow \infty]{\text{a.s.}} \Phi_{\mathbf{x}}^* g,$
2. $\|\theta\|_{\mathcal{Y}}^2 \xrightarrow[D \rightarrow \infty]{\text{a.s.}} \|g\|_{L^2(\hat{\mathcal{X}}, d\mu; \mathcal{Y}')}^2.$

Proof. Proof of proposition 29 item 1: since $\omega_1, \dots, \omega_D$ are i.i.d. random vectors, for all $\mathbf{y} \in \mathcal{Y}$ and for all $\mathbf{y}' \in \mathcal{Y}'$, $\langle \mathbf{y}, B(\cdot) \mathbf{y}' \rangle \in L^2(\hat{\mathcal{X}}, d\mu)$ and $g \in L^2(\hat{\mathcal{X}}, d\mu; \mathcal{Y})$, from the strong law of large numbers

$$\begin{aligned} \tilde{\Phi}_{\tau:D}(\mathbf{x})^* \theta &= \frac{1}{D} \sum_{j=1}^D \overline{(\mathbf{x}, \omega_j)} B(\omega_j) g(\omega_j), \quad \omega_j \sim \mu \\ &\xrightarrow[D \rightarrow \infty]{\text{a.s.}} \int_{\hat{\mathcal{X}}} \overline{(\mathbf{x}, \omega)} B(\omega) g(\omega) d\mu(\omega) = (Wg)(\mathbf{x}) := \Phi_{\mathbf{x}}^* g. \end{aligned}$$

Proof of proposition 29 item 2: again, since $\omega_1, \dots, \omega_D$ 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

$$\begin{aligned} \|\theta\|_{\mathcal{Y}}^2 &= \sum_{j=1}^D \|g(\omega_j)\|_{\mathcal{Y}}^2, \quad \omega_j \sim \mu \\ &\xrightarrow[D \rightarrow \infty]{\text{a.s.}} \int_{\hat{\mathcal{X}}} \|g(\omega)\|_{\mathcal{Y}}^2 d\mu(\omega) := \|g\|_{L^2(\hat{\mathcal{X}}, d\mu; \mathcal{Y}')}^2 \end{aligned}$$

□

Hence the sequence of function $\tilde{f}_D := \tilde{\Phi}_{\tau:D}(\cdot)^* \theta$ converges almost surely to a function $f \in \mathcal{H}_{(\mathcal{A}, d\mu)} \implies \mathcal{H}_{\mathcal{K}}$. Therefore, in light of eq. (21), it is possible to define an approximate feature map of an Operator-Valued Kernel from its regularization properties in the [vv-RKHS](#). Otherwise corollary 30 exhibit a construction of an [ORFF](#) directly from an [OVK](#).

Corollary 30. *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*

$$\tilde{\Phi}_{1:D}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D (x, \omega_j) B(\omega_j)^*, \quad \omega_j \sim \mu, \quad (26)$$

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

Proof. Find $(A(\omega), d_\mu)$ from proposition 23 and apply proposition 27. \square

We write $\tilde{\Phi}_{1:D}(x)^* \tilde{\Phi}_{1:D}(x) \approx K(x, z)$ when $\tilde{\Phi}_{1:D}(x)^* \tilde{\Phi}_{1:D}(x) \xrightarrow{\text{a.s.}} K(x, z)$ 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 ϕ_x i.e. $\tilde{\Phi}_{1:D}(x) \approx \phi_x$, when for all $y \in \mathcal{Y}$, $\langle y, K(x, z)y' \rangle = \langle \phi_x y, \phi_z y' \rangle \approx \langle \tilde{\Phi}_{1:D}(x)y, \tilde{\Phi}_{1:D}(z)y' \rangle := \tilde{K}(x, z)$ where ϕ_x is defined in the sense of proposition 26.

Remark 31. *We find a decomposition such that for all $j = 1, \dots, D$, $A(\omega_j) = B(\omega_j)B(\omega_j)^*$ either by exhibiting an analytic closed-form or using a numerical decomposition.*

Corollary 30 allows us to define algorithm 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, ω) 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 $\omega_j, j = 1, \dots, D$ from the probability distribution μ ;
 - 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 on*

$\mathcal{X} = \mathbb{R}^d$ invariant w.r.t. the group operation $+$. Then the function $\tilde{\Phi}_{\tau,D}$ defined as follow is an Operator-valued Random Fourier Feature of K_0 .

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

Proof. The (Pontryagin) dual of $\mathcal{X} = \mathbb{R}^d$ is $\hat{\mathcal{X}} = \mathbb{R}^d$, and the duality pairing is $\langle x - z, \omega \rangle = \exp(i \langle x - z, \omega \rangle)$. The kernel approximation yields:

$$\begin{aligned} \tilde{K}(x, z) &= \tilde{\Phi}_{\tau,D}(x)^* \tilde{\Phi}_{\tau,D}(z) \\ &= \frac{1}{D} \sum_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle & \sin \langle x, \omega_j \rangle \end{pmatrix} \begin{pmatrix} \cos \langle z, \omega_j \rangle \\ \sin \langle z, \omega_j \rangle \end{pmatrix} A(\omega_j) \\ &= \frac{1}{D} \sum_{j=1}^D \cos \langle x - z, \omega_j \rangle A(\omega_j) \\ &\xrightarrow[D \rightarrow \infty]{\text{a.s.}} \mathbf{E}_{\mu} [\cos \langle x - z, \omega \rangle A(\omega)] \end{aligned}$$

in the weak operator topology. Since for all $x \in \mathcal{X}$, $\sin(x, \cdot)$ is an odd function and $A(\cdot) p_{\mu}(\cdot)$ is even,

$$\mathbf{E}_{\mu} [\cos \langle x - z, \omega \rangle A(\omega)] = \mathbf{E}_{\mu} [\exp(-i \langle x - z, \omega \rangle) A(\omega)] = K(x, z).$$

Hence $\tilde{K}(x, z) \xrightarrow[D \rightarrow \infty]{\text{a.s.}} K(x, z)$. \square

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

Example 2 (Matrix-valued kernel on the skewed multiplicative group). In the following, $K(x, z) = K_{1-c}(x \odot z)$ is a \mathbb{R}^p -Mercer matrix-valued kernel on $\mathcal{X} = (-c; +\infty)^d$ invariant w.r.t. the group operation⁷ \odot . Then the function $\tilde{\Phi}_{\tau,D}$ defined as follow is an Operator-valued Random Fourier Feature of K_{1-c} .

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

Proof. The dual of $\mathcal{X} = (-c; +\infty)^d$ is $\hat{\mathcal{X}} = \mathbb{R}^d$, and the duality pairing is $\langle x \odot z^{-1}, \omega \rangle = \exp(i \langle \log(x \odot z^{-1} + c), \omega \rangle)$ (see Li, Ionescu, and Sminchisescu [14]). Following the proof of example 1, we have

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

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

⁷ The group operation \odot is defined in ??.

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

4.1 BACKGROUND

4.2 THE NYSTRÖM METHOD

4.3 SUB-SAMPLING THE DATA

4.4 CONCLUSIONS



Part III

FINAL WORDS

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.

CONCLUSIONS

Part IV

APPENDIX



OPERATOR-VALUED FUNCTIONS AND INTEGRATION

BIBLIOGRAPHY

- [1] Erik M Alfsen. “A simplified constructive proof of the existence and uniqueness of Haar measure.” In: *Mathematica Scandinavica* 12.1 (1964), pp. 106–116.
- [2] 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.
- [3] 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.
- [4] L. Baldassarre, L. Rosasco, A. Barla, and A. Verri. “Multi-output learning via spectral filtering.” In: *Machine Learning* 87.3 (2012), pp. 259–301.
- [5] A. Caponnetto, C. A. Micchelli, M., and Y. Ying. “Universal Multi-Task Kernels.” In: *Journal of Machine Learning Research* 9 (2008), pp. 1615–1646.
- [6] 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.
- [7] 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.
- [8] John B Conway. *A course in functional analysis*. Vol. 96. Springer Science & Business Media, 2013.
- [9] 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.
- [10] T. Evgeniou, C. A. Micchelli, and M. Pontil. “Learning Multiple Tasks with kernel methods.” In: *JMLR* 6 (2005), pp. 615–637.
- [11] Gerald B Folland. *A course in abstract harmonic analysis*. CRC press, 1994.
- [12] E. Fuselier. “Refined Error Estimates for Matrix-Valued Radial Basis Functions.” PhD thesis. Texas A&M University, 2006.

- [13] 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](https://doi.org/10.1007/978-3-642-15986-2_27). URL: http://dx.doi.org/10.1007/978-3-642-15986-2_27.
- [14] 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.
- [15] 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.
- [16] 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.
- [17] C. A. Micchelli and M. A. Pontil. "On Learning Vector-Valued Functions." In: *Neural Computation* 17 (2005), pp. 177–204.
- [18] M. Micheli and J. Glaunes. *Matrix-valued kernels for shape deformation analysis*. Tech. rep. Arxiv report, 2013.
- [19] 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.
- [20] N. Wahlström, M. Kok, T.B. Schön, and Fredrik Gustafsson. "Modeling magnetic fields using Gaussian processes." In: *in Proc. of the 38th ICASSP*. 2013.
- [21] 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

COLOPHON

This document was typeset using the typographical look-and-feel `classicthesis` developed by André Miede. The style was inspired by Robert Bringhurst's seminal book on typography "*The Elements of Typographic Style*". `classicthesis` is available for both \LaTeX and \LyX :

<https://bitbucket.org/amiede/classicthesis/>

Happy users of `classicthesis` usually send a real postcard to the author, a collection of postcards received so far is featured here:

<http://postcards.miede.de/>

Final Version as of November 1, 2016 (`classicthesis` version 0.1).