

NNT: 2017SACLE024

Thèse de Doctorat soumise en vue de l'obtention du titre de

**DOCTEUR DE
L'UNIVERSITÉ PARIS-SACLAY**

Préparée à

**L'UNIVERSITÉ D'ÉVRY-VAL-D'ESSONNE
&
TÉLÉCOM PARISTECH**

Présentée par
M. Romain Brault¹

Intitulée

**LARGE-SCALE OPERATOR-VALUED
KERNEL REGRESSION**

Sous-titrée
Data Are Not Reals!

École Doctorale — 580

Sciences et Technologies de l'Information et de la Communication (STIC)
Laboratoire Informatique, Biologie Intégrative et Systèmes Complexes (IBISC)
Spécialité de mathématiques et d'informatique

Sous supervision de Professeur (Prof.) FLORENCE D'ALCHÉ-BUC²

Présentée et soutenue à l'Université d'Évry-Val-d'Essonne, le 3 Juillet 2017,

devant le jury composé de

Prof.	Florence	D'ALCHÉ-BUC	Télécom-ParisTech	Directeur,
Dr.	Aurélien	BELLET	INRIA, Lille	Examinateur
Prof.	Jean-Marc	DELOSME	Université d'Évry-Val-d'Essonne	Examinateur,
Prof.	Paul	HONEINE	Université de Rouen Normandie	Rapporteur,
Assoc. Prof.	Hachem	KADRI	Université de Aix-Marseille	Examinateur,
Prof.	Liva	RALAIVOLA	Université de Aix-Marseille	Rapporteur,
Assoc. Prof.	Zoltán	SZABÓ	École Polytechnique	Examinateur,
Assoc. Prof.	Marie	SZAFRANSKI	ENSIIE, Évry	Examinateur.

¹ Courriel: roman.brault@ibisc.fr

² Courriel: florence.dalche@telecom-paristech.fr

Engineer (Eng.) ENSIIE, Romain Brault: *Data Are Not Reals!* Large-scale Operator-Valued Kernel Regression, © Romain Brault, October 6, 2017

SUPERVISOR:

Professor (Prof.) Florence d'Alché-Buc

LOCATION:

46, Rue Barrault, 75013 — Paris, France

DIGITALY SIGNED DOCUMENT:

```
gpg --keyserver keys.gnupg.net \
--recv-keys A276D73294A106E2544FFF9E3E5B5D0B181C5E04
gpg --verify ThesisRomainBrault.pdf.asc ThesisRomainBrault.pdf
```

ABSTRACT

In this thesis we study scalable methods to perform regression with *Operator-Valued Kernels* in order to learn *vector-valued functions*.

When data present structure, or relations between them or their different components, a common approach is to treat the data as a vector living in an appropriate Hilbert space rather than a collection of real numbers. This representation allows to take into account the structure of the data by defining an appropriate space embedding the underlying structure. Thus many problems in machine learning can be cast into learning vector-valued functions. *Operator-Valued Kernels* and *vector-valued Reproducing Kernel Hilbert Spaces* provide a theoretical and practical framework to address the issue of learning vector-valued functions by naturally extending the well-known framework of scalar-valued kernels. In the context of scalar-valued functions learning, a scalar-valued kernel can be seen a similarity measure between two data points. A solution of the learning problem has the form of a linear combination of these similarities with respect to weights (to determine), in order to have the best “fit” of the data. When dealing with Operator-Valued Kernels, the evaluation of the kernel is no longer a scalar similarity, but a function (an operator) acting on vectors. A solution is then a linear combination of operators with respect to vector weights.

Although Operator-Valued Kernels generalize strictly scalar-valued kernels, large scale applications are usually not affordable with these tools that require an important computational power along with a large memory capacity. In this thesis, we propose and study scalable methods to perform regression with *Operator-Valued Kernels*. To achieve this goal, we extend Random Fourier Features, an approximation technique originally introduced for scalar-valued kernels, to *Operator-Valued Kernels*. The idea is to take advantage of an approximated operator-valued feature map in order to come up with a linear model in a finite dimensional space.

First we develop a general framework devoted to the approximation of shift-invariant Mercer kernels on Locally Compact Abelian groups and study their properties along with the complexity of the algorithms based on them. Second we show theoretical guarantees by bounding the error due to the approximation, with high probability. Third, we study various applications of Operator Random Fourier Features to different tasks of Machine learning such as multi-class

classification, multi-task learning, time series modeling, functional regression and anomaly detection. We also compare the proposed framework with other state of the art methods. Fourth, we conclude by drawing short-term and mid-term perspectives.



"We have at our command computers with adequate data-handling ability and with sufficient computational speed to make use of machine-learning techniques, but our knowledge of the basic principles of these techniques is still rudimentary.

Lacking such knowledge, it is necessary to specify methods of problem solution in minute and exact detail, a time-consuming and costly procedure.

Programming computers to learn from experience should eventually eliminate the need for much of this detailed programming effort."

— Arthur Samuel [145]

ACKNOWLEDGEMENTS

First of all I would like to express my sincere gratitude to my advisor Prof. Florence d'Alché-Buc for guiding me, the continuous support during this Ph. D thesis and her insightful remarks.

Besides my advisor, I would also like to thank my thesis jury Prof. Liva Ralaivola and Prof. Paul Honeine for their time.

This thesis would not have been possible without my collaborators Dr. Markus Heinonen and Dr. Maxime Sangnier, whose help have been precious to develop fully my ideas and publish. I would like to thank Dr. Nicolas Goix, Dr. Nicolas Drougard and Maël Chiapino with whom it has been a pleasure to collaborate on a paper. Big thank to Alexandre Gramfort for helping me through the deployment of operalib in the context of the mission Paris-Saclay Center for Data Science.

I would also like to acknowledge the Université d'Évry-val-d'Essonne for hosting and funding¹ me for the first part of my thesis and Télécom ParisTech for the second part.

This manuscript is dedicated to my parents and family for backing me up during this three years and especially my grand parents for welcoming me any week-end to their place to rest.

I am also thankful to all my colleagues at Evry: Néhémy Lim, Yves-Stan Le Corne, Adel Mezime, Sébastien Maurais, Vincent Chau, Céline Brouard, Arnaud Fouchet, Frédéric Papadopoulos, Markus Heinonen, Laurent Poligny and Jean-Baptiste Meunier and also

¹ PhD grant numbered 76391.

my colleague of Télécom-ParisTech: Maxime Sangnier, Paul Lagrée, Claire Vernade, Moussab Djerab, Alexandre Garcia, Nicolas Goix, Maël Chiapino, Ana Korba, Albert Thomas, Adil Salim and Jean Lafont for the great conversations at work as well as the pub!



CONTENTS

I INTRODUCTION	1
1 OUTLINE AND MOTIVATIONS	3
1.1 Motivation	4
1.2 Outline	4
2 ON LEARNING EFFICIENTLY SCALAR-VALUED FUNCTIONS	7
2.1 About statistical learning	8
2.1.1 Introduction to kernel methods	11
2.1.2 Towards large scale learning with kernels	14
3 BACKGROUND	21
3.1 Notations	22
3.1.1 Algebraic structures	22
3.1.2 Topology and continuity	22
3.1.3 Measure theory	23
3.1.4 Vector spaces, linear operators and matrices	24
3.2 Elements of abstract harmonic analysis	26
3.2.1 Locally compact Abelian groups	26
3.2.2 The Haar measure	29
3.2.3 Even and odd functions	30
3.2.4 Characters	30
3.2.5 The Fourier Transform	33
3.2.6 Representations of Groups	35
3.3 On Operator-Valued Kernels	36
3.3.1 Definitions and properties	36
3.3.2 Shift-Invariant OVK on LCA groups	43
3.3.3 Examples of Operator-Valued Kernels	44
3.3.4 Some use of Operator-valued kernels	48
II CONTRIBUTIONS	51
4 OPERATOR-VALUED RANDOM FOURIER FEATURES	53
4.1 Motivation	54
4.2 Theoretical study	54
4.2.1 Sufficient conditions of existence	56
4.2.2 Examples of spectral decomposition	62
4.2.3 Functional Fourier feature map	67
4.3 Operator-valued Random Fourier Features	69
4.3.1 Building Operator-valued Random Fourier Features	69
4.3.2 From Operator Random Fourier Feature maps to OVKs	71
4.3.3 Examples of Operator Random Fourier Feature maps	75

4.3.4	Regularization property	77
4.4	Conclusions	79
5	BOUNDED THE ERROR OF THE ORFF APPROXIMATION	81
5.1	Convergence with high probability of the ORFF estimator	82
5.1.1	Random Fourier Features in the scalar case and decomposable OVK	84
5.1.2	Uniform convergence of ORFF approximation on LCA groups	85
5.1.3	Dealing with infinite dimensional operators . .	90
5.1.4	Variance of the ORFF approximation	91
5.1.5	Application on decomposable, curl-free and divergence-free OVK	93
5.2	Conclusions	94
6	LEARNING WITH FEATURE MAPS	95
6.1	Learning with OVK	96
6.1.1	Supervised learning within VV-RKHS	96
6.1.2	Learning with Operator Random Fourier Feature maps	100
6.2	Solving ORFF-based regression	103
6.2.1	Gradient descent methods	106
6.2.2	Complexity analysis	108
6.3	Efficient learning with ORFF	109
6.3.1	Case of study: the decosubmposable kernel . .	110
6.3.2	Linear operators in matrix form	113
6.3.3	Curl-free kernel	117
6.3.4	Divergence-free kernel	117
6.3.5	Iterative solvers using linear operators	118
6.4	Experiments	118
6.4.1	Learning with ORFF vs learning with OVK . .	118
6.5	Conclusion	124
7	CONSISTENCY AND GENERALIZATION BOUND FOR ORFF	125
7.1	Generalization bound	126
7.1.1	Generalization by bounding the function space complexity	127
7.1.2	Algorithm stability	130
7.2	Consistency of learning with ORFF	131
7.3	Discussion	135
8	APPLICATION TO TIME SERIES MODELLING	137
8.1	Introduction	138
8.2	Operator-Valued Kernels for Vector Autoregression .	138
8.3	Operator-Valued Random Fourier Features	141
8.3.1	Numerical Performance	143
8.3.2	Simulated data	143
8.3.3	Influence of the number of random features . .	144
8.3.4	Real datasets	145
8.4	Discussion	146

III CONCLUSION AND WORK IN PROGRESS	147
9 WORK IN PROGRESS	149
9.1 Learning function-valued functions	150
9.1.1 Quantile regression	150
9.1.2 Functional output data	151
9.1.3 ORFF for functional output data	151
9.1.4 Many quantile regression	155
9.1.5 One-class SVM revisited	157
9.2 Operalib	160
10 CONCLUSION	163
10.1 Contributions	164
10.2 Perspectives	165
IV APPENDIX	167
A PROOFS OF THEOREMS	169
A.1 Proof of the error bound with high probability of the ORFF estimator	170
A.1.1 Epsilon-net	170
A.1.2 Bounding the Lipschitz constant	171
A.1.3 Bounding the error on a given anchor point .	172
A.1.4 Union Bound and examples	176
A.2 Proof of the ORFF estimator variance bound	181
B MISCELLANEOUS	183
B.1 Learning with semi-supervision	184
B.1.1 Representer theorem and feature equivalence .	184
B.1.2 Gradients	189
B.1.3 Complexity	191
C RELEVANT PIECE OF CODE	197
c.1 Python code for figure 3.1	198
c.2 Python code for figure 5.1	199
c.3 Python code for figure 5.3	201
c.4 Python code for figure 5.4	203
c.5 Python code for figure 5.5	206
c.6 Python (tensorflow) code for continuous quantile re- gression	209
D ONE CLASS SPLITTING CRITERIA FOR RANDOM FORESTS	213
D.1 Background on decision trees	217
D.2 Adaptation to the one-class setting	218
D.2.1 One-class splitting criterion	218
D.2.2 Prediction: scoring function of the forest . . .	222
D.2.3 OneClassRF: a Generic One-Class Random For- est algorithm	223
D.3 Benchmarks	224
D.3.1 Default parameters of OneClassRF	224
D.3.2 Results	224
D.4 Theoretical analysis	226

D.4.1	Underlying model	226
D.4.2	Adaptive approach	228
D.5	Conclusion	229
D.6	Further insights on the algorithm	230
D.6.1	Interpretation of parameter gamma	230
D.6.2	Alternative scoring functions	230
D.6.3	Alternative stopping criteria	231
D.6.4	Variable importance	231
D.7	Hyper-parameters of tested algorithms	231
D.8	Description of the datasets	233
D.9	Further details on benchmarks and outlier detection results	234

BIBLIOGRAPHY	245
--------------	-----

LIST OF FIGURES

Figure 2.1	Borel’s strong law of large numbers.	10
Figure 2.2	Separation of nested circles with linear classifier	11
Figure 2.3	A scalar-valued feature map	15
Figure 3.1	Riesz map, dual spaces and adjoints.	25
Figure 3.2	Synthetic 2D curl-free field	46
Figure 3.3	Synthetic 2D divergence-free field	47
Figure 4.1	Relationships between feature-maps.	68
Figure 4.2	Approximation of a function in a vv-RKHS using different realizations of Operator Random Fourier Feature	72
Figure 5.1	ORFF reconstruction error	82
Figure 5.2	decomposable ORFF variance bound	92
Figure 5.3	Curl-free ORFF variance bound	92
Figure 6.1	ORFF equivalence theorem.	104
Figure 6.2	ORFF equivalence theorem with overfitting. . . .	105
Figure 6.3	Efficient decomposable Gaussian ORFF	115
Figure 6.4	Efficient curl-free Gaussian ORFF	117
Figure 6.5	Efficient divergence-free Gaussian ORFF	118
Figure 6.6	Prediction Error in percent on the MNIST dataset versus D, the number of Fourier features	120
Figure 6.7	Empirical comparison between curl-free ORFF, curl-free OVK, independent ORFF, independent OVK on a synthetic vector field regression task.	121
Figure 6.8	Decomposable kernel on the third dataset: R^2 score vs number of data in the train set (N) .	122
Figure 6.9	Decomposable kernel on the third dataset: R^2 score vs number of data in the train set (N) for different number for different number of random samples (D).	123
Figure 9.1	Learning a continuous quantile function with ORFF regression.	152
Figure 9.2	Learning many quantile with joint OVK regression.	153
Figure 9.3	Learning a continuous quantile function with ORFF regression.	156
Figure 9.4	Continuous OCSVM: proportion of inlier with respect to nu	159
Figure 9.5	Continuous OCSVM for outlier detection . . .	161
Figure B.1	ORFF equivalence theorem (semi-supervised)	190

Figure B.2	ORFF multiclass semi-supervised	192
Figure D.1	Outliers distribution G in the naive and adaptive approach.	218
Figure D.2	Adaptative splitting criteria	220
Figure D.3	One Class Random Forest level-sets	222
Figure D.4	Illustration of the standard splitting criterion on two modes when the proportion γ varies. .	232
Figure D.5	Performances of the novelty detection algorithms	235
Figure D.6	Performances of the outlier detection algorithms	235
Figure D.7	ROC and PR curves for ONECLASSRF (novelty detection framework)	236
Figure D.8	ROC and PR curves for ONECLASSRF (outlier detection framework)	236
Figure D.9	ROC and PR curves for IFOREST (novelty detection framework)	236
Figure D.10	ROC and PR curves for IFOREST (outlier detection framework)	237
Figure D.11	ROC and PR curves for OCRLSAMPLING (novelty detection framework)	237
Figure D.12	ROC and PR curves for OCRLSAMPLING (outlier detection framework)	237
Figure D.13	ROC and PR curves for OCSM (novelty detection framework)	238
Figure D.14	ROC and PR curves for OCSM (outlier detection framework)	238
Figure D.15	ROC and PR curves for LOF (novelty detection framework)	238
Figure D.16	ROC and PR curves for LOF (outlier detection framework)	239
Figure D.17	ROC and PR curves for Orca (novelty detection framework)	239
Figure D.18	ROC and PR curves for Orca (outlier detection framework)	239
Figure D.19	ROC and PR curves for LSAD (novelty detection framework)	240
Figure D.20	ROC and PR curves for LSAD (outlier detection framework)	240
Figure D.21	ROC and PR curves for RFC (novelty detection framework)	240
Figure D.22	ROC and PR curves for RFC (outlier detection framework)	241

LIST OF TABLES

Table 3.1	Mathematical symbols and their signification (part 1).	27
Table 3.3	Mathematical symbols and their signification (part 2).	28
Table 3.5	Classification of Fourier Transforms in terms of their domain and transform domain.	33
Table 6.1	Time complexity of efficient linear-operators for different ORFF.	115
Table 6.3	Efficient linear-operators for different ORFF.	116
Table 6.5	Error (% of nMSE) on SARCOS dataset.	124
Table 8.1	Sequential SCV-MSE and computation times for VAR(1), ORFFVAR and OKVAR on synthetic data (Settings 1, 2 and 3).	144
Table 8.2	SVC-MSE with respect to D the number of random features for ORFFVAR.	145
Table 8.3	SCV-MSE and computation times for ORFF- VAR, VAR(1) and OKVAR on real datasets.	146
Table D.1	Original datasets characteristics	223
Table D.3	Results for the novelty detection setting.	225
Table D.4	Results for the outlier detection setting	242

ACRONYMS

ACML	Asian Conference in Machine Learning
AUC	Area Under the Curve
c. f.	confer
cum.	cumulative
ECML	European Conference in Machine Learning
e. g.	exempli gratia
FT	Fourier Transform
i. e.	id est
IForest	Isolation Forest
i. i. d.	independent identically distributed
KDD	The Association for Computing Machinery's Special Interest Group on Knowledge Discovery and Data Mining
L-BFGS-B	Limited-memory BroydenFletcherGoldfarb-Shanno algorithm for Bound constraint optimization
LCA	Locally Compact Abelian
LOF	Local Outlier Factor
LSAD	Least Squares Anomaly Detection
MGF	Moment Generating Function
N. A.	Not Available
NORMA	Naive Online regularized Risk Minimization Algorithm
OCRFsampling	One-Class Random Forest Sampling
OCSM	One-Class Support Vector Machine
OKVAR	Operator-Valued Kernel-Based Vector Autoregressive
OneClassRF	One-Class Random Forest
ONORMA	Operator-valued Naive Online regularized Risk Minimization Algorithm
ORFF	Operator-valued Random Fourier Feature
ORFFVAR	Operator-valued Random Fourier Feature Vector Autoregressive
OVK	Operator-Valued Kernel

PD	Positive definite
p. d. f	probability density function
POVM	Positive Operator-Valued Measure
PR	Precision Recall
PSD	Positive Semi-definite
RF	Random Forest
RFC	Random Forest Clustering
RFF	Random Fourier Feature
RKHS	Reproducing Kernel Hilbert Space
ROC	Receiver Operating Characteristic
r. v.	random variable
SCV	Sequential cross-validation
SCV-MSE	Sequential cross-validation Mean Squared Error
SPSD	Symmetric Positive Semi-definite
SVM	Support Vector Machine
UCI	University of California Irvine
VAR	Vector Autoregressive
VC-dimension	Vapnik-Chernonenkis dimension
VV-RKHS	Vector Valued Reproducing Kernel Hilbert Space
w. r. t.	with respect to



Part I
INTRODUCTION

1

OUTLINE AND MOTIVATIONS

In this chapter we present our motivations as well as the structure of the present manuscript.

Contents

1.1	Motivation	4
1.2	Outline	4

1.1 MOTIVATION

This thesis is dedicated to the definition of a general and flexible approach to learn vector-valued functions together with an efficient implementation of the learning algorithms. To achieve this goal, we study shallow architectures, namely the product of a (nonlinear) operator-valued feature $\tilde{\Phi}(x)$ and a parameter vector θ such that $\tilde{f}(x) = \tilde{\Phi}(x)^*\theta$, and combine two appealing methodologies: Operator-Valued Kernel Regression and Random Fourier Features.

Operator-Valued Kernels [5, 34, 41, 84, 113] extend the classic scalar-valued kernels to functions with values in some *output* Hilbert space. As in the scalar case, Operator-Valued Kernels (OVKs) are used to build Reproducing Kernel Hilbert Spaces (RKHS) in which representer theorems apply as for ridge regression or other appropriate loss functional. In these cases, learning a model in the RKHS boils down to learning a function of the form $f(x) = \sum_{i=1}^N K(x, x_i)\alpha_i$ where x_1, \dots, x_N are the training input data and each $\alpha_i, i = 1, \dots, N$ is a vector of the output space \mathcal{Y} , and each $K(x, x_i)$ is an operator on vectors of \mathcal{Y} .

However, OVKs suffer from the same drawbacks as classic (scalar-valued) kernel machines: they scale poorly to large datasets because they are exceedingly demanding in terms of memory and computations. We propose to approximate OVKs by extending a methodology called Random Fourier Features (RFFs) [12, 94, 139, 144, 164, 167, 191] so far developed to speed up scalar-valued kernel machines. The RFF approach linearizes a shift-invariant kernel model by generating explicitly an approximated feature map $\tilde{\varphi}$. RFFs has been shown to be efficient on large datasets and has been further improved by efficient matrix computations such as [94, “FastFood”] and [61, “SORF”], which are considered as the best large scale implementations of kernel methods, along with Nyström approaches proposed in Drineas and Mahoney [55]. Moreover thanks to RFFs, kernel methods have been proved to be competitive with deep architectures [51, 108, 192].

1.2 OUTLINE

Chapter 2. In this introductory chapter we recall some elements of the statistical learning theory started by Vapnik [177]. Then we recall kernel methods [9] which are used to construct spaces of *scalar-valued* functions (called RKHSs) that are used model and learn non linear dependencies from the data. We finish by a literature review on large-scale implementations of kernel methods based on random Fourier features [139] and the Nyström method [185].

Chapter 3. In this chapter, to conclude the introduction, we develop briefly the mathematical tools used throughout this manuscript. We give a full table of notations, and present elements of functional analysis [93] and abstract harmonic analysis [65]. Then we turn our attention to the case where the functions we want to learn are not real-valued, but vector-valued. To learn vector-valued functions we define Operator-Valued Kernels [41, 113] that generalize the scalar-valued kernel presented in [Chapter 2](#). We conclude by giving a non-exhaustive list of Operator-Valued Kernels along with the context in which they have been used.

Chapter 4. In this first contribution chapter we present a generalization of the RFF framework introduced in [Chapter 2](#) [29]. This is based on an operator-valued Bochner theorem proposed by Carmeli et al. [41]. We use this theorem to show how to construct an Operator-valued Random Fourier Feature (ORFF) from an OVK. Conversely we also show that it is possible to construct an ORFF from the regularization properties it induces rather than from an OVK. We give various examples of ORFF maps such as an ORFF map for the decomposable kernel, the curl-free kernel and the divergence-free kernel.

Chapter 5. In this contribution chapter we refine the bound on the OVK approximation with ORFF we first proposed in [29] and presented in [28]. It generalizes the proof technique of Rahimi and Recht [139] to OVK on LCA groups thanks to the recent results of Koltchinskii [92], Minsker [121], Sutherland and Schneider [167], and Tropp [175]. As a Bernstein bound it depends on the variance of the estimator for which we derive an “upper bound”.

Chapter 6. This contribution chapter focus on explaining how to define an efficient implementation and algorithm to train an ORFF model. First we recall the supervised ridge regression with OVK and the celebrated representer theorem [182]. Then we show under which conditions learning with an ORFF is equivalent to learn with a kernel approximation. Eventually we give the gradient for the ridge regression problem, useful to find an optimal solution with gradient descent algorithms, as well as a closed form algorithm. We conclude by showing how viewing ORFFs as linear operators rather than matrices yields a more efficient implementation and finish with some numerical applications on toy and real-world datasets.

Chapter 7. This contribution chapter deals with a generalization bound for the a regression problem with ORFF based on the results of Maurer [112] and Rahimi and Recht [140]. We also discuss the case of Ridge regression presented in [Chapter 6](#).

Chapter 8. This contribution chapter shows how to use the ORFF methodology for non-linear vector autoregression. It is an instantiation of the ORFF framework to $\mathcal{X} = \mathcal{Y} = (\mathbb{R}^d, +)$. We also give a generalization of a stochastic gradient descent [51] to ORFF. This is a joint work with Néhémy Lim and Florence d’Alché-Buc and has been published at a workshop of ECML. It is based on the previous work of Lim et al. [101] for time series vector autoregression with operator-valued kernels [30].

Chapter 9. To conclude our work we present some work in progress. We show practical applications of operator-valued kernels acting on an infinite dimensional space \mathcal{Y} . We give two examples. First we show how to generalize many quantile regression to learn a continuous function of the quantiles on the data. Second we apply the same methodology to the One-Class Support Vector Machine (OCSM) algorithm in order to learn a continuous function of all the level sets. We conclude by presenting Operalib, a python library developed during this thesis which aims at implementing OVK-based algorithms in the spirit of Scikit-learn [132].

2

ON LEARNING EFFICIENTLY SCALAR-VALUED FUNCTIONS

“For such a model there is no need to ask the question “Is the model true?”.
true?”.

If “truth” is to be the “whole truth” the answer must be “No”.
The only question of interest is “Is the model illuminating and
useful?”.
— George Box [27]

In this chapter we recall some elements of the statistical learning theory started by Vapnik [177]. Then we recall kernel methods [9] which are used to construct spaces of real-valued functions (called RKHS) that are used model and learn non linear dependencies from the data. We finish by a litterature review on large-scale implementation of kernel methods based on random Fourier features [139] and the Nyström method [185].

Contents

2.1	About statistical learning	8
2.1.1	Introduction to kernel methods	11
2.1.2	Towards large scale learning with kernels .	14

2.1 ABOUT STATISTICAL LEARNING

We focus on the context of supervised learning. Supervised learning aims at building a function that predicts an output from a given input, by exploiting a “training set” composed of pairs of observed inputs/outputs. Denote \mathcal{X} , an *input space* and \mathcal{Y} , the *output space*. In this chapter, $\mathcal{Y} \subseteq \mathbb{R}$. When $\mathcal{Y} = \{1, \dots, C\}$, we talk about *supervised classification*. When $\mathcal{Y} = \mathbb{R}$, supervised learning corresponds to usual *regression*. We are given an independent identically distributed (i. i. d.) sample of size N of training data $s = (x_i, y_i)_{i=1}^N$, drawn from an unknown but fixed joint probability law $\Pr(X, Y)$. We call *learning algorithm*, a function \mathcal{A} that takes a class of functions \mathcal{F} , a training sample s and returns a function in \mathcal{F} . The learning algorithm can be studied through many angles, from a computational point of view to a statistical point of view.

From a limited number of observations, we wish to build a function that captures the relationship between the two random variables X and Y . More specifically, we search for a function f in some class of functions, denoted \mathcal{F} and called the *hypothesis class* such that the $f \in \mathcal{F}$ makes good predictions for the pair (X, Y) distributed according $\Pr(X, Y)$. To convert this abstract goal into a mathematical definition, we define a local loss function $L : \mathcal{X} \times \mathcal{F} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ that evaluates the capacity of a function f to predict the outcome y from an input x .

Hence, the goal of supervised learning is to find a function $f \in \mathcal{F}$ that minimizes the following criterion, called the *true risk* associated to L :

$$\mathfrak{R}(f) = \mathbf{E}_{\Pr}[L(X, f, Y)], \quad (2.1)$$

using the training dataset. However, this definition comes with an important issue: we do not know $\Pr(X, Y)$ and thus we cannot compute this risk nor minimize it. A first proposition is to replace this true risk by its empirical counterpart, the *empirical risk*, i. e. the empirical mean of the loss computed on the training data:

$$\mathfrak{R}_{\text{emp}}(f, s) = \frac{1}{N} \sum_{i=1}^N L(x_i, f, y_i).$$

Since the training data are usually supposed to be i. i. d., the celebrated strong law of large numbers tells us that for any given function f in \mathcal{F} , the *empirical risk* converges almost surely to the true risk.

Intuitively the empirical risk measures the performance of a model on the training data, while the true risk measures the performance of a model with respect to all the possible experiments (even the ones that are not present in the training set). Although the convergence of

the empirical risk to the true risk is guaranteed by the strong law of large numbers, for a given value of N , the function produced by minimization of the empirical risk may suffer from overfitting, i. e. being too much adapted to the training data and having a poor behavior on new unseen data.

Generalization error bounds, first introduced by the seminal work of Vapnik [178] in the context of supervised binary classification and then largely studied in wider contexts (see for instance, Mohri, Rostamizadeh, and Talwalkar [122]), provide a tool to understand how the difference between the true risk and the empirical risk behaves given N , the size of the sample used to compute the empirical risk, and d , a measure of the capacity the hypothesis class. These bounds usually take the following form. For any $\delta \in (0, 1)$, with probability $1 - \delta$, the following holds for any function $f \in \mathcal{F}$ of capacity $|\mathcal{F}| \in \mathbb{R}$:

$$\mathfrak{R}(f) \leq \mathfrak{R}_{\text{emp}}(f, s) + C(\delta, N, |\mathcal{F}|)$$

Especially for the functions of interest f_s returned by a learning algorithm, we have

$$\mathfrak{R}(f_s) \leq \mathfrak{R}_{\text{emp}}(f_s, s) + C(\delta, N, |\mathcal{F}|).$$

Usually it is expected from the quantity $C(\delta, N, |\mathcal{F}|)$ to increase with the capacity of the class of functions $|\mathcal{F}|$, and to decrease when the number of points N increases. This suggests to control the complexity of the hypothesis class while minimizing the empirical risk. In other words, if the class of functions \mathcal{F} is not too big, we expect that a low empirical risk implies a low true risk in particular when the number of training points N is large. Also when δ goes to zero, it is expected for $C(\delta, N, |\mathcal{F}|)$ to go to infinity since $1 - \delta$ is the probability of the bound to be valid¹.

Most of the approaches in machine learning, and specifically in supervised learning, are based on regularizing approaches: in this case, learning algorithms minimize the empirical loss while controlling a penalty term on the model f . In [Subsection 2.1.1](#), we will choose an hypothesis class as an Hilbert space where the penalty can be expressed as the ℓ_2 norm in this Hilbert space.

There is a crucial difference between the strong law of large numbers and the generalization property of a learning algorithm. The strong law of large numbers holds *after* a model f has been selected and fixed in \mathcal{F} . Thus minimizing the empirical risk does not yield *ipso facto* a model that minimizes the true risk (which measures the adequation of the model on unseen data). This can be illustrated by an intuitive example adapted from Cornuéjols and Miclet [47, page 64] and the infinite monkey theorem.

¹ We give examples of such bounds in [Chapter 7](#).

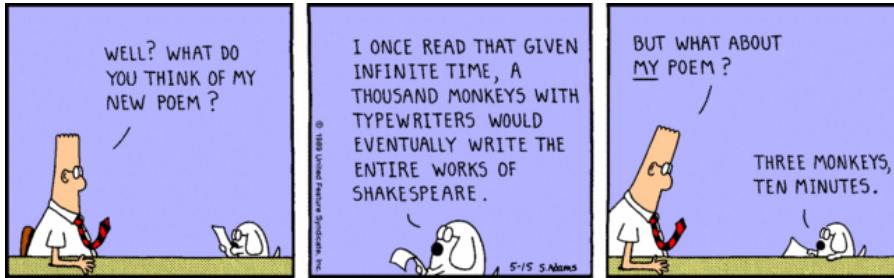


Figure 2.1: Borel's strong law of large numbers.

Example 2.1 Suppose we have a recruiter (a learning algorithm) whose task is to select the best students from a pool of candidates (the class of functions). Given ten students the recruiter makes them pass a test with N questions. If the exam is well constructed and there are enough questions the recruiter should be able to retrieve the best student.

Now suppose that ten million monkeys $\gg N$ take the test and answer randomly to the questions. Then with high probability a monkey will score better or as well as the best student (strong law of large numbers). Can we say then that the recruiter has identified the best student?

Intuitively we see that when the capacity of the class of function grows (the number of students and random monkeys), the performance of the best element a posteriori (minimizing the empirical risk) is not linked to the future performance (minimizing the true risk). In the present example we see that the capacity of the class of function is too large with respect to the number of data and thus presents a risk of overfitting.

On the contrary the generalization property ensures that the difference between the empirical risk and the true risk is controlled because the bound does not depend on a single fixed model, but on the whole class of functions. In this case if there are too many random monkeys, $C(\delta, N, |\mathcal{F}|)$ will blow-up, resulting in a poor generalization property.

A slightly stronger requirement is the consistency of learning algorithm. Given a loss function L and a class of function \mathcal{F} there exists a optimal solutions that minimize the true risk.

$$f_* \in \arg \min_{f \in \mathcal{F}} \mathfrak{R}(f).$$

The excess risk is defined as the difference between the empirical risk of a model returned by a learning algorithm and f_* . A learning algorithm is said to be consistent when it is possible to bound the excess risk uniformly over all the solutions returned by a learning algorithm.

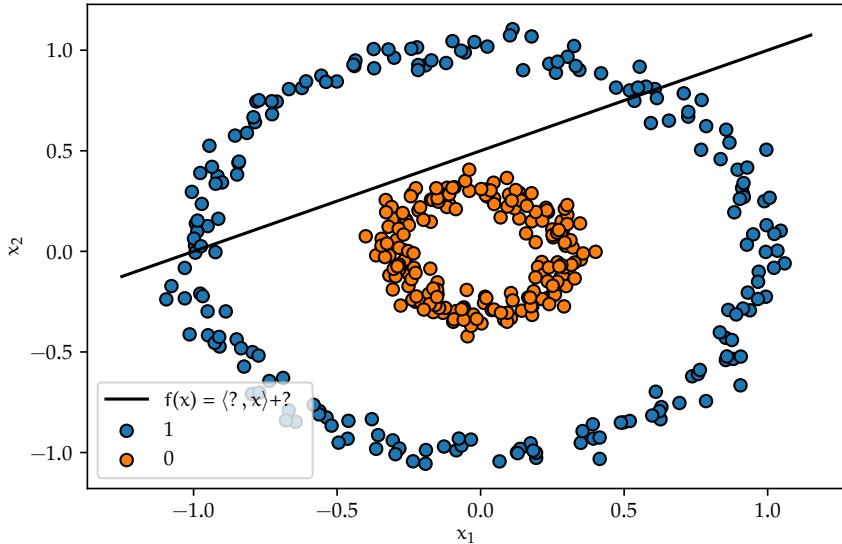


Figure 2.2: It is impossible to find a linear classifier that splits perfectly two nested circles.

2.1.1 Introduction to kernel methods

2.1.1.1 Kernels and Reproducing Kernel Hilbert Spaces

A fair simple choice for \mathcal{F} is the set of all linear functions. In this case we focus on defining learning algorithm picking up the “best” function(s) in the class

$$\mathcal{F}_{\text{lin.}} = \{ f \mid f(x) = \langle w, x \rangle + b, \forall w \in \mathbb{R}^d, \forall x \in \mathbb{R}^d, \forall b \in \mathbb{R} \}.$$

Although this class of functions has been well studied and has good generalization properties (as long as the norm of w is not too big), it has a rather low capacity. For instance in \mathbb{R}^2 it is impossible to separate two nested circles with a line (see Figure 2.2). On the other hand if one considers the class of functions of all functions $\{ f \mid f : \mathcal{X} \rightarrow \mathbb{R} \}$, this space contains too many functions for any algorithm to be able to find a solution to the minimization of the empirical risk. The idea of kernel methods [9, 20, 24, 89, 154] is to work in a subset of the set of all functions, namely a Reproducing Kernel Hilbert Space (RKHS), associated to a well chosen positive semi-definite and symmetric function (*a kernel*).

Definition 2.1 (Positive Definite kernels). Let \mathcal{X} be a locally compact second countable topological space. A kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is said to

be Positive Semi-definite (PSD) if for any $(x_1, \dots, x_N) \in \mathcal{X}^N$, the (Gram) matrix

$$\mathbf{K} = \left(k(x_i, x_j) \right)_{i=1, j=1}^{i=N, j=N} \in \mathcal{M}_{N,N}(\mathbb{R})$$

is Symmetric Positive Semi-definite (SPSD)².

The following proposition gives necessary and sufficient conditions to obtain a SPSD matrix:

Proposition 2.1 (SPSD matrix). \mathbf{K} is SPSD if and only if it is symmetric and one of the following assertions holds:

- The eigenvalues of \mathbf{K} are non-negative
- for any column vector $c = (c_1, \dots, c_N)^T \in \mathcal{M}_{N,1}(\mathbb{R})$,

$$c^T \mathbf{K} c = \sum_{i,j=1}^N c_i K_{ij} c_j \geq 0$$

One of the most important property of PD kernels [122] is that a PD kernel defines a unique RKHS. Note that the converse is also true.

Theorem 2.1 (Aronszajn [9]). Suppose k is a symmetric, positive definite kernel on a set \mathcal{X} . Then there is a unique Hilbert space of functions \mathcal{H} on \mathcal{X} for which k is a reproducing kernel, i.e.

$$\forall x \in \mathcal{X}, k(\cdot, x) \in \mathcal{H} \quad (2.2a)$$

$$\forall h \in \mathcal{H}, \forall x \in \mathcal{X}, h(x) = \langle h, k(\cdot, x) \rangle_{\mathcal{H}}. \quad (2.2b)$$

\mathcal{H} is called a reproducing kernel Hilbert space (Reproducing Kernel Hilbert Space) associated to k , and will be denoted, \mathcal{H}_k .

Another way to use Aronszajn's results is to state the feature map property for the PSD kernels.

Proposition 2.2 (Feature map). Suppose k is a symmetric, positive definite kernel on a set \mathcal{X} . Then, there exists a Hilbert space \mathcal{H} and a mapping φ from \mathcal{X} to \mathcal{H} such that:

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

The mapping φ is called a feature map and \mathcal{H} , a feature space.

Remark 2.1 Aronszajn's theorem tells us that there always exists at least one feature map, the so-called canonical feature map and the feature space associate, the Reproducing Kernel Hilbert Space \mathcal{H}_k

$$\varphi(x) = k(\cdot, x)$$

and $\mathcal{H} = \mathcal{H}_k$. However there exists several pairs of feature maps and features spaces for a given kernel k .

2 Note that for historical reasons valid kernels are called "Positive Definite kernels", although for any sequences of points the corresponding Gram matrix needs only to be (symmetric) Positive Semi-Definite [67].

2.1.1.2 Learning in Reproducing Kernel Hilbert Spaces

Back to learning and minimizing the empirical risk, a fair question is how to pick-up functions that minimize the empirical risk, in a space \mathcal{H}_k with infinite cardinality in polynomial time? The answer comes from the regularization and interpolation theory. To limit the size of the space in which we search for the function minimizing the empirical risk we add a regularization term to the empirical risk.

$$\begin{aligned}\mathfrak{R}_\lambda(f, s) &= \mathfrak{R}_{\text{emp}}(f, s) + \frac{\lambda}{2} \|f\|_{\mathcal{H}_k}^2 \\ &= \frac{1}{N} \sum_{i=1}^N L(x_i, f, y_i) + \frac{\lambda}{2} \|f\|_{\mathcal{H}_k}^2\end{aligned}$$

and we minimize \mathfrak{R}_λ instead of $\mathfrak{R}_{\text{emp}}$. Then the representer theorem (also called minimal norm interpolation theorem) states the following.

Theorem 2.2 (Representer theorem, Wahba [182]). *If f_s is a solution of*

$$\arg \min_{f \in \mathcal{H}_k} \mathfrak{R}_\lambda(f, s),$$

where $\lambda > 0$ then $f_s = \sum_{i=1}^N k(\cdot, x_i) \alpha_i$.

We note the vector $\alpha = (\alpha_i)_{i=1}^N$ and the matrix $\mathbf{K} = (k(x_i, x_k))_{i,k=1}^N$. Because of the representer theorem, stating that a solution of the empirical risk minimization is a linear combination of kernel evaluations weighted by a vector α , with mild abuse of notation we identify the function $f \in \mathcal{H}_k$ with the vector α . Thus we rewrite the loss $L(x, f, y)$ as $L(x, \alpha, y)$. Then we can rewrite

$$\mathfrak{R}_\lambda(\alpha, s) = \frac{1}{N} \sum_{i=1}^N L(x_i, \alpha, y_i) + \frac{\lambda}{2} \langle \alpha, \mathbf{K} \alpha \rangle_2,$$

and $f(x_i) = (\mathbf{K} \alpha)_i$ for any x_i in the training set. For instance if we choose $L(x, f, y) = \frac{1}{2} |f(x) - y|^2$ to be the least square loss, then

$$L(x_i, \alpha, y_i) = \frac{1}{2} |(\mathbf{K} \alpha)_i - y_i|^2.$$

In this case L is convex in α , thus it is possible to derive a polynomial time (in N) algorithm minimizing \mathfrak{R}_λ for the least square loss, which is called *kernel Ridge regression*:

$$\mathfrak{R}_\lambda(\alpha, s) = \frac{1}{2N} \|\mathbf{K} \alpha - (y_i)_{i=1}^N\|_2^2 + \frac{\lambda}{2} \langle \alpha, \mathbf{K} \alpha \rangle_2. \quad (2.3)$$

As a result of the representer theorem we see that we search a minimizer over $\alpha \in \mathbb{R}^N$ instead of $f \in \mathcal{H}_k$. By strict convexity and coercivity of \mathfrak{R}_λ , and because $\mathbf{K} + \lambda I_N$ is invertible³ for any $\lambda > 0$, the unique

³ Note that although $\mathbf{K} + \lambda I_N$ is always invertible if $\lambda > 0$, choosing a too small value of λ can lead to an ill-conditioned system if the eigenvalues of $\mathbf{K} + \lambda I_N$ are too small.

solution is $\alpha_s = \arg \min_{\alpha \in \mathbb{R}^N} \mathfrak{R}_\lambda(\alpha, s) = (\mathbf{K}/N + \lambda I_N)^{-1}(y_i)_{i=1}^N$. This is an $O(N^3)$ algorithm.

Another way of describing positive definite kernels and RKHS consists in defining a *feature map* $\varphi : \mathcal{X} \rightarrow \mathcal{H}$ where \mathcal{H} is a Hilbert space. Then any function in \mathcal{H}_k can be written $f(x) = \langle \varphi(x), \theta \rangle_{\mathcal{H}}$. In a nutshell the function φ is called feature map because it “extracts characteristic elements from a vector”. Usually a feature map takes a vector in an input space with low dimension and maps it to a potentially infinite dimensional Hilbert space. Put it differently, any function in \mathcal{H}_k is the composition of linear functional θ^* with a non linear feature map φ . Thus if the feature map φ is fixed (which is equivalent to fixing the kernel), it is possible to “learn” with a linear class of functions $\theta \in \mathcal{H}$ (see [Figure 2.3](#)). If we note

$$\varphi = \begin{pmatrix} \varphi(x_1) & \dots & \varphi(x_N) \end{pmatrix}$$

the “matrix” where each column represents the feature map evaluated at the point x_i with $1 \leq i \leq N$, the regularized risk minimization with the least square loss reads

$$\mathfrak{R}_\lambda(\theta, s) = \frac{1}{2N} \|\varphi^\top \theta - (y_i)_{i=1}^N\|_2^2 + \frac{\lambda}{2} \|\theta\|_2^2.$$

and the unique solution is $\theta_s = (\varphi \varphi^\top / N + \lambda I_{\mathcal{H}})^{-1} \varphi$. This is an $O_t(\dim(\mathcal{H})^2(N + \dim \mathcal{H}))$.

time complexity algorithm. This algorithm seems more appealing than its kernel counterpart when many data are given since once the space \mathcal{H} has been fixed, the algorithm is linear in the number of training points. However many questions remains. First although it is possible to design a feature map *ex nihilo*, can we design systematically a feature map from a kernel? For some kernels (e.g. the Gaussian kernel) it is well known that the Hilbert space corresponding to it has dimension $\dim(\mathcal{H}) = \infty$. Is it possible to find an approximation of the kernel such that $\dim(\mathcal{H}) < \infty$? If such a construction is possible and we know that N training data are available, is it possible to have a sufficiently good approximation ⁴ with $\dim(\mathcal{H}) \ll N$?

2.1.2 Towards large scale learning with kernels

Motivated by large scale applications, different methodologies have been proposed to approximate kernels and feature maps. This subsection briefly reminds the main approaches based on Random Fourier Features and Nyström techniques. Notice that another line of research

⁴ When $\dim(\mathcal{H}) \geq N$ then it is better to use the kernel algorithm than the feature algorithm. This is called the kernel trick.

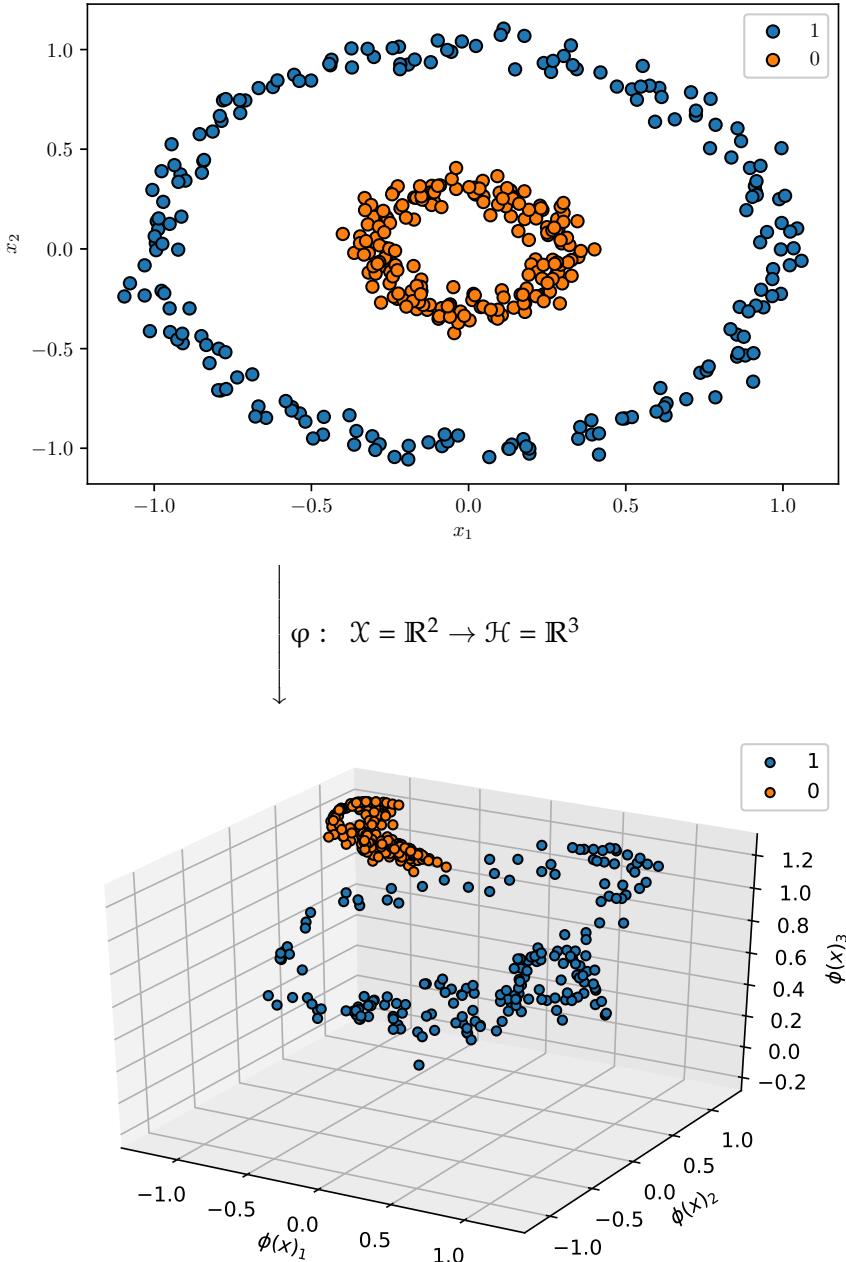


Figure 2.3: We map the two circles in \mathbb{R}^2 to \mathbb{R}^3 . In \mathbb{R}^3 it is now possible to separate the circles with a linear functional: a plane. We used the feature map

$$\varphi(x) = 0.82 \begin{pmatrix} \cos(1.76x_1 + 2.24x_2 + 2.75) \\ \cos(0.40x_1 + 1.87x_2 + 5.6) \\ \cos(0.98x_1 - 0.98x_2 + 6.05) \end{pmatrix}.$$

Here $\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ has been chosen as a realization of an RFF map (see Equation 2.5). A “cleaner” feature map adapted to this problem could have been

$$\varphi(x) = \begin{pmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{pmatrix}.$$

concerns online learning method such as NORMA developed in [90], later extended to the operator-valued kernel case by Audiffren and Kadri [10]. We start with the seminal work of Rahimi and Recht [139] who show that given a continuous shift-invariant kernel ($\forall x, z, t \in \mathcal{X}$, $k(x + t, z + t) = k(x, z)$), it is possible to obtain a feature map called RFF that approximate the given kernel.

2.1.2.1 Random Fourier Features map

The Random Fourier Features methodology introduced by Rahimi and Recht [139] provides a way to scale up kernel methods when kernels are Mercer and *translation-invariant*. We view the input space \mathcal{X} as a group endowed with the addition. Extensions to other group laws such as Li, Ionescu, and Sminchisescu [98] are described in [Subsubsection 4.2.2.2](#) within the general framework of operator-valued kernels.

Denote $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ a positive definite kernel on $\mathcal{X} = \mathbb{R}^d$. A kernel k is said to be *shift-invariant* or *translation-invariant* for the addition if for all $(x, z, t) \in (\mathbb{R}^d)^3$ we have $k(x + t, z + t) = k(x, z)$. Then, we define $k_0 : \mathbb{R}^d \rightarrow \mathbb{R}$ the function such that $k(x, z) = k_0(x - z)$. k_0 is called the *signature* of kernel k . Bochner's theorem [65] is the theoretical result that leads to the Random Fourier Features.

Theorem 2.3 (Bochner's theorem). *Any continuous positive definite function is the Fourier Transform of a bounded non-negative Borel measure.*

It implies that any positive definite, continuous and shift-invariant kernel k , has a continuous and positive semi-definite signature k_0 , which is the Fourier Transform \mathcal{F} of a non-negative measure μ . Hence we have $k(x, z) = k_0(x - z) = \int_{\mathbb{R}^d} \exp(-i\langle \omega, x - z \rangle) d\mu(\omega) = \mathcal{F}[k_0](\omega)$. Moreover $\mu = \mathcal{F}^{-1}[k_0]$. Without loss of generality, we assume that μ is a probability measure, i.e. $\int_{\mathbb{R}^d} d\mu(\omega) = 1$ by renormalizing the kernel since

$$\int_{\mathbb{R}^d} d\mu(\omega) = \int_{\mathbb{R}^d} \exp(-i\langle \omega, 0 \rangle) d\mu(\omega) = k_0(0).$$

and we can write the above equation as an expectation over μ . For all $x, z \in \mathbb{R}^d$

$$k_0(x - z) = \mathbf{E}_\mu [\exp(-i\langle \omega, x - z \rangle)].$$

Eventually, if k is real valued we only write the real part,

$$\begin{aligned} k(x, z) &= \mathbf{E}_\mu [\cos\langle \omega, x - z \rangle] \\ &= \mathbf{E}_\mu [\cos\langle \omega, z \rangle \cos\langle \omega, x \rangle + \sin\langle \omega, z \rangle \sin\langle \omega, x \rangle]. \end{aligned}$$

Let $\bigoplus_{j=1}^D x_j$ denote the D -length column vector obtained by stacking vectors $x_j \in \mathbb{R}^d$. The feature map $\tilde{\varphi} : \mathbb{R}^d \rightarrow \mathbb{R}^{2D}$ defined as

$$\tilde{\varphi}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos\langle x, \omega_j \rangle \\ \sin\langle x, \omega_j \rangle \end{pmatrix}, \quad \omega_j \sim \mathcal{F}^{-1}[k_0] \text{ i.i.d.} \quad (2.4)$$

is called a *Random Fourier Features* (map). Each $\omega_j, j = 1, \dots, D$ is independently and identically sampled from the inverse Fourier transform μ of k_0 . This Random Fourier Features map provides the following Monte-Carlo estimator of the kernel: $\tilde{k}(x, z) = \tilde{\varphi}(x)^* \tilde{\varphi}(z)$. Using trigonometric identities, Rahimi and Recht [139] showed that the same feature map can also be written

$$\tilde{\varphi}(x) = \sqrt{\frac{2}{D}} \bigoplus_{j=1}^D \left(\cos(\langle x, \omega_j \rangle + b_j) \right), \quad (2.5)$$

where $\omega_j \sim \mathcal{F}^{-1}[k_0]$, $b_j \sim \mathcal{U}(0, 2\pi)$ i. i. d.. The feature map defined by Equation 2.4 and Equation 2.5 have been compared in Sutherland and Schneider [167] where they give the condition under which Equation 2.4 has lower variance than Equation 2.5. For instance for the Gaussian kernel, Equation 2.4 has always lower variance. In practice, Equation 2.5 is easier to program. In this manuscript we focus on random Fourier feature of the form given in Equation 2.4.

The dimension D governs the precision of this approximation, whose uniform convergence towards the target kernel (as defined in ??) can be found in Rahimi and Recht [139] and in more recent papers with some refinements proposed in Sutherland and Schneider [167] and Sriperumbudur and Szabo [164]. Finally, it is important to notice that Random Fourier Features approach *only* requires two steps before the application of a learning algorithm: (1) define the inverse Fourier transform of the given shift-invariant kernel, (2) compute the randomized feature map using the spectral distribution μ . Rahimi and Recht [139] show that for the Gaussian kernel $k_0(x - z) = \exp(-\gamma \|x - z\|_2^2)$, the spectral distribution μ is a Gaussian distribution. For the Laplacian kernel $k_0(x - z) = \exp(-\gamma \|x - z\|_1)$, the spectral distribution is a Cauchy distribution.

We now focus on another famous way of obtaining feature maps for any scalar valued kernel called the Nyström method.

2.1.2.2 Nyström approximation

To overcome the bottleneck of Gram matrix computations in kernel methods, Williams and Seeger [185] have proposed to generate a low-rank matrix approximation of the Gram matrix using a subset of its columns. Since this feature map is based on a decomposition of the Gram matrix, the feature map resulting from the Nyström method is data dependent. Let $k : \mathcal{X}^2 \rightarrow \mathbb{R}$ be any scalar-valued kernel and let

$$\mathbf{s} = (x_i)_{i=1}^N$$

be the training data. We note a subsample of the training data

$$\mathbf{s}_M = (x_i)_{i=1}^M$$

where $M \leq N$ and s_M is a subsequence of s . Then construct the Gram matrix K_M on the subsequence s_M . Namely

$$K_M = \left(k(x_i, x_j) \right)_{i,j=1}^M.$$

Then perform the singular-valued decomposition $K_M = U \Lambda U^\top$. The Nyström feature map is given by

$$\tilde{\varphi}(x) = \Lambda^{-1/2} U^\top \left(\bigoplus_{i=1}^M k(x, x_i) \right).$$

Here M plays the same role as D in the RFF case: it controls the quality of the approximation. Let K be the full Gram matrix on the training data s , let

$$K_b = \left(k(x_i, x_j) \right)_{i=1, j=1}^{N, M}.$$

Then it is easy to verify that $\varphi^\top \varphi = K_b K_M^\dagger K_b^\top \approx K$, where K_M^\dagger is the pseudo-inverse of K_M and the quantity $K_b K_M^\dagger K_b^\top$ is a low rank approximation of the Gram matrix K .

2.1.2.3 Random features vs Nyström method

The main conceptual difference between the Nyström features and the Random Fourier Feature is that the Nyström construction is data dependent, while the RFF is not. The advantage of random Fourier feature lies in their fast construction. For N data in \mathbb{R}^d , it costs $O(N D d)$ to featurize all the data. For the Nyström features it costs $O(M^2(M + d))$. Moreover if one desires to add a new feature, the RFF methodology is as simple as drawing a new random vector $\omega \sim \mathcal{F}^{-1}[k_0]$, compute $\cos(\langle \omega, x \rangle + b)$, where $b \sim \mathcal{U}(0, 2\pi)$ and concatenate it the existing feature. For the Nyström features one needs to recompute the singular value decomposition of the new augmented Gram matrix K_{M+1} .

To analyse the RFF and Nyström features authors usually study the approximation error of the approximate Gram matrix and the target kernel $\|\varphi^\top \varphi - K\|$ (see [55, 142, 190]) or the supremum of the error between the approximated kernel and the true kernel over a compact subset X of the support if k : $\sup_{(x,z) \in \mathcal{C} \subseteq X^2} |\tilde{\varphi}(x)^\top \tilde{\varphi}(z) - k(x, z)|$ (see Bach [12], Rahimi and Recht [139], Rudi, Camoriano, and Rosasco [144], and Sutherland and Schneider [167]). Because Bartlett and Mendelson [17] showed that for generalization error to be below $\epsilon \in \mathbb{R}_{>0}$ for kernel methods is $O(N^{-1/2})$, the number of samples M or D required to reach some approximation error below ϵ should not grow faster than $O(M^{-1/2})$ for the Nyström method or $O(D^{-1/2})$ for the RFF method to match kernel learning. Concerning the Nyström method, Yang et al. [190] suggest

that the number of samples M is reduced to $O(M^{-1})$ to reach an error below ϵ when the gap between the eigenvalues of \mathbf{K} is large enough. As a result in this specific case, one should sample $M = O(\sqrt{N})$ Nyström features to ensure good generalization. On the other hand Rahimi and Recht [140] reported that the generalization performance of RFF learning is $O(N^{-1/2} + D^{-1/2})$, which indicates that $D = O(N)$ features should be sampled to generalize well. As a result the complexity of learning with the RFF seems not to decrease. However the bounds of Rahimi and Recht [140] are suboptimal and very recently (end of 2016) Rudi, Camoriano, and Rosasco [144] proved that in the case of ridge regression (Equation 2.3), the generalization error is $O(N^{-1/2} + D^{-1})$ meaning that $D = O(\sqrt{N})$ random features are required for good generalization with RFFs. We refer the interested reader to Yang et al. [190] for an empirical comparison between the Nyström method and the RFF method.

2.1.2.4 *Extensions of the RFF method*

The seminal idea of Rahimi and Recht [139] has opened a large literature on random features. Nowadays, many classes of kernels other than translation invariant are now proved to have an efficient random feature representation. Kar and Karnick [87] proposed random feature maps for dot product kernels (rotation invariant) and Hamid et al. [77] improved the rate of convergence of the approximation error for such kernels by noticing that feature maps for dot product kernels are usually low rank and may not utilize the capacity of the projected feature space efficiently. Pham and Pagh [135] proposed fast random feature maps for polynomial kernels.

Li, Ionescu, and Sminchisescu [98] generalized the original RFF of Rahimi and Recht [139]. Instead of computing feature maps for shift-invariant kernels on the additive group $(\mathbb{R}^d, +)$, they used the generalized Fourier transform on any locally compact abelian group to derive random features on the multiplicative group $(\mathbb{R}_{>0}^d, *)$. In the same spirit Yang et al. [189] noticed that an theorem equivalent to Bochner's theorem exists on the semi-group $(\mathbb{R}_+^d, +)$. From this they derived "Random Laplace" features and used them to approximate kernels adapted to learn on histograms.

To speed-up the convergence rate of the random features approximation, Yang et al. [188] proposed to sample the random variable from a quasi Monte-Carlo sequence instead of i.i.d. random variables. Le, Sarlós, and Smola [94] proposed the "Fastfood" algorithm to reduce the complexity of computing a RFF –using structured matrices and a fast Walsh-Hadamard transform– from $O_t(Dd)$ to $O_t(D \log(d))$. More recently Felix et al. [61] proposed also an algorithm "SORF" to compute Gaussian RFF in $O_t(D \log(d))$ but with bet-

ter convergence rates than “Fastfood” [94]. Mukuta and Harada [124] proposed a data dependent feature map (comparable to the Nyström method) by estimating the distribution of the input data, and then finding the eigenfunction decomposition of Mercer’s integral operator associated to the kernel.

In the context of large scale learning and deep learning, Lu et al. [108] showed that RFFs can achieve performances comparable to deep-learning methods by combining multiple kernel learning and composition of kernels along with a scalable parallel implementation. Dai et al. [51] and Xie, Liang, and Song [186] combined RFFs and stochastic gradient descent to define an online learning algorithm called “Doubly stochastic gradient descent” adapted to large scale learning. Yang et al. [192] proposed and studied the idea of replacing the last fully interconnected layer of a deep convolutional neural network [95] by the “Fastfood” implementation of RFFs.

Eventually Yang et al. [191] introduced the algorithm “À la Carte”, based on “Fastfood” which is able to learn the spectral distribution corresponding to a kernel rather than defining it from the kernel. Very recently Kawaguchi, Xie, and Song [88] proposed to use semi-random features which are a tradeoff between the random features based on kernel methods (e.g. RFFs) and the trainable layer in deep learning.



3

BACKGROUND

In this chapter we introduce briefly the mathematical tools used throughout this manuscript. We give a full table of notations, and present elements of functional analysis [93] and abstract harmonic analysis [65]. Then we turn our attention to the case when the functions we want to learn are not real-valued, but vector-valued. To learn vector-valued functions we define Operator-Valued Kernels that generalize the scalar-valued kernel presented in [Chapter 2](#). We conclude by giving a non-exhaustive list of Operator-Valued Kernels and in which context they have been used.

Contents

3.1	Notations	22
3.1.1	Algebraic structures	22
3.1.2	Topology and continuity	22
3.1.3	Measure theory	23
3.1.4	Vector spaces, linear operators and matrices	24
3.2	Elements of abstract harmonic analysis	26
3.2.1	Locally compact Abelian groups	26
3.2.2	The Haar measure	29
3.2.3	Even and odd functions	30
3.2.4	Characters	30
3.2.5	The Fourier Transform	33
3.2.6	Representations of Groups	35
3.3	On Operator-Valued Kernels	36
3.3.1	Definitions and properties	36
3.3.2	Shift-Invariant OVK on LCA groups	43
3.3.3	Examples of Operator-Valued Kernels . . .	44
3.3.4	Some use of Operator-valued kernels . . .	48

3.1 NOTATIONS

In this section we summarize briefly important notions used throughout this document. It is mainly based on books and lecture notes of Cotaescu [49] and Kurdila and Zabarankin [93].

3.1.1 Algebraic structures

¹ Commutative.

We note \mathbb{K} any Abelian¹ field and call its elements scalars. \mathbb{R} is the Abelian field of real numbers and \mathbb{C} is the Abelian field of complex numbers. The unit pure imaginary number $\sqrt{-1} \in \mathbb{C}$ is denoted i and the Euler constant $\exp(1) \in \mathbb{R}$ is denoted e . \mathbb{N} represents the set of natural numbers and \mathbb{N}_n , $n \in \mathbb{N}$ the set of natural numbers smaller or equal to n . For any space \mathcal{S} , \mathcal{S}^d , $d \in \mathbb{N}$ represents the Cartesian product space $\mathcal{S}^d = \mathcal{S} \times \dots \times \mathcal{S}$. For any two algebraic structures \mathcal{S} and \mathcal{S}' we write $\mathcal{S} \cong \mathcal{S}'$ if there exist an isomorphism between these two structures. If $a + ib = x \in \mathbb{C}$ then $\bar{x} = a - ib \in \mathbb{C}$ denotes the complex conjugate. By extension if $x \in \mathbb{R}$, $\bar{x} = x \in \mathbb{R}$.

3.1.2 Topology and continuity

In order to define a proper notion of continuity, we focus on topological spaces. A topological space is a pair of sets $(\mathcal{X}, \mathcal{T}_x)$ where \mathcal{X} describes the points considered, and \mathcal{T}_x describes the possible neighbourhoods. The standard axioms of topology suppose that $\mathcal{T}_x \subseteq \mathcal{P}(\mathcal{X})$ is a collection of subsets of \mathcal{X} such that the empty set and \mathcal{X} itself belongs to \mathcal{T}_x , any (finite or infinite) union of members of \mathcal{T}_x still belongs to \mathcal{T}_x and the intersection of any finite number of members of \mathcal{T}_x still belongs to \mathcal{T}_x . The elements of \mathcal{T}_x are called open sets and the collection \mathcal{T}_x is a topology on \mathcal{X} . If $(\mathcal{X}, \mathcal{T}_x)$ and $(\mathcal{Y}, \mathcal{T}_y)$ are topological spaces, a function f is said to be continuous if for every open set $\mathcal{V} \in \mathcal{T}_y$, the inverse image $f^{-1}(\mathcal{V}) = \{x \in \mathcal{X} \mid f(x) \in \mathcal{V}\}$ is an open subset of \mathcal{T}_x . Since the notion of continuity depends on open sets, it depends on the topology of the spaces \mathcal{X} and \mathcal{Y} .

If $(\mathcal{X}, \mathcal{T}_x)$ is a topological space and x is a point in \mathcal{X} , a neighbourhood of x is a subset \mathcal{V} of \mathcal{X} that includes an *open* set \mathcal{U} containing x . A topological space \mathcal{X} is said to be Hausdorff (T2) when all distinct points in \mathcal{X} are pairwise neighbourhood-separable. i.e. if there exists a neighbourhood \mathcal{U} of x and a neighbourhood \mathcal{V} of y such that \mathcal{U} and \mathcal{V} are disjoint. It implies the uniqueness of limits of sequences and existence of nets used throughout this thesis. Therefore in the whole document we always assume that a topological space \mathcal{X} is Haussdorff.

A topological space is said to be second countable if it has a countable base. Every second-countable space is separable and Lindelöf²

(The reverse implications do not hold). A space is metrisable if and only if it is second countable.

² Every open cover has a countable subcover.

A topological space is said to be separable if there exists a sequence $(x_n)_{n \in \mathbb{N}^*}$ of elements of \mathcal{X} such that every nonempty open subsets of the space contains at least one element of the sequence. Separability plays an important role in numerical analysis because many theorems have only constructive proofs for separable spaces. Such constructive proofs can be turned into algorithms which is the primary goal of this work. In this document we also assume that any topological space is separable if there is no specific mention of the contrary. Moreover we recall that a Hilbert space is separable if and only if it has a countable orthonormal basis (Hence separable Hilbert spaces are second countable). Hence an operator between two separable Hilbert spaces can be written as an infinite dimensional matrix. In some cases we also introduce *Polish spaces* which are separable topological spaces \mathcal{X} that have at least one metric d such that (\mathcal{X}, d) is complete. Then d induces the topology \mathcal{T}_d of \mathcal{X} . As metrisable spaces, Polish spaces are always second countable. Moreover every second countable locally compact Hausdorff space is a Polish space and every separable Banach space is a Polish space.

If \mathcal{X} and \mathcal{Y} are two topological spaces, we denote by $\mathcal{F}(\mathcal{X}; \mathcal{Y})$ the topological 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, endowed with the product topology (topology of pointwise convergence).

3.1.3 Measure theory

A σ -algebra on \mathcal{X} is a set $\mathcal{M} \subseteq \mathcal{P}(\mathcal{X})$ of subsets of \mathcal{X} , containing the empty set, which is closed under taking complements and countable unions. A pair $(\mathcal{X}, \mathcal{M})$ where \mathcal{X} is a set and \mathcal{M} is a σ -algebra is called a measure space. The Borel σ -algebra $\mathcal{B}(\mathcal{X})$ is a σ -algebra generated by the open sets of \mathcal{X} . A measure on a measurable space $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ is a map $\mu : \mathcal{B}(\mathcal{X}) \rightarrow \overline{\mathbb{R}}_+$ which is zero on the empty set and countably additive, i. e. for any subset $(Z_n)_{n \in \mathbb{N}}$ is a sequence of pairwise disjoint measurable sets,

$$\mu \left(\bigcup_{n \in \mathbb{N}} Z_n \right) = \sum_{n \in \mathbb{N}} \mu(Z_n).$$

We note $\mathcal{N}(m, \sigma)$ the Gaussian distribution with mean $m \in \mathbb{R}$ and variance $\sigma^2 \in \mathbb{R}$. $\mathcal{U}(a, b)$ is the uniform distribution with support (a, b) and $\mathcal{S}(m, \sigma)$ is the hyperbolic secant distribution with mean m and variance σ^2 .

3.1.4 Vector spaces, linear operators and matrices

Given any vector space \mathcal{H} over an Abelian field \mathbb{K} , the (continuous) dual space¹ \mathcal{H}^* is defined as the set of all *continuous* linear functionals $x^* : \mathcal{H} \rightarrow \mathbb{K}$. When \mathcal{H} is a vector space, there is a natural duality pairing between \mathcal{H}^* and \mathcal{H} defined for all $x^* \in \mathcal{H}^*$ and all $z \in \mathcal{H}$ as $(x^*, z)_{\mathcal{H}^*, \mathcal{H}} = x^*(z) = x^*z$. The duality paring $(\cdot, \cdot)_{\mathcal{H}^*, \mathcal{H}}$ is then a bilinear form.

Let \mathcal{H}_1 and \mathcal{H}_2 be two vector spaces. We call operator any linear function from \mathcal{H}_1 to \mathcal{H}_2 . The transpose (or dual) of an operator $W : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ is defined as $W^T : \mathcal{H}_2^* \rightarrow \mathcal{H}_1^*$ such that $W^T : x^* \mapsto x^*(W)$. It is characterized by the relation $(x^*, Wz)_{\mathcal{H}_2^*, \mathcal{H}_1} = (W^T x^*, z)_{\mathcal{H}_1^*, \mathcal{H}_1}$ for all $x^* \in \mathcal{H}_2^*$ and all $z \in \mathcal{H}_1$. An operator is called self-dual when $W^T = W$.

Let \mathcal{H}_1 and \mathcal{H}_2 be two vector space. We set $\mathcal{L}(\mathcal{H}_1; \mathcal{H}_2)$ to be the space of *bounded* (linear) operators from \mathcal{H}_1 to \mathcal{H}_2 . The vector space \mathcal{H}_1 is called the domain, noted Dom and \mathcal{H}_2 the codomain. We use the shortcut notation $\mathcal{L}(\mathcal{H}) = \mathcal{L}(\mathcal{H}; \mathcal{H})$. Interestingly if \mathcal{H}_1 and \mathcal{H}_2 are normed vector spaces, they can be viewed as topological vector spaces, and the notion of continuity coincides with that of boundedness. We recall that the norm of a linear operator is given by

$$\|W\|_{\mathcal{H}_1, \mathcal{H}_2} = \sup_{x \neq 0} \frac{\|Wx\|_{\mathcal{H}_2}}{\|x\|_{\mathcal{H}_1}}.$$

If $W \in \mathcal{L}(\mathcal{H}_1, \mathcal{H}_2)$

$$\text{Ker } W = \{x \in \text{Dom}(W) \mid Wx = 0\}$$

denotes the kernel (nullspace), which is a vector subspace of the domain and

$$\text{Im } W = \{y \in \mathcal{H}_2 \mid y = Wx, x \in \text{Dom}(W)\}$$

the image (range) which is a vector subspace of the codomain \mathcal{H}_2 .

If \mathcal{H} is an Hilbert space on a field \mathbb{K} we denote its scalar product by $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ and its norm by $\|\cdot\|_{\mathcal{H}}$. When the base field of \mathcal{H} is \mathbb{R} , $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is a *bilinear* form. When the base field of \mathcal{H} is \mathbb{C} , $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is a *sesquilinear* form.

¹ The continuous dual space is also called topological dual space. This must be differentiate from the *algebraic* dual space, which is the space of linear functionals from the original vector-space to its base field. Hence the continuous dual space is a subset of the algebraic dual space. The continuous and the algebraic dual space only match when considering finite dimensional vector-spaces

Let \mathcal{H} be a Hilbert space. From Riesz's representation theorem, there is a unique isometric isomorphism $\iota_R : \mathcal{H} \rightarrow \mathcal{H}^*$ such that for any x and $y \in \mathcal{H}$, $(\iota_R(x), y)_{\mathcal{H}^*, \mathcal{H}} = \langle x, y \rangle_{\mathcal{H}}$ and $\|\iota_R(x)\|_{\mathcal{H}^*} = \|x\|_{\mathcal{H}}$. The Riesz map ι_R is self-dual, thus if \mathcal{H} is a Hilbert space, \mathcal{H} is reflexive. i.e. $\mathcal{H}^{**} \cong \mathcal{H}$. When the base field of \mathcal{H} is \mathbb{C} , then the Riesz map ι_R is an *anti-linear* form since $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is sesquilinear and $(\cdot, \cdot)_{\mathcal{H}^*, \mathcal{H}}$ is bilinear. In the same way when the base field of \mathcal{H} is \mathbb{R} then ι_R is *linear* since both $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ and $(\cdot, \cdot)_{\mathcal{H}^*, \mathcal{H}}$ are bilinear. If \mathcal{H} is a Hilbert space we make the dual space \mathcal{H}^* a Hilbert space by endowing it with the inner product $\langle x^*, z^* \rangle_{\mathcal{H}^*} = \langle \iota_R^{-1}(x^*), \iota_R^{-1}(z^*) \rangle_{\mathcal{H}}$ for all $x^*, z^* \in \mathcal{H}^*$.

Let \mathcal{H}_1 and \mathcal{H}_2 be two Hilbert spaces. The adjoint of an operator $W : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ is the unique mapping $W^* : \mathcal{H}_2 \rightarrow \mathcal{H}_1$ such that $\langle W^*x, z \rangle_{\mathcal{H}_1} = \langle x, Wz \rangle_{\mathcal{H}_2}$ for all $x \in \text{Dom}(W^*)$, $z \in \text{Dom}(W)$. Its existence is guaranteed by Riesz's representation theorem. An operator $W : \text{Dom}(W) \subseteq \mathcal{H} \rightarrow \mathcal{H}$ is said to be symmetric when $W^* = W$, and self-adjoint when W is bounded, symmetric, $\text{Dom}(W^*) = \text{Dom}(W)$ and $\text{Dom}(W)$ is dense in \mathcal{H} . If W is bounded, symmetric and $\text{Dom}(W) = \mathcal{H}$ then W is self-adjoint. Notice that the transpose is linked to the adjoint by the relation $W^* = \iota_R^{-1}W^\top\iota_R$. When \mathcal{H} is a Hilbert space, if $x \in \mathcal{H}$, we always define $x^* \in \mathcal{H}^*$ to be

$$x^* = \iota_R(x) = \langle x, \cdot \rangle_{\mathcal{H}}.$$

$$\begin{array}{ccc} \mathcal{H}_2 & \xrightarrow{W^*} & \mathcal{H}_1 \\ \downarrow \iota_R & & \iota_R^{-1} \swarrow \quad \searrow \iota_R \\ \mathcal{H}_2^* & \xrightarrow{W^\top} & \mathcal{H}_1^* \end{array}$$

Figure 3.1: Riesz map, dual spaces and adjoints.

Let \mathcal{H} be a *separable* Hilbert space and let $(e_i)_{i \in \mathbb{N}^*}$ be a basis of \mathcal{H} . We call $(e_i^*)_{i \in \mathbb{N}^*}$ the dual basis of \mathcal{H} , the basis of \mathcal{H}^* such that for all $i, j \in \mathbb{N}^*$, $e_i^*(e_j) = \langle e_i, e_j \rangle_{\mathcal{H}} = \delta_{ij}$. In the whole document we consider that \mathcal{H}^* is always equipped with the dual basis of \mathcal{H} . For a vector $x \in \mathcal{H}$ with a basis $(e_i)_{i \in \mathbb{N}^*}$ we write $x_i = e_i^*(x)$. For a linear operator $W : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ where \mathcal{H}_1 and \mathcal{H}_2 are Hilbert spaces with respective basis $(e_i)_{i \in \mathbb{N}^*}$ and $(e'_j)_{j \in \mathbb{N}^*}$, we note $W_i = We_i$ and $W_{ij} = e_j^*(We_i)$. Eventually given two separable Hilbert spaces \mathcal{H}_1 and \mathcal{H}_2 , an operator $W : \mathcal{H}_1 \rightarrow \mathcal{H}_2$, $(e_i)_{i \in \mathbb{N}^*}$ a basis of \mathcal{H}_1 and $(e'_i)_{i \in \mathbb{N}^*}$ a basis of \mathcal{H}_2 we have

$$(W^\top)_{ij} = e_j^* W^\top e'_i = e_j^* e'^*_i W = e'^*_i We_j = W_{ji}.$$

We call matrix M of size $(m, n) \in \mathbb{N}^2$ on an Abelian field \mathbb{K} a collection of elements $M = (m_{ij})_{1 \leq i \leq m, 1 \leq j \leq n}$, $m_{ij} \in \mathbb{K}$. We note $\mathcal{M}_{m,n}(\mathbb{K})$ the vector space of all matrices. If \mathcal{H}_1 and \mathcal{H}_2 are two separable Hilbert spaces on an Abelian field \mathbb{K} , any linear operator $L \in \mathcal{L}(\mathcal{H}_1; \mathcal{H}_2)$ can be viewed as a (potentially infinite) matrix. Let $n = \dim(\mathcal{H}_1)$, $m = \dim(\mathcal{H}_2)$ and let $B = (e_i)_{i=1}^n$ and $C = (e'_j)_{j=1}^m$ be the respective bases of \mathcal{H}_1 and \mathcal{H}_2 . We note $\text{mat}_{B,C} : \mathcal{L}(\mathcal{H}_1; \mathcal{H}_2) \rightarrow \mathcal{M}_{m,n}(\mathbb{K})$ such that $M = \text{mat}_{B,C}(L) = (e_j'^* L e_i)_{1 \leq i \leq n, 1 \leq j \leq m} \in \mathcal{M}_{m,n}(\mathbb{K})$. Let $M_1 \in \mathcal{M}_{m,n}(\mathbb{K})$ and $M_2 \in \mathcal{M}_{n,l}(\mathbb{K})$. The product between two matrices is written $M_1 M_2 \in \mathcal{M}_{m,l}(\mathbb{K})$ and obey $(M_1 M_2)_{ij} = \sum_{k=1}^n M_{ik} M_{kj}$. Given two linear operator $L_1 \in \mathcal{L}(\mathcal{H}_1; \mathcal{H}_2)$ and $L_2 \in \mathcal{L}(\mathcal{H}_2; \mathcal{H}_3)$ we have $L_1 L_2 \in \mathcal{L}(\mathcal{H}_1; \mathcal{H}_3)$ and i

$$\text{mat}_{B,D}(L_1 L_2) = \text{mat}_{B,C}(L_1) \text{mat}_{C,D}(L_2).$$

The operator $\text{mat}_{B,C}$ is a vector space isomorphism allowing us to identify $\mathcal{L}(\mathcal{H}_1; \mathcal{H}_2)$ with $\mathcal{M}_{mn}(\mathbb{K})$ where $n = \dim(\mathcal{H}_1)$ and $m = \dim(\mathcal{H}_2)$. All these notations are summarized in [Tables 3.1](#) and [3.3](#).

3.2 ELEMENTS OF ABSTRACT HARMONIC ANALYSIS

3.2.1 Locally compact Abelian groups

Definition 3.1 (Locally Compact Abelian (LCA) group.). A group X endowed with a binary operation \star is said to be a Locally Compact Abelian group if X is a topological commutative group w.r.t. \star for which every point has a compact neighborhood and is Hausdorff (T2).

Moreover given a element z of a LCA group X , we define the set $z \star X = X \star z = \{z \star x \mid \forall x \in X\}$ and the set $X^{-1} = \{x^{-1} \mid \forall x \in X\}$. We also note e the neutral element of X such that $x \star e = e \star x = e$ for all $x \in X$. Throughout this thesis we focus on positive-definite functions. Let \mathcal{Y} be a complex separable Hilbert space. A function $f : X \rightarrow \mathcal{Y}$ is positive definite if for all $N \in \mathbb{N}$ and all $y \in \mathcal{Y}$,

$$\sum_{i,j=1}^N \langle y_i, f(x_j^{-1} \star x_i) y_j \rangle_{\mathcal{Y}} \geq 0 \quad (3.1)$$

for all sequences $(y_i)_{i \in \mathbb{N}_N^*} \in \mathcal{Y}^N$ and all sequences $(x_i)_{i \in \mathbb{N}_N^*} \in X^N$. If \mathcal{Y} is real we add the assumption that $f(x^{-1}) = f(x)^*$ for all $x \in X$. A consequence is that a positive-definite function is bounded, as shown by Falb [60], $\|f(x)\|_{\mathcal{Y}, \mathcal{Y}} \leq 2\|f(e)\|_{\mathcal{Y}, \mathcal{Y}}$ for all $x \in X$, however positive-definite functions are not necessarily continuous. This motivates the introduction of functions of positive type which are nothing but continuous positive-definite function.

Table 3.1: Mathematical symbols and their signification (part 1).

Symbol	Meaning
$:=$	Equal by definition.
\mathbb{N}	The semi-group of natural numbers.
\mathbb{K}	Any non-discrete Abelian field endowed with an absolute value. Elements of \mathbb{K} are called scalars.
\mathbb{R}	The Abelian field of real numbers.
\mathbb{C}	The Abelian field of complex numbers.
\mathbb{U}	The circle group of complex numbers with unit module.
$i \in \mathbb{C}$	Unit pure imaginary number $i^2 := -1$.
$e \in \mathbb{R}$	Euler constant.
$e \in \mathcal{X}$	The neutral element of the group \mathcal{X} .
δ_{ij}	Kronecker delta function. $\delta_{ij} = 0$ if $i \neq j$, 1 otherwise.
$\langle \cdot, \cdot \rangle_2$	Euclidean inner product.
$\ \cdot\ _2$	Euclidean norm.
\mathcal{X}	Input space.
$\hat{\mathcal{X}}$	The Pontryagin dual of \mathcal{X} when \mathcal{X} is a LCA group.
\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} .
$\mathcal{F}(\mathcal{X}; \mathcal{Y})$	Topological vector space of functions from \mathcal{X} to \mathcal{Y} .
$\mathcal{C}(\mathcal{X}; \mathcal{Y})$	The topological vector subspace of \mathcal{F} of continuous functions 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} .
$\ \cdot\ _{\mathcal{Y}, \mathcal{Y}'}$	The operator norm $\ \Gamma\ _{\mathcal{Y}, \mathcal{Y}'} = \sup_{\ \mathbf{y}\ _{\mathcal{Y}}=1} \ \Gamma \mathbf{y}\ _{\mathcal{Y}'}$ for all $\Gamma \in \mathcal{L}(\mathcal{Y}, \mathcal{Y}')$
$\mathcal{M}_{m,n}(\mathbb{K})$	The set of matrices of size (m, n) .
$\mathcal{L}(\mathcal{Y})$	The set of bounded linear operator from a Hilbert space \mathcal{Y} 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 a topological space \mathcal{X} .
$\mu(\mathcal{X})$	A scalar positive measure of \mathcal{X} .
$\text{Leb}(\mathcal{X})$	The Lebesgue measure of \mathcal{X} .
$\text{Haar}(\mathcal{X})$	A Haar measure of \mathcal{X} .

Table 3.3: Mathematical symbols and their signification (part 2).

Symbol	Meaning
$\Pr_{\mu, \rho}(\mathcal{X})$	A probability measure of \mathcal{X} whose Radon-Nikodym derivative (density) with respect to the measure μ is ρ .
$\mathcal{F}[\cdot]$	The Fourier Transform operator.
$\mathcal{F}^{-1}[\cdot]$	The Inverse Fourier Transform operator.
ess sup	The essential supremum.
$L^p(\mathcal{X}, \mu)$	The Banach space of $ \cdot ^p$ -integrable function from $(\mathcal{X}, \mathcal{B}(\mathcal{X}), \mu)$ to \mathbb{C} for $p \in \mathbb{R}_+$.
$L^p(\mathcal{X}, \mu; \mathcal{Y})$	The Banach space of $\ \cdot\ _{\mathcal{Y}}^p$ (Bochner)-integrable function from $(\mathcal{X}, \mathcal{B}(\mathcal{X}), \mu)$ to \mathcal{Y} for $p \in \mathbb{R}_+$. $L^p(\mathcal{X}, \mu, \mathbb{R}) := L^p(\mathcal{X}, \mu)$.
$\bigoplus_{j=1}^D x_i$	The direct sum of $D \in \mathbb{N}$ vectors x_i 's in the Hilbert spaces \mathcal{H}_i . By definition $\langle \bigoplus_{j=1}^D x_j, \bigoplus_{j=1}^D z_j \rangle = \sum_{j=1}^D \langle x_j, z_j \rangle_{\mathcal{H}_i}$ [11].
$\ \cdot\ _p$	The $L^p(\mathcal{X}, \mu, \mathcal{Y})$ norm. $\ f\ _p^p := \int_{\mathcal{X}} \ f(x)\ _{\mathcal{Y}}^p d\mu(x)$. When $\mathcal{X} = \mathbb{N}^*$, $\mathcal{Y} \subseteq \mathbb{R}$ and μ is the counting measure and $p = 2$ it coincide with the Euclidean norm $\ \cdot\ _2$ for finite dimensional vectors.
$\ \cdot\ _\infty$	The uniform norm $\ f\ _\infty = \text{ess sup} \{ \ f(x)\ _{\mathcal{Y}} \mid x \in \mathcal{X} \} = \lim_{p \rightarrow \infty} \ f\ _p$.
T^\top	The transpose operator of a linear operator.
$*$	The adjoint operator of a linear operator.
$ \Gamma $	The absolute value of the linear operator $\Gamma \in \mathcal{L}(\mathcal{Y})$, i.e. $ \Gamma ^2 = \Gamma^* \Gamma$.
$\text{Tr}[\Gamma]$	The trace of a linear operator $\Gamma \in \mathcal{L}(\mathcal{Y})$.
$\sigma(\Gamma)$	The spectrum of the bounded linear operator $\Gamma \in \mathcal{L}(\mathcal{Y})$ where \mathcal{Y} is a Hilbert space, i.e. $\sigma(\Gamma) = \{\lambda \in \mathbb{C} \mid \nexists s, s(\lambda e - \Gamma) = e\}$.
$\lambda_i(\Gamma)$	The i -th eigenvalue of $\Gamma \in \mathcal{L}(\mathcal{Y})$, ranked by increasing modulus, where \mathcal{Y} is a <i>separable</i> Hilbert space and $i \in \mathbb{N}^*$.
$\rho(\Gamma)$	The spectral radius of the linear operator Γ i.e. $\rho(\Gamma) = \sup \{ \lambda \mid \lambda \in \sigma(\Gamma) \}$.
$\ \cdot\ _{\sigma, p}$	The Schatten p -norm, $\ \Gamma\ _{\sigma, p}^p = \text{Tr}[\ \Gamma\ ^p]$ for $\Gamma \in \mathcal{L}(\mathcal{Y})$, where \mathcal{Y} is a Hilbert space. Note that $\ \Gamma\ _{\sigma, \infty} = \rho(\Gamma) \leq \ \Gamma\ _{\mathcal{Y}, \mathcal{Y}}$.
\succcurlyeq	“Greater than” in the Loewner partial order of operators. $\Gamma_1 \succcurlyeq \Gamma_2$ if $\sigma(\Gamma_1 - \Gamma_2) \subseteq \mathbb{R}_+$.
$\overline{\mathbb{R}}$	The one point compactification of the real line $\mathbb{R} \cup \{\infty\}$.
\cong	Given two sets \mathcal{X} and \mathcal{Y} , $\mathcal{X} \cong \mathcal{Y}$ if there exists an isomorphism $\varphi : \mathcal{X} \rightarrow \mathcal{Y}$.

3.2.2 The Haar measure

Measures on topological spaces which appear in practice often satisfy the following regularity properties.

Definition 3.2 (Radon measure). A Radon measure $\mu = \mathbf{Rad}$ on a topological measurable space \mathcal{X} is a measure on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ which satisfies the following properties.

1. The measure **Rad** is finite on every compact set.

$$\mathbf{Rad}(K) < \infty, \text{ for any compact set } K \in \mathcal{B}(\mathcal{X}).$$

2. The measure **Rad** is outer regular on any Borel sets E .

$$\mathbf{Rad}(E) = \inf \{ \mathbf{Rad}(U) \mid E \subseteq U \}, \text{ for any open set } U.$$

3. The measure **Rad** is inner regular on open sets E .

$$\mathbf{Rad}(E) = \sup \{ \mathbf{Rad}(K) \mid K \subseteq E \}, \text{ for any compact set } K.$$

When dealing with topological groups it is natural to look for measures which are invariant under translation. There exists, up to a positive multiplicative constant, a unique countably additive, nontrivial measure **Haar** on any LCA group. For more details and constructive proofs see Alfsen [4], Conway [46], and Folland [65].

Definition 3.3 (The Haar measure). A Haar measure $\mu = \mathbf{Haar}$ on a LCA group $\mathcal{X} = (G, \star)$ is a Radon measure on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ which is non-zero on non-empty open sets and is invariant under translation. Namely

1. if $Z \subseteq \mathcal{X}$ is open and not empty, then $\mathbf{Haar}(Z) > 0$.
2. For all $Z \in \mathcal{B}(\mathcal{X})$ and $x \in \mathcal{X}$, $\mathbf{Haar}(x \star Z) = \mathbf{Haar}(Z)$.

Such a measure on a LCA group \mathcal{X} is called a Haar measure³. An immediate consequence of the invariance is that for any $s \in \mathcal{X}$,

$$\int_{\mathcal{X}} f(s \star x) d\mathbf{Haar}(x) = \int_{\mathcal{X}} f(x) d\mathbf{Haar}(x).$$

It can be shown that $\mathbf{Haar}(U) > 0$ for every non-empty open subset U . In particular, if \mathcal{X} is compact then $\mathbf{Haar}(\mathcal{X})$ is finite and positive, so we can uniquely specify a Haar measure on \mathcal{X} by adding the normalization condition $\mathbf{Haar}(\mathcal{X}) = 1$. We call measured space the space $(\mathcal{X}, \mathcal{B}(\mathcal{X}), \mathbf{Haar})$. In other words the (topological) space \mathcal{X} endowed with its Borel σ -algebra $\mathcal{B}(\mathcal{X})$ and a measure **Haar**. If $\mathbf{Haar}(\mathcal{X}) = 1$ then the space $(\mathcal{X}, \mathcal{B}(\mathcal{X}), \mathbf{Haar})$ is called a probability space. Last but not least, on the additive group $(\mathbb{R}, +)$, the Lebesgue measure noted **Leb** is a valid Haar measure. For a concise introduction and important properties we refer the reader to the lecture of Tornier [173].

³ 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

3.2.3 Even and odd functions

Let \mathcal{X} be a LCA group and \mathbb{K} be a field viewed as an additive group. We say that a function $f : \mathcal{X} \rightarrow \mathbb{K}$ 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 3.4 (Even and odd operator-valued function on a LCA group). Let \mathcal{X} be a measured LCA group and \mathcal{Y} be a Hilbert space, and $\mathcal{L}(\mathcal{Y})$ the space of bounded linear operators from \mathcal{Y} to itself viewed as an additive group. A function $f : \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ is (weakly) even if for all $x \in \mathcal{X}$ and all $y, y' \in \mathcal{Y}$,

$$\langle y, f(x^{-1})y' \rangle_{\mathcal{Y}} = \langle y, f(x)y' \rangle_{\mathcal{Y}} \quad (3.2)$$

and (weakly) odd if

$$\langle y, f(x^{-1})y' \rangle_{\mathcal{Y}} = -\langle y, f(x)y' \rangle_{\mathcal{Y}} \quad (3.3)$$

It is easy to check that if f is odd then $\int_{\mathcal{X}} \langle y, f(x)y' \rangle_{\mathcal{Y}} d\text{Haar}(x) = 0$.

Proof

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

□

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 $h(x) = \langle y, f(x)g(x)y' \rangle$. Then we have

$$\begin{aligned} h(x^{-1}) &= \langle y, f(x^{-1})g(x^{-1})y' \rangle_{\mathcal{Y}} = \langle y, f(x)(-g(x))y' \rangle_{\mathcal{Y}} \\ &= -h(x). \end{aligned} \quad (3.4)$$

3.2.4 Characters

Locally Compact Abelian (LCA) groups are central to the general definition of Fourier Transform which is related to the concept of Pontryagin duality [65]. 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 $\omega : \mathcal{X} \rightarrow \mathbb{U}$ 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 $\widehat{\mathcal{X}}$. The dual group of an LCA group is an LCA group such that we can endow $\widehat{\mathcal{X}}$ with a “dual” Haar measure noted $\widehat{\text{Haar}}$. Then 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 $\widehat{\mathcal{X}} \cong \mathcal{X}$. 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

$$\omega(x) = (x, \omega) = (\omega, x) = x(\omega), \quad (3.5)$$

where $x \in \mathcal{X} \cong \widehat{\mathcal{X}}$ and $\omega \in \widehat{\mathcal{X}}$. The form (\cdot, \cdot) defined in [Equation 3.5](#) is called (duality) pairing. Another important property involves the complex conjugate of the pairing which is defined as

$$\overline{(x, \omega)} = (x^{-1}, \omega) = (x, \omega^{-1}). \quad (3.6)$$

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,

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

The following example shows how to determine the (Pontryagin) dual of a LCA group.

Example 3.1 (Folland [65]). *On the additive group $\mathcal{X} = (\mathbb{R}, +)$ we have $\widehat{\mathbb{R}} \cong \mathbb{R}$ with the duality pairing $(x, \omega) = \exp(ix\omega)$ for all $x \in \mathbb{R}$ and all $\omega \in \mathbb{R}$. The Haar measure on \mathcal{X} is the Lebesgue measure.*

Proof If $\omega \in \widehat{\mathbb{R}}$ then $\omega(0) = 1$ since ω is an homeomorphism from \mathbb{R} to \mathbb{U} . Therefore there exists $a > 0$ such that $\int_0^a \omega(t)d\text{Leb}(t) \neq 0$. Setting $A\omega = \int_0^a \omega(t)d\text{Leb}(t)$ we have

$$(A\omega)(x) = \int_0^a \omega(x+t)d\text{Leb}(t) = \int_x^{a+x} \omega(t)d\text{Leb}(t).$$

so ω is differentiable and

$$\omega'(x) = A^{-1}(\omega(a+x) - \omega(x)) = c\omega(x) \quad \text{where} \quad c = A^{-1}(\omega(a) - 1).$$

It follow that $\omega(x) = e^{cx}$, and since $|\omega| = 1$, one can take $c = i\xi$ for some $\xi \in \mathbb{R}$. Hence we can identify ω with ξ and $\widehat{\mathbb{R}}$ with \mathbb{R} since ξ uniquely determines ω , thus we identify $\omega = \xi$. \square

We also especially mention the duality pairing associated to the skewed multiplicative LCA product group. This group together with the operation \odot has been proposed by Li, Ionescu, and Sminchisescu [98] to handle histograms features especially useful in image recognition applications. Let $\mathcal{X} = (-c_k; +\infty)_{k=1}^d$, where $c_k \in \mathbb{R}_+$, endowed with the group operation \odot defined component-wise for all $x, z \in \mathcal{X}$ as follow.

$$x \odot z := ((x_k + c_k)(z_k + c_k) - c_k)_{k=1}^d.$$

Example 3.2 (Li, Ionescu, and Sminchisescu [98]). On the skewed multiplicative group $\mathcal{X} = ((-c, +\infty), \odot)$ we have $\widehat{(-, +\infty)} \cong \mathbb{R}$, with duality pairing $(x, \omega) = \exp(i \log(x+c)\omega)$ for all $x \in \mathcal{X}$ and all $\omega \in \widehat{\mathcal{X}}$. The Haar measure on \mathcal{X} is given for all $\mathcal{Z} \in \mathcal{B}(\mathcal{X})$ by $\text{Haar}(\mathcal{Z}) = \int_{\mathcal{Z}} (z+c)^{-1} d\text{Leb}(z)$.

Proof Let $a, b \in (-c, +\infty)$ and $\mu([a, b]) = \int_a^b (z+c)^{-1} d\text{Leb}(z)$. Then for all $d \in (-c, +\infty)$

$$\begin{aligned} \mu([d \odot a, d \odot b]) &= \int_{(d+c)(a+c)-c}^{(d+c)(b+c)-c} (z+c)^{-1} d\text{Leb}(z) \\ &= \log(d+c)(b+c) - \log(d+c)(a+c) \\ &= \log(b+c) - \log(a+c) \\ &= \int_a^b (z+c)^{-1} d\text{Leb}(z) = \mu([a, b]). \end{aligned}$$

Thus μ is translation invariant, making $\text{Haar} = \lambda\mu$ a valid Haar measure on \mathcal{X} for any multiplicative constant $\lambda \in \mathbb{R}_*$. Let $(x, \omega) = \exp(i \log(x+c)\omega)$ for all $x \in \mathcal{X}$ and all $\omega \in \widehat{\mathcal{X}}$. We have for all $z \in \mathcal{X}$

$$\begin{aligned} (x \odot z, \omega) &= \exp(i \log((x+c)(z+c))\omega) \\ &= \exp(i \log(x+c)\omega) \exp(i \log(z+c)\omega) \\ &= (x, \omega)(z, \omega) \end{aligned}$$

Thus $\omega(x \odot z) := \omega(x)\omega(z)$, which defines a valid pairing, therefore we can identify $\widehat{\mathcal{X}} = \widehat{(-c, +\infty)} \cong \mathbb{R}$ where \mathbb{R} is the additive group endowed with the Haar measure being the Lebesgue measure. \square

It is easy to extend the Pontryagin dual of groups to dual groups, as well as defining the pairing on the dual group using the following proposition [65]

Proposition 3.1 (Folland [65]). Let $(\mathcal{X}_i)_{i \in \mathbb{N}}$ be a collection of LCA groups. Then

$$\left(\widehat{\prod_{i \in \mathbb{N}} \mathcal{X}_i} \right) \cong \prod_{i \in \mathbb{N}} \widehat{\mathcal{X}_i}$$

Proof Each $\omega = (\omega_1, \dots, \omega_N) \in \prod_{i=1}^N \mathcal{X}_i$ defines a character on $\prod_{i=1}^N \mathcal{X}_i$ by

$$((x_1, \dots, x_N), (\omega_1, \dots, \omega_N)) = (x_1, \omega_1) \cdots (x_N, \omega_N).$$

Moreover, every character ω on $\prod_{i=1}^N \mathcal{X}_i$ is of this form, where ω_i is defined by

$$(x_i, \omega_i) = ((e_1, \dots, e_{i-1}, x_j, e_{i+1}, \dots, e_N), \omega),$$

where e_i 's denotes the neutral elements of the LCA group \mathcal{X}_i . \square

Hence $\widehat{\mathbb{R}^d} \cong \mathbb{R}^d$ with duality pairing

$$(x, \omega) = \exp \left(i \sum_{k=1}^d x_k \omega_k \right),$$

hence $h_\omega(x) = \sum_{k=1}^d \omega_k x_k = \langle x, \omega \rangle_2$. For the skewed multiplicative group $(-c_k; +\infty)_{k=1}^d \cong \mathbb{R}^d$ and the duality pairing is defined by

$$(x, \omega) = \exp \left(i \sum_{k=1}^d \log(x_k + c_k) \omega_k \right).$$

Hence $h_\omega(x) = \sum_{k=1}^d \log(x_k + c_k) \omega_k = \langle \log(x + c), \omega \rangle_2$. Eventually the natural Haar measure on a product group is the product measure. e.g. for $\mathcal{X} = \mathbb{R}^d$, the Haar measure on \mathbb{R}^d is the d -th power of the Lebesgue measure on \mathbb{R} . Table 3.5 provides an explicit list of pairings for various groups based on \mathbb{R}^d or its subsets. The interested reader can refer to Folland [65] for a more detailed construction of LCA, Pontryagin duality and Fourier Transforms on LCA.

Table 3.5: Classification of Fourier Transforms in terms of their domain and transform domain.

$\mathcal{X} =$	$\widehat{\mathcal{X}} \cong$	Operation	Pairing
\mathbb{R}^d	\mathbb{R}^d	+	$(x, \omega) = \exp(i \langle x, \omega \rangle_2)$
$\mathbb{R}_{*,+}^d$	\mathbb{R}^d	.	$(x, \omega) = \exp(i \langle \log(x), \omega \rangle_2)$
$(-c; +\infty)^d$	\mathbb{R}^d	\odot	$(x, \omega) = \exp(i \langle \log(x + c), \omega \rangle_2)$

3.2.5 The Fourier Transform

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

$$\forall \omega \in \widehat{\mathcal{X}}, \quad \mathcal{F}[f](\omega) = \int_{\mathcal{X}} \overline{(x, \omega)} f(x) d\text{Haar}(x).$$

The Inverse Fourier Transform (IFT) of a function $g \in L^1(\widehat{\mathcal{X}}, \widehat{\text{Haar}}; \mathcal{Y})$ is noted $\mathcal{F}^{-1}[g]$ defined by

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

We also define the flip operator \mathcal{R} by $(\mathcal{R}f)(x) := f(x^{-1})$.

Theorem 3.1 (Fourier inversion (Folland [65])). *Given a measure **Haar** defined on \mathcal{X} , there exists a unique suitably normalized dual measure $\widehat{\text{Haar}}$ on $\widehat{\mathcal{X}}$ such that for all $f \in L^1(\mathcal{X}, \text{Haar}; \mathbb{Y})$ and if $\mathcal{F}[f] \in L^1(\widehat{\mathcal{X}}, \widehat{\text{Haar}}; \mathbb{Y})$ we have*

$$f(x) = \int_{\widehat{\mathcal{X}}} (x, \omega) \mathcal{F}[f](\omega) d\widehat{\text{Haar}}(\omega), \quad \text{for Haar-almost all } x \in \mathcal{X}. \quad (3.7)$$

i.e. such that $(\mathcal{R}\mathcal{F}\mathcal{F}[f])(x) = \mathcal{F}^{-1}\mathcal{F}[f](x) = f(x)$ for **Haar**-almost all $x \in \mathcal{X}$. If f is continuous this relation holds for all $x \in \mathcal{X}$.

Proof The proof is based on Bochner's theorem and the Pontryagin duality theorem. We refer the reader to Folland [65, theorem 4.22 page 105 and theorem 4.33 page 111] for the full proof. \square

Thus when a Haar measure **Haar** on \mathcal{X} is given, the measure on $\widehat{\mathcal{X}}$ that makes Theorem 3.1 true is called the dual measure of **Haar**, noted $\widehat{\text{Haar}}$. Let $c \in \mathbb{R}_*$. If $c\text{Haar}$ is the measure on \mathcal{X} , then $c^{-1}\widehat{\text{Haar}}$ is the dual measure on $\widehat{\mathcal{X}}$. Hence one must replace $\widehat{\text{Haar}}$ by $c^{-1}\widehat{\text{Haar}}$ in the inversion formula to compensate. Therefore, we always take the Haar measure $\widehat{\text{Haar}}$ on $\widehat{\mathcal{X}}$ to be the dual of the given Haar measure **Haar** on \mathcal{X} . Whenever $\widehat{\text{Haar}} = \text{Haar}$ we say that the Haar measure is self-dual. Moreover if $\widehat{\text{Haar}}$ is normalized, the Fourier Transform on

$$L^1(\mathcal{X}, \text{Haar}; \mathbb{Y}) \cap L^2(\mathcal{X}, \text{Haar}; \mathbb{Y})$$

extends uniquely to a unitary isomorphism from $L^2(\mathcal{X}, \text{Haar}, \mathbb{Y})$ onto $L^2(\widehat{\mathcal{X}}, \widehat{\text{Haar}}; \mathbb{Y})$ (Plancherel theorem). For the familiar case of a scalar-valued function f on the LCA group $(\mathbb{R}^d, +)$, we have for all $\omega \in \widehat{\mathcal{X}} = \mathbb{R}^d$

$$\begin{aligned} \mathcal{F}[f](\omega) &= \int_{\mathcal{X}} \overline{(x, \omega)} f(x) d\text{Haar}(x) \\ &= \int_{\mathbb{R}^d} \exp(-i\langle x, \omega \rangle_2) f(x) d\text{Leb}(x), \end{aligned} \quad (3.8)$$

the Haar measure being here the Lebesgue measure. Notice that the normalization factor of $\widehat{\text{Haar}}$ on $\widehat{\mathcal{X}}$ depends on the measure **Haar** on \mathcal{X} and the duality pairing. For instance let $\mathcal{X} = (\mathbb{R}^d, +)$. In Example 3.1 we showed that $\widehat{\mathcal{X}} \cong \mathbb{R}^d$ with pairing $(x, \omega) = \exp(i\langle x, \omega \rangle_2)$, for all $x \in \mathcal{X}$ and $\omega \in \widehat{\mathcal{X}}$. If one endows \mathcal{X} with the Lebesgue measure as the Haar measure, the Haar measure on the dual is defined for all $Z \in \mathcal{B}(\mathbb{R}^d)$ by

$$\text{Haar}(Z) = \text{Leb}(Z), \quad \text{and} \quad \widehat{\text{Haar}}(Z) = \frac{1}{(2\pi)^d} \text{Leb}(Z),$$

in order to have $\mathcal{F}^{-1}\mathcal{F}[f] = f$. If one uses the cleaner equivalent pairing $(x, \omega) = \exp(2i\pi\langle x, \omega \rangle_2)$ rather than $(x, \omega) = \exp(i\langle x, \omega \rangle_2)$, then

$$\widehat{\text{Haar}}(Z) = \text{Leb}(Z).$$

The pairing $(x, \omega) = \exp(2i\pi\langle x, \omega \rangle_2)$ looks more attractive in theory since it limits the messy factor outside the integral sign and makes the Haar measure self-dual. However it is of lesser use in practice since it yields additional unnecessary computation when evaluating the pairing. Hence for symmetry reason on $(\mathbb{R}^d, +)$ and reduce computations we settle with the Haar measure on \mathbb{R}^d groups (additive and multiplicative) defined as

$$\widehat{\text{Haar}}(\mathcal{Z}) = \text{Haar}(\mathcal{Z}) = \frac{1}{\sqrt{2\pi}^d} \mathbf{Leb}(\mathcal{Z}).$$

We conclude this subsection by recalling the injectivity property of the Fourier Transform.

Corollary 3.1 (Fourier Transform injectivity (Folland [65])). *Given μ and ν two measures, if $\mathcal{F}[\mu] = \mathcal{F}[\nu]$ then $\mu = \nu$. Moreover given two functions f and $g \in L^1(\mathcal{X}, \text{Haar}; \mathcal{Y})$ if $\mathcal{F}[f] = \mathcal{F}[g]$ then $f = g$*

Proof We refer the reader to the proof of Folland [65, corollary 4.34 page 112]. \square

3.2.6 Representations of Groups

Representations of groups are convenient tools that allows group-theoretic problems to be replaced by linear algebra problems. Let $\text{Gl}(\mathcal{H})$ 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 :

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

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 $\{0\}$ and \mathcal{H} are the only two stable invariant subspaces under all operators $\pi(x)$ for all $x \in \mathcal{X}$. i.e. for all $\mathcal{U} \subset \mathcal{H}$, $\mathcal{U} \neq \{0\}$,

$$\{ \pi(x)v \mid \forall x \in \mathcal{X}, \forall v \in \mathcal{U} \} = \mathcal{U}.$$

To study LCA groups we also introduce the left regular representation of \mathcal{X} acting on a Hilbert space of function $\mathcal{H} \subset \mathcal{F}(\mathcal{X}; \mathcal{Y})$. 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 $\lambda_x = 1 \iff x = e$.

3.3 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. The fundations of the general theory of scalar kernels is mostly due to Aronszajn [9] and provides a unifying point of view for the study of an important class of Hilbert spaces of real or complex valued functions. It has been first applied in the theory of partial differential equation. The theory of Operator-Valued Kernels (OVKs) which extends the scalar-valued kernel was first developped by Pedrick [133] in his Ph. D Thesis. Since then it has been successfully applied to machine learning by many authors. In particular we introduce the notion of Operator-Valued Kernels following the propositions of Carmeli, De Vito, and Toigo [40], Carmeli et al. [41], and Micchelli and Pontil [113].

3.3.1 Definitions and properties

In machine learning the goal is often to find a function f belonging to a class (space) of functions $\mathcal{F}(\mathcal{X}; \mathcal{Y})$ that minimizes a criterion called the true risk (see [Section 2.1](#)). The class of functions we consider are functions living in a Hilbert space $\mathcal{H} \subset \mathcal{F}(\mathcal{X}; \mathcal{Y})$. The completeness allows to consider sequences of functions $f_n \in \mathcal{H}$ where the limit $f_n \rightarrow f$ is in \mathcal{H} . Moreover the existence of an inner product gives rise to a norm and also makes \mathcal{H} a metric space.

Among all these functions $f \in \mathcal{H}$, we consider a subset of functions $f \in \mathcal{H}_K \subset \mathcal{H}$ such that the evaluation map $ev_x : f \mapsto f(x)$ is bounded for all x . i. e. such that $\|ev_x\|_{\mathcal{H}_K} \leq C_x \in \mathbb{R}$ for all x . For scalar valued kernel the evaluation map is a linear functional. Thus by Riesz's representation theorem there is an isomorphism between evaluating a function at a point and an inner product: $f(x) = ev_x f = \langle K_x, f \rangle_K$. From this we deduce the reproducing property $K(x, z) = \langle K_x, K_z \rangle_K$ which is the cornerstone of many proofs in machine learning and functional analysis. When dealing with vector-valued functions, the evaluation map ev_x is no longer a linear functional, since it is vector-valued.

However, inspired by the theory of scalar valued kernel, many authors showed that if the evaluation map of functions with values in a Hilbert space \mathcal{Y} is bounded, a similar reproducing property can be obtained; namely $\langle y', K(x, z)y \rangle = \langle K_x y', K_z y \rangle_K$ for all $y, y' \in \mathcal{Y}$. This motivates the following definition of a Vector Valued Reproducing Kernel Hilbert Space (VV-RKHS).

Definition 3.5 (Vector Valued Reproducing Kernel Hilbert Space [40, 113]). Let \mathcal{Y} be a (real or complex) Hilbert space. A Vector Valued Reproducing Kernel Hilbert Space on a locally compact second countable topological space \mathcal{X} 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 $C_x \in \mathbb{R}$ such that for all $f \in \mathcal{H}$

$$\|f(x)\|_{\mathcal{Y}} \leq C_x \|f\|_{\mathcal{H}}. \quad (3.9)$$

Throughout this section we show that a VV-RKHS defines a unique positive-definite function called Operator-Valued Kernel (OVK) and conversely an OVK uniquely defines a VV-RKHS. The bijection between OVKs and VV-RKHSs has been first proved by Senkene and Tempel'man [152] in 1973. In this introduction to OVKs we follow the definitions and most recent proofs of Carmeli et al. [41].

Definition 3.6 (Positive-definite Operator-Valued Kernel acting on a complex Hilbert 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 a positive-definite Operator-Valued Kernel if

$$\sum_{i,j=1}^N \langle K(x_i, x_j)y_j, y_i \rangle_{\mathcal{Y}} \geq 0, \quad (3.10)$$

for all $N \in \mathbb{N}$, for all sequences of points $(x_i)_{i=1}^N$ in \mathcal{X}^N and all sequences of points $(y_i)_{i=1}^N$ in \mathcal{Y}^N .

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

Definition 3.7 (Positive-definite Operator-Valued Kernel 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 a positive-definite Operator-Valued Kernel kernel if

$$K(x, z) = K(z, x)^* \quad (3.11)$$

and

$$\sum_{i,j=1}^N \langle K(x_i, x_j) y_j, y_i \rangle_{\mathcal{Y}} \geq 0, \quad (3.12)$$

for all $N \in \mathbb{N}$, for all sequences of points $(x_i)_{i=1}^N$ in \mathcal{X}^N , and all sequences of points $(y_i)_{i=1}^N$ in \mathcal{Y}^N .

As in the scalar case any Vector Valued Reproducing Kernel Hilbert Space defines a unique positive-definite Operator-Valued Kernel and conversely a positive-definite Operator-Valued Kernel defines a unique Vector Valued Reproducing Kernel Hilbert Space.

Proposition 3.2 (Carmeli, De Vito, and Toigo [40]). *Given a Vector Valued Reproducing Kernel Hilbert Space there is a unique positive-definite Operator-Valued Kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$.*

Proof Given $x \in \mathcal{X}$, [Equation 3.9](#) ensures that the evaluation map at x defined as

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

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

$$K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}) \quad K(x, z) = ev_x ev_z^*.$$

Since for all $(x_i)_{i=1}^N$ in \mathcal{X}^N and all $(y_i)_{i=1}^N$ in \mathcal{Y}^N ,

$$\begin{aligned} \sum_{i,j=1}^N \langle K(x_i, x_j) y_j, y_i \rangle_{\mathcal{Y}} &= \sum_{i,j=1}^N \langle ev_{x_j}^* y_j, ev_{x_i}^* y_i \rangle_{\mathcal{Y}} \\ &= \left\langle \sum_{i=1}^N ev_{x_i}^* y_i, \sum_{i=1}^N ev_{x_i}^* y_i \right\rangle_{\mathcal{Y}} \\ &= \left\| \sum_{i=1}^N ev_{x_i}^* y_i \right\|_{\mathcal{Y}} \geq 0, \end{aligned}$$

the map K is positive-definite. □

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 $K_x = ev_x^*$. As a consequence we have that

$$K(x, z)y = ev_x ev_z^* y = K_x^* K_z y = (K_z y)(x). \quad (3.13)$$

Some direct consequences follow 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 $x \in \mathcal{X}$, $\text{ev}_x^*y = K_x y = K(\cdot, x)y$ such that

$$\langle f(x), y \rangle_{\mathcal{Y}} = \langle f, K(\cdot, x)y \rangle_{\mathcal{H}} = \langle K_x^*f, y \rangle_{\mathcal{Y}}. \quad (3.14)$$

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 comes from the fact that $\|K_x\|_{\mathcal{Y}, \mathcal{H}} = \|K_x^*\|_{\mathcal{H}, \mathcal{Y}} = \sqrt{\|K(x, x)\|_{\mathcal{Y}, \mathcal{Y}}}$ and the operator norm is sub-multiplicative.

Additionally given a positive-definite Operator-Valued Kernel, it defines a unique VV-RKHS.

Proposition 3.3 (Carmeli, De Vito, and Toigo [40]). *Given a positive-definite Operator-Valued Kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$, there is a unique Vector Valued Reproducing Kernel Hilbert Space \mathcal{H} on \mathcal{X} with reproducing kernel K .*

Proof Let $K_{x,y} = K(\cdot, x)y \in \mathcal{F}(\mathcal{X}; \mathcal{Y})$ and 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

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

such 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 a positive-definite Operator-Valued Kernel, we have that $\langle f, f \rangle_{\mathcal{H}_0} \geq 0$ for all $f \in \mathcal{H}_0$. Because the sesquilinear form is positive if \mathcal{Y} is a complex Hilbert space, it is also Hermitian. If \mathcal{Y} is a real Hilbert space, by assumption $K(x, z) = K(z, x)^*$, making $\langle \cdot, \cdot \rangle_{\mathcal{H}_0}$ an Hermitian sesquilinear form. 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

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

which implies that

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

Therefore if $\langle f, f \rangle_{\mathcal{H}_0} = 0$ then $f = 0$. Eventually we deduce that $\langle \cdot, \cdot \rangle_{\mathcal{H}_0}$ is an inner 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. 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 $f = 0$. Besides, as 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 $K_x^* f = (Wf)(x) = f(x) = ev_x f$, showing that \mathcal{H} is a Vector Valued Reproducing Kernel Hilbert Space with reproducing kernel

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

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

The above theorem also holds if \mathcal{Y} is a real Hilbert space provided we add the assumption that $K(x, z)$ is self-adjoint i.e. $K(x, z) = K(z, x)^*$ for all $x, z \in \mathcal{X}$. Then $K(x, z)$ still defines a valid symmetric bilinear form on \mathcal{Y} when \mathcal{Y} is a real Hilbert space.

Since an positive-definite Operator-Valued Kernel defines a unique Vector Valued Reproducing Kernel Hilbert Space (VV-RKHS) and conversely a VV-RKHS defines a unique Operator-Valued Kernel, we denotes 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 positive-definite Operator-Valued Kernels or reproducing Operator-Valued Kernels as Operator-Valued Kernels whether they act on complex or real Hilbert spaces. As a consequence, given K an Operator-Valued Kernel, define $K_x = K(\cdot, x)$ we have

$$K(x, z) = K_x^* K_z \quad \forall x, z \in \mathcal{X}, \tag{3.15a}$$

$$\mathcal{H}_K = \overline{\text{span}} \{ K_{x,y} \mid \forall x \in \mathcal{X}, \forall y \in \mathcal{Y} \}. \tag{3.15b}$$

Where $\overline{\text{span}}$ is the closed span of a given set. Another way to describe functions of \mathcal{H}_K consists in using a suitable feature map.

Proposition 3.4 (Feature Operator (Carmeli et al. [41])). 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

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

Then

$$\|f\|_K = \inf \{ \|g\|_{\mathcal{H}} \mid \forall g \in \mathcal{H}, Wg = f \}. \quad (3.16)$$

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

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

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

Let $\mathcal{H}_K = Im W$ 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 . We call this new partial isometry (again) 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 $Ker W$ is closed, W is bounded, and $\Phi(x)$ is bounded by definition such that the evaluation map

$$ev_x = (W\Phi(x))^*$$

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

$$K(x, z) = ev_x ev_z^* = (W\Phi(x))^* (W\Phi(z)) = \Phi(x)^* W^* W\Phi(z) = \Phi(x)^* \Phi(z),$$

Since $W^* W$ is the identity on $Im \Phi(z)$. Hence \mathcal{H}_K is a VV-RKHS (see proof of [Proposition 3.2](#)). \square

We call Φ a *feature map*, W a *feature operator* and \mathcal{H} a *feature space*. Since W is an 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} . Notice that W is a partial isometry, meaning that there can exist multiple functions $g \in \mathcal{H}$, the redescription space, such that $Wg = f$ where f is a function of the VV-RKHS \mathcal{H}_K . However [Equation 3.16](#) shows that there is a unique function $g \in \mathcal{H}$ such that $Wg = f$, and $\|g\|_{\mathcal{H}} = \|f\|_{\mathcal{H}_K}$. Among all functions $g \in \mathcal{H}$ such that $Wg = f$, the only one making the norm in the VV-RKHS and the redescription space is the one with minimal norm.

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 3.8 (\mathcal{Y} -Mercer kernel (Carmeli et al. [41])). 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 characterizes \mathcal{Y} -Mercer kernel in terms of the properties of a kernel rather than properties of the VV-RKHS.

Proposition 3.5 (Characterization of \mathcal{Y} -Mercer kernel (Carmeli et al. [41])). Let K be a reproducing kernel. The kernel K is Mercer if and only if the function $x \mapsto \|K(x, x)\|_{\mathcal{Y}, \mathcal{Y}}$ 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 [Equation 3.15b](#)). In addition for all $f \in \mathcal{H}_K$, $\|K_x^* f\| = \|f(x)\| \leq \|f\|_\infty$. Hence there exists a constant $M \in \mathbb{R}_+$ such that for all $x \in \mathcal{X}$, $\|K_x\| \leq M$. Therefore from [Equation 3.15a](#), 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

The next lemma shows that when \mathcal{X} and \mathcal{Y} are separable and \mathcal{H}_K is a space of continuous functions then \mathcal{H}_K is separable. It is worth mentioning that when the Hilbert space \mathcal{H}_K is separable, it admits a countable orthonormal basis.

Lemma 3.1 (Separable VV-RKHS (Carmeli, De Vito, and Toigo [40])). 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} ensure 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 $f(x) = ev_x f = 0$. Since f is continuous and \mathcal{X}_0 is dense in \mathcal{X} , for all $x \in \mathcal{X}$, $f(x) = 0$ thus $f = 0$. \square

Since a \mathcal{Y} -Mercer kernel K defines a VV-RKHS \mathcal{H}_K of continuous functions, \mathcal{H}_K is separable when \mathcal{X} and \mathcal{Y} are separable.

Proposition 3.6 (Separable VV-RKHS for \mathcal{Y} -Mercer kernel (Carmeli, De Vito, and Toigo [40])). *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 is separable.*

Proof From [Proposition 3.5](#) K is a \mathcal{Y} -Mercer kernel if and only if $\mathcal{H}_K \subset \mathcal{C}(\mathcal{X}; \mathcal{Y})$. Applying [Lemma 3.1](#) of Carmeli, De Vito, and Toigo [40], we have that \mathcal{H}_K is separable. \square

Thus since \mathcal{H}_K is also a Hilbert space and is separable it is second countable (i.e. it has a countable orthonormal basis). 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 .

3.3.2 Shift-Invariant OVK on LCA groups

The main subjects of interest of [Chapter 4](#) 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 3.9 (Shift-invariant OVK). *A reproducing Operator-Valued Kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ is called shift-invariant⁴ if for all $x, z, t \in \mathcal{X}$,*

$$K(x * t, z * t) = K(x, z). \quad (3.17)$$

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 binary group operation $*$.

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 $f \in \mathcal{H}_K$,

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

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 let us see a group as a linear operator which are well studied mathematical objects.

⁴ Also referred to as translation-invariant OVK.

Proposition 3.7 (Kernel signature (Carmeli et al. [41])). *Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ be a reproducing kernel. The following conditions are equivalents.*

1. K is a positive-definite shift-invariant Operator-Valued Kernel.
2. There is a positive-definite function $K_e : \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ such that $K(x, z) = K_e(z^{-1} * x)$.

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

$$K(x, z) = K_e^* \lambda_{x^{-1} * z} K_e \quad \forall (x, z) \in \mathcal{X}^2, \quad (3.18a)$$

$$\|K(x, x)\| = \|K_e(e)\| \quad \forall x \in \mathcal{X}. \quad (3.18b)$$

Proof Assume [Proposition 3.7 item 1](#) holds true. Given $x, z \in \mathcal{X}$, [Equation 3.13](#) and [Equation 3.17](#) yields

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

Since K is a reproducing kernel, K_e is of completely positive type, so that [Proposition 3.7 item 2](#) holds true. Besides if [Proposition 3.7 item 2](#) holds true obviously the definition of a reproducing kernel ([definition 3.6](#)) is fulfilled so that [Proposition 3.7 item 1](#) holds true.

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

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

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

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

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 [Equation 3.15b](#)), 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 [Equation 3.13](#) since for all $x \in \mathcal{X}$ and all $y \in \mathcal{Y}$

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

Moreover notice that shift-invariant \mathcal{Y} -Mercer kernels are directly linked to functions of positive type (see [Equation 3.1](#)), since shift-invariant \mathcal{Y} -Mercer kernels are nothing but functions whose signature is of positive type (continuous positive-definite functions).

3.3.3 Examples of Operator-Valued Kernels

In this subsection we list some Operator-Valued Kernels (OVKs) that have been used successfully in the litterature. We do not recall the proof that the following kernels are well defined and refer the interested reader to the respective authors original work.

OVKs 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.

Proposition 3.8 (Decomposable kernel (Micchelli and Pontil [113])). *Let Γ be a non-negative operator of $\mathcal{L}_+(\mathcal{Y})$. K is said to be a decomposable kernel⁵ if for all $(x, z) \in \mathcal{X}^2$,*

$$K(x, z) := k(x, z)\Gamma,$$

where k is a scalar kernel.

When $\mathcal{Y} = \mathbb{R}^p$, the operator Γ can be represented by a matrix which can be 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 [5], Baldassarre et al. [14], and Evgeniou, Micchelli, and Pontil [59] that Γ 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, Γ can be set to the empirical covariance of the output training data or learned with one of the algorithms proposed in the literature [54, 101, 158]. 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 of continuous functions $C(\mathcal{X}, \mathcal{Y})$. The conditions for a kernel to be universal have been discussed in Caponnetto et al. [39] and Carmeli et al. [41]. In particular they show that a decomposable kernel is universal provided that the scalar kernel k is universal and the operator Γ is injective. Given $(e_k)_{k=1}^p$ a basis of \mathcal{Y} , we recall here how the matrix Γ acts as a regularizer between the components of the outputs $f_k = \langle f(\cdot), e_k \rangle_{\mathcal{Y}}$ of a function $f \in \mathcal{H}_K$.

Proposition 3.9 (Kernels and Regularizers (Álvarez, Rosasco, and Lawrence [5])). *Let $K(x, z) := k(x, z)\Gamma$ for all $x, z \in \mathcal{X}$ be a decomposable kernel where Γ is a matrix of size $p \times p$. Then for all $f \in \mathcal{H}_K$,*

$$\|f\|_K = \sum_{i,j=1}^p \left(\Gamma^\dagger \right)_{ij} \langle f_i, f_j \rangle_K \quad (3.19)$$

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 .

We prove a generalized version of Proposition 3.9 to any Operator-Valued Kernel in Subsection 4.3.4.

⁵ Some authors also refer to as separable kernels.

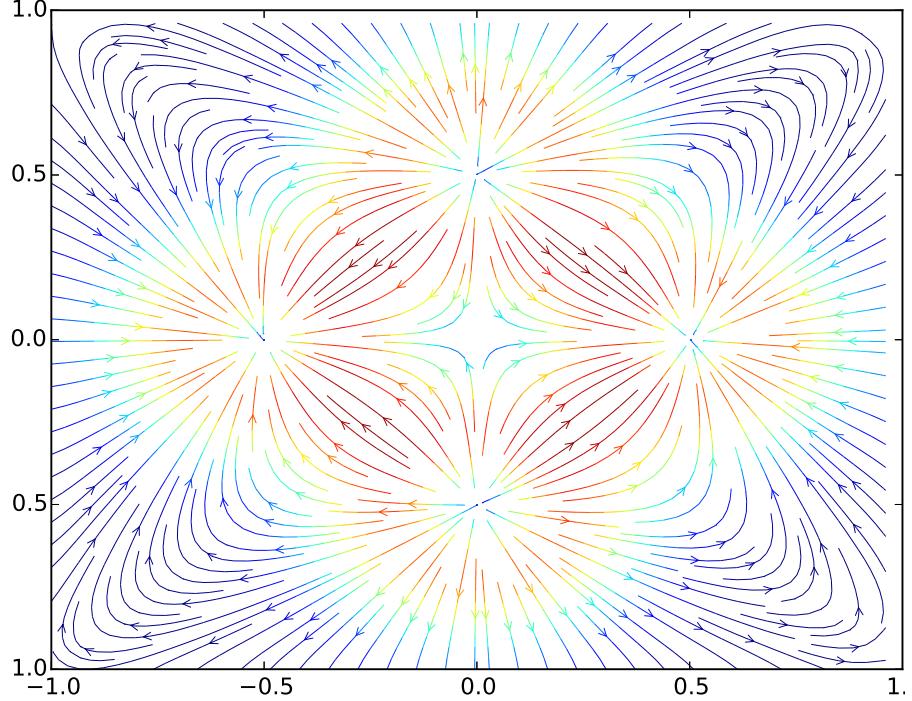


Figure 3.2: Synthetic 2D curl-free field .

Curl-free and divergence-free kernels provide an interesting application of operator-valued kernels [15, 109, 115] to *vector field* learning, for which input and output spaces have the same dimensions ($d = p$). Applications cover shape deformation analysis [115] and magnetic fields approximations [183]. These kernels discussed in [68] allow encoding input-dependent similarities between vector-fields. An illustration of a synthetic 2D curl-free and divergence free fields are given respectively in Figure 3.2 and Figure 3.3. To obtain the curl-free field we took the gradient of a mixture of five two dimensional Gaussians (since the gradient of a potential is always curl-free). We generated the divergence-free field by taking the orthogonal of the curl-free field.

Proposition 3.10 (Curl-free and Div-free kernel (Macedo and Castro [109])). 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

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

⁶ See Subsection 6.2.1 for a formal definition of the operator ∇ .

Although taken separately these kernels are not universal, a convex combination of the curl-free and divergence-free kernels allows to

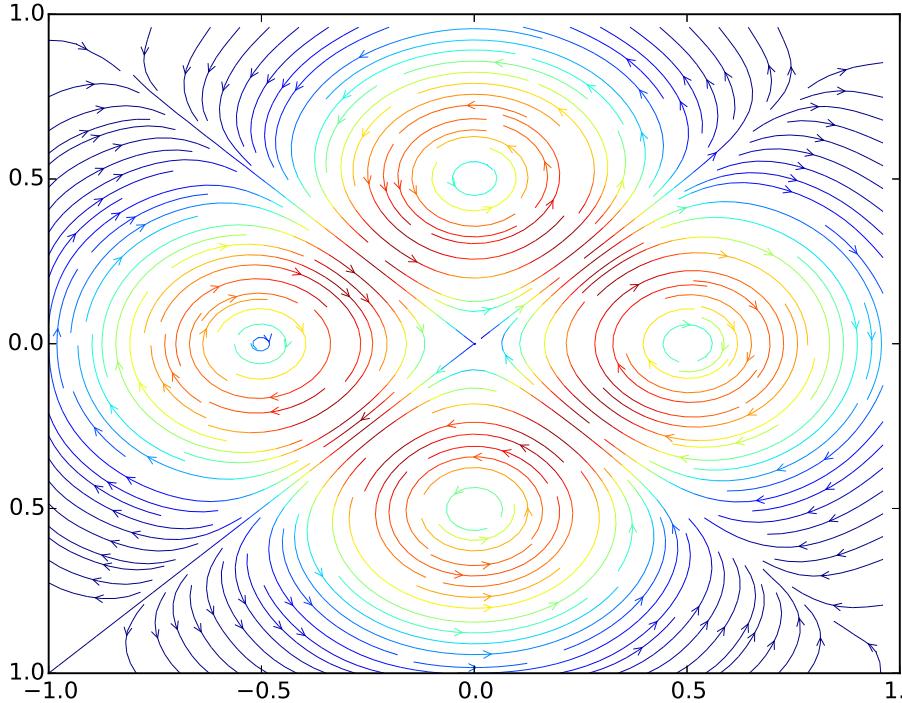


Figure 3.3: Synthetic 2D divergence-free field .

learn any vector field that satisfies the Helmholtz decomposition theorem [15, 109]. The next class of kernels we present are transformable kernels, whose action on each coordinate of an output vector is determined by “views” of an input data.

Proposition 3.11 (Transformable kernel (Caponnetto et al. [39])). *Let $k : \mathcal{X}' \times \mathcal{X}' \rightarrow \mathbb{R}$ be a scalar-valued kernels and ψ_1, \dots, ψ_p be a collection functions from $\mathcal{X} \rightarrow \mathcal{X}'$. Then the transformable kernel is defined for all $(i, j) \in (\mathbb{N}_p^*)^2$ as*

$$K(x, z)_{ij} = \langle e_i, K(x, z)e_j \rangle_y = k(\psi_i(x), \psi_j(z)),$$

for all $x, z \in \mathcal{X}$.

Transformable kernels have been successfully used for network inference from time series by means of autoregressive models (Lim et al. [102, 103]), and by Vazquez and Walter [179] for cokriging the multi-output version of kriging⁷, which takes into account the correlations between the outputs.

⁷ Gaussian process regression.

We also introduce an example of Operator-Valued Kernel acting on a function space which found applications in Kadri et al. [86].

Proposition 3.12 (Hilbert Schmidt Integral kernel (Kadri et al. [86])). *Let $k_{\mathcal{X}} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a scalar valued kernel acting on the inputs and $k_{\mathcal{T}} : \mathcal{T} \times \mathcal{T} \rightarrow \mathbb{R}$ be a scalar valued kernel acting on the outputs.*

Define the integral operator $L_{\mathcal{T}}g = \int_{\mathcal{Y}} k_{\mathcal{T}}(\cdot, t)g(t)d\mu(t)$. Then the Hilbert Schmidt Integral kernel is defined as

$$K : \begin{cases} \mathcal{X} \times \mathcal{X} & \rightarrow \mathcal{L}(\mathcal{Y}) \\ (x, z) & \mapsto k_{\mathcal{X}}(x, z)L_{\mathcal{T}}. \end{cases}$$

This kernel is useful to learn functions f that are function valued. In other words, $f \in \mathcal{F}(\mathcal{X}; \mathcal{F}(\mathcal{T}; \mathbb{R}))$ and the Operator-Valued Kernel K act on a function g in the following way.

$$K(x, z)g = k_{\mathcal{X}}(x, z) \int_{\mathcal{Y}} k_{\mathcal{Y}}(\cdot, t)g(t)d\mu(t). \quad (3.20)$$

In Kadri et al. [86], the authors studied the case where $\mathcal{T} = \mathbb{R}$ with $k_{\mathcal{Y}}(t, s) = \exp(-|t - s|)$ and applied it to speech inversion (Which and how Human articulators are activated from an audible speech signal). Notice that the Hilbert Schmidt integral kernel is a particular case of decomposable kernel, where $\mathcal{Y} = \mathcal{F}(\mathcal{T}; \mathbb{R})$.

3.3.4 Some use of Operator-valued kernels

We give here a non exhaustive list of works concerning Operator-Valued Kernels. A good review of Operator-Valued Kernels has been conducted in Álvarez, Rosasco, and Lawrence [5]. For a theoretical introduction to OVKs the interested reader can refer to the papers Caponnetto et al. [39], Carmeli, De Vito, and Toigo [40], and Carmeli et al. [41]. Generalization bounds for OVK have been studied in Kadri et al. [86], Maurer [112], Sangnier, Fercoq, and Buc [146], and Sindhwani, Minh, and Lozano [158].

Operator-valued Kernel Regression has first been studied in the context of Ridge Regression and Multi-task learning by Micchelli and Pontil [113].

Baldassarre et al. [15] and Macedo and Castro [109] showed the interest of spectral algorithms in Ridge regression and introduced vector field learning as a new multiple output task in Machine Learning community. Wahlström et al. [183] applied vector field learning with OVK-based Gaussian processes to the reconstruction of magnetic fields (which are curl-free).

Multi-task regression [114] and structured multi-class classification [54, 119, 123] are undoubtedly the first target applications for working in Vector Valued Reproducing Kernel Hilbert Space. Operator-Valued Kernels have been shown useful to provide a general framework for structured output prediction [34, 35] with a link to Output Kernel Regression [83]. Beyond structured classification, other

various applications such as link prediction, drug activity prediction or recently metabolite identification [36] and image colorization [76] have been developed.

The works of Kadri et al. [84] and Kadri et al. [86] have been the precursors of regression with functional values, opening a new avenue of applications. Appropriate algorithms devoted to on-line learning have been also derived by Audiffren and Kadri [10].

Kernel learning was addressed at least in two ways: first with using Multiple Kernel Learning in Kadri et al. [85] and second, using various penalties, smooth ones in Ciliberto et al. [43] and Dinuzzo et al. [54] for decomposable kernels and non smooth ones in Lim et al. [103] using proximal methods in the case of decomposable and transformable kernels.

Dynamical modeling was tackled in the context of multivariate time series modelling in Lim et al. [102, 103] and Sindhiani, Minh, and Lozano [158] and as a generalization of Recursive Least Square Algorithm in Amblard and Kadri [6].

Sangnier, Fercoq, and Buc [146] recently explored the minimization of a pinball loss under regularizing constraints induced by a well chosen decomposable kernel in order to handle joint quantile regression.



Part II
CONTRIBUTIONS

4

OPERATOR-VALUED RANDOM FOURIER FEATURES

In this first contribution chapter we present a generalization of the RFF framework introduced in [Chapter 2 \[29\]](#). This is based on an operator-valued Bochner theorem proposed by Carmeli et al. [41]. We use this theorem to construct an Operator-valued Random Fourier Feature (ORFF) from an OVK. Conversely we also show that it is possible to construct an ORFF from the regularization properties it induces rather than from an OVK. We give various examples of ORFF maps such as an ORFF map for the decomposable kernel, the curl-free kernel and the divergence-free kernel.

Contents

4.1	Motivation	54
4.2	Theoretical study	54
4.2.1	Sufficient conditions of existence	56
4.2.2	Examples of spectral decomposition	62
4.2.3	Functional Fourier feature map	67
4.3	Operator-valued Random Fourier Features	69
4.3.1	Building Operator-valued Random Fourier Features	69
4.3.2	From Operator Random Fourier Feature maps to OVKs	71
4.3.3	Examples of Operator Random Fourier Feature maps	75
4.3.4	Regularization property	77
4.4	Conclusions	79

4.1 MOTIVATION

Random Fourier Features have been proved useful to implement efficiently kernel methods in the scalar case. In this work, we propose to extend Random Fourier Feature methodology in order to approximate OVKs. As in the scalar case, we are mainly interested on explicit approximated feature maps because they open the door to learning linear models. Our final goal is to come up with the definition of $\tilde{\Phi} : \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}, \mathcal{H})$ for some Hilbert spaces \mathcal{Y} and \mathcal{H} , a feature map of some approximation \tilde{K} of a given OVK, K . This chapter is devoted to the construction of these approximations based on Random Fourier Feature principles. It is followed by a non asymptotical study of the error of approximation (Chapter 5) and the development of learning tools based on Operator Random Fourier Feature maps with practical and theoretical insights.

We present in this chapter a construction methodology devoted to shift-invariant \mathcal{Y} -Mercer operator-valued kernels defined on any Locally Compact Abelian (LCA) group, noted $(\mathcal{X}, *)$, for some operation noted $*$. This allows us to use the general context of Pontryagin duality for Fourier Transform of functions on LCA groups. Building upon a generalization of the celebrated 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 RFF for scalar-valued kernels and derive a general methodology to build Operator-valued Random Fourier Feature (ORFF) when operator-valued kernels are shift-invariant according to the chosen group operation. Elements of this chapter have been developed in Brault, Heinonen, and Buc [29].

We present a construction of feature maps called Operator-valued Random Fourier Feature (ORFF), such that $f : x \mapsto \Phi(x)^*\theta$ 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 $\tilde{K}(x, z) = \tilde{\Phi}(x)^*\tilde{\Phi}(z)$ with high probability on *compact* subsets of the LCA \mathcal{X} .

4.2 THEORETICAL STUDY

The following proposition of Neeb [126] and Zhang, Xu, and Zhang [193] extends Bochner's theorem to any shift-invariant \mathcal{Y} -Mercer kernel.

Proposition 4.1 (Operator-valued Bochner's theorem [126, 193]). If a 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 $\widehat{Q} : \mathcal{B}(\mathcal{X}) \rightarrow \mathcal{L}_+(\mathcal{Y})$ such that for all $x, z \in \mathcal{X}$,

$$K(x, z) = \int_{\widehat{\mathcal{X}}} \overline{(x * z^{-1}, \omega)} d\widehat{Q}(\omega), \quad (4.1)$$

where \widehat{Q} belongs to the set of all the projection-valued measures of bounded variation on the σ -algebra of Borel subsets of $\widehat{\mathcal{X}}$. Conversely, from any positive operator-valued measure M , a shift-invariant kernel K can be defined by [Equation 4.1](#).

Although this theorem is central to the spectral decomposition of shift-invariant \mathcal{Y} -Mercer OVK, the following results proved by Carmeli et al. [41] provides more 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 $(A, \widehat{\mu})$ where A is an operator-valued function on $\widehat{\mathcal{X}}$ and $\widehat{\mu}$ is a *real-valued* positive measure on $\widehat{\mathcal{X}}$ (instead of a operator-valued measure as in [Proposition 4.1](#)). 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 has such a spectral decomposition when \mathcal{Y} is finite dimensional or \mathcal{X} is compact.

Proposition 4.2 (Carmeli et al. [41]). Let $\widehat{\mu}$ be a positive measure on $\mathcal{B}(\widehat{\mathcal{X}})$ and $A : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$ such that $\langle A(\cdot)y, y' \rangle \in L^1(\mathcal{X}, \widehat{\mu})$ for all $y, y' \in \mathcal{Y}$ and $A(\omega) \succcurlyeq 0$ for $\widehat{\mu}$ -almost all $\omega \in \widehat{\mathcal{X}}$. Then, for all $\delta \in \mathcal{X}$,

$$K_e(\delta) = \int_{\widehat{\mathcal{X}}} \overline{(\delta, \omega)} A(\omega) d\widehat{\mu}(\omega) \quad (4.2)$$

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

$$(Wg)(x) = \int_{\widehat{\mathcal{X}}} \overline{(x, \omega)} B(\omega) g(\omega) d\widehat{\mu}(\omega), \quad (4.3)$$

Where $B(\omega)B(\omega)^* = A(\omega)$ and both integrals converge 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 $(A, \widehat{\mu})$.

When $p = 1$ one can always assume A is reduced to the scalar 1, $\widehat{\mu}$ is still a bounded positive measure and we retrieve the Bochner theorem applied to the scalar case ([Theorem 2.3](#)).

[Proposition 4.2](#) shows that a pair $(A, \hat{\mu})$ entirely characterizes an OVK. Namely a given measure $\hat{\mu}$ and a function A such that $\langle y', A(\cdot)y \rangle \in L^1(\mathcal{X}, \hat{\mu})$ for all $y, y' \in \mathcal{Y}$ and $A(\omega) \succcurlyeq 0$ for $\hat{\mu}$ -almost all ω , give rise to an OVK. Since $(A, \hat{\mu})$ determines a unique kernel we can write $\mathcal{H}_{(A, \hat{\mu})} \Rightarrow \mathcal{H}_K$ where K is defined as in [Equation 4.2](#). However the converse is not true: Given a \mathcal{Y} -Mercer shift invariant Operator-Valued Kernel, there exist infinitely many pairs $(A, \hat{\mu})$ that characterize an OVK.

The main difference between [Equation 4.1](#) and [Equation 4.2](#) 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 $\hat{\mu}$; and that this operator-valued density is not unique.

Finally [Proposition 4.2](#) does not provide any *constructive* way to obtain the pair $(A, \hat{\mu})$ that characterizes an OVK. The following [Sub-section 4.2.1](#) 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 $(C, \widehat{\text{Haar}})$ from it. Additionally, we show that we can always extract a scalar-valued *probability* density function from C such that we obtain a pair $(A, \Pr_{\widehat{\text{Haar}}, \rho})$ where $\Pr_{\widehat{\text{Haar}}, \rho}$ is a *probability* distribution absolutely continuous with respect to $\widehat{\text{Haar}}$ and with associated probability density function (p. d. f) ρ . Thus for all $Z \subset \mathcal{B}(\widehat{\mathcal{X}})$,

$$\Pr_{\widehat{\text{Haar}}, \rho}(Z) = \int_Z \rho(\omega) d\widehat{\text{Haar}}(\omega).$$

When the reference measure $\widehat{\text{Haar}}$ is the Lebesgue measure, we note

$$\begin{aligned} \Pr_{\widehat{\text{Haar}}, \rho} &= \Pr_{\text{Leb}, \rho} \\ &= \Pr_\rho. \end{aligned}$$

For any function $f : \mathcal{X} \times \widehat{\mathcal{X}} \times \mathcal{Y} \rightarrow \mathbb{R}$, we also use the notation

$$\begin{aligned} \mathbb{E}_{\widehat{\text{Haar}}, \rho}[f(x, \omega, y)] &= \mathbb{E}_{\omega \sim \Pr_{\widehat{\text{Haar}}, \rho}}[f(x, \omega, y)] \\ &= \int_{\widehat{\mathcal{X}}} f(x, \omega, y) d\Pr_{\widehat{\text{Haar}}, \rho}(\omega) \\ &= \int_{\widehat{\mathcal{X}}} f(x, \omega, y) \rho(\omega) d\widehat{\text{Haar}}(\omega). \end{aligned}$$

where the two last equalities hold by the transfer theorem and the fact that $\Pr_{\widehat{\text{Haar}}, \rho}$ has density ρ .

4.2.1 Sufficient conditions of existence

While [Proposition 4.2](#) 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 $(A, \mathbf{Pr}_{\hat{\mu}, \rho})$ from the kernel signature K_e . Proposition 14 in Carmeli et al. [41] gives the solution, and also provides a sufficient condition for [Proposition 4.2](#) to apply.

Proposition 4.3 (Carmeli et al. [41]). *Let K be a shift-invariant \mathcal{Y} -Mercer kernel of signature K_e . Suppose that for all $z \in \mathcal{X}$ and for all $y, y' \in \mathcal{Y}$, the function*

$$\langle K_e(\cdot)y, y' \rangle_y \in L^1(\mathcal{X}, \mathbf{Haar}),$$

where \mathcal{X} is endowed with the group law \star . Denote $C : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$, the function defined for all $\omega \in \widehat{\mathcal{X}}$ that satisfies for all y, y' in \mathcal{Y} :

$$\begin{aligned} \langle y', C(\omega)y \rangle_y &= \int_{\mathcal{X}} (\delta, \omega) \langle y', K_e(\delta)y \rangle_y d\mathbf{Haar}(\delta) \\ &= \mathcal{F}^{-1} [\langle y', K_e(\cdot)y \rangle_y](\omega). \end{aligned} \quad (4.4)$$

Then

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

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

We found that there has been confusion in the literature whether a kernel is the Fourier Transform or Inverse Fourier Transform of a measure. However [Lemma 4.1](#) clarifies the relation between the Fourier Transform and Inverse Fourier Transform for a translation invariant Operator-Valued Kernel. Here we show that in the real scalar case the Fourier Transform and Inverse Fourier Transform of a shift-invariant kernel are the same as well as in the operator-valued case.

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 to simplify the definition of ORFF.

Lemma 4.1 *Let K_e be the signature of a shift-invariant \mathcal{Y} -Mercer kernel such that for all $y, y' \in \mathcal{Y}$, $\langle y', K_e(\cdot)y \rangle_y \in L^1(\mathcal{X}, \mathbf{Haar})$ and let*

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

Then

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

3. $K_e(\delta)$ is self-adjoint and K_e is even.

Proof For any function f on (\mathcal{X}, \star) define the flip operator \mathcal{R} by

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

For any shift invariant \mathbb{Y} -Mercer kernel and for all $\delta \in \mathcal{X}$, $K_e(\delta) = K_e(\delta^{-1})^*$. Indeed from the definition of a shift-invariant kernel,

$$K_e(\delta^{-1}) = K(\delta^{-1}, e) = K(e, \delta) = K(\delta, e)^* = K_e(\delta)^*.$$

Item 1: taking the Fourier Transform yields,

$$\begin{aligned} \langle y', C(\omega)y \rangle_y &= \mathcal{F}^{-1} [\langle y', K_e(\cdot)y \rangle_y](\omega) \\ &= \mathcal{F}^{-1} [\langle y', (\mathcal{R}K_e(\cdot))^*y \rangle_y](\omega) \\ &= \mathcal{F}^{-1} [\langle \mathcal{R}K_e(\cdot)y', y \rangle_y](\omega) \\ &= \mathcal{F}^{-1} [\mathcal{R}\langle K_e(\cdot)y', y \rangle_y](\omega) \\ &= \mathcal{R}\mathcal{F}^{-1} [\langle K_e(\cdot)y', y \rangle_y](\omega) \\ &= \mathcal{R}\langle C(\cdot)y', y \rangle_y(\omega) \\ &= \langle y', C(\omega^{-1})^*y \rangle_y. \end{aligned}$$

Hence $C(\omega) = C(\omega^{-1})^*$. Suppose that \mathbb{Y} is a complex Hilbert space. Since for all $\omega \in \widehat{\mathcal{X}}$, $C(\omega)$ is bounded and non-negative so $C(\omega)$ is self-adjoint. Besides we have $C(\omega) = C(\omega^{-1})^*$ so C must be even. Suppose that \mathbb{Y} is a real Hilbert space. The Fourier Transform of a real valued function obeys $\mathcal{F}[f](\omega) = \overline{\mathcal{F}[f](\omega^{-1})}$. Therefore since $C(\omega)$ is non-negative for all $\omega \in \widehat{\mathcal{X}}$,

$$\begin{aligned} \langle y', C(\omega)y \rangle &= \overline{\langle y', C(\omega^{-1})y \rangle} = \langle y, C(\omega^{-1})^*y' \rangle \\ &= \langle y, C(\omega)y' \rangle. \end{aligned}$$

Hence $C(\omega)$ is self-adjoint and thus C is even.

Item 2: simply, for all $y, y' \in \mathbb{Y}$, $\langle y, C(\omega^{-1})y' \rangle = \langle y', C(\omega)y \rangle$ thus

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

Item 3: from Item 2 we have $\mathcal{F}^{-1} [\langle y', K_e(\cdot)y \rangle] = \mathcal{F}^{-1} \mathcal{R} \langle y', K_e(\cdot)y \rangle$. 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

While [Proposition 4.3](#) gives an explicit form of the operator $C(\omega)$ defined as the Fourier Transform of the kernel K , it is not really convenient to work with the Haar measure $\widehat{\text{Haar}}$ on $\mathcal{B}(\widehat{\mathcal{X}})$. However it is easily possible to turn $\widehat{\text{Haar}}$ into a probability measure to allow efficient (Monte-Carlo) integration over an infinite domain.

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 function from $C(\cdot)$.

Proposition 4.4 (Shift-invariant \mathcal{Y} -Mercer kernel spectral decomposition). *Let K_e be the signature of a shift-invariant \mathcal{Y} -Mercer kernel. If for all $y, y' \in \mathcal{Y}$, $\langle K_e(\cdot)y, y' \rangle \in L^1(\mathcal{X}, \text{Haar})$ then there exists a positive probability measure $\widehat{\text{Pr}_{\text{Haar}, \rho}}$ and an operator-valued function A such that for all $y, y' \in \mathcal{Y}$,*

$$\langle y', K_e(\delta)y \rangle = E_{\widehat{\text{Haar}}, \rho} \left[\overline{(\delta, \omega)} \langle y', A(\omega)y \rangle \right], \quad (4.5)$$

with

$$\langle y', A(\omega)y \rangle \rho(\omega) = \mathcal{F}[\langle y', K_e(\cdot)y \rangle](\omega). \quad (4.6)$$

Moreover

1. for all $y, y' \in \mathcal{Y}$, $\langle A(\cdot)y, y' \rangle \in L^1(\widehat{\mathcal{X}}, \widehat{\text{Pr}_{\text{Haar}, \rho}})$,
2. $A(\omega)$ is non-negative for $\widehat{\text{Pr}_{\text{Haar}, \rho}}$ -almost all $\omega \in \widehat{\mathcal{X}}$,
3. $A(\cdot)$ and $\rho(\cdot)$ are even functions.

Proof This is a simple consequence of [Proposition 4.3](#) and [Lemma 4.1](#). 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 .

$$\begin{aligned} \langle y', K_e(\delta)y \rangle(\omega) &= \int_{\widehat{\mathcal{X}}} \overline{(\delta, \omega)} \langle y', C(\omega)y \rangle d\widehat{\text{Haar}}(\omega) \\ &= \int_{\widehat{\mathcal{X}}} \overline{(\delta, \omega)} \left\langle y', \frac{1}{\rho(\omega)} C(\omega)y \right\rangle \rho(\omega) d\widehat{\text{Haar}}(\omega). \end{aligned}$$

It is always possible to choose $\rho(\omega)$ such that $\int_{\widehat{\mathcal{X}}} \rho(\omega) d\widehat{\text{Haar}}(\omega) = 1$. For instance choose

$$\rho(\omega) = \frac{\|C(\omega)\|_{y, y}}{\int_{\widehat{\mathcal{X}}} \|C(\omega)\|_{y, y} d\widehat{\text{Haar}}(\omega)}$$

Since for all $y, y' \in \mathcal{Y}$, $\langle y', C(\cdot)y \rangle \in L^1(\widehat{\mathcal{X}}, \widehat{\text{Haar}})$ and \mathcal{Y} is a separable Hilbert space, by Pettis measurability theorem, $\int_{\widehat{\mathcal{X}}} \|C(\omega)\|_{y, y} d\widehat{\text{Haar}}(\omega)$ is finite and so is $\|C(\omega)\|_{y, y}$ for all $\omega \in \widehat{\mathcal{X}}$. Therefore $\rho(\omega)$ is the density of a probability measure $\widehat{\text{Pr}_{\text{Haar}, \rho}}$, i.e. conclude by taking

$$\widehat{\text{Pr}_{\text{Haar}, \rho}}(\mathcal{Z}) = \int_{\mathcal{Z}} \rho(\omega) d\widehat{\text{Haar}}(\omega),$$

for all $\mathcal{Z} \in \mathcal{B}(\widehat{\mathcal{X}})$. □

In the case where $\mathcal{Y} = \mathbb{R}^p$, we rewrite [Equation 4.6](#) coefficient-wise by choosing an orthonormal basis $(e_j)_{j \in \mathbb{N}_p^*}$ of \mathbb{R}^p .

$$A(\omega)_{ij}\rho(\omega) = \mathcal{F}[K_e(\cdot)_{ij}](\omega). \quad (4.7)$$

It follows that for all i and j in \mathbb{N}_p^* ,

$$K_e(x * z^{-1})_{ij} = \mathcal{F}[A(\cdot)_{ij}\rho(\cdot)](x * z^{-1}) \quad (4.8)$$

Remark 4.1 Note that although the Fourier Transform of K_e yields a unique operator-valued function $C(\cdot)$, the decomposition of $C(\cdot)$ into $A(\cdot)\rho(\cdot)$ is again not unique. The choice of the decomposition may be justified by the computational cost.

Another difficulty arises from the fact that the quantity

$$\sup_{\omega \in \widehat{\mathcal{X}}} \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}}$$

obtained in [Proposition 4.4](#) might not be bounded. Later, when we will focus on Monte-Carlo approximation of these integrals, we will have to take care of the unboundedness of $\|A(\cdot)\|_{\mathcal{Y}, \mathcal{Y}}$ that forbids the use of the most simple concentrations inequalities that require the boundedness of the random variable to be controlled. Therefore in the context of Operator-Valued Kernel concentration inequalities for unbounded random operators should be used.

However, as pointed out by Minh [118], under some condition on the trace of $K_e(\delta)$, it is possible to turn $A(\cdot)$ into a bounded random operator for all ω in $\widehat{\mathcal{X}}$. The idea is to define a sum measure $\rho = \sum_{j \in \mathbb{N}^*} \rho_{e_j}$, which gives rise to a bounded operator $A(\omega)$ and is independent of the $\{e_j\}_{j \in \mathbb{N}^*}$ base, instead of constructing a measure from the operator norm as in [Proposition 4.4](#). Additionally with such construction the measure associated to $A(\cdot)$ is *independent* from the basis of \mathcal{Y} . We present this result and in this proof we relax the assumptions of Minh [118] which requires $\int_{\mathcal{X}} |\text{Tr } K_e(\delta)| d\text{Haar}(\delta)$ to be well defined. We only require $\text{Tr } K_e(e)$ to be well defined.

Proposition 4.5 (Bounded shift-invariant \mathcal{Y} -Mercer kernel spectral decomposition (adaptation of a result from Minh [118])). Let K_e be the signature of a shift-invariant \mathcal{Y} -Mercer kernel. If for all y and y' in \mathcal{Y} , $\langle K_e(\cdot)y, y' \rangle \in L^1(\mathcal{X}, \text{Haar})$ and $\text{Tr } K_e(e) \in \mathbb{R}$, then

$$\langle y', K_e(\delta)y \rangle = E_{\widehat{\text{Haar}}, \rho_{\text{Tr}}} \left[\overline{(\delta, \omega)} \langle y, A_{\text{Tr}}(\omega)y' \rangle \right]. \quad (4.9)$$

with

$$\langle y', C(\cdot)y \rangle = \mathcal{F}[\langle y', K_e(\cdot)y \rangle] \quad (4.10a)$$

$$c_{\text{Tr}} = \text{Tr}[K_e(e)] \quad (4.10b)$$

$$A_{\text{Tr}}(\omega) = c_{\text{Tr}} \text{Tr}[C(\omega)]^{-1} C(\omega) \quad (4.10c)$$

$$\rho_{\text{Tr}}(\omega) = c_{\text{Tr}}^{-1} \text{Tr}[C(\omega)]. \quad (4.10d)$$

Moreover

1. For all $y, y' \in \mathcal{Y}$, $\langle y, A_{\text{Tr}}(\cdot)y' \rangle \in L^1(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho_{\text{Tr}}})$.
2. $A_{\text{Tr}}(\omega)$ is non-negative for all $\omega \in \widehat{\mathcal{X}}$,
3. $\text{ess sup}_{\omega \in \widehat{\mathcal{X}}} \|A_{\text{Tr}}(\omega)\|_{y, y} \leq c_{\text{Tr}}$,
4. $A_{\text{Tr}}(\cdot)$ and ρ_{Tr} are even functions.

Proof Let $\{e_j\}_{j \in \mathbb{N}^*}$ be an orthonormal basis of \mathcal{Y} . Notice that

$$\begin{aligned} \int_{\widehat{\mathcal{X}}} \langle e_j, C(\omega)e_j \rangle d\widehat{\text{Haar}}(\omega) &= \int_{\widehat{\mathcal{X}}} \underbrace{\overline{(e_j, \omega)}}_{=1} \langle e_j, C(\omega)e_j \rangle d\widehat{\text{Haar}}(\omega) \\ &= \langle e_j, K_e(e)e_j \rangle. \end{aligned}$$

Since $C(\omega)$ is non-negative, all the $\langle e_j, C(\omega)e_j \rangle$. Thus using the monotone convergence theorem,

$$\begin{aligned} \int_{\widehat{\mathcal{X}}} \text{Tr}[C(\omega)] d\widehat{\text{Haar}}(\omega) &= \int_{\widehat{\mathcal{X}}} \sum_{j \in \mathbb{N}^*} \langle e_j, C(\omega)e_j \rangle d\widehat{\text{Haar}}(\omega) \\ &= \sum_{k \in \mathbb{N}^*} \langle e_j, K_e(e)e_j \rangle \\ &= \text{Tr}[K_e(e)] = c_{\text{Tr}} < \infty. \end{aligned}$$

Let $A_{\text{Tr}}(\omega)$ and $\rho_{\text{Tr}}(\omega)$ be defined as in [Equation 4.10c](#) and [Equation 4.10d](#), respectively. By definition, $\int_{\widehat{\mathcal{X}}} \rho_{\text{Tr}}(\omega) d\widehat{\text{Haar}}(\omega) = 1$ and $A_{\text{Tr}}(\omega)\rho_{\text{Tr}}(\omega) = C(\omega)$. Now it remains to check the finiteness of $\text{Tr}[C(\omega)]$ for all $\omega \in \widehat{\mathcal{X}}$. Since for all $\omega \in \widehat{\mathcal{X}}$, $\text{Tr}[C(\omega)] \geq 0$,

$$\text{Tr}[C(\omega)] \leq \int_{\widehat{\mathcal{X}}} \text{Tr}[C(\omega)] d\widehat{\text{Haar}}(\omega) = \text{Tr}[K_e(e)] < \infty.$$

Since $\text{Tr}[C(\omega)]$ is positive and its integral is finite, ρ_{Tr} is a probability density function. In particular $C(\omega)$ is self-adjoint operator thus $\|C(\omega)\|_{\sigma, \infty} = \|C(\omega)\|_{y, y}$ for all $\omega \in \widehat{\mathcal{X}}$. Thus the Schatten norms $\|\cdot\|_{\sigma, p}$ verifies $\text{Tr}[\|\cdot\|] = \|\cdot\|_{\sigma, 1} \geq \|\cdot\|_{\sigma, p} \geq \|\cdot\|_{\sigma, q} \geq \|\cdot\|_{\sigma, \infty} = \|\cdot\|_{y, y}$ for all $p, q \in \mathbb{N}^*$ such that $1 \leq p \leq q \leq \infty$. Therefore since for all $\omega \in \widehat{\mathcal{X}}$, $C(\omega)$ is non-negative, we have for $\Pr_{\widehat{\text{Haar}}, \rho}$ -almost all $\omega \in \widehat{\mathcal{X}}$,

$$\begin{aligned} \|A_{\text{Tr}}(\omega)\|_{y, y} &= c_{\text{Tr}} \text{Tr}[C(\omega)]^{-1} \|C(\omega)\|_{\sigma, \infty} \\ &\leq c_{\text{Tr}} \text{Tr}[C(\omega)]^{-1} \|C(\omega)\|_{\sigma, 1} \\ &= c_{\text{Tr}} \text{Tr}[C(\omega)]^{-1} \text{Tr}[|C(\omega)|] \\ &= c_{\text{Tr}} \text{Tr}[C(\omega)]^{-1} \text{Tr}[C(\omega)] \\ &\leq c_{\text{Tr}} < \infty. \end{aligned}$$

Thus $\text{ess sup}_{\omega \in \widehat{\mathcal{X}}} \|A(\omega)\|_{y, y} \leq c_{\text{Tr}} < \infty$. As C is an even function, so are A_{Tr} and ρ_{Tr} . Eventually $\langle y', C(\cdot)y \rangle$ is in $L^1(\widehat{\mathcal{X}}, \widehat{\text{Haar}})$, thus $\langle y, A_{\text{Tr}}(\cdot)\rho_{\text{Tr}}(\cdot)y' \rangle$ is in $L^1(\widehat{\mathcal{X}}, \widehat{\text{Haar}})$, hence $\langle y, A_{\text{Tr}}(\cdot)y' \rangle \in L^1(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho_{\text{Tr}}})$. Since the trace is independent of the basis of \mathcal{Y} , so is ρ_{Tr} . \square

If \mathcal{Y} is finite dimensional then $\text{Tr}[K_e(e)]$ is well defined hence [Proposition 4.5](#) is valid as long as $K_e(\cdot)_{ij} \in L^1(\mathcal{X}, \text{Haar})$ for all $i, j \in \mathbb{N}_p^*$, where p is the dimension of \mathcal{Y} .

4.2.2 Examples of spectral decomposition

In this section we give examples of spectral decomposition for various \mathcal{Y} -Mercer kernels, based on [Proposition 4.4](#) and [Proposition 4.5](#).

4.2.2.1 Gaussian decomposable kernel

Recall that a decomposable \mathbb{R}^p -Mercer introduced in the Background section has the form $K(x, z) = k(x, z)\Gamma$, where $k(x, z)$ is a scalar Mercer kernel and $\Gamma \in \mathcal{L}(\mathbb{R}^p)$ is a non-negative operator. Let us focus on $K_e^{dec, gauss}(\cdot) = k_e^{gauss}(\cdot)\Gamma$, the Gaussian decomposable kernel where $K_e^{dec, gauss}$ and k_e^{gauss} are respectively the signature of K and k on the additive group $\mathcal{X} = (\mathbb{R}^d, +)$ – i.e. $\delta = x - z$ and $e = 0$. The well known Gaussian kernel is defined for all $\delta \in \mathbb{R}^d$ as follows

$$k_0^{gauss}(\delta) = \exp\left(-\frac{1}{2\sigma^2}\|\delta\|_2^2\right)$$

where $\sigma \in \mathbb{R}_+$ is an hyperparameter corresponding to the bandwidth of the kernel. The –Pontryagin– dual group of $\mathcal{X} = (\mathbb{R}^d, +)$ is $\widehat{\mathcal{X}} \cong (\mathbb{R}^d, +)$ with the pairing

$$(\delta, \omega) = \exp(i\langle \delta, \omega \rangle)$$

where δ and $\omega \in \mathbb{R}^d$. In this case the Haar measures on \mathcal{X} and $\widehat{\mathcal{X}}$ are in both cases the Lebesgue measure. However in order to have the property that $\mathcal{F}^{-1}[\mathcal{F}[f]] = f$ and $\mathcal{F}^{-1}[f] = \mathcal{R}\mathcal{F}[f]$ one must normalize both measures by $\sqrt{2\pi}^{-d}$, i.e. for all $\mathcal{Z} \in \mathcal{B}(\mathbb{R}^d)$,

$$\begin{aligned} \sqrt{2\pi}^d \mathbf{Haar}(\mathcal{Z}) &= \mathbf{Leb}(\mathcal{Z}) \text{ and} \\ \sqrt{2\pi}^d \widehat{\mathbf{Haar}}(\mathcal{Z}) &= \mathbf{Leb}(\mathcal{Z}). \end{aligned}$$

Then the Fourier Transform on $(\mathbb{R}^d, +)$ is

$$\begin{aligned} \mathcal{F}[f](\omega) &= \int_{\mathbb{R}^d} \exp(-i\langle \delta, \omega \rangle) f(\delta) d\mathbf{Haar}(\delta) \\ &= \int_{\mathbb{R}^d} \exp(-i\langle \delta, \omega \rangle) f(\delta) \frac{d\mathbf{Leb}(\delta)}{\sqrt{2\pi}^d}. \end{aligned}$$

Since $k_0^{gauss} \in L^1$ and Γ is bounded, it is possible to apply [Proposition 4.4](#), and obtain for all y and $y' \in \mathcal{Y}$,

$$\begin{aligned} \langle y', C^{dec, gauss}(\omega)y \rangle &= \mathcal{F}\left[\langle y', K_0^{dec, gauss}(\cdot)y \rangle\right](\omega) \\ &= \mathcal{F}[k_0^{gauss}](\omega) \langle y', \Gamma y \rangle. \end{aligned}$$

Thus

$$C^{dec, gauss}(\omega) = \int_{\mathbb{R}^d} \exp\left(-i\langle \omega, \delta \rangle - \frac{\|\delta\|_2^2}{2\sigma^2}\right) \frac{d\mathbf{Leb}(\delta)}{\sqrt{2\pi}^d} \Gamma.$$

Hence

$$C^{\text{dec, gauss}}(\omega) = \underbrace{\frac{1}{\sqrt{2\pi\frac{1}{\sigma^2}^d}} \exp\left(-\frac{\sigma^2}{2}\|\omega\|_2^2\right)}_{\rho(\cdot)=\mathcal{N}(0, \sigma^{-2}I_d)\sqrt{2\pi}^d} \underbrace{\Gamma}_{A(\cdot)=\Gamma}.$$

Therefore the canonical decomposition of $C^{\text{dec, gauss}}$ is $A^{\text{dec, gauss}}(\omega) = \Gamma$ and $\rho^{\text{dec, gauss}} = \mathcal{N}(0, \sigma^{-2}I_d)\sqrt{2\pi}^d$, where \mathcal{N} is the Gaussian probability distribution. Note that this decomposition is done with respect to the *normalized* Lebesgue measure $\widehat{\text{Haar}}$, meaning that for all $Z \in \mathcal{B}(\widehat{\mathcal{X}})$,

$$\begin{aligned} \mathbf{Pr}_{\widehat{\text{Haar}}, \mathcal{N}(0, \sigma^{-2}I_d)\sqrt{2\pi}^d}(Z) &= \int_Z \mathcal{N}(0, \sigma^{-2}I_d)\sqrt{2\pi}^d d\widehat{\text{Haar}}(\omega) \\ &= \int_{\widehat{\mathcal{X}}} \mathcal{N}(0, \sigma^{-2}I_d) d\mathbf{Leb}(\omega) \\ &= \mathbf{Pr}_{\mathcal{N}(0, \sigma^{-2}I_d)}(Z). \end{aligned}$$

Thus, the same decomposition with respect to the usual –non-normalized– Lebesgue measure \mathbf{Leb} yields

$$A^{\text{dec, gauss}}(\cdot) = \Gamma \quad (4.11a)$$

$$\rho^{\text{dec, gauss}} = \mathcal{N}(0, \sigma^{-2}I_d). \quad (4.11b)$$

If Γ is a trace class operator, applying [Proposition 4.5](#) yields the same decomposition since $\mathbf{Tr}[K_0^{\text{dec, gauss}}(0)] = \mathbf{Tr}[\Gamma]$ and

$$\mathbf{Tr}[C^{\text{dec, gauss}}(\cdot)] = \mathcal{N}(0, \sigma^{-2}I_d)\sqrt{2\pi}^d \mathbf{Tr}[\Gamma].$$

4.2.2.2 Skewed- χ^2 decomposable kernel

The skewed- χ^2 scalar kernel [98], useful for image processing, is defined on the LCA group $\mathcal{X} = (-c_k; +\infty)_{k=1}^d$, with $c_k \in \mathbb{R}_+$ and endowed with the group operation \odot . Let $(e_k)_{k=1}^d$ be the standard basis of \mathcal{X} and $_k : x \mapsto \langle x, e_k \rangle$. The operator $\odot : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{X}$ is defined by

$$x \odot z = ((x_k + c_k)(z_k + c_k) - c_k)_{k=1}^d.$$

The identity element e is $(1 - c_k)_{k=1}^d$ since $(1 - c) \odot x = x$. Thus the inverse element x^{-1} is $((x_k + c_k)^{-1} - c_k)_{k=1}^d$. The skewed- χ^2 scalar kernel reads

$$k_{1-c}^{\text{skewed}}(\delta) = \prod_{k=1}^d \frac{2}{\sqrt{\delta_k + c_k} + \sqrt{\frac{1}{\delta_k + c_k}}}. \quad (4.12)$$

The dual of \mathcal{X} is $\widehat{\mathcal{X}} \cong \mathbb{R}^d$ with the pairing

$$(\delta, \omega) = \prod_{k=1}^d \exp(i \log(\delta_k + c_k) \omega_k).$$

The Haar measure are defined for all $\mathcal{Z} \in \mathcal{B}((-c; +\infty)^d)$ and all $\widehat{\mathcal{Z}} \in \mathcal{B}(\mathbb{R}^d)$ by

$$\begin{aligned}\sqrt{2\pi}^d \mathbf{Haar}(\mathcal{Z}) &= \int_{\mathcal{Z}} \prod_{k=1}^d \frac{1}{z_k + c_k} d\mathbf{Leb}(z) \\ \sqrt{2\pi}^d \widehat{\mathbf{Haar}}(\widehat{\mathcal{Z}}) &= \mathbf{Leb}(\widehat{\mathcal{Z}}).\end{aligned}$$

Thus the Fourier Transform is

$$\mathcal{F}[f](\omega) = \int_{(-c; +\infty)^d} \prod_{k=1}^d \frac{\exp(-i \log(\delta_k + c_k) \omega_k)}{\delta_k + c_k} f(\delta) \frac{d\mathbf{Leb}(\delta)}{\sqrt{2\pi}^d}.$$

Then, applying Fubini's theorem over product space, and the fact that each dimension is independent

$$\mathcal{F}[k_0^{skewed}](\omega) = \prod_{k=1}^d \int_{-\infty}^{+\infty} \frac{2 \exp(-i \log(\delta_k + c_k) \omega_k)}{(\delta_k + c_k) \left(\sqrt{\delta_k + c_k} + \sqrt{\frac{1}{\delta_k + c_k}} \right)} \frac{d\mathbf{Leb}(\delta_k)}{\sqrt{2\pi}^d}.$$

Making the change of variable $t_k = (\delta_k + c_k)^{-1}$ yields

$$\begin{aligned}\mathcal{F}[k_0^{skewed}](\omega) &= \prod_{k=1}^d \int_{-\infty}^{+\infty} \frac{2 \exp(-it_k \omega_k)}{\exp(\frac{1}{2}t_k) + \exp(-\frac{1}{2}t_k)} \frac{d\mathbf{Leb}(t_k)}{\sqrt{2\pi}^d} \\ &= \sqrt{2\pi}^d \prod_{k=1}^d \operatorname{sech}(\pi \omega_k).\end{aligned}$$

Since $k_{1-c}^{skewed} \in L^1$ and Γ is bounded, it is possible to apply [Proposition 4.4](#), and obtain

$$\begin{aligned}C^{\text{dec}, \text{skewed}}(\omega) &= \mathcal{F}[k_{1-c}^{skewed}](\omega) \Gamma \\ &= \sqrt{2\pi}^d \underbrace{\prod_{k=1}^d \operatorname{sech}(\pi \omega_k)}_{\rho(\cdot) = \mathcal{S}(0, 2^{-1})^d \sqrt{2\pi}^d} \underbrace{\Gamma}_{A(\cdot)}.\end{aligned}$$

Hence the decomposition with respect to the usual –non-normalized– Lebesgue measure \mathbf{Leb} yields

$$A^{\text{dec}, \text{skewed}}(\cdot) = \Gamma \tag{4.13a}$$

$$\rho^{\text{dec}, \text{skewed}} = \mathcal{S}(0, 2^{-1})^d. \tag{4.13b}$$

4.2.2.3 Curl-free Gaussian kernel

The curl-free Gaussian kernel is defined as $K_0^{\text{curl,gauss}} = -\nabla\nabla^\top k_0^{\text{gauss}}$. Here $\mathcal{X} = (\mathbb{R}^d, +)$ so the setting is the same than [Subsubsection 4.2.2.1](#).

$$\begin{aligned} C^{\text{curl,gauss}}(\omega)_{ij} &= \mathcal{F}\left[K_{1-c}^{\text{curl,gauss}}(\cdot)_{ij}\right](\omega) \\ &= \mathcal{F}\left[-\frac{\partial^2}{\partial\delta_i\partial\delta_j}k_0^{\text{gauss}}\right](\omega) \\ &= -(i\omega_i)(i\omega_j)\mathcal{F}[k_0^{\text{gauss}}](\omega) \\ &= \omega_i\omega_j\mathcal{F}[k_0^{\text{gauss}}](\omega) \\ &= \sqrt{2\pi\frac{1}{\sigma^2}} \exp\left(-\frac{\sigma^2}{2}\|\omega\|_2^2\right) \sqrt{2\pi}^d \omega_i\omega_j. \end{aligned}$$

Hence

$$C^{\text{curl,gauss}}(\omega) = \underbrace{\frac{1}{\sqrt{2\pi\frac{1}{\sigma^2}}} \exp\left(-\frac{\sigma^2}{2}\|\omega\|_2^2\right) \sqrt{2\pi}^d}_{\mu(\cdot)=\mathcal{N}(0,\sigma^{-2}I_d)\sqrt{2\pi}^d} \underbrace{\omega\omega^\top}_{A(\omega)=\omega\omega^\top}.$$

Here a canonical decomposition is $A^{\text{curl,gauss}}(\omega) = \omega\omega^\top$ for all $\omega \in \mathbb{R}^d$ and $\rho^{\text{curl,gauss}} = \mathcal{N}(0, \sigma^{-2}I_d)\sqrt{2\pi}^d$ with respect to the normalized Lebesgue measure **Leb**. Again the decomposition with respect to the usual –non-normalized– Lebesgue measure is for all $\omega \in \mathbb{R}^d$

$$A^{\text{curl,gauss}}(\omega) = \omega\omega^\top \quad (4.14a)$$

$$\rho^{\text{curl,gauss}} = \mathcal{N}(0, \sigma^{-2}I_d). \quad (4.14b)$$

Notice that in this case $\|A^{\text{curl,gauss}}(\cdot)\|_{\mathbb{R}^d, \mathbb{R}^d}$ is not bounded. However applying [Proposition 4.5](#) yields a different decomposition where the quantity $\|A_{\text{Tr}}^{\text{curl,gauss}}(\cdot)\|_{\mathbb{R}^d, \mathbb{R}^d}$ is bounded. First we have for all $\delta \in \mathbb{R}^d$ and for all $i, j \in \mathbb{N}_d^*$

$$\frac{\partial^2}{\partial\delta_i\partial\delta_j}k_0^{\text{gauss}}(\delta) = \frac{\exp\left(-\frac{1}{2\sigma^2}\|\delta\|_2^2\right)}{\sigma^2} \begin{cases} \frac{\delta_i\delta_j}{\sigma^2} & \text{if } i \neq j \\ \left(1 - \frac{\delta_i\delta_j}{\sigma^2}\right) & \text{otherwise.} \end{cases}$$

Hence

$$-\nabla\nabla^\top k_0^{\text{gauss}}(\delta) = \left(I_d - \frac{\delta\delta^\top}{\sigma^2}\right) \frac{\exp\left(-\frac{1}{2\sigma^2}\|\delta\|_2^2\right)}{\sigma^2}.$$

Thus

$$\begin{aligned} \text{Tr}\left[K_0^{\text{curl,gauss}}(0)\right] &= \text{Tr}\left[\nabla\nabla^\top k_0^{\text{gauss}}(0)\right] \\ &= d\sigma^{-2} \end{aligned}$$

and

$$\text{Tr}[C(\omega)] = \|\omega\|_2^2 \mathcal{N}(0, \sigma^{-2} I_d) \sqrt{2\pi}^d.$$

Apply [Proposition 4.5](#) to obtain the decomposition $A_{\text{Tr}}^{\text{curl}, \text{gauss}}(\omega) = \omega\omega^\top \|\omega\|_2^{-2}$ and the measure $\rho_{\text{Tr}}^{\text{curl}, \text{gauss}}(\omega) = \sigma^2 d^{-1} \|\omega\|_2^2 \mathcal{N}(0, \sigma^{-2}) \sqrt{2\pi}^d$ for all $\omega \in \mathbb{R}^d$, with respect to the normalized Lebesgue measure. Therefore the decomposition with respect to the usual non-normalized Lebesgue measure is

$$A_{\text{Tr}}^{\text{curl}, \text{gauss}}(\omega) = \frac{\omega\omega^\top}{\|\omega\|_2^2} \quad (4.15a)$$

$$\rho_{\text{Tr}}^{\text{curl}, \text{gauss}}(\omega) = \frac{\sigma^2}{d} \|\omega\|_2^2 \mathcal{N}(0, \sigma^{-2})(\omega). \quad (4.15b)$$

This example also illustrates that there exists many decompositions of $C(\omega)$ into $(A(\omega), \widehat{\text{Pr}_{\text{Haar}, \rho}}(\omega))$.

4.2.2.4 Divergence-free kernel

The divergence-free Gaussian kernel is defined as $K_0^{\text{div}, \text{gauss}} = (\nabla\nabla^\top - \Delta)k_0^{\text{gauss}}$ on the group $\mathcal{X} = (\mathbb{R}^d, +)$. The setting is the same than [Subsubsection 4.2.2.1](#). Hence

$$\begin{aligned} C^{\text{div}, \text{gauss}}(\omega)_{ij} &= \mathcal{F}\left[K_0^{\text{div}, \text{gauss}}(\cdot)_{ij}\right](\omega) \\ &= \mathcal{F}\left[\frac{\partial^2}{\partial \delta_i \partial \delta_j} k_0^{\text{gauss}} - \delta_{i=j} \sum_{k=1}^d \frac{\partial^2}{\partial \delta_k \partial \delta_k} k_0^{\text{gauss}}\right](\omega) \\ &= \left(-(i\omega_i)(i\omega_j) - \delta_{i=j} \sum_{k=1}^d (i\omega_k)^2\right) \mathcal{F}[k_0^{\text{gauss}}] \\ &= \left(\delta_{i=j} \sum_{k=1}^d \omega_k^2 - \omega_i \omega_j\right) \mathcal{F}[k_0^{\text{gauss}}](\omega). \end{aligned}$$

Hence

$$C^{\text{div}, \text{gauss}}(\omega) = \underbrace{\frac{1}{\sqrt{2\pi \frac{1}{\sigma^2}}} \exp\left(-\frac{\sigma^2}{2} \|\omega\|_2^2\right)}_{\rho(\cdot) = \mathcal{N}(0, \sigma^{-2} I_d) \sqrt{2\pi}^d} \underbrace{\sqrt{2\pi}^d \left(I_d \|\omega\|_2^2 - \omega\omega^\top\right)}_{A(\omega) = I_d \|\omega\|_2^2 - \omega\omega^\top}.$$

Thus the canonical decomposition with respect to the normalized Lebesgue measure is $A^{\text{div}, \text{gauss}}(\omega) = I_d \|\omega\|_2^2 - \omega\omega^\top$ and the measure

$$\rho^{\text{div}, \text{gauss}} = \mathcal{N}(0, \sigma^{-2} I_d) \sqrt{2\pi}^d.$$

The canonical decomposition with respect to the usual Lebesgue measure is

$$A^{\text{div}, \text{gauss}}(\omega) = I_d \|\omega\|_2^2 - \omega\omega^\top \quad (4.16a)$$

$$\rho^{\text{div}, \text{gauss}} = \mathcal{N}(0, \sigma^{-2} I_d). \quad (4.16b)$$

To obtain the bounded decomposition, again, apply [Proposition 4.5](#). For all $\delta \in \mathbb{R}^d$,

$$\sum_{k=1}^d \frac{\partial^2}{\partial \delta_k \partial \delta_k} k_0^{gauss}(\delta) = \left(d - \frac{\|\delta\|_2^2}{\sigma^2} \right) \frac{\exp\left(-\frac{1}{2\sigma^2}\|\delta\|_2^2\right)}{\sigma^2}.$$

Thus overall,

$$K_0^{div, gauss}(\delta) = \left(\frac{\delta \delta^\top}{\sigma^2} + \left((d-1) - \frac{\|\delta\|_2^2}{\sigma^2} \right) I_d \right) \frac{\exp\left(-\frac{1}{2\sigma^2}\|\delta\|_2^2\right)}{\sigma^2}.$$

Eventually $\text{Tr}[K_0^{div, gauss}(0)] = \text{Tr}[(\nabla \nabla^\top - \Delta) k_0^{gauss}(0)] = d(d-1)\sigma^{-2}$ and $\text{Tr}[C(\omega)] = (d-1)\|\omega\|_2^2 \mathcal{N}(0, \sigma^2 I_d) \sqrt{2\pi}^d$. As a result the decomposition with respect to the normalized Lebesgue measure is $A_{\text{Tr}}^{div, gauss}(\omega) = (I_d - \omega \omega^\top \|\omega\|_2^{-2})$ and $\rho_{\text{Tr}}^{div, gauss}(\omega) = d^{-1}\sigma^2 \|\omega\|_2^2 \mathcal{N}(0, \sigma^2 I_d) \sqrt{2\pi}^d$. The decomposition with respect to the normalized Lebesgue measure being

$$A_{\text{Tr}}^{div, gauss}(\omega) = I_d - \frac{\omega \omega^\top}{\|\omega\|_2^2} \quad (4.17a)$$

$$\rho_{\text{Tr}}^{div, gauss} = \frac{\sigma^2}{d} \|\omega\|_2^2 \mathcal{N}(0, \sigma^{-2} I_d). \quad (4.17b)$$

4.2.3 Functional Fourier feature map

We introduce a *functional* feature map, we call *Fourier Feature map*, defined by the following proposition as a direct consequence of [Proposition 4.2](#).

Proposition 4.6 (Functional Fourier feature map). *Let \mathcal{Y} and \mathcal{Y}' be two Hilbert spaces. If there exist an operator-valued function $B : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y}, \mathcal{Y}')$ such that for all $y, y' \in \mathcal{Y}$,*

$$\langle y, B(\omega)B(\omega)^*y' \rangle_{\mathcal{Y}} = \langle y', A(\omega)y \rangle_{\mathcal{Y}}$$

$\widehat{\mu}$ -almost everywhere and $\langle y', A(\cdot)y \rangle \in L^1(\widehat{\mathcal{X}}, \widehat{\mu})$ then the operator Φ_x defined for all y in \mathcal{Y} by

$$(\Phi_x y)(\omega) = (x, \omega)B(\omega)^*y, \quad (4.18)$$

is a feature map⁸ of some shift-invariant \mathcal{Y} -Mercer kernel K .

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

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

⁸ 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.

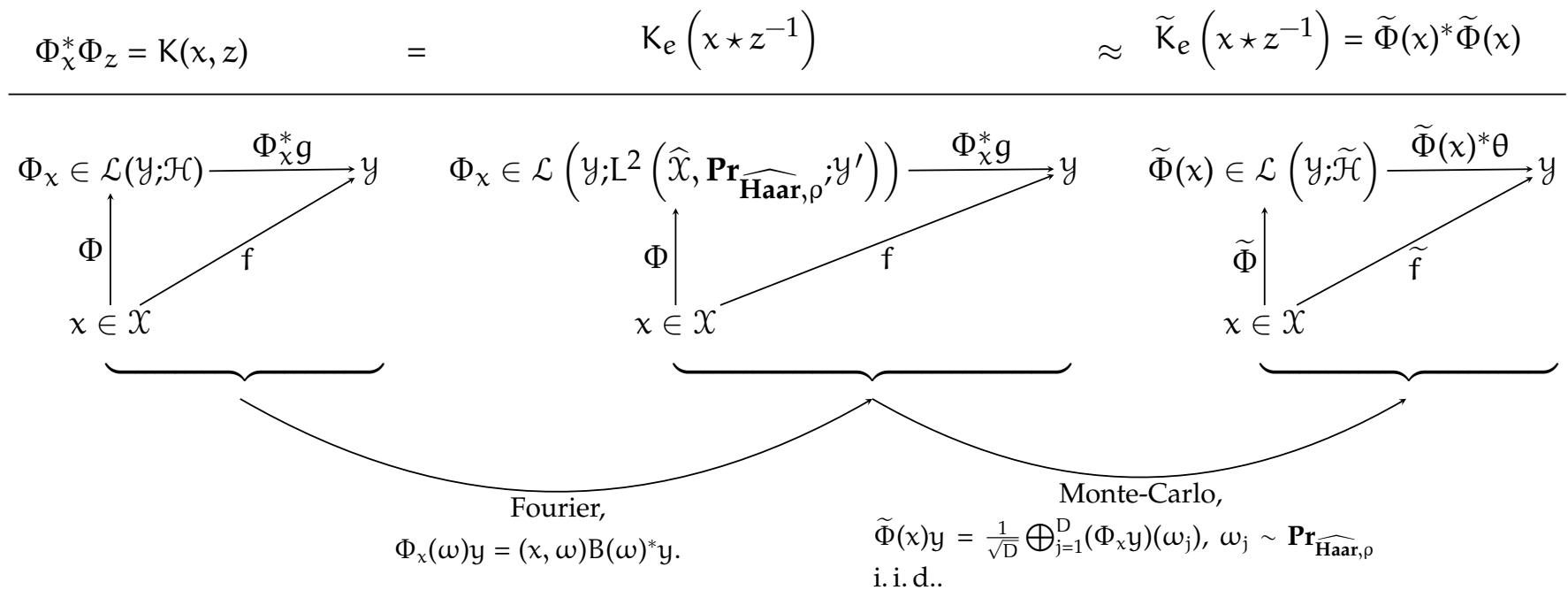


Figure 4.1: Relationships between feature-maps. For any realization of $\omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho}$ i.i.d., $\widetilde{\mathcal{H}} = \bigoplus_{j=1}^D \mathcal{Y}'$.

which defines a \mathcal{Y} -Mercer according to [Proposition 4.2](#) of Carmeli et al. [[41](#)]. \square

With this notation we have $\Phi : \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}; L^2(\widehat{\mathcal{X}}, \widehat{\mu}; \mathcal{Y}'))$ such that $\Phi_x \in \mathcal{L}(\mathcal{Y}; L^2(\widehat{\mathcal{X}}, \widehat{\mu}; \mathcal{Y}'))$ where $\Phi_x := \Phi(x)$.

4.3 OPERATOR-VALUED RANDOM FOURIER FEATURES

4.3.1 Building Operator-valued Random Fourier Features

As shown in [Propositions 4.4](#) and [4.5](#) it is always possible to find a pair $(A, \widehat{\Pr_{\text{Haar}, \rho}})$ from a shift invariant \mathcal{Y} -Mercer Operator-Valued Kernel K_e such that $\widehat{\Pr_{\text{Haar}, \rho}}$ is a probability measure, i.e. $\int_{\widehat{\mathcal{X}}} \rho d\widehat{\Pr_{\text{Haar}}} = 1$ where ρ is the density of $\widehat{\Pr_{\text{Haar}, \rho}}$ and $K_e(\delta) = E_\rho(\delta, \omega)A(\omega)$. In order to obtain an approximation of K from a decomposition $(A, \widehat{\Pr_{\text{Haar}, \rho}})$ we turn our attention to a Monte-Carlo estimation of the expectations in [Equation 4.9](#) and [Equation 4.5](#) characterizing a \mathcal{Y} -Mercer shift-invariant Operator-Valued Kernel.

Proposition 4.7 *Let $K(x, z)$ be a shift-invariant \mathcal{Y} -Mercer kernel with signature K_e such that for all $y, y' \in \mathcal{Y}$, $\langle y', K_e(\cdot)y \rangle \in L^1(\mathcal{X}, \text{Haar})$. Then one can find a pair $(A, \widehat{\Pr_{\text{Haar}, \rho}})$ that satisfies [Proposition 4.4](#). i.e. for $\widehat{\Pr_{\text{Haar}, \rho}}$ -almost all ω , and all $y, y' \in \mathcal{Y}$,*

$$\langle y, A(\omega)y' \rangle \rho(\omega) = \mathbb{F}[\langle y', K_e(\cdot)y \rangle](\omega).$$

If $(\omega_j)_{j=1}^D$ be a sequence of $D \in \mathbb{N}^*$ i. i. d. random variables following the law $\widehat{\Pr_{\text{Haar}, \rho}}$ then the operator-valued function \tilde{K} defined for $(x, z) \in \mathcal{X} \times \mathcal{X}$ as

$$\tilde{K}(x, z) = \frac{1}{D} \sum_{j=1}^D \overline{(x * z^{-1}, \omega_j)} A(\omega_j)$$

is an approximation⁹ of K .

Proof From the strong law of large numbers

$$\frac{1}{D} \sum_{j=1}^D \overline{(x * z^{-1}, \omega_j)} A(\omega_j) \xrightarrow[D \rightarrow \infty]{a.s.} E_{\widehat{\Pr_{\text{Haar}, \rho}}}[(\overline{(x * z^{-1}, \omega_j)}) A(\omega)]$$

⁹ i.e. it satisfies for all $x, z \in \mathcal{X}$,
 $\tilde{K}(x, z) \xrightarrow[D \rightarrow \infty]{a.s.} K(x, z)$ in the weak operator topology,
where K is a \mathcal{Y} -Mercer OVK.

where the integral converges in the weak operator topology. Then by [Proposition 4.4](#),

$$E_{\widehat{\Pr_{\text{Haar}, \rho}}}[(\overline{(x * z^{-1}, \omega_j)}) A(\omega)] = K_e(x * z^{-1}).$$

\square

Now, for efficient computations as motivated in the introduction, we are interested in finding an approximated *feature map* instead of a kernel approximation. Indeed, an approximated feature map will allow to build linear models in regression tasks. The following proposition deals with the feature map construction.

Proposition 4.8 *Assume the same conditions as Proposition 4.7. Moreover, if one can define $B : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y}', \mathcal{Y})$ such that for $\Pr_{\widehat{\text{Haar}}, \rho}$ -almost all ω , and all $u, v \in \mathcal{Y}$,*

$$\langle u, B(\omega)B(\omega)^*v \rangle \rho(\omega) = \langle u, A(\omega)v \rangle \rho(\omega) = \mathcal{F}[\langle u, K_e(\cdot)v \rangle](\omega),$$

then the function $\tilde{\Phi} : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y}, \bigoplus_{j=1}^D \mathcal{Y}')$ defined for all $y \in \mathcal{Y}$ as follows:

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D (x, \omega_j) B(\omega_j)^* y, \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho} \text{ i. i. d.},$$

is an approximated feature map¹⁰ for the kernel K .

Proof Let $(\omega_j)_{j=1}^D$ be a sequence of $D \in \mathbb{N}^*$ i. i. d. random variables following the law $\Pr_{\widehat{\text{Haar}}, \rho}$. For all $x, z \in \mathcal{X}$ and all $u, v \in \mathcal{Y}$,

$$\begin{aligned} & \left\langle \tilde{\Phi}(x)u, \tilde{\Phi}(z)v \right\rangle_{\bigoplus_{j=1}^D \mathcal{Y}'} \\ &= \frac{1}{D} \left\langle \bigoplus_{j=1}^D ((x, \omega_j) B(\omega_j)^* u), \bigoplus_{j=1}^D ((z, \omega_j) B(\omega_j)^* v) \right\rangle \end{aligned}$$

By definition of the inner product in direct sum of Hilbert spaces,

$$\begin{aligned} & \frac{1}{D} \left\langle \bigoplus_{j=1}^D ((x, \omega_j) B(\omega_j)^* u), \bigoplus_{j=1}^D ((z, \omega_j) B(\omega_j)^* v) \right\rangle \\ &= \frac{1}{D} \sum_{j=1}^D \left\langle u, \overline{(x, \omega_j)} B(\omega_j)(z, \omega_j) B(\omega_j)^* v \right\rangle_{\mathcal{Y}} \\ &= \left\langle u, \left(\frac{1}{D} \sum_{j=1}^D \overline{(x \star z^{-1}, \omega_j)} A(\omega_j) \right) v \right\rangle_{\mathcal{Y}}, \end{aligned}$$

Eventually apply Proposition 4.7 to obtain the convergence of the Monte-Carlo plug-in estimator to the true kernel K . \square

Remark 4.2 We find a decomposition such that $A(\omega_j) = B(\omega_j)B(\omega_j)^*$ for all $j \in \mathbb{N}_D^*$ either by exhibiting a closed-form or using a numerical decomposition.

Notice that an ORFF map as defined in Proposition 4.8 is also the Monte-Carlo sampling of the corresponding functional Fourier feature map $\Phi_x : \mathcal{Y} \rightarrow L^2(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho}; \mathcal{Y}')$ as defined in Proposition 4.6. Indeed, for all $y \in \mathcal{Y}$ and all $x \in \mathcal{X}$,

$$\tilde{\Phi}(x)y = \bigoplus_{j=1}^D (\Phi_x y)(\omega_j), \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho} \text{ i. i. d..}$$

¹⁰ i. e. it satisfies $\tilde{\Phi}(x)^* \tilde{\Phi}(z) \xrightarrow[D \rightarrow \infty]{a.s.} K(x, z)$ in the weak operator topology, where K is a \mathcal{Y} -Mercer OVK

Proposition 4.8 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 shift-invariant \mathcal{Y} -Mercer kernel such that $\forall y, y' \in \mathcal{Y}, \langle y', K_e(\cdot)y \rangle \in L^1(\mathbb{R}^d, \text{Haar})$ and D the number of features.
- Output:** A random feature $\tilde{\Phi}(x)$ such that $\tilde{\Phi}(x)^* \tilde{\Phi}(z) \approx K(x, z)$
- 1 Define the pairing (x, ω) from the LCA group $(\mathcal{X}, *)$;
 - 2 Find a decomposition $(A, \Pr_{\widehat{\text{Haar}}, \rho})$ and B such that

$$\begin{aligned} B(\omega)B(\omega)^*\rho(\omega) &= A(\omega)\rho(\omega) \\ &= \mathcal{F}^{-1}[K_e](\omega); \end{aligned}$$

- 3 Draw D i. i. d. realizations $(\omega_j)_{j=1}^D$ from the probability distribution $\Pr_{\widehat{\text{Haar}}, \rho}$;
 - 4 **return** $\begin{cases} \tilde{\Phi}(x) \in \mathcal{L}(\mathcal{Y}, \tilde{\mathcal{H}}) & : y \mapsto \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D (x, \omega_j) B(\omega_j)^* y; \\ \tilde{\Phi}(x)^* \in \mathcal{L}(\tilde{\mathcal{H}}, \mathcal{Y}) & : \theta \mapsto \frac{1}{\sqrt{D}} \sum_{j=1}^D (x, \omega_j) B(\omega_j) \theta_j \end{cases}$
-

We give a numerical illustration of different \tilde{K} built from different i. i. d. samples $(\omega_1, \dots, \omega_D)$. In [Figure 4.2](#), we represent the approximation of a reference function (black line) defined as $(y_1, y_2)^T = f(x_i) = \sum_{j=1}^{250} K_{ij} u_j$ where $u_j \sim \mathcal{N}(0, I_2)$ and K is a Gaussian decomposable kernel defined as

$$K_{ij} = \exp\left(-\frac{(x_i - x_j)^2}{2(0.1)^2}\right) \Gamma, \quad \text{for } i, j \in \mathbb{N}_{250}^*.$$

We took $\Gamma = .5I_2 + .51_2$ such that the outputs y_1 and y_2 share some similarities. We generated 250 points $(x_i)_{i=1}^N$, equally separated on the segment $(-1; 1)$. Then we computed an approximate kernel matrix $\tilde{K} \approx K$ for 25 increasing values of D ranging from 1 to 10^4 . The top row of the two graphs in [Figure 4.2](#) shows that the more the number of features increases the closer the model $\tilde{f}(x_i) = \sum_{j=1}^{250} \tilde{K}_{ij} u_j$ is to f . The bottom row shows the same experiment but for a different realization of \tilde{K} . When D is small the curves of the bottom and top rows are very dissimilar –and sine wave like– while they both converge to f when D increases.

4.3.2 From Operator Random Fourier Feature maps to OVKs

It is also interesting to notice that we can go the other way and define from the general form of an Operator-valued Random Fourier Feature, an operator-valued kernel.

Proposition 4.9 (Operator Random Fourier Feature map). *Let \mathcal{Y} and \mathcal{Y}' be two Hilbert spaces. If one defines an operator-valued function on the*

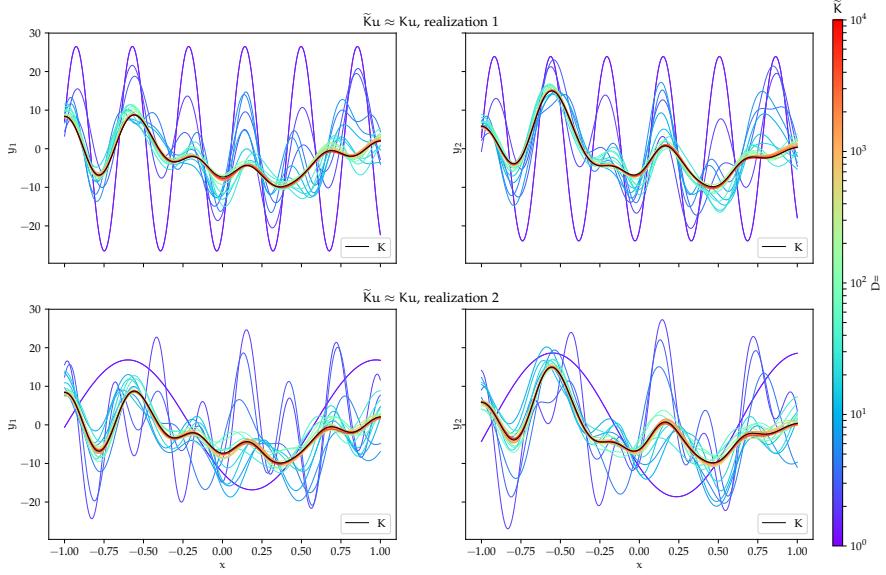


Figure 4.2: Approximation of a function in a VV-RKHS using different realizations of Operator Random Fourier Feature. Top row and bottom row correspond to two different realizations of \tilde{K} , which are different Operator-Valued Kernel. However when D tends to infinity, the different realizations of \tilde{K} yield the same OVK.

dual of a LCA group \mathcal{X} , $B : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y}, \mathcal{Y}')$, and a probability measure $\Pr_{\widehat{\text{Haar}}, \rho}$ on $B(\widehat{\mathcal{X}})$, such that for all $y \in \mathcal{Y}$ and all $y' \in \mathcal{Y}'$, $\langle y, B(\cdot)y' \rangle \in L^2(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho})$, then the operator-valued function

$$\tilde{\Phi} : \mathcal{X} \rightarrow \mathcal{L}\left(y, \bigoplus_{j=1}^D y'\right)$$

defined for all $x \in \mathcal{X}$ and for all $y \in \mathcal{Y}$ by

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D (x, \omega_j) B(\omega_j)^* y, \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho}, \text{ i.i.d.}, \quad (4.19)$$

¹¹ i.e. it satisfies
 $\tilde{\Phi}(x)^* \tilde{\Phi}(z) \xrightarrow[D \rightarrow \infty]{a.s.} K(x, z)$ in the weak
operator topology,
where K is a
 \mathcal{Y} -Mercer OVK

is an approximated feature map of some \mathcal{Y} -Mercer operator-valued kernel¹¹.

Proof Similarly to the proof of [Proposition 4.8](#), we have

$$\left\langle \tilde{\Phi}(x)y, \tilde{\Phi}(z)y' \right\rangle_{\bigoplus_{j=1}^D \mathcal{Y}'} = \left\langle y, \frac{1}{D} \sum_{j=1}^D \overline{(x * z^{-1}, \omega_j)} A(\omega_j) y' \right\rangle$$

and by the strong law of large numbers,

$$\frac{1}{D} \sum_{j=1}^D \overline{(x * z^{-1}, \omega_j)} A(\omega_j) \xrightarrow[D \rightarrow \infty]{a.s.} E_{\widehat{\text{Haar}}, \rho}[(\overline{x * z^{-1}, \omega}) A(\omega)]$$

in the weak operator topology. Now from [Proposition 4.2](#) with $\widehat{\mu} = \Pr_{\widehat{\text{Haar}}, \rho}$ we obtain $E_{\widehat{\text{Haar}}, \rho}[(\overline{x * z^{-1}, \omega}) A(\omega)] = K_e(x * z^{-1})$, K_e being the signature of some shift-invariant \mathcal{Y} -Mercer kernel. \square

The difference between [Proposition 4.9](#) and [Proposition 4.8](#) is that in [Proposition 4.9](#) we do not assume that $\Lambda(\omega)$ and $\Pr_{\widehat{\mathcal{H}}^{\text{haar}}, \rho}$ have been obtained from [Proposition 4.4](#). We conclude by showing that any realization of an approximate feature map gives a proper operator valued kernel. Hence we can always view $\tilde{K}(x, z) = \tilde{\Phi}(x)^* \tilde{\Phi}(z)$ —where $\tilde{\Phi}$ is defined as in [Proposition 4.7](#) (construction from an OVK) or [Proposition 4.9](#)— as a \mathcal{Y} -Mercer and thus apply all the classic results of the Operator-Valued Kernel theory on \tilde{K} .

Proposition 4.10 *Let $\omega \in \widehat{\mathcal{X}}^D$. If for all $y, y' \in \mathcal{Y}$*

$$\begin{aligned} \langle y', \tilde{K}_e(x * z^{-1}) y \rangle_{\mathcal{Y}} &= \langle \tilde{\Phi}(x)y', \tilde{\Phi}(z)y \rangle_{\widetilde{\mathcal{H}}} \\ &= \left\langle y', \frac{1}{D} \sum_{j=1}^D \overline{(x * z^{-1}, \omega_j)} B(\omega_j) B(\omega_j)^* y \right\rangle_{\mathcal{Y}}, \end{aligned}$$

for all $x, z \in \mathcal{X}$, then \tilde{K} is a shift-invariant \mathcal{Y} -Mercer Operator-Valued Kernel.

Proof Apply [Proposition 3.4](#) to $\tilde{\Phi}$ considering the Hilbert space $\widetilde{\mathcal{H}}$ to show that \tilde{K} is an OVK. Then [Proposition 3.7](#) shows that \tilde{K} is shift-invariant since $\tilde{K}(x, z) = \tilde{K}_e(x * z^{-1})$. Since $B(\omega)$ is a bounded operator, \tilde{K} is \mathcal{Y} -Mercer because all the functions in the sum are continuous. \square

Note that the above theorem does not consider the ω_j 's as random variables and therefore does not show the convergence of the kernel \tilde{K} to some target kernel K . However it shows that any realization of \tilde{K} when ω_j 's are random variables yields a valid \mathcal{Y} -Mercer operator-valued kernel.

Indeed, as a result of [Proposition 4.10](#), in the same way we defined an ORFF, we can define an approximate feature operator \tilde{W} which maps $\widetilde{\mathcal{H}}$ onto $\mathcal{H}_{\tilde{K}}$, where

$$\tilde{K}(x, z) = \tilde{\Phi}(x)^* \tilde{\Phi}(z), \quad \text{for all } x, z \in \mathcal{X}.$$

Definition 4.1 (Random Fourier feature operator). Let $\omega = (\omega_j)_{j=1}^D \in \widehat{\mathcal{X}}^D$ and let

$$\tilde{K}_e = \frac{1}{D} \sum_{j=1}^D \overline{(\cdot, \omega_j)} B(\omega_j) B(\omega_j)^*.$$

We call random Fourier feature operator the linear application $\tilde{W} : \widetilde{\mathcal{H}} \rightarrow \mathcal{H}_{\tilde{K}}$ defined as

$$(\tilde{W}\theta)(x) := \tilde{\Phi}(x)^* \theta = \frac{1}{\sqrt{D}} \sum_{j=1}^D \overline{(x, \omega_j)} B(\omega_j) \theta_j$$

where $\theta = \bigoplus_{j=1}^D \theta_j \in \widetilde{\mathcal{H}}$. Then from [Proposition 3.4](#),

$$(Ker \tilde{W})^\perp = \overline{\text{span}} \left\{ \tilde{\Phi}(x)y \mid \forall x \in \mathcal{X}, \forall y \in \mathcal{Y} \right\} \subseteq \widetilde{\mathcal{H}}.$$

The random Fourier feature operator is useful to show the relations between the random Fourier feature map with the functional feature map defined in [Proposition 4.6](#). The relationship between the generic feature map (defined for all Operator-Valued Kernel) the functional feature map (defining a shift-invariant \mathcal{Y} -Mercer Operator-Valued Kernel) and the random Fourier feature map is presented in [Figure 4.1](#).

Proposition 4.11 *For any $g \in \mathcal{H} = L^2(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho}; \mathcal{Y}')$, let*

$$\theta := \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D g(\omega_j), \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho} \text{ i.i.d. .}$$

Then

1. $(\widetilde{W}\theta)(x) = \widetilde{\Phi}(x)^* \theta \xrightarrow[D \rightarrow \infty]{a.s.} \Phi_x^* g = (Wg)(x)$,
2. $\|\theta\|_{\mathcal{H}}^2 \xrightarrow[D \rightarrow \infty]{a.s.} \|g\|_{\mathcal{H}}^2$,

Proof (of Proposition 4.11 item 1) Since $(\omega_j)_{j=1}^D$ are i.i.d. random vectors, for all $y \in \mathcal{Y}$ and for all $y' \in \mathcal{Y}'$, $\langle y, B(\cdot)y' \rangle \in L^2(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho})$ and $g \in L^2(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho}; \mathcal{Y}')$, from the strong law of large numbers

$$\begin{aligned} (\widetilde{W}\theta)(x) &= \widetilde{\Phi}(x)^* \theta \\ &= \frac{1}{D} \sum_{j=1}^D \overline{(x, \omega_j)} B(\omega_j) g(\omega_j), \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho} \text{ i.i.d.} \\ &\xrightarrow[D \rightarrow \infty]{a.s.} \int_{\widehat{\mathcal{X}}} \overline{(x, \omega)} B(\omega) g(\omega) d\Pr_{\widehat{\text{Haar}}, \rho}(\omega) \\ &= (Wg)(x) := \Phi_x^* g. \end{aligned} \quad \square$$

Proof (of Proposition 4.11 item 2) Again, since $(\omega_j)_{j=1}^D$ are i.i.d. random vectors and

$$g \in L^2(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho}; \mathcal{Y}'),$$

from the strong law of large numbers

$$\begin{aligned} \|\theta\|_{\mathcal{H}}^2 &= \frac{1}{D} \sum_{j=1}^D \|g(\omega_j)\|_{\mathcal{Y}'}^2, \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho} \text{ i.i.d.} \\ &\xrightarrow[D \rightarrow \infty]{a.s.} \int_{\widehat{\mathcal{X}}} \|g(\omega)\|_{\mathcal{Y}'}^2 d\Pr_{\widehat{\text{Haar}}, \rho}(\omega) \\ &= \|g\|_{L^2(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho}; \mathcal{Y}')}^2. \end{aligned} \quad \square$$

We write $\tilde{\Phi}(x)^*\tilde{\Phi}(x) \approx K(x, z)$ when $\tilde{\Phi}(x)^*\tilde{\Phi}(x) \xrightarrow{a.s.} K(x, z)$ in the weak operator topology when D tends to infinity. With mild abuse of notation we say that $\tilde{\Phi}(x)$ is an approximate feature map of the functional feature map Φ_x i. e. $\tilde{\Phi}(x) \approx \Phi_x$, when for all $y', y \in \mathcal{Y}$,

$$\begin{aligned}\langle y, K(x, z)y' \rangle_{\mathcal{Y}} &= \langle \Phi_x y, \Phi_z y' \rangle_{L^2(\widehat{\mathcal{X}}, \Pr_{\text{Haar}, \rho}; \mathcal{Y}')} \\ &\approx \langle \tilde{\Phi}(x)y, \tilde{\Phi}(x)y' \rangle_{\widetilde{\mathcal{H}}} := \langle y, \tilde{K}(x, z)y' \rangle_{\mathcal{Y}}\end{aligned}$$

where Φ_x is defined in the sense of [Proposition 4.6](#).

4.3.3 Examples of Operator Random Fourier Feature maps

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

Example 4.1 (Matrix-valued kernel on the additive group). *In the following let $K(x, z) = K_0(x - z)$ be a \mathcal{Y} -Mercer matrix-valued kernel on $\mathcal{X} = \mathbb{R}^d$, invariant w. r. t. the group operation $+$. Then the function $\tilde{\Phi}$ defined as follow is an Operator-valued Random Fourier Feature of K_0 .*

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle_2 B(\omega_j)^* y \\ \sin \langle x, \omega_j \rangle_2 B(\omega_j)^* y \end{pmatrix}, \quad \omega_j \sim \Pr_{\text{Haar}, \rho} \text{ i. i. d..}$$

for all $y \in \mathcal{Y}$.

Proof The (Pontryagin) dual of $\mathcal{X} = \mathbb{R}^d$ is $\widehat{\mathcal{X}} \cong \mathbb{R}^d$, and the duality pairing is $(x - z, \omega) = \exp(i\langle x - z, \omega \rangle_2)$. The kernel approximation yields

$$\begin{aligned}\tilde{K}(x, z) &= \tilde{\Phi}(x)^*\tilde{\Phi}(z) \\ &= \frac{1}{D} \sum_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle_2 & \sin \langle x, \omega_j \rangle_2 \end{pmatrix} \begin{pmatrix} \cos \langle z, \omega_j \rangle_2 \\ \sin \langle z, \omega_j \rangle_2 \end{pmatrix} A(\omega_j) \\ &= \frac{1}{D} \sum_{j=1}^D \cos \langle x - z, \omega_j \rangle_2 A(\omega_j) \\ &\xrightarrow[D \rightarrow \infty]{a.s.} \mathbf{E}_{\rho} [\cos \langle x - z, \omega \rangle_2 A(\omega)]\end{aligned}$$

in the weak operator topology. Since for all $x \in \mathcal{X}$, $\sin \langle x, \cdot \rangle_2$ is an odd function and $A(\cdot)\rho(\cdot)$ is even,

$$\mathbf{E}_{\rho} [\cos \langle x - z, \omega \rangle_2 A(\omega)] = \mathbf{E}_{\rho} [\exp(-i\langle x - z, \omega \rangle_2) A(\omega)] = K(x, z).$$

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

In particular we deduce the following feature maps for the kernels proposed in [Subsection 4.2.2](#).

- For the decomposable Gaussian kernel $K_0^{\text{dec}, \text{gauss}}(\delta) = k_0^{\text{gauss}}(\delta)\Gamma$ for all $\delta \in \mathbb{R}^d$, let $BB^* = \Gamma$. A bounded –and unbounded– ORFF map is

$$\begin{aligned}\tilde{\Phi}(x)y &= \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle_2 B^* y \\ \sin \langle x, \omega_j \rangle_2 B^* y \end{pmatrix} \\ &= (\tilde{\varphi}(x) \otimes B^*)y,\end{aligned}$$

where $\omega_j \sim \mathbf{Pr}_{\mathcal{N}(0, \sigma^{-2} I_d)}$ i. i. d. and $\tilde{\varphi}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle_2 \\ \sin \langle x, \omega_j \rangle_2 \end{pmatrix}$ is a scalar RFF map [139].

- For the curl-free Gaussian kernel, $K_0^{\text{curl}, \text{gauss}} = -\nabla \nabla^\top k_0^{\text{gauss}}$ an unbounded ORFF map is

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle_2 \omega_j^\top y \\ \sin \langle x, \omega_j \rangle_2 \omega_j^\top y \end{pmatrix}, \quad (4.20)$$

$\omega_j \sim \mathbf{Pr}_{\mathcal{N}(0, \sigma^{-2} I_d)}$ i. i. d. and a bounded ORFF map is

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle_2 \frac{\omega_j^\top}{\|\omega_j\|} y \\ \sin \langle x, \omega_j \rangle_2 \frac{\omega_j^\top}{\|\omega_j\|} y \end{pmatrix}, \quad \omega_j \sim \mathbf{Pr}_\rho \text{ i. i. d.}$$

where $\rho(\omega) = \frac{\sigma^2 \|\omega\|^2}{d} \mathcal{N}(0, \sigma^{-2} I_d)(\omega)$ for all $\omega \in \mathbb{R}^d$.

- For the divergence-free Gaussian kernel $K_0^{\text{div}, \text{gauss}}(x, z) = (\nabla \nabla^\top - \Delta I_d) k_0^{\text{gauss}}(x, z)$ an unbounded ORFF map is

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle_2 B(\omega_j)^\top y \\ \sin \langle x, \omega_j \rangle_2 B(\omega_j)^\top y \end{pmatrix} \quad (4.21)$$

where $\omega_j \sim \mathbf{Pr}_\rho$ i. i. d. and $B(\omega) = (\|\omega\| I_d - \omega \omega^\top)$ and $\rho = \mathcal{N}(0, \sigma^{-2} I_d)$ for all $\omega \in \mathbb{R}^d$. A bounded ORFF map is

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle_2 B(\omega_j)^\top y \\ \sin \langle x, \omega_j \rangle_2 B(\omega_j)^\top y \end{pmatrix}, \quad \omega_j \sim \mathbf{Pr}_\rho \text{ i. i. d.},$$

where $B(\omega) = \left(I_d - \frac{\omega \omega^\top}{\|\omega\|^2} \right)$ and $\rho(\omega) = \frac{\sigma^2 \|\omega\|^2}{d} \mathcal{N}(0, \sigma^{-2} I_d)$ for all $\omega \in \mathbb{R}^d$.

The second example extends scalar-valued Random Fourier Features on the skewed multiplicative group –described in [Subsection 3.2.4](#) and [Subsubsection 4.2.2.2](#)– to the operator-valued case.

Example 4.2 (Matrix-valued kernel on the skewed multiplicative group). In the following, $K(x, z) = K_{1-c}(x \odot z^{-1})$ is a \mathcal{Y} -Mercer matrix-valued kernel on $\mathcal{X} = (-c; +\infty)^d$ invariant w.r.t. the group operation¹² \odot . Then the function $\tilde{\Phi}$ defined as follow is an Operator-valued Random Fourier Feature of K_{1-c} .

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle \log(x+c), \omega_j \rangle_2 B(\omega_j)^* y \\ \sin \langle \log(x+c), \omega_j \rangle_2 B(\omega_j)^* y \end{pmatrix},$$

$\omega_j \sim \widehat{\text{Pr}_{\text{Haar}, \rho}}$ i.i.d., for all $y \in \mathcal{Y}$.

Proof The dual of $\mathcal{X} = (-c; +\infty)^d$ is $\widehat{\mathcal{X}} \cong \mathbb{R}^d$, and the duality pairing is $(x \odot z^{-1}, \omega) = \exp(i \langle \log(x \odot z^{-1} + c), \omega \rangle_2)$. Following the proof of Example 4.1, we have

$$\tilde{K}(x, z) = \frac{1}{D} \sum_{j=1}^D \cos \left\langle \log \left(\frac{x+c}{z+c} \right), \omega_j \right\rangle_2 A(\omega_j).$$

which converges almost surely to

$$E_\rho [\exp(-i \langle \log(x \odot z^{-1} + c), \omega \rangle_2) A(\omega)] = E_\rho [\overline{(x \odot z^{-1}, \omega)} A(\omega)] = K(x, z)$$

when D tends to infinity, in the weak operator topology. \square

- For the skewed- χ^2 decomposable kernel defined as $K_{1-c}^{\text{dec, skewed}}(\delta) = k_{1-c}^{\text{skewed}}(\delta)\Gamma$ for all $\delta \in \mathcal{X}$, let $BB^* = \Gamma$. A bounded –and unbounded– ORFF map is

$$\begin{aligned} \tilde{\Phi}(x)y &= \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle \log(x+c), \omega_j \rangle_2 B^* y \\ \sin \langle \log(x+c), \omega_j \rangle_2 B^* y \end{pmatrix}, \quad \omega_j \sim \text{Pr}_\rho \text{ i.i.d.} \\ &= (\tilde{\Phi}(x) \otimes B^*)y, \end{aligned}$$

where $\rho = \mathcal{S}(0, 2^{-1})$ and $\tilde{\Phi}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle \log(x+c), \omega_j \rangle_2 \\ \sin \langle \log(x+c), \omega_j \rangle_2 \end{pmatrix}$ is a scalar RFF map [98].

4.3.4 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 $L^2(\widehat{\mathcal{X}}, \widehat{\text{Haar}}; \mathcal{Y}')$ induced by $\|\cdot\|_K$?

¹² The group operation \odot is defined in Subsubsection 4.2.2.2.

Proposition 4.12 *Let K be a shift-invariant \mathcal{Y} -Mercer Kernel such that for all y, y' in \mathcal{Y} , $\langle y', K_e(\cdot)y \rangle_y \in L^1(\mathcal{X}, \text{Haar})$. Then for all $f \in \mathcal{H}_K$*

$$\|f\|_K^2 = \int_{\widehat{\mathcal{X}}} \frac{\langle \mathcal{F}[f](\omega), A(\omega)^\dagger \mathcal{F}[f](\omega) \rangle_y}{\rho(\omega)} d\widehat{\text{Haar}}(\omega). \quad (4.22)$$

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

Proof We first show how the Fourier Transform relates to the feature operator. Since \mathcal{H}_K is embedded into $\mathcal{H} = L^2(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho}; \mathcal{Y}')$ by means of the feature operator W , we have for all $f \in \mathcal{H}_K$, for all $f \in \mathcal{H}$ and for all $x \in \mathcal{X}$

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

By injectivity of the Fourier Transform, $\mathcal{F}^{-1}[f](\omega) = \rho(\omega)B(\omega)g(\omega)$. From Proposition 3.4 we have

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

The pseudo inverse of the operator $B(\omega)$ – noted $B(\omega)^\dagger$ – is the unique solution of the system $\mathcal{F}^{-1}[f](\omega) = \rho(\omega)B(\omega)g(\omega)$ w.r.t. $g(\omega)$ with minimal norm¹. Eventually,

$$\|f\|_K^2 = \int_{\widehat{\mathcal{X}}} \frac{\|B(\omega)^\dagger \mathcal{F}^{-1}[f](\omega)\|_{\mathcal{Y}'}^2}{\rho(\omega)^2} d\Pr_{\widehat{\text{Haar}}, \rho}(\omega)$$

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

$$\begin{aligned} \|f\|_K^2 &= \int_{\widehat{\mathcal{X}}} \frac{\|\mathcal{R}[B(\cdot)^\dagger \rho(\cdot)](\omega) \mathcal{F}[f](\omega)\|_{\mathcal{Y}'}^2}{\rho(\omega)^2} d\widehat{\text{Haar}}(\omega) \\ &= \int_{\widehat{\mathcal{X}}} \frac{\|B(\omega)^\dagger \rho(\omega) \mathcal{F}[f](\omega)\|_{\mathcal{Y}'}^2}{\rho(\omega)^2} d\widehat{\text{Haar}}(\omega) \\ &= \int_{\widehat{\mathcal{X}}} \frac{\langle B(\omega)^\dagger \mathcal{F}[f](\omega), B(\omega)^\dagger \mathcal{F}[f](\omega) \rangle_{\mathcal{Y}'} d\widehat{\text{Haar}}(\omega)}{\rho(\omega)} \\ &= \int_{\widehat{\mathcal{X}}} \frac{\langle \mathcal{F}[f](\omega), A(\omega)^\dagger \mathcal{F}[f](\omega) \rangle_{\mathcal{Y}'} d\widehat{\text{Haar}}(\omega)}{\rho(\omega)}. \end{aligned}$$

□

¹ Note that since $B(\omega)$ is bounded the pseudo inverse of $B(\omega)$ is well defined for $\widehat{\text{Haar}}$ -almost all ω

Note that if $K(x, z) = k(x, z)$ is a scalar kernel then for all ω in $\widehat{\mathcal{X}}$, $A(\omega) = 1$. Therefore we recover the well known result for kernels that is for any $f \in \mathcal{H}_k$ we have $\|f\|_k = \int_{\widehat{\mathcal{X}}} \mathcal{F}[k_e](\omega)^{-1} \mathcal{F}[f](\omega)^2 d\widehat{\text{Haar}}(\omega)$ [162, 180, 190]. Eventually from this last equation we also recover [Proposition 3.9](#) for decomposable kernels. If $A(\omega) = \Gamma \in \mathcal{L}_+(\mathbb{R}^p)$,

$$\|f\|_k = \sum_{i,j=1}^p \left(\Gamma^\dagger \right)_{ij} \langle f_i, f_j \rangle_k \quad (4.23)$$

Algorithm 2: Construction of ORFF

Input :

- The pairing (x, ω) of the LCA group $(\mathcal{X}, *)$.
- A probability measure $\Pr_{\widehat{\text{Haar}}, \rho}$ with density ρ w. r. t. the haar measure $\widehat{\text{Haar}}$ on $\widehat{\mathcal{X}}$.
- An operator-valued function $B : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y}, \mathcal{Y}')$ such that for all $y, y' \in \mathcal{Y}, \langle y', B(\cdot)B(\cdot)^*y \rangle \in L^1(\widehat{\mathcal{X}}, \Pr_{\widehat{\text{Haar}}, \rho})$.
- D the number of features.

Output: A random feature $\tilde{\Phi}(x)$ such that $\tilde{\Phi}(x)^* \tilde{\Phi}(z) \approx K(x, z)$.

1 Draw D random vectors $(\omega_j)_{j=1}^D$ i. i. d. from the probability law

$\Pr_{\widehat{\text{Haar}}, \rho}$;

2 **return** $\begin{cases} \tilde{\Phi}(x) \in \mathcal{L}(\mathcal{Y}, \widetilde{\mathcal{H}}) & : y \mapsto \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D (x, \omega_j) B(\omega_j)^* y; \\ \tilde{\Phi}(x)^* \in \mathcal{L}(\widetilde{\mathcal{H}}, \mathcal{Y}) & : \theta \mapsto \frac{1}{\sqrt{D}} \sum_{j=1}^D (x, \omega_j) B(\omega_j) \theta_j \end{cases}$

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 $A(\omega) = B(\omega)B(\omega)^*$. Moreover, in light of [Equation 4.22](#) the regularization property of the kernel with respect to the Fourier Transform, it is also possible to define an approximate feature map of an Operator-Valued Kernel from its regularization properties in the VV-RKHS as proposed in [Algorithm 2](#).

4.4 CONCLUSIONS

We have presented a generic framework that generalizes scalar-valued RFFs presented in [Chapter 2](#). We first showed how to construct an ORFF from an OVK in the very general case when \mathcal{X} is an LCA

group and \mathcal{Y} is an infinite dimensional space. Then, conversely, how to use the regularization properties ([Proposition 4.12](#)) to construct an ORFF without defining an OVK.



5

BOUNDING THE ERROR OF THE ORFF APPROXIMATION

In this chapter we refine the bound on the OVK approximation with ORFF we first proposed in [29] and presented in [28]. It generalizes the proof technique of Rahimi and Recht [139] to OVK on LCA groups thanks to the recent results of Koltchinskii [92], Minsker [121], Sutherland and Schneider [167], and Tropp [175]. As a Bernstein bound it depends on the variance of the estimator for which we derive an “upper bound”.

Contents

5.1	Convergence with high probability of the ORFF estimator	82
5.1.1	Random Fourier Features in the scalar case and decomposable OVK	84
5.1.2	Uniform convergence of ORFF approximation on LCA groups	85
5.1.3	Dealing with infinite dimensional operators	90
5.1.4	Variance of the ORFF approximation	91
5.1.5	Application on decomposable, curl-free and divergence-free OVK	93
5.2	Conclusions	94

5.1 CONVERGENCE WITH HIGH PROBABILITY OF THE ORFF ESTIMATOR

We are now interested in a non-asymptotic analysis of the ORFF approximation of shift-invariant \mathcal{Y} -Mercer kernels on LCA group \mathcal{X} endowed with the operation group \star where \mathcal{X} is a Banach space (The more general case where \mathcal{X} is a Polish space is discussed in [Appendix A.1](#)). For a given D , we study how close is the approximation $\tilde{K}(x, z) = \tilde{\Phi}(x)^* \tilde{\Phi}(z)$ to the target kernel $K(x, z)$ for any x, z in \mathcal{X} .

If $A \in \mathcal{L}_+(\mathcal{Y})$ we denote $\|A\|_{\mathcal{Y}, \mathcal{Y}}$ its operator norm (the induced norm). For x and z in some non-empty compact $\mathcal{C} \subset \mathbb{R}^d$, we consider: $F(x \star z^{-1}) = \tilde{K}(x, z) - K(x, z)$ and study how the uniform norm

$$\|\tilde{K} - K\|_{\mathcal{C} \times \mathcal{C}} = \sup_{(x, z) \in \mathcal{C} \times \mathcal{C}} \|\tilde{K}(x, z) - K(x, z)\|_{\mathcal{Y}, \mathcal{Y}} \quad (5.1)$$

behaves according to D . All along this document we denote $\delta = x \star z^{-1}$ for all x and $z \in \mathcal{X}$. [Figure 5.1](#) empirically shows convergence of three different OVK approximations for x, z sampled from the compact $[-1, 1]^4$ and using an increasing number of sample points D . The log-log plot shows that all three kernels have the same convergence rate, up to a multiplicative factor.

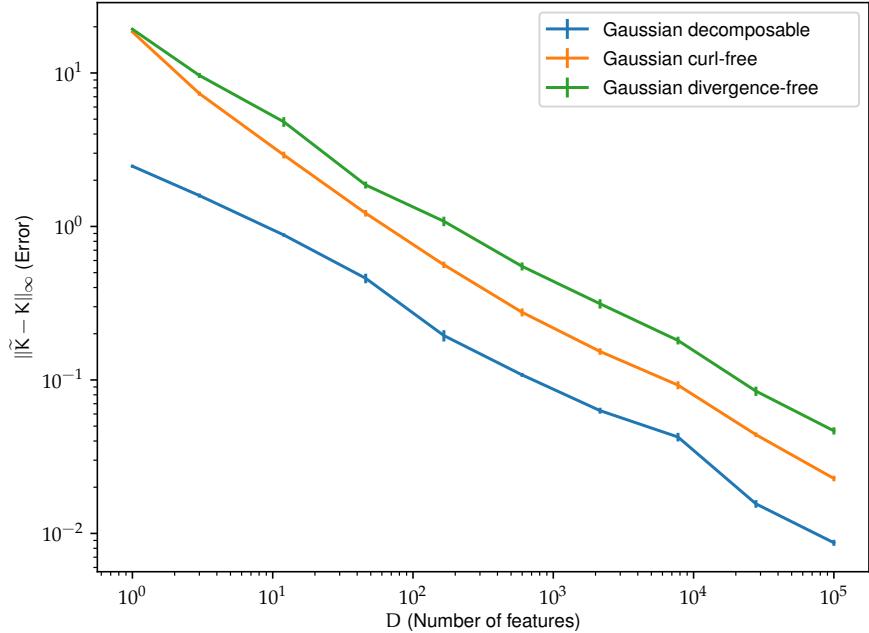


Figure 5.1: Error reconstructing the target operator-valued kernel K with ORFF approximation \tilde{K} for the decomposable, curl-free and divergence-free kernel.

In order to bound the error with high probability, we turn to concentration inequalities devoted to random matrices [25]. The concentration phenomenon can be summarized by the following sentence of Ledoux [96]. “A random variable that depends (in a smooth way) on the influence of many random variables (but not too much on any of them) is essentially constant”.

A typical application is the study of the deviation of the empirical mean of independent identically distributed random variables to their expectation. This means that given an error ϵ between the kernel approximation \tilde{K} and the true kernel K , if we are given enough samples to construct \tilde{K} , the probability of measuring an error greater than ϵ is essentially zero (it drops at an exponential rate with respect to the number of samples D). To measure the error between the kernel approximation and the true kernel at a given point, many metrics are possible. For instance, any matrix norm such as the Hilbert-Schmidt norm, trace norm, the operator norm or Schatten norms. In this work we focus on measuring the error in terms of operator norm. For all $x, z \in \mathcal{X}$ we look for a bound on

$$\begin{aligned} & \Pr_{\rho} \left\{ (\omega_j)_{j=1}^D \mid \left\| \tilde{K}(x, z) - K(x, z) \right\|_{y, y} \geq \epsilon \right\} \\ &= \Pr_{\rho} \left\{ (\omega_j)_{j=1}^D \mid \sup_{0 \neq y \in \mathcal{Y}} \frac{\left\| (\tilde{K}(x, z) - K(x, z))y \right\|_y}{\|y\|_y} \geq \epsilon \right\} \end{aligned}$$

In other words, given any vector $y \in \mathcal{Y}$ we study how the residual operator $\tilde{K} - K$ is able to send y to zero. We believe that this way of measuring the “error” to be more intuitive. Moreover, on contrary to an error measure with the Hilbert-Schmidt norm, the operator norm error does not grow linearly with the dimension of the output space as the Hilbert-Schmidt norm does. On the other hand the Hilbert-Schmidt norm makes the studied random variables Hilbert space valued, for which it is much easier to derive concentration inequalities [125, 136, 161]. Note that in the scalar case ($A(\omega) = 1$) the Hilbert-Schmidt norm error and the operator norm are the same and measure the deviation between \tilde{K} and K as the absolute value of their difference.

A raw concentration inequality of the kernel estimator gives the error on one point. If one is interesting in bounding the maximum error over N points, applying a union bound on all the point would yield a bound that grows linearly with N . This would suggest that when the number of points increase, even if all of them are concentrated in a small subset of \mathcal{X} , we should draw increasingly more features to have an error below ϵ with high probability. However if we restrict ourselves to study the error on a compact subset of \mathcal{X} (and in practice data points lies often in a closed bounded subset of \mathbb{R}^d), we can cover

this compact subset by a finite number of closed balls and apply the concentration inequality and the union bound only on the center of each ball. Then if the function $\|\tilde{K}_e - K_e\|$ is smooth enough on each ball (i.e. Lipschitz) we can guarantee with high probability that the error between the centers of the balls will not be too high. Eventually we obtain a bound in the worst case scenario on all the points in a subset \mathcal{C} of \mathcal{X} . This bound depends on the covering number $\mathcal{N}(\mathcal{C}, r)$ of \mathcal{X} with ball of radius r . When \mathcal{X} is a Banach space, the covering number is proportional to the diameter of $\mathcal{C} \subseteq \mathcal{X}$.

Prior to the presentation of general results, we briefly recall the uniform convergence of RFF approximation for a scalar shift invariant kernel on the additive LCA group \mathbb{R}^d and introduce a direct corollary about decomposable shift-invariant OVK on the LCA group $(\mathbb{R}^d, +)$.

5.1.1 Random Fourier Features in the scalar case and decomposable OVK

Rahimi and Recht [139] proved the uniform convergence of Random Fourier Feature (RFF) approximation for a scalar shift-invariant kernel on the LCA group \mathbb{R}^d endowed with the group operation $\star = +$. In the case of the shift-invariant decomposable OVK, an upper bound on the error can be obtained as a direct consequence of the result in the scalar case obtained by Rahimi and Recht [139] and other authors [164, 167].

Theorem 5.1 (Uniform error bound for RFF, Rahimi and Recht [139]). *Let \mathcal{C} be a compact subset of \mathbb{R}^d of diameter $|\mathcal{C}|$. Let k be a shift invariant kernel, differentiable with a bounded second derivative and \Pr_ρ its normalized Inverse Fourier Transform such that it defines a probability measure. Let*

$$\tilde{k} = \sum_{j=1}^D \cos \langle \cdot, \omega_j \rangle \approx k(x, z) \text{ and } \sigma^2 = E_\rho \|\omega\|_2^2.$$

Then we have

$$\Pr_\rho \left\{ (\omega_j)_{j=1}^D \mid \|\tilde{k} - k\|_{\mathcal{C} \times \mathcal{C}} \geq \epsilon \right\} \leq 2^8 \left(\frac{\sigma |\mathcal{C}|}{\epsilon} \right)^2 \exp \left(-\frac{\epsilon^2 D}{4(d+2)} \right)$$

From [Theorem 5.1](#), we can deduce the following corollary about the uniform convergence of the ORFF approximation of the decomposable kernel. We recall that for a given pair x, z in \mathcal{C} , $\tilde{K}(x, z) = \tilde{\Phi}(x)^* \tilde{\Phi}(z) = \Gamma \tilde{k}(x, z)$ and $K_0(x - z) = \widehat{\Gamma E_{\text{Haar}, \rho} [\tilde{k}(x, z)]}$.

Corollary 5.1 (Uniform error bound for decomposable ORFF). *Let \mathcal{C} be a compact subset of \mathbb{R}^d of diameter $|\mathcal{C}|$. Let K be a decomposable kernel*

built from a positive operator self-adjoint Γ , and k a shift invariant kernel with bounded second derivative such that

$$\tilde{K} = \sum_{j=1}^D \cos \langle \cdot, \omega_j \rangle \Gamma \approx K \text{ and } \sigma^2 = \mathbf{E}_\rho \|\omega\|_2^2.$$

Then

$$\begin{aligned} \Pr_\rho \left\{ (\omega_j)_{j=1}^D \mid \left\| \tilde{K} - K \right\|_{\mathcal{C} \times \mathcal{C}} \geq \epsilon \right\} \\ \leq 2^8 \left(\frac{\sigma \|\Gamma\|_{\mathcal{Y}, \mathcal{Y}} |\mathcal{C}|}{\epsilon} \right)^2 \exp \left(-\frac{\epsilon^2 D}{4 \|\Gamma\|_2^2 (d+2)} \right) \end{aligned}$$

Proof The proof directly extends [Theorem 5.1](#) given by [139]. Let \tilde{k} be the Random Fourier approximation for the scalar-valued kernel k . Since

$$\sup_{(x,z) \in \mathcal{C} \times \mathcal{C}} \left\| \tilde{K}(x, z) - K(x, z) \right\|_{\mathcal{Y}, \mathcal{Y}} = \sup_{(x,z) \in \mathcal{C} \times \mathcal{C}} \|\Gamma\|_{\mathcal{Y}, \mathcal{Y}} \left| \tilde{k}(x, z) - k(x, z) \right|,$$

taking $\epsilon' = \|\Gamma\|_{\mathcal{Y}, \mathcal{Y}} \epsilon$ gives the following result for all positive ϵ' :

$$\begin{aligned} \Pr_\rho \left\{ (\omega_j)_{j=1}^D \mid \sup_{x, z \in \mathcal{C}} \left\| \Gamma \left(\tilde{k}(x, z) - k(x, z) \right) \right\|_{\mathcal{Y}, \mathcal{Y}} \geq \epsilon' \right\} \\ \leq 2^8 \left(\frac{\sigma \|\Gamma\|_{\mathcal{Y}, \mathcal{Y}} |\mathcal{C}|}{\epsilon'} \right)^2 \exp \left(-\frac{(\epsilon')^2 D}{4 \|\Gamma\|_{\mathcal{Y}, \mathcal{Y}}^2 (d+2)} \right) \end{aligned}$$

which concludes the proof. \square

Note that a similar corollary could have been obtained for the recent result of Sutherland and Schneider [167] who refined the bound proposed by Rahimi and Recht by using a Bernstein concentration inequality instead of the Hoeffding inequality. More recently Sriperumbudur and Szabo [164] showed an optimal bound for Random Fourier Feature. The improvement of Sriperumbudur and Szabo [164] is mainly in the constant factors where the bound does not depend linearly on the diameter $|\mathcal{C}|$ of \mathcal{C} but exhibit a logarithmic dependency $\log(|\mathcal{C}|)$, hence requiring significantly less random features to reach a desired uniform error with high probability. Moreover, Sutherland and Schneider [167] also considered a bound on the expected max error $\mathbf{E}_{\widehat{\text{Haar}}, \rho} \left\| \tilde{K} - K \right\|_\infty$, which is obtained using Dudley's entropy integral [25, 56] as a bound on the supremum of an empirical process by the covering number of the indexing set. This useful theorem is also part of the proof of Sriperumbudur and Szabo [164].

5.1.2 Uniform convergence of ORFF approximation on LCA groups

In this analysis, we assume that \mathcal{Y} is finite dimensional, in [Subsection 5.1.3](#), we discuss how the proof could be extended to infinite

dimensional output Hilbert spaces. We propose a bound for Operator-valued Random Fourier Feature approximation in the general case. It relies on two main ideas:

1. a matrix-Bernstein concentration inequality for random matrices need to be used instead of concentration inequality for scalar random variables,
2. a general theorem, valid for random matrices with bounded norms (such as decomposable kernel ORFF approximation) as well as unbounded norms (such as the ORFF approximation we proposed for curl and divergence-free kernels, for which the norm behave as subexponential random variables).

Before introducing the new theorem, we give the definition of the Orlicz norm which gives a proxy-bound on the norm of subexponential random variables.

Definition 5.1 (Orlicz norm [176]). Let $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ be a non-decreasing convex function with $\psi(0) = 0$. For a random variable X on a measured space $(\Omega, \mathcal{T}(\Omega), \mu)$, the quantity

$$\|X\|_\psi = \inf \{ C > 0 \mid E_\mu[\psi(|X|/C)] \leq 1 \}.$$

is called the Orlicz norm of X .

Here, the function ψ is chosen as $\psi(u) = \psi_\alpha(u)$ where $\psi_\alpha(u) := e^{u^\alpha} - 1$. When $\alpha = 1$, a random variable with finite Orlicz norm is called a *subexponential variable* because its tails decrease at an exponential rate. Let X be a self-adjoint random operator. Given a scalar-valued measure μ , we call *variance* of an operator X the quantity $\text{Var}_\mu[X] = E_\mu[X - E_\mu[X]]^2$. With this convention if X is a $p \times p$ Hermitian matrix,

$$\text{Var}_\mu[X]_{\ell m} = \sum_{r=1}^p \text{Cov}_\mu[X_{\ell r}, X_{rm}].$$

Among the possible concentration inequalities adapted to random operators [92, 97, 121, 136, 175], we focus on the results of Minsker [121] and Tropp [175], for their robustness to high or potentially infinite dimension of the output space \mathcal{Y} . To guarantee a good scaling with the dimension of \mathcal{Y} we introduce the notion of intrinsic dimension (or effective rank) of an operator.

Definition 5.2 Let A be a trace class operator acting on a Hilbert space \mathcal{Y} . We call *intrinsic dimension* the quantity

$$\text{IntDim}(A) = \frac{\text{Tr}[A]}{\|A\|_{\mathcal{Y}, \mathcal{Y}}}.$$

Indeed the bound proposed in our first publication at ACML [29] based on Koltchinskii [92] depends on p while the present bound depends on the intrinsic dimension of the variance of $A(\omega)$ which is always smaller than p when the operator $A(\omega)$ is Hilbert-Schmidt ($p \leq \infty$).

Corollary 5.2 Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ be a shift-invariant \mathcal{Y} -Mercer kernel, where \mathcal{Y} is a finite dimensional Hilbert space of dimension p and \mathcal{X} a finite dimensional Banach space of dimension d . Moreover, let \mathcal{C} be a closed ball of \mathcal{X} centred at the origin of diameter $|\mathcal{C}|$, $A : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$ and $\Pr_{\widehat{\text{Haar}}, \rho}$ a pair such that

$$\tilde{K}_e = \sum_{j=1}^D \cos(\cdot, \omega_j) A(\omega_j) \approx K_e, \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho} \text{ i.i.d.}$$

Let $\mathcal{D}_{\mathcal{C}} = \mathcal{C} \star \mathcal{C}^{-1}$ and

$$V(\delta) \succcurlyeq \text{Var}_{\widehat{\text{Haar}}, \rho} \tilde{K}_e(\delta), \quad \text{for all } \delta \in \mathcal{D}_{\mathcal{C}}$$

and H_{ω} be the Lipschitz constant of the function $h : x \mapsto (x, \omega)$. If the three following constants exist

$$m \geq \int_{\widehat{\mathcal{X}}} H_{\omega} \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} d\Pr_{\widehat{\text{Haar}}, \rho}(\omega) < \infty$$

and

$$u \geq 4 \left(\| \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} \|_{\Psi_1} + \sup_{\delta \in \mathcal{D}_e} \|K_e(\delta)\|_{\mathcal{Y}, \mathcal{Y}} \right) < \infty$$

and

$$v \geq \sup_{\delta \in \mathcal{D}_e} D \|V(\delta)\|_{\mathcal{Y}, \mathcal{Y}} < \infty.$$

Define $p_{\text{int}} \geq \sup_{\delta \in \mathcal{D}_e} \text{IntDim}(V(\delta))$, then for all $0 < \epsilon \leq m|\mathcal{C}|$,

$$\begin{aligned} & \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \|\tilde{K} - K\|_{\mathcal{C} \times \mathcal{C}} \geq \epsilon \right\} \\ & \leq 8\sqrt{2} \left(\frac{m|\mathcal{C}|}{\epsilon} \right) (p_{\text{int}} r_{v/D}(\epsilon))^{\frac{1}{d+1}} \begin{cases} \exp \left(-D \frac{\epsilon^2}{8v(d+1)(1+\frac{1}{p})} \right), & \epsilon \leq \frac{v}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8u(d+1)K(v,p)} \right), & \text{otherwise,} \end{cases} \end{aligned}$$

where $K(v, p) = \log(16\sqrt{2}p) + \log\left(\frac{u^2}{v}\right)$ and $r_{v/D}(\epsilon) = 1 + \frac{3}{\epsilon^2 \log^2(1+D\epsilon/v)}$.

Sketch of the proof In the following, let $F(\delta) = F(x \star z^{-1}) = \tilde{K}(x, z) - K(x, z)$. Let $\mathcal{D}_{\mathcal{C}} = \mathcal{C} \star \mathcal{C}^{-1} = \{x \star z^{-1} \mid x, z \in \mathcal{C}\}$. Since \mathcal{C} is supposed compact, so is $\mathcal{D}_{\mathcal{C}}$. Its diameter is at most $2|\mathcal{C}|$ where $|\mathcal{C}|$ is the diameter of \mathcal{C} . Since \mathcal{C} is supposed to be a closed ball of a Banach space it is then possible to find an ϵ -net covering $\mathcal{D}_{\mathcal{C}}$ with at most $T = (4|\mathcal{C}|/r)^d$ balls of radius r [50]. We call δ_i for $i \in \{1, \dots, T\}$ the center of the i -th ball, called anchors of the ϵ -net. Denote L_F the Lipschitz constant of F . We introduce the following lemma proved in [139].

Lemma 5.1 *For all $\delta \in \mathcal{D}_C$, if*

$$L_F \leq \frac{\epsilon}{2r} \quad (5.2)$$

and

$$\|\mathbb{F}(\delta_i)\|_{y,y} \leq \frac{\epsilon}{2}, \quad \text{for all } i \in \mathbb{N}_T^* \quad (5.3)$$

then $\|\mathbb{F}(\delta)\|_{y,y} \leq \epsilon$ for all $\delta \in \mathcal{D}_C$.

To apply the lemma, we must check assumptions Equation 5.2 and Equation 5.3.

Sketch of the proof (Equation 5.2) **Lemma 5.2** *Let $H_\omega \in \mathbb{R}_+$ be the Lipschitz constant of $h_\omega(\cdot)$ and assume that*

$$\int_{\widehat{\mathcal{X}}} H_\omega \|A(\omega)\|_{y,y} d\Pr_{\widehat{\text{Haar}},\rho}(\omega) < \infty.$$

Then the operator-valued function $K_e : \mathcal{X} \rightarrow \mathcal{L}(Y)$ is Lipschitz with

$$\|K_e(x) - K_e(z)\|_{y,y} \leq d_{\mathcal{X}}(x,z) \int_{\widehat{\mathcal{X}}} H_\omega \|A(\omega)\|_{y,y} d\Pr_{\widehat{\text{Haar}},\rho}(\omega) \quad (5.4)$$

In the same way, considering $\tilde{K}_e(\delta) = \frac{1}{D} \sum_{j=1}^D \cos h_{\omega_j}(\delta) A(\omega_j)$, where $\omega_j \sim \Pr_{\widehat{\text{Haar}},\rho}$, we can show that \tilde{K}_e is Lipschitz with

$$\|\tilde{K}_e(x) - \tilde{K}_e(z)\|_{y,y} \leq d_{\mathcal{X}}(x,z) \frac{1}{D} \sum_{j=1}^D H_{\omega_j} \|A(\omega_j)\|_{y,y}.$$

Taking the expectation yields

$$\mathbb{E}_{\widehat{\text{Haar}},\rho}[L_F] = 2 \int_{\widehat{\mathcal{X}}} H_\omega \|A(\omega)\|_{y,y} d\Pr_{\widehat{\text{Haar}},\rho}$$

Thus by Markov's inequality,

$$\begin{aligned} \Pr_{\widehat{\text{Haar}},\rho}\{(\omega_j)_{j=1}^D \mid L_F \geq \epsilon\} &\leq \frac{\mathbb{E}_{\widehat{\text{Haar}},\rho}[L_F]}{\epsilon} \\ &\leq \frac{2}{\epsilon} \int_{\widehat{\mathcal{X}}} H_\omega \|A(\omega)\|_{y,y} d\Pr_{\widehat{\text{Haar}},\rho}. \end{aligned} \quad (5.5)$$

Sketch of the proof (Equation 5.3) *To obtain a bound on the anchors we apply theorem 4 of Koltchinskii [92]. We suppose the existence of the two constants*

$$v_i = D \mathbb{V}\text{ar}_{\widehat{\text{Haar}},\rho}[\tilde{K}(\delta_i)]$$

and

$$u_i = 4 \left(\|\|A(\omega)\|_{y,y}\|_{\Psi_1} + \|K_e(\delta_i)\|_{y,y} \right)$$

Then $\forall i \in \{1, \dots, T\}$,

$$\Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \|F(\delta_i)\|_{y,y} \geq \epsilon \right\} \\ \leq \begin{cases} 4 \text{IntDim}(v_i) \exp \left(-D \frac{\epsilon^2}{2\|v_i\|_{y,y} (1+\frac{1}{p})} \right) r_{v_i/D}(\epsilon), & \epsilon \leq \frac{\|v_i\|_{y,y}}{2u_i} \frac{1+1/p}{K(v_i,p)} \\ 4 \text{IntDim}(v_i) \exp \left(-D \frac{\epsilon}{4u_i K(v_i,p)} \right) r_{v_i/D}(\epsilon), & \text{otherwise.} \end{cases}$$

where

$$K(v_i, p) = \log \left(16\sqrt{2}p \right) + \log \left(\frac{u_i^2}{\|v_i\|_{y,y}} \right)$$

and

$$r_{v_i/D} = 1 + \frac{3}{\epsilon^2 \log^2(1 + D\epsilon/\|v_i\|_{y,y})}.$$

Combining [Equation 5.2](#) and [Equation 5.3](#). Now applying the lemma and taking the union bound over the centers of the ϵ -net yields

$$\Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \|\tilde{K} - K\|_{C \times C} \geq \epsilon \right\} \\ \leq 4 \left(\frac{rm}{\epsilon} + p_{\text{int}} \left(\frac{2|C|}{r} \right)^d r_{v/D}(\epsilon) \right. \\ \left. \begin{cases} \exp \left(-D \frac{\epsilon^2}{8v(1+\frac{1}{p})} \right), & \epsilon \leq \frac{v}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8u K(v,p)} \right), & \text{otherwise.} \end{cases} \right)$$

The right hand side of the equation has the form $ar + br^{-d}$ with

$$a = \frac{m}{\epsilon}$$

and

$$b = p_{\text{int}}(2|C|)^d r_{v/D}(\epsilon) \begin{cases} \exp \left(-D \frac{\epsilon^2}{8v(1+\frac{1}{p})} \right), & \epsilon \leq \frac{v}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8u K(v,p)} \right), & \text{otherwise.} \end{cases}$$

Following [\[118, 139, 167\]](#), we optimize over r . It is a convex continuous function on \mathbb{R}_+ and achieve the minimum value

$$r_* = a^{\frac{d}{d+1}} b^{\frac{1}{d+1}} \left(d^{\frac{1}{d+1}} + d^{-\frac{d}{d+1}} \right),$$

hence

$$\Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \|\tilde{K} - K\|_{C \times C} \geq \epsilon \right\} \\ \leq 8\sqrt{2} \left(\frac{m|C|}{\epsilon} \right) (p_{\text{int}} r_{v/D}(\epsilon))^{\frac{1}{d+1}} \begin{cases} \exp \left(-D \frac{\epsilon^2}{8v(d+1)(1+\frac{1}{p})} \right), & \epsilon \leq \frac{v}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8u(d+1)K(v,p)} \right), & \text{otherwise,} \end{cases}$$

which concludes the sketch of the proof. \square

We give a comprehensive full proof of the theorem in [Appendix A.1](#). It follows the usual scheme derived in Rahimi and Recht [[139](#)] and Sutherland and Schneider [[167](#)] and involves Bernstein concentration inequality for unbounded symmetric matrices ([Theorem A.3](#)).

5.1.3 Dealing with infinite dimensional operators

We studied the concentration of ORFFs under the assumption that \mathcal{Y} is finite dimensional. Indeed a d term characterizing the dimension of the input space \mathcal{X} appears in the bound proposed in [Corollary 5.2](#), and when d tends to infinity, the exponential part goes to zero so that the probability is bounded by a constant greater than one. Unfortunately, considering unbounded random operators [[121](#)] does not give any tighter solution.

In our first bound presented at ACML, we presented a bound based on a matrix concentration inequality for unbounded random variable. Compared to this previous bound, [Corollary 5.2](#) does not depend on the dimensionality p of the output space \mathcal{Y} but on the intrinsic dimension of the operator $A(\omega)$. However to remove the dependency in p in the exponential part, we must turn our attention to operator concentration inequalities for bounded random variable. To the best of our knowledge we are not aware of concentration inequalities working for “unbounded” operator- valued random variables acting on infinite dimensional spaces. Following the same proof than [Corollary 5.2](#) we obtain

Corollary 5.3 *Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ be a shift-invariant \mathcal{Y} -Mercer kernel, where \mathcal{Y} is a Hilbert space and \mathcal{X} a finite dimensional Banach space of dimension D . Moreover, let \mathcal{C} be a closed ball of \mathcal{X} centered at the origin of diameter $|\mathcal{C}|$, subset of \mathcal{X} , $A : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$ and $\Pr_{\widehat{\text{Haar}}, p}$ a pair such that*

$$\tilde{K}_e = \sum_{j=1}^D \cos(\cdot, \omega_j) A(\omega_j) \approx K_e, \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, p} \text{ i.i.d.}$$

where $A(\omega_j)$ is a Hilbert-Schmidt operator for all $j \in \mathbb{N}_D^*$. Let $\mathcal{D}_{\mathcal{C}} = \mathcal{C} * \mathcal{C}^{-1}$ and

$$V(\delta) \succcurlyeq \text{Var}_{\widehat{\text{Haar}}, p} \tilde{K}_e(\delta), \quad \text{for all } \delta \in \mathcal{D}_{\mathcal{C}}$$

and H_ω be the Lipschitz constant of the function $h : x \mapsto (x, \omega)$. If the three following constants exists

$$m \geq \int_{\widehat{\mathcal{X}}} H_\omega \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} d\Pr_{\widehat{\text{Haar}}, p}(\omega) < \infty$$

and

$$u \geq \text{ess sup}_{\omega \in \widehat{\mathcal{X}}} \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} + \sup_{\delta \in \mathcal{D}_{\mathcal{C}}} \|K_e(\delta)\|_{\mathcal{Y}, \mathcal{Y}} < \infty$$

and

$$v \geq \sup_{\delta \in \mathcal{D}_e} D \|V(\delta)\|_{\mathcal{Y}, \mathcal{Y}} < \infty.$$

define $p_{int} \geq \sup_{\delta \in \mathcal{D}_e} IntDim(V(\delta))$ then for all $\sqrt{\frac{v}{D}} + \frac{u}{3D} < \epsilon < m|\mathcal{C}|$,

$$\begin{aligned} & \Pr_{\widehat{\text{Haar}}, p} \left\{ (\omega_j)_{j=1}^D \mid \sup_{\delta \in \mathcal{D}_e} \|F(\delta)\|_{\mathcal{Y}, \mathcal{Y}} \geq \epsilon \right\} \\ & \leq 8\sqrt{2} \left(\frac{m|\mathcal{C}|}{\epsilon} \right) p_{int}^{\frac{1}{d+1}} \exp(-D\psi_{v,d,u}(\epsilon)) \end{aligned}$$

where $\psi_{v,d,u}(\epsilon) = \frac{\epsilon^2}{2(d+1)(v+u\epsilon/3)}$.

Again a full comprehensive proof is given in [Appendix A.1](#) of the appendix. Notice that in this result, The dimension $p = \dim \mathcal{Y}$ does not appear. Only the intrinsic dimension of the variance of the estimator. Moreover when d is large, the term $p_{int}^{\frac{1}{d+1}}$ goes to one, so that the impact of the intrinsic dimension on the bound vanish when the dimension of the input space is large.

5.1.4 Variance of the ORFF approximation

We now provide a bound on the norm of the variance of \tilde{K} , required to apply [Corollaries 5.2](#) and [5.3](#). This is an extension of the proof of Sutherland and Schneider [167] to the operator-valued case, and we recover their results in the scalar case when $A(\omega) = 1$. An illustration of the bound is provided in [Figures 5.2](#) and [5.3](#) for the decomposable and the curl-free OVK.

Proposition 5.1 (Bounding the variance of \tilde{K}). *Let K be a shift invariant \mathcal{Y} -Mercer kernel on a second countable LCA topological space \mathcal{X} . Let $A : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$ and $\Pr_{\widehat{\text{Haar}}, p}$ a pair such that*

$$\tilde{K}_e = \sum_{j=1}^D \cos(\cdot, \omega_j) A(\omega_j) \approx K_e, \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, p} \text{ i.i.d.}$$

Then,

$$\begin{aligned} \text{Var}_{\widehat{\text{Haar}}, p} [\tilde{K}_e(\delta)] & \leq \frac{1}{2D} \left((K_e(2\delta) + K_e(e)) E_{\widehat{\text{Haar}}, p} [A(\omega)] \right. \\ & \quad \left. - 2K_e(\delta)^2 + \text{Var}_{\widehat{\text{Haar}}, p} [A(\omega)] \right) \end{aligned}$$

Proof It relies on the i.i.d. property of the random vectors ω_j and trigonometric identities (see the proof in [Proposition A.3](#) of the appendix). \square

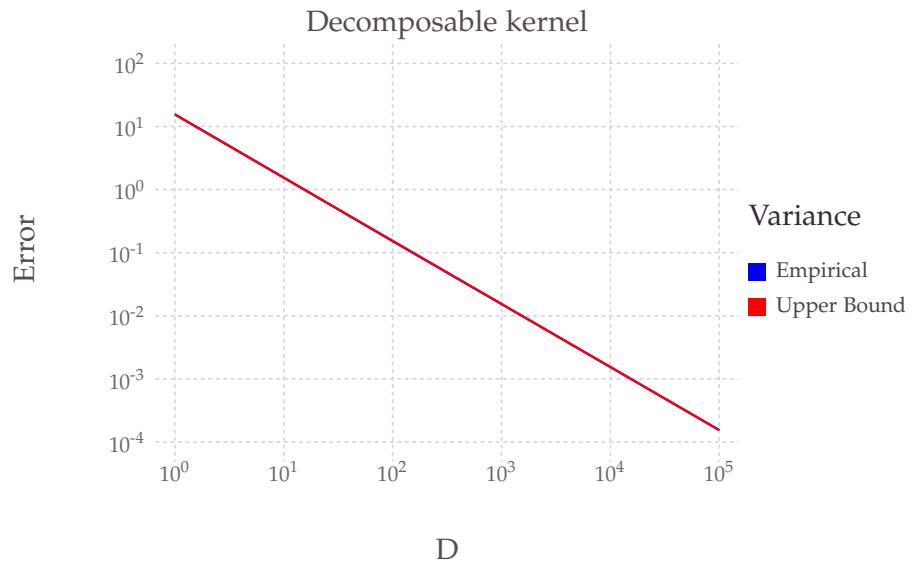


Figure 5.2: Comparison between an empirical bound on the norm of the variance of the decomposable ORFF obtained and the theoretical bound proposed in [Proposition A.3](#) versus D

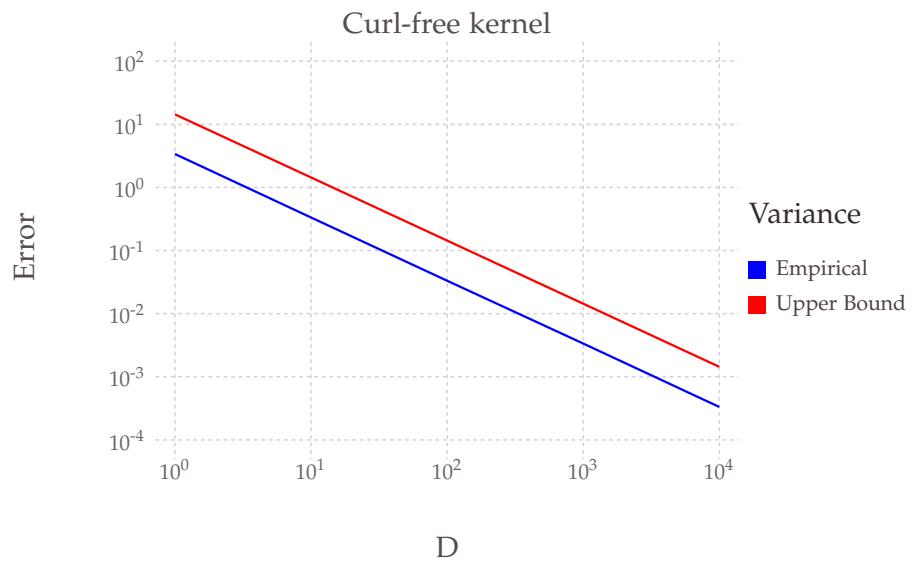


Figure 5.3: Comparison between an empirical bound on the norm of the variance of the curl-free ORFF obtained and the theoretical bound proposed in [Proposition A.3](#) versus D

5.1.5 Application on decomposable, curl-free and divergence-free OVK

First, the two following examples discuss the form of H_ω for the additive group and the skewed-multiplicative group. Here we view $\mathcal{X} = \mathbb{R}^d$ as a Banach space endowed with the Euclidean norm. Thus the Lipschitz constant H_ω is bounded by the supremum of the norm of the gradient of h_ω .

Example 5.1 (Additive group). *On the additive group, $h_\omega(\delta) = \langle \omega, \delta \rangle$. Hence $H_\omega = \|\omega\|_2$.*

Example 5.2 (Skewed-multiplicative group). *On the skewed multiplicative group, $h_\omega(\delta) = \langle \omega, \log(\delta + c) \rangle$. Therefore*

$$\sup_{\delta \in \mathcal{C}} \|\nabla h_\omega(\delta)\|_2 = \sup_{\delta \in \mathcal{C}} \|\omega / (\delta + c)\|_2.$$

Eventually \mathcal{C} is compact subset of \mathcal{X} and finite dimensional thus \mathcal{C} is closed and bounded. Thus $H_\omega = \|\omega\|_2 / (\min_{\delta \in \mathcal{C}} \|\delta\|_2 + c)$.

Now we compute upper bounds on the norm of the variance and Orlicz norm of the three ORFFs we took as examples.

5.1.5.1 Decomposable kernel

Notice that in the case of the Gaussian decomposable kernel, i.e. $A(\omega) = A$, $e = 0$, $K_0(\delta) = Ak_0(\delta)$, $k_0(\delta) \geq 0$ and $k_0(\delta) = 1$, then we have

$$D\|\mathbf{Var}_\mu [\tilde{K}_0(\delta)]\|_{y,y} \leq (1 + k_0(2\delta))\|A\|_{y,y}/2 + k_0(\delta)^2.$$

5.1.5.2 Curl-free and divergence-free kernels:

recall that in this case $p = d$. For the (Gaussian) curl-free kernel, $A(\omega) = \omega\omega^*$ where $\omega \in \mathbb{R}^d \sim \mathcal{N}(0, \sigma^{-2}I_d)$ thus $E_\mu[A(\omega)] = I_d/\sigma^2$ and $\mathbf{Var}_\mu[A(\omega)] = (d+1)I_d/\sigma^4$. Hence,

$$D\|\mathbf{Var}_\mu [\tilde{K}_0(\delta)]\|_{y,y} \leq \frac{1}{2} \left\| \frac{1}{\sigma^2} K_0(2\delta) - 2K_0(\delta)^2 \right\|_{y,y} + \frac{(d+1)}{\sigma^4}.$$

This bound is illustrated by [Figure 5.1](#) B, for a given datapoint. Eventually for the Gaussian divergence-free kernel, $A(\omega) = I\|\omega\|_2^2 - \omega\omega^*$, thus $E_\mu[A(\omega)] = I_d(d-1)/\sigma^2$ and $\mathbf{Var}_\mu[A(\omega)] = d(4d-3)I_d/\sigma^4$. Hence,

$$D\|\mathbf{Var}_\mu [\tilde{K}_0(\delta)]\|_{y,y} \leq \frac{1}{2} \left\| \frac{(d-1)}{\sigma^2} K_0(2\delta) - 2K_0(\delta)^2 \right\|_{y,y} + \frac{d(4d-3)}{\sigma^4}.$$

To conclude, we ensure that the random variable $\|A(\omega)\|_{y,y}$ has a finite Orlicz norm with $\psi = \psi_1$ in these three cases.

5.1.5.3 Computing the Orlicz norm

For a random variable with strictly monotonic moment generating function (MGF), one can characterize its inverse ψ_1 Orlicz norm by taking the functional inverse of the MGF evaluated at 2 (see Lemma A.3 of the appendix). In other words $\|X\|_{\psi_1}^{-1} = \text{MGF}(x)_X^{-1}(2)$. For the Gaussian curl-free and divergence-free kernel,

$$\|\mathcal{A}^{\text{div}}(\omega)\|_{y,y} = \|\mathcal{A}^{\text{curl}}(\omega)\|_{y,y} = \|\omega\|_2^2,$$

where $\omega \sim \mathcal{N}(0, I_d/\sigma^2)$, hence $\|\mathcal{A}(\omega)\|_2 \sim \Gamma(p/2, 2/\sigma^2)$. The MGF of this gamma distribution is $\text{MGF}(x)(t) = (1 - 2t/\sigma^2)^{-(p/2)}$. Eventually

$$\|\|\mathcal{A}^{\text{div}}(\omega)\|_{y,y}\|_{\psi_1}^{-1} = \|\|\mathcal{A}^{\text{curl}}(\omega)\|_{y,y}\|_{\psi_1}^{-1} = \frac{\sigma^2}{2} \left(1 - 4^{-\frac{1}{p}}\right).$$

5.2 CONCLUSIONS

In this chapter we have seen how to bound $\|\tilde{K} - K\|$ in the operator norm with high probability (Section 5.1). We studied the case of unbounded finite dimensional OVKs and bounded potentially infinite dimensional OVKs. The current lack of concentration inequalities working for both unbounded and infinite dimensional with the operator norm (Banach space) in the literature prevents us to unify these bounds.



6

LEARNING WITH FEATURE MAPS

This contribution chapter focuses on explaining how to define an efficient implementation and algorithm to train an ORFF model. First we recall the supervised ridge regression with OVK and the celebrated representer theorem [182]. Then we show under which condition learning with an ORFF is equivalent to learning with a kernel approximation. Eventually we give the gradient for the Ridge regression problem, in order to find an optimal solution with gradient descent algorithms, as well as a closed form algorithm. We conclude by showing how viewing ORFFs as linear operators rather than matrices yields a more efficient implementation and finish with some numerical applications on toy and real-world datasets.

Contents

6.1	Learning with OVK	96
6.1.1	Supervised learning within VV-RKHS	96
6.1.2	Learning with Operator Random Fourier Feature maps	100
6.2	Solving ORFF-based regression	103
6.2.1	Gradient descent methods	106
6.2.2	Complexity analysis	108
6.3	Efficient learning with ORFF	109
6.3.1	Case of study: the decosubmposable kernel	110
6.3.2	Linear operators in matrix form	113
6.3.3	Curl-free kernel	117
6.3.4	Divergence-free kernel	117
6.3.5	Iterative solvers using linear operators	118
6.4	Experiments	118
6.4.1	Learning with ORFF vs learning with OVK	118
6.5	Conclusion	124

6.1 LEARNING WITH OVK

Before focusing on learning function with an ORFF model, we briefly review the context of supervised learning in VV-RKHS.

6.1.1 Supervised learning within VV-RKHS

6.1.1.1 General results

Let $\mathbf{s} = (x_i, y_i)_{i=1}^N \in (\mathcal{X} \times \mathcal{Y})^N$ be a sequence of training samples. Given a local loss function $L : \mathcal{X} \times \mathcal{F} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ such that L is proper, convex and lower semi-continuous in \mathcal{F} , we are interested in finding a *vector-valued function* $f_s : \mathcal{X} \rightarrow \mathcal{Y}$, that lives in a VV-RKHS and minimizes a tradeoff between a data fitting term and a regularization term to prevent from overfitting. Namely finding $f_s \in \mathcal{H}_K$ such that

$$f_s = \arg \min_{f \in \mathcal{H}_K} \frac{1}{N} \sum_{i=1}^N L(x_i, f, y_i) + \frac{\lambda}{2} \|f\|_K^2 \quad (6.1)$$

¹³ *Tychonov regularization.* where $\lambda \in \mathbb{R}_+$ is a regularization¹³ parameter. We recall that the quantity

$$\mathfrak{R}_{\text{emp}}(f, \mathbf{s}) = \frac{1}{N} \sum_{i=1}^N L(x_i, f, y_i), \quad \forall f \in \mathcal{H}_K, \forall \mathbf{s} \in (\mathcal{X} \times \mathcal{Y})^N.$$

is called the empirical risk of the model $f \in \mathcal{H}_K$ according the local loss L . A common choice for L is the quadratic loss $L : (x, f, y) \mapsto \|f(x) - y\|_y^2$. We introduce a corollary from Mazur and Schauder proposed in 1936 showing that [Equation 6.1](#) –and [Equation 6.3](#)– attains a unique mimimizer (see Górniewicz [74] and Kurdila and Zabarankin [93]).

Theorem 6.1 (Mazur-Schauder). *Let \mathcal{H} be a Hilbert space and $\mathfrak{R} : \mathcal{H} \rightarrow \overline{\mathbb{R}}$ be a proper, convex, lower semi-continuous and coercive function. Then \mathfrak{R} is bounded from below and attains a minimizer. Moreover if \mathfrak{R} is strictly convex the minimizer is unique.*

This is easily verified for Ridge regression. Define

$$\mathfrak{R}_\lambda(f, \mathbf{s}) = \frac{1}{N} \sum_{i=1}^N \|f(x_i) - y_i\|_y^2 + \frac{\lambda}{2} \|f\|_K^2, \quad (6.2)$$

¹⁴ *Reminder, if $f \in \mathcal{H}_K$, $ev_x : f \mapsto f(x)$ is continuous, see [Proposition 3.2](#).* where $f \in \mathcal{H}_K$ and $\lambda \in \mathbb{R}_{>0}$. \mathfrak{R}_λ is continuous¹⁴ and strictly convex. Additionally \mathfrak{R}_λ is coercive since $\|f\|_K$ is coercive, $\lambda \in \mathbb{R}_{>0}$, and all the summands of \mathfrak{R}_λ are positive. Hence for all positive λ , $f_s = \arg \min_{f \in \mathcal{H}_K} \mathfrak{R}_\lambda(f)$ exists, is unique and attained.

Remark 6.1 (Kadri et al. [86]). We consider the optimization problem proposed in [Equation 6.2](#) where $L : (x_i, f, y_i) \mapsto \|f(x_i) - y_i\|_y^2$. If given a training sample s , we have

$$\frac{1}{N} \sum_{i=1}^N \|y_i\|_y^2 \leq \sigma_y^2,$$

then $\lambda \|f_s\|_K \leq 2\sigma_y^2$. Indeed, since \mathcal{H}_K is a Hilbert space, $0 \in \mathcal{H}_K$, thus

$$\begin{aligned} \frac{\lambda}{2} \|f_s\|_K^2 &\leq \frac{1}{N} \sum_{i=1}^N L(x_i, f_s, y_i) + \frac{\lambda}{2} \|f_s\|_K^2 \\ &\leq \frac{1}{N} \sum_{i=1}^N L(x_i, 0, y_i) \leq \sigma_y^2, \quad \text{by optimality of } f_s. \end{aligned}$$

Since for all $x \in \mathcal{X}$, $\|f(x)\|_y \leq \sqrt{\|\mathbf{K}(x, x)\|_{y,y}} \|f\|_K$, the maximum value that the solution $\|f_s(x)\|_y$ of [Equation 6.2](#) can reach is $\sigma_y \sqrt{\frac{2\|\mathbf{K}(x, x)\|_{y,y}}{\lambda}}$. Thus when solving a Ridge regression problem, given a shift-invariant kernel K_e , one should choose

$$0 < \lambda \leq 2\|K_e(e)\|_{y,y} \frac{\sigma_y^2}{C^2}.$$

with $C \in \mathbb{R}_{>0}$ to have a chance to fit all the y_i with norm $\|y_i\|_y \leq C$ in the train set.

6.1.1.2 Representer theorem and Feature equivalence

Regression in Vector Valued Reproducing Kernel Hilbert Space has been well studied [5, 8, 35, 86, 113, 116, 120, 146], and a cornerstone of learning in VV-RKHS is the representer theorem¹⁵, which allows to replace the search of a minimizer in a infinite dimensional VV-RKHS by a finite number of parameters $(u_i)_{i=1}^N$, $u_i \in \mathcal{Y}$.

In the following we suppose we are given a cost function $c : \mathcal{Y} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$, such that $c(f(x), y)$ returns the error of the prediction $f(x)$ w.r.t. the ground truth y . A loss function of a model f with respect to an example $(x, y) \in \mathcal{X} \times \mathcal{Y}$ can be naturally defined from a cost function as $L(x, f, y) = c(f(x), y)$. Conceptually the function c evaluates the quality of the prediction versus its ground truth $y \in \mathcal{Y}$ while the loss function L evaluates the quality of the model f at a training point $(x, y) \in \mathcal{X} \times \mathcal{Y}$.

¹⁵ Sometimes referred to as minimal norm interpolation theorem.

Theorem 6.2 (Representer theorem). Let K be a \mathcal{Y} -Mercer Operator-Valued Kernel and \mathcal{H}_K its corresponding \mathcal{Y} -Reproducing Kernel Hilbert space.

Let $c : \mathcal{Y} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ be a cost function such that $L(x, f, y) = c(f(x), y)$ is a proper convex lower semi-continuous function in f for all $x \in \mathcal{X}$ and all $y \in \mathcal{Y}$.

Eventually let $\lambda \in \mathbb{R}_{>0}$ be the Tychonov regularization hyperparameter. The solution $f_s \in \mathcal{H}_K$ of the regularized optimization problem

$$f_s = \arg \min_{f \in \mathcal{H}_K} \frac{1}{N} \sum_{i=1}^N c(f(x_i), y_i) + \frac{\lambda}{2} \|f\|_K^2 \quad (6.3)$$

has the form $f_s = \sum_{j=1}^N K(\cdot, x_j) u_{s,j}$ where $(u_s)_j \in \mathcal{Y}$ and

$$u_s = \arg \min_{u \in \bigoplus_{i=1}^N \mathcal{Y}} \frac{1}{N} \sum_{i=1}^N c \left(\sum_{j=1}^N K(x_i, x_j) u_j, y_i \right) + \frac{\lambda}{2} \sum_{j,k=1}^N u_j^* K(x_j, x_k) u_k. \quad (6.4)$$

The first representer theorem was introduced by Wahba [182] in the case where $\mathcal{Y} = \mathbb{R}$. The extension to an arbitrary Hilbert space \mathcal{Y} has been proved by many authors in different forms [34, 86, 113]. The idea behind the representer theorem is that even though we minimize over the whole space \mathcal{H}_K , when $\lambda > 0$, the solution of Equation 6.3 falls inevitably into the set

$$\mathcal{H}_{K,s} = \left\{ \sum_{j=1}^N K_{x_j} u_j \mid \forall (u_i)_{i=1}^N \in \mathcal{Y}^N \right\}.$$

Therefore the result can be expressed as a finite linear combination of basis functions of the form $K(\cdot, x_k)$. Remark that we can perform the kernel expansion of $f_s = \sum_{j=1}^N K(\cdot, x_j) u_{s,j}$ even though $\lambda = 0$. However f_s is no longer the solution of Equation 6.3 over the whole space \mathcal{H}_K but a projection on the subspace $\mathcal{H}_{K,s}$. While this is in general not a problem for practical applications, it might raise issues for further theoretical investigations. In particular, it makes it difficult to perform theoretical comparison of the “exact” solution of Equation 6.3 with respect to the ORFF approximation solution given in Theorem 6.3.

Proof Since $f(x) = K_x^* f$ (see Equation 3.14), the optimization problem reads

$$f_s = \arg \min_{f \in \mathcal{H}_K} \frac{1}{N} \sum_{i=1}^N c(K_{x_i}^* f, y_i) + \frac{\lambda}{2} \|f\|_K^2$$

Let $W_s : \mathcal{H}_K \rightarrow \bigoplus_{i=1}^N \mathcal{Y}$ be the restriction¹ linear operator defined as

$$W_s f = \bigoplus_{i=1}^N K_{x_i}^* f,$$

¹ W_s is sometimes called the sampling or evaluation operator as in Minh, Bazzani, and Murino [120]. However we prefer calling it “restriction operator” as in Rosasco, Belkin, and Vito [142] since $W_s f$ is the restriction of f to the points in s .

with $K_{x_i}^* : \mathcal{H}_K \rightarrow \mathcal{Y}$ and $K_{x_i} : \mathcal{Y} \rightarrow \mathcal{H}_K$. Let $Y = \bigoplus_{i=1}^N y_i \in \mathcal{Y}^N$. We have

$$\langle Y, W_s f \rangle_{\bigoplus_{i=1}^N \mathcal{Y}} = \sum_{i=1}^N \langle y_i, K_{x_i}^* f \rangle_{\mathcal{Y}} = \sum_{i=1}^N \langle K_{x_i} y_i, f \rangle_{\mathcal{H}_K}.$$

Thus the adjoint operator $W_s^* : \bigoplus_{i=1}^N \mathcal{Y} \rightarrow \mathcal{H}_K$ is

$$W_s^* Y = \sum_{i=1}^N K_{x_i} y_i,$$

and the operator $W_s^* W_s : \mathcal{H}_K \rightarrow \mathcal{H}_K$ is

$$W_s^* W_s f = \sum_{i=1}^N K_{x_i} K_{x_i}^* f.$$

Let

$$\mathfrak{R}_\lambda(f, s) = \underbrace{\frac{1}{N} \sum_{i=1}^N c(f(x_i), y_i)}_{=\mathfrak{R}_c} + \frac{\lambda}{2} \|f\|_K^2$$

To ensure that \mathfrak{R}_λ has a global minimizer we need the following technical lemma (which is a consequence of the Hahn-Banach theorem for lower-semicontinuous functional, see Kurdila and Zabarankin [93]).

Lemma 6.1 *Let \mathfrak{R} be a proper, convex, lower semi-continuous functional, defined on a Hilbert space \mathcal{H} . If \mathfrak{R} is strongly convex, then \mathfrak{R} is coercive.*

Proof Consider the convex function $G(f) := \mathfrak{R}(f) - \lambda \|f\|^2$, for some $\lambda > 0$. Since \mathfrak{R} is by assumption proper, lower semi-continuous and strongly convex with parameter λ , G is proper, lower semi-continuous and convex. Thus Hahn-Banach theorem apply, stating that G is bounded by below by an affine functional. i. e. there exists f_0 and $f_1 \in \mathcal{H}$ such that

$$G(f) \geq G(f_0) + \langle f - f_0, f_1 \rangle, \quad \text{for all } f \in \mathcal{H}.$$

Then substitute the definition of G to obtain

$$\mathfrak{R}(f) \geq \mathfrak{R}(f_0) + \lambda (\|f\| - \|f_0\|) + \langle f - f_0, f_1 \rangle.$$

By the Cauchy-Schwartz inequality, $\langle f, f_1 \rangle \geq -\|f\|\|f_1\|$, thus

$$\mathfrak{R}(f) \geq \mathfrak{R}(f_0) + \lambda (\|f\| - \|f_0\|) - \|f\|\|f_1\| - \langle f_0, f_1 \rangle,$$

which tends to infinity as f tends to infinity. Hence \mathfrak{R} is coercive \square

Since c is proper, lower semi-continuous and convex by assumption, thus the term \mathfrak{R}_c is also proper, lower semi-continuous and convex. Moreover the term $\frac{\lambda}{2} \|f\|_K^2$ is strongly convex. Thus \mathfrak{R}_λ is strongly convex. Apply

[Lemma 6.1](#) to obtain the coercivity of \mathfrak{R}_λ , and then [Theorem 6.1](#) to show that \mathfrak{R}_λ has a unique minimizer and is attained. Then let

$$\mathcal{H}_{K,s} = \left\{ \sum_{j=1}^N K_{x_j} u_j \mid \forall (u_i)_{i=1}^N \in \mathcal{Y}^N \right\}.$$

¹⁶ $\mathcal{H}_{K,s}^\perp \oplus \mathcal{H}_{K,s} = \mathcal{H}_K$ because W_s is bounded.

For $f \in \mathcal{H}_{K,s}^\perp$, the operator W_s satisfies

$$\langle Y, W_s f \rangle_{\bigoplus_{i=1}^N \mathcal{Y}} = \langle \underbrace{f}_{\in \mathcal{H}_{K,s}^\perp}, \underbrace{\sum_{i=1}^N K_{x_i} y_i}_{\in \mathcal{H}_{K,s}} \rangle_{\mathcal{H}_K} = 0$$

for all sequences $(y_i)_{i=1}^N$, since $y_i \in \mathcal{Y}$. Hence,

$$(f(x_i))_{i=1}^N = 0 \quad (6.5)$$

Now for an arbitrary $f \in \mathcal{H}_K$, consider the orthogonal decomposition $f = f^\perp + f^\parallel$, where $f^\perp \in \mathcal{H}_{K,s}^\perp$ and $f^\parallel \in \mathcal{H}_{K,s}$. Then since $\|f^\perp + f^\parallel\|_{\mathcal{H}_K}^2 = \|f^\perp\|_{\mathcal{H}_K}^2 + \|f^\parallel\|_{\mathcal{H}_K}^2$, [Equation 6.5](#) and ?? shows that if $\lambda > 0$, clearly then

$$\mathfrak{R}_\lambda(f, s) = \mathfrak{R}_\lambda(f^\perp + f^\parallel, s) \geq \mathfrak{R}_\lambda(f^\parallel, s)$$

The last inequality holds only when $\|f^\perp\|_{\mathcal{H}_K} = 0$, that is when $f^\perp = 0$. As a result since the minimizer of \mathfrak{R}_λ is unique and attained, it must lies in $\mathcal{H}_{K,s}$. \square

The representer theorem shows that minimizing a functional in a VV-RKHS yields a solution which depends on all the points in the training set. Assuming that for all $x_i, x \in \mathcal{X}$ and for all $u_i \in \mathcal{Y}$ it takes time P , to compute $K(x_i, x)u_i$, making a prediction using the representer theorem takes NP . If $\mathcal{Y} = \mathbb{R}^p$, we can view $K(x_i, x)$ as a matrix and then $P = O(p^2)$ thus making a prediction cost $O(Np^2)$ operations.

6.1.2 Learning with Operator Random Fourier Feature maps

Instead of learning a model f that depends on all the points of the training set, we would like to learn a parametric model of the form $\tilde{f}(x) = \Phi(x)^*\theta$, where θ lives in some space $\widetilde{\mathcal{H}}$. We are interested in finding a parameter vector θ_s such that

$$\theta_s = \arg \min_{\theta \in \widetilde{\mathcal{H}}} \frac{1}{N} \sum_{i=1}^N c(\tilde{\Phi}(x_i)^*\theta, y_i) + \frac{\lambda}{2} \|\theta\|_{\widetilde{\mathcal{H}}}^2 \quad (6.6)$$

The following theorem states that when $\lambda > 0$ then learning with a feature map is equivalent to learn with a kernel. Moreover if $f_s \in \mathcal{H}_K$ is a solution of [Equation 6.7](#) and $\theta_s \in \mathcal{H}$ is the solution of [Equation 6.8](#), then $f_s = \Phi(\cdot)^*\theta_s$. This equivalence could have been obtained

by means of Lagrange duality. However in this proof we do not use such tool: we only focus on the representer theorem and the fact that there exist a partial isometry W between the VV-RKHS and a feature space \mathcal{H} . We show that if θ_s is a solution of [Equation 6.8](#), then theta belongs to $(\text{Ker } W)^\perp$, thus there is an isometry between $\theta_s \in \mathcal{H}$ and $\mathcal{H}_{\tilde{K}}$: namely W .

Theorem 6.3 (Feature equivalence). *Let \tilde{K} be an Operator-Valued Kernel such that for all $x, z \in \mathcal{X}$, $\tilde{\Phi}(x)^* \tilde{\Phi}(z) = \tilde{K}(x, z)$ where \tilde{K} is a \mathcal{Y} -Mercer OVK and $\mathcal{H}_{\tilde{K}}$ its corresponding \mathcal{Y} -Reproducing kernel Hilbert space.*

Let $c : \mathcal{Y} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ be a cost function such that $L(x, \tilde{f}, y) = c(\tilde{f}(x), y)$ is a proper convex lower semi-continuous function in $\tilde{f} \in \mathcal{H}_{\tilde{K}}$ for all $x \in \mathcal{X}$ and all $y \in \mathcal{Y}$.

Eventually let $\lambda \in \mathbb{R}_{>0} \mathbb{R}_+$ be the Tychonov regularization hyperparameter. The solution $f_s \in \mathcal{H}_{\tilde{K}}$ of the regularized optimization problem

$$\tilde{f}_s = \arg \min_{\tilde{f} \in \mathcal{H}_{\tilde{K}}} \frac{1}{N} \sum_{i=1}^N c(\tilde{f}(x_i), y_i) + \frac{\lambda}{2} \|\tilde{f}\|_{\tilde{K}}^2 \quad (6.7)$$

has the form $\tilde{f}_s = \tilde{\Phi}(\cdot)^ \theta_s$, where $\theta_s \in (\text{Ker } \tilde{W})^\perp$ and*

$$\theta_s = \arg \min_{\theta \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N c(\tilde{\Phi}(x_i)^* \theta, y_i) + \frac{\lambda}{2} \|\theta\|_{\mathcal{H}}^2 \quad (6.8)$$

Proof Since \tilde{K} is an operator-valued kernel, from [Theorem 6.2](#), [Equation 6.7](#) has a solution of the form

$$\begin{aligned} \tilde{f}_s &= \sum_{i=1}^N \tilde{K}(\cdot, x_i) u_i, \quad u_i \in \mathcal{Y}, x_i \in \mathcal{X} \\ &= \sum_{i=1}^N \tilde{\Phi}(\cdot)^* \tilde{\Phi}(x_i) u_i = \tilde{\Phi}(\cdot)^* \underbrace{\left(\sum_{i=1}^N \tilde{\Phi}(x_i) u_i \right)}_{=\theta \in (\text{Ker } \tilde{W})^\perp \subset \mathcal{H}}. \end{aligned}$$

Let

$$\theta_s = \arg \min_{\theta \in (\text{Ker } \tilde{W})^\perp} \frac{1}{N} \sum_{i=1}^N c(\tilde{\Phi}(x_i)^* \theta, y_i) + \frac{\lambda}{2} \|\tilde{\Phi}(\cdot)^* \theta\|_{\tilde{K}}^2. \quad (6.9)$$

Since $\theta \in (\text{Ker } \tilde{W})^\perp$ and W is an isometry from $(\text{Ker } \tilde{W})^\perp \subset \mathcal{H}$ onto $\mathcal{H}_{\tilde{K}}$, we have $\|\tilde{\Phi}(\cdot)^* \theta\|_{\tilde{K}}^2 = \|\theta\|_{\mathcal{H}}^2$. Hence

$$\theta_s = \arg \min_{\theta \in (\text{Ker } \tilde{W})^\perp} \frac{1}{N} \sum_{i=1}^N c(\tilde{\Phi}(x_i)^* \theta, y_i) + \frac{\lambda}{2} \|\theta\|_{\mathcal{H}}^2$$

Finding a minimizer θ_s over $(\text{Ker } \tilde{W})^\perp$ is not the same as finding a minimizer over $\tilde{\mathcal{H}}$. Although in both cases Mazur-Schauder's theorem guarantees that the respective minimizers are unique, they might not be the same. Since \tilde{W} is bounded, $\text{Ker } \tilde{W}$ is closed, so that we can perform the decomposition $\tilde{\mathcal{H}} = (\text{Ker } \tilde{W})^\perp \oplus (\text{Ker } \tilde{W})$. Then clearly by linearity of W and the fact that for all $\theta^{\parallel} \in \text{Ker } \tilde{W}$, $\tilde{W}\theta^{\parallel} = 0$, if $\lambda > 0$ we have

$$\theta_s = \arg \min_{\theta \in \tilde{\mathcal{H}}} \frac{1}{N} \sum_{i=1}^N c(\tilde{\Phi}(x_i)^* \theta, y_i) + \frac{\lambda}{2} \|\theta\|_{\tilde{\mathcal{H}}}^2$$

Thus

$$\begin{aligned} \theta_s = & \arg \min_{\substack{\theta^\perp \in (\text{Ker } \tilde{W})^\perp, \\ \theta^{\parallel} \in \text{Ker } \tilde{W}}} \frac{1}{N} \sum_{i=1}^N c\left(\left(\tilde{W}\theta^\perp\right)(x) + \underbrace{\left(\tilde{W}\theta^{\parallel}\right)(x)}_{=0 \text{ for all } \theta^{\parallel}}, y_i\right) \\ & + \frac{\lambda}{2} \|\theta^\perp\|_{\tilde{\mathcal{H}}}^2 + \underbrace{\frac{\lambda}{2} \|\theta^{\parallel}\|_{\tilde{\mathcal{H}}}^2}_{=0 \text{ only if } \theta^{\parallel}=0} \end{aligned}$$

Thus

$$\theta_s = \arg \min_{\theta^\perp \in (\text{Ker } \tilde{W})^\perp} \frac{1}{N} \sum_{i=1}^N c\left(\left(\tilde{W}\theta^\perp\right)(x), y_i\right) + \frac{\lambda}{2} \|\theta^\perp\|_{\tilde{\mathcal{H}}}^2.$$

Hence minimizing over $(\text{Ker } \tilde{W})^\perp$ or $\tilde{\mathcal{H}}$ is the same when $\lambda > 0$. Eventually,

$$\theta_s = \arg \min_{\theta \in \tilde{\mathcal{H}}} \frac{1}{N} \sum_{i=1}^N c(\tilde{\Phi}(x_i)^* \theta, y_i) + \frac{\lambda}{2} \|\theta\|_{\tilde{\mathcal{H}}}^2.$$

□

In the aforementioned theorem, we use the notation \tilde{K} and $\tilde{\Phi}$ because our main subject of interest is the ORFF map. However this theorem works for *any* feature maps $\Phi(x) \in \mathcal{L}(\mathcal{Y}, \mathcal{H})$ even when \mathcal{H} is infinite dimensional.² This shows that when $\lambda > 0$ the solution of [Equation 6.4](#) with the approximated kernel $K(x, z) \approx \tilde{K}(x, z) = \tilde{\Phi}(x)^* \tilde{\Phi}(z)$ is the same than the solution of [Equation 6.6](#) up to a linear transformation. Namely, if u_s is the solution of [Equation 6.4](#), θ_s is the solution of [Equation 6.9](#) and $\lambda > 0$ we have

$$\theta_s = \sum_{i=1}^N \tilde{\Phi}(x_i)(u_s)_i \in (\text{Ker } W)^\perp \subseteq \tilde{\mathcal{H}}.$$

² If $\Phi(x) : \mathcal{L}(\mathcal{Y}, \mathcal{H})$ and $\dim(\mathcal{H}) = \infty$, the decomposition $\mathcal{H} = (\text{Ker } W) \oplus (\text{Ker } W)^\perp$ holds since \mathcal{H} is a Hilbert space and W is a bounded operator.

If $\lambda_K = 0$ we can still find a solution u_s of [Equation 6.4](#). By construction of the kernel expansion, we have $u_s \in (\text{Ker } W)^\perp$. However looking at the proof of [Theorem 6.3](#) we see that θ_s might *not* belong to $(\text{Ker } W)^\perp$. We can compute a residual vector

$$r_s = \sum_{i=1}^N \tilde{\Phi}(x_i)(u_s)_i - \theta_s.$$

Since $\sum_{j=1}^N \tilde{\Phi}(x_j) \in (\text{Ker } W)^\perp$ by construction, if $r_s = 0$, it means that λ_K is large enough for both representer theorem and ORFF representer theorem to apply. If $r_s \neq 0$ but $\tilde{\Phi}(\cdot)^* r_s = 0$ it means that both θ_s and $\sum_{j=1}^N \tilde{\Phi}(x_j) u_s$ are in $(\text{Ker } W)^\perp$, thus the representer theorem fails to find the “true” solution over the whole space $\mathcal{H}_{\tilde{K}}$ but returns a projection onto $\mathcal{H}_{\tilde{K},s}$ of the solution. If $r_s \neq 0$ and $\tilde{\Phi}(\cdot)^* r_s \neq 0$ means that θ_s is *not* in $(\text{Ker } W)^\perp$, thus the feature equivalence theorem fails to apply. Since $r_s = \sum_{i=1}^N \tilde{\Phi}(x_i)(u_s)_i - \theta_s^\perp - \theta_s^\parallel$ and $\sum_{i=1}^N \tilde{\Phi}(x_i)(u_s)_i$ is in $(\text{Ker } W)^\perp$, with mild abuse of notation we write $r_s = \theta^\parallel$. This remark is illustrated in [Figures 6.1](#) and [6.2](#).

In [Figure 6.1](#), we generated the data from a sine wave to which we add some Gaussian noise. We learned a Gaussian kernel based RFF model (blue curve) and a kernel model (yellow curve) where the kernel is obtained from the RFF map. The left column represents the fit of the model to the points for four different values of λ (top to bottom: 10^{-2} , 10^{-5} , $10e^{-10}$, 0). The middle column shows if the RFF solution θ_s is in $(\text{Ker } \tilde{W})^\perp$. This is true for all values of λ . The right column shows that even though θ_s is in $(\text{Ker } \tilde{W})^\perp$, when $\lambda \rightarrow 0$ learning with RFF is different from learning with the kernel constructed from the RFF maps since the coefficients of θ^\parallel are all different from 0.

[Figure 6.2](#) is the same setting than [Figure 6.1](#) except that we decreased the scale parameter σ of the Gaussian kernel to make it overfit, and emphasize that when $\lambda = 0$, θ_s might not belong to $(\text{Ker } \tilde{W})^\perp$, as represented on the middle column.

6.2 SOLVING ORFF-BASED REGRESSION

In order to find a solution to [Equation 6.9](#), we first turn our attention to gradient descent methods. In the following we let

$$\mathfrak{R}_\lambda(\theta) = \frac{1}{N} \sum_{i=1}^N c \left(\tilde{\Phi}(x_i)^* \theta, y_i \right) + \frac{\lambda}{2} \|\theta\|_{\mathcal{H}}^2 \quad (6.10)$$

Then, we study the complexity in time of the proposed algorithm.

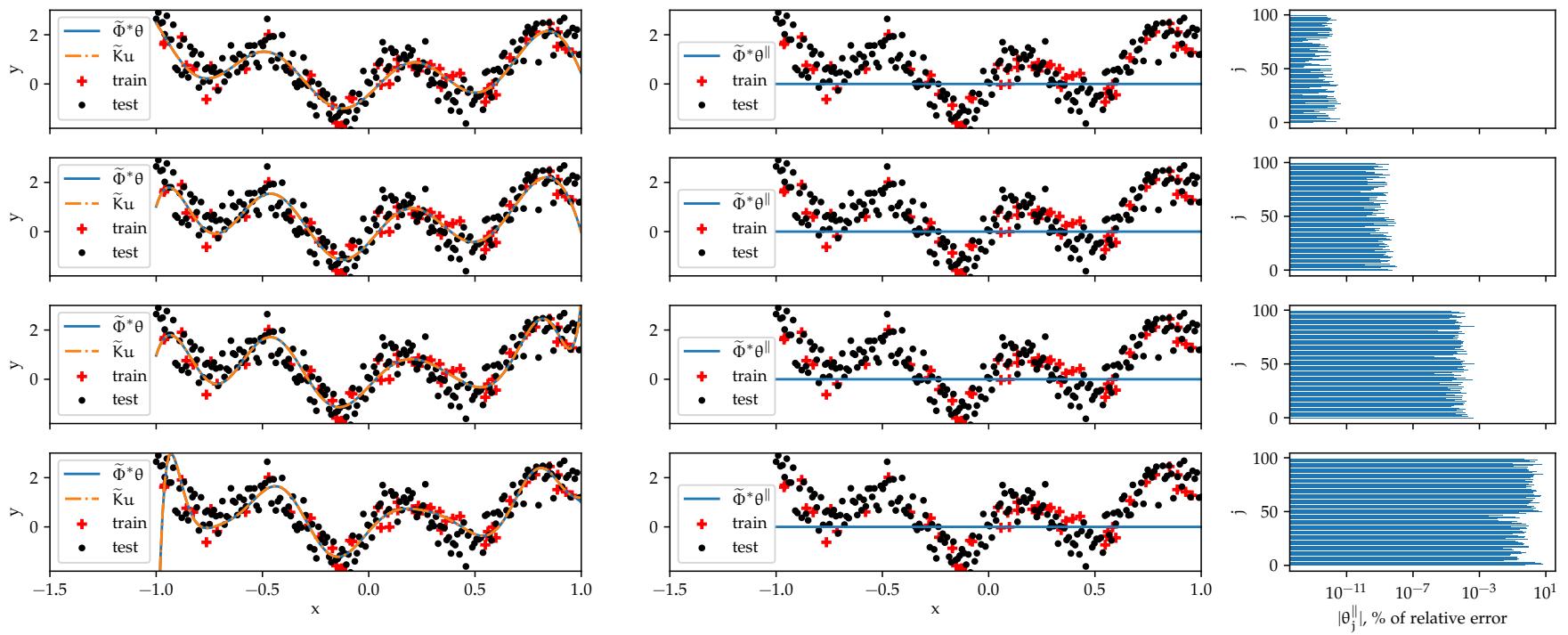


Figure 6.1: ORFF equivalence theorem.

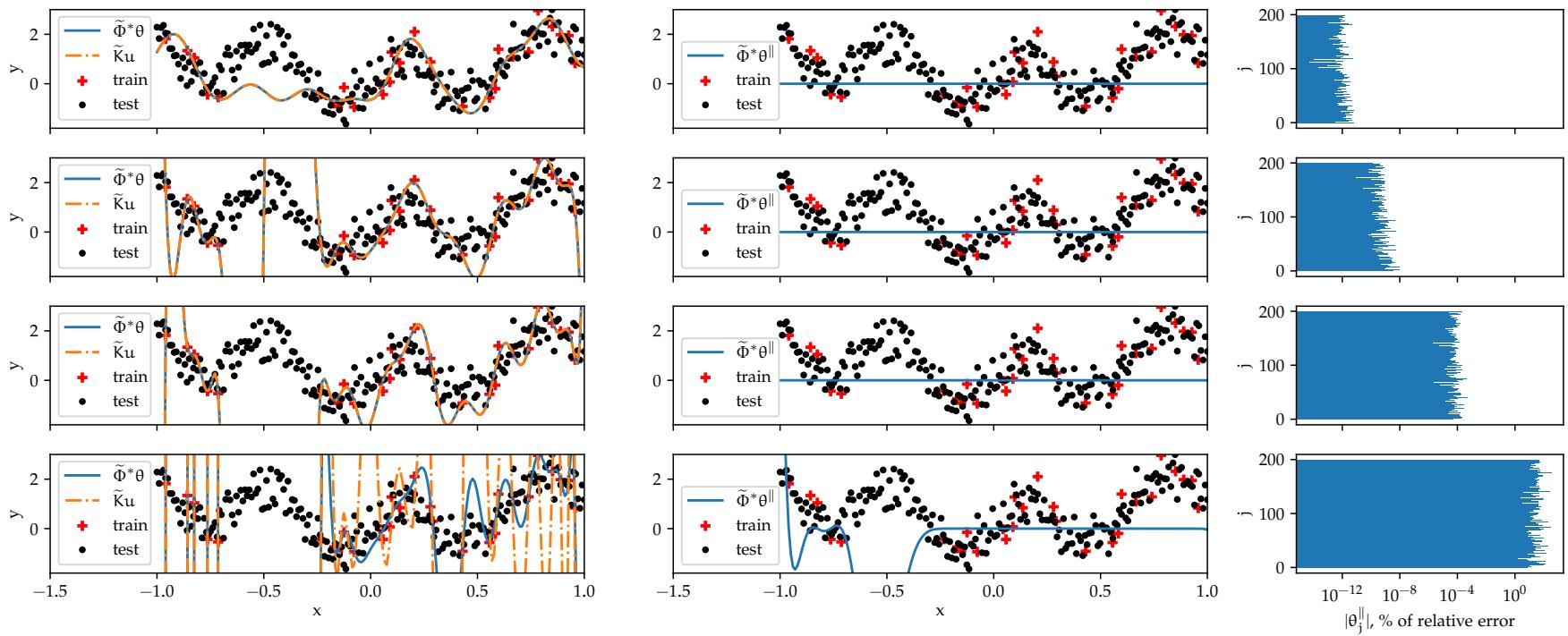


Figure 6.2: ORFF equivalence theorem with overfitting.

6.2.1 Gradient descent methods

Since the solution of [Equation 6.9](#) is unique when $\lambda > 0$, a sufficient and necessary condition is that the gradient of \mathfrak{R}_λ at the minimizer θ_s is zero. We use the Frechet derivative, the strongest notion of derivative in Banach spaces¹⁷ [46, 93] which directly generalizes the notion of gradient to Banach spaces. A function $f : \mathcal{H}_0 \rightarrow \mathcal{H}_1$ is called Frechet differentiable at $\theta_0 \in \mathcal{H}_0$ if there exists a bounded linear operator $A \in \mathcal{L}(\mathcal{H}_0, \mathcal{H}_1)$ such that

$$\lim_{\|h\|_{\mathcal{H}_0} \rightarrow 0} \frac{\|f(\theta_0 + h) - f(\theta_0) - Ah\|_{\mathcal{H}_1}}{\|h\|_{\mathcal{H}_0}} = 0$$

We write

$$(D_F f)(\theta_0) = \left. \frac{\partial f(\theta)}{\partial \theta} \right|_{\theta=\theta_0} = A$$

and call it Frechet derivative of f with respect to θ at θ_0 . With mild abuse of notation we write

$$\left. \frac{\partial f(\theta)}{\partial \theta} \right|_{\theta=\theta_0} = \frac{\partial f(\theta_0)}{\partial \theta_0}.$$

The chain rule is valid in this context [93, theorem 4.1.1 page 140]. Namely, let \mathcal{H}_0 , \mathcal{H}_1 and \mathcal{H}_2 be three Hilbert spaces. If a function $f : \mathcal{H}_0 \rightarrow \mathcal{H}_1$ is Frechet differentiable at θ and $g : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ is Frechet differentiable at $f(\theta)$ then $g \circ f$ is Frechet differentiable at θ and for all $h \in \mathcal{H}_0$

$$\frac{\partial}{\partial \theta} (g \circ f)(\theta) \circ h = \frac{\partial g(f(\theta))}{\partial f(\theta)} \circ \frac{\partial f(\theta)}{\partial \theta} \circ h,$$

or equivalently,

$$D_F(g \circ f)(\theta) \circ h = (D_F g)(f(\theta)) \circ (D_F f)(\theta) \circ h.$$

If $f : \mathcal{H} \rightarrow \mathbb{R}$ then $(D_F f)(\theta_0) \in \mathcal{H}^*$ for all $\theta_0 \in \mathcal{H}$, and by Riesz's representation theorem we define the gradient of f noted $\nabla_\theta f(\theta) \in \mathcal{H}$ as the vector in \mathcal{H} such that

$$\langle \nabla_\theta f(\theta), h \rangle_{\mathcal{H}} = (D_F f)(\theta) \circ h = \frac{\partial f(\theta)}{\partial \theta} \circ h.$$

For a function $f : \mathcal{H}_0 \rightarrow \mathcal{H}_1$ we note the jacobian of f as $J_\theta f(\theta) = \frac{\partial f(\theta)}{\partial \theta}$. In this context if $f : \mathcal{H}_0 \rightarrow \mathcal{H}_1$ and $g : \mathcal{H}_1 \rightarrow \mathbb{R}$ the chain rule reads for all $h \in \mathcal{H}_0$

$$\frac{\partial}{\partial \theta} (g \circ f)(\theta) \circ h = \frac{\partial g(f(\theta))}{\partial f(\theta)} \circ J_\theta f(\theta) \circ h.$$

By Riesz's representation theorem,

$$\begin{aligned} \langle \nabla_\theta (g \circ f)(\theta), h \rangle_{\mathcal{H}_0} &= \langle \nabla_{f(\theta)} g(f(\theta)), J_\theta f(\theta) h \rangle_{\mathcal{H}_0} \\ &= \langle (J_\theta f(\theta))^* \nabla_{f(\theta)} g(f(\theta)), h \rangle_{\mathcal{H}_0} \end{aligned}$$

¹⁷ Here we view the Hilbert space \mathcal{H} (feature space) as a reflexive Banach space.

Hence

$$\nabla_{\theta}(g \circ f)(\theta) = (J_{\theta} f(\theta))^* \nabla_{f(\theta)} g(f(\theta)).$$

Thus by linearity and applying the chaine rule to [Equation 6.9](#) we have

$$\begin{aligned}\nabla_{\theta} c\left(\tilde{\Phi}(x_i)^* \theta, y_i\right) &= \tilde{\Phi}(x_i) V^* \left(\frac{\partial}{\partial y} c(y, y_i) \Big|_{y=\tilde{\Phi}(x_i)^* \theta} \right)^*, \\ \nabla_{\theta} \|\theta\|_{\tilde{\mathcal{H}}}^2 &= 2\theta.\end{aligned}$$

Provided that $c(y, y_i)$ is Frechet differentiable w.r.t. y , for all y and $y_i \in \mathcal{Y}$ we have $\nabla_{\theta} \mathfrak{R}_{\lambda}(\theta) \in \tilde{\mathcal{H}}$ and

$$\nabla_{\theta} \mathfrak{R}_{\lambda}(\theta) = \frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) \left(\frac{\partial}{\partial y} c(y, y_i) \Big|_{y=\tilde{\Phi}(x_i)^* \theta} \right)^* + \lambda \theta \quad (6.11)$$

Example 6.1 (Naive closed form for the squared error cost). Consider the cost function defined for all $y, y' \in \mathcal{Y}$ by $c(y, y') = \frac{1}{2} \|y - y'\|_2^2$. Then

$$\left(\frac{\partial}{\partial y} c(y, y_i) \Big|_{y=\tilde{\Phi}(x_i)^* \theta} \right)^* = (\tilde{\Phi}(x_i)^* \theta - y_i).$$

Thus, since the optimal solution θ_s verifies $\nabla_{\theta_s} \mathfrak{R}_{\lambda}(\theta_s) = 0$ we have

$$\frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) (\tilde{\Phi}(x_i)^* \theta_s - y_i) + \lambda \theta_s = 0.$$

Therefore,

$$\left(\frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) \tilde{\Phi}(x_i)^* + \lambda I_{\tilde{\mathcal{H}}} \right) \theta_s = \frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) y_i. \quad (6.12)$$

Suppose that $\mathcal{Y} \subseteq \mathbb{R}^p$, and for all $x \in \mathcal{X}$, $\tilde{\Phi}(x) : \mathbb{R}^r \rightarrow \mathbb{R}^p$ where all spaces are endowed with the euclidean inner product. From this we can derive [Algorithm 3](#) which returns the closed form solution of [Equation 6.10](#) for $c(y, y') = \frac{1}{2} \|y - y'\|_2^2$.

If one considers a Mahalanobis inner product, evaluation of operators has to be done with extra care since the adjoint operator is *not* the classic conjugate transpose of the operator (see [Remark 6.2](#)). Indeed let $x, z \in \mathcal{Y} = \mathbb{C}^p$ endowed with its standard basis B , and $\langle x, y \rangle_{\mathcal{Y}} = \langle x, \Sigma^{-1}z \rangle_2$ where Σ is some symmetric positive-definite operator w.r.t. the basis B . Some simple calculations shows that given an operator $A \in \mathcal{L}(\mathcal{Y})$,

$$(A^*)_{ij} := \langle e_j, \Sigma^{-1} A^* e_i \rangle_2 = \overline{\langle \Sigma^{-1} A^* e_i, e_j \rangle_2} = \overline{\langle e_i, A \Sigma^{-1} e_j \rangle_2} := \overline{(\Sigma A \Sigma^{-1})_{ji}}$$

Thus $A^* = \Sigma^{-1} \bar{A}^T \Sigma$.

Remark 6.2 Notice that the evaluation of each operator $\nabla_\theta \mathfrak{R}_\lambda(\theta)$, V^* , $\tilde{\Phi}(x_i)^*$'s depends on the inner product of the respective spaces in which they are defined. Namely \mathcal{Y} , and \mathcal{H} . For instance if one chooses $\mathcal{H} = \bigoplus_{j=1}^D \mathcal{Y}'$, $\mathcal{Y}' = \mathbb{R}^{u'}$ endowed with the Euclidean inner product $\langle \theta', \theta \rangle_{\mathcal{Y}'} = \langle \theta', \theta \rangle_2$, \mathcal{Y} endowed with a Mahalanobis inner product $\langle u', u \rangle_{\mathcal{U}} = \langle u', \Sigma^{-1}u \rangle_2$ where Σ is some symmetric positive definite operator, then for all $x \in \mathcal{X}$,

$$\tilde{\Phi}(x)_{ij} = \langle e_j \tilde{\Phi}(x) e_i \rangle_2 = \langle e_i, \Sigma^{-1} \Sigma \tilde{\Phi}(x)^* e_j \rangle_2 = (\Sigma \tilde{\Phi}(x)^*)_{ji}.$$

Thus $\tilde{\Phi}(x)^* = \Sigma^{-1} \tilde{\Phi}(x)^\top$ and then [Equation 6.12](#) reads

$$\left(\frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) \Sigma^{-1} \tilde{\Phi}(x_i)^\top + \lambda I_{\mathcal{H}} \right) \theta_s = \frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) y_i.$$

6.2.2 Complexity analysis

[Algorithm 3](#) constitutes our first step toward large-scale learning with Operator-Valued Kernels. We can easily compute the time complexity of [Algorithm 3](#) when all the operators act on finite dimensional Hilbert spaces. Suppose that $p = \dim(\mathcal{Y}) < \infty$ and for all $x \in \mathcal{X}$, $\tilde{\Phi}(x) : \mathcal{Y} \rightarrow \mathcal{H}$ where $r = \dim(\mathcal{H}) < \infty$ is the dimension of the re-description space $\mathcal{H} = \mathbb{R}^r$. Since p and $r < \infty$, we view the operators $\tilde{\Phi}(x)$ and $I_{\mathcal{H}}$ as matrices. Step 1 costs $O_t(Nr^2p)$. Steps 2 costs $O_t(Nrp)$. For step 3, the naive inversion of the operator costs $O_t(r^3)$. Eventually the overall complexity of [Algorithm 3](#) is

$$O_t(r^2(Np + r)),$$

while the space complexity is $O_s(r^2)$. This complexity is to compare with the kernelized solution. Let

$$\mathbf{K} : \begin{cases} \mathcal{Y}^N \rightarrow \mathcal{Y}^N \\ u \mapsto \bigoplus_{i=1}^{N+u} \sum_{j=1}^{N+u} K(x_i, x_j) u_j \end{cases}$$

When $\mathcal{Y} = \mathbb{R}$,

$$\mathbf{K} = \begin{pmatrix} K(x_1, x_1) & \dots & K(x_1, x_{N+u}) \\ \vdots & \ddots & \vdots \\ K(x_{N+u}, x_1) & \dots & K(x_{N+u}, x_{N+u}) \end{pmatrix}$$

is called the Gram matrix of K . When $\mathcal{Y} = \mathbb{R}^p$, \mathbf{K} is a matrix-valued Gram matrix of size $pN \times pN$ where each entry $\mathbf{K}_{ij} \in \mathcal{M}_{p,p}(\mathbb{R})$. Then the equivalent kernelized solution u_s of [Theorem 6.2](#) is

$$\left(\frac{1}{N} \mathbf{K} + \lambda I_{\bigoplus_{i=1}^N \mathcal{Y}} \right) u_s = \bigoplus_{i=1}^N y_i.$$

Algorithm 3: Naive closed form for the squared error cost.**Input :**

- $s = (x_i, y_i)_{i=1}^N \in (\mathcal{X} \times \mathbb{R}^p)^N$ a sequence of supervised training points,
- $\tilde{\Phi}(x_i) \in \mathcal{L}(\mathbb{R}^p, \mathbb{R}^r)$ a feature map defined for all $x_i \in \mathcal{X}$,
- $\lambda \in \mathbb{R}_{>0}$ the Tychonov regularization term,

Output: A model

$$h : \begin{cases} \mathcal{X} \rightarrow \mathbb{R}^p \\ x \mapsto \tilde{\Phi}(x)^T \theta_s, \end{cases}$$

such that θ_s minimizes [Equation 6.10](#), where $c(y, y') = \|y - y'\|_2^2$
and \mathbb{R}^r and \mathbb{R}^p are Hilbert spaces endowed with the
Euclidean inner product.

- 1 $P \leftarrow \frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) \tilde{\Phi}(x_i)^T \in \mathcal{L}(\mathbb{R}^r, \mathbb{R}^r);$
- 2 $Y \leftarrow \frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) y_i \in \mathbb{R}^r;$
- 3 $\theta_s \leftarrow \text{solve}_\theta((P + \lambda I_r) \theta = Y);$
- 4 **return** $h : x \mapsto \tilde{\Phi}(x)^T \theta_s;$

which has time complexity $O(N^3 p^3)$ and space complexity $O_s(N^2 p^2)$. Suppose we are given a generic ORFF map (see [Subsection 4.3.3](#)). Then $r = 2Dp$, where D is the number of samples. Hence [Algorithm 3](#) is better than its kernelized counterpart when $r = 2Dp$ is small compared to Np . Thus, roughly speaking it is better to use [Algorithm 3](#) when the number of features, r , required is small compared to the number of training points. Notice that [Algorithm 3](#) has a linear complexity with respect to the number of supervised training points N , thus it is more suitable for large scale learning provided that D does not grow linearly with N .

Yet naive learning with [Algorithm 3](#) by viewing all the operators as matrices is still problematic. Indeed learning p independent models with scalar Random Fourier Features would cost $O_t(D^2 p^3 (N + D))$ since $r = 2Dp$. This means that learning a vector-valued function has increased the (expected) complexity from p to p^3 . However in some cases we can drastically reduce the complexity by viewing the feature-maps as linear operators rather than matrices.

6.3 EFFICIENT LEARNING WITH ORFF

When developing [Algorithm 3](#) we considered that the feature map $\tilde{\Phi}(x)$ was a matrix from \mathbb{R}^p to \mathbb{R}^r for all $x \in \mathcal{X}$, and therefore that computing $\tilde{\Phi}(x)\tilde{\varphi}(z)^T$ has a time complexity of $O(r^2 p)$. While this holds

true in the most generic scenario, in many cases the feature maps present some structure or sparsity allowing to reduce the computational cost of evaluating the feature map. We focus on the Operator-valued Random Fourier Feature given by [Algorithm 1](#), developped in [Section 4.3](#) and [Subsection 4.3.3](#) and treat the decomposable kernel, the curl-free kernel and the divergence-free kernel as an example. We recall that if $\mathcal{Y}' = \mathbb{R}^{p'}$ and $\mathcal{Y} = \mathbb{R}^p$, then $\tilde{\mathcal{H}} = \mathbb{R}^{2Dp'}$ thus the Operator-valued Random Fourier Features given in [Chapter 4](#) have the form

$$\begin{cases} \tilde{\Phi}(x) \in \mathcal{L}\left(\mathbb{R}^p, \mathbb{R}^{2Dp'}\right) & : y \mapsto \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D (x, \omega_j) B(\omega_j)^T y \\ \tilde{\Phi}(x)^T \in \mathcal{L}\left(\mathbb{R}^{2Dp'}, \mathbb{R}^p\right) & : \theta \mapsto \frac{1}{\sqrt{D}} \sum_{j=1}^D (x, \omega_j) B(\omega_j) \theta_j \end{cases},$$

where $\omega_j \sim \Pr_{\widehat{\text{Haar}}, p}$ i. i. d. and $B(\omega_j) \in \mathcal{L}\left(\mathbb{R}^p, \mathbb{R}^{p'}\right)$ for all $\omega_j \in \widehat{\mathcal{X}}$. Hence the Operator-valued Random Fourier Feature can be seen as the block matrix

$$\tilde{\Phi}(x) = \begin{pmatrix} \cos(x, \omega_1) B(\omega_1)^T \\ \sin(x, \omega_1) B(\omega_1)^T \\ \vdots \\ \cos(x, \omega_D) B(\omega_D)^T \\ \sin(x, \omega_D) B(\omega_D)^T \end{pmatrix} \in \mathcal{M}_{2Dp', p}(\mathbb{R}), \quad (6.13)$$

$\omega_j \sim \Pr_{\widehat{\text{Haar}}, p}$ i. i. d..

6.3.1 Case of study: the decosubmposable kernel

Throughout this section we show how the mathematical formulation relates to a concrete (Python) implementation. We propose a Python implementation based on NumPy [127], SciPy [81] and Scikit-learn [132]. Following [Equation 6.13](#), the feature map associated to the decomposable kernel would be

$$\tilde{\Phi}(x) = \frac{1}{\sqrt{D}} \begin{pmatrix} \cos(x, \omega_1) B^T \\ \sin(x, \omega_1) B^T \\ \vdots \\ \cos(x, \omega_D) B^T \\ \sin(x, \omega_D) B^T \end{pmatrix} = \underbrace{\frac{1}{\sqrt{D}} \begin{pmatrix} \cos(x, \omega_1) \\ \sin(x, \omega_1) \\ \vdots \\ \cos(x, \omega_D) \\ \sin(x, \omega_D) \end{pmatrix}}_{\tilde{\varphi}(x)} \otimes B^T,$$

$\omega_j \sim \Pr_{\widehat{\text{Haar}}, p}$ i. i. d., which would lead to the following naive python implementation for the Gaussian (RBF) kernel of parameter γ , whose associated spectral distribution is $\Pr_p = \mathcal{N}(0, 2\gamma)$.

```

def NaiveDecomposableGaussianORFF(X, A, gamma=1.,
                                    D=100, eps=1e-5, random_state=0):
    """Return the Naive ORFF map associated with the data X.

    Parameters
    -----
    X : {array-like}, shape = [n_samples, n_features]
        Samples.
    A : {array-like}, shape = [n_targets, n_targets]
        Operator of the Decomposable kernel (positive semi-definite)
    gamma : {float},
        Gamma parameter of the RBF kernel.
    D : {integer}
        Number of random features.
    eps : {float}
        Cutoff threshold for the singular values of A.
    random_state : {integer}
        Seed of the generator.

    Returns
    -----
    \tilde{\Phi}(X) : array
    """
    # Decompose A=BB^ transpose
    u, s, v = svd(A, full_matrices=False, compute_uv=True)
    B = dot(diag(sqrt(s[s > eps])), v[s > eps, :])

    # Sample a RFF from the scalar Gaussian kernel
    phi_s = RBFSampler(gamma=gamma, n_components=D, random_state=random_state)
    phiX = phi_s.fit_transform(X)

    # Create the ORFF linear operator
    return matrix(kron(phiX, B))

```

Let $\theta \in \mathbb{R}^{2Dp'}$ and $y \in \mathbb{R}$. With such implementation evaluating a matrix vector product such as $\tilde{\Phi}(x)^T\theta$ or $\tilde{\Phi}(x)y$ have $O_t(2Dp'p)$ time complexity and $O_s(2Dp'p)$ of space complexity, which is utterly inefficient. Indeed, recall that if $B \in \mathcal{M}_{p,p'}(\mathbb{R}^{p'})$ is matrix, the operator $\tilde{\Phi}(x)$ corresponding to the decomposable kernel is

$$\begin{aligned}\tilde{\Phi}(x)y &= \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos\langle x, \omega_j \rangle B^T y \\ \sin\langle x, \omega_j \rangle B^T y \end{pmatrix} \\ &= \left(\frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos\langle x, \omega_j \rangle \\ \sin\langle x, \omega_j \rangle \end{pmatrix} \right) \otimes (B^T y)\end{aligned}$$

and

$$\begin{aligned}\tilde{\Phi}(x)^T\theta &= \frac{1}{\sqrt{D}} \sum_{j=1}^D \cos\langle x, \omega_j \rangle B\theta_j + \sin\langle x, \omega_j \rangle B\theta_j \\ &= B \left(\frac{1}{\sqrt{D}} \sum_{j=1}^D (\cos\langle x, \omega_j \rangle + \sin\langle x, \omega_j \rangle) \theta_j \right).\end{aligned}\tag{6.14}$$

Which requires only evaluation of B on y and can be implemented easily in Python thanks to SciPy's LinearOperator. Note that the computation of these expressions can be fully vectorized¹⁸ using the vec-

¹⁸ See Walt, Colbert, and Varoquaux [184].

torization property of the Kronecker product. In the following we consider $\Theta \in \mathcal{M}_{2D,u'}(\mathbb{R})$ and the operator $\text{vec} : \mathcal{M}_{p',2D}(\mathbb{R}) \rightarrow \mathbb{R}^{2Dp'}$ which turns a matrix into a vector (i.e. $\theta_{p'i+j} = \text{vec}(\Theta_{ij})$, $i \in \mathbb{N}_{(2D-1)}$ and $j \in \mathbb{N}_{p'}^*$). Then

$$(\tilde{\varphi}(x) \otimes B^T)^T \theta = (\tilde{\varphi}(x)^T \otimes B) \text{vec}(\Theta) = \text{vec}(B\Theta\tilde{\varphi}(x)).$$

with this trick, many authors [29, 142, 158] notice that the decomposable kernel usually yields a Stein equation [134]. Indeed rewriting step 3 of [Algorithm 3](#) gives a system to solve of the form

$$\tilde{\varphi}(X)\tilde{\varphi}(X)^T\Theta B^T B + \lambda\Theta - Y = 0 \Leftrightarrow (\tilde{\varphi}(X)\tilde{\varphi}(X)^T \otimes B^T B + \lambda I_{2Dp'}) \theta - Y = 0.$$

¹⁹ See for instance Sleijpen, Sonneveld, and Van Gijzen [160].

Many solvers exist to solve efficiently this kind of systems¹⁹, but most of them share the particularity that they are not just restricted to handle Stein equations. Broadly speaking, iterative solvers (or matrix-free solvers) are designed to solve any system of equation of the form $PX = C$, where P is a linear operator (not a matrix). This is exactly our case where $\tilde{\varphi}(x) \otimes B^T$ is the matrix form of the operator $\Theta \mapsto \text{vec}(B\Theta\tilde{\varphi}X)$.

This leads us to the following (more efficient) Python implementation of the Decomposable ORFF “operator” to be fed to a matrix-free solvers.

```
def EfficientDecomposableGaussianORFF(X, A, gamma=1.,
                                         D=100, eps=1e-5, random_state=0):
    """Return the efficient ORFF map associated with the data X.

    Parameters
    -----
    X : {array-like}, shape = [n_samples, n_features]
        Samples.
    A : {array-like}, shape = [n_targets, n_targets]
        Operator of the Decomposable kernel (positive semi-definite)
    gamma : {float},
        Gamma parameter of the RBF kernel.
    D : {integer}
        Number of random features.
    eps : {float}
        Cutoff threshold for the singular values of A.
    random_state : {integer}
        Seed of the generator.

    Returns
    -----
    \tilde{\Phi}(X) : Linear Operator, callable
    """
    # Decompose A=BB^T
    u, s, v = svd(A, full_matrices=False, compute_uv=True)
    B = dot(diag(sqrt(s[s > eps])), v[s > eps, :])

    # Sample a RFF from the scalar Gaussian kernel
    phi_s = RBFSampler(gamma=gamma, n_components=D, random_state=random_state)
    phiX = phi_s.fit_transform(X)

    # Create the ORFF linear operator
    cshape = (D, B.shape[0])
```

```

rshape = (X.shape[0], B.shape[1])
return LinearOperator((phiX.shape[0] * B.shape[1], D * B.shape[0]),
                      matvec=lambda b: dot(phiX, dot(b.reshape(cshape),
                        B)),
                      rmatvec=lambda r: dot(phiX.T, dot(r.reshape(rshape),
                        B.T)))

```

6.3.2 Linear operators in matrix form

For convenience we give the operators corresponding to the decomposable, curl-free and divergence-free kernels in matrix form. Let $(x_i)_{i=1}^N$, $N \in \mathbb{N}^*$, x_i 's in \mathbb{R}^d , $d \leq \infty$ be a sequence of points in \mathbb{R}^d . We note

$$X = \begin{pmatrix} x_1 & \dots & x_N \end{pmatrix} \in \mathcal{M}_{d,N}$$

the data matrix where each column represents a data point³. Naturally if $\tilde{\Phi}(x) : \mathbb{R}^u \rightarrow \mathbb{R}^{r_1}$ and $\tilde{\varphi}(x) : \mathbb{R} \rightarrow \mathbb{R}^{r_2}$, for all $x \in \mathbb{R}^d$ we define

$$\tilde{\Phi}(X) = \begin{pmatrix} \tilde{\Phi}(x_1) & \dots & \tilde{\Phi}(x_N) \end{pmatrix} \in \mathcal{M}_{r_1, N_u}$$

and

$$\tilde{\varphi}(X) = \begin{pmatrix} \tilde{\varphi}(x_1) & \dots & \tilde{\varphi}(x_N) \end{pmatrix} \in \mathcal{M}_{r_2, N}$$

and

$$Y = \begin{pmatrix} y_1 & \dots & y_N \end{pmatrix} \in \mathcal{M}_{u,N}.$$

Given a matrix $X \in \mathcal{M}_{m,n}(\mathbb{R})$, we note $X_{\bullet i}$ the *column* vector corresponding to the i -th column of the matrix X and $X_{i\bullet}$ the *row* vector (covector) corresponding to the i -th line of the matrix X . With these notations, if $X \in \mathcal{M}_{m,n}$ and $Z \in \mathcal{M}_{n,m'}$, $X_{i\bullet} Z_{\bullet j} \in \mathbb{R}$ is the inner product between the i -th row of X and the j -th column of Z and $X_{\bullet i} Z_{j\bullet} \in \mathcal{M}_{m,m'}(\mathbb{R})$ is the outer product between the i -th column of X and j -th row of Z .

For the curl-free and divergence-free kernel given in [Subsection 4.3.3](#) we recall the unbounded ORFF maps are respectively for all $y \in \mathcal{Y}$

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle_2 \omega_j^\top y \\ \sin \langle x, \omega_j \rangle_2 \omega_j^\top y \end{pmatrix},$$

³ In many programming language, such as Python, C, C++ or Java each data point is traditionally represented by a row in the data matrix (row major formulation). While this is more natural when parsing a data file, it is less common in mathematical formulations. In this document we adopt the *column major* formulation used by Matlab, Fortran or Julia. Moreover although C++ is commonly row major, some libraries such as Eigen are column major. When dealing with row major formulation, one should “transpose” all the equations given in [Table 6.3](#).

and

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \left(\cos \langle x, \omega_j \rangle_2 \left(\|\omega_j\|_2 I_d - \frac{\omega_j \omega_j^\top}{\|\omega_j\|_2} \right) y + \sin \langle x, \omega_j \rangle_2 \left(\|\omega_j\|_2 I_d - \frac{\omega_j \omega_j^\top}{\|\omega_j\|_2} \right) y \right),$$

where $\omega_j \sim \text{Pr}_{\mathcal{N}(0, \sigma^{-2} I_d)}$. To avoid complex index notations we decompose the feature maps $\tilde{\Phi}(X)$ into two sub feature maps $\tilde{\Phi}^c$ and $\tilde{\Phi}^s$ corresponding to the cosine part and the sine part of each feature map. Namely, for the curl-free kernel, for all $y \in \mathcal{Y}$

$$\tilde{\Phi}(x)y = \begin{cases} \tilde{\Phi}^c(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \left(\cos \langle x, \omega_j \rangle_2 \omega_j^\top y \right), \\ \tilde{\Phi}^s(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \left(\sin \langle x, \omega_j \rangle_2 \omega_j^\top y \right). \end{cases}$$

In the same way, for the divergence-free kernel,

$$\begin{aligned} & \tilde{\Phi}(x)y \\ &= \begin{cases} \tilde{\Phi}^c(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \left(\cos \langle x, \omega_j \rangle_2 \left(\|\omega_j\|_2 I_d - \frac{\omega_j \omega_j^\top}{\|\omega_j\|_2} \right) y \right), \\ \tilde{\Phi}^s(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \left(\sin \langle x, \omega_j \rangle_2 \left(\|\omega_j\|_2 I_d - \frac{\omega_j \omega_j^\top}{\|\omega_j\|_2} \right) y \right). \end{cases} \end{aligned}$$

We also introduce $\tilde{\Phi}^e$, $e \in \{s, c\}$ which denotes either $\tilde{\Phi}^s$ or $\tilde{\Phi}^c$. This equivalent formulation allows us to keep the notation “lighter” and closer to a proper Python/Matlab implementation with vectorization. With these notations, a summary of efficient linear operators in matrix form is given in [Table 6.3](#). The complexity of evaluating all this operators is given in [Table 6.1](#).

It is worth mentioning that the same strategy can be applied in many different language. For instance in C++, the library Eigen [75] allows to wrap a sparse matrix with a custom type, where the user overloads the transpose and dot product operator (as in Python). Then the custom user operator behaves as a (sparse) matrix –see https://eigen.tuxfamily.org/dox/group__MatrixfreeSolverExample.html. With this implementation the time complexity of $\tilde{\Phi}(x)^\top \theta$ and $\tilde{\Phi}(x)y$ falls down to $O_t(Dp' + p'p)$ and the same holds for space complexity.

A quick experiment shows the advantage of seeing the decomposable kernel as a linear operator rather than a matrix. We draw $N = 100$ points $(x_i)_{i=1}^N$ in the interval $(0, 1)^{20}$ and use a decomposable kernel with matrix $\Gamma = BB^\top \in \mathcal{M}_{p,p}(\mathbb{R})$ where $B \in \mathcal{M}_{p,p}(\mathbb{R})$ is a random matrix with coefficients drawn uniformly in $(0, 1)$. We compute $\tilde{\Phi}(x)^\top \theta$ for all x_i 's, where $\theta \in \mathcal{M}_{2D,1}(\mathbb{R})$, $D = 100$, with the implementation `EfficientDecomposableGaussianORFF`, [Equation 6.14](#), and `NaiveDecomposableGaussianORFF`, [Equation 6.13](#). The coefficients of θ were

Table 6.1: Complexity of efficient linear-operator (in matrix form) for different Feature maps given in [Table 6.3](#).

Kernel	$\tilde{\Phi}(X)^*$	$\tilde{\Phi}(X)$
Decomposable	$O((p'D + p'p)N)$	$O((pN + p'p)D)$
Curl-free	$O(pND)$	$O(pND)$
Divergence-free	$O((p^2 + pN)D)$	$O((p^2 + pN)D)$

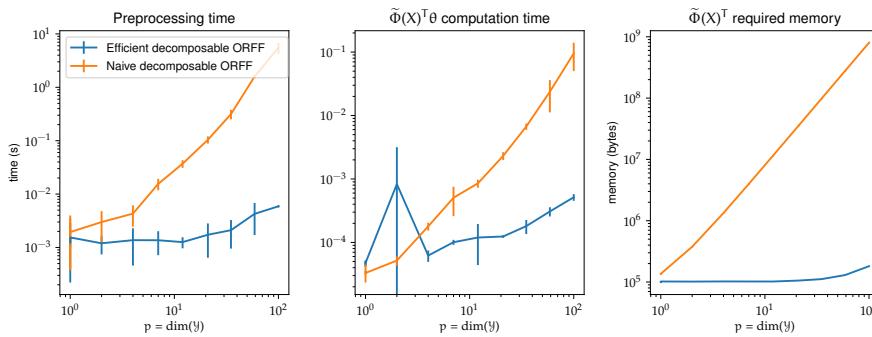


Figure 6.3: Efficient decomposable Gaussian ORFF (lower is better).

Table 6.3: Efficient linear-operator (in matrix form) for different Feature maps.

Kernel	$\tilde{\Phi}(X)^*$	$\tilde{\Phi}(X)$
Decomposable ¹	$\Theta \mapsto B(\Theta \tilde{\varphi}(X))$	$Y \mapsto B^T(Y \tilde{\varphi}(X)^T)$
Gaussian curl-free ²	$\Theta^c, \Theta^s \mapsto \sum_{j=1}^D \omega_j (\Theta_j^c \tilde{\varphi}^c(X)_{j\bullet} + \Theta_j^s \tilde{\varphi}^s(X)_{j\bullet})$	$Y \mapsto \Theta_j^e = \omega_j^T (Y \tilde{\varphi}^e(X)_{\bullet j}^T)$
Gaussian divergence-free ^{2,3}	$\Theta^c, \Theta^s \mapsto \sum_{j=1}^D (B(\omega_j) \Theta_{\bullet j}^c) \tilde{\varphi}^c(X)_{j\bullet} + (B(\omega_j) \Theta_{\bullet j}^s) \tilde{\varphi}^s(X)_{j\bullet}$	$Y \mapsto \Theta_{\bullet j}^e = B(\omega_j) (Y \tilde{\varphi}^e(X)_{\bullet j}^T)$

¹ Where $\tilde{\varphi}(X) = (\tilde{\varphi}(X_{\bullet 1}) \dots \tilde{\varphi}(X_{\bullet N})) \in \mathcal{M}_{r,N}$ is any design matrix, with scalar feature map $\tilde{\varphi} : \mathbb{R}^d \rightarrow \mathbb{R}^r$ such that $\tilde{\varphi}(x)^* \tilde{\varphi}(z) = k(x, z) \in \mathbb{R}$ for all $x, z \in \mathcal{X}$. The input data $X \in \mathcal{M}_{d,N}(\mathbb{R})$, the output data $U \in \mathcal{M}_{p,N}(\mathbb{R})$, the parameter matrices Θ^c and $\Theta^s \in \mathcal{M}_{p',r}(\mathbb{R})$ and the decomposable operator $B \in \mathcal{M}_{p,p'}(\mathbb{R})$.

² Where $\tilde{\varphi}^c(X)_{ji} = \cos(\omega_j, x_i)$ and $\tilde{\varphi}^s(X)_{ji} = \sin(\omega_j, x_i)$, $j \in \mathbb{N}_D^*$ and $i \in \mathbb{N}_N^*$. Thus $\tilde{\varphi}^c(X) \in \mathcal{M}_{D,N}(\mathbb{R})$ and $\tilde{\varphi}^s(X) \in \mathcal{M}_{D,N}(\mathbb{R})$. The input data $X \in \mathcal{M}_{d,N}(\mathbb{R})$, the output data $U \in \mathcal{M}_{d,N}(\mathbb{R})$, the parameter matrices Θ^c and $\Theta^s \in \mathbb{R}^D$, $\omega_j \sim \text{Pr}_{\mathcal{N}(0, \sigma^{-2} I_d)}$ i.i.d. for all $j \in \mathbb{N}_D^*$.

Eventually $e \in \{s, c\}$, namely $\Theta^c = (\Theta_1^{e=c} \dots \Theta_D^{e=c})^T$ and $\Theta^s = (\Theta_1^{e=s} \dots \Theta_D^{e=s})^T$.

³ Here, Θ^c and $\Theta^s \in \mathcal{M}_{d,D}(\mathbb{R})$ thus $\Theta^c = (\Theta_{\bullet 1}^{e=c} \dots \Theta_{\bullet D}^{e=c})$, $\Theta^s = (\Theta_{\bullet 1}^{e=s} \dots \Theta_{\bullet D}^{e=s})$ and $B(\omega) = \left(\|\omega\|_2 I_d - \frac{\omega \omega^T}{\|\omega\|_2} \right) \in \mathcal{M}_{d,d}$.

drawn at random uniformly in $(0, 1)$. We report the execution time in [Figure 6.3](#) for different values of p , $1 \leq p \leq 100$. The left plot reports the execution time in seconds of the construction of the feature. The middle plot reports the execution time of $\tilde{\Phi}(x)^T \theta$, and the right plot the memory used in bytes to store $\tilde{\Phi}(x)$ for all x_i 's. We averaged the results over ten runs. Full code is given in [Appendix C.3](#).

6.3.3 Curl-free kernel

We use the unbounded ORFF map presented in [Equation 4.20](#). We draw $N = 1000$ points $(x_i)_{i=1}^N$ in the interval $(0, 1)^p$ and use a curl-free kernel. We compute $\tilde{\Phi}(x)^T \theta$ for all x_i 's, where $\theta \in \mathcal{M}_{2D,1}(\mathbb{R})$, $D = 500$, with the matrix implementation and the `LinearOperator` implementation. The coefficients of θ were drawn at random uniformly in $(0, 1)$. We report the execution time in [Figure 6.4](#) for different values of p , $1 \leq p \leq 100$. The left plot reports the execution time in seconds of the

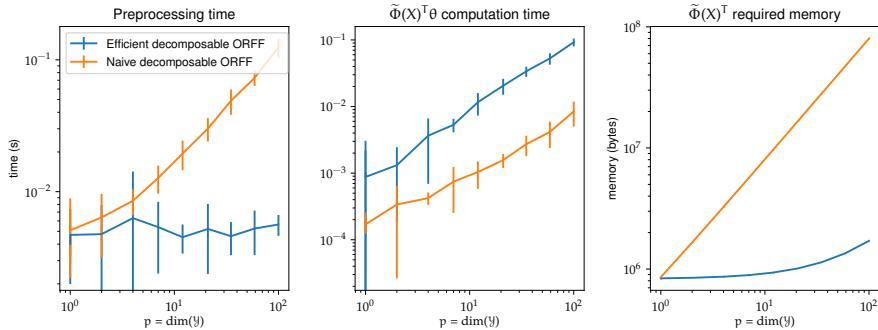


Figure 6.4: Efficient curl-free Gaussian ORFF (lower is better).

construction of the features. The middle plot reports the execution time of $\tilde{\Phi}(x)^T \theta$, and the right plot the memory used in bytes to store $\tilde{\Phi}(x)$ for all x_i 's. We averaged the results over fifty runs. Full code is given in [Appendix C.4](#). As we can see the linear-operator implementation is one order of magnitude slower than its matrix counterpart. However it uses considerably less memory.

6.3.4 Divergence-free kernel

We use the unbounded ORFF map presented in [Equation 4.21](#). We draw $N = 100$ points $(x_i)_{i=1}^N$ in the interval $(0, 1)^p$ and use a curl-free kernel. We compute $\Phi(x)^T \theta$ for all x_i 's, where $\theta \in \mathcal{M}_{2Dp,1}(\mathbb{R})$, $D = 100$, with the matrix implementation and the `LinearOperator` implementation. The coefficients of θ were drawn at random uniformly in $(0, 1)$. We report the execution time in [Figure 6.4](#) for different values of p , $1 \leq p \leq 100$. The left plot reports the execution time in seconds of the construction of the feature. The middle plot reports

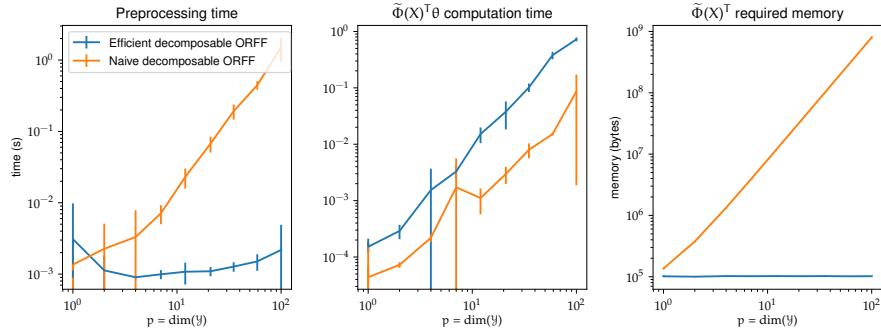


Figure 6.5: Efficient divergence-free Gaussian ORFF (lower is better).

the execution time of $\tilde{\Phi}(x)^T \theta$, and the right plot the memory used in bytes to store $\tilde{\Phi}(x)$ for all x_i 's. We averaged the results over ten runs. Full code is given in [Appendix C.5](#). We draw the same conclusions as the curl-free kernel.

6.3.5 Iterative solvers using linear operators

We have shown in this section that viewing the Operator-valued Random Fourier Feature maps as linear operator rather than staking matrices leads to more efficient computations.

6.4 EXPERIMENTS

We present a set of experiments to complete the theoretical contribution and illustrate the behavior of ORFF-regression. First we study how well the ORFF regression recover the result of operator-valued kernel regression. Second we show the advantages of ORFF regression over independent RFF regression. A code implementing ORFF is available in Operalib, a framework for machine learning with OVKs, at <https://github.com/operalib/operalib>

6.4.1 Learning with ORFF vs learning with OVK

6.4.1.1 Datasets

The *first dataset* considered is the handwritten digits recognition dataset MNIST available at <http://yann.lecun.com/exdb/mnist>. We select a training set of 12,000 images and a test set of 10,000 images. The inputs are images represented as a vector $x_i \in [0, 255]^{784}$ and the targets $y_i \in \mathbb{N}_9$ are integers between 0 and 9.

First we scaled the inputs such that they take values in $[-1, 1]^{784}$. Then we binarize the targets such that each number is represented by

a unique binary vector of dimension 10. The vector y_i is zero everywhere except on the dimension corresponding to the class where it is one. For instance the class 4 is encoded

$$(0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0)^\top.$$

To predict classes, we use the simplex coding method presented in Mroueh et al. [123]. The intuition behind simplex coding is to project the binarized labels of dimension p onto the most separated vectors on the hypersphere of dimension $p - 1$. For ORFF we can encode directly this projection in the B matrix of the decomposable kernel $K_0(\delta) = BB^*k_0(\delta)$ where k_0 is a Gaussian kernel. The matrix B is computed via the recursion

$$B_{p+1} = \begin{pmatrix} 1 & u^\top \\ 0_{p-1} & \sqrt{1-p^{-2}}B_p \end{pmatrix}, \quad B_2 = \begin{pmatrix} 1 & -1 \end{pmatrix},$$

where $u = (-p^{-2} \ \dots \ -p^{-2})^\top \in \mathbb{R}^{p-1}$ and $0_{p-1} = (0 \ \dots \ 0)^\top \in \mathbb{R}^{p-1}$. For Operator-Valued Kernels we project the binarized targets on the simplex as a preprocessing step, before learning with the decomposable $K_0(\delta) = I_p k_0(\delta)$, where k_0 is a scalar Gaussian kernel.

The *second dataset* is a simulated five dimensional (5D) vector field with structure. We generate a scalar field as a random function $f : [-1, 1]^5 \rightarrow \mathbb{R}$, where $\tilde{f}(x) = \tilde{\varphi}(x)^*\theta$ where θ is a random matrix with each entry following a standard normal distribution, $\tilde{\varphi}$ is a scalar Gaussian RFF with bandwidth $\sigma = 0.4$. The input data x are generated from a uniform probability distribution. We take the gradient of \tilde{f} to generate the curl-free 5D vector field.

The *third dataset* is a synthetic of data from $\mathbb{R}^{20} \rightarrow \mathbb{R}^4$ as described in Audiffren and Kadri [10]. In this dataset, inputs (x_1, \dots, x_{20}) are generated independently and uniformly over $[0, 1]$ and the different outputs are computed as follows. Let

$$\varphi(x) = (x_1^2, x_4^2, x_1x_2, x_3x_5, x_2, x_4, 1)$$

and (w_i) denotes the i.i.d. copies of a seven dimensional Gaussian distribution with zero mean and covariance $\Sigma \in \mathcal{M}_{7,7}(\mathbb{R})$ such that

$$\Sigma = \text{Diag} \left(0.5 \ 0.25 \ 0.1 \ 0.05 \ 0.15 \ 0.1 \ 0.15 \right)$$

Then, the outputs of the different tasks are generated as $y_i = w_i \varphi(x)$. We use this dataset with $p = 4, 10^5$ instances and for the train set and also 10^5 instances for the test set.

6.4.1.2 Results

Performance of ORFF regression on the first dataset. We trained both ORFF and OVK models on MNIST dataset with a decomposable Gaussian kernel with signature

$$K_0(\delta) = \exp(-\|\delta\|/(2\sigma^2)) \Gamma.$$

To apply [Algorithm 3](#) after noticing that in the case of the decomposable kernel with $\lambda_M = 0$, it boils down to a Stein equation [29, section 5.1], we use an off-the-shelf solver⁴ able to handle Stein's equation. For both methods we choose $\sigma = 20$ and use a 2-fold cross validation on the training set to select the optimal λ . First, [Figure 6.6](#) compares the running time between OVK and ORFF models using $D = 1000$ Fourier features against the number of datapoints N . The log-log plot shows ORFF scaling better than the OVK w. r. t. the number of points. Second, [Figure 6.6](#) shows the test prediction error versus the number of ORFFs D , when using $N = 1000$ training points. As expected, the ORFF model converges toward the OVK model when the number of features increases.

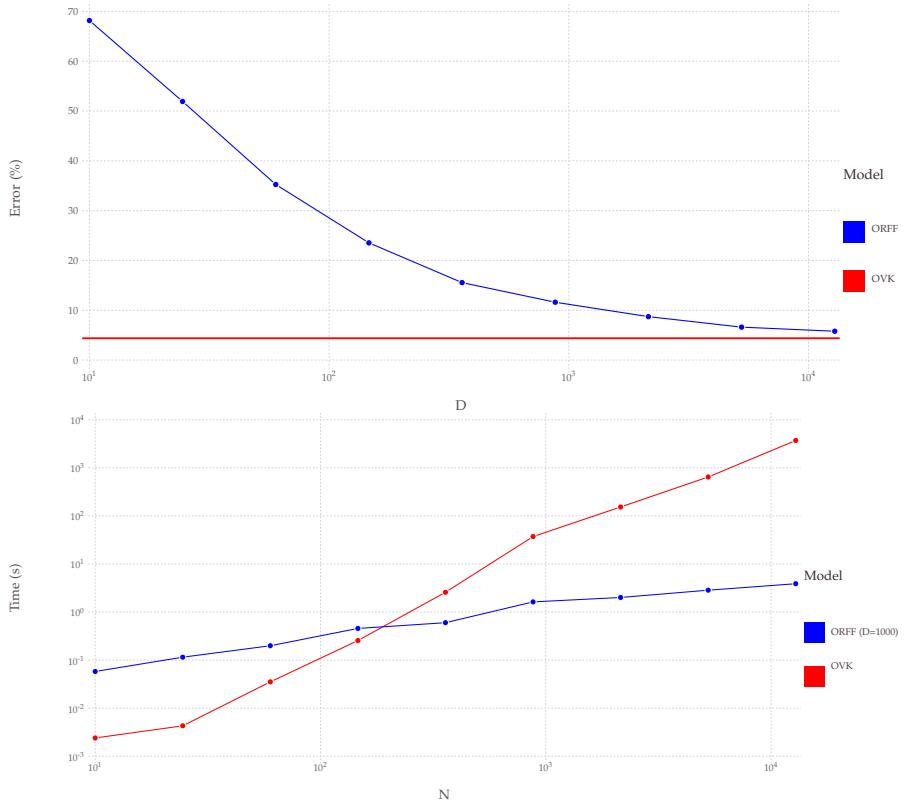


Figure 6.6: Empirical comparison of ORFF and OVK regression on MNIST dataset and empirical behavior of ORFF regression versus D and N .

⁴ Available at <http://ta.twi.tudelft.nl/nw/users/gijzen/IDR.html>

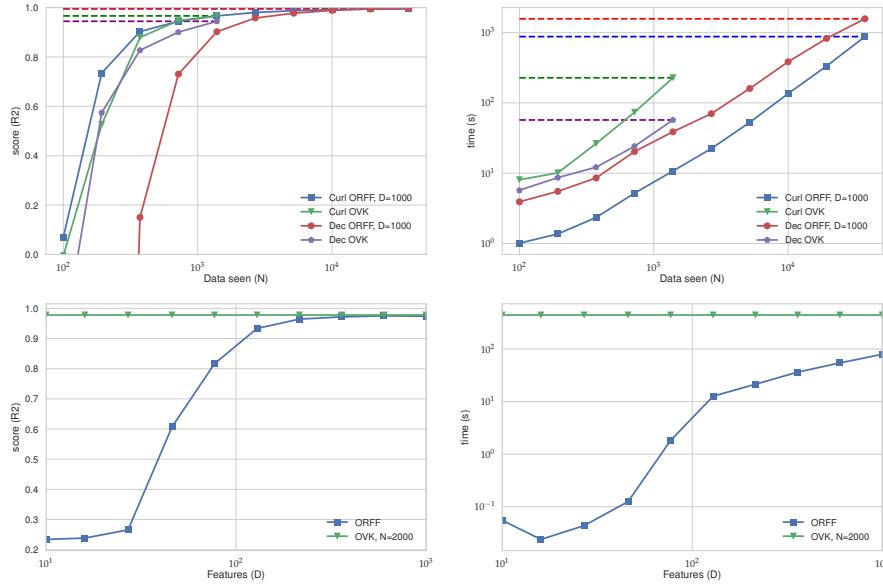


Figure 6.7: Empirical comparison between curl-free ORFF, curl-free OVK, independent ORFF, independent OVK on a synthetic vector field regression task.

Performance of ORFF regression on the second dataset. We perform a similar experiment on the second dataset (5D-vector field with structure). We use a Gaussian curl-free kernel with bandwidth equal to the median of the pairwise distances and tune the hyperparameter λ on a grid. Here we optimize Equation 6.9 using Scipy’s L-BFGS-B [37] solver⁵ with the gradients given in Equation 6.11 and the efficient linear operator described in Subsection 6.3.2. Figure 6.7 (bottom row) reports the R^2 score on the test set versus the number of curl-ORFF D with a comparison with curl-OKV. In this experiment, we see that curl-ORFF can even be better than curl-OKV, suggesting that ORFF might play an additional regularizing role. It also shows the computation time of curl-ORFF and curl-OKV. We see that OKV regression does not scale with large datasets, while ORFF regression does. When $N > 10^4$, OKV regression exceeds memory capacity.

Structured prediction vs Independent (RFF) prediction. On the second dataset, Figure 6.7 (top row) compares R^2 score and time of ORFF regression using the trivial identity decomposable kernel, e.g. independent RFFs, to curl-free ORFF regression. Curl-free ORFF outperforms independent RFFs, as expected, since the dataset involves structured outputs.

Impact of the number of random features (D). In this setting we solved the optimisation problem for both ORFF and OKV using a L-BFGS-B. Figure 6.8 top row shows that for a fixed number of instance

⁵ Available at <http://docs.scipy.org/doc/scipy/reference/optimize.html>

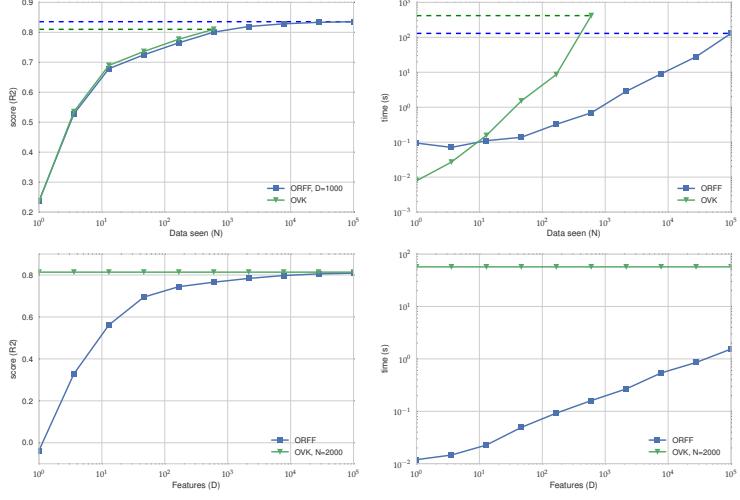


Figure 6.8: Decomposable kernel on the third dataset: R^2 score vs number of data in the train set (N)

in the train set, OVK performs better than ORFF in terms of accuracy (R^2). However ORFF scales better than OVK w. r. t. the number of data. ORFF is able to process more data than OVK in the same time and thus reach a better accuracy for a given amount of time. Bottom row shows that ORFF tends to reach OVK's accuracy for a fixed number of data when the number of features increase.

Multitask learning. In this experiment we are interested in multitask learning with operator-valued random Fourier features, and see whether the approximation of a joint OVK performs better than an independent OVK. In this setting we assume that for each entry $x_i \in \mathbb{R}^d$ we only have access to one observation $y_i \in \mathbb{R}$ corresponding to a task t_i . We used the SARCOS dataset, taken from <http://www.gaussianprocess.org/gpml/data/> website. This is an inverse dynamics problem, i. e. we have to predict the 7 joint torques given the joint positions, velocities and accelerations. Hence, we have to solve a regression problem with 21 inputs and 7 outputs which is a very nonlinear function. It has 45K inputs data. Suppose that we are given a collection of inputs data $x_1, \dots, x_N \in \mathbb{R}^{21}$ and a collection of output data $((y_1, t_1), \dots, (y_N, t_N)) \in (\mathbb{R} \times \mathbb{N}_T)^N$ where T is the number of tasks. We consider the following multitask loss function

$$L(h(x), (y, t)) = \frac{1}{2} (\langle h(x), e_t \rangle_2 - y)^2,$$

This loss function is adapted to datasets where the number of data per tasks is unbalanced (i. e. for one input data we observe the value

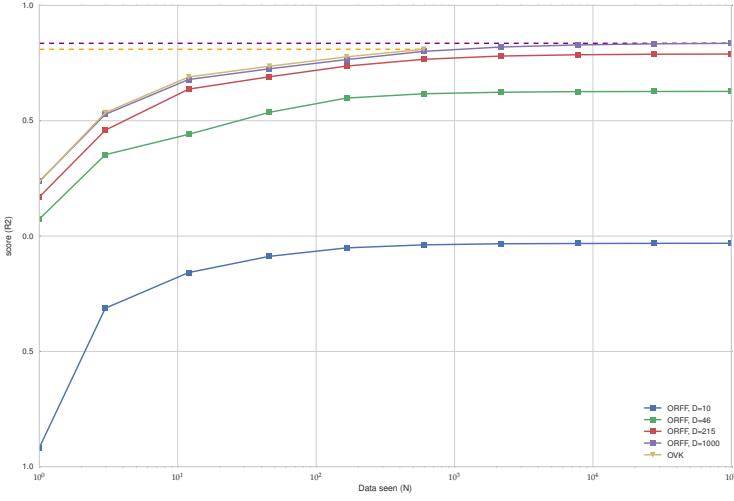


Figure 6.9: Decomposable kernel on the third dataset: R^2 score vs number of data in the train set (N) for different number of random samples (D).

of only one task and not all the tasks.). We optimise the regularized risk

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N L(h(x_i), (y_i, t_i)) \\ + \frac{\lambda}{2N} \|h\|_{\mathcal{H}}^2 = \frac{1}{2N} \sum_{i=1}^N (\langle h(x_i), e_{t_i} \rangle - y_i)^2 + \frac{\lambda}{2N} \|h\|_{\mathcal{H}}^2 \end{aligned}$$

We used a model h based on the decomposable kernel

$$h(x) = (\varphi(x)^T \otimes B)\theta$$

we chose B such that $BB^T = A$, where A is the inverse graph Laplacian L of the similarities between the tasks, parametrized by an hyperparameter $\gamma \in \mathbb{R}_+$. Namely:

$$L_{kl} = \exp \left(-\gamma \sqrt{\sum_{i=1}^N (y_i^k - y_i^l)^2} \right).$$

We draw N data randomly for each task, hence creating a dataset of $N \times 7$ data and computed the nMSE on the proposed test set (4.5K points). We repeated the experiments 80 times to avoid randomness. We choose $D = \frac{\max(N, 500)}{2}$ features, and optimized the problem with a second order batch gradient. Table 6.5 shows that using the ORFF approximation of an operator-valued kernel with a good prior on the data improves the performances w.r.t. the independent ORFF. However the advantage seems to be less important the more data are available.

N	Independant (%)	Laplacian (%)	p-value	T
50×7	23.138 ± 0.577	22.254 ± 0.536	2.68%	4(s)
100×7	16.191 ± 0.221	15.568 ± 0.187	< 0.1%	16(s)
150×7	13.821 ± 0.115	13.459 ± 0.106	< 0.1%	13(s)
200×7	12.713 ± 0.0978	12.554 ± 0.0838	1.52%	12(s)
400×7	10.785 ± 0.0579	10.651 ± 0.0466	< 0.1%	10(s)
800×7	7.512 ± 0.0344	7.512 ± 0.0344	100%	15(s)
1600×7	6.486 ± 0.0242	6.486 ± 0.0242	100%	20(s)
3200×7	5.658 ± 0.0187	5.658 ± 0.0187	100%	20(s)

Table 6.5: Error (%) of nMSE on SARCOS dataset.

6.5 CONCLUSION

ORFF approximations open the door to the literature of efficient learning with linear models: the feature map can be seen as a function that linearize non linear functions by embedding them in a high dimensional feature space in which we can learn with linear models. Indeed, we described in this chapter how learning a non linear, non parametric model with OVKs is converted into learning a linear and parametric model based on ORFF.

The complexity in time of these approximations together with the linear learning algorithm make this implementation scalable with the data size and thus appealing compared to OVK regression as shown in numerical experiments. Further work concerns generalization bounds and consistency for ORFF-regression. Finally this work opens the door to building deeper architectures by stacking vector-valued functions while keeping a kernel view for large datasets.



7

CONSISTENCY AND GENERALIZATION BOUND FOR ORFF

This short chapter deals with a generalization bound for a regression problem with ORFF based on the results of Maurer [112] and Rahimi and Recht [140]. We also discuss the case of Ridge regression presented in [Chapter 6](#).

Contents

7.1	Generalization bound	126
7.1.1	Generalization by bounding the function space complexity	127
7.1.2	Algorithm stability	130
7.2	Consistency of learning with ORFF	131
7.3	Discussion	135

7.1 GENERALIZATION BOUND

In this section, we are interested in finding a function $f_* : \mathcal{X} \rightarrow \mathcal{Y}$, where \mathcal{X} is a Polish space and \mathcal{Y} a separable Hilbert space such that for all x_i in \mathcal{X} and all y_i in \mathcal{Y} that minimizes a criterion. In statistical supervised learning, we consider a training set sequence $s = (x_i, y_i)_{i=1}^N \in (\mathcal{X} \times \mathcal{Y})^N$, $N \in \mathbb{N}^*$ drawn i.i.d. from an unknown probability law \mathbf{Pr} . We suppose we are given a cost function $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, such that $c(f(x), y)$ returns the error of the prediction $f(x)$ w.r.t. the ground truth y . We define the true risk as the sum of the cost over all possible training examples drawn from a latent probability law \mathbf{Pr} ,

$$\mathfrak{R}(f) = \int_{\mathcal{X} \times \mathcal{Y}} L(x, f, y) d\mathbf{Pr}(x, y) = \int_{\mathcal{X} \times \mathcal{Y}} c(f(x), y) d\mathbf{Pr}(x, y)$$

Thus given a class of functions \mathcal{F} , the goal of a learning algorithm is to find an optimal model f_* that minimizes the true risk. Namely

$$f_* \in \arg \min_{f \in \mathcal{F}} \mathfrak{R}(f) = \arg \min_{f \in \mathcal{F}} \int_{\mathcal{X} \times \mathcal{Y}} c(f(x), y) d\mathbf{Pr}(x, y).$$

Since in practice we do not have access to the joint probability law of (X, Y) , we define its empirical counterpart as the empirical mean estimate, where the sequence $s = (x_i, y_i)_{i=1}^N$ is made of $(\mathcal{X} \times \mathcal{Y})$ -valued random vectors drawn i.i.d. from some law \mathbf{Pr} . The empirical risk then reads

$$\mathfrak{R}_{\text{emp}}(f, s) = \frac{1}{N} \sum_{i=1}^N c(f(x_i), y_i), \quad (x_i, y_i) \sim \mathbf{Pr} \text{ i.i.d.}$$

As a result, in practice we seek a function f_s such that

$$f_s \in \arg \min_{f \in \mathcal{F}} \mathfrak{R}_{\text{emp}}(f, s) = \arg \min_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N c(f(x_i), y_i). \quad (7.1)$$

The basic requirement for any learning algorithm is the generalization property: the empirical error must be a good proxy of the expected error, that is the difference between the two must be “small” when N is large. A generalization bound allows to study, for any $f \in \mathcal{F}$ the difference between its true risk $\mathfrak{R}(f)$ and its empirical risk, $\mathfrak{R}(f, s)$. This quantifies the impact of having a limited number of observations. Generalization (upper) bounds [177] involve two components: one being the empirical risk and the other depends on the dataset size as well as some capacity notion that reflects the richness of the family of functions \mathcal{F} considered. First generalization bounds proved by Vapnik and Chervonenkis involve the dimension of Vapnik-Chervonenkis dimension of \mathcal{F} . In practice, generalization bounds suggest that when learning a function from a finite dataset,

it is necessary to control the size (richness) of the class of functions \mathcal{F} . Hence, a regularizer is added to the data-fitting term in order to maintain the solution f_s of [Equation 7.1](#) unique and belong to a ball of \mathcal{F} . As a result if \mathcal{F} is a Banach space, it is common to find f_s such that

$$f_s = \arg \min_{f \in \mathcal{F}} \mathfrak{R}_{\text{emp}}(f, s) + \frac{\lambda}{2} \|f\|_{\mathcal{F}}^2.$$

(Tychonov regularization) or

$$f_s = \begin{cases} \arg \min_{f \in \mathcal{F}} \mathfrak{R}_{\text{emp}}(f, s) \\ \text{subject to } \|f\|_{\mathcal{F}} < C \in \mathbb{R}_{>0} \end{cases}$$

(Ivanov regularization) or

$$f_s = \begin{cases} \arg \min_{f \in \mathcal{F}} \mathfrak{R}_{\text{emp}}(f, s) \\ \text{subject to } \|f\|_{\infty} < C \in \mathbb{R}_{>0}. \end{cases}$$

7.1.1 Generalization by bounding the function space complexity

In the following we consider functions living in a Vector Valued Reproducing Kernel Hilbert Space, with kernel K (or \tilde{K}).

Proposition 7.1 (Bartlett and Mendelson [17] and Maurer [112]).
Suppose that $f \in \mathcal{H}_K$ a VV-RKHS where

$$\sup_{x \in \mathcal{X}} \text{Tr}[K(x, x)] < T$$

and $\|f\|_{\mathcal{H}_K} < B$. Moreover let $c : \mathcal{Y} \rightarrow [0, C]$ be a L-Lipschitz cost function and \mathcal{Y} a separable Hilbert space. Then if we are given N i. i. d. random variables with values in \mathcal{X} (training samples, noted s), then we have with at least probability $1 - \delta$, $\delta \in (0, 1)$ over the drawn training samples s that for any $f \in \mathcal{H}_K$,

$$\mathfrak{R}(f) \leq \mathfrak{R}_{\text{emp}}(f, s) + 2\sqrt{\frac{2}{N}} \left(LB T^{1/2} + C \sqrt{\ln(2/\delta)} \right). \quad (7.2)$$

The following proof is due to Maurer [112] generalizing the work of Bartlett and Mendelson [17, section 4.3]: we do not claim any originality for this proof.

Proof First let us introduce the notion of Rademacher complexity of a class of functions \mathcal{F} . We recall that the probability mass function of a uniformly distributed Rademacher random variable is given for any $k \in \{-1, 1\}$ by

$$f(k) = \begin{cases} 1/2 & \text{if } k = -1 \\ 1/2 & \text{otherwise.} \end{cases}$$

Definition 7.1 (Bartlett and Mendelson [17]). Let \mathcal{X} be any set. Let $\epsilon_1, \dots, \epsilon_N$ be N independent Rademacher random variables, identically uniformly distributed on $\{-1; 1\}$. For any class of functions $\mathcal{F} : \mathcal{X} \rightarrow \mathbb{R}$, then for all $x_1, \dots, x_N \in \mathcal{X}$ the quantity

$$\mathcal{R}_N(\mathcal{F}) := E \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^N \epsilon_i f(x_i) \mid x_1, \dots, x_N \right]$$

is called Rademacher complexity of the class \mathcal{F} .

In a few words the Rademacher complexity measures the richness of a class a function by its capacity to be correlated to noise. In generalization bounds, the Rademacher complexity of a class of functions often involves a composition between a target function to be learn and a cost function, part of the risk we want to minimize. The idea is to bound the Rademacher complexity with a term that does not depends on the cost function, but only on the target function.

Proposition 7.2 (Maurer [112]). Let \mathcal{X} be any set and (x_1, \dots, x_N) in \mathcal{X}^N and let \mathcal{F} be a class of functions $f : \mathcal{X} \rightarrow \mathcal{Y}$ and for $i=1, \dots, N$, each function $h_i : \mathcal{Y} \rightarrow \mathbb{R}$ be a L -Lipschitz function, where \mathcal{Y} is a separable Hilbert space endowed with Euclidean inner product. Then

$$\begin{aligned} & E \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^N \epsilon_i h_i(f(x_i)) \mid x_1, \dots, x_N \right] \\ & \leq \sqrt{2L} E \left[\sup_{f \in \mathcal{F}} \sum_{i=1, k}^{i=N} \epsilon_{ik} f_k(x_i) \mid x_1, \dots, x_N \right], \end{aligned}$$

where ϵ_{ik} is a doubly indexed independent Rademacher sequence and $f_k(x_i)$ is the k -th component of $f(x_i)$. We use the shortcut notation $\sum_{i=1, k}^N$ which stands for $\sum_{i=1}^N \sum_k$.

From now on, we consider functions $f \in \mathcal{H}_K$ a Vector Valued Reproducing Kernel Hilbert Space. Then there exists an induced feature-map $\Phi : \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}, \mathcal{H})$ such that for all $y, y' \in \mathcal{Y}$ the kernel is given by

$$\langle y, K(x, z)y' \rangle = \langle \Phi_x y, \Phi_z y' \rangle.$$

We say that the feature space \mathcal{H} is embedded into the RKHS \mathcal{H}_K by means of the feature operator $(W\theta)(x) := (\Phi_x^* \theta)$. Indeed W defines a partial isometry between \mathcal{H} and \mathcal{H}_K . Suppose that \mathcal{Y} is a separable Hilbert space and let the class of \mathcal{Y} -valued functions \mathcal{F} be

$$\mathcal{F} = \{ f \mid f : x \mapsto (W\theta)(x), \|\theta\|_{\mathcal{H}} < B \} \subset \mathcal{H}_K.$$

Let $c_{y_i} = c(\cdot - y_i)$, for all $i \in \mathbb{N}_N$. Then from Proposition 7.2 and if K is trace class, we have

$$\begin{aligned} E \sup_{\|\theta\|_{\mathcal{H}} < B} \sum_{i=1}^N \epsilon_i c_{y_i}(\Phi_{x_i}^* \theta) & \leq \sqrt{2L} E \sup_{\|\theta\|_{\mathcal{H}} < B} \sum_{i=1, k}^{i=N} \epsilon_{ik} \langle \Phi_{x_i}^* \theta, e_k \rangle \\ & = \sqrt{2L} E \sup_{\|\theta\|_{\mathcal{H}} < B} \left\langle \theta, \sum_{i=1, k}^{i=N} \epsilon_{ik} \Phi_{x_i} e_k \right\rangle_y. \end{aligned}$$

Thus

$$\begin{aligned}
E \sup_{\|\theta\|_{\mathcal{H}} < B} \sum_{i=1}^N \epsilon_i c_{y_i}(\Phi_{x_i}^* \theta) &\leq \sqrt{2LB} E \left\| \sum_{i=1, k}^{i=N} \epsilon_{ik} \Phi_{x_i} e_k \right\|_y \\
&\leq \sqrt{2LB} \sqrt{\sum_{i=1, k}^{i=N} \|\Phi_{x_i} e_k\|_y^2} \\
&\leq \sqrt{2LB} \sqrt{\sum_{i=1}^N \text{Tr}[K(x_i, x_i)]} \\
&\leq \sqrt{2LB} \sqrt{N} \sqrt{\sup_{x \in \mathcal{X}} \text{Tr}[K(x, x)]}.
\end{aligned} \tag{7.3}$$

Then we apply the following theorem (Theorem 7.1) from Bartlett and Mendelson [17] and Maurer [112] to conclude.

Theorem 7.1 Let \mathcal{X} be any set, \mathcal{F} a class of functions $f : \mathcal{X} \rightarrow [0, C]$ and let X_1, \dots, X_N be a sequence of i.i.d. random variables with value in \mathcal{X} . Then for $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have for all $f \in \mathcal{F}$ that

$$Ef(X) \leq \frac{1}{N} \sum_{i=1}^N f(X_i) + \frac{2}{N} \mathcal{R}_N(\mathcal{F}) + C \sqrt{\frac{8 \ln(2/\delta)}{N}} \tag{7.4}$$

Conclude by plugging Equation 7.3 in Theorem 7.1. \square

As an example, let us consider Equation 6.1, which is a solution of the regularized empirical risk, and Algorithm 3. We first list the following assumptions useful in the rest of the section. Let $s = (x_i, y_i)_{i=1}^N \in \mathcal{X}^N \times \mathcal{Y}^N$ be the training samples.

Assumption 7.1 There exists a positive constant $\kappa \in \mathbb{R}_{\geq 0}$ such that

$$\max_{i \in \mathbb{N}_N^*} \left\| \tilde{K}(x_i, x_i) \right\|_{y, y} < \kappa.$$

Assumption 7.2 There exists a positive constant $T \in \mathbb{R}_{\geq 0}$ such that

$$\max_{i \in \mathbb{N}_N^*} \text{Tr}[\tilde{K}(x_i, x_i)] < T.$$

Assumption 7.3 There exists a positive constant $C \in \mathbb{R}_{\geq 0}$ such that

$$\max_{i \in \mathbb{N}_N^*} \|y_i\|_y \leq C.$$

Assumption 7.4 Given a Loss function L , there exists a positive constant $\xi \in \mathbb{R}_{\geq 0}$ such that for all $x \in \mathcal{X}$, for all $y \in \mathcal{Y}$ and for any $s \in \mathcal{X}^N \times \mathcal{Y}^N$,

$$L(x, f_s, y) \leq \xi.$$

Under [assumption 7.3](#), from [Remark 6.1](#), we know that $\|f\|_{\mathcal{H}_K} \leq \sqrt{\frac{2}{\lambda}}\sigma_y$, where

$$\frac{1}{N} \sum_{i=1}^N \|y_i\|_y^2 \leq \sigma_y^2 \leq C^2.$$

Thus we see straight away that it is possible to choose $B = \sqrt{\frac{2}{\lambda}}C$. Let $\kappa = \|\tilde{K}_e(e)\|_{y,y}$ the Lipschitz constant of the least square loss $c(f_s(x), y) = \frac{1}{2}\|f_s(x) - y\|_y^2$ with respect to $f_s(x)$ is $L = \max\left(\sqrt{\frac{2\kappa}{\lambda}}C, C\right)$ and the loss takes values in $[0, \frac{1}{2}L^2]$. Hence under assumption that $\lambda < 2\kappa$ and [assumption 7.2](#), and [assumption 7.3](#), [Equation 7.2](#) applies especially that for any $f_s \in \mathcal{H}_K$, solution of [Algorithm 3](#),

$$\mathfrak{R}(f_s) \leq \mathfrak{R}_{\text{emp}}(f_s, s) + 8 \frac{C^2}{\lambda} \sqrt{\frac{\kappa}{N}} \left(T^{1/2} + \sqrt{\frac{\kappa \ln(1/\delta)}{2}} \right). \quad (7.5)$$

This bound is to be compared to the results of Kadri et al. [86] in the context of β -stability.

7.1.2 Algorithm stability

²⁰ Other methods using covering numbers [171, 194] or VC-dimension [178] have also been used as a proxy on the complexity of the hypothesis space.

The approach to generalization bounds presented in [Proposition 7.1](#) is based on controlling the complexity of the hypothesis space²⁰ using Rademacher complexity. On the other hand, the idea of stability is that a reliable algorithm should not change its solution too much if we modify slightly the training data. Given a training sequence

$$s = ((x_1, y_1), \dots, (x_N, y_N)) \in (\mathcal{X} \times \mathcal{Y})^N,$$

we note $s^{\setminus i}$ the training sequence

$$s^{\setminus i} = ((x_1, y_1), \dots, (x_{i-1}, y_{i-1}), (x_{i+1}, y_{i+1}), \dots, (x_N, y_N)) \in (\mathcal{X} \times \mathcal{Y})^N,$$

the subsequence of s from which we removed the i -th element.

Definition 7.2 (Uniform stability Bousquet and Elisseeff [26, definition 6]). A learning algorithm $s \mapsto f_s$ has uniform stability β with respect to the loss function L if the following holds

$$\forall i \in \mathbb{N}_N^*, \forall s \in (\mathcal{X} \times \mathcal{Y})^N \sup_{x \in \mathcal{X}, y \in \mathcal{Y}} |L(x, f_s, y) - L(x, f_{s^{\setminus i}}, y)| \leq \beta.$$

As shown by Bousquet and Elisseeff [26], algorithm stability has direct link with generalization. Indeed if an algorithm has β -stability, and a “bounded” loss for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ ([assumption 7.4](#)), it is possible to exhibit a generalization bound.

Theorem 7.2 (Bousquet and Elisseeff [26, theorem 12]). Let $s \mapsto f_s$ be a learning algorithm with uniform stability β with respect to a loss L that satisfies [assumption 7.4](#). Then $\forall N \in \mathbb{N}^*, \forall \delta \in (0, 1)$, the following bound holds with probability at least $1 - \delta$ over the i.i.d. drawn training samples s .

$$\mathfrak{R}(f_s) \leq \mathfrak{R}_{\text{emp}}(f_s, s) + 2\beta + (4N\beta + \xi)\sqrt{\frac{\ln(1/\delta)}{2N}}.$$

In their original paper on learning function-valued output data, Kadri et al. [86] showed that under [assumption 7.1](#), [assumption 7.3](#), and provided that K is weakly measurable, the algorithm is β -stable with $\beta = \frac{\sigma^2 \kappa^2}{2\lambda N}$. Moreover [assumption 7.4](#) holds with $\xi = \sigma^2/2$, where $\sigma = \sigma_y(1 + \kappa/\sqrt{\lambda})$. Thus another generalization bound for [Algorithm 3](#) is

$$\begin{aligned} \mathfrak{R}(f_s) &\leq \mathfrak{R}_{\text{emp}}(f_s, s) + \frac{\kappa^2 C^2 \left(1 + \frac{\kappa}{\sqrt{\lambda}}\right)^2}{\lambda N} \\ &\quad + C^2 \left(1 + \frac{\kappa}{\sqrt{\lambda}}\right)^2 \left(\frac{4\kappa^2}{\lambda} + 1\right) \sqrt{\frac{\ln(1/\delta)}{2N}}. \end{aligned} \tag{7.6}$$

Although both bounds have a convergence rate in $O(N^{1/2})$, an importance difference between the bound [Equation 7.5](#) and the bound [Equation 7.6](#) is that in [Equation 7.6](#) C , κ and λ play a role, while [Equation 7.5](#) add also the trace constant T . This means that [Equation 7.5](#) is less general than [Equation 7.5](#) because when \mathcal{Y} is infinite dimensional, κ is always well defined, while the trace T can be possibly infinite. On the other hand [Equation 7.5](#) is a simpler bound, with a better behaviour in λ when $\lambda < 1$. Indeed [Equation 7.5](#) is a bound in $O(\lambda^{-1})$ while [Equation 7.6](#) is in $O((\lambda\sqrt{\lambda})^{-1})$. Thus the choice between [Equation 7.6](#) and [Equation 7.5](#) has to be done according the the kind of OVK used, as well as the regularization parameter λ .

7.2 CONSISTENCY OF LEARNING WITH ORFF

In this section we are interested by measuring how $\mathfrak{R}(\tilde{f}_s)$ is close to the smallest true risk achieved in the function class \mathcal{F} . The quantity of interest is:

$$\mathfrak{R}(\tilde{f}_s) - \min_{f \in \mathcal{F}} \mathfrak{R}(f).$$

In other words, we quantify the difference between the risk of the optimal solution belonging to a given class of functions \mathcal{F} , and the risk given a solution f_s returned by some learning algorithm. Here to derive a consistency result, we study an algorithm slightly different from [Algorithm 3](#). Given a loss function $L : \mathcal{X} \times \mathcal{F} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ and its canonical cost function $c(f(x), y) := L(x, f, y)$ such that c is Lipschitz

in its first argument. We consider learning with an ORFF $\tilde{\Phi}(x) : \mathcal{Y} \rightarrow \bigoplus_{j=1}^D \mathcal{Y}'$ thanks to the algorithm

$$\theta_s = \begin{cases} \arg \min_{\theta \in \bigoplus_{j=1}^D \mathcal{Y}'} & \frac{1}{N} \sum_{i=1}^N c(\tilde{\Phi}(x_i)^* \theta) \\ \text{subject to} & \max_{j \in \mathbb{N}_D^*} \|\theta_j\|_{\mathcal{Y}} \leq \frac{B}{D}. \end{cases} \quad (7.7)$$

Then the associated output function return is $\tilde{f}_s = \tilde{\Phi}(\cdot)^* \theta_s$. We suppose that the operator $A(\omega)$ used in the construction of $\tilde{\Phi}(x)$ has bounded trace $\Pr_{\rho, \widehat{\text{Haar}}}$ -almost everywhere.

Proposition 7.3 *Let $\Phi_x = (x, \cdot)B(\cdot)$ be a Fourier feature such that there exists a constant $T \in \mathbb{R}_+$ such that*

$$\operatorname{ess \sup}_{\omega \in \widehat{\mathcal{X}}} \operatorname{Tr}[A(\omega)] < T$$

and a constant $u \in \mathbb{R}_+$ such that

$$\operatorname{ess \sup}_{\omega \in \widehat{\mathcal{X}}} \sqrt{\|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}}^2} < u.$$

where $A(\omega) = B(\omega)B(\omega)^*$. Let ρ be the density of a probability distribution with respect to the Haar measure **Haar** and define the set

$$\mathcal{F} = \left\{ f \mid f : x \mapsto \int_{\widehat{\mathcal{X}}} \Phi_x(\omega) \theta(\omega)^* d\widehat{\text{Haar}}(\omega), \|\theta(\omega)\|_{\mathcal{Y}} < B\rho(\omega) \right\} \subseteq \mathcal{H}_K.$$

Eventually let $c : \mathcal{Y}^2 \rightarrow [0, C]$ be a cost function L-Lipschitz in its first argument. Then for any $\delta \in (0, 1)$, given a training sequence $s = (x_i, y_i) \in (\mathcal{X} \times \mathcal{Y})^N$ drawn i.i.d., if \tilde{f}_s is given by [Equation 7.7](#) then we have

$$\begin{aligned} \mathfrak{R}(\tilde{f}_s) - \min_{f \in \mathcal{F}} \mathfrak{R}(f) &\leq 4 \underbrace{\sqrt{\frac{2}{N} \left(LB T^{1/2} + C \sqrt{\ln(2/\delta)} \right)}}_{\text{Estimation error.}} \\ &\quad + \underbrace{\frac{uLB}{\sqrt{D}} \left(1 + \sqrt{2 \ln(1/\delta)} \right)}_{\text{Approximation error.}}. \end{aligned}$$

with probability $1 - 2\delta$ over the training sequence and the random vectors $(\omega_j)_{j=1}^D$.

Proof We follow the proof idea of Rahimi and Recht [\[140\]](#) in the scalar case and adapt it to the vector-valued case in the light of the results of Maurer [\[112\]](#). We first define the two following sets.

$$\mathcal{F} = \left\{ f \mid f : x \mapsto \int_{\widehat{\mathcal{X}}} \Phi_x(\omega)^* \theta(\omega) d\widehat{\text{Haar}}(\omega), \|\theta(\omega)\|_{\mathcal{Y}} < B\rho(\omega) \right\}$$

and

$$\widetilde{\mathcal{F}} = \left\{ f \mid f : x \mapsto \sum_{j=1}^D \Phi_x(\omega_j)^* \theta_j, \forall j \in \mathbb{N}_D^*, \|\theta_j\|_{\mathcal{Y}} < \frac{B}{D} \right\}.$$

Proposition 7.4 (Existence of an approximate function). Let μ be a measure on \mathcal{X} , and f_* a function in \mathcal{F} . Moreover let $\text{ess sup}_{\omega \in \widehat{\mathcal{X}}} \|B(\omega)\|_{y,y}^2 \leq u$. If $(\omega_j)_{j=1}^D$ are drawn i.i.d. from a probability distribution of density ρ w.r.t. $\widehat{\text{Haar}}$, then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over $(\omega_j)_{j=1}^D$, there exists a function \tilde{f} in $\widetilde{\mathcal{F}}$ such that

$$\sqrt{\int_{\mathcal{X}} \left\| \tilde{f}(x) - f_*(x) \right\|_y^2 d\mu(x)} \leq \frac{uB}{\sqrt{D}} \left(1 + \sqrt{2 \ln(1/\delta)} \right).$$

Proof Since $f_* \in \mathcal{F}$, we can write $f_*(x) = \int_{\mathcal{X}} \overline{(x, \omega)} B(\omega) \theta(\omega) d\widehat{\text{Haar}}(\omega)$. Construct the functions $f_j = (\cdot, \omega_j) B(\omega_j) \beta_j$ with $\beta_j := \frac{\theta(\omega_j)}{\rho(\omega)}$, so that $E_{\rho, \widehat{\text{Haar}}} f_j = f_*$ pointwise. Let

$$\tilde{f}(x) = \sum_{j=1}^D \Phi_x(\omega_j)^* \frac{\beta_j}{D}$$

be the sample average of these functions. Then, $\tilde{f} \in \widetilde{\mathcal{F}}$ because $\|\beta_j\|_y / D < B/D$. Also, under the inner product $\int_{\mathcal{X}} \langle f(x), g(x) \rangle_y d\mu(x)$, we have that

$$\left\| \overline{(\cdot, \omega_j)} B(\omega_j) \beta_j \right\|_{L^2(\mathcal{X}, \mu; y)} \leq \text{ess sup}_{\omega \in \widehat{\mathcal{X}}} \|B(\omega)\|_{y,y} B.$$

We introduce the following technical lemma of Rahimi and Recht [140] for concentration of random variable in Hilbert spaces (similar to Pinelis [136]).

Lemma 7.1 Let X_1, \dots, X_D be i.i.d. random variables with values in a ball of radius M centered at the origin in a Hilbert space \mathcal{H} . Denote the sample average $\bar{X} = \frac{1}{D} \sum_{j=1}^D X_j$. Then for any $\delta \in (0, 1)$ with probability $1 - \delta$,

$$\|E\bar{X} - \bar{X}\|_y \leq \frac{M}{\sqrt{D}} \left(1 + \sqrt{2 \ln(1/\delta)} \right).$$

Eventually apply Lemma 7.1 to f_1, \dots, f_D under the canonical inner product of the vector valued function space $L^2(\mathcal{X}, \mu; \mathcal{Y})$ to conclude the proof. \square

Proposition 7.5 (Bound on the approximation error). Let $L(x, f, y)$ be a loss function and $c_y(f(x)) = L(x, f, y)$ be a L -Lipschitz cost function for all $y \in \mathcal{Y}$. Let f_* be a function in \mathcal{F} . Suppose there exists a constant $u \in \mathbb{R}_+$ such that

$$\text{ess sup}_{\omega \in \widehat{\mathcal{X}}} \sqrt{\|A(\omega)\|_{y,y}} \leq u.$$

If $(\omega_j)_{j=1}^D$ are i.i.d. random variables drawn from a probability distribution of density ρ , then for any $\delta \in (0, 1)$ there exists, with probability $1 - \delta$ over $(\omega_j)_{j=1}^D$, a function $\tilde{f} \in \widetilde{\mathcal{F}}$ such that

$$\mathfrak{R}(\tilde{f}) \leq \mathfrak{R}(f_*) + \frac{uLB}{\sqrt{D}} \left(1 + \sqrt{2 \ln(1/\delta)} \right).$$

Proof Given any functions f and g in \mathcal{F} , the Lipschitz hypothesis on c_y , followed by the concavity of the square root (Jensen's inequality) gives

$$\begin{aligned}\mathfrak{R}(f) - \mathfrak{R}(g) &= E_\mu c_y(f(x)) - c_y(g(x)) \\ &\leq E_\mu |c_y(f(x)) - c_y(g(x))| \\ &\leq L E_\mu \|f(x) - g(x)\|_y \\ &\leq L \sqrt{E_\mu \|f(x) - g(x)\|_y^2}.\end{aligned}$$

apply [Proposition 7.4](#) to conclude. \square

Proposition 7.6 (Bound on the estimation error). Let $c_y : \mathcal{Y} \rightarrow [0; C]$ be a L -Lipschitz cost function for all $y \in \mathcal{Y}$. Let $(\omega_j)_{j=1}^D$ be D fixed vectors in $\widehat{\mathcal{X}}$. If $s = (x_i, y_i)_{i=1}^N \in (\mathcal{X} \times \mathcal{Y})^N$ are i. i. d. random variables then for all $\delta \in (0, 1)$, then it holds with probability $1 - \delta$ for all $\tilde{f} \in \widetilde{\mathcal{F}}$ that

$$\mathfrak{R}(\tilde{f}) \leq \mathfrak{R}_{emp}(\tilde{f}, s) + 2\sqrt{\frac{2}{N}} (LB T^{1/2} + C \sqrt{\ln(2/\delta)}).$$

where $\text{Tr}[\tilde{K}_e(e)] < T \in \mathbb{R}_+$.

Proof Since $\tilde{f} \in \widetilde{\mathcal{F}}$ and for all $j \in \mathbb{N}_D^*$, $\|\theta_j\|_y < B/D$, Thus $\|\theta\|_{\bigoplus_{j=1}^D y} < B/\sqrt{D}$. Moreover, if we define $\tilde{\Phi}(x) = \bigoplus_{j=1}^D \Phi_x(\omega_j)$, it gives birth to a RKHS with kernel $D\Phi_x^*\Phi_z$ for all $x, z \in \mathcal{X}$. Thus with arguments similar to [Equation 7.3](#), noticing that the terms in \sqrt{D} cancels out, we obtain a bound on the Rademacher complexity

$$\mathfrak{R}_N(\widetilde{\mathcal{F}}) \leq \sqrt{2BL} \sqrt{N \text{Tr}[\tilde{K}_e(e)]}.$$

Eventually apply [Theorem 7.1](#). \square

We are now ready to prove the main claim. Let f_* be a minimizer of \mathfrak{R} over \mathcal{F} , \tilde{f} a minimizer of \mathfrak{R}_{emp} over $\widetilde{\mathcal{F}}$ and \tilde{f}_* a minimizer of \mathfrak{R} over $\widetilde{\mathcal{F}}$. Then

$$\mathfrak{R}(\tilde{f}) - \mathfrak{R}(f_*) = \mathfrak{R}(\tilde{f}) - \mathfrak{R}(\tilde{f}_*) + \mathfrak{R}(\tilde{f}_*) - \mathfrak{R}(f_*). \quad (7.8)$$

The first difference in the right hand side of the equation is the estimation error. By [Proposition 7.6](#), with probability $1 - \delta$, $\mathfrak{R}(\tilde{f}_*) - \mathfrak{R}_{emp}(\tilde{f}_*, s) \leq \epsilon_{est}$ and simultaneously, $\mathfrak{R}(\tilde{f}) - \mathfrak{R}_{emp}(\tilde{f}, s) \leq \epsilon_{est}$. By optimality of \tilde{f} , $\mathfrak{R}_{emp}(\tilde{f}, s) \leq \mathfrak{R}(\tilde{f}_*)$. Combining these facts, with probability $1 - \delta$,

$$\mathfrak{R}(\tilde{f}) - \mathfrak{R}(\tilde{f}_*) \leq 4\sqrt{\frac{2}{N}} (LB T^{1/2} + C \sqrt{\ln(2/\delta)}) = 2\epsilon_{est}.$$

Applying [Proposition 7.5](#) yields

$$\mathfrak{R}(\tilde{f}_*) - \mathfrak{R}(f_*) \leq \frac{uLB}{\sqrt{D}} (1 + \sqrt{2 \ln(1/\delta)}) = \epsilon_{app}.$$

Conclude by the union bound with probability $1 - 2\delta$ [Equation 7.8](#) is bounded by above by $2\epsilon_{est} + \epsilon_{app}$. Notice that $\text{Tr}[\tilde{K}_e(e)] = \frac{1}{D} \sum_{j=1}^D A(\omega_j)$. Thus if we have $\text{ess sup}_{x \in \widehat{\mathcal{X}}} \text{Tr}[A(\omega)] < \infty$, $\text{Tr}[\tilde{K}_e(e)] \leq \text{ess sup}_{x \in \widehat{\mathcal{X}}} \text{Tr}[A(\omega)]$. \square

7.3 DISCUSSION

In this chapter we reviewed two ways of obtaining generalization bounds (see [Section 7.1](#) and [Subsection 7.1.2](#)) for OVKs by bounding the function class complexity (Maurer [112]) or using algorithm stability arguments (Kadri et al. [86]). Then we used the results of Maurer [112] to prove the consistency of the algorithm obtained by minimizing [Equation 7.7](#), which is a variant of [Algorithm 3](#), where we replace the Tychonov regularizer by a projection in a $\|\cdot\|_\infty$ ball. This bound generalizes the work of Rahimi and Recht [140] to vector-valued learning.

Notice that we cannot directly derive a consistency bound from [Proposition 7.3](#) to [Algorithm 3](#). Indeed with arguments similar to [Remark 6.1](#), we can show that $\tilde{f}_s = \tilde{\Phi}(x)^*\theta$ has a parameter vector θ such that $\|\theta_j\|_y < \sqrt{\frac{2}{\lambda D}}\sigma_y$, where $\sigma_y^2 = \frac{1}{N} \sum_{i=1}^N \|y_i\|_y^2$. Thus if \tilde{f}_s is a solution of [Algorithm 3](#), we do not have $\tilde{f}_s \in \tilde{\mathcal{F}}$, i.e. the Tychonov regularization is not “powerful” enough to guarantee that \tilde{f}_s belongs to $\tilde{\mathcal{F}}$. One could argue that we could choose $\lambda = O(\sqrt{D})$ to obtain consistency with Tychonov regularization, however this makes little sense since in this case if $D \rightarrow \infty$ then $\lambda \rightarrow \infty$ the [Algorithm 3](#) will always return $\tilde{f}_s = 0$.

While the bound in [Proposition 7.3](#) shows the consistency of learning with ORFF it still has low and possibly suboptimal rate. Moreover it does not allow to derive a number of features D smaller than the number of data since both of them decrease the error in $O(D^{-1/2})$ (respectively $O(N^{-1/2})$) as in the reference bound for scalar-valued random features by Rahimi and Recht [140]. In the scalar-valued kernel literature, recent work of Bach [13] with much more involved analysis, gives similar results to Rahimi and Recht [140] in the case of Tichonov regularization. Moreover it suggests that the number of features D to guarantee an error below some constant is linked to the decrease rate of the eigenvalues of the Mercer decomposition of scalar-valued kernel k . If the eigenvalues decrease in $O(m^{-2s})$ then the error is in $O(\log(D)^s D^{-s})$. Lastly the new results of Rudi, Camoriano, and Rosasco [144] show that for scalar-valued kernels, the kernel ridge regression algorithm (which is [Algorithm 3](#) with $A = 1$) generalizes optimality with a number of features $D = O(\sqrt{N})$. Thus the time complexity required for optimal generalization with RFFs in the case of kernel ridge regression is $O(ND^2) = O(N^2)$ and the space complexity is in $O(N^{1.5})$, if the random features are all stored and not computed, on the fly, in an online fashion²¹.

²¹ See [Subsection 6.2.2](#).



8

APPLICATION TO TIME SERIES MODELLING

This chapter shows how to use the ORFF methodology to non-linear vector autoregression. It is an instantiation of the ORFF framework to $\mathcal{X} = \mathcal{Y} = (\mathbb{R}^d, +)$. We also give a generalization of a stochastic gradient descent [51] to ORFF. This is a joint work with Néhémy Lim and Florence d’Alché-Buc and has been published at a workshop of ECML. It is based on the previous work Lim et al. [101] for time series vector autoregression with operator-valued kernels [30].

Contents

8.1	Introduction	138
8.2	Operator-Valued Kernels for Vector Autoregression	138
8.3	Operator-Valued Random Fourier Features	141
8.3.1	Numerical Performance	143
8.3.2	Simulated data	143
8.3.3	Influence of the number of random features	144
8.3.4	Real datasets	145
8.4	Discussion	146

8.1 INTRODUCTION

Time series are ubiquitous in various fields such as climate, biomedical signal processing, videos understanding to name but a few. When linear models are not appropriate, a generic nonparametric approach to modelling is relevant. In this work we build on a recent work about Vector Autoregressive models using Operator-Valued Kernels [100, 101]. Vector autoregression is addressed in a Vector Valued Reproducing Kernel Hilbert Space with the important property to allow for couplings between outputs. Given a d -dimensional time series of N data points $\{x_1, \dots, x_N\}$, autoregressive models based on operator-valued kernels have the form $\hat{x}_{t+1} = h(x_t) = \sum_{\ell=1}^{N-1} K(x_t, x_\ell) c_\ell$ where coefficients $c_\ell \in \mathbb{R}^d$, $\ell = 1, \dots, N - 1$ are the model parameters. A naive approach for training such a model requires a memory complexity $O(N^2 d^2)$, which makes the method prohibitive for large-scale problems.

To scale up standard algorithms, we define an approximated operator-valued feature map $\tilde{\Phi} : \mathbb{R}^d \rightarrow \mathbb{R}^D$ that allows to approximate the aforementioned model h in the RKHS by the following function

$$\tilde{h}(x_t) = \tilde{\Phi}(x_t)^* \theta \approx h(x_t).$$

The features maps are matrices of size $D \times d$ where D controls the quality of the approximation, d is the dimension of the inputs and θ is here the parameter vector to learn. This formulation allows to reduce the memory complexity to $O((N - 1)D + (N - 1)d)$ which is now linear w.r.t. the number of data points (see [Section 6.1](#)). The principle used for building the feature map extends the idea of scalar Random Fourier Features to the operator-valued case [139, 167].

8.2 OPERATOR-VALUED KERNELS FOR VECTOR AUTOREGRESSION

Assume that we observe a dynamical system composed of $d \in \mathbb{N}^*$ state variables at $N \in \mathbb{N}^*$ evenly-spaced time points. The resulting discrete multivariate time series is denoted by $x_{1:N} = (x_\ell)_{\ell=1}^N$ where $x_\ell \in \mathbb{R}^d$ denotes the state of the system at time t_ℓ , $\ell \in \mathbb{N}_N^*$. It is generally assumed that the evolution of the state of the system is governed by a function h , such that $x_t = h(x_{t-p}, \dots, x_{t-1}) + u_t$ where t is a discrete measure of time and u_t is a zero-mean noise random vector. Then h is usually referred to as a vector autoregressive model of order p . In the remainder of the chapter, we consider first-order vector autoregressive models, that is $p = 1$. In a supervised learning set-

ting, the vector autoregression problem consists in learning a model $\hat{h} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ from a given training set

$$\mathbf{s} = ((x_1, x_2), \dots, (x_{N-1}, x_N)) \in (\mathbb{R}^d \times \mathbb{R}^d)^N.$$

In the literature, a standard approach to vector autoregressive modelling is to fit a VAR model. The VAR(1) model reads $h(x_t) = Ax_t$ where A is an $d \times d$ matrix whose structure encodes the temporal relationships among the d state variables.

However, due to their intrinsically linear nature, VAR models fail to capture the nonlinearities underlying realistic dynamical systems. In this chapter we builds upon the work of Lim et al. [101] where the authors introduced a family of nonparametric nonlinear autoregressive models called OKVAR. OKVAR models rely on the theory of operator-valued kernels [133, 152], which provides a versatile framework for learning vector-valued functions [5, 41, 113]. Those models can be regarded as natural extensions of VAR models to the nonlinear case.

Next, we recall key elements of the theory of VV-RKHS of functions from \mathbb{R}^d to \mathbb{R}^d (see [Section 4.2](#) for the detailed construction). We first introduced the matrix-valued kernel which is an instance of OVKs.

Definition 8.1 (Matrix-valued kernels). A function $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ is said to be a positive $\mathbb{R}^{d \times d}$ -valued kernel if:

1. $\forall x, z \in \mathbb{R}^d, K(x, z) = K(z, x)^*$,
2. $\forall N \in \mathbb{N}, \forall ((x_i, y_i))_{i=1}^N \in (\mathbb{R}^d \times \mathbb{R}^d)^N, \sum_{i,j=1}^N y_i^* K(x_i, x_j) y_j \geq 0$.

Furthermore, for a given $\mathbb{R}^{d \times d}$ -valued kernel K , we associate K with a unique VV-RKHS $(\mathcal{H}_K, \langle \cdot, \cdot \rangle_{\mathcal{H}_K})$ of functions from \mathbb{R}^d to \mathbb{R}^d . The precise construction of \mathcal{H}_K can be found in [Section 3.3](#). In this section, we assume that all functions $h \in \mathcal{H}_K$ are continuous. Then K is called an \mathbb{R}^d -Mercer kernel (see [definition 3.8](#)).

Similarly to the case of scalar-valued kernels, working within the framework of VV-RKHS allows to take advantage of representer theorems ([Theorem 6.2](#)) for a class of regularized loss functions such as ridge regression. More precisely, we consider h , a nonparametric vector autoregressive model of the following form assuming we have observed N data points. Given x_t the state vector at time t , we have $\hat{x}_{t+1} = \sum_{\ell=1}^{N-1} K(x_t, x_\ell) c_\ell$ where $x_{1:N} = (x_i)_{i=1}^N \in (\mathbb{R}^d)^N$ is the observed time series, $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ is a matrix-valued kernel and $(c_i)_{i=1}^{N-1} \in (\mathbb{R}^d)^{N-1}$ are the model parameters. We call OKVAR any model of the above form. In Lim et al. [101], the authors developed a

²² See for instance Parikh and Boyd [130] about proximal algorithms and Fercq and Peter [62], Fercq and Richtárik [63], and Richtárik and Takáč [141] for proximal block coordinate descent.

family of OKVAR models based on appropriate choices of kernels to address the problem of network inference where both the parameters $c_\ell, \ell \in \mathbb{N}_{N-1}^*$ and the OVK itself are learned using a proximal block coordinate descent algorithm ²² under sparsity constraints. In the following, we will not consider the kernel learning problem and will use a simple ridge loss. We will also illustrate our approach to a well known class of OVK, called *decomposable* or *separable* matrix-valued kernels [39, 113], and instance of Decomposable OVK that were originally introduced to solve multi-task learning problems [59]. Other kernels may also be considered as developed in Subsection 3.3.3.

Proposition 8.1 (Decomposable matrix-valued kernels). *Let the function $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be a scalar-valued kernel and $\Gamma \in \mathbb{R}^{d \times d}$ a positive semidefinite matrix of size $d \times d$. Then function $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ defined for all $(x, z) \in \mathbb{R}^d \times \mathbb{R}^d$ as $K(x, z) = k(x, z)\Gamma$ is a decomposable matrix-valued kernel.*

A common choice for the scalar-valued kernel is the Gaussian kernel

$$k_{\text{Gauss}}(x, z) = \exp\left(-\frac{1}{2\sigma^2}\|x - z\|_2^2\right)$$

for any $x, z \in \mathbb{R}^d$ and $\sigma \in \mathbb{R}_+$. Notice that k_{Gauss} can equivalently be written with an hyperparameter $\gamma \in \mathbb{R}_+$:

$$k_{\text{Gauss}}(x, z) = \exp(-\gamma\|x - z\|_2^2),$$

with $\sigma = (2\gamma)^{-1/2}$. The corresponding decomposable kernel is referred to as K_{dec} and is as $K_{\text{dec}}(x, z) = k_{\text{Gauss}}(x, z)\Gamma$ with Γ a positive semidefinite matrix.

While the model parameters c_ℓ 's are estimated under sparsity constraints in Lim et al. [101], here we consider the classic kernel ridge regression setting where the loss function to minimize is

$$\mathfrak{R}_\lambda(h, s) = \frac{1}{N-1} \sum_{\ell=2}^N \|h(x_{\ell-1}) - x_\ell\|_2^2 + \lambda \|h\|_{\mathcal{H}_K}^2 \quad (8.1)$$

with $\lambda \geq 0$ and $\|h\|_{\mathcal{H}_K}^2 = \sum_{t, \ell=1}^{N-1} c_t^* K(x_t, x_\ell) c_\ell$. The optimization problem is solved using a L-BFGS-B [37] which is well suited for optimization problems with a large number of parameters, and is widely used as a training algorithm on small/medium-scale problems. However, like standard kernel methods, OKVAR suffers from unfavourable computational complexity both in time and memory since it needs to store the full Gram matrix, preventing its ability to scale to large data sets and making it really slow on medium scale problem. We argue that this obstacle can be effectively overcome: in the following we develop a method to scale up OKVAR to successfully tackle medium/large scale autoregression problems.

8.3 OPERATOR-VALUED RANDOM FOURIER FEATURES

We now introduce our methodology to approximate OVKs. Given a shift-invariant kernel $K(x, z) = K_0(x - z)$, we approximate K by finding an explicit feature map such that $\tilde{\Phi}(x)^* \tilde{\Phi}(z) \approx K_0(x - z)$. The idea is to use a generalization of Bochner's theorem for the OVK family that states that any translation-invariant OVK can be written as the Fourier transform of a positive operator-valued measure. More precisely, we build on the following proposition first proved in [41]. More details can be found in [Section 4.2](#).

In the following, suppose that $K_0 = k_0(\cdot)A$ is a decomposable kernel. Decomposable kernels belong to the family of translation-invariant OVKs. From [Proposition 4.3](#) we see that $C(\omega)_{ij} = \mathcal{F}^{-1}[k_0(\cdot)](\omega)A_{ij}$. We decompose A as $A = BB^*$, note that A does not depend on ω , and we denote $\bigoplus_{j=1}^D z_j$ the Dm -long column vector obtained by stacking vectors $z_j \in \mathbb{R}^m$. Then we define an approximate feature map for K_0 , called Operator-valued Random Fourier Feature (ORFF) map [30] as follows (see [Subsection 4.2.2](#) and [Section 4.3](#)). For all $x \in \mathbb{R}^d$,

$$\tilde{\Phi}^{dec}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle B^* \\ \sin \langle x, \omega_j \rangle B^* \end{pmatrix}, \quad \omega_j \sim \mathcal{F}^{-1}[k_0],$$

which can also be expressed as a Kronecker product \otimes of a scalar feature map with a matrix (see [Subsection 6.3.1](#)): $\tilde{\Phi}^{dec}(x) = \tilde{\varphi}(x) \otimes B^*$ where

$$\tilde{\varphi}(x) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos \langle x, \omega_j \rangle \\ \sin \langle x, \omega_j \rangle \end{pmatrix}, \quad \omega_j \sim \mathcal{F}^{-1}[k_0]$$

is a scalar-valued feature map. In particular, if k_0 is a Gaussian kernel with bandwidth σ^2 , then $\mathcal{F}^{-1}[k_0] = \mathcal{N}(0, 1/\sigma^2)$ as proven in Rahimi and Recht [139]. More examples on different OVK can be found in [Subsection 4.2.2](#) as well as a proof of the uniform convergence of the kernel approximation in [Section 5.1](#) defined by $\tilde{K}(x, z) = \tilde{\Phi}(x)^* \tilde{\Phi}(z)$ towards the true kernel. In the case of vector autoregression, we consider a model \tilde{h} of the form: $\hat{x}_{t+1} = \tilde{\Phi}(x_t)^* \theta$. That model is referred to as ORFFVAR in the remainder of the section. Now, given the operator-valued feature map, we get a linear model, and we want to minimize the regularized risk

$$\mathfrak{R}_\lambda(\theta, s) = \frac{1}{N-1} \sum_{\ell=2}^N \|(\tilde{\varphi}(x_{\ell-1})^* \otimes B)\theta - x_\ell\|_2^2 + \lambda \|\theta\|_2^2$$

with $\lambda > 0$ instead of [Equation 8.1](#) (see [Theorem 6.3](#)). In their paper Brault, Lim, and Buc [30] proposed to formulate the learning

problem as a Stein equation when dealing with decomposable kernels, and then used an appropriate solver [163]. We opted here for a more general algorithm, which is a variant of the doubly stochastic gradient descent [51]. In a few words, this algorithm is a stochastic gradient descent that takes advantage of the feature representation of the kernel allowing the number of features to grow along with the number of points. Dai et al. [51] show that the number of iterations needed for achieving a desired accuracy ϵ using a stochastic approximation is $\Omega(1/\epsilon)$, making it competitive compared to other stochastic methods for kernels such as NORMA [90] and its OVK adaptation ONORMA [10]. We propose here in [Algorithm 4](#), an extension of the doubly stochastic gradient descent of Dai et al. [51] to OVKs. Additionally we consider a batch approach w. r. t. the data and the features, and make it possible to “cap” the maximum number of features. The inputs of the algorithm are: \mathcal{X} the input data, \mathcal{Y} the targets, K_e the OVK used for learning, γ_t the learning rate (see Dai et al. [51] for a discussion on the selection of a proper learning rate), T the number of iterations, n the size of data batch, b the size of the feature batch, and D the maximum number of features. Note that if K_0 is a scalar kernel, $D = T$, $b = 1$ and $n = 1$, we recover the algorithm formulated in Dai et al. [51].

Algorithm 4: Block-coordinate mini-batch doubly SGD.

Data: $\mathcal{X}, \mathcal{Y}, K_e, \gamma_t, \lambda, T, n, D, b$
Result: Find θ

- 1 Let $D_b = D/b$ and find (ω, x) , $B(\omega)$ and $\mu(\omega)$ from K_e ;
- 2 **for** $i = 1$ **to** D_b **do**
- 3 $\theta_{i,.}^1 = 0$;
- 4 **end**
- 5 **for** $t = 1$ **to** T **do**
- 6 $\mathcal{A}_t = \mathcal{X}_t \times \mathcal{Y}_t$, a random subsample of n data from $\mathcal{X} \times \mathcal{Y}$;
- 7 $h(\mathcal{X}_t) = \text{predict}(\mathcal{X}_t, \theta^t, K_e)$; // Make a prediction.
- 8 $\Omega_i \sim \mu$ with seed i , where $i = ((t - 1) \bmod D_b) + 1$;
 // Sample b features from μ .
- 9 **for** $\omega \in \Omega_i$ // Update the parameters from the
 gradient.
- 10 **do**
- 11 $\theta_{i,\omega}^{t+1} =$
 $\theta_{i,\omega}^t - \gamma_t \left(\frac{1}{|\mathcal{A}_t|} \sum_{(x,y) \in \mathcal{A}_t} \frac{B(\omega)^*(\omega, x)(h(x) - y)}{\sqrt{D}} + \lambda \theta_{i,\omega}^t \right)$;
- 12 **end**
- 13 **end**
- 14 **return** θ^{t+1}

In addition, the convergence of the algorithm can be speeded-up by preconditioning by the Hessian of the system. An experimental C++ code is available at <https://github.com/RomainBrault/0V2SGD>.

Algorithm 5: $h(\mathcal{X}) = \text{predict}(\mathcal{X}, \theta, K_e)$

Data: \mathcal{X}, θ, K_0

- 1 Find $(\omega, x), B(\omega)$ and $\mu(\omega)$ from K_0 ;
- 2 $f(\mathcal{X}) = 0$;
- 3 **for** $x \in \mathcal{X}$ **do**
- 4 **for** $i = 1$ to D **do**
- 5 $\Omega_i \sim \mu(\omega)$ with seed i ;
- 6 **for** $\omega \in \Omega_i$ **do**
- 7 $h(x) = h(x) + \overline{(\omega, x)} B(\omega) \theta_{i,\omega}$;
- 8 **end**
- 9 **end**
- 10 **end**
- 11 **return** $h(\mathcal{X})$

8.3.1 Numerical Performance

We now apply [Algorithm 4](#) to toy and real datasets.

8.3.2 Simulated data

To assess the performance of our models, we start our investigation by generating discrete d -dimensional time series $(x_t)_{t \geq 1}$ as follows

$$\begin{cases} x_1 \sim \mathcal{N}(0, \Sigma_x) \\ x_{t+1} = h(x_t) + u_{t+1}, \quad \forall t > 0. \end{cases} \quad (8.2)$$

where the residuals are homoscedastic and distributed according to $u_t \sim \mathcal{N}(0, \Sigma_u)$. We study two different kinds of noise: an isotropic noise with covariance $\Sigma_u = \sigma_u^2 I_d$ and an anisotropic noise with Toeplitz structure $\Sigma_{u,ij} = \nu^{|i-j|}$, where ν lives in $(0, 1)$. We generated $N = 1000$ data points and used a Sequential cross-validation (SCV) with time windows $N_t = N/2$ to measure the Mean Squared Error SCV-MSE of the different models. Next, we compare the performances of VAR(1), OKVAR and ORFFVAR through three scenarios. Across the simulations, the topological structures of the underlying dynamical systems are encoded by a matrix A of size 5×5 . All entries of A are set to zero except for the diagonal where all coefficients are equal to 0.9 for Settings 1 and 3 and 0.5 for Setting 2. Then five off-diagonal coefficients are drawn randomly from $\mathcal{N}(0, 0.3)$ for Settings

1 and 3 and $\mathcal{N}(0, 0.5)$ for Setting 2. We check that all the eigenvalues of A are less than one to ensure the stability of the system. More specifically, we picked the following values of parameters for each scenario.

- **Setting 1: Linear model.**: $h(x_t) = Ax_t$, $\nu = 0.9$ and $\sigma_u = 0.9$,
- **Setting 2: Exponential model.**: $h(x_t) = A \exp(x_t)$ where \exp is the element-wise exponential function, $\nu = 0.09$ and $\sigma_u = 0.09$,
- **Setting 3: Sine model.**: $h(x_t) = A \sin(x_t)$ where \sin is the element-wise sine function, $\nu = 0.9$ and $\sigma_u = 0.009$.

ORFFVAR is instantiated with $D = 25$ random features in presence of a white noise while we set $D = 50$ in case of a Toeplitz noise. We summarize the computational efficiency and the statistical accuracy of the models in Table 8.1. Throughout all the experiments, we set B as the identity matrix of size $d \times d$. This reflects the absence of a prior on the structure of the data. A further study on the influence of the choice of B can be found in Álvarez, Rosasco, and Lawrence [5] and Propositions 3.9 and 4.12.

In Setting 1, we observe that OKVAR does not provide any advantage over VAR(1) as expected since the data were generated according to a linear VAR(1) model. Note that OKVAR takes orders of magnitude more time to achieve the same performance as VAR(1) while ORFFVAR performs equally well with a competitive timing. In nonlinear scenarios (Settings 2 and 3), OKVAR and ORFFVAR consistently outperform VAR(1). Noticeably, ORFFVAR reaches the accuracy of OKVAR with the computation time of VAR(1).

model	Setting	noise	SVC-MSE	1 variance	1 time	SVC-MSE	2 variance	2 time	SVC-MSE	3 variance	3 time
VAR(1)	White	0.914979	0.572485	0.002467(s)	0.001275	0.000994	0.002346(s)	0.009534	0.006003	0.001697(s)	
	Toeplitz	1.091096	1.267880	0.004822(s)	0.017014	0.013498	0.002050(s)	0.116901	0.127396	0.001702(s)	
ORFFVAR	White	0.919663	0.572936	0.000994(s)	0.001003	0.000647	0.001284(s)	0.009536	0.005998	0.002377(s)	
	Toeplitz	1.097183	1.268978	0.001022(s)	0.012635	0.008837	0.012144(s)	0.116964	0.127395	0.000934(s)	
OKVAR	White	0.958790	0.591934	0.104706(s)	0.001100	0.000731	0.027099(s)	0.009227	0.005717	0.014458(s)	
	Toeplitz	1.410969	1.312243	0.289046(s)	0.013854	0.010977	1.856988(s)	0.160133	0.136570	0.019170(s)	

Table 8.1: Sequential SCV-MSE and computation times for VAR(1), ORFFVAR and OKVAR on synthetic data (Settings 1, 2 and 3).

8.3.3 Influence of the number of random features

Here, we investigate the impact of D , the number of random features for ORFFVAR. To this end, we generated $N = 10000$ data points following Equation 8.2, with exponential nonlinearities and white noise as in Setting 2. We performed a sequential cross-validation on a window of $N/2$ data. As expected the error decreases with the number

of random features D ([Table 8.2](#)). For the same computation time ($D = 25$) as VAR(1), ORFFVAR achieves an SCV-MSE that is twice as small.

model	D = 1	D = 5	D = 10	D = 25	D = 50	D = 100	VAR(1)
SVC-MSE	0.005342	0.001111	0.000991	0.000962	0.000949	0.000944	0.001660
variance	0.008639	0.000793	0.000660	0.000618	0.000608	0.000605	0.001363
time	0.001191(s)	0.002384(s)	0.003614(s)	0.018469(s)	0.038229(s)	0.069294(s)	0.019634(s)

Table 8.2: SVC-MSE with respect to D the number of random features for ORFFVAR.

8.3.4 Real datasets

We now investigate three real datasets. The performances of the models on those datasets are recorded in [Table 8.3](#). Throughout the experiments, the hyperparameters are set as follows: the bandwidth of the Gaussian kernel σ is chosen as the median of the Euclidean pairwise distances and the regularization parameter λ was tuned on a grid. The number of random features D and the parameters in [Algorithm 4](#) were picked so as to reach the level of accuracy of OKVAR/VAR.

MACRODATA This dataset is part of the Python library [Statmodels](#)¹. It contains 204 US macroeconomic data points collected on the period 1959–2009. Each data point represents 12 economic features. No pre-processing is applied before learning. We measure SCV-MSE using a window of 25 years (50 points). We instantiated [Algorithm 4](#) as follows: $\gamma_t = 1$, $\lambda = 10^{-3}$, $D = 100$, $T = 2$ and $b = 50$ for ORFF and $\lambda = 0.00025$ and $\sigma = 11.18$ for OKVAR.

GESTURE PHASE. This dataset² is constructed using features extracted from seven videos with people gesticulating. We present the results for videos 1 and 4, consisting in 1069 data points and 31 features. Data are normalized prior to learning. We measure SCV-MSE using a time window of 200 points. We implemented ORFFVAR with $\gamma_t = 1$, $\lambda = 10^{-3}$, $D = 100$, $T = 2$ and $b = 50$.

CLIMATE. This dataset [106] contains monthly meteorological measurements of 18 variables (temperature, CO₂ concentration, ...) collected at 135 different locations throughout the USA and recorded over 13 years, thus resulting in 135 time series of dimension 18 and length 156. Data are standardized at each station. A unique model is learned for all stations. SCV-MSE is measured on a window of 1872 points, corresponding to the data of all the 135 stations over one year.

1 <https://github.com/statsmodels/statsmodels>

2 <https://archive.ics.uci.edu/ml/datasets/Gesture+Phase+Segmentation>

Specifically, we set the parameters of ORFFVAR as follows: $\gamma_t = 1$, $\lambda = 10^{-6}$, $D = 100$, $T = 1$ and $b = 100$.

HEART. The dataset is a multivariate time-serie recorded from a patient in the sleep laboratory of the Beth Israel Hospital in Boston, Massachusetts³. The attributes are the heart rate, the chest volume (respiration force) and the blood oxygen concentration. The time-serie contains 17000 points recorded at 2Hz during roughly 4 hours 30 minutes. We used a window of 240 points for the sequential cross-validation (corresponding to 2 minutes of observations).

Dataset	N	d	ORFFVAR			VAR(1)			OKVAR		
			SCV-MSE	variance	time	SCV-MSE	variance	time	SCV-MSE	variance	time
Macrodata	#203	#12	445.9	84.5	0.014(s)	449.1	1021	0.0005(s)	499.8	793.0	0.641(s)
Gesture phase 1	#1743	#31	0.741	2.999	0.009(s)	0.980	3.370	0.0014(s)	N. A.	N. A.	N. A.
Gesture phase 4	#1069	#31	0.473	2.406	0.061(s)	0.768	6.49	0.0075(s)	N. A.	N. A.	N. A.
Climate	#19375	#18	0.237	0.2128	0.396(s)	0.266	0.218	0.0124(s)	N. A.	N. A.	N. A.
Heart	#16999	#3	0.262	1.020	0.011(s)	0.259	1.040	0.0010(s)	N. A.	N. A.	N. A.

Table 8.3: SCV-MSE and computation times for ORFFVAR, VAR(1) and OKVAR on real datasets.

8.4 DISCUSSION

Operator-Valued Random Fourier Feature provides a way to approximate OVK and in the context of time series, allows for nonlinear Vector Autoregressive models that can be efficiently learned both in terms of computing time and memory. We illustrate the approach with a simple family of Operator-valued kernels, the so-called decomposable kernels but other kernels may be used. While we focused on first-order autoregressive models, we will consider extensions of our models for higher orders. In this work, the kernel hyperparameter B is given prior to learning, however it would be interesting to learn B as in OKVAR. Thus, a promising perspective is to use these models in tasks such as network inference and search for causality graphs among the state variables for large-scale time series [100, 101].



³ <http://www-psych.stanford.edu/~andreas/Time-Series/SantaFe.html>

Part III

CONCLUSION AND WORK IN PROGRESS

9

WORK IN PROGRESS

To conclude our work we present some work in progress. We show practical applications of operator-valued kernels acting on an infinite dimensional space \mathcal{Y} . We give two examples. First we show how to generalize many quantile regression to learn a continuous function of the quantiles on the data. Second we apply the same methodology to the one-class SVM algorithm in order to learn a continuous function of all the level sets. We conclude by presenting Operalib, a python library developed during this thesis which aims at implementing OVK-based algorithms in the spirit of Scikit-learn [132].

Contents

9.1	Learning function-valued functions	150
9.1.1	Quantile regression	150
9.1.2	Functional output data	151
9.1.3	ORFF for functional output data	151
9.1.4	Many quantile regression	155
9.1.5	One-class SVM revisited	157
9.2	Operalib	160

9.1 LEARNING FUNCTION-VALUED FUNCTIONS

In this section we show how to use OVK in hand with the ORFF framework to learn function-valued functions. We focus on two application cases: quantile regression and one-class classification. This section is rather an informal (but detailed) discussion on ideas that we plan to improve for future publications.

9.1.1 Quantile regression

This introduction to quantile regression is adapted from the paper of Sangnier, Fercoq, and Buc [146]. As we have seen in the introductory [Chapter 2](#), a standard task in Machine Learning is to estimate the conditional expectation $f(x) = E_{\Pr}[Y|X = x]$, where $(X, Y) \sim \Pr$ with some function belonging to a hypothesis space $f \in \mathcal{F}$. Yet, many sensitive applications need more than the expected valued of the relationship between random variables. To control the “quality” of the predicted value from an input x , fields such as economics, medicine, physics or social science require to have access to the different quantile to model the distribution around the mean $f(x) \in \mathbb{R}$ and strengthen their analysis.

Here we are interested in learning and predicting simultaneously *all* the quantiles on the compact $[0, 1]$, of the scalar-valued random variable $Y|X$. We place ourselves in the setting of conditional quantile regression by minimization of the pinball loss [91]. For $\tau \in [0, 1]$ the pinball loss reads

$$L_\tau(x, f, y) = \max(\tau(f(x) - y), (\tau - 1)(f(x) - y)).$$

In a nutshell, this loss has been introduced by noticing that finding the optimal location parameter $\mu = f(x)$ in the ℓ_1 loss $L(x, f, y) = |f(x) - y|$ yields an estimator of the unconditional median [91]. Recently Sangnier, Fercoq, and Buc [146] proposed to learn simultaneously many quantiles by minimizing the multi-quantile loss function. Given a vector of quantiles $\boldsymbol{\tau} = (\tau_1, \dots, \tau_p) \in [0, 1]^p$

$$L_{\boldsymbol{\tau}}(x, f, y) = \sum_{i=1}^p \max(\tau_i(f(x)_i - y), (\tau_i - 1)(f(x)_i - y)).$$

We see that now it is necessary for $f(x) \in \mathbb{R}^p$ to be vector-valued. In this work we push further the idea by considering that $f(x)$ is a function of an arbitrary quantile $\tau \in [0, 1]$. Thus we view f as a vector-valued function $f : \mathbb{R} \rightarrow ([0, 1] \rightarrow \mathbb{R})$. For the sake of simplicity we note $f(x) = f_x$ and introduce the generalized pinball loss

$$L(x, f, y) = \int_{[0,1]} \max(\tau(f_x(\tau) - y), (\tau - 1)(f_x(\tau) - y)) d\tau. \quad (9.1)$$

9.1.2 Functional output data

Pioneer work on learning function-valued function has been done by Kadri et al. [86]. Inspired by them we develop an ORFF methodology to learn functional data where the outputs are functions that we suppose living in a RKHS.

Namely, we suppose that the image of a function f , has values $f(x) \in \mathcal{H}_{k_T}$ in a RKHS, where $k_T : T^2 \rightarrow \mathbb{R}$ is a scalar-valued kernel and \mathcal{H}_{k_T} is the corresponding RKHS. From this hypothesis we see that

$$f_x(\tau) = \langle f(x), k_T(\cdot, \tau) \rangle_{\mathcal{H}_{k_T}}$$

If we add the second hypothesis that $f \in \mathcal{H}_K$, where \mathcal{H}_K is a *Vector Valued Reproducing Kernel Hilbert Space* for some Operator-Valued Kernel K , Carmeli et al. [41] showed in example 6 page 17-18 that in this case the operator K is given by

$$K = \begin{cases} \mathcal{X} \times \mathcal{X} & \rightarrow \mathcal{L}(\mathcal{H}_{k_T}) \\ x, z & \mapsto k_X(x, z) I_{\mathcal{H}_{k_T}}, \end{cases} \quad (9.2)$$

where $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is another scalar-valued kernel. Moreover Carmeli et al. [41] showed in example 7 page 18-19 that the VV-RKHS induced by K is the same RKHS than the one induced by the kernel K' defined as follow for some measure μ with support T .

$$K' = \begin{cases} \mathcal{X} \times \mathcal{X} & \rightarrow \mathcal{L}(L^2(T, \mu)) \\ x, z & \mapsto (g \mapsto k_X(x, z) \int_T k_T(\cdot, \tau) g(\tau) d\mu(\tau)). \end{cases}$$

This is exactly the decomposable kernel introduced in [Proposition 3.12](#) in [Chapter 2](#). Because the RKHSs induced by K and K' is the same, we can either view its elements as functions from \mathcal{X} into \mathcal{H}_{k_T} (through \mathcal{H}_K) or as functions from \mathcal{X} into $L^2(T, \mu)$ (through $\mathcal{H}_{K'}$).

9.1.3 ORFF for functional output data

Because $\mathcal{Y} = \mathcal{H}_{k_T}$ is a proper (infinite dimensional) Hilbert space, we can apply the ORFF methodology. Let k_X be a scalar Mercer kernel and $\mathcal{X} = \mathbb{R}$. Then by [Proposition 4.8](#) applied to the decomposable kernel (see [Subsection 4.3.3](#)) we have the following approximate feature map for K defined in [Equation 9.2](#):

$$\tilde{\Phi}(x)y = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos(x\omega_j) B^*y \\ \sin(x\omega_j) B^*y \end{pmatrix}, \quad \omega_j \sim \mathcal{F}[k_X] \text{ i.i.d.}$$

where $BB^* = I_{\mathcal{H}_{k_T}}$ and $y \in \mathcal{H}_{k_T}$. At this point we could choose $B = I_{\mathcal{H}_{k_T}}$. However this is not really useful since it would make the

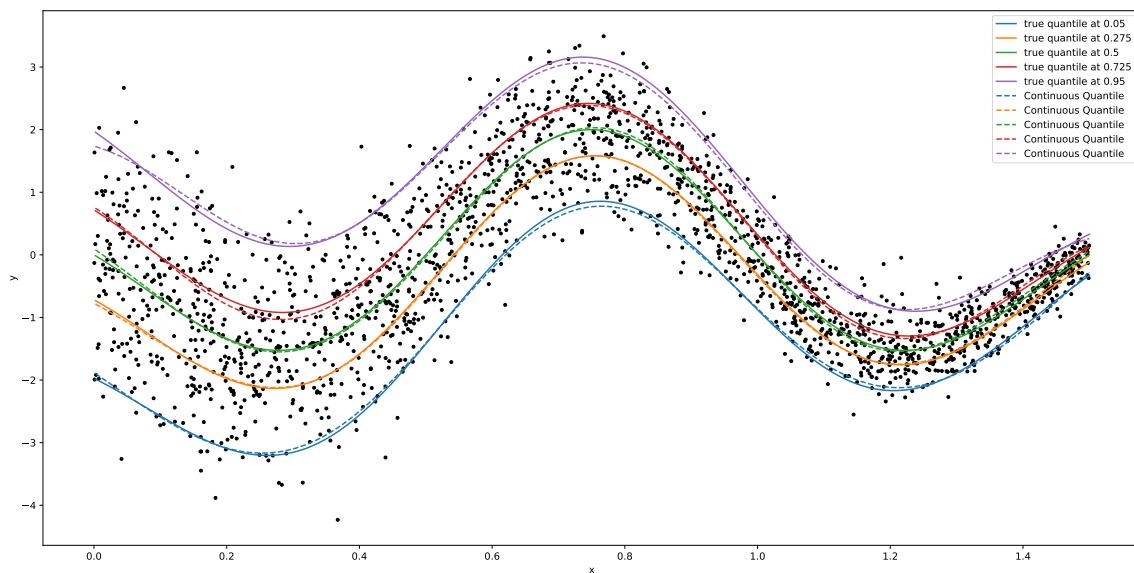


Figure 9.1: Learning a continuous quantile function with ORFF regression.

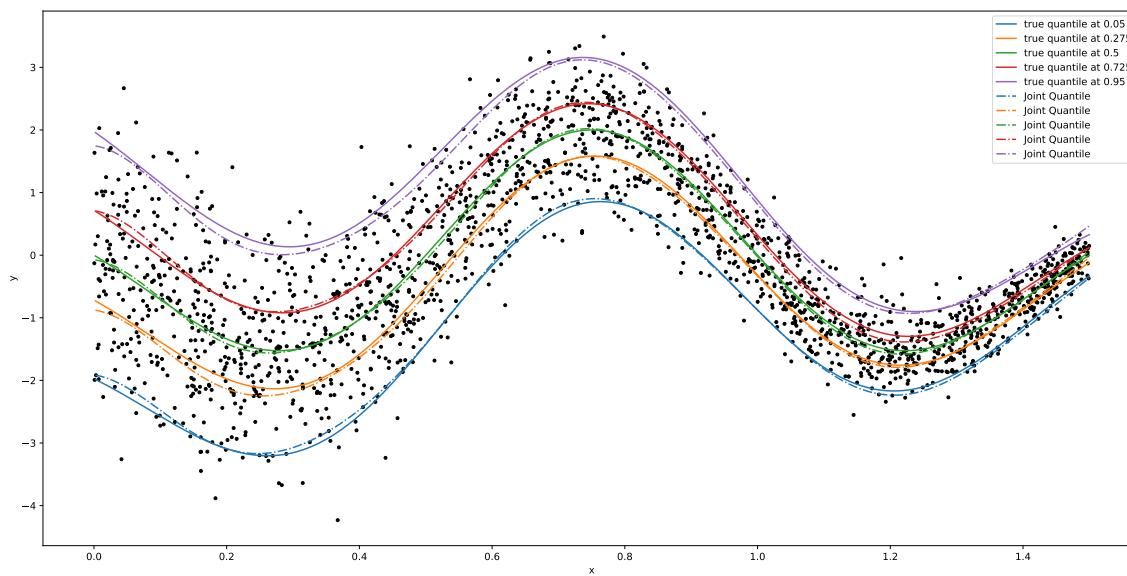


Figure 9.2: Learning many quantile with joint OVK regression.

redescription space $\tilde{\mathcal{H}} = \bigoplus_{j=1}^D \mathcal{H}_{k_j}$, which is a direct sum of infinite dimensional RKHS. Yet since \mathcal{H}_{k_j} is a RKHS, according to [Proposition 3.4](#) it is possible to define a feature operator $W : \mathcal{H} \rightarrow \mathcal{H}_{k_j}$ such that $(Wg)(\tau) = \Phi_\tau^* g$. Moreover $W^* W$ is the identity on $\text{Im } \Phi_\tau$ which is here \mathcal{H}_{k_j} . (see the proof of [Proposition 3.4](#) and Carmeli et al. [41]). Thus we can choose $\Phi_\tau = \Phi(\tau)$ to be the functional Fourier feature map associated to k_j defined in [Proposition 4.6](#). Then we have $BB^* = I_{\mathcal{H}_{k_j}} = WW^*$. Thus we can choose $B = W = \Phi(\cdot)^*$ and the approximate feature map reads

$$\tilde{\Phi}(x) \in \mathcal{L} \left(\mathcal{H}_{k_j}; \bigoplus_{j=1}^D L^2 \left(\mathcal{T}, \widehat{\Pr_{\text{Haar}, \rho}} \right) \right)$$

and

$$(\tilde{\Phi}(x)g)(\tau) = \frac{1}{\sqrt{D}} \bigoplus_{j=1}^D \begin{pmatrix} \cos(x\omega_j)W^*g \\ \sin(x\omega_j)W^*g \end{pmatrix}, \quad \omega_j \sim \mathcal{F}[k_j] \text{ i.i.d.}$$

Then it is easy to verify that the adjoint operator is given by

$$\begin{aligned} \left(\tilde{\Phi}(x)^* \theta \right) (\tau) &= \frac{1}{\sqrt{DD'}} \sum_{j=1}^D \left(\cos(x\omega_j) \right. \\ &\quad \left. + \sin(x\omega_j) \right) \left(\bigoplus_{k=1}^{D'} \begin{pmatrix} \cos(\tau\omega'_k) \\ \sin(\tau\omega'_k) \end{pmatrix} \right)^* \theta_j \\ &= \frac{1}{\sqrt{DD'}} \sum_{j=1}^D \left(\cos(x\omega_j) + \sin(x\omega_j) \right) \theta_{jk} \left(\cos(\tau\omega'_k) \right. \\ &\quad \left. + \sin(\tau\omega'_k) \right), \\ \omega_j &\sim \mathcal{F}[k_j] \text{ i.i.d. and } \omega'_k \sim \mathcal{F}[k_j] \text{ i.i.d..} \end{aligned}$$

where $\theta_k \in \mathbb{R}^{D'}$, for all $k \in \mathbb{N}_D^*$ and $\theta_{jk} \in \mathbb{R}$ for all $j \in \mathbb{N}_D^*$ and all $k \in \mathbb{N}_{D'}^*$. The above equations can be rewritten in matrix form which results in the following conjecture.

Conjecture 9.1 *If $\tilde{\varphi}_X$ is an RFF for k_X such that $\tilde{\varphi}(x) \in \mathbb{R}^D$ and $\tilde{\varphi}_{\mathcal{T}}$ is an RFF for $k_{\mathcal{T}}$, such that $\tilde{\varphi}(\tau) \in \mathbb{R}^{D'}$ then an ORFF map for*

$$K(x, z) = k_X(x, z) I_{\mathcal{H}_{k_j}}$$

is given for all $x \in \mathbb{R}$, all $\tau \in \mathbb{R}$ and all $\Theta \in \mathcal{M}_{D, D'}(\mathbb{R})$ by

$$\left(\tilde{\Phi}_K(x)^* \Theta \right) (\tau) = \tilde{\varphi}_X(x)^* \Theta \tilde{\varphi}_{\mathcal{T}}(\tau)$$

and

$$\left(\tilde{\Phi}_K(x)G \right) (\tau) = \tilde{\varphi}_X(x) \tilde{\varphi}_{\mathcal{T}}(\tau)^* G,$$

where $g \in \mathbb{R}^{D'}$.

Moreover if one defines $\tilde{\Phi}_K(x, \tau) = \left(\tilde{\Phi}_K(x)^* \Theta \right) (\tau)$ one have of course

$$\tilde{\Phi}_K(x, \tau)^* \Theta = \tilde{\varphi}_X(x)^* \Theta \tilde{\varphi}_{\mathcal{T}}(\tau)$$

9.1.4 Many quantile regression

From the loss defined in [Equation 9.1](#) we defined the regularized risk using the “continuous” pinball loss for the quantile regression problem. For all $f \in \mathcal{H}_K$,

$$\mathfrak{R}_\lambda(f, s) = \frac{1}{N} \sum_{i=1}^N \int_{[0,1]} \begin{cases} \tau (f_{x_i}(\tau) - y_i) & \text{if } f_{x_i}(\tau) \geq y_i \\ (1-\tau)(y_i - f_{x_i}(\tau)) & \text{otherwise} \end{cases} + \lambda \|f\|_K^2.$$

The issue with the above risk is that the different quantile for a given point $x \in \mathbb{R}$ may cross (see Sangnier, Fercoq, and Buc [[146](#)]). To avoid this to happen we need to force the function $f_x(\tau)$ to be *increasing* in τ for any $x \in \mathbb{R}$. Because a decreasing function has a negative derivative we can add a penalty term to the risk to avoid $f_x(\tau)$ to be decreasing in τ .

$$\Omega_{\text{cross}}(f_{x_i}) = -\min \left(\frac{\partial f_{x_i}}{\partial \tau}(\tau), 0 \right)$$

Thus the regularized risk with the no crossing constraint is

$$\mathfrak{R}_{\lambda_1, \lambda_2}(f, s) = \frac{1}{N} \sum_{i=1}^N \int_{[0,1]} \begin{cases} \tau (f_{x_i}(\tau) - y_i) & \text{if } f_{x_i}(\tau) \geq y_i \\ (1-\tau)(y_i - f_{x_i}(\tau)) & \text{otherwise} \end{cases} - \lambda_1 \min \left(\frac{\partial f_{x_i}}{\partial \tau}(\tau), 0 \right) + \lambda_2 \|f\|_K^2.$$

Eventually we replace the integral by a Monte-Carlo sampling with the uniform law $\mathcal{U}([0, 1])$ and plug in the approximate function of f using the ORFF map proposed in [conjecture 9.1](#). The final regularized risk to be minimized reads

$$\mathfrak{R}_{\lambda_1, \lambda_2}(\Theta, s) = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \begin{cases} \tau_t (\tilde{f}_{x_i}(\tau_t) - y_i) & \text{if } \tilde{f}_{x_i}(\tau_t) \geq y_i \\ (1-\tau_t)(y_i - \tilde{f}_{x_i}(\tau_t)) & \text{otherwise} \end{cases} - \lambda_1 \min \left(\frac{\partial \tilde{f}_{x_i}}{\partial \tau}(\tau_t), 0 \right) + \lambda_2 \|\Theta\|_{\text{fro}}^2.$$

where $\tilde{f}_x(\tau) = \tilde{\varphi}_x(x)^* \Theta \tilde{\varphi}_T(\tau)$, $\tau_t \sim \mathcal{U}([0, 1])$ and

$$\begin{aligned} \frac{\partial \tilde{f}_x}{\partial \tau}(\tau) &= \tilde{\varphi}_x(x)^* \Theta \frac{\partial \tilde{\varphi}_T}{\partial \tau}(\tau) \\ &= \tilde{\varphi}_x(x)^* \Theta \bigoplus_{k=1}^{D'} \begin{pmatrix} -\omega'_k \sin(\omega'_k \tau) \\ \omega'_k \cos(\omega'_k \tau) \end{pmatrix}, \quad \omega'_k \sim \mathcal{F}[k_T] \text{ i. i. d..} \end{aligned}$$

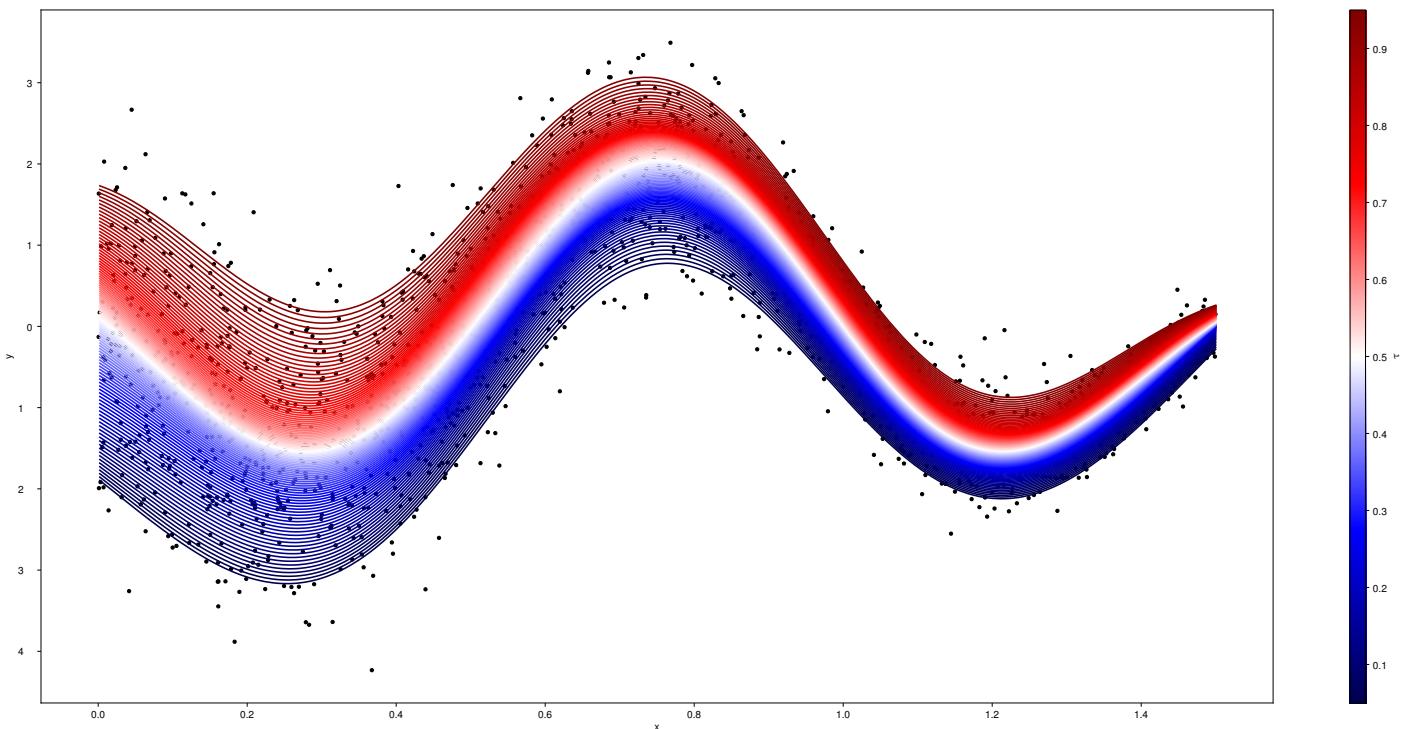


Figure 9.3: Learning a continuous quantile function with ORFF regression.

9.1.4.1 Some results

We minimized the quantity $\mathfrak{R}_{\lambda_1, \lambda_2}(\Theta, s)$ on a toy dataset: a sine wave with some heteroscedastic noise. First we compared our methodology to the joint quantile regression proposed in Sangnier, Fercoq, and Buc [146]. We generate $N = 2500$ for the train set and $N' = 1000$ points for the test set and use a Gaussian kernel for both k_X and k_T . We choosed $\sigma_X = 0.25$ and σ_T has been set to be the median of the pairwise distance of the τ_t 's drawn randomly from $\mathcal{U}([0, 1])$. Notice that $\mathfrak{R}_{\lambda_1, \lambda_2}(\Theta, s)$ is convex in Θ . To avoid computing complex gradients and by lack of time, we used Tensorflow [1] to perform a gradient descent (with RMSProp [172]) with automatic symbolic differentiation. [Figure 9.1](#) show the result for the quantile at 0.05, 0.275, 0.5, 0.775 and 0.95 using the ORFF methodology. [Figure 9.2](#) shows the joint quantile regression of Sangnier, Fercoq, and Buc [146] on the same dataset. Not only our method matched the the performances of Sangnier, Fercoq, and Buc [146]¹ but we cutted down the computation time from circa 1330 seconds to circa 30 seconds (training and testing). Moreover on contrary to Sangnier, Fercoq, and Buc [146] we have access to all the quantile of the model (see [Figure 9.3](#)).

9.1.5 One-class SVM revisited

We also propose an extension of the celebrated One-Class Support Vector Machine (OCSM) such that it is possible to learn jointly all the level sets. One-class classification, also known as unary classification, tries to identify objects of a specific class amongst all objects, by learning from a training set containing only the objects of that class. In this framework, we assume that we only observe examples of one class (referred to as the inlier class). The second class is called outlier class. We turn our attention to the OCSM of Schölkopf et al. [147] which extends the Support Vector Machine (SVM) methodology [48, 155] to handle training using only inliers.

We recall that given an hyperparameter $\nu \in [0, 1]$ that controls the proportion of inlier, given as scalar kernel k , the OCSM problem reads

$$\arg \min_{f \in \mathcal{H}_k, \tau \in \mathbb{R}} \frac{1}{2} \|f\|_{\mathcal{H}_k}^2 - \nu \tau + \frac{1}{N} \sum_{i=1}^N \max(\tau - f(x_i), 0)$$

The decision function is then

$$h(x, \tau) = \mathbb{1}_{[\tau, \infty)}(f(x)).$$

¹ We reported an error computed with the pinball loss on the test set of 0.818 for our method and 0.817 for joint regression (note that we don't report here an average on many experiments to avoid randomness introduced by the random features, but the results seems robust in practice.)

As in [Subsection 9.1.1](#) we can rewrite the optimization problem as an integral over all the value of v and suppose that f is function-valued (a function of v). Moreover τ must also change its value according to v . Thus given a kernel $k_{\mathcal{X}}$ on the inputs $x \in \mathbb{R}^d$ with its approximate feature map $\tilde{\varphi}_{\mathcal{X}}$ and a kernel $k_{\mathcal{T}}$ on the level sets with its approximate feature map $\tilde{\varphi}_{\mathcal{T}}$, we define the continuous one-class SVM problem as

$$\begin{aligned} & \arg \min_{f \in \mathcal{H}_K, \tau \in \mathcal{H}_{k_{\mathcal{T}}}} \frac{1}{N} \sum_{i=1}^N \int_{[0,1]} \max(0, \tau(v) - f_{x_i}(v)) dv \\ & + \frac{1}{2} \int_{[0,1]} v \|f.(v)\|_{\mathcal{H}_{k_{\mathcal{X}}}}^2 dv - \int_{[0,1]} v \tau(v) dv. \end{aligned}$$

Again we can compute the integral by Monte-Carlo sampling and replace f and τ by their respective approximation. Notice that the RKHS of τ should match the RKHS of the output space of f_K . Hence

$$\begin{aligned} & \arg \min_{\Theta \in \mathcal{M}_{D,D'}(\mathbb{R}), \tau \in \mathbb{R}^D} \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \max(0, \tilde{\tau}(v_t) - \tilde{f}_{x_i}(v_t)) \\ & + \frac{1}{2T} \sum_{t=1}^T v_t \|\tilde{f}_(v_t)\|_2^2 - \frac{1}{T} \sum_{t=1}^T v_t \tilde{\tau}(v_t), \end{aligned}$$

where $v_t \sim \mathcal{U}([0, 1])$ i. i. d., $\tilde{\tau}(v) = \tilde{\varphi}_{\mathcal{T}}(v)$ and $\tilde{f}_x(v) = \tilde{\varphi}(x)^* \Theta \tilde{\varphi}_{\mathcal{T}}(v)$. We also deduce that $f.(v) = \Theta \tilde{\varphi}_{\mathcal{T}}(v)$. Here the natural decision function is

$$h(x, v) = \mathbb{1}_{[\tilde{\tau}(v), \infty)} (\tilde{f}_x(v)). \quad (9.3)$$

9.1.5.1 Proof of concept

First we ensure that the variable v is a good proxy for the proportion of inlier. For this we generate a dataset of points in $\mathcal{X} = \mathbb{R}^2$ from a mixture of three Gaussians. One Gaussian is located at $\mu_1 = (0, 0)$, the second at $\mu_2 = (5, 5)$ and the third at $\mu_3 = (10, 10)$. Each Gaussian has unit variance and we draw 250 points from the first and third one and 100 from the second one, so that we have 600 points in the dataset. We take $k_{\mathcal{X}}$ as a Gaussian kernel with scale parameter $\gamma_{\mathcal{X}} = 2$ and $k_{\mathcal{T}}$ another Gaussian kernel with scale parameter $\gamma_{\mathcal{T}}$ (see [Subsubsection 4.2.2.2](#)). After training we apply the decision function to the train test to which we add 100 points generated from a uniform distribution to model the novelty detection setting. We refer to this new augmented set as the test set. In [Figure 9.4](#) we show the proportion of inlier with respect to v . The top figure shows the result for a model trained with $\gamma_{\mathcal{T}} = 100$ the middle figure for $\gamma_{\mathcal{T}} = 1$ and the bottom figure for $\gamma_{\mathcal{T}} = 0.01$. We see that when $\gamma_{\mathcal{T}} = 0.01$, the proportion of inlier on the train and test set does not follow the theoretical black curve, because the algorithm regularizes too much between the level

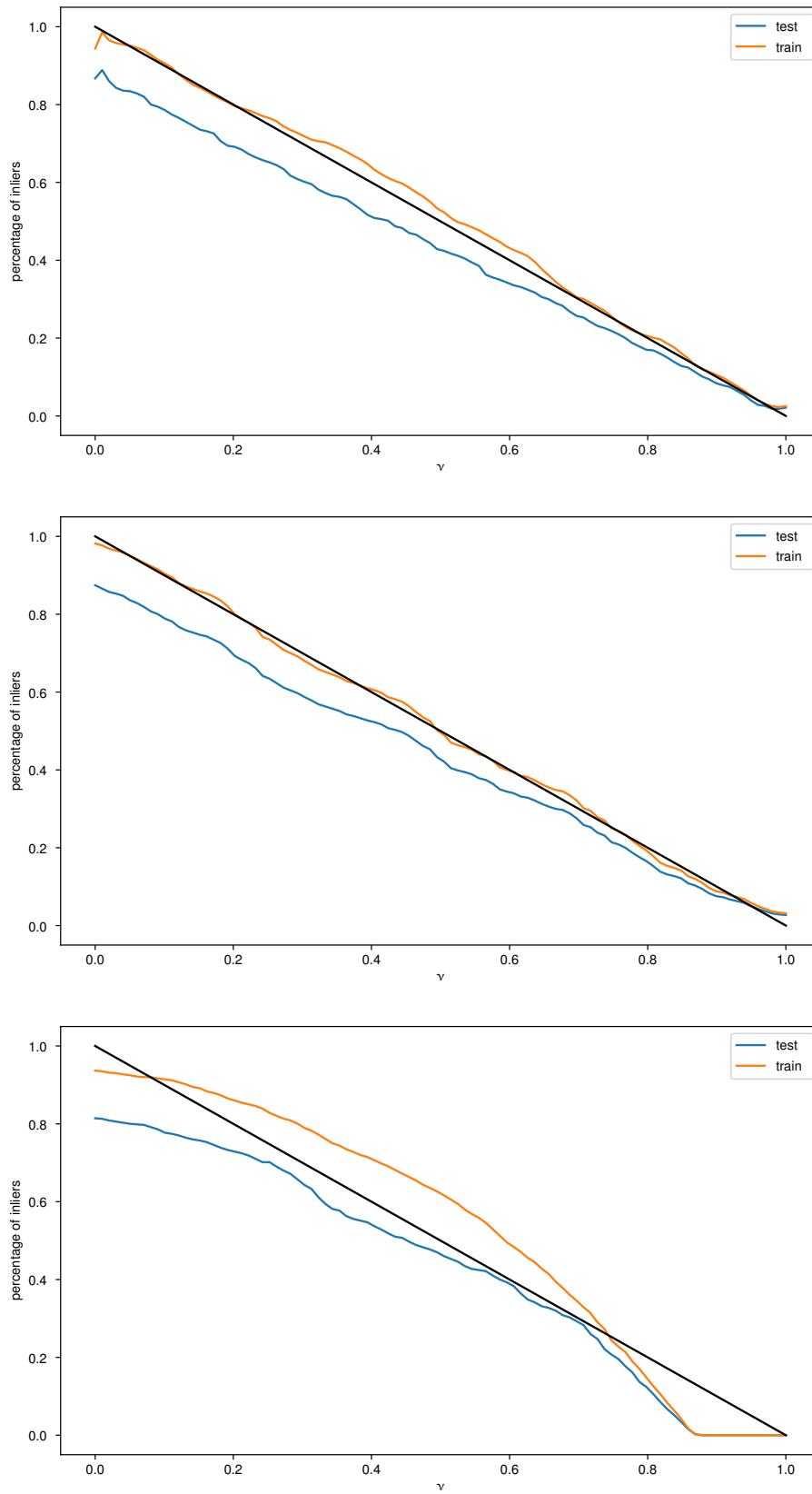


Figure 9.4: Continuous OCSVVM: proportion of inlier with respect to γ . Top figure corresponds to $\gamma_T = 100$, middle to 1 and bottom to 0.01.

sets. When $\gamma_T = 1$ or 100 the proportion of inliers almost follows the theoretical black curve. We see that when $\gamma_T = 100$ the curves are less stable than when $\gamma_T = 0$ especially around $v = 0$. The gap between the train curve (orange) and test curve (blue) corresponds to the “novel” points that are not distributed according to the mixture of Gaussians.

We can give an interpretation to k_T and k_X . k_X control the complexity of the boundary of each level set and k_T tells how much each level set differ from the neighbour level sets.

We propose a second experiment in a context of outlier detection. This time the train set is polluted with outliers. We replicate the example given in the documentations of Scikit-Learn at http://scikit-learn.org/stable/auto_examples/covariance/plot_outlier_detection.html We compare our method to three other well known outlier detection methods from the literature: Isolation Forest [105], OCSM [147] and a Robust Covariance estimator [38, 132] on Figure 9.5. Our method achieves the state of the art on this simple example which is encouraging. However the computation time of our continuous OCSM is higher than the other methods. It took circa 0.25 second for the OCSM. 10 seconds for our method, 5 seconds for isolation forest and 0.1 second for the Robust covariance estimator. This can be due to the implementation since we used a (sub-optimal) hand-crafted full gradient descent. Notice that however our method is able to retrieve all the level sets after training, not only the one presented in Figure 9.5. When one is interested in a specific level set or range of level set one could sample the v_t from another distribution than the uniform distribution $\mathcal{U}[0, 1]$ to give more importance to the desired range of level sets.

9.2 OPERALIB

During this Thesis we started the development of a library named “Operalib” implementing various machine learning algorithms based on operator-valued kernels. We are grateful to Alexandre Gramfort (LTCI, Télécom ParisTech) who served as a technical mentor at the beginning of this software development and provided many advices. Operator-valued kernels defines a framework allowing learning vector/function/structured output. To install the library it should be as simple as

Listing 9.1: Installation of Operalib.

```
1 pip install operalib
```

The library currently features:

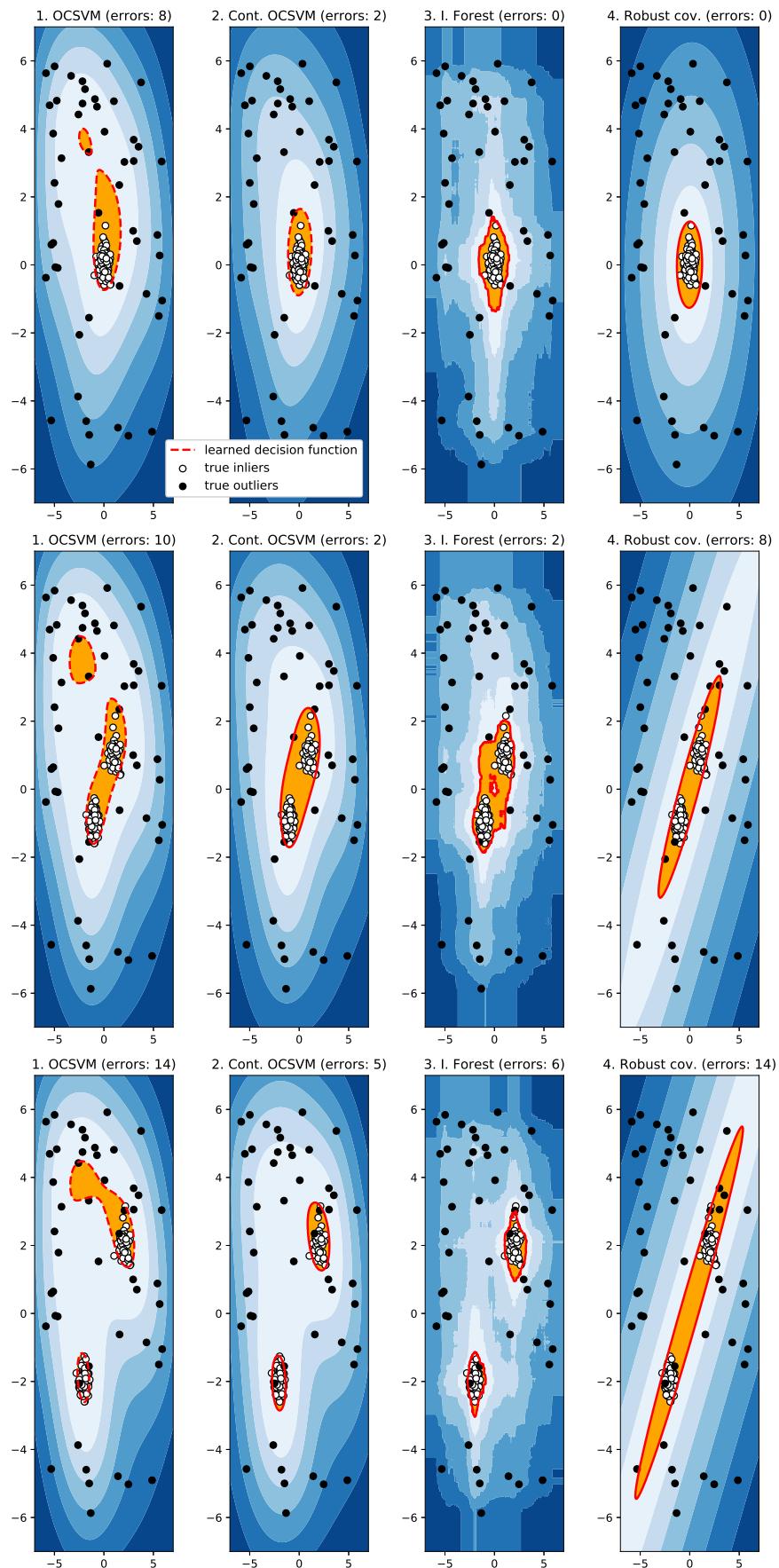


Figure 9.5: Continuous OCSVM for outlier detection.

- Quantile regression [146],
- ONORMA [10],
- semi-supervised Ridge regression [35],
- some elements of the ORFF framework [29].

The algorithms work for a selection of popular operator-valued kernels such that the matrix-valued decomposable kernel, the curl-free kernel and the divergence-free kernel. The library is structured so that it is easy for the user to define its own operator-valued kernel and plug it to the existing optimisation algorithms, while keeping efficient computations thanks to the methodology presented in [Section 6.3](#) (i.e. by seeing operator-valued kernels as operators along with matrix-free solver rather than plain matrices). We designed the library in order to have a close compatibility with Scikit-learn. Code and documentation are publicly available at <https://github.com/operalib/operalib>. In a near future we plan to add the family of works of Brouard, d’Alché-Buc, and Szafranski [35] around Input Output Kernel Regression, the work of Lim et al. [103] about the two learning algorithms defined in Lim et al. [103]: a sparse learning of OVK and a learning algorithm for both kernel and weights with a block coordinate descent scheme and a proximal gradient method to deal with non-smooth constraints. Development for modeling time series will be also included. We hope to expand with more algorithms from various authors of the OVK community and welcome any new contributor!



10

CONCLUSION

To conclude this work we would like to summarize our contributions, and show how they answered the initial question of large-scale learning with Operator-Valued Kernels. Then we finish with some short and mid term perspectives.

Contents

10.1 Contributions	164
10.2 Perspectives	165

In Machine Learning, many algorithms focus on learning functions that model dependencies between inputs and outputs where outputs are real numbers. However in numerous application fields such as biology, economics, physics, etc. the output data are not reals: they can be a collection of reals, present complex structures or can even be functions. To overcome this difficulty and take into account the structure of the data, a common approach is to see them as *vectors* of some Hilbert space. From this observation, in this thesis, we took interest in vector-valued functions. Looking at the literature we focused on mathematical objects called Operator-Valued Kernels to learn such functions.

10.1 CONTRIBUTIONS

OVKs naturally extend the celebrated kernel methods used to learn scalar-valued functions, to the case of learning vector-valued functions. Yet, although OVKs are appealing from a theoretical aspect, these methods scale poorly in terms of computation time when the number of data is high. Indeed, in order to evaluate a *function* on a unknown point with an Operator-Valued Kernel, it requires to evaluate an Operator-Valued Kernel on all the point in the given dataset. Hence naive learning with kernels usually scales cubically in time with the number of data. In the context of large-scale learning such scaling is not acceptable. Through this work we propose a methodology to tackle this difficulty.

Enlightened by the literature on large-scale learning with *scalar*-valued kernel, in particular the work of Rahimi and Recht [139], we propose to replace an OVK by a random feature map that we called Operator-valued Random Fourier Feature. Our contributions start with the formal mathematical construction of this feature from an OVK. Then we show that it is also possible to obtain a kernel from an ORFF. Eventually we analyse the regularization properties in terms of Fourier Transform of γ -Mercer kernels. Then we moved on giving a bound on the error due to the random approximation of the OVK with high probability. We showed that it is possible to bound the error even though the ORFF estimator of an OVK is not a bounded random variable. Moreover we also give a bound when the dimension of the output data infinite.

After ensuring that an ORFF is a good approximation of a kernel, we moved on giving a framework for supervised learning with Operator-Valued Kernels. We showed that learning with a feature map is equivalent to learning with the reconstructed OVK under some mild conditions. Then we focused on an efficient implementation of ORFF by viewing them as linear operators rather than matrices

and using matrix-free (iterative) solvers and concluded with some numerical experiments. Eventually we gave a generalization bound for ORFF learning that suggests that the number of features sampled in an ORFF should be proportional to the number of data. We concluded our contribution by applying the ORFF framework to learning vector-valued time series.

10.2 PERSPECTIVES

To start with the theoretical perspectives, following Rahimi and Recht we gave a generalization bound for ORFF kernel ridge that suggests that the number of features to draw is proportional to the number of data. However new results of Rudi, Camoriano, and Rosasco [144] suggest that the number of feature should be proportional to the *square root* of the number of data. In a future work, we shall investigate this result and extend it to ORFF.

On the methodological perspectives we gave an intuition on how Operator-Valued Kernels can be used to learn outputs that are functions. We used the ORFF framework to speed up quantile regression and at the same time obtain the full quantile function. We applied the same methodology to the anomaly detection setting and showed that it is possible to learn jointly all the level sets of a distribution with an extension of a One-Class Support Vector Machine. We are convinced that this will open the door to many new applications. Given a problem with some hyperparameters, the combination of ORFF and Operator-Valued Kernels allow to learn functions of the hyperparameters.

Another nice extension would be to be able to learn the structure of an ORFF i. e. the spectral distribution and the operator from the data, as in Yang et al. [191] so that we avoid to inject directly ourselves a prior on the data by the mean of an Operator-Valued Kernel.

On the implementation level, we are really enthusiastic about Operalib, a library for learning with Operator-Valued Kernels started during this thesis as a project of Paris-Saclay Center for Data Science¹, and will extend the library with other OVK-based algorithms. Moreover much work is remaining to do concerning the implementation of efficient algorithms based on (O)RFFs. We could extend the Multiple Kernel learning setting to OVKs to see if we can match the performances of Deep Neural Networks as in Lu et al. [108]. We could also improve the Doubly Stochastic Gradient descent in the light of the recent results of Rudi, Camoriano, and Rosasco [144] on generalization.

¹ As a collaboration with Alexandre Gramfort

Eventually during these three years, we have witnessed the rise of deep-learning methods with neural networks. As pointed out by many authors, random features share deep connections with neural networks: an ORFF-based shallow architecture can be seen as a one-layer neural architecture. Conversely, a neural network can be seen as a compositional feature map. As in the work of Yang et al. [192] we could replace the last layer of a convolutional neural network [95] with an ORFF map in order to open these architectures to the setting offered by OVK to deal with structured and functional outputs.

Part IV
APPENDIX

A

PROOFS OF THEOREMS

In this appendix we detail the proofs of [Corollary 5.2](#) and [Corollary 5.3](#). These two corollaries applying on compact subsets of Banach spaces are the consequences of more generic propositions ([Proposition A.1](#) and [Proposition A.2](#)) working on any compact subsets of Polish spaces. Eventually we give a proof on the variance bound given in [Proposition A.3](#).

Contents

A.1	Proof of the error bound with high probability of the ORFF estimator	170
A.1.1	Epsilon-net	170
A.1.2	Bounding the Lipschitz constant	171
A.1.3	Bounding the error on a given anchor point	172
A.1.4	Union Bound and examples	176
A.2	Proof of the ORFF estimator variance bound	181

A.1 PROOF OF THE ERROR BOUND WITH HIGH PROBABILITY OF THE ORFF ESTIMATOR

We recall the notations $\delta = x \star z^{-1}$, for all $x, z \in \mathcal{X}$, $\tilde{K}(x, z) = \tilde{\Phi}(x)^* \tilde{\Phi}(z)$, $\tilde{K}^j(x, z) = \Phi_x(\omega_j)^* \Phi_z(\omega_j)$, where $\omega_j \sim \widehat{\text{Pr}_{\text{Haar}, \rho}}$ and $K_e(\delta) = K(x, z)$ and $\tilde{K}_e(\delta) = \tilde{K}(x, z)$. For the sake of readability, we use throughout the proof the quantities

$$\begin{aligned} F(\delta) &:= \tilde{K}(x, z) - K(x, z) \\ F^j(\delta) &:= \frac{1}{D} (\tilde{K}^j(x, z) - K(x, z)). \end{aligned}$$

We also view \mathcal{X} as a metric space endowed with the distance $d_{\mathcal{X}} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$. Compared to the scalar case, the proof follows the same scheme as the one described in [139, 167], but we consider an operator norm as measure of the error and therefore concentration inequality dealing with these operator norm. The main feature of [Proposition A.1](#) is that it covers the case of bounded ORFF as well as unbounded ORFF. In the case of bounded ORFF, a Bernstein inequality for matrix concentration such that the one proved in Mackey et al. [110, Corollary 5.2] or the formulation of Tropp [174] recalled in Koltchinskii [92] is suitable. However some kernels like the curl and the divergence-free kernels do not have obvious bounded $\|F^j\|_{\mathcal{Y}, \mathcal{Y}}$ but exhibit F^j with subexponential tails. Therefore, we use an operator Bernstein concentration inequality adapted for random matrices with subexponential norms.

A.1.1 *Epsilon*-net

Let $\mathcal{C} \subseteq \mathcal{X}$ be a compact subset of \mathcal{X} . Let $\mathcal{D}_{\mathcal{C}} = \{x \star z^{-1} \mid x, z \in \mathcal{C}\}$ with diameter at most $2|\mathcal{C}|$ where $|\mathcal{C}|$ is the diameter of \mathcal{C} . Since \mathcal{C} is supposed compact, so is $\mathcal{D}_{\mathcal{C}}$. Since $\mathcal{D}_{\mathcal{C}}$ is also a metric space it is well known that a compact metric space is totally bounded. Thus it is possible to find a finite ϵ -net covering $\mathcal{D}_{\mathcal{C}}$. We call $T = N(\mathcal{D}_{\mathcal{C}}, r)$ the number of closed balls of radius r required to cover $\mathcal{D}_{\mathcal{C}}$. For instance if $\mathcal{D}_{\mathcal{C}}$ is a subspace finite dimensional Banach space with diameter at most $2|\mathcal{C}|$ it is possible to cover the space with at most $T = (4|\mathcal{C}|/r)^d$ balls of radius r (see Cucker and Smale [50, proposition 5]).

Let us call $\delta_i, i = 1, \dots, T$ the center of the i -th ball, also called anchor of the ϵ -net. Denote L_F the Lipschitz constant of F . Let $\|\cdot\|_{\mathcal{Y}, \mathcal{Y}}$ be the operator norm on $\mathcal{L}(\mathcal{Y})$ (largest eigenvalue). We introduce the following technical lemma.

Lemma A.1 $\forall \delta \in \mathcal{D}_{\mathcal{C}}$, if

$$L_F \leq \frac{\epsilon}{2r} \tag{A.1}$$

and

$$\|\mathbb{F}(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}} \leq \frac{\epsilon}{2}, \quad \text{for all } i \in \mathbb{N}_T^* \quad (A.2)$$

then $\|\mathbb{F}(\delta)\|_{\mathcal{Y}, \mathcal{Y}} \leq \epsilon$.

Proof

$$\begin{aligned} \|\mathbb{F}(\delta)\|_{\mathcal{Y}, \mathcal{Y}} &= \|\mathbb{F}(\delta) - \mathbb{F}(\delta_i) + \mathbb{F}(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}} \\ &\leq \|\mathbb{F}(\delta) - \mathbb{F}(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}} + \|\mathbb{F}(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}} \end{aligned}$$

for all $0 < i < T$. Using the Lipschitz continuity of \mathbb{F} we have

$$\|\mathbb{F}(\delta) - \mathbb{F}(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}} \leq d_{\mathcal{X}}(\delta, \delta_i)L_{\mathbb{F}} \leq rL_{\mathbb{F}}$$

hence

$$\|\mathbb{F}(\delta)\|_{\mathcal{Y}, \mathcal{Y}} \leq rL_{\mathbb{F}} + \|\mathbb{F}(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}} = \frac{r\epsilon}{2r} + \frac{\epsilon}{2} = \epsilon.$$

To apply the lemma, we must bound the Lipschitz constant of the operator-valued function \mathbb{F} (Equation A.1) and $\|\mathbb{F}(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}}$, for all $i = 1, \dots, T$ as well (Equation A.2).

A.1.2 Bounding the Lipschitz constant

This proof is a slight generalization of Minh [118] to arbitrary metric spaces. It differ from our first approach [29], based on the proof of Sutherland and Schneider [167] which was only valid for a finite dimensional input space \mathcal{X} and imposed a twice differentiability condition on the considered kernel.

Lemma A.2 *Let $H_{\omega} \in \mathbb{R}_+$ be the Lipschitz constant of $h_{\omega}(\cdot)$ and assume that*

$$\int_{\widehat{\mathcal{X}}} H_{\omega} \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} d\Pr_{\widehat{\text{Haar}}, \rho}(\omega) < \infty.$$

Then the operator-valued function $K_e : \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ is Lipschitz with

$$\|K_e(x) - K_e(z)\|_{\mathcal{Y}, \mathcal{Y}} \leq d_{\mathcal{X}}(x, z) \int_{\widehat{\mathcal{X}}} H_{\omega} \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} d\Pr_{\widehat{\text{Haar}}, \rho}(\omega) \quad (A.3)$$

Proof We use the fact that the cosine function is Lipschitz with constant 1 and h_{ω} Lipschitz with constant H_{ω} . For all $x, z \in \mathcal{X}$ we have

$$\begin{aligned} \|\tilde{K}_e(x) - K_e(z)\|_{\mathcal{Y}, \mathcal{Y}} &= \left\| \int_{\widehat{\mathcal{X}}} (\cos h_{\omega}(x) - \cos h_{\omega}(z)) A(\omega) d\Pr_{\widehat{\text{Haar}}, \rho} \right\|_{\mathcal{Y}, \mathcal{Y}} \\ &\leq \int_{\widehat{\mathcal{X}}} |\cos h_{\omega}(x) - \cos h_{\omega}(z)| \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} d\Pr_{\widehat{\text{Haar}}, \rho} \\ &\leq \int_{\widehat{\mathcal{X}}} |h_{\omega}(x) - h_{\omega}(z)| \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} d\Pr_{\widehat{\text{Haar}}, \rho} \\ &\leq d_{\mathcal{X}}(x, z) \int_{\widehat{\mathcal{X}}} H_{\omega} \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} d\Pr_{\widehat{\text{Haar}}, \rho} \end{aligned}$$

In the same way, considering $\tilde{K}_e(\delta) = \frac{1}{D} \sum_{j=1}^D \cos h_{\omega_j}(\delta) A(\omega_j)$, where $\omega_j \sim \widehat{\Pr}_{\text{Haar}, \rho}$, we can show that \tilde{K}_e is Lipschitz with

$$\|\tilde{K}_e(x) - \tilde{K}_e(z)\|_{y,y} \leq d_X(x, z) \frac{1}{D} \sum_{j=1}^D H_{\omega_j} \|A(\omega_j)\|_{y,y}.$$

Combining the Lipschitz continuity of \tilde{K}_e and \tilde{K} (Lemma A.2) we obtain

$$\begin{aligned} \|F(x) - F(z)\|_{y,y} &= \|\tilde{K}_e(x) - \tilde{K}_e(z) + K_e(z)\|_{y,y} \\ &\leq \|\tilde{K}_e(x) - \tilde{K}_e(z)\|_{y,y} + \|K_e(x) - K_e(z)\|_{y,y} \\ &\leq d_X(x, z) \left(\int_{\widehat{\mathcal{X}}} H_{\omega} \|A(\omega)\|_{y,y} d\Pr_{\text{Haar}, \rho} \right. \\ &\quad \left. + \frac{1}{D} \sum_{j=1}^D H_{\omega_j} \|A(\omega_j)\|_{y,y} \right) \end{aligned}$$

Taking the expectation yields

$$E_{\widehat{\Pr}_{\text{Haar}, \rho}} [L_F] = 2 \int_{\widehat{\mathcal{X}}} H_{\omega} \|A(\omega)\|_{y,y} d\Pr_{\text{Haar}, \rho} (\omega)$$

Thus by Markov's inequality,

$$\begin{aligned} \Pr_{\widehat{\Pr}_{\text{Haar}, \rho}} \{ (\omega_j)_{j=1}^D \mid L_F \geq \epsilon \} &\leq \frac{E_{\widehat{\Pr}_{\text{Haar}, \rho}} [L_F]}{\epsilon} \\ &\leq \frac{2}{\epsilon} \int_{\widehat{\mathcal{X}}} H_{\omega} \|A(\omega)\|_{y,y} d\Pr_{\text{Haar}, \rho}. \end{aligned} \tag{A.4}$$

A.1.3 Bounding F on a given anchor point δ_i

To bound $\|F(\delta_i)\|_{y,y}$, Hoeffding inequality devoted to matrix concentration [110] can be applied. We prefer here to turn to tighter and refined inequalities such as Matrix Bernstein inequalities (Sutherland and Schneider [167] also pointed that for the scalar case). The first non-commutative (matrix) concentration inequalities are due to the pioneer work of Ahlswede and Winter [3], using bound on the moment generating function. This gave rise to many applications Koltchinskii [92], Oliveira [128], and Tropp [174] ranging from analysis of randomized optimization algorithm to analysis of random graphs and generalization bounds useful in machine learning. The following inequality has been proposed in [92].

Theorem A.1 (Bounded non-commutative Bernstein). *From Theorem 3 of Koltchinskii [92], consider a sequence $(X_j)_{j=1}^D$ of D independent Hermitian $p \times p$ random matrices acting on a finite dimensional Hilbert space \mathcal{Y} that satisfy $E X_j = 0$, and suppose that there exist some constant*

$U \geq \|X_j\|_{y,y}$ for each index j . Denote the proxy bound on the matrix variance

$$V \succcurlyeq \sum_{j=1}^D E X_j^2.$$

Then, for all $\epsilon \geq 0$,

$$\Pr \left\{ \left\| \sum_{j=1}^D X_j \right\|_{y,y} \geq \epsilon \right\} \leq p \exp \left(-\frac{\epsilon^2}{2\|V\|_{y,y} + 2U\epsilon/3} \right)$$

This bound we used in our original paper [29] has the default to grow linearly with the dimension p of the output space y . However if the evaluation of the operator-valued kernel at two points yields a low-rank matrix, this bound could be improved since only a few principal dimensions are relevant. Moreover this bound cannot be used when dealing with operator-valued kernel acting on infinite dimensional Hilbert spaces. Recent results of Minsker [121] consider the notion of intrinsic dimension to avoid this “curse of dimensionality”.

Definition A.1 Let A be a trace class operator acting on a Hilbert space y . We call intrinsic dimension the quantity

$$\text{IntDim}(A) = \frac{\text{Tr}[A]}{\|A\|_{y,y}}.$$

When A is approximately low-rank (i.e. many eigenvalues are small), or go quickly to zero, the intrinsic dimension can be much lower than the dimensionality. Indeed,

$$1 \leq \text{IntDim}(A) \leq \text{Rank}(A) \leq \dim(A).$$

Theorem A.2 (Bounded non-commutative Bernstein with intrinsic dimension [121, 175]). Consider a sequence $(X_j)_{j=1}^D$ of D independent Hilbert-Schmidt self-adjoint random operators acting on a separable Hilbert y space that satisfy $EX_j = 0$ for all $j \in \mathbb{N}_D^*$. Suppose that there exist some constant $U \geq 2\|X_j\|_{y,y}$ almost surely for all $j \in \mathbb{N}_D^*$. Define a semi-definite upper bound for the the operator-valued variance

$$V \succcurlyeq \sum_{j=1}^D E X_j^2.$$

Then for all $\epsilon \geq \sqrt{\|V\|_{y,y}} + U/3$,

$$\Pr \left\{ \left\| \sum_{j=1}^D X_j \right\|_{y,y} \geq \epsilon \right\} \leq 4 \text{IntDim}(V) \exp(-\psi_{V,U}(\epsilon))$$

$$\text{where } \psi_{V,U}(\epsilon) = \frac{\epsilon^2}{2\|V\|_{y,y} + 2U\epsilon/3}$$

Essentially, compared to [Theorem A.1](#), [Theorem A.2](#) replace the dimension of \mathcal{Y} by four times the intrinsic dimension of the variance of the matrix valued random variable. The concentration inequality is restricted to the case where $\epsilon \geq \sqrt{\|V\|_{\mathcal{Y},\mathcal{Y}}} + U/3$ since the probability is vacuous on the contrary. The assumption that X_j 's are Hilbert-Schmidt operators comes from the fact that the product of two such operator yields a trace-class operator, for which the intrinsic dimension is well defined.

However, to cover the general case including unbounded ORFFs like curl and divergence-free ORFFs, we choose a version of Bernstein matrix concentration inequality proposed in [\[92\]](#) that allows to consider matrices that are not uniformly bounded but have subexponential tails. In the following we use the notion of Orlicz norm to bound random variable by their tail behavior rather than their value.

Definition A.2 (Orlicz norm). We follow the definition given by Koltchinskii [\[92\]](#). Let $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ be a non-decreasing convex function with $\psi(0) = 0$. For a random variable X on a measured space $(\Omega, \mathcal{T}(\Omega), \mu)$

$$\|X\|_\psi := \inf \{ C > 0 \mid \mathbf{E}[\psi(|X|/C)] \leq 1 \}.$$

For the sake of simplicity, we now fix $\psi(t) = \psi_1(t) = \exp(t) - 1$. Although the Orlicz norm should be adapted to the tail of the distribution of the random operator we want to quantify to obtain the sharpest bounds. We also introduce two technical lemmas related to Orlicz norm. The first one relates the ψ_1 -Orlicz norm to the moment generating function (MGF).

Lemma A.3 Let X be a random variable with a strictly monotonic moment-generating function. We have $\|X\|_{\psi_1}^{-1} = MGF_{|X|}^{-1}(2)$.

Proof We have

$$\begin{aligned} \|X\|_{\psi_1} &= \inf \{ C > 0 \mid \mathbf{E}[\exp(|X|/C)] \leq 2 \} \\ &= \frac{1}{\sup \{ C > 0 \mid MGF_{|X|}(C) \leq 2 \}} \end{aligned}$$

X has strictly monotonic moment-generating thus $C^{-1} = MGF_{|X|}^{-1}(2)$. Hence

$$\|X\|_{\psi_1}^{-1} = MGF_{|X|}^{-1}(2).$$

The second lemma gives the Orlicz norm of a positive constant.

Lemma A.4 If $a \in \mathbb{R}_+$ then $\|a\|_{\psi_1} = \frac{a}{\ln(2)} < 2a$.

Proof We consider a as a positive constant random variable, whose Moment Generating Function (MGF) is

$$MGF_a(t) = \exp(at).$$

From [Lemma A.3](#), $\|a\|_{\psi_1} = \frac{1}{MGF_a^{-1}(2)}$. Then $MGF_a^{-1}(2) = \frac{\ln(2)}{|a|}$, $a \neq 0$. If $a = 0$ then $\|a\|_{\psi_1} = 0$ by definition of a norm. Thus $\|a\|_{\psi_1} = \frac{a}{\ln(2)}$. \square

We now turn our attention to Minsker [121]’s theorem to for unbounded random variables.

Theorem A.3 (Unbounded non-commutative Bernstein with intrinsic dimension). Consider a sequence $(X_j)_{j=1}^D$ of D independent self-adjoint random operators acting on a finite dimensional Hilbert space \mathcal{Y} of dimension p that satisfy $\mathbf{E} X_j = 0$ for all $j \in \mathbb{N}_D^*$. Suppose that there exist some constant $U \geq \|\|X_j\|_{\mathcal{Y}, \mathcal{Y}}\|_{\psi}$ for all $j \in \mathbb{N}_D^*$. Define a semi-definite upper bound for the the operator-valued variance

$$V \succcurlyeq \sum_{j=1}^D \mathbf{E} X_j^2.$$

Then for all $\epsilon > 0$,

$$\begin{aligned} & \Pr \left\{ \left\| \sum_{j=1}^D X_j \right\|_{\mathcal{Y}, \mathcal{Y}} \geq \epsilon \right\} \\ & \leq \begin{cases} 2 \text{IntDim}(V) \exp \left(-\frac{\epsilon^2}{2\|V\|_{\mathcal{Y}, \mathcal{Y}}(1+\frac{1}{p})} \right) r_V(\epsilon), & \epsilon \leq \frac{\|V\|_{\mathcal{Y}, \mathcal{Y}}^{1+1/p}}{2U K(V, p)} \\ 2 \text{IntDim}(V) \exp \left(-\frac{\epsilon}{4U K(V, p)} \right) r_V(\epsilon), & \text{otherwise.} \end{cases} \end{aligned}$$

where $K(V, p) = \log(16\sqrt{2}p) + \log\left(\frac{DU^2}{\|V\|_{\mathcal{Y}, \mathcal{Y}}}\right)$ and $r_V(\epsilon) = 1 + \frac{3}{\epsilon^2 \log^2(1+\epsilon/\|V\|_{\mathcal{Y}, \mathcal{Y}})}$

Let $\psi = \psi_1$. To use [Theorem A.3](#), we set $X_j = F^j(\delta_i)$. We have indeed $\widehat{\mathbf{E}_{\text{Haar}, \rho}}[F^j(\delta_i)] = 0$ since $\tilde{K}(\delta_i)$ is the Monte-Carlo approximation of $K_e(\delta_i)$ and the matrices $F^j(\delta_i)$ are self-adjoint. We assume we can bound all the Orlicz norms of the $F^j(\delta_i) = \frac{1}{D}(\tilde{K}^j(\delta_i) - K_e(\delta_i))$. In the following we use constants u_i such that $u_i = DU$. Using [Lemma A.4](#) and the sub-additivity of the $\|\cdot\|_{\mathcal{Y}, \mathcal{Y}}$ and $\|\cdot\|_{\psi_1}$ norm,

$$\begin{aligned} u_i &= 2D \max_{1 \leq j \leq D} \|\|F^j(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}}\|_{\psi_1} \\ &\leq 2 \max_{1 \leq j \leq D} \|\|\tilde{K}^j(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}}\|_{\psi_1} + 2\|\|K_e(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}}\|_{\psi_1} \\ &< 4 \max_{1 \leq j \leq D} \|\|A(\omega_j)\|_{\mathcal{Y}, \mathcal{Y}}\|_{\psi_1} + 4\|\|K_e(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}}\|_{\psi_1} \\ &= 4 \left(\|\|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}}\|_{\psi_1} + \|K_e(\delta_i)\|_{\mathcal{Y}, \mathcal{Y}} \right) \end{aligned}$$

In the same way we defined the constants $v_i = DV$,

$$\begin{aligned} v_i &= D \sum_{j=1}^D \widehat{\mathbf{E}_{\text{Haar}, \rho}} F^j(\delta_i)^2 \\ &= D \widehat{\mathbf{Var}}_{\text{Haar}, \rho} [\tilde{K}(\delta_i)] \end{aligned}$$

Then applying [Theorem A.3](#), we get for all $i \in \mathbb{N}_{\mathcal{N}(\mathcal{D}_e, r)}^*$ (i is the index of each anchor)

$$\begin{aligned} & \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \|F(\delta_i)\|_{y,y} \geq \epsilon \right\} \\ & \leq \begin{cases} 4 \text{IntDim}(v_i) \exp \left(-D \frac{\epsilon^2}{2\|v_i\|_{y,y}(1+\frac{1}{p})} \right) r_{v_i/D}(\epsilon), & \epsilon \leq \frac{\|v_i\|_{y,y}}{2u_i} \frac{1+1/p}{K(v_i,p)} \\ 4 \text{IntDim}(v_i) \exp \left(-D \frac{\epsilon}{4u_i K(v_i,p)} \right) r_{v_i/D}(\epsilon), & \text{otherwise.} \end{cases} \end{aligned}$$

with

$$K(v_i, p) = \log \left(16\sqrt{2}p \right) + \log \left(\frac{u_i^2}{\|v_i\|_{y,y}} \right)$$

and

$$r_{v_i/D} = 1 + \frac{3}{\epsilon^2 \log^2(1 + D\epsilon/\|v_i\|_{y,y})}.$$

To unify the bound on each anchor we define two constant

$$u = 4 \left(\|\|A(\omega)\|_{y,y}\|_{\psi_1} + \sup_{\delta \in \mathcal{D}_e} \|K_e(\delta)\|_{y,y} \right) \geq \max_{i=1,\dots,T} u_i$$

and

$$v = \sup_{\delta \in \mathcal{D}_e} D\text{Var}_{\widehat{\text{Haar}}, \rho} [\tilde{K}_e(\delta)] \geq \max_{i=1,\dots,T} v_i.$$

A.1.4 Union Bound and examples

Taking the union bound over the anchors yields

$$\begin{aligned} & \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \bigcup_{i=1}^{\mathcal{N}(\mathcal{D}_e, r)} \|F(\delta_i)\|_{y,y} \geq \epsilon \right\} \\ & \leq 4\mathcal{N}(\mathcal{D}_e, r) r_{v/D}(\epsilon) \text{IntDim}(v) \\ & \quad \begin{cases} \exp \left(-D \frac{\epsilon^2}{2\|v\|_{y,y}(1+\frac{1}{p})} \right), & \epsilon \leq \frac{\|v\|_{y,y}}{2u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{4u K(v,p)} \right), & \text{otherwise.} \end{cases} \end{aligned} \tag{A.5}$$

Hence combining [Equation A.4](#) and [Equation A.5](#) gives and summing up the hypothesis yields the following proposition

Proposition A.1 *Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ be a shift-invariant \mathcal{Y} -Mercer kernel, where \mathcal{Y} is a finite dimensional Hilbert space of dimension p and \mathcal{X} a metric space. Moreover, let \mathcal{C} be a compact subset of \mathcal{X} , $A : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$ and $\Pr_{\widehat{\text{Haar}}, \rho}$ a pair such that*

$$\tilde{K}_e = \sum_{j=1}^D \cos(\cdot, \omega_j) A(\omega_j) \approx K_e, \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho} \text{ i.i.d.}$$

Let

$$V(\delta) \succcurlyeq \mathbf{Var}_{\widehat{\text{Haar}}, \rho} \tilde{K}_e(\delta), \quad \text{for all } \delta \in \mathcal{D}_e$$

and H_ω be the Lipschitz constant of the function $h : x \mapsto (x, \omega)$. If the three following constant exists

$$m \geq \int_{\widehat{\mathcal{X}}} H_\omega \|A(\omega)\|_{y,y} d\Pr_{\widehat{\text{Haar}}, \rho} < \infty$$

and

$$u \geq 4 \left(\| \|A(\omega)\|_{y,y}\|_{\psi_1} + \sup_{\delta \in \mathcal{D}_e} \|K_e(\delta)\|_{y,y} \right) < \infty$$

and

$$v \geq \sup_{\delta \in \mathcal{D}_e} D \|V(\delta)\|_{y,y} < \infty.$$

Define $p_{\text{int}} \geq \sup_{\delta \in \mathcal{D}_e} \text{IntDim}(V(\delta))$ then for all $r \in \mathbb{R}_+^*$ and all $\epsilon \in \mathbb{R}_+^*$,

$$\begin{aligned} & \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \| \tilde{K} - K \|_{\mathcal{C} \times \mathcal{C}} \geq \epsilon \right\} \\ & \leq 4 \left(\frac{rm}{\epsilon} + p_{\text{int}} N(\mathcal{D}_e, r) r_{v/D}(\epsilon) \right. \\ & \quad \left. \begin{cases} \exp \left(-D \frac{\epsilon^2}{8v(1+\frac{1}{p})} \right), & \epsilon \leq \frac{v}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8uK(v,p)} \right), & \text{otherwise.} \end{cases} \right) \end{aligned}$$

where

$$K(v, p) = \log \left(16\sqrt{2}p \right) + \log \left(\frac{u^2}{\|v\|_{y,y}} \right)$$

and

$$r_{v/D}(\epsilon) = 1 + \frac{3}{\epsilon^2 \log^2(1 + D\epsilon/\|v\|_{y,y})}.$$

Proof Let $m = \int_{\widehat{\mathcal{X}}} H_\omega \|A(\omega)\|_{y,y} d\Pr_{\widehat{\text{Haar}}, \rho}$. From [Lemma A.2](#),

$$\Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid L_F \geq \frac{\epsilon}{2r} \right\} \leq \frac{4rm}{\epsilon}.$$

Thus from [Lemma A.1](#), for all $r \in \mathbb{R}_+^*$,

$$\begin{aligned} & \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \sup_{\delta \in \mathcal{D}_e} \|F(\delta)\|_{y,y} \geq \epsilon \right\} \\ & \leq \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid L_F \geq \frac{\epsilon}{2r} \right\} \\ & \quad + \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \bigcup_{i=1}^{N(\mathcal{D}_e, r)} \|F(\delta_i)\|_{y,y} \geq \epsilon \right\} \\ & = 4 \frac{rm}{\epsilon} + 4N(\mathcal{D}_e, r) r_{v/D}(\epsilon) \text{IntDim}(v) \\ & \quad \begin{cases} \exp \left(-D \frac{\epsilon^2}{8\|v\|_{y,y}(1+\frac{1}{p})} \right), & \epsilon \leq \frac{\|v\|_{y,y}}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8uK(v,p)} \right), & \text{otherwise.} \end{cases} \end{aligned}$$

With slight modification we can obtain a second inequality for the case where the random operators $A(\omega_j)$ are bounded almost surely. This second bound with more restrictions on A has the advantage of working in infinite dimension as long as $A(\omega_j)$ is a Hilbert-Schmidt operator.

Proposition A.2 *Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ be a shift-invariant \mathcal{Y} -Mercer kernel, where \mathcal{Y} is a Hilbert space and \mathcal{X} a metric space. Moreover, let \mathcal{C} be a compact subset of \mathcal{X} , $A : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$ and $\Pr_{\widehat{\text{Haar}}, \rho}$ a pair such that*

$$\tilde{K}_e = \sum_{j=1}^D \cos(\cdot, \omega_j) A(\omega_j) \approx K_e, \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho} \text{ i.i.d.}$$

where $A(\omega_j)$ is a Hilbert-Schmidt operator for all $j \in \mathbb{N}_D^*$. Let $\mathcal{D}_e = \mathcal{C} * \mathcal{C}^{-1}$ and

$$V(\delta) \succcurlyeq \text{Var}_{\widehat{\text{Haar}}, \rho} \tilde{K}_e(\delta), \quad \text{for all } \delta \in \mathcal{D}_e$$

and H_ω be the Lipschitz constant of the function $h : x \mapsto (x, \omega)$. If the three following constant exists

$$m \geq \int_{\widehat{\mathcal{X}}} H_\omega \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} d\Pr_{\widehat{\text{Haar}}, \rho} < \infty$$

and

$$u \geq \text{ess sup}_{\omega \in \widehat{\mathcal{X}}} \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} + \sup_{\delta \in \mathcal{D}_e} \|K_e(\delta)\|_{\mathcal{Y}, \mathcal{Y}} < \infty$$

and

$$v \geq \sup_{\delta \in \mathcal{D}_e} D \|V(\delta)\|_{\mathcal{Y}, \mathcal{Y}} < \infty.$$

define $p_{\text{int}} \geq \sup_{\delta \in \mathcal{D}_e} \text{IntDim}(V(\delta))$ then for all $r \in \mathbb{R}_+^*$ and all $\epsilon > \sqrt{\frac{v}{D}} + \frac{1}{3D} u$,

$$\begin{aligned} & \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \sup_{\delta \in \mathcal{D}_e} \|F(\delta)\|_{\mathcal{Y}, \mathcal{Y}} \geq \epsilon \right\} \\ & \leq 4 \left(\frac{rm}{\epsilon} + p_{\text{int}} N(\mathcal{D}_e, r) \exp(-D \psi_{v, u}(\epsilon)) \right) \end{aligned}$$

where $\psi_{v, u}(\epsilon) = \frac{\epsilon^2}{2(v+u\epsilon/3)}$.

When the covering number $N(\mathcal{D}_e, r)$ of the metric space \mathcal{D}_e has an analytical form, it is possible to optimize the bound over the radius r of the covering balls. As an example, we refine [Proposition A.1](#) and [Proposition A.2](#) in the case where \mathcal{C} is a finite dimensional Banach space.

Corollary A.1 Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ be a shift-invariant \mathcal{Y} -Mercer kernel, where \mathcal{Y} is a finite dimensional Hilbert space of dimension p and \mathcal{X} a finite dimensional Banach space of dimension d . Moreover, let \mathcal{C} be a closed ball of \mathcal{X} centered at the origin of diameter $|\mathcal{C}|$, $A : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$ and $\Pr_{\widehat{\text{Haar}}, \rho}$ a pair such that

$$\tilde{K}_e = \sum_{j=1}^D \cos(\cdot, \omega_j) A(\omega_j) \approx K_e, \quad \omega_j \sim \Pr_{\widehat{\text{Haar}}, \rho} \text{ i.i.d.}$$

Let $\mathcal{D}_{\mathcal{C}} = \mathcal{C} \star \mathcal{C}^{-1}$ and

$$V(\delta) \succcurlyeq \text{Var}_{\widehat{\text{Haar}}, \rho} \tilde{K}_e(\delta), \quad \text{for all } \delta \in \mathcal{D}_{\mathcal{C}}$$

Let H_ω be the Lipschitz constant of $h_\omega : x \mapsto (x, \omega)$. If the three following constant exists

$$m \geq \int_{\widehat{\mathcal{X}}} H_\omega \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} d\Pr_{\widehat{\text{Haar}}, \rho} < \infty$$

and

$$u \geq 4 \left(\|\|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}}\|_{\psi_1} + \sup_{\delta \in \mathcal{D}_{\mathcal{C}}} \|K_e(\delta)\|_{\mathcal{Y}, \mathcal{Y}} \right) < \infty$$

and

$$v \geq \sup_{\delta \in \mathcal{D}_{\mathcal{C}}} D \|V(\delta)\|_{\mathcal{Y}, \mathcal{Y}} < \infty.$$

Define $p_{\text{int}} \geq \sup_{\delta \in \mathcal{D}_{\mathcal{C}}} \text{IntDim}(V(\delta))$, then for all $0 < \epsilon \leq m|C|$,

$$\begin{aligned} & \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \|\tilde{K} - K\|_{\mathcal{C} \times \mathcal{C}} \geq \epsilon \right\} \\ & \leq 8\sqrt{2} \left(\frac{m|C|}{\epsilon} \right) (p_{\text{int}} r_{v/D}(\epsilon))^{\frac{1}{d+1}} \begin{cases} \exp \left(-D \frac{\epsilon^2}{8v(d+1)(1+\frac{1}{p})} \right), & \epsilon \leq \frac{v}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8u(d+1)K(v,p)} \right), & \text{otherwise,} \end{cases} \end{aligned}$$

where $K(v, p) = \log(16\sqrt{2}p) + \log\left(\frac{u^2}{v}\right)$ and $r_{v/D}(\epsilon) = 1 + \frac{3}{\epsilon^2 \log^2(1+D\epsilon/v)}$.

Proof As we have seen in Appendix A.1.1, suppose that \mathcal{X} is a finite dimensional Banach space. Let $\mathcal{C} \subset \mathcal{X}$ be a closed ball centered at the origin of diameter $|\mathcal{C}| = C$ then the difference ball centered at the origin

$$\begin{aligned} \mathcal{D}_{\mathcal{C}} &= \mathcal{C} \star \mathcal{C}^{-1} \\ &= \left\{ x \star z^{-1} \mid \|x\|_{\mathcal{X}} \leq C/2, \|z\|_{\mathcal{X}} \leq C/2, (x, z) \in \mathcal{X}^2 \right\} \subset \mathcal{X} \end{aligned}$$

is closed and bounded, so compact and has diameter $|C| = 2C$. It is possible to cover it with

$$\mathcal{N}(\mathcal{D}_{\mathcal{C}}, r) = \left(\frac{2|C|}{r} \right)^d$$

closed balls of radius r . Plugging back into Equation A.5 yields

$$\begin{aligned} & \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \|\tilde{K} - K\|_{C \times C} \geq \epsilon \right\} \\ & \leq 4 \left(\frac{rm}{\epsilon} + p_{\text{int}} \left(\frac{2|C|}{r} \right)^d r_{v/D}(\epsilon) \right. \\ & \quad \left. \begin{cases} \exp \left(-D \frac{\epsilon^2}{8v(1+\frac{1}{p})} \right), & \epsilon \leq \frac{v}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8uK(v,p)} \right), & \text{otherwise.} \end{cases} \right) \end{aligned}$$

The right hand side of the equation has the form $ar + br^{-d}$ with

$$a = \frac{m}{\epsilon}$$

and

$$b = p_{\text{int}}(2|C|)^d r_{v/D}(\epsilon) \begin{cases} \exp \left(-D \frac{\epsilon^2}{8v(1+\frac{1}{p})} \right), & \epsilon \leq \frac{v}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8uK(v,p)} \right), & \text{otherwise.} \end{cases}$$

Following [118, 139, 167], we optimize over r . It is a convex continuous function on \mathbb{R}_+ and achieve minimum at

$$r = \left(\frac{bd}{a} \right)^{\frac{1}{d+1}}$$

and the minimum value is

$$r_* = a^{\frac{d}{d+1}} b^{\frac{1}{d+1}} \left(d^{\frac{1}{d+1}} + d^{-\frac{d}{d+1}} \right),$$

hence

$$\begin{aligned} & \Pr_{\widehat{\text{Haar}}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \|\tilde{K} - K\|_{C \times C} \geq \epsilon \right\} \\ & \leq C_d \left(\frac{2m|C|}{\epsilon} \right)^{\frac{d}{d+1}} (p_{\text{int}} r_{v/D}(\epsilon))^{\frac{1}{d+1}} \begin{cases} \exp \left(-D \frac{\epsilon^2}{8v(d+1)(1+\frac{1}{p})} \right), & \epsilon \leq \frac{v}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8u(d+1)K(v,p)} \right), & \text{otherwise,} \end{cases} \\ & \leq 8\sqrt{2} \left(\frac{m|C|}{\epsilon} \right) (p_{\text{int}} r_{v/D}(\epsilon))^{\frac{1}{d+1}} \begin{cases} \exp \left(-D \frac{\epsilon^2}{8v(d+1)(1+\frac{1}{p})} \right), & \epsilon \leq \frac{v}{u} \frac{1+1/p}{K(v,p)} \\ \exp \left(-D \frac{\epsilon}{8u(d+1)K(v,p)} \right), & \text{otherwise,} \end{cases} \end{aligned}$$

where $C_d = 4 \left(d^{\frac{1}{d+1}} + d^{-\frac{d}{d+1}} \right)$. Eventually when X is a Banach space, the Lipschitz constant of h_ω is the supremum of the gradient

$$H_\omega = \sup_{\delta \in \mathcal{D}_e} \|(\nabla h_\omega)(\delta)\|_{\widehat{X}}.$$

Following the same proof technique we obtain the second bound for bounded ORFF.

Corollary A.2 Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y})$ be a shift-invariant \mathcal{Y} -Mercer kernel, where \mathcal{Y} is a Hilbert space and \mathcal{X} a finite dimensional Banach space of dimension D . Moreover, let \mathcal{C} be a closed ball of \mathcal{X} centered at the origin of diameter $|\mathcal{C}|$, subset of \mathcal{X} , $A : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$ and $\widehat{\text{Pr}}_{\text{Haar}, \rho}$ a pair such that

$$\tilde{K}_e = \sum_{j=1}^D \cos(\cdot, \omega_j) A(\omega_j) \approx K_e, \quad \omega_j \sim \widehat{\text{Pr}}_{\text{Haar}, \rho} \text{ i.i.d.}$$

where $A(\omega_j)$ is a Hilbert-Schmidt operator for all $j \in \mathbb{N}_D^*$. Let $\mathcal{D}_e = \mathcal{C} * \mathcal{C}^{-1}$ and

$$V(\delta) \succcurlyeq \widehat{\text{Var}}_{\text{Haar}, \rho} \tilde{K}_e(\delta), \quad \text{for all } \delta \in \mathcal{D}_e$$

and H_ω be the Lipschitz constant of the function $h : x \mapsto (x, \omega)$. If the three following constant exists

$$m \geq \int_{\widehat{\mathcal{X}}} H_\omega \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} d\widehat{\text{Pr}}_{\text{Haar}, \rho} < \infty$$

and

$$u \geq \text{ess sup}_{\omega \in \widehat{\mathcal{X}}} \|A(\omega)\|_{\mathcal{Y}, \mathcal{Y}} + \sup_{\delta \in \mathcal{D}_e} \|K_e(\delta)\|_{\mathcal{Y}, \mathcal{Y}} < \infty$$

and

$$v \geq \sup_{\delta \in \mathcal{D}_e} D \|V(\delta)\|_{\mathcal{Y}, \mathcal{Y}} < \infty.$$

define $p_{\text{int}} \geq \sup_{\delta \in \mathcal{D}_e} \text{IntDim}(V(\delta))$ then for all $\sqrt{\frac{v}{D}} + \frac{u}{3D} < \epsilon < m|\mathcal{C}|$,

$$\begin{aligned} & \widehat{\text{Pr}}_{\text{Haar}, \rho} \left\{ (\omega_j)_{j=1}^D \mid \sup_{\delta \in \mathcal{D}_e} \|F(\delta)\|_{\mathcal{Y}, \mathcal{Y}} \geq \epsilon \right\} \\ & \leq 8\sqrt{2} \left(\frac{m|\mathcal{C}|}{\epsilon} \right) p_{\text{int}}^{\frac{1}{d+1}} \exp(-D\psi_{v, d, u}(\epsilon)) \end{aligned}$$

where $\psi_{v, d, u}(\epsilon) = \frac{\epsilon^2}{2(d+1)(v+u\epsilon/3)}$.

A.2 PROOF OF THE ORFF ESTIMATOR VARIANCE BOUND

We use the notations $\delta = x * z^{-1}$ for all $x, z \in \mathcal{X}$, $\tilde{K}(x, z) = \tilde{\Phi}(x)^* \tilde{\Phi}(z)$, $\tilde{K}^j(x, z) = \Phi_x(\omega_j)^* \Phi_z(\omega_j)$ and $K_e(\delta) = K_e(x, z)$.

Proposition A.3 (Bounding the variance of \tilde{K}). Let K be a shift invariant \mathcal{Y} -Mercer kernel on a second countable LCA topological space \mathcal{X} . Let $A : \widehat{\mathcal{X}} \rightarrow \mathcal{L}(\mathcal{Y})$ and $\widehat{\text{Pr}}_{\text{Haar}, \rho}$ a pair such that

$$\tilde{K}_e = \sum_{j=1}^D \cos(\cdot, \omega_j) A(\omega_j) \approx K_e, \quad \omega_j \sim \widehat{\text{Pr}}_{\text{Haar}, \rho} \text{ i.i.d.}$$

Then,

$$\begin{aligned} \widehat{\text{Var}}_{\text{Haar}, \rho} [\tilde{K}_e(\delta)] & \leq \frac{1}{2D} \left((K_e(2\delta) + K_e(\delta)) E_{\widehat{\text{Pr}}_{\text{Haar}, \rho}} [A(\omega)] \right. \\ & \quad \left. - 2K_e(\delta)^2 + \widehat{\text{Var}}_{\text{Haar}, \rho} [A(\omega)] \right) \end{aligned}$$

Proof Let $\delta \in \mathcal{D}_e$ be a constant. From the definition of the variance of a random variable and using the fact that the $(\omega_j)_{j=1}^D$ are i.i.d. random variables,

$$\begin{aligned}\text{Var}_{\widehat{\text{Haar}}, \rho} [\tilde{K}_e(\delta)] &= E_{\widehat{\text{Haar}}, \rho} \left[\frac{1}{D} \sum_{j=1}^D \tilde{K}_e^j(\delta) - K_e(\delta) \right]^2 \\ &= \frac{1}{D^2} E_{\widehat{\text{Haar}}, \rho} \left[\sum_{j=1}^D \tilde{K}_e^j(\delta) - K_e(\delta) \right]^2 \\ &= \frac{1}{D} E_{\widehat{\text{Haar}}, \rho} [\tilde{K}_e^j(\delta)^2 - \tilde{K}_e^j(\delta)K_e(\delta) - K_e(\delta)\tilde{K}_e^j(\delta) \\ &\quad + K_e(\delta)^2]\end{aligned}$$

From the definition of \tilde{K}_e^j , $E_{\widehat{\text{Haar}}, \rho} \tilde{K}_e^j(\delta) = K_e(\delta)$, which leads to

$$\text{Var}_{\widehat{\text{Haar}}, \rho} [\tilde{K}_e(\delta)] = \frac{1}{D} E_{\widehat{\text{Haar}}, \rho} [\tilde{K}_e^j(\delta)^2 - K_e(\delta)^2]$$

A trigonometric identity gives us $(\cos(\delta, \omega))^2 = \frac{1}{2} (\cos(2\delta, \omega) + \cos(e, \omega))$. Thus

$$\begin{aligned}\text{Var}_{\widehat{\text{Haar}}, \rho} [\tilde{K}_e(\delta)] &= \frac{1}{2D} E_{\widehat{\text{Haar}}, \rho} [(\cos(2\delta, \omega) + \cos(e, \omega)) A(\omega)^2 \\ &\quad - 2K_e(\delta)^2].\end{aligned}$$

Also,

$$\begin{aligned}E_{\widehat{\text{Haar}}, \rho} [\cos(2\delta, \omega)A(\omega)^2] &= E_{\widehat{\text{Haar}}, \rho} [\cos(2\delta, \omega)A(\omega)] E_{\widehat{\text{Haar}}, \rho} [A(\omega)] \\ &\quad + \text{Cov}_{\widehat{\text{Haar}}, \rho} [\cos(2\delta, \omega)A(\omega), A(\omega)] \\ &= K_e(2\delta) E_{\widehat{\text{Haar}}, \rho} [A(\omega)] \\ &\quad + \text{Cov}_{\widehat{\text{Haar}}, \rho} [\cos(2\delta, \omega)A(\omega), A(\omega)]\end{aligned}$$

Similarly we obtain

$$\begin{aligned}E_{\widehat{\text{Haar}}, \rho} [\cos(e, \omega)A(\omega)^2] &= K_e(e) E_{\widehat{\text{Haar}}, \rho} [A(\omega)] \\ &\quad + \text{Cov}_{\widehat{\text{Haar}}, \rho} [\cos(e, \omega)A(\omega), A(\omega)]\end{aligned}$$

Therefore

$$\begin{aligned}\text{Var}_{\widehat{\text{Haar}}, \rho} [\tilde{K}_e(\delta)] &= \frac{1}{2D} \left((K_e(2\delta) + K_e(e)) E_{\widehat{\text{Haar}}, \rho} [A(\omega)] - 2K_e(\delta)^2 \right. \\ &\quad \left. + \text{Cov}_{\widehat{\text{Haar}}, \rho} [(\cos(2\delta, \omega) + \cos(e, \omega)) A(\omega), A(\omega)] \right) \\ &= \frac{1}{2D} \left((K_e(2\delta) + K_e(e)) E_{\widehat{\text{Haar}}, \rho} [A(\omega)] - 2K_e(\delta)^2 \right. \\ &\quad \left. + \text{Cov}_{\widehat{\text{Haar}}, \rho} [(\cos(\delta, \omega))^2 A(\omega), A(\omega)] \right) \\ &\leq \frac{1}{2D} \left((K_e(2\delta) + K_e(e)) E_{\widehat{\text{Haar}}, \rho} [A(\omega)] - 2K_e(\delta)^2 \right. \\ &\quad \left. + \text{Var}_{\widehat{\text{Haar}}, \rho} [A(\omega)] \right)\end{aligned}$$



B

MISCELLANEOUS

In this appendix chapter we present how to use ORFF framework in the context of semi-supervised learning. We use the setting of Minh, Bazzani, and Murino [119, 120]

Contents

B.1	Learning with semi-supervision	184
B.1.1	Representer theorem and feature equivalence	184
B.1.2	Gradients	189
B.1.3	Complexity	191

B.1 LEARNING WITH SEMI-SUPERVISION

We present here a direct extension of [Chapter 6](#) to semi-supervised learning with ORFF. Semi-supervised learning in vv-RKHS has been first presented at the same time by Brouard, d'Alché-Buc, and Szafranski [34] and Minh and Sindhwani [117], and then more deeply developed in Brouard, d'Alché-Buc, and Szafranski [35] and Minh, Bazzani, and Murino [120]. We have chosen to adopt here the presentation of Minh, Bazzani, and Murino [120] slightly more general, encompassing Vector-valued Manifold Regularization [19, 34, 119] and Co-regularized Multi-view Learning [31, 143, 159, 166].

B.1.1 Representer theorem and feature equivalence

We suppose that we are given a training sample $\mathbf{u} = (x_i)_{i=1}^{N+U} \in \mathcal{X}^U$ of unlabeled examples. We note $\mathbf{z} \in (\mathcal{X} \times \mathcal{Y})^N \times \mathcal{X}^U$ the sequence $\mathbf{z} = \mathbf{s}\mathbf{u}$ concatenating both labeled (\mathbf{s}) and unlabeled (\mathbf{u}) training examples.

Theorem B.1 (Representer theorem, Minh, Bazzani, and Murino [120]). *Let K be a \mathcal{U} -Mercer Operator-Valued Kernel and \mathcal{H}_K its corresponding \mathcal{U} -Reproducing Kernel Hilbert space.*

Let $V : \mathcal{U} \rightarrow \mathcal{Y}$ be a bounded linear operator and let $c : \mathcal{Y} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ be a cost function such that $L(x, f, y) = c(Vf(x), y)$ is a proper convex lower semi-continuous function in f for all $x \in \mathcal{X}$ and all $y \in \mathcal{Y}$.

Eventually let $\lambda_K \in \mathbb{R}_{>0}$ and $\lambda_M \in \mathbb{R}_+$ be two regularization hyperparameters and $(M_{ik})_{i,k=1}^{N+U}$ be a sequence of data dependent bounded linear operators in $\mathcal{L}(\mathcal{U})$, such that

$$\sum_{i,j=1}^{N+U} \langle u_i, M_{ik} u_k \rangle \geq 0, \quad \forall (u_i)_{i=1}^{N+U} \in \mathcal{U}^{N+U} \text{ and } M_{ik} = M_{ki}^*.$$

The solution $f_z \in \mathcal{H}_K$ of the regularized optimization problem

$$\begin{aligned} f_z = \arg \min_{f \in \mathcal{H}_K} & \frac{1}{N} \sum_{i=1}^N c(Vf(x_i), y_i) + \frac{\lambda_K}{2} \|f\|_K^2 \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \langle f(x_i), M_{ik} f(x_k) \rangle_U \end{aligned} \tag{B.1}$$

has the form $f_z = \sum_{j=1}^{N+U} K(\cdot, x_j) u_{z,j}$ where $u_{z,j} \in \mathcal{U}$ and

$$\begin{aligned} u_z = \arg \min_{u \in \bigoplus_{i=1}^{N+U} \mathcal{U}} & \frac{1}{N} \sum_{i=1}^N c \left(V \sum_{k=1}^{N+U} K(x_i, x_k) u_k, y_i \right) \\ & + \frac{\lambda_K}{2} \sum_{k=1}^{N+U} u_k^* K(x_i, x_k) u_k \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \left\langle \sum_{j=1}^{N+U} K(x_i, x_j) u_j, M_{ik} \sum_{j=1}^{N+U} K(x_k, x_j) u_j \right\rangle_u . \end{aligned} \quad (B.2)$$

We present here the proof of the formulation proposed by Minh, Bazzaani, and Murino [120]. In the mean time we clarify some elements of the proof. Indeed the existence of a global minimizer is not trivial and we must invoke the Mazur-Schauder theorem. Moreover the coercivity of the objective function required by the Mazur-Schauder theorem is not obvious when we do not require the cost function to take only positive values. However a corollary of Hahn-Banach theorem linking strong convexity to coercivity gives the solution.

Proof Since $f(x) = K_x^* f$ (see [Equation 3.14](#)), the optimization problem reads

$$\begin{aligned} f_z = \arg \min_{f \in \mathcal{H}_K} & \frac{1}{N} \sum_{i=1}^N c(VK_{x_i}^* f, y_i) + \frac{\lambda_K}{2} \|f\|_K^2 \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \langle K_{x_i}^* f, M_{ik} K_{x_k}^* f \rangle_u \end{aligned}$$

Let $W_{V,s} : \mathcal{H}_K \rightarrow \bigoplus_{i=1}^N \mathcal{Y}$ be the restriction linear operator defined as

$$W_{V,s} f = \bigoplus_{i=1}^N VK_{x_i}^* f,$$

with $VK_{x_i}^* : \mathcal{H}_K \rightarrow \mathcal{Y}$ and $K_{x_i} V^* : \mathcal{Y} \rightarrow \mathcal{H}_K$. Let $Y = \bigoplus_{i=1}^N y_i \in \mathcal{Y}^N$. We have

$$\langle Y, W_{V,s} f \rangle_{\bigoplus_{i=1}^N \mathcal{Y}} = \sum_{i=1}^N \langle y_i, VK_{x_i}^* f \rangle_{\mathcal{Y}} = \sum_{i=1}^N \langle K_{x_i} V^* y_i, f \rangle_{\mathcal{H}_K}.$$

Thus the adjoint operator $W_{V,s}^* : \bigoplus_{i=1}^N \mathcal{Y} \rightarrow \mathcal{H}_K$ is

$$W_{V,s}^* Y = \sum_{i=1}^N K_{x_i} V^* y_i,$$

and the operator $W_{V,s}^* W_{V,s} : \mathcal{H}_K \rightarrow \mathcal{H}_K$ is

$$W_{V,s}^* W_{V,s} f = \sum_{i=1}^N K_{x_i} V^* VK_{x_i}^* f$$

where $V^*V \in \mathcal{L}(\mathcal{U})$. Let

$$\begin{aligned} J_{\lambda_K}(f) &= \underbrace{\frac{1}{N} \sum_{i=1}^N c(Vf(x_i), y_i) + \frac{\lambda_K}{2} \|f\|_K^2}_{=J_c} \\ &\quad + \underbrace{\frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \langle f(x_i), M_{ik}f(x_k) \rangle_U}_{=J_M} \end{aligned}$$

Since c is proper, lower semi-continuous and convex by assumption, thus the term J_c is also proper, lower semi-continuous and convex. Moreover the term J_M is always positive for any $f \in \mathcal{H}_K$ and $\frac{\lambda_K}{2} \|f\|_K^2$ is strongly convex. Thus J_{λ_K} is strongly convex. Apply [Lemma 6.1](#) to obtain the coercivity of J_{λ_K} , and then [Theorem 6.1](#) to show that J_{λ_K} has a unique minimizer and is attained. Then let

$$\mathcal{H}_{K,z} = \left\{ \sum_{j=1}^{N+U} K_{x_j} u_j \mid \forall (u_i)_{i=1}^{N+U} \in \mathcal{U}^{N+U} \right\}.$$

For $f \in \mathcal{H}_{K,z}^\perp$, the operator $W_{V,s}$ satisfies

$$\langle Y, W_{V,s} f \rangle_{\bigoplus_{i=1}^N \mathcal{Y}} = \langle \underbrace{f}_{\in \mathcal{H}_{K,z}^\perp}, \underbrace{\sum_{i=1}^{N+U} K_{x_i} V^* y_i}_{\in \mathcal{H}_{K,z}} \rangle_{\mathcal{H}_K} = 0$$

for all sequences $(y_i)_{i=1}^N$, since $V^* y_i \in \mathcal{U}$. Hence,

$$(Vf(x_i))_{i=1}^N = 0 \tag{B.3}$$

In the same way,

$$\sum_{i=1}^{N+U} \langle K_{x_i}^* f, u_i \rangle_{\mathcal{U}} = \langle \underbrace{f}_{\in \mathcal{H}_{K,z}^\perp}, \underbrace{\sum_{j=1}^{N+U} K_{x_j} u_j}_{\in \mathcal{H}_{K,z}} \rangle_{\mathcal{H}_K} = 0.$$

for all sequences $(u_i)_{i=1}^{N+U} \in \mathcal{U}^{N+U}$. As a result,

$$(f(x_i))_{i=1}^{U+N} = 0. \tag{B.4}$$

Now for an arbitrary $f \in \mathcal{H}_K$, consider the orthogonal decomposition $f = f^\perp + f^\parallel$, where $f^\perp \in \mathcal{H}_{K,z}^\perp$ and $f^\parallel \in \mathcal{H}_{K,z}$. Then since $\|f^\perp + f^\parallel\|_{\mathcal{H}_K}^2 = \|f^\perp\|_{\mathcal{H}_K}^2 + \|f^\parallel\|_{\mathcal{H}_K}^2$, [Equation B.3](#) and [Equation B.4](#) shows that if $\lambda_K > 0$, clearly then

$$J_{\lambda_K}(f) = J_{\lambda_K}(f^\perp + f^\parallel) \geq J_{\lambda_K}(f^\parallel)$$

The last inequality holds only when $\|f^\perp\|_{\mathcal{H}_K} = 0$, that is when $f^\perp = 0$. As a result since the minimizer of J_{λ_K} is unique and attained, it must lies in $\mathcal{H}_{K,z}$. \square

Theorem B.2 (Feature equivalence). Let \tilde{K} be an Operator-Valued Kernel such that for all $x, z \in \mathcal{X}$, $\tilde{\Phi}(x)^* \tilde{\Phi}(z) = \tilde{K}(x, z)$ where \tilde{K} is a \mathcal{U} -Mercer OVK and $\mathcal{H}_{\tilde{K}}$ its corresponding \mathcal{U} -Reproducing kernel Hilbert space.

Let $V : \mathcal{U} \rightarrow \mathcal{Y}$ be a bounded linear operator and let $c : \mathcal{Y} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$ be a cost function such that $L(x, \tilde{f}, y) = c(V\tilde{f}(x), y)$ is a proper convex lower semi-continuous function in $\tilde{f} \in \mathcal{H}_{\tilde{K}}$ for all $x \in \mathcal{X}$ and all $y \in \mathcal{Y}$.

Eventually let $\lambda_K \in \mathbb{R}_{>0}$ and $\lambda_M \in \mathbb{R}_+$ be two regularization hyperparameters and $(M_{ik})_{i,k=1}^{N+U}$ be a sequence of data dependent bounded linear operators in $\mathcal{L}(\mathcal{U})$, such that

$$\sum_{i,j=1}^{N+U} \langle u_i, M_{ik} u_k \rangle \geq 0, \quad \forall (u_i)_{i=1}^{N+U} \in \mathcal{U}^{N+U} \text{ and } M_{ik} = M_{ki}^*.$$

The solution $\tilde{f}_z \in \mathcal{H}_{\tilde{K}}$ of the regularized optimization problem

$$\begin{aligned} \tilde{f}_z &= \arg \min_{\tilde{f} \in \mathcal{H}_{\tilde{K}}} \frac{1}{N} \sum_{i=1}^N c(V\tilde{f}(x_i), y_i) + \frac{\lambda_K}{2} \|\tilde{f}\|_{\tilde{K}}^2 \\ &\quad + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \langle \tilde{f}(x_i), M_{ik} \tilde{f}(x_k) \rangle_U \end{aligned} \tag{B.5}$$

has the form $\tilde{f}_z = \tilde{\Phi}(\cdot)^* \theta_z$, where $\theta_z \in (\text{Ker } \tilde{W})^\perp$ and

$$\begin{aligned} \theta_z &= \arg \min_{\theta \in \tilde{\mathcal{H}}} \frac{1}{N} \sum_{i=1}^N c(V\tilde{\Phi}(x_i)^* \theta, y_i) + \frac{\lambda_K}{2} \|\theta\|_{\tilde{\mathcal{H}}}^2 \\ &\quad + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \langle \theta, \tilde{\Phi}(x_i) M_{ik} \tilde{\Phi}(x_k)^* \theta \rangle_{\tilde{\mathcal{H}}}. \end{aligned} \tag{B.6}$$

Proof Since \tilde{K} is an operator-valued kernel, from [Theorem B.1](#), [Equation B.5](#) has a solution of the form

$$\begin{aligned} \tilde{f}_z &= \sum_{i=1}^{N+U} \tilde{K}(\cdot, x_i) u_i, \quad u_i \in \mathcal{U}, x_i \in \mathcal{X} \\ &= \sum_{i=1}^N \tilde{\Phi}(\cdot)^* \tilde{\Phi}(x_i) u_i = \tilde{\Phi}(\cdot)^* \underbrace{\left(\sum_{i=1}^{N+U} \tilde{\Phi}(x_i) u_i \right)}_{=\theta \in (\text{Ker } \tilde{W})^\perp \subset \tilde{\mathcal{H}}}. \end{aligned}$$

Let

$$\begin{aligned} \theta_z &= \arg \min_{\theta \in (\text{Ker } \tilde{W})^\perp} \frac{1}{N} \sum_{i=1}^N c(V\tilde{\Phi}(x_i)^* \theta, y_i) + \frac{\lambda_K}{2} \|\tilde{\Phi}(\cdot)^* \theta\|_{\tilde{K}}^2 \\ &\quad + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \langle \tilde{\Phi}(x_i)^* \theta, M_{ik} \tilde{\Phi}(x_k)^* \theta \rangle_U. \end{aligned}$$

Since $\theta \in (\text{Ker } \widetilde{W})^\perp$ and W is an isometry from $(\text{Ker } \widetilde{W})^\perp \subset \widetilde{\mathcal{H}}$ onto $\mathcal{H}_{\widetilde{K}}$, we have $\left\| \widetilde{\Phi}(\cdot)^* \theta \right\|_{\widetilde{\mathcal{H}}}^2 = \|\theta\|_{\widetilde{\mathcal{H}}}^2$. Hence

$$\begin{aligned}\theta_z = & \arg \min_{\theta \in (\text{Ker } \widetilde{W})^\perp} \frac{1}{N} \sum_{i=1}^N c \left(V \widetilde{\Phi}(x_i)^* \theta, y_i \right) + \frac{\lambda_K}{2} \|\theta\|_{\widetilde{\mathcal{H}}}^2 \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \langle \widetilde{\Phi}(x_i)^* \theta, M_{ik} \widetilde{\Phi}(x_k)^* \theta \rangle_U.\end{aligned}$$

Finding a minimizer θ_z over $(\text{Ker } \widetilde{W})^\perp$ is not the same as finding a minimizer over $\widetilde{\mathcal{H}}$. Although in both cases Mazur-Schauder's theorem guarantees that the respective minimizers are unique, they might not be the same. Since \widetilde{W} is bounded, $\text{Ker } \widetilde{W}$ is closed, so that we can perform the decomposition $\widetilde{\mathcal{H}} = (\text{Ker } \widetilde{W})^\perp \oplus (\text{Ker } \widetilde{W})$. Then clearly by linearity of W and the fact that for all $\theta^\parallel \in \text{Ker } \widetilde{W}$, $\widetilde{W}\theta^\parallel = 0$, if $\lambda > 0$ we have

$$\begin{aligned}\theta_z = & \arg \min_{\theta \in \widetilde{\mathcal{H}}} \frac{1}{N} \sum_{i=1}^N c \left(V \widetilde{\Phi}(x_i)^* \theta, y_i \right) + \frac{\lambda_K}{2} \|\theta\|_{\widetilde{\mathcal{H}}}^2 \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \langle \widetilde{\Phi}(x_i)^* \theta, M_{ik} \widetilde{\Phi}(x_k)^* \theta \rangle_U\end{aligned}$$

Thus

$$\begin{aligned}\theta_z = & \arg \min_{\substack{\theta^\perp \in (\text{Ker } \widetilde{W})^\perp, \\ \theta^\parallel \in \text{Ker } \widetilde{W}}} \frac{1}{N} \sum_{i=1}^N c \left(V \left(\widetilde{W}\theta^\perp \right) (x_i) + \underbrace{V \left(\widetilde{W}\theta^\parallel \right) (x_i)}_{=0 \text{ for all } \theta^\parallel}, y_i \right) \\ & + \frac{\lambda_K}{2} \|\theta^\perp\|_{\widetilde{\mathcal{H}}}^2 + \underbrace{\frac{\lambda_K}{2} \|\theta^\parallel\|_{\widetilde{\mathcal{H}}}^2}_{=0 \text{ only if } \theta^\parallel=0} \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \langle \widetilde{\Phi}(x_i)^* \theta^\perp, M_{ik} \left(\widetilde{W}\theta^\perp \right) (x_k) \rangle_U \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \underbrace{\left\langle \left(\widetilde{W}\theta^\parallel \right) (x_i), M_{ik} \left(\widetilde{W}\theta^\perp \right) (x_k) \right\rangle_U}_{=0 \text{ for all } \theta^\parallel} \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \underbrace{\left\langle \left(\widetilde{W}\theta^\perp \right) (x_i), M_{ik} \underbrace{\left(\widetilde{W}\theta^\parallel \right) (x_k)}_{=0 \text{ for all } \theta^\parallel} \right\rangle_U}_{=0 \text{ for all } \theta^\parallel} \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \underbrace{\left\langle \left(\widetilde{W}\theta^\parallel \right) (x_i), M_{ik} \underbrace{\left(\widetilde{W}\theta^\parallel \right) (x_k)}_{=0 \text{ for all } \theta^\parallel} \right\rangle_U}.\end{aligned}$$

Thus

$$\begin{aligned}\theta_z = \arg \min_{\theta^\perp \in (\text{Ker } \widetilde{W})^\perp} & \frac{1}{N} \sum_{i=1}^N c \left(V \left(\widetilde{W} \theta^\perp \right) (x), y_i \right) + \frac{\lambda_K}{2} \|\theta^\perp\|_{\widetilde{\mathcal{H}}}^2 \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \left\langle \widetilde{\Phi}(x_i)^* \theta^\perp, M_{ik} \left(\widetilde{W} \theta^\perp \right) (x_k) \right\rangle_u.\end{aligned}$$

Hence minimizing over $(\text{Ker } \widetilde{W})^\perp$ or $\widetilde{\mathcal{H}}$ is the same when $\lambda_K > 0$. Eventually,

$$\begin{aligned}\theta_z = \arg \min_{\theta \in \widetilde{\mathcal{H}}} & \frac{1}{N} \sum_{i=1}^N c \left(V \widetilde{\Phi}(x_i)^* \theta, y_i \right) + \frac{\lambda_K}{2} \|\theta\|_{\widetilde{\mathcal{H}}}^2 \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \left\langle \widetilde{\Phi}(x_i)^* \theta, M_{ik} \widetilde{\Phi}(x_k)^* \theta \right\rangle_u \\ = \arg \min_{\theta \in \widetilde{\mathcal{H}}} & \frac{1}{N} \sum_{i=1}^N c \left(V \widetilde{\Phi}(x_i)^* \theta, y_i \right) + \frac{\lambda_K}{2} \|\theta\|_{\widetilde{\mathcal{H}}}^2 \\ & + \frac{\lambda_M}{2} \sum_{i,k=1}^{N+U} \left\langle \theta, \widetilde{\Phi}(x_i) M_{ik} \widetilde{\Phi}(x_k)^* \theta \right\rangle_{\widetilde{\mathcal{H}}}.\end{aligned}$$

This theorem is illustrated by [Figure B.1](#). We use the classic two moons dataset¹. We first perform an unsupervised spectral clustering step [181] and construct the matrix where C_{ik} is 1 if x_i and x_k are in the same cluster, 0 otherwise. Then we take the inverse Laplacian of this matrix and use it as the data dependent operator M . E

B.1.2 Gradients

By linearity and applying the chaine rule to [Equation B.6](#) and since $M_{ik}^* = M_{ki}$ for all $i, k \in \mathbb{N}_{N+U}^*$, we have

$$\begin{aligned}\nabla_\theta c \left(V \widetilde{\Phi}(x_i)^* \theta, y_i \right) &= \widetilde{\Phi}(x_i) V^* \left(\frac{\partial}{\partial y} c(y, y_i) \Big|_{y=V \widetilde{\Phi}(x_i)^* \theta} \right)^*, \\ \nabla_\theta \left\langle \widetilde{\Phi}(x_i)^* \theta, M_{ik} \widetilde{\Phi}(x_k)^* \theta \right\rangle_u &= \widetilde{\Phi}(x_i) (M_{ik} + M_{ki}^*) \widetilde{\Phi}(x_k)^* \theta, \\ \nabla_\theta \|\theta\|_{\widetilde{\mathcal{H}}}^2 &= 2\theta.\end{aligned}$$

Provided that $c(y, y_i)$ is Frechet differentiable w.r.t. y , for all y and $y_i \in \mathcal{Y}$ we have $\nabla_\theta J_{\lambda_K}(\theta) \in \widetilde{\mathcal{H}}$ and

$$\begin{aligned}\nabla_\theta J_{\lambda_K}(\theta) = \frac{1}{N} \sum_{i=1}^N & \widetilde{\Phi}(x_i) V^* \left(\frac{\partial}{\partial y} c(y, y_i) \Big|_{y=V \widetilde{\Phi}(x_i)^* \theta} \right)^* \\ & + \lambda_K \theta + \lambda_M \sum_{i,k=1}^{N+U} \widetilde{\Phi}(x_i) M_{ik} \widetilde{\Phi}(x_k)^* \theta\end{aligned}$$

¹ Available at http://scikit-learn.org/stable/modules/generated/sklearn.datasets.make_moons.html.

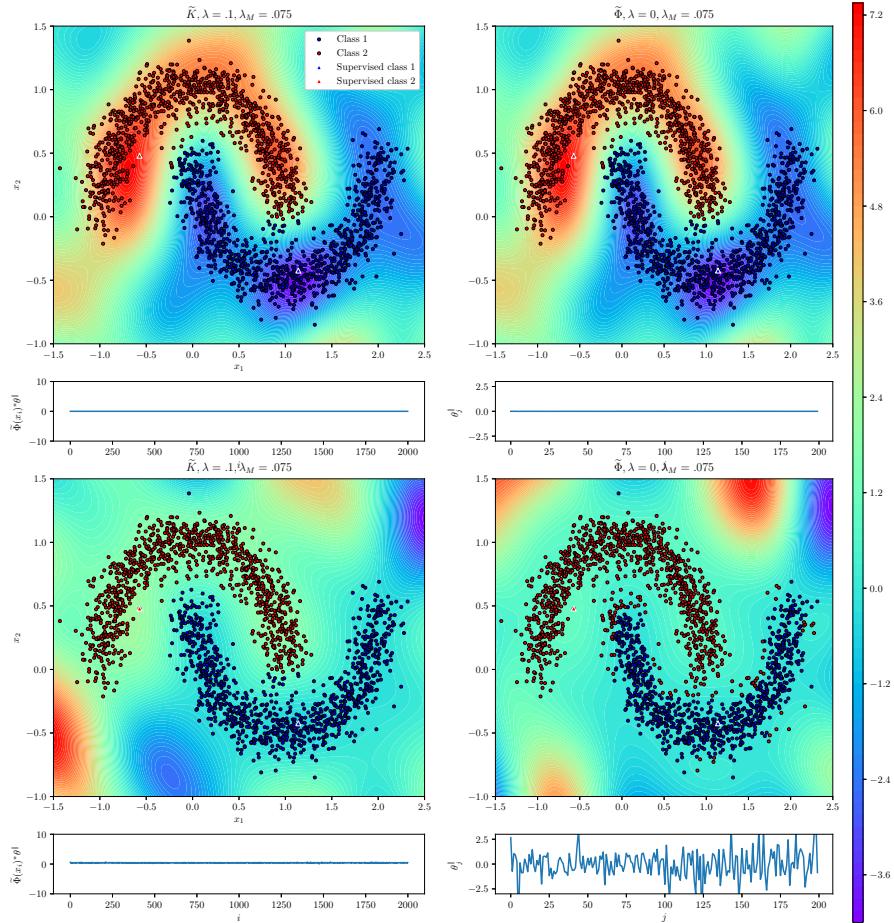


Figure B.1: ORFF equivalence theorem (semi-supervised). Each row compares the scalar ORFF $\tilde{\Phi}$ method constructed from a Gaussian with the kernel method where $\tilde{K} = \tilde{\Phi}^T \tilde{\Phi}$. The top row corresponds to the case $\lambda_K = 0.1$ and $\lambda_M = 0.075$. Since $\lambda_K > 0$, the solution with \tilde{K} and $\tilde{\Phi}$ are exactly the same (Theorem B.2 applies) and we see that $\theta^\parallel = 0$. The bottom row corresponds to the case $\lambda_K = 0$ and $\lambda_M = 0.075$. Here the solution with \tilde{K} and $\tilde{\Phi}$ doesn't match (Theorem B.2 fails to apply since $\lambda_K = 0$). Moreover we can see that $\theta^\parallel \neq 0$ and $\tilde{\Phi}(x)^* \theta \neq 0$, thus θ is not in $(\text{Ker } W)^\perp$.

Therefore after factorization, considering $\lambda_K > 0$,

$$\begin{aligned}\nabla_{\theta} J_{\lambda_K}(\theta) &= \frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) V^* \left(\frac{\partial}{\partial y} c(y, y_i) \Big|_{y=V\tilde{\Phi}(x_i)^*\theta} \right)^* \\ &\quad + \lambda_K \left(I_{\tilde{\mathcal{H}}} + \frac{\lambda_M}{\lambda_K} \sum_{i,k=1}^{N+U} \tilde{\Phi}(x_i) M_{ik} \tilde{\Phi}(x_k)^* \right) \theta\end{aligned}$$

We note the quantity

$$\tilde{\mathbf{M}}_{(\lambda_K, \lambda_M)} = I_{\tilde{\mathcal{H}}} + \frac{\lambda_M}{\lambda_K} \sum_{i,k=1}^{N+U} \tilde{\Phi}(x_i) M_{ik} \tilde{\Phi}(x_k)^* \in \mathcal{L}(\tilde{\mathcal{H}}) \quad (\text{B.7})$$

so that

$$\nabla_{\theta} J_{\lambda_K}(\theta) = \frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) V^* \left(\frac{\partial}{\partial y} c(y, y_i) \Big|_{y=V\tilde{\Phi}(x_i)^*\theta} \right)^* + \lambda_K \tilde{\mathbf{M}}_{(\lambda_K, \lambda_M)} \theta \quad (\text{B.8})$$

Example B.1 (Naive closed form for the squared error cost). Consider the cost function defined for all $y, y' \in \mathcal{Y}$ by $c(y, y') = \frac{1}{2} \|y - y'\|_2^2$. Then

$$\left(\frac{\partial}{\partial y} c(y, y_i) \Big|_{y=V\tilde{\Phi}(x_i)^*\theta} \right)^* = (V\tilde{\Phi}(x_i)^*\theta - y_i).$$

Thus, since the optimal solution θ_z verifies $\nabla_{\theta z} J_{\lambda_K}(\theta_z) = 0$ we have

$$\frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) V^* (V\tilde{\Phi}(x_i)^*\theta_z - y_i) + \lambda_K \tilde{\mathbf{M}}_{(\lambda_K, \lambda_M)} \theta_z = 0.$$

Therefore,

$$\left(\frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) V^* V\tilde{\Phi}(x_i)^* + \lambda_K \tilde{\mathbf{M}}_{(\lambda_K, \lambda_M)} \right) \theta_z = \frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) V^* y_i$$

Suppose that $\mathcal{Y} \subseteq \mathbb{R}^p$, $V : \mathcal{U} \rightarrow \mathcal{Y}$ where $\mathcal{U} \subseteq \mathbb{R}^u$ and for all $x \in \mathcal{X}$, $\tilde{\Phi}(x) : \mathbb{R}^r \rightarrow \mathbb{R}^u$ where all spaces are endowed with the euclidean inner product. From this we can derive [Algorithm 6](#) which returns the closed form solution of [Equation 6.10](#) for $c(y, y') = \frac{1}{2} \|y - y'\|_2^2$.

b.1.3 Complexity

Suppose that $u = \dim(\mathcal{U}) < +\infty$ and $u' = \dim(\mathcal{U}') < \infty$ and for all $x \in \mathcal{X}$, $\tilde{\Phi}(x) : \mathcal{U}' \rightarrow \tilde{\mathcal{H}}$ where $r = \dim(\tilde{\mathcal{H}}) < \infty$ is the dimension of the redescription space $\tilde{\mathcal{H}} = \mathbb{R}^r$. Since u, u' , and $r < \infty$, we view the operators $\tilde{\Phi}(x)$, V and $\tilde{\mathbf{M}}_{(\lambda_K, \lambda_M)}$ as matrices. Computing V^*V cost $O_t(u^2p)$. Step 1 costs $O_t(r^2u + ru^2)$. Steps 5 (optional) has the same

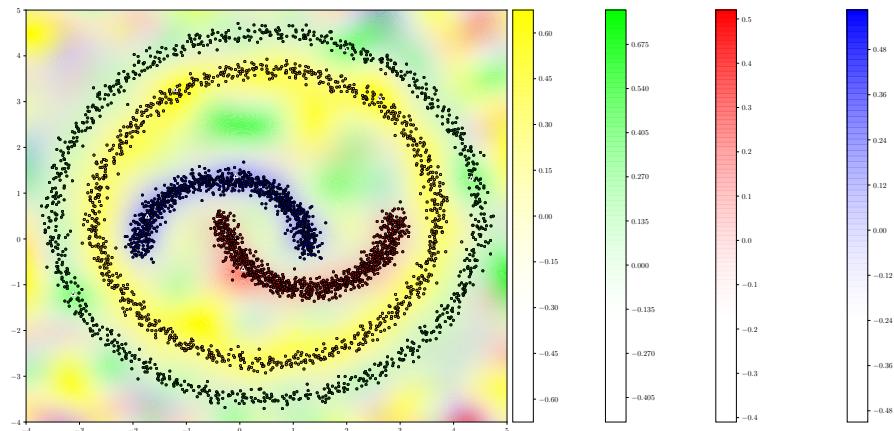


Figure B.2: Semi-supervised learning on a multiclass nested circle and two moons dataset. We performed an unsupervised spectral clustering [181] step to construc the matrices M_{ik} 's and, choose the matrix B of the ORFF using the simplex coding [123].

Algorithm 6: Naive closed form for the squared error cost.**Input :**

- $\mathbf{s} = (x_i, y_i)_{i=1}^N \in (\mathcal{X} \times \mathbb{R}^p)^N$ a sequence of supervised training points,
- $\mathbf{u} = (x_i)_{i=N+1}^{N+U} \in \mathcal{X}^U$ a sequence of unsupervised training points,
- $\tilde{\Phi}(x_i) \in \mathcal{L}(\mathbb{R}^u, \mathbb{R}^r)$ a feature map defined for all $x_i \in \mathcal{X}$,
- $(M_{ik})_{i,k=1}^{N+U}$ a sequence of data dependent operators (see [Theorem B.2](#)),
- $V \in \mathcal{L}(\mathbb{R}^u, \mathbb{R}^p)$ a combination operator,
- $\lambda_K \in \mathbb{R}_{>0}$ the Tychonov regularization term,
- $\lambda_M \in \mathbb{R}_+$ the manifold regularization term.

Output: A model

$$h : \begin{cases} \mathcal{X} \rightarrow \mathbb{R}^p \\ x \mapsto \tilde{\Phi}(x)^T \theta_z, \end{cases}$$

such that θ_z minimize [Equation 6.10](#), where $c(y, y') = \|y - y'\|_2^2$
and \mathbb{R}^u , \mathbb{R}^r and \mathbb{R}^p are Hilbert spaces endowed with
the euclidean inner product.

```

1  $P \leftarrow \frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) V^T V \tilde{\Phi}(x_i)^T \in \mathcal{L}(\mathbb{R}^r, \mathbb{R}^r);$ 
2 if  $\lambda_M = 0$  then
3    $\tilde{M}_{(\lambda_K, \lambda_M)} \leftarrow I_r \in \mathcal{L}(\mathbb{R}^r, \mathbb{R}^r);$ 
4 else
5    $\tilde{M}_{(\lambda_K, \lambda_M)} \leftarrow \left( I_r + \frac{\lambda_M}{\lambda_K} \sum_{i,k=1}^{N+U} \tilde{\Phi}(x_i) M_{ik} \tilde{\Phi}(x_k)^T \right) \in \mathcal{L}(\mathbb{R}^r, \mathbb{R}^r);$ 
6 end
7  $Y \leftarrow \frac{1}{N} \sum_{i=1}^N \tilde{\Phi}(x_i) V^T y_i \in \mathbb{R}^r;$ 
8  $\theta_z \leftarrow \text{solve}_\theta \left( (P + \lambda_K \tilde{M}_{(\lambda_K, \lambda_M)}) \theta = Y \right) \in \mathbb{R}^r;$ 
9 return  $h : x \mapsto \tilde{\Phi}(x)^T \theta_z;$ 

```

cost except that the sum is done over all pair of $N + U$ points thus it costs $O_t((N + U)^2(r^2u + ru^2))$. Steps 7 costs $O_t(N(ru + up))$. For step 8, the naive inversion of the operator costs $O_t(r^3)$. Eventually the overall complexity of [Algorithm 6](#) is

$$O_t \left(ru(r + u) \begin{cases} (N + U)^2 & \text{if } \lambda_M > 0 \\ N & \text{if } \lambda_M = 0 \end{cases} + r^3 + Nu(r + p) \right),$$

while the space complexity is $O_s(r^2)$. This complexity is to compare with the kernelized solution proposed by Minh, Bazzani, and Murino [120]. Let

$$\mathbf{K} : \begin{cases} \mathcal{U}^{N+U} \rightarrow \mathcal{U}^{N+U} \\ u \mapsto \bigoplus_{i=1}^{N+U} \sum_{j=1}^{N+U} K(x_i, x_j) u_j \end{cases}$$

and

$$\mathbf{M} : \begin{cases} \mathcal{U}^{N+U} \rightarrow \mathcal{U}^{N+U} \\ u \mapsto \bigoplus_{i=1}^{N+U} \sum_{k=1}^{N+U} M_{ik} u_k. \end{cases}$$

When $\mathcal{U} = \mathbb{R}$,

$$\mathbf{K} = \begin{pmatrix} K(x_1, x_1) & \dots & K(x_1, x_{N+U}) \\ \vdots & \ddots & \vdots \\ K(x_{N+U}, x_1) & \dots & K(x_{N+U}, x_{N+U}) \end{pmatrix}$$

is called the Gram matrix of K . When $\mathcal{U} = \mathbb{R}^p$, \mathbf{K} is a matrix-valued Gram matrix of size $u(N+U) \times u(N+U)$ where each entry $K_{ij} \in \mathcal{M}_{u,u}(\mathbb{R})$. When $\mathcal{U} = \mathbb{R}^u$, \mathbf{M} can also be seen as a matrix-valued matrix where each entry is $M_{ik} \in \mathcal{M}_{u,u}(\mathbb{R})$. We also introduce the matrices $\mathbf{V}^\top \mathbf{V} := I_{N+U} \otimes (V^\top V)$ and

$$\mathbf{P} : \begin{cases} \mathcal{U}^{N+U} \rightarrow \mathcal{U}^{N+U} \\ u \mapsto \left(\bigoplus_{j=1}^N u_j \right) \oplus \left(\bigoplus_{j=N+1}^{N+U} 0 \right) \end{cases}$$

The operator \mathbf{P} is a projection that sets all the terms u_j , $N < j \leq N+U$ of u to zero. When $\mathcal{U} = \mathbb{R}^u$ it can also be seen as the block matrix of size $u(N+U) \times u(N+U)$ and

$$\mathbf{P} = \begin{pmatrix} & & 0 & \dots & 0 \\ & I_u \otimes I_N & \vdots & \ddots & \vdots \\ & & 0 & \dots & 0 \\ 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 & \dots & 0 \end{pmatrix}$$

Then the equivalent kernelized solution u_z of [Theorem B.1](#) is given by Minh, Bazzani, and Murino [120]

$$\left(\frac{1}{N} \mathbf{V}^\top \mathbf{V} \mathbf{P} \mathbf{K} + \lambda_M \mathbf{M} \mathbf{K} + \lambda_K I_{\bigoplus_{i=1}^{N+U} \mathcal{U}} \right) u_z = \left(\bigoplus_{i=1}^N V^\top y_i \right) \oplus \left(\bigoplus_{i=N+1}^{N+U} 0 \right).$$

which has time complexity $O_t(((N+U)u)^3 + Nup)$ and space complexity $O_s(((N+U)u)^2)$. Notice that computing the data dependent norm

(manifold regularization) is expensive. Indeed when $\lambda_M = 0$, [Algorithm 6](#) has a linear complexity with respect to the number of supervised training points N while the complexity becomes quadratic in the number of supervised and unsupervised training points $N + U$ when $\lambda_M > 0$.

C

RELEVANT PIECE OF CODE

In the following, we give minimal short samples of Python code showing how to implement efficient ORFF. Each section represent an independent snippet of code, and a simple copy-paste in a python editor should generate the corresponding figure presented in this manuscript.

Contents

c.1	Python code for figure 3.1	198
c.2	Python code for figure 5.1	199
c.3	Python code for figure 5.3	201
c.4	Python code for figure 5.4	203
c.5	Python code for figure 5.5	206
c.6	Python (tensorflow) code for continuous quantile regression	209

C.1 PYTHON CODE FOR FIGURE 3.1

```

r"""Plot figure: Different outcomes of a Gaussian kernel approximation."""

import matplotlib
import numpy as np
import matplotlib.pyplot as plt

from sklearn.metrics import pairwise_kernels

def phi(x, w, D):
    r"""RFF map."""
    Z = np.dot(x, w)
    return np.hstack((np.cos(Z), np.sin(Z))) / np.sqrt(D)

def createColorbar(lwr, upr, fig, axes):
    r"""Create colorbar for multiple Pyplot plot."""
    cax = fig.add_axes([.92, 0.1, 0.01, 0.8])
    norm = matplotlib.colors.LogNorm(vmin=lwr, vmax=upr, clip=False)
    c = matplotlib.colorbar.ColorbarBase(cax, cmap=plt.get_cmap('rainbow'),
                                         norm=norm, label='D=')
    plt.title(r'$\widetilde{K}$')
    return c

def main():
    r"""Plot figure: Different outcomes of a Gaussian kernel approximation."""
    T = 25 # Number of curves

    cm_subsection = np.linspace(0, 1, T + 1)
    colors = [matplotlib.cm.rainbow(x) for x in cm_subsection]

    d = 1 # Dimension of the input
    N = 250 # Number of points per curves

    # Generate N data in (-1, 1) and exact Gram matrix
    np.random.seed(0)
    X = np.linspace(-1, 1, N).reshape(N, d)
    K = pairwise_kernels(X, metric='rbf', gamma=1. / (2. * .1 ** 2))

    # A Matrix for the decomposable kernel. Link the outputs to some mean value
    c = np.random.randn(N, 2)
    A = .5 * np.eye(2) + .5 * np.ones((2, 2))

    plt.close()
    plt.rc('text', usetex=True)
    plt.rc('font', family='serif')
    f, axes = plt.subplots(2, 2, figsize=(12, 8), sharex=True, sharey=True)

    # For each curve with different D
    for k, D in enumerate(np.logspace(0, 4, T)):
        D = int(D)
        np.random.seed(0)

        w = np.random.randn(d, D) / .1
        phiX = phi(X, w, D)
        Kt = np.dot(phiX, phiX.T)

        # Generate outputs with the exact Gram matrix
        pred = np.dot(np.dot(Kt, c), A)
        axes[0, 0].plot(X, pred[:, 0], c=colors[k], lw=.5, linestyle='--')
        axes[0, 0].set_ylabel(r'$y_1$')
        axes[0, 1].plot(X, pred[:, 1], c=colors[k], lw=.5, linestyle='--')
        axes[0, 1].set_ylabel(r'$y_2$')

```

```

# Generate outputs with the a realization of the random Gram matrix
w = np.random.randn(d, D) / .1
phiX = phi(X, w, D)
Kt = np.dot(phiX, phiX.T)

pred = np.dot(np.dot(Kt, c), A)
axes[1, 0].plot(X, pred[:, 0], c=colors[k], lw=.5, linestyle='--')
axes[1, 0].set_xlabel(r'$x$')
axes[1, 0].set_ylabel(r'$y_1$')
axes[1, 1].plot(X, pred[:, 1], c=colors[k], lw=.5, linestyle='--')
axes[1, 1].set_xlabel(r'$x$')
axes[1, 1].set_ylabel(r'$y_2$')

axes[0, 0].plot(X, np.dot(np.dot(K, c), A)[:, 0], c='k', lw=.5, label='K')
axes[0, 1].plot(X, np.dot(np.dot(K, c), A)[:, 1], c='k', lw=.5, label='K')
axes[1, 0].plot(X, np.dot(np.dot(K, c), A)[:, 0], c='k', lw=.5, label='K')
axes[1, 1].plot(X, np.dot(np.dot(K, c), A)[:, 1], c='k', lw=.5, label='K')

axes[0, 0].set_title(r'$\widetilde{K} \approx K$, realization 1', x=1.1)
axes[1, 0].set_title(r'$\widetilde{K} \approx K$, realization 2', x=1.1)

for xx in axes.ravel():
    xx.legend(loc=4)

createColorbar(1, D, f, axes)
plt.savefig('not_Mercer.pgf', bbox_inches='tight')

if __name__ == "__main__":
    main()

```

C.2 PYTHON CODE FOR FIGURE 5.1

```

"""Plot figure: ORFF Representer theorem pt.1."""

import numpy as np
import matplotlib.pyplot as plt
import matplotlib

def phi(x, w, D):
    """RFF map."""
    Z = np.dot(x, w)
    return np.hstack((np.cos(Z), np.sin(Z))) / np.sqrt(D)

def main():
    """Plot figure: ORFF Representer theorem pt.1."""
    d = 1 # dimensionality of the inputs
    D = 50 # number of random features

    N = 50
    Nt = 200

    # N training points in (0,1)
    np.random.seed(0)
    x = 2 * np.random.rand(N, d) - 1
    y = np.sin(10 * x)
    y += .5 * np.random.randn(y.shape[0], y.shape[1]) + 2. * x ** 2

    # Nt testing points in (0,1)
    xt = np.linspace(-1, 1, Nt).reshape((-1, 1))
    yt = np.sin(10 * xt) + 2. * xt ** 2
    yt += .5 * np.random.randn(yt.shape[0], yt.shape[1])

```

```

sigma = .3
w = np.random.randn(d, D) / sigma # Realization of  $(\omega_j)_{j=1}^D$ 

phiX = phi(x, w, D) # Train RFF
phiXt = phi(xt, w, D) # Test RFF

# Create plot
plt.close()
plt.rc('text', usetex=True)
plt.rc('font', family='serif')
f, axis = plt.subplots(4, 3, gridspec_kw={'width_ratios': [3, 3, 1.5]}, figsize=(16, 6), sharex='col', sharey='col')
f.subplots_adjust(hspace=.25)
formatter = matplotlib.ticker.ScalarFormatter()
formatter.set_powerlimits((-3, 4))

# For different hyperparameters \lambda
for k, lbda in enumerate([1e-2, 5e-6, 1e-10, 0]):
    # Train with ORFF with kernel approximation (dual)
    ck = np.linalg.lstsq(np.dot(phiX, phiX.T) + lbda * np.eye(N),
                         y, rcond=-1)[0]
    # Train with ORFF without kernel approximation (primal)
    c = np.linalg.lstsq(np.dot(phiX.T, phiX) + lbda * np.eye(2 * D),
                        np.dot(phiX.T, y), rcond=-1)[0]
    cc = np.sum((phi(x, w, D) * ck), axis=0)
    # Link dual coefficient with primal coefficients
    cr = (cc - c.ravel()) / np.linalg.norm(c) * 100
    err = np.array([np.linalg.norm(np.dot(phiXt, c) - yt) ** 2 / Nt,
                    np.linalg.norm(np.dot(np.dot(phiXt,
                                                phiX.T),
                                         ck) - yt) ** 2 / Nt,
                    np.linalg.norm(np.dot(phiXt, cr)) ** 2 / Nt,
                    np.linalg.norm(cr)])
    # Plot
    lmin = -1.8
    lmax = 3.
    axis[k, 0].set_xlim([-1.5, 1])
    axis[k, 0].set_ylim([lmin, lmax])
    axis[k, 0].plot(xt, np.dot(phiXt, c),
                     label=r'$\widetilde{\Phi}^* \theta$')
    axis[k, 0].plot(xt, np.dot(np.dot(phiXt, phiX.T), ck),
                     label=r'$\widetilde{K}u$', linestyle='-.')
    axis[k, 0].scatter(x, y, c='r', marker='+', label='train', lw=2)
    axis[k, 0].scatter(xt, yt, c='k', marker='.', label='test')
    axis[k, 0].legend(loc=3)
    axis[k, 0].set_ylabel('y')
    if k == 3:
        axis[k, 0].set_xlabel('x')

    lmin = -1.8
    lmax = 3.
    pred = np.dot(phi(xt, w, D), cr)
    axis[k, 1].set_xlim([-1.5, 1])
    axis[k, 1].set_ylim([lmin, lmax])
    axis[k, 1].plot(xt, pred,
                     label=r'$\widetilde{\Phi}^* \theta_{parallel}$')
    axis[k, 1].scatter(x, y, c='r', marker='+', label='train', lw=2)
    axis[k, 1].scatter(xt, yt, c='k', marker='.', label='test')
    axis[k, 1].legend(loc=3)
    if k == 3:
        axis[k, 1].set_xlabel('x')

    xs = np.arange(cr.size)
    axis[k, 2].barh(xs, np.abs(cr), edgecolor="none", log=True)
    axis[k, 2].set_ylabel(r'$j$')

```

```

    if k == 3:
        axis[k, 2].set_xlabel(
            r'$|\theta^{\parallel}|_j|$, \% of relative error')
    plt.savefig('representer.pgf', bbox_inches='tight')

    return err

if __name__ == "__main__":
    main()

```

C.3 PYTHON CODE FOR FIGURE 5.3

"""Efficient implementation of the Gaussian ORFF decomposable kernel."""

```

from time import time

from pypmlel.asizeof import asizeof

from numpy.linalg import svd
from numpy.random import rand, seed
from numpy import (dot, diag, sqrt, kron, zeros,
                  logspace, log10, matrix, eye, int, float)
from scipy.sparse.linalg import LinearOperator
from sklearn.kernel_approximation import RBFSampler
from matplotlib.pyplot import savefig, subplots, tight_layout

def NaiveDecomposableGaussianORFF(X, A, gamma=1.,
                                    D=100, eps=1e-5, random_state=0):
    r"""Return the Naive ORFF map associated with the data X.

    Parameters
    -----
    X : {array-like}, shape = [n_samples, n_features]
        Samples.
    A : {array-like}, shape = [n_targets, n_targets]
        Operator of the Decomposable kernel (positive semi-definite)
    gamma : {float},
        Gamma parameter of the RBF kernel.
    D : {integer},
        Number of random features.
    eps : {float},
        Cutoff threshold for the singular values of A.
    random_state : {integer},
        Seed of the generator.

    Returns
    -----
    \tilde{\Phi}(X) : array
    """
    # Decompose A=BB^T
    u, s, v = svd(A, full_matrices=False, compute_uv=True)
    B = dot(diag(sqrt(s[s > eps])), v[s > eps, :])

    # Sample a RFF from the scalar Gaussian kernel
    phi_s = RBFSampler(gamma=gamma, n_components=D, random_state=random_state)
    phiX = phi_s.fit_transform(X)

    # Create the ORFF linear operator
    return matrix(kron(phiX, B))

def EfficientDecomposableGaussianORFF(X, A, gamma=1.,
                                       D=100, eps=1e-5, random_state=0):
    r"""Return the Efficient ORFF map associated with the data X.

```

```

Parameters
-----
X : {array-like}, shape = [n_samples, n_features]
    Samples.
A : {array-like}, shape = [n_targets, n_targets]
    Operator of the Decomposable kernel (positive semi-definite)
gamma : {float},
    Gamma parameter of the RBF kernel.
D : {integer},
    Number of random features.
eps : {float},
    Cutoff threshold for the singular values of A.
random_state : {integer},
    Seed of the generator.

Returns
-----
\tilde{\Phi}(X) : Linear Operator, callable
"""

# Decompose A=BB^T
u, s, v = svd(A, full_matrices=False, compute_uv=True)
B = dot(diag(sqrt(s[s > eps])), v[s > eps, :])

# Sample a RFF from the scalar Gaussian kernel
phi_s = RBFSampler(gamma=gamma, n_components=D, random_state=random_state)
phiX = phi_s.fit_transform(X)

# Create the ORFF linear operator
cshape = (D, B.shape[0])
rshape = (X.shape[0], B.shape[1])
return LinearOperator((phiX.shape[0] * B.shape[1], D * B.shape[0]),
                      matvec=lambda b: dot(phiX, dot(b.reshape(cshape),
                                                     B)),
                      rmatvec=lambda r: dot(phiX.T, dot(r.reshape(rshape),
                                                       B.T)),
                      dtype=float)

def main():
    r"""Plot figure: Efficient decomposable gaussian ORFF."""
    N = 100 # Number of points
    pmax = 100 # Maximum output dimension
    d = 20 # Input dimension
    D = 100 # Number of random features

    seed(0)
    X = rand(N, d)

    R, T = 10, 10
    time_Efficient, mem_Efficient = zeros((R, T, 2)), zeros((R, T))
    time_naive, mem_naive = zeros((R, T, 2)), zeros((R, T))

    for i, p in enumerate(logspace(0, log10(pmax), T)):
        A = rand(int(p), int(p))
        A = dot(A.T, A) + eye(int(p))

        # Perform \Phi(X)^T \theta with Efficient implementation
        for j in range(R):
            start = time()
            phiX1 = EfficientDecomposableGaussianORFF(X, A, D)
            time_Efficient[j, i, 0] = time() - start
            theta = rand(phiX1.shape[1], 1)
            start = time()
            phiX1 * theta
            time_Efficient[j, i, 1] = time() - start

```

```

mem_Efficient[j, i] = asizeof(phiX1, code=True)

# Perform \Phi(X)^T \theta with naive implementation
for j in range(R):
    start = time()
    phiX2 = NaiveDecomposableGaussianORFF(X, A, D)
    time_naive[j, i, 0] = time() - start
    theta = rand(phiX2.shape[1], 1)
    start = time()
    phiX2 * theta
    time_naive[j, i, 1] = time() - start
    mem_naive[j, i] = asizeof(phiX2, code=True)

# Plot
f, axes = subplots(1, 3, figsize=(10, 4), sharex=True, sharey=False)
axes[0].errorbar(logspace(0, log10(pmax), T).astype(int),
                  time_Efficient[:, :, 0].mean(axis=0),
                  time_Efficient[:, :, 0].std(axis=0),
                  label='Efficient decomposable ORFF')
axes[0].errorbar(logspace(0, log10(pmax), T).astype(int),
                  time_naive[:, :, 0].mean(axis=0),
                  time_naive[:, :, 0].std(axis=0),
                  label='Naive decomposable ORFF')
axes[1].errorbar(logspace(0, log10(pmax), T).astype(int),
                  time_Efficient[:, :, 1].mean(axis=0),
                  time_Efficient[:, :, 1].std(axis=0),
                  label='Efficient decomposable ORFF')
axes[1].errorbar(logspace(0, log10(pmax), T).astype(int),
                  time_naive[:, :, 1].mean(axis=0),
                  time_naive[:, :, 1].std(axis=0),
                  label='Naive decomposable ORFF')
axes[2].errorbar(logspace(0, log10(pmax), T).astype(int),
                  mem_Efficient[:, :, 0].mean(axis=0),
                  mem_Efficient[:, :, 0].std(axis=0),
                  label='Efficient decomposable ORFF')
axes[2].errorbar(logspace(0, log10(pmax), T).astype(int),
                  mem_naive[:, :, 0].mean(axis=0),
                  mem_naive[:, :, 0].std(axis=0),
                  label='Naive decomposable ORFF')

axes[0].set_xscale('log')
axes[0].set_yscale('log')
axes[1].set_xscale('log')
axes[1].set_yscale('log')
axes[2].set_xscale('log')
axes[2].set_yscale('log')
axes[0].set_xlabel(r'$p=\dim(\mathcal{Y})$')
axes[1].set_xlabel(r'$p=\dim(\mathcal{Y})$')
axes[2].set_xlabel(r'$p=\dim(\mathcal{Y})$')
axes[0].set_ylabel('time (s)')
axes[2].set_ylabel('memory (bytes)')
axes[0].set_title('Preprocessing time')
axes[1].set_title(r'$\widetilde{\Phi}(X)^T \theta$ computation time')
axes[2].set_title(r'$\widetilde{\Phi}(X)^T$ required memory')
axes[0].legend(loc=2)
tight_layout()
savefig('efficient_decomposable_gaussian.pgf', bbox_inches='tight')

if __name__ == "__main__":
    main()

```

C.4 PYTHON CODE FOR FIGURE 5.4

```

r"""Efficient implementation of the Gaussian curl-free kernel."""

from time import time

```

```

from pympler.asizeof import asizeof

from numpy.random import rand, seed
from numpy import dot, zeros, logspace, log10, matrix, int, float
from scipy.sparse.linalg import LinearOperator
from sklearn.kernel_approximation import RBFSampler
from matplotlib.pyplot import savefig, subplots, tight_layout


def NaiveCurlFreeGaussianORFF(X, gamma=1.,
                               D=100, eps=1e-5, random_state=0):
    r"""Return the Naive ORFF map associated with the data X.

    Parameters
    -----
    X : {array-like}, shape = [n_samples, n_features]
        Samples.
    gamma : {float},
        Gamma parameter of the RBF kernel.
    D : {integer},
        Number of random features.
    eps : {float},
        Cutoff threshold for the singular values of A.
    random_state : {integer},
        Seed of the generator.

    Returns
    -----
    \tilde{\Phi}(X) : array
    """
    phi_s = RBFSampler(gamma=gamma, n_components=D,
                        random_state=random_state)
    phiX = phi_s.fit_transform(X)
    phiX = (phiX.reshape((phiX.shape[0], 1, phiX.shape[1])) *
            phi_s.random_weights_.reshape((1, -1, phiX.shape[1])))

    return matrix(phiX.reshape((-1, phiX.shape[2])))


def EfficientCurlFreeGaussianORFF(X, gamma=1.,
                                   D=100, eps=1e-5, random_state=0):
    r"""Return the Efficient ORFF map associated with the data X.

    Parameters
    -----
    X : {array-like}, shape = [n_samples, n_features]
        Samples.
    gamma : {float},
        Gamma parameter of the RBF kernel.
    D : {integer},
        Number of random features.
    eps : {float},
        Cutoff threshold for the singular values of A.
    random_state : {integer},
        Seed of the generator.

    Returns
    -----
    \tilde{\Phi}(X) : array
    """
    phi_s = RBFSampler(gamma=gamma, n_components=D,
                        random_state=random_state)
    phiX = phi_s.fit_transform(X)

    return LinearOperator((phiX.shape[0] * X.shape[1], phiX.shape[1]),

```



```

        label='Naive decomposable ORFF')
    axes[2].errorbar(logspace(0, log10(dmax), T).astype(int),
                      mem_Efficient[:, :].mean(axis=0),
                      mem_Efficient[:, :].std(axis=0),
                      label='Efficient decomposable ORFF')
    axes[2].errorbar(logspace(0, log10(dmax), T).astype(int),
                      mem_naive[:, :].mean(axis=0),
                      mem_naive[:, :].std(axis=0),
                      label='Naive decomposable ORFF')
axes[0].set_xscale('log')
axes[0].set_yscale('log')
axes[1].set_xscale('log')
axes[1].set_yscale('log')
axes[2].set_xscale('log')
axes[2].set_yscale('log')
axes[0].set_xlabel(r'$p=\dim(\mathcal{Y})$')
axes[1].set_xlabel(r'$p=\dim(\mathcal{Y})$')
axes[2].set_xlabel(r'$p=\dim(\mathcal{Y})$')
axes[0].set_ylabel(r'time (s)')
axes[2].set_ylabel(r'memory (bytes)')
axes[0].set_title(r'Preprocessing time')
axes[1].set_title(r'$\widetilde{\Phi}(X)^T \theta$ computation time')
axes[2].set_title(r'$\widetilde{\Phi}(X)^T$ required memory')
axes[0].legend(loc=2)
tight_layout()
savefig('efficient_curlfree_gaussian.pdf', bbox_inches='tight')

if __name__ == "__main__":
    main()

```

C.5 PYTHON CODE FOR FIGURE 5.5

```

r"""Efficient implementation of the Gaussian divergence-free kernel."""

from time import time

from pympler.asizeof import asizeof

from numpy.random import rand, seed
from numpy.linalg import norm
from numpy import dot, zeros, logspace, log10, matrix, int, eye, float
from scipy.sparse.linalg import LinearOperator
from sklearn.kernel_approximation import RBFSampler
from matplotlib.pyplot import savefig, subplots, tight_layout


def _rebase(phiX, W, Wn):
    return (phiX.reshape((phiX.shape[0], 1, 1, phiX.shape[1])) *
            (eye(W.shape[1]).reshape(1, W.shape[1], W.shape[1], 1) * Wn -
             W * W.reshape(1, 1, W.shape[1], phiX.shape[1]) / Wn)).reshape(
                (-1, W.shape[1] * Wn.shape[3]))


def NaiveDivergenceFreeGaussianORFF(X, gamma=1.,
                                     D=100, eps=1e-5, random_state=0):
    r"""Return the Naive ORFF map associated with the data X.

    Parameters
    -----
    X : {array-like}, shape = [n_samples, n_features]
        Samples.
    gamma : {float},
        Gamma parameter of the RBF kernel.
    D : {integer},
        Number of random features.
    """

```

```

eps : {float},
      Cutoff threshold for the singular values of A.
random_state : {integer},
      Seed of the generator.

>Returns
-----
\tilde{\Phi}(X) : array
"""

phi_s = RBFSampler(gamma=gamma, n_components=D,
                     random_state=random_state)

phiX = _rebase(phi_s.fit_transform(X),
                phi_s.random_weights_.reshape((1, -1, 1, D)),
                norm(phi_s.random_weights_, axis=0).reshape((1, 1, 1, -1)))

return matrix(phiX)

def EfficientDivergenceFreeGaussianORFF(X, gamma=1.,
                                         D=100, eps=1e-5, random_state=0):
    r"""Return the Efficient ORFF map associated with the data X.

Parameters
-----
X : {array-like}, shape = [n_samples, n_features]
    Samples.
gamma : {float},
    Gamma parameter of the RBF kernel.
D : {integer},
    Number of random features.
eps : {float},
    Cutoff threshold for the singular values of A.
random_state : {integer},
    Seed of the generator.

>Returns
-----
\tilde{\Phi}(X) : array
"""

phi_s = RBFSampler(gamma=gamma, n_components=D,
                     random_state=random_state)
phiX = phi_s.fit_transform(X)
W = phi_s.random_weights_.reshape((1, -1, 1, phiX.shape[1]))
Wn = norm(phi_s.random_weights_, axis=0).reshape((1, 1, 1, -1))
return LinearOperator((phiX.shape[0] * X.shape[1],
                      phiX.shape[1] * X.shape[1]),
                      matvec=lambda b: dot(_rebase(phiX, W, Wn), b),
                      rmatvec=lambda r: dot(_rebase(phiX, W, Wn).T, r),
                      dtype=float)

def main():
    r"""Plot figure: Efficient decomposable gaussian ORFF."""
    N = 100 # Number of points
    dmax = 100 # Input dimension
    D = 100 # Number of random features

    seed(0)

    R, T = 10, 10
    time_Efficient, mem_Efficient = zeros((R, T, 2)), zeros((R, T))
    time_naive, mem_naive = zeros((R, T, 2)), zeros((R, T))

    for i, d in enumerate(logspace(0, log10(dmax), T)):
        X = rand(N, int(d))

```

```

# Perform  $\Phi(X)^T \theta$  with Efficient implementation
for j in range(R):
    start = time()
    phiX1 = EfficientDivergenceFreeGaussianORFF(X, D)
    time_Efficient[j, i, 0] = time() - start
    theta = rand(phiX1.shape[1], 1)
    start = time()
    phiX1 * theta
    time_Efficient[j, i, 1] = time() - start
    mem_Efficient[j, i] = asizeof(phiX1, code=True)

# Perform  $\Phi(X)^T \theta$  with naive implementation
for j in range(R):
    start = time()
    phiX2 = NaiveDivergenceFreeGaussianORFF(X, D)
    time_naive[j, i, 0] = time() - start
    theta = rand(phiX2.shape[1], 1)
    start = time()
    phiX2 * theta
    time_naive[j, i, 1] = time() - start
    mem_naive[j, i] = asizeof(phiX2, code=True)

# Plot
f, axes = subplots(1, 3, figsize=(10, 4), sharex=True, sharey=False)
axes[0].errorbar(logspace(0, log10(dmax), T).astype(int),
                  time_Efficient[:, :, 0].mean(axis=0),
                  time_Efficient[:, :, 0].std(axis=0),
                  label='Efficient decomposable ORFF')
axes[0].errorbar(logspace(0, log10(dmax), T).astype(int),
                  time_naive[:, :, 0].mean(axis=0),
                  time_naive[:, :, 0].std(axis=0),
                  label='Naive decomposable ORFF')
axes[1].errorbar(logspace(0, log10(dmax), T).astype(int),
                  time_Efficient[:, :, 1].mean(axis=0),
                  time_Efficient[:, :, 1].std(axis=0),
                  label='Efficient decomposable ORFF')
axes[1].errorbar(logspace(0, log10(dmax), T).astype(int),
                  time_naive[:, :, 1].mean(axis=0),
                  time_naive[:, :, 1].std(axis=0),
                  label='Naive decomposable ORFF')
axes[2].errorbar(logspace(0, log10(dmax), T).astype(int),
                  mem_Efficient[:, :].mean(axis=0),
                  mem_Efficient[:, :].std(axis=0),
                  label='Efficient decomposable ORFF')
axes[2].errorbar(logspace(0, log10(dmax), T).astype(int),
                  mem_naive[:, :].mean(axis=0),
                  mem_naive[:, :].std(axis=0),
                  label='Naive decomposable ORFF')
axes[0].set_xscale('log')
axes[0].set_yscale('log')
axes[1].set_xscale('log')
axes[1].set_yscale('log')
axes[2].set_xscale('log')
axes[2].set_yscale('log')
axes[0].set_xlabel(r'$p=\dim(\mathcal{Y})$')
axes[1].set_xlabel(r'$p=\dim(\mathcal{Y})$')
axes[2].set_xlabel(r'$p=\dim(\mathcal{Y})$')
axes[0].set_ylabel(r'time (s)')
axes[2].set_ylabel(r'memory (bytes)')
axes[0].set_title(r'Preprocessing time')
axes[1].set_title(r'$\widetilde{\Phi}(X)^T \theta$ computation time')
axes[2].set_title(r'$\widetilde{\Phi}(X)^T$ required memory')
axes[0].legend(loc=2)
tight_layout()
savefig('efficient_divfree_gaussian.pdf', bbox_inches='tight')

```

```

if __name__ == "__main__":
    main()

C.6 PYTHON (TENSORFLOW) CODE FOR CONTINUOUS QUAN-
TILE REGRESSION

import tensorflow as tf
import numpy as np
import matplotlib.pyplot as plt
from time import time
from operalib import toy_data_quantile
from sklearn.metrics.pairwise import euclidean_distances
from operalib import toy_data_quantile

def phi(X, omegas, D):
    Z = tf.matmul(X, omegas)
    phiX = tf.concat([tf.cos(Z), tf.sin(Z)], 1) / np.sqrt(D)
    return phiX

def phigrad(X, omegas, D):
    Z = tf.matmul(X, omegas)
    Zc = tf.cos(Z)
    Zs = tf.sin(Z)
    phiX = tf.concat([Zc, Zs], 1) / np.sqrt(D)
    phiXg = tf.concat([-omegas * Zs, omegas * Zc], 1) / np.sqrt(D)
    return phiX, phiXg

def model(phiX, phit, theta):
    return tf.matmul(tf.matmul(phiX, theta), tf.transpose(phit))

def loss(theta, X, t, y, omegas1, omegas2, D1, D2, lbda1, lbda2):
    phiX = phi(X, omegas1, D1)
    phit, phitg = phigrad(t, omegas2, D2)

    gxp = tf.matmul(phiX, theta)
    pred = tf.matmul(gxp, tf.transpose(phit))

    pin = tf.reduce_mean(tf.reduce_mean(tf.where(tf.greater_equal(pred, y),
                                                (pred - y) * tf.transpose(t),
                                                (y - pred) *
                                                (1 - tf.transpose(t))), 1)))
    cross = tf.reduce_mean(tf.reduce_mean(tf.maximum(tf.matmul(gxp,
                                                               tf.transpose(phitg)), 0), 1)))
    reg = tf.nn.l2_loss(theta) / (D1 * D2)
    return pin + lbda1 * cross + lbda2 * reg

def main():
    np.random.seed(0)

    print("Creating dataset...")
    N = 2000
    Nt = 1000
    d = 1
    p = 1
    probs = np.linspace(0.05, 0.95, 5) # Quantile levels of interest
    x_train, y_train, z_train = toy_data_quantile(N)
    x_test, y_test, z_test = toy_data_quantile(Nt, probs=probs)

```

```

ctype = tf.float32

lbda1 = 1
lbda2 = 1e-5
ts = 100
D1 = 100
D2 = 100
sigma1 = .25
tn = np.random.rand(ts)
sigma2 = np.median(euclidean_distances(tn.reshape(-1, 1)))
with tf.device('/gpu:0'):
    X = tf.placeholder(ctype, [N, d], name='input_batch')
    y = tf.placeholder(ctype, [N, p], name='input_batch')
    t = tf.placeholder(ctype, [ts, 1])

    omegas1 = tf.Variable(tf.random_normal([d, D1],
                                           mean=0, stddev=1 / sigma1,
                                           dtype=ctype), trainable=False)
    omegas2 = tf.Variable(tf.random_normal([d, D2],
                                           mean=0, stddev=1 / sigma2,
                                           dtype=ctype), trainable=False)
    theta = tf.Variable(tf.random_normal([2 * D1, 2 * D2],
                                         mean=0, stddev=1,
                                         dtype=ctype), trainable=True)
    ls = loss(theta, X, t, y, omegas1, omegas2, D1, D2, lbda1, lbda2)

    opt = tf.train.RMSPropOptimizer(.00075).minimize(ls)
    test_X = tf.placeholder(ctype, [Nt, d], name='input_batch')
    test_t = tf.placeholder(ctype, [None, 1], name='input_batch')
    phitest_X = phi(test_X, omegas1, D1)
    phitest_t, phigtest_t = phigrad(test_t, omegas2, D2)

config = tf.ConfigProto(allow_soft_placement=True)
start = time()
tt5 = 1 - np.linspace(0.05, 0.95, 5).reshape(-1, 1)
with tf.Session(config=config) as session:
    init = tf.global_variables_initializer()
    session.run(init)

    for i in range(0, 1000):
        session.run(opt,
                    feed_dict={X: np.asarray(x_train, dtype=np.float32),
                               y: np.asarray(y_train.reshape(-1, 1),
                                             dtype=np.float32),
                               t: np.asarray(tn.reshape(-1, 1),
                                             dtype=np.float32)})}

    if i % 100 == 0:
        print(session.run(ls,
                          feed_dict={X: np.asarray(x_train,
                                                    dtype=np.float32),
                                     y: np.asarray(y_train.reshape(-1,
                                                               1),
                                                               dtype=np.float32),
                                     t: np.asarray(tn.reshape(-1, 1),
                                                               dtype=np.float32)}))

    pred_test = session.run(model(phitest_X, phitest_t, theta),
                            feed_dict={test_X:
                                       np.asarray(x_test,
                                                 dtype=np.float32),
                                       test_t:
                                       np.asarray(tt5, dtype=np.float32)}))

session.close()
print(time() - start)
print('done')

plt.figure(figsize=(20, 10))

```

```
plt.gca().set_prop_cycle(None)
for i, q in enumerate(z_test):
    plt.plot(x_test, q, '--', label='true quantile at ' + str(probs[i]))
plt.gca().set_prop_cycle(None)
plt.plot(x_test, pred_test, '--', label='Continuous Quantile')
plt.scatter(x_train.ravel(), y_train.ravel(), marker='.', c='k')
plt.xlabel('x')
plt.ylabel('y')
plt.legend()
plt.tight_layout()
plt.show()

if __name__ == "__main__":
    main()
```

*
*
*

D

ONE CLASS SPLITTING CRITERIA FOR RANDOM FORESTS

This thesis has also been the opportunity to be part of a collaborative work on anomaly detection with random forests, with other fellow Ph. D. students of Télécom ParisTech. The following paper Goix et al. [73] is based on a original idea of Nicolas Goix and a joint work with Nicolas Drougard and Maël Chiapino. It is currently under review at ECML. Our original paper can be found at <https://arxiv.org/pdf/1611.01971.pdf>.

Contents

d.1	Background on decision trees	217
d.2	Adaptation to the one-class setting	218
d.2.1	One-class splitting criterion	218
d.2.2	Prediction: scoring function of the forest .	222
d.2.3	OneClassRF: a Generic One-Class Random Forest algorithm	223
d.3	Benchmarks	224
d.3.1	Default parameters of OneClassRF	224
d.3.2	Results	224
d.4	Theoretical analysis	226
d.4.1	Underlying model	226
d.4.2	Adaptive approach	228
d.5	Conclusion	229
d.6	Further insights on the algorithm	230
d.6.1	Interpretation of parameter gamma	230
d.6.2	Alternative scoring functions	230
d.6.3	Alternative stopping criteria	231
d.6.4	Variable importance	231
d.7	Hyper-parameters of tested algorithms	231
d.8	Description of the datasets	233
d.9	Further details on benchmarks and outlier detection results	234

Anomalies, novelties or outliers are usually assumed to lie in low probability regions of the data generating process. This assumption drives many statistical anomaly detection methods. Parametric techniques [16, 58] suppose that the inliers are generated by a distribution belonging to some specific parametric model a priori known. Here and hereafter, we denote by inliers the “not abnormal” data, and by outliers/anomalies/novelties the data from the abnormal class. Classical non-parametric approaches are based on density (level set) estimation [33, 138, 147, 151], on dimensionality reduction [2, 157] or on decision trees [105, 156]. Relevant overviews of current research on anomaly detection can be found in Chandola, Banerjee, and Kumar [42], Hodge and Austin [79], Markou and Singh [111], and Patcha and Park [131].

The algorithm proposed in this paper lies in the *novelty detection* setting, also called *one-class classification*. In this framework, we assume that we only observe examples of one class (referred to as the normal class, or inlier class). The second –hidden– class is called the abnormal class, or outlier class. The goal is to identify characteristics of the inlier class, such as its support or some density level sets with levels close to zero. This setup is for instance used in some –non-parametric– kernel methods such as OCSM [147], which extends the SVM methodology [48, 155] to handle training using only inliers. Recently, LSAD [138], a kernel method similarly extends a multi-class probabilistic classifier [165] to the one-class setting.

Random Forests (RFs) are strong machine learning tools [32], comparing well with state-of-the-art methods such as SVM or boosting algorithms [66], and used in a wide range of domains [53, 69, 168]. These estimators fit a number of decision tree classifiers on different random sub-samples of the dataset. Each tree is built recursively, according to a splitting criterion based on some impurity measure of a node. The prediction is done by an average over each tree prediction. In classification the averaging is based on a majority vote. Practical and theoretical insights on RFs are given in Biau, Devroye, and Lugosi [21], Biau and Scornet [22], Genuer, Poggi, and Tuleau [70], and Louppe [107].

Yet few attempts have been made to transfer the idea of RFs to one-class classification [52, 105, 156]. In Liu, Ting, and Zhou [105], the novel concept of *isolation* is introduced. The Isolation Forest algorithm isolates anomalies, instead of profiling the inlier behavior which is the usual approach. It avoids adapting splitting rules to the one-class setting by using extremely randomized trees, also named extra trees [71]: isolation trees are built completely randomly, without any splitting rule. Therefore, Isolation Forest is not really based

on RFs, the base estimators being extra trees instead of classical decision trees. Isolation Forest performs very well in practice with low memory and time complexities. In Désir et al. [52] and Shi and Horvath [156], outliers are generated to artificially form a second class. In Désir et al. [52] the authors propose a technique to reduce the number of outliers needed by shrinking the dimension of the input space. The outliers are then generated from the reduced space using a distribution complementary to the inlier distribution. Thus their algorithm artificially generates a second class, to use classical RFs. In Shi and Horvath [156], two different outliers generating processes are compared. In the first one, an artificial second class is created by randomly sampling from the product of empirical marginal –inlier– distributions. In the second one outliers are uniformly generated from the hyper-rectangle that contains the observed data. The first option is claimed to work best in practice, which can be understood from the curse of dimensionality argument: in large dimension [170], when the outliers distribution is not tightly defined around the target set, the chance for an outlier to be in the target set becomes very small, so that a huge number of outliers is needed.

Looking beyond the RF literature, Scott and Nowak [151] proposes a methodology to build dyadic decision trees to estimate minimum-volume sets [57, 137]. This is done by reformulating their structural risk minimization problem to be able to use the algorithm in Blanchard, Schäfer, and Rozenholc [23]. While this methodology can also be used for non-dyadic trees pruning (assuming such a tree has been previously constructed, e. g. using some greedy heuristic), it does not allow to grow such trees. Also, the theoretical guarantees derived there relies on the dyadic structure assumption. In the same spirit, Clémençon and Robbiano [44] proposes to use the two-class splitting criterion defined in Clémençon and Vayatis [45]. This two-class splitting rule aims at producing oriented decision trees with a “left-to-right” structure to address the bipartite ranking task. Extension to the one-class setting is done by assuming a uniform distribution for the outlier class. Consistency and rate bounds relies also on this left-to-right structure. building process to a recursive optimization procedure, thus allowing Thus, these two references [44, 151] impose constraints on the tree structure (designed to allow a statistical study) which differs then significantly from the general structure of the base estimators in RF. The price to pay is the flexibility of the model, and its ability to capture complex broader patterns or structural characteristics from the data.

In this paper, we make the choice to stick to the RF framework. We do not assume any structure for the binary decision trees. The price to pay is the lack of statistical guarantees –the consistency of

RFs has only been proved recently [149] and in the context of regression additive models. The gain is that we preserve the flexibility and strength of RFs, the algorithm presented here being able to compete well with state-of-the-art anomaly detection algorithms. Besides, we do not assume any –fixed in advance– outlier distribution as in Cléménçon and Robbiano [44], but define it in an adaptive way during the tree building process.

To the best of our knowledge, no algorithm structurally extends (without second class sampling and without alternative base estimators) RFs to one-class classification. The main purpose of this work is to introduce such a methodology. It builds on a natural adaptation of two-class Gini-based criterion specially designed for splitting criteria to the one-class setting, as well as an adaptation of the two-class majority vote.

The basic underlying idea is the following. To split a node without second class examples (outliers), we proceed as follows. Each time we look for the best split for a node t , we simply replace (in the two-class *impurity decrease* to be maximized going to the left child node t_L by the proportion expectation $\text{Leb}(t_L)/\text{Leb}(t)$ (idem for the right node), $\text{Leb}(t)$ being the Lebesgue measure, i.e. the volume of the rectangular cell corresponding to node t . It ensures that one child node manages to capture the maximum number of observations with a minimal volume, while the other child looks for the opposite.

This simple idea corresponds to an adaptive modeling of the outlier distribution. The proportion expectation mentioned above is weighted proportionally to the number of inliers in node t . Thus, the resulting outlier distribution is tightly concentrated around the inliers, the latter concentrates outside, closely around but also inside the support of the normal distribution. Besides, and this attests the consistency of our approach with the two-class framework, it turns out that the one-class model promoted here corresponds to the asymptotic behavior of an adaptive outliers generating methodology.

This paper is structured as follows. Appendix D.1 provides the reader with necessary background, to address Appendix D.2 which proposes an adaptation of RFs to the one-class setting and describes a generic one-class random forest algorithm. The latter is compared empirically with state-of-the-art anomaly detection methods in Appendix D.3. Finally a theoretical justification of the one-class criterion is given in Appendix D.4.

D.1 BACKGROUND ON DECISION TREES

Let us denote by $\mathcal{X} \subset \mathbb{R}^d$ the d -dimensional hyper-rectangle containing all the observations. Consider a binary tree on \mathcal{X} whose node values are subsets of \mathcal{X} , iteratively produced by splitting \mathcal{X} into two disjoint subsets. Each internal node t with value \mathcal{X}_t is labeled with a split feature m_t and split value c_t (along that feature), in such a way that it divides \mathcal{X}_t into two disjoint spaces $\mathcal{X}_{t_L} := \{x \in \mathcal{X}_t \mid x_{m_t} < c_t\}$ and $\mathcal{X}_{t_R} := \{x \in \mathcal{X}_t \mid x_{m_t} \geq c_t\}$, where t_L (respectively t_R) denotes the left (respectively right) children of node t , and x_j denotes the j th coordinate of vector x . Such a binary tree is grown from a sample X_1, \dots, X_n ($\forall i \in \mathbb{N}_n^*, X_i \in \mathcal{X}$) and its finite depth is determined either by a fixed maximum depth value or by a stopping criterion evaluated on the nodes (e.g. based on an impurity measure). The external nodes (the *leaves*) form a partition of \mathcal{X} .

In a supervised classification setting, these binary trees are called *classification trees* and prediction is made by assigning to each sample $x \in \mathcal{X}$ the majority class of the leaves containing x . This is called the *majority vote*. Classification trees are usually built using an impurity measure $i(t)$ whose decrease is maximized at each split of a node t , yielding an optimal split (m_t^*, c_t^*) . The decrease of impurity (also called *goodness of split*) $\Delta i(t, t_L, t_R)$ w.r.t. the split (m_t, c_t) and corresponding to the partition $\mathcal{X}_t = \mathcal{X}_{t_L} \sqcup \mathcal{X}_{t_R}$ of the node t is defined as

$$\Delta i(t, t_L, t_R) = i(t) - p_L i(t_L) - p_R i(t_R), \quad (\text{D.1})$$

where $p_L = p_L(t)$ (respectively $p_R = p_R(t)$) is the proportion of samples from \mathcal{X}_t going to \mathcal{X}_{t_L} (respectively to \mathcal{X}_{t_R}). The impurity measure $i(t)$ reflects the goodness of node t : the smaller $i(t)$, the purer the node t and the better the prediction by majority vote on this node. Usual choices for $i(t)$ are the Gini index [72] or the Shannon entropy [153]. To produce a randomized tree, these optimization steps are usually partially randomized (conditionally on the data, splits (m_t^*, c_t^*) 's become random variables). A classification tree can even be grown totally randomly [71]. In a two-class classification setup, the Gini index is

$$i_G(t) = 2 \left(\frac{n_t}{n_t + n'_t} \right) \left(\frac{n'_t}{n_t + n'_t} \right) \quad (\text{D.2})$$

where n_t (respectively n'_t) stands for the number of observations with label 0 (respectively 1) in node t . The Gini index is maximal when $n_t/(n_t + n'_t) = n'_t/(n_t + n'_t) = 0.5$, namely when the conditional probability to have label 0 given that we are in node t is the same as to have label 0 unconditionally: the node t does not discriminate at all between the two classes. For a node t , maximizing the impurity decrease [Equation D.1](#) is equivalent to minimizing $p_L i(t_L) + p_R i(t_R)$.

Since $p_L = (n_{t_L} + n'_{t_L})/(n_t + n'_t)$ and $p_R = (n_{t_R} + n'_{t_R})/(n_t + n'_t)$, and the quantity $(n_t + n'_t)$ being constant in the optimization problem, this is equivalent to minimizing the following proxy of the impurity decrease,

$$I(t_L, t_R) = (n_{t_L} + n'_{t_L})i(t_L) + (n_{t_R} + n'_{t_R})i(t_R). \quad (\text{D.3})$$

Note that with the Gini index $i_G(t)$ given in [Equation D.2](#), the corresponding proxy of the impurity decrease is

$$I_G(t_L, t_R) = \frac{n_{t_L}n'_{t_L}}{n_{t_L} + n'_{t_L}} + \frac{n_{t_R}n'_{t_R}}{n_{t_R} + n'_{t_R}}. \quad (\text{D.4})$$

In the one-class setting, no label is available, hence the impurity measure $i(t)$ does not apply to this setup. The standard splitting criterion which consists in minimizing the latter cannot be used anymore.

D.2 ADAPTATION TO THE ONE-CLASS SETTING

The two reasons why RFs do not apply to one-class classification are that the standard splitting criterion does not apply to this setup, as well as the majority vote. In this section, we propose a one-class splitting criterion and a one-class version of the majority vote.

D.2.1 One-class splitting criterion

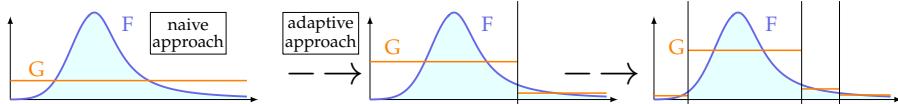


Figure D.1: Outliers distribution G in the naive and adaptive approach. In the naive approach, G does not depend on the tree and is constant on the input space. In the adaptive approach the distribution depends on the inlier distribution F through the tree. The outliers density is constant and equal to the average of F on each node before splitting it.

As one does not observe the second-class (outliers), n'_t needs to be defined. In the naive approach below, it is defined as $n'_t := n' \mathbf{Leb}(\mathcal{X}_t)/\mathbf{Leb}(\mathcal{X})$, where n' is the assumed total number of –hidden– outliers. In the adaptive approach hereafter, it is defined as $n'_t := \gamma n_t$, with typically $\gamma = 1$. Thus, the class ratio $\gamma_t := n'_t/n_t$ is well defined in both approaches and in the naive approach, goes to 0 when $\mathbf{Leb}(\mathcal{X}_t) \rightarrow 0$ while it is maintained constant to γ in the adaptive one.

D.2.1.1 Naive approach

A naive approach to extend the Gini splitting criterion to the one-class setting is to assume a uniform distribution for the second class

(outliers), and to replace their number n'_t in node t by the expectation $n' \mathbf{Leb}(\mathcal{X}_t) / \mathbf{Leb}(\mathcal{X})$, where n' denotes the total number of outliers (for instance, it can be chosen as a proportion of the number of inliers). The problem with this approach appears when the dimension is *not small*. As mentioned in the introduction (curse of dimensionality), when actually generating n' uniform outliers on \mathcal{X} , the probability that a node (sufficiently small to yield a good precision) contains at least one of them is very close to zero. That is why data-dependent distributions for the outlier class are often considered [52, 156]. Taking the expectation $n' \mathbf{Leb}(\mathcal{X}_t) / \mathbf{Leb}(\mathcal{X})$ to replace the number of points in node t does not solve the curse of dimensionality mentioned in the introduction: the volume proportion $L_t := \mathbf{Leb}(\mathcal{X}_t) / \mathbf{Leb}(\mathcal{X})$ is very close to 0 for nodes t deep in the tree, especially in large dimension. In addition, we typically grow trees on sub-samples of the input data, meaning that even the root node of the trees may be very small compared to the hyper-rectangle containing all the input data. An other problem is that the Gini splitting criterion is skew-sensitive [64], and has here to be apply on nodes t with $0 \simeq n'_t \ll n_t$. When trying empirically this approach, we observe that splitting such nodes produces a child containing (almost) all the data (see [Appendix D.4](#)).

Example D.1 *To illustrate the fact that the volume proportion*

$$L_t := \frac{\mathbf{Leb}(\mathcal{X}_t)}{\mathbf{Leb}(\mathcal{X})}$$

becomes very close to zero in large dimension for lots of nodes t (in particular the leaves), suppose for the sake of simplicity that the input space is $\mathcal{X} = [0, 1]^d$. Suppose that we are looking for a rough precision of $1/2^3 = 0.125$ in each dimension, i.e. a unit cube precision of 2^{-3d} . To achieve such a precision, the splitting criterion has to be used on nodes/cells t of volume of order 2^{-3d} , namely with $L_t = 1/2^{3d}$. Note that if we decide to choose n' to be 2^{3d} times larger than the number of inliers in order that $n'L_t$ is not negligible w.r.t. the number of inliers, the same –reversed– problem of unbalanced classes appears on nodes with small depth.

D.2.1.2 Adaptive approach

Our solution is to remove the uniform assumption on the outliers, and to choose their distribution adaptively in such a way it is tightly concentrated around the inlier distribution. Formally, the idea is to maintain constant the class ratio $\gamma_t := n'_t / n_t$ on each node t : before looking for the best split, we update the number of outliers to be equal (up to a scaling constant γ) to the number of inliers, $n'_t = \gamma n_t$, i.e. $\gamma_t \equiv \gamma$. These –hidden– outliers are uniformly distributed on node t . The parameter γ is typically set to $\gamma = 1$, see suppl. [Appendix D.6.1](#) for a discussion on the relevance of this choice (in a nutshell, γ has an influence on optimal splits).

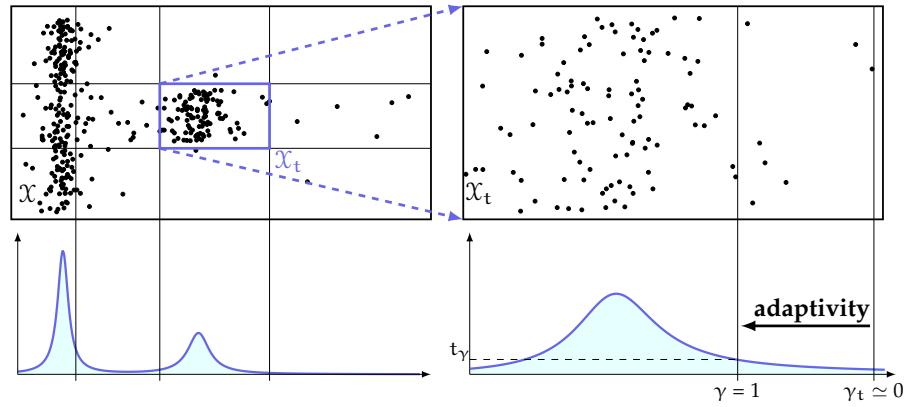


Figure D.2: The left part represents the dataset under study and the underlying density. The node \mathcal{X}_t obtained after some splits is illustrated in the right part of this figure: without the proposed adaptive approach, the class ratio γ_t becomes too small and yields poor splits (all the data are in the “inlier side” of the split, which thus does not discriminate at all). Contrariwise, setting γ to one, i.e. using the adaptive approach, is far preferable.

With this methodology, one cannot derive a one-class version of the Gini index [Equation D.2](#), but we can define a one-class version of the proxy of the impurity decrease [Equation D.4](#), by simply replacing n'_{t_L} (respectively n'_{t_R}) by $n'_t \lambda_L$ (resp. $n'_t \lambda_R$), where $\lambda_L := \text{Leb}(\mathcal{X}_{t_L})/\text{Leb}(\mathcal{X}_t)$ and $\lambda_R := \text{Leb}(\mathcal{X}_{t_R})/\text{Leb}(\mathcal{X}_t)$ are the volume proportion of the two child nodes

$$I_G^{\text{OC-ad}}(t_L, t_R) = \frac{n_{t_L} \gamma n_t \lambda_L}{n_{t_L} + \gamma n_t \lambda_L} + \frac{n_{t_R} \gamma n_t \lambda_R}{n_{t_R} + \gamma n_t \lambda_R}. \quad (\text{D.5})$$

Minimization of the one-class Gini improvement proxy [Equation D.5](#) is illustrated in [Figure D.2](#). Note that $n'_t \lambda_L$ (resp. $n'_t \lambda_R$) is the expectation of the number of uniform observations (on \mathcal{X}_t) among n'_t (fixed to $n'_t = \gamma n_t$) falling into the left (respectively right) node.

Choosing the split minimizing $I_G^{\text{OC-ad}}(t_L, t_R)$ at each step of the tree building process, corresponds to generating $n'_t = \gamma n_t$ outliers each time the best split has to be chosen for node t , and then using the classical two-class Gini proxy [Equation D.4](#). The only difference is that n'_{t_L} and n'_{t_R} are replaced by their expectations $n'_t \lambda_{t_L}$ and $n'_t \lambda_{t_R}$ in our method.

d.2.1.3 Resulting outlier distribution

[Figure D.1](#) shows the corresponding outlier density G (we drop the dependence in the number of splits to keep the notations uncluttered). Note that G is a piece-wise constant approximation of the inlier distribution F . Considering the Neyman-Pearson test $X \sim F$ versus $X \sim G$ instead of $X \sim F$ versus $X \sim \mathcal{U}$ may seem surprising at first sight. Let

us try to give some intuition on why this works in practice. First, there exists (at each step) $\epsilon > 0$ such that $G > \epsilon$ on the entire input space, since the density G is constant on each node and equal to the average of F on this node *before splitting it*. If the average of F was estimated to be zero (no inlier in the node), the node would obviously not have been split, from where the existence of ϵ . Thus, at each step, one can also view G as a piece-wise approximation of $F_\epsilon := (1 - \epsilon)F + \epsilon\mathcal{U}$, which is a mixture of F and the uniform distribution. (ϵ depending on the step/number of splits) Yet, one can easily show that optimal tests for the Neyman-Pearson problem $H_0 : X \sim F$ vs. $H_1 : X \sim F_\epsilon$ are identical to the optimal tests for $H_0 : X \sim F$ vs. $H_1 : X \sim \mathcal{U}$, since the corresponding likelihood ratios are related by a monotone transformation, see Scott and Blanchard [150] for instance (in fact, this reference shows that these two problems are even equivalent in terms of consistency and rates of convergence of the learning rules). An other intuitive justification is as follows. In the first step, the algorithm tries to discriminate F from \mathcal{U} . When going deeper in the tree, splits manage to discriminate F from a (more and more accurate) approximation of F . Asymptotically, splits become irrelevant since they are trying to discriminate F from itself (a perfect approximation, $\epsilon \rightarrow 0$).

Remark D.1 (Consistency with the two-class framework) Consider the following method to generate outliers –tightly concentrated around the support of the inlier distribution. Sample uniformly $n' = \gamma n$ outliers on the rectangular cell containing all the inliers. Split this root node using classical two-class impurity criterion (e.g. minimizing [Equation D.4](#)). Apply recursively the three following steps: for each node t , remove the potential outliers inside \mathcal{X}_t , re-sample $n'_t = \gamma n_t$ uniform outliers on \mathcal{X}_t , and use the latter to find the best split using [Equation D.4](#). Then, each optimization problem [Equation D.4](#) we have solved is equivalent (in expectation) to its one-class version [Equation D.5](#). In other words, by generating outliers adaptively, we can recover (in average) a tree grown using the one-class impurity, from a tree grown using the two-class impurity.

Remark D.2 (Extension to other impurity criteria) Our extension to the one-class setting also applies to other impurity criteria. For instance, in the case of the Shannon entropy defined in the two-class setup by

$$i_S(t) = \frac{n_t}{n_t + n'_t} \log_2 \frac{n_t + n'_t}{n_t} + \frac{n'_t}{n_t + n'_t} \log_2 \frac{n_t + n'_t}{n'_t},$$

the one-class impurity improvement proxy becomes

$$I_S^{OC-ad}(t_L, t_R) = n_{t_L} \log_2 \frac{n_{t_L} + \gamma n_t \lambda_L}{n_{t_L}} + n_{t_R} \log_2 \frac{n_{t_R} + \gamma n_t \lambda_R}{n_{t_R}}.$$

D.2.2 Prediction: scoring function of the forest

Now that RFs can be grown in the one-class setting using the one-class splitting criterion, the forest has to return a prediction adapted to this framework. In other words we also need to extend the concept of majority vote.

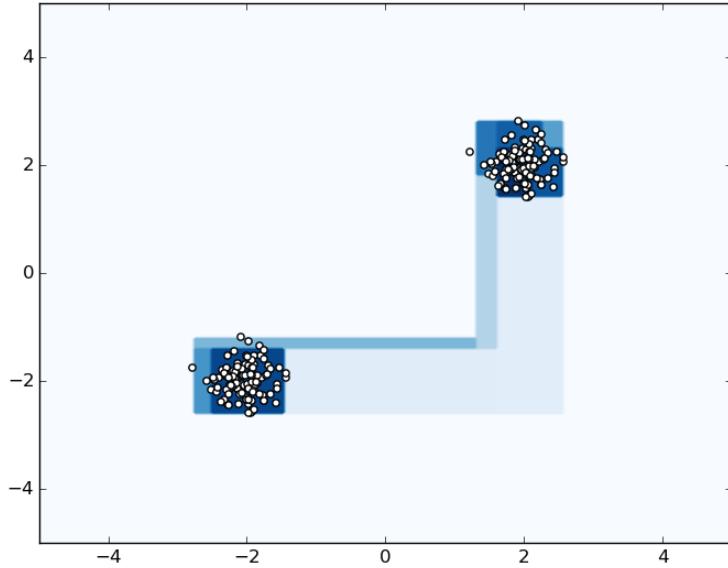


Figure D.3: OneClassRF with one tree, level-sets of the scoring function.

Most usual one-class (or more generally anomaly detection) algorithms actually provide more than just a level-set estimate or a predicted label for any new observation, abnormal versus normal. Instead, they return a real valued function, termed *scoring function*, defining a preorder/ranking on the input space. Such a function $s : \mathbb{R}^d \rightarrow \mathbb{R}$ allows to rank any observations according to their supposed “degree of abnormality”. Thresholding it provides level-set estimates, as well as a decision rule that splits the input space into inlier/normal and outlier/abnormal regions. The scoring function $s(x)$ we use is the one defined in Liu, Ting, and Zhou [105] in view of its established high performance. It is a decreasing function of the average depth of the leaves containing x in the forest. An average term is added to each node containing more than one sample, say containing N samples. This term $c(N)$ is the average depth of an extremely randomized tree [71] (i.e. built without minimizing any criterion, by randomly choosing one feature and one uniform value over this feature to split on) on N samples. Formally,

$$\log_2 s(x) = - \left(\sum_{t \text{ leaves}} \mathbb{1}_{\{x \in t\}} d_t + c(n_t) \right) / c(n), \quad (\text{D.6})$$

where d_t is the depth of node t , and $c(n) = 2H(n - 1) - 2(n - 1)/n$, $H(i)$ being the harmonic number. Alternative scoring functions can be defined for this one-class setting (see [Appendix D.6.2](#)).

D.2.3 OneClassRF: a Generic One-Class Random Forest algorithm

Let us summarize the One Class Random Forest algorithm, based on generic RFs [32]. It has 6 parameters, namely `max_samples`, `max_features_tree`, `max_features_node`, `gamma`, `max_depth`, `n_trees`.

Table D.1: Original datasets characteristics

Datasets	nb of samples	nb of features	anomaly class	
adult	48842	6	class '> 50K'	(23.9%)
annthyroid	7200	6	classes $\neq 3$	(7.42%)
arrhythmia	452	164	classes $\neq 1$ (features 10-14 removed)	(45.8%)
forestcover	286048	10	class 4 (versus class 2)	(0.96%)
http	567498	3	attack	(0.39%)
ionosphere	351	32	bad	(35.9%)
pendigits	10992	16	class 4	(10.4%)
pima	768	8	pos (class 1)	(34.9%)
shuttle	85849	9	classes $\neq 1$ (class 4 removed)	(7.17%)
smtp	95156	3	attack	(0.03%)
spambase	4601	57	spam	(39.4%)
wilt	4839	5	class 'w' (diseased trees)	(5.39%)

Each tree is classically grown on a random subset of both the input samples and the input features [78, 129]. This random subset is a sub-sample of size `max_samples`, with `max_features_tree` variables chosen at random without replacement (replacement is only done after the tree is grown). The tree is built by minimizing [Equation D.5](#) for each split, using parameter γ (recall that $n'_t := \gamma n_t$), until either the maximal depth `max_depth` is achieved or the node contains only one point. Minimizing [Equation D.5](#) is done as introduced in Amit and Geman [7]: at each node, we search the best split over a random selection of features with fixed size `max_features_node`. The forest is composed of a number `n_trees` of trees. The predicted score of a point x is given by $s(x)$, with s defined by [Equation D.6](#). Remarks on alternative stopping criteria and variable importances are available in [Appendix D.6.3](#).

[Figure D.3](#) represents the level sets of the scoring function produced by `ONECLASSRF`, with only one tree of maximal depth 4, without sub-sampling, and using the Gini-based one-class splitting criterion with $\gamma = 1$.

D.3 BENCHMARKS

In this section, we compare the `ONECLASSRF` algorithm described above to seven state-of-art anomaly detection algorithms: the `IFOREST` algorithm [105], a one-class RFs algorithm based on sampling a second class `OCRFsAMPLING` [52], One-Class Support Vector Machine (`OCSM`) [147], Local Outlier Factor (`LOF`) [33], `Orca` [18], Least Squares Anomaly Detection (`LSAD`) [138], Random Forest Clustering (`RFC`) [156].

D.3.1 Default parameters of `OneClassRF`

The default parameters taken for our algorithm are the followings.

- `max_samples` is fixed to 20% of the training sample size (with a minimum of 100);
- `max_features_tree` is fixed to 50% of the total number of features with a minimum of 5 (i.e. each tree is built on 50% of the total number of features);
- `max_features_node` is fixed to 5;
- γ is fixed to 1;
- `max_depth` is fixed to \log_2 (logarithm in base 2) of the training sample size as in Liu, Ting, and Zhou [105];
- `n_trees` is fixed to 100 as in the previous reference.

The other algorithms in the benchmark are trained with their recommended (default) hyper-parameters as seen in their respective paper or author's implementation. See [Appendix D.7](#) for details. The characteristics of the twelve reference datasets considered here are summarized in [Table D.1](#). They are all available on the UCI repository [99] and the preprocessing is done as usually in the literature (see [Appendix D.8](#)).

D.3.2 Results

All the code is available at <https://github.com/ngoix/OCRF>. The experiments are performed in the novelty detection framework, where the training set consists of inliers only. No significance level test are

Table D.3: Results for the novelty detection setting.

Datasets	ONECLASSRF		IFOREST		OCRFsAMPLING		OCSM		LOF		Orca		LSAD		RFC	
	ROC	PR	ROC	PR	ROC	PR	ROC	PR	ROC	PR	ROC	PR	ROC	PR	ROC	PR
AUC																
adult	0.665	0.278	0.661	0.227	N. A.	N. A.	0.638	0.201	0.615	0.188	0.606	0.218	0.647	0.258	N. A.	N. A.
annthyroid	0.936	0.468	0.913	0.456	0.918	0.532	0.706	0.242	0.832	0.446	0.587	0.181	0.810	0.327	N. A.	N. A.
arrhythmia	0.684	0.510	0.763	0.492	0.639	0.249	0.922	0.639	0.761	0.473	0.720	0.466	0.778	0.514	0.716	0.299
forestcover	0.968	0.457	0.863	0.046	N. A.	N. A.	N. A.	N. A.	0.990	0.795	0.946	0.558	0.952	0.166	N. A.	N. A.
http	0.999	0.838	0.994	0.197	N. A.	0.999	0.812	0.981	0.537	N. A.	N. A.					
ionosphere	0.909	0.643	0.902	0.535	0.859	0.609	0.973	0.849	0.959	0.807	0.928	0.910	0.978	0.893	0.950	0.754
pendigits	0.960	0.559	0.810	0.197	0.968	0.694	0.603	0.110	0.983	0.827	0.993	0.925	0.983	0.752	N. A.	N. A.
pima	0.719	0.247	0.726	0.183	0.759	0.266	0.716	0.237	0.700	0.152	0.588	0.175	0.713	0.216	0.506	0.090
shuttle	0.999	0.998	0.996	0.973	N. A.	N. A.	0.992	0.924	0.999	0.995	0.890	0.782	0.996	0.956	N. A.	N. A.
smtp	0.922	0.499	0.907	0.005	N. A.	N. A.	0.881	0.656	0.924	0.149	0.782	0.142	0.877	0.381	N. A.	N. A.
spambase	0.850	0.373	0.824	0.372	0.797	0.485	0.737	0.208	0.746	0.160	0.631	0.252	0.806	0.330	0.723	0.151
wilt	0.593	0.070	0.491	0.045	0.442	0.038	0.323	0.036	0.697	0.092	0.441	0.030	0.677	0.074	0.896	0.631
average	0.850	0.495	0.821	0.311	0.769	0.410	0.749	0.410	0.837	0.462	0.759	0.454	0.850	0.450	0.758	0.385
cum. train time	61s		68s		N. A.		N. A.		N. A.		2232s		73s		N. A.	

given, but experiments for each algorithm are repeated 10 times on random training and testing datasets are performed, yielding averaged ROC and PR curves whose AUCs are summarized in [Table D.3](#) (higher is better). The training time of each algorithm has been limited (for each experiment among the 10 performed for each dataset) to 30 minutes, where N. A. indicates that the algorithm could not finish training within the allowed time limit. In average on all the datasets, our proposed algorithm **ONECLASSRF** achieves both best AUC ROC and AUC PR scores (with LSAD for AUC ROC). It also achieves the lowest cumulative training time. For further insights on the benchmarks c.f. [Appendix D.6](#). It appears that **ONECLASSRF** has the best performance on five datasets in terms of ROC AUCs, and is also the best in average. Computation times (training plus testing) of **ONECLASSRF** are also very competitive. Experiments in an outlier detection framework (the training set is polluted by outliers) have also been made (see [Appendix D.9](#)). The anomaly rate is arbitrarily bounded to 10% max (before splitting data into training and testing sets).

D.4 THEORETICAL ANALYSIS

This section aims at recovering [Equation D.5](#) from a natural modeling of the one-class framework, along with a theoretical study of the problem raised by the naive approach.

D.4.1 *Underlying model*

In order to generalize the two-class framework to the one-class one, we need to consider the population versions associated to empirical quantities [Equation D.1](#), [Equation D.2](#) and [Equation D.3](#), as well as the underlying model assumption. The latter can be described as follows.

D.4.1.1 *Existing Two-Class Model (n, α)*.

We consider a r.v. $X : \Omega \rightarrow \mathbb{R}^d$ w.r.t. a probability space $(\Omega, \mathcal{F}, \Pr)$. The law of X depends on another r.v. $y \in \{0, 1\}$, verifying $\Pr\{y = 1\} = 1 - \Pr\{y = 0\} = \alpha$. We assume that conditionally on $y = 0$, X follows a law F , and conditionally on $y = 1$ a law G ;

$$\begin{aligned} X \mid y = 0 &\sim F, & \Pr\{y = 0\} &= 1 - \alpha, \\ X \mid y = 1 &\sim G, & \Pr\{y = 1\} &= \alpha. \end{aligned}$$

Then, considering

$$p(t_L|t) = \Pr\{X \in \mathcal{X}_{t_L} \mid X \in \mathcal{X}_t\},$$

and

$$p(t_R|t) = \Pr\{X \in \mathcal{X}_{t_R} \mid X \in \mathcal{X}_t\},$$

the population version (probabilistic version) of [Equation D.1](#) is

$$\Delta i^{theo}(t, t_L, t_R) = i^{theo}(t) - p(t_L|t)i^{theo}(t_L) - p(t_R|t)i^{theo}(t_R)$$

It can be used with the Gini index i_G^{theo} ,

$$i_G^{theo}(t) = 2\Pr\{y = 0 | X \in \mathcal{X}_t\}\Pr\{y = 1 | X \in \mathcal{X}_t\} \quad (\text{D.8})$$

which is the population version of [Equation D.2](#).

D.4.1.2 One-Class-Model (n, α).

We model the one-class framework as follows. Among the n i.i.d. observations, we only observe those with $y = 0$ (the inliers), namely N realizations of $(X | y = 0)$, where N is itself a realization of a r.v. N of law $N \sim \text{Bin}(n, (1 - \alpha))$. Here and hereafter, $\text{Bin}(n, p)$ denotes the binomial distribution with parameters (n, p) . As outliers are not observed, it is natural to assume that G follows a uniform distribution on the hyper-rectangle \mathcal{X} containing all the observations, so that G has a constant density $g(\cdot) = 1/\text{Leb}(\mathcal{X})$ on \mathcal{X} . Note that this assumption *will be removed* in the adaptive approach described below – which aims at maintaining a non-negligible proportion of (hidden) outliers in every nodes.

Let us define $L_t = \text{Leb}(\mathcal{X}_t)/\text{Leb}(\mathcal{X})$. Then, $\Pr\{X \in \mathcal{X}_t | y = 1\} = \Pr\{y = 1\}\Pr\{X \in \mathcal{X}_t | y = 1\} = \alpha L_t$. Replacing $\Pr\{X \in \mathcal{X}_t | y = 0\}$ by its empirical version n_t/n in [Equation D.8](#), we obtain the one-class empirical Gini index

$$i_G^{OC}(t) = \frac{n_t \alpha n L_t}{(n_t + \alpha n L_t)^2}. \quad (\text{D.9})$$

This one-class index can be seen as a *semi-empirical* version of [Equation D.8](#), in the sense that it is obtained by considering empirical quantities for the (observed) inlier behavior and population quantities for the (non-observed) outlier behavior. Now, maximizing the population version of the impurity decrease $\Delta i_G^{theo}(t, t_L, t_R)$ as defined in [Equation D.7](#) is equivalent to minimizing

$$p(t_L|t)i_G^{theo}(t_L) + p(t_R|t)i_G^{theo}(t_R). \quad (\text{D.10})$$

Considering semi-empirical versions of $p(t_L|t)$ and $p(t_R|t)$, as for [Equation D.9](#), gives $p_n(t_L|t) = (n_{t_L} + \alpha n L_{t_L})/(n_t + \alpha n L_t)$ and $p_n(t_R|t) = (n_{t_R} + \alpha n L_{t_R})/(n_t + \alpha n L_t)$. Then, the semi-empirical version of [Equation D.10](#) is

$$p_n(t_L|t)i_G^{OC}(t_L) + p_n(t_R|t)i_G^{OC}(t_R) = \frac{1}{(n_t + \alpha n L_t)} \left(\frac{n_{t_L} \alpha n L_{t_L}}{n_{t_L} + \alpha n L_{t_L}} + \frac{n_{t_R} \alpha n L_{t_R}}{n_{t_R} + \alpha n L_{t_R}} \right) \quad (\text{D.11})$$

where $1/(n_t + \alpha n L_t)$ is constant when the split varies. This means that finding the split minimizing [Equation D.11](#) is equivalent to finding the split minimizing

$$I_G^{OC}(t_L, t_R) = \frac{n_{t_L} \alpha n L_{t_L}}{n_{t_L} + \alpha n L_{t_L}} + \frac{n_{t_R} \alpha n L_{t_R}}{n_{t_R} + \alpha n L_{t_R}}. \quad (\text{D.12})$$

Note that [Equation D.12](#) can be obtained from the two-class impurity decrease [Equation D.4](#) as described in the naive approach paragraph in [Appendix D.2](#). In other words, it is the naive one-class version of [Equation D.4](#).

Remark D.3 (Direct link with the two-class framework). *The two-class proxy of the Gini impurity decrease [Equation D.4](#) is recovered from [Equation D.12](#) by replacing $\alpha n L_{t_L}$ (resp. $\alpha n L_{t_R}$) by n'_{t_L} (respectively n'_{t_R}), the number of second class instances in t_L (respectively in t_R). When generating αn of them uniformly on \mathcal{X} , $\alpha n L_t$ is the expectation of n'_t .*

As detailed in [Appendix D.2.1](#), this approach suffers from the curse of dimensionality. We can summarize the problem as follows. Note that when setting $n'_t := \alpha n L_t$, the class ratio $\gamma_t = n'_t / n_t$ is then equal to

$$\gamma_t = \alpha n L_t / n_t. \quad (\text{D.13})$$

This class ratio is close to 0 for lots of nodes t , which makes the Gini criterion unable to discriminate accurately between the -hidden- outliers and the inliers. Minimizing this criterion produces splits corresponding to $\gamma_t \simeq 0$ in [Figure D.2](#): one of the two child nodes, say t_L contains almost all the data.

D.4.2 Adaptive approach

The solution presented [Appendix D.2](#) is to remove the uniform assumption for the outlier class. From the theoretical point of view, the idea is to choose in an adaptive way (w.r.t. the volume of \mathcal{X}_t) the number αn , which can be interpreted as the number of (hidden) outliers. α . Doing so, we aim at avoiding $\alpha n L_t \ll n_t$ when L_t is too small. Namely, with γ_t defined in [Equation D.13](#), we aim at avoiding $\gamma_t \simeq 0$ when $L_t \simeq 0$. The idea is to consider $\alpha(L_t)$ and $n(L_t)$ such that $\alpha(L_t) \rightarrow 1$, $n(L_t) \rightarrow \infty$ when $L_t \rightarrow 0$. We then define the one-class adaptive proxy of the impurity decrease by

$$I_G^{OC-\text{ad}}(t_L, t_R) = \frac{n_{t_L} \alpha(L_t) n(L_t) L_{t_L}}{n_{t_L} + \alpha(L_t) n(L_t) L_{t_L}} + \frac{n_{t_R} \alpha(L_t) n(L_t) L_{t_R}}{n_{t_R} + \alpha(L_t) \cdot n(L_t) \cdot L_{t_R}}. \quad (\text{D.14})$$

In other words, instead of considering one general model One-Class-Model(n , α) defined in [Appendix D.4.1](#), we adapt it to each node

t , considering One-Class-Model($n(L_t)$, $\alpha(L_t)$) before searching the best split. We still consider the N inliers as a realization of this model. When growing the tree, using One-Class-Model($n(L_t)$, $\alpha(L_t)$) allows to maintain a non-negligible expected proportion of outliers in the node to be split, despite L_t becomes close to zero. Of course, constraints have to be imposed to ensure consistency between these models. Recalling that the number N of inliers is a realization of N following a Binomial distribution with parameters $(n, 1 - \alpha)$, a first natural constraint on $(n(L_t), \alpha(L_t))$ is

$$(1 - \alpha)n = (1 - \alpha(L_t)) \cdot n(L_t), \quad \text{for all } t, \quad (\text{D.15})$$

so that the expectation of N remains unchanged.

Remark D.4 In our adaptive model One-Class-Model($n(L_t)$, $\alpha(L_t)$) which varies when we grow the tree, let us denote by $N(L_t) \sim \text{Bin}(n(L_t), 1 - \alpha(L_t))$ the r.v. ruling the number of inliers. The number of inliers N is still viewed as a realization of it. Note that the distribution of $N(L_t)$ converges in distribution to $\mathcal{P}((1 - \alpha)n)$ a Poisson distribution with parameter $(1 - \alpha)n$ when $L_t \rightarrow 0$, while the distribution $\text{Bin}(n(L_t), \alpha(L_t))$ of the r.v. $n(L_t) - N(L_t)$ ruling the number of (hidden) outliers goes to infinity almost surely. In other words, the asymptotic model (when $L_t \rightarrow 0$) consists in assuming that the number of inliers N we observed is a realization of $N_\infty \sim \mathcal{P}((1 - \alpha)n)$, and that an infinite number of outliers have been hidden.

A second natural constraint on $(\alpha(L_t), n(L_t))$ is related to the class ratio γ_t . As explained in [Appendix D.2.1](#), we do not want γ_t to go to zero when L_t does. Let us say we want γ_t to be constant for all node t , equal to $\gamma > 0$. From the constraint $\gamma_t = \gamma$ and [Equation D.13](#), we get

$$\alpha(L_t)n(L_t)L_t = \gamma n_t := n'_t. \quad (\text{D.16})$$

The constant γ is a parameter ruling the expected proportion of outliers in each node. Typically, $\gamma = 1$ so that there is as much expected uniform (hidden) outliers than inliers at each time we want to find the best split minimizing [Equation D.14](#). [Equation D.15](#) and [Equation D.16](#) allow to explicitly determine $\alpha(L_t)$ and $n(L_t)$: $\alpha(L_t) = n'_t / ((1 - \alpha)nL_t + n'_t)$ and $n(L_t) = ((1 - \alpha)nL_t + n'_t) / L_t$. Regarding [Equation D.14](#), $\alpha(L_t)n(L_t)L_{t_L} = \frac{n'_t}{L_t}L_{t_L} = n'_t \frac{\text{Leb}(\mathcal{X}_{t_L})}{\text{Leb}(\mathcal{X}_t)}$ by [Equation D.16](#) and $\alpha(L_t)n(L_t)L_{t_R} = n'_t \frac{\text{Leb}(\mathcal{X}_{t_R})}{\text{Leb}(\mathcal{X}_t)}$, so that we recover [Equation D.5](#).

D.5 CONCLUSION

Through a natural adaptation of both (two-class) splitting criteria and majority vote, this paper introduces a methodology to structurally

extend RFs to the one-class setting. Our one-class splitting criteria correspond to the asymptotic behavior of an adaptive outliers generating methodology, so that consistency with two-class RFs seems respected. While no statistical guaranties have been derived in this paper, a strong empirical performance attests the relevance of this methodology.

D.6 FURTHER INSIGHTS ON THE ALGORITHM

D.6.1 Interpretation of parameter γ

In order for the splitting criterion [Equation D.5](#) to perform well, n'_t is expected to be of the same order of magnitude as the number of inliers n_t . If $\gamma = n'_t/n_t \ll 1$, the split puts every inliers on the same side, even the ones which are far in the tail of the distribution, thus widely over-estimating the support of inliers. If $\gamma \gg 1$, the opposite effect happens, yielding an estimate of a t -level set with t not close enough to 0. [Figure D.2](#) illustrates the splitting criterion when γ varies. It clearly shows that there is a link between parameter γ and the level t_γ of the induced level-set estimate. But from the theory, an explicit relation between γ and t_γ is hard to derive. By default we set γ to 1. One could object that in some situations, it is useful to randomize this parameter. For instance, in the case of a bi-modal distribution for the inlier/normal behavior, one split of the tree needs to separate two clusters, in order for the level set estimate to distinguish between the two modes. As illustrated in [Figure D.4](#), it can only occur if n'_t is large with respect to n_t ($\gamma \gg 1$). However, the randomization of γ is somehow included in the randomization of each tree, thanks to the sub-sampling inherent to RFs. Moreover, small clusters tend to vanish when the sub-sample size is sufficiently small: a small sub-sampling size is used by Liu, Ting, and Zhou [105] to isolate outliers even when they form clusters.

D.6.2 Alternative scoring functions

Although we use the scoring function defined in [Equation D.6](#) because of its established high performance [105], other scoring functions can be defined. A natural idea to adapt the majority vote to the one-class setting is to change the single vote of a leaf node t into the fraction $\frac{n_t}{\text{Leb}(\mathcal{X}_t)}$, the forest output being the average of the latter quantity over the forest, $s(x) = \sum_{t \text{ leaves}} \mathbb{1}_{\{x \in t\}} \frac{n_t}{\text{Leb}(\mathcal{X}_t)}$. In such a case, each tree of the forest yields a piece-wise density estimate on its induced partition. The output produced by the forest is then a *step-wise density estimate*. We could also think about the *local density of a typical cell*. For each point x of the input space, it returns the average number of observations in the leaves containing x , divided by the average

volume of such leaves. The output of **OneClassRF** is then the scoring function $s(x) = \left(\sum_{t \text{ leaves}} \mathbb{1}_{\{x \in t\}} n_t \right) \left(\sum_{t \text{ leaves}} \mathbb{1}_{\{x \in t\}} \text{Leb}(\mathcal{X}_t) \right)^{-1}$, where the sums are over each leave of each tree in the forest. This score can be interpreted as the local density of a “typical” cell (typical among those usually containing x).

D.6.3 Alternative stopping criteria

Other stopping criteria than a maximal depth may be considered. We could stop splitting a node t when it contains less than n_min observations, or when the quantity $n_t / \text{Leb}(\mathcal{X}_t)$ is large enough (all the points in the cell \mathcal{X}_t are likely to be inliers) or close enough to 0 (all the points in the cell \mathcal{X}_t are likely to be outliers). These options are not discussed in this work.

D.6.4 Variable importance

In the multiclass setting, Breiman [32] proposed to evaluate the importance of a feature $j \in \{1, \dots, d\}$ for prediction by adding up the weighted impurity decreases for all nodes t where X_j is used, averaged over all the trees. The analogue quantity can be computed with respect to the one-class impurity decrease proxy. In our one-class setting, this quantity represents the size of the tail of X_j , and can be interpreted as the capacity of feature j to discriminate between inliers/outliers.

D.7 HYPER-PARAMETERS OF TESTED ALGORITHMS

Overall we chose to train the different algorithms with their (default) hyper-parameters as seen in their respective paper or author’s implementation. Indeed, since we are in an unsupervised setting, there is no trivial way to select/learn the hyperparameters of the different algorithm in the training phase – the labels are not supposed to be available. Hence the more realistic way to test the algorithms is to use their recommended/default hyperparameters.

The OCSM algorithm uses default parameters: `kernel='rbf'` with `tol=1e-3`, `nu=0.5`, `shrinking=True` and `gamma=1/n_features`, where `tol` is the tolerance for stopping criterion.

The LOF algorithm uses default parameters: `n_neighbors=5` with the `leaf_size=30` and `metric='minkowski'` and `contamination=0.1` and `algorithm='auto'`, where the algorithm parameters stipulates how to compute the nearest neighbors (either ball-tree, kd-tree or brute-force).

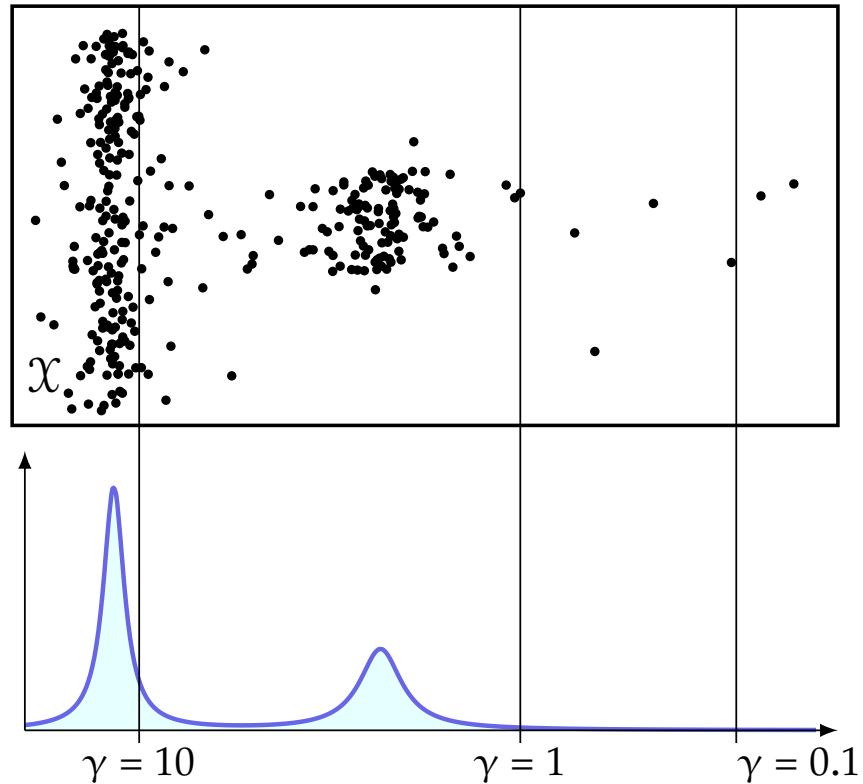


Figure D.4: Illustration of the standard splitting criterion on two modes when the proportion γ varies.

The IFOREST algorithm uses default parameters: `n_estimators=100` and `max_samples=min(256, n_samples)` and `max_features=1` and setting `bootstrap=false`, where `bootstrap` states whether samples are drawn with replacement.

The OCRFsAMPLING algorithm uses default parameters: the number of dimensions for the Random Subspace Method `krsm=-1`, the number of features randomly selected at each node during the induction of the tree `krfs=-1`, `n_tree=100`, the factor controlling the extension of the outlier domain used to sample outliers according to the volume of the hyper-box surrounding the target data `alpha=1.2`, the factor controlling the number of outlier data generated according to the number of target data `beta=10`, whether outliers are generated from uniform distribution `optimize=0` and eventually whether data outside target bounds are considered as outlier data `rejectOutOfBounds=0`.

The *Orca* algorithm uses default parameter `k=5` (number of nearest neighbors) as well as `N=n/8` (how many anomalies are to be reported). The last setting, set up in the empirical evaluation of iForest in Liu,

Ting, and Zhou [104], allows a better computation time without impacting Orca’s performance.

The RFC algorithm uses default parameters: `no.forest=25` with the number of trees `no.trees=3000`, the `Addcl1` Random Forest dissimilarity `addcl1=T`, `addcl2=F` use the importance measure `imp=T`, the data generating process `oob.prox1=T`, the number of features sampled at each split `mtry1=3`.

The LSAD algorithm uses default parameters: the maximum number of samples per kernel `n_kernels_max=500`, the center of each kernel (the center of the random sample subset by default) `kernel_pos='None'`, the kernel scale parameter (using the pairwise median trick by default [80]) `gamma='None'`, the regularization parameter `rho=0.1`.

D.8 DESCRIPTION OF THE DATASETS

The characteristics of the twelve reference datasets considered here are summarized in Table D.1. They are all available on the UCI repository [99] and the preprocessing is done in a classical way. In anomaly detection, we typically have data from two class (inliers/outliers) – in novelty detection, the second class is unavailable in training in outlier detection, training data are polluted by second class (anonymous) examples. The classical approach to adapt multi-class data to this framework is to set classes forming the outlier class, while the other classes form the inlier class.

We removed all categorial attributes. Indeed, our method is designed to handle data whose distribution is absolutely continuous w.r.t. the Lebesgue measure. The `http` and `smtp` datasets belong to the KDD Cup ‘99 dataset [82, 169], which consist of a wide variety of hand-injected attacks (anomalies) in a closed network (normal/inlier background). They are classically obtained as described in Yamanishi et al. [187]. These two datasets are available on the *scikit-learn* library [132]. The `shuttle` dataset is the fusion of the training and testing datasets available in the UCI repository. As in Liu, Ting, and Zhou [105], we use instances from all different classes but class 4. In the `forestcover` data, the inliers are the instances from class 2 while instances from class 4 are anomalies (as in Liu, Ting, and Zhou [105]). The `ionosphere` dataset differentiates “good” from “bad” radars, considered here as abnormal. A “good” radar shows evidence of some type of structure in the ionosphere. A “bad” radar does not, its signal passing through the ionosphere. The `spambase` dataset consists of spam or non-spam emails. The former constitute our anomaly class. The `annthyroid` medical dataset on hypothyroidism contains one nor-

mal class and two abnormal ones, which form our outliers. The *arrhythmia* dataset reflects the presence and absence (class 1) of cardiac arrhythmia. The number of attributes being large considering the sample size, we removed attributes containing missing data. Besides, we removed attributes taking less than 10 different values, the latter breaking too strongly our absolutely continuous assumption (w.r.t. to **Leb**). The *pendigits* dataset contains 10 classes corresponding to the digits from 0 to 9, examples being handwriting samples. As in Schubert et al. [148], the outliers are chosen to be those from class 4. The *pima* dataset consists of medical data on diabetes. Patients suffering from diabetes (inlier class) were considered outliers. The *wild* dataset involves detecting diseased trees in Quickbird imagery. Diseased trees (class ‘w’) is our outlier class. In the *adult* dataset, the goal is to predict whether income exceeds \$ 50K/year based on census data. We only keep the 6 continuous attributes.

D.9 FURTHER DETAILS ON BENCHMARKS AND OUTLIER DETECTION RESULTS

[Figure D.5](#) shows that the amount of time to train¹ and test any dataset takes less than one minute with **ONECLASSRF**, whereas some algorithms have far higher computation times (OCRFSAMPLING, OCSM, LOF and Orca have computation times higher than 30 minutes in some datasets). Our approach yields results similar to quite new algorithms such as IFOREST and LSAD.

In this section present experiments in the outlier detections setting. For each algorithm, 10 experiments on random training and testing datasets are performed. Averaged ROC and PR curves AUC are summarized in [Table D.4](#). For the experiments made in an unsupervised framework (meaning that the training set is polluted by outliers), the anomaly rate is arbitrarily bounded to 10% max (before splitting data into training and testing sets).

¹ For **ONECLASSRF**, Orca and RFC, testing and training time cannot be isolated because of algorithms implementation: for these algorithms, the sum of the training and testing times are displayed in [Figure D.5](#) and [Figure D.6](#).

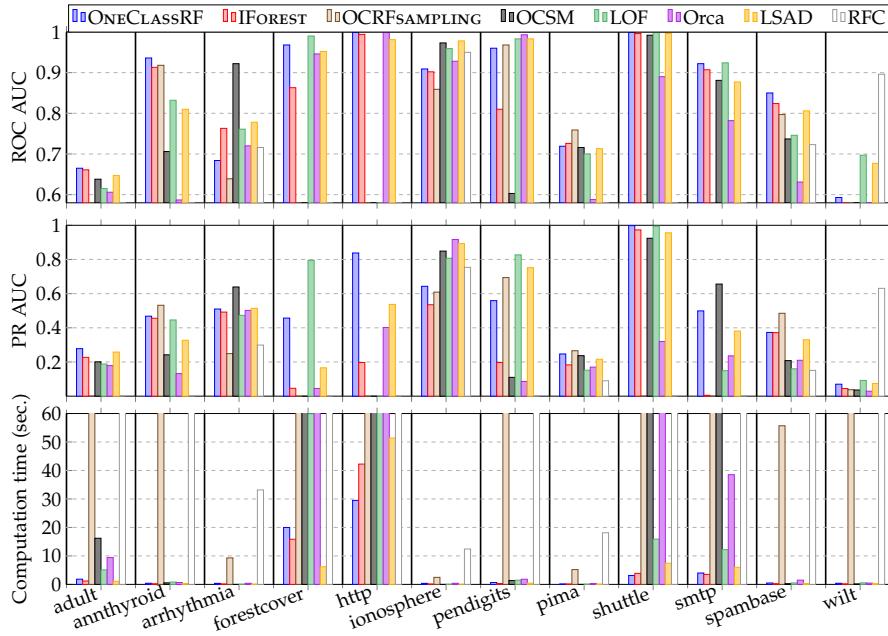


Figure D.5: Performances of the algorithms on each dataset in the novelty detection framework: ROC AUCs are displayed on the top, PR AUCs in the middle and training times on the bottom, for each dataset and algorithm. The x-axis represents the datasets.

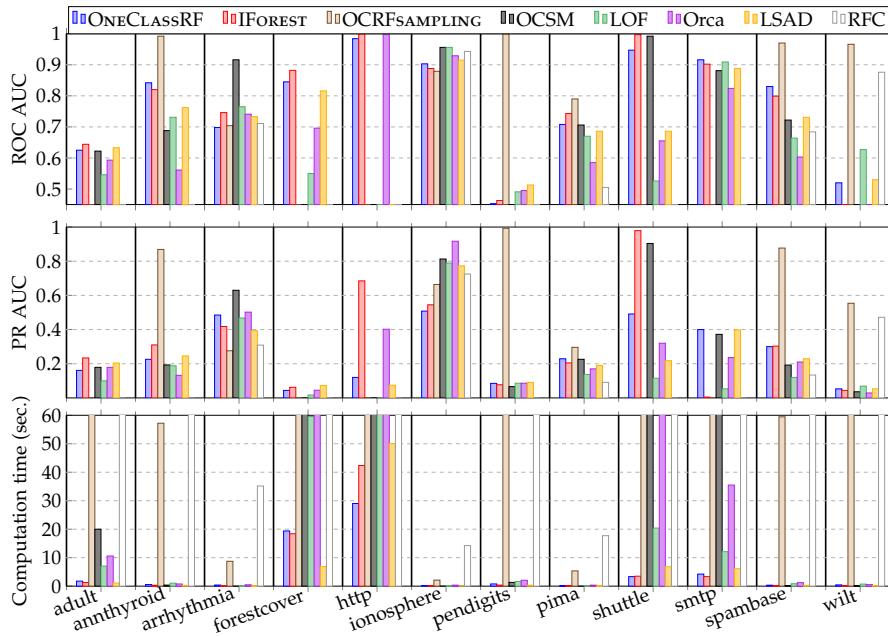


Figure D.6: Performances of the algorithms on each dataset in the outlier detection framework: ROC AUCs are on the top, PR AUCs in the middle and processing times are displayed below (for each dataset and algorithm). The x-axis represents the datasets.

Figure D.7: ROC and PR curves for OneClassRF (novelty detection framework)

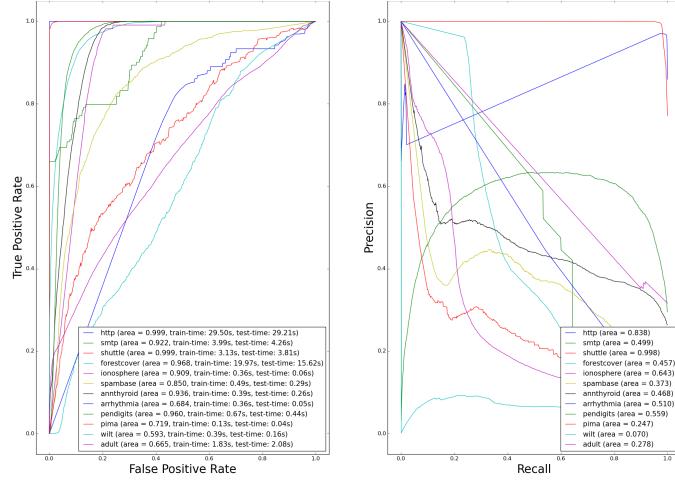


Figure D.8: ROC and PR curves for OneClassRF (outlier detection framework)

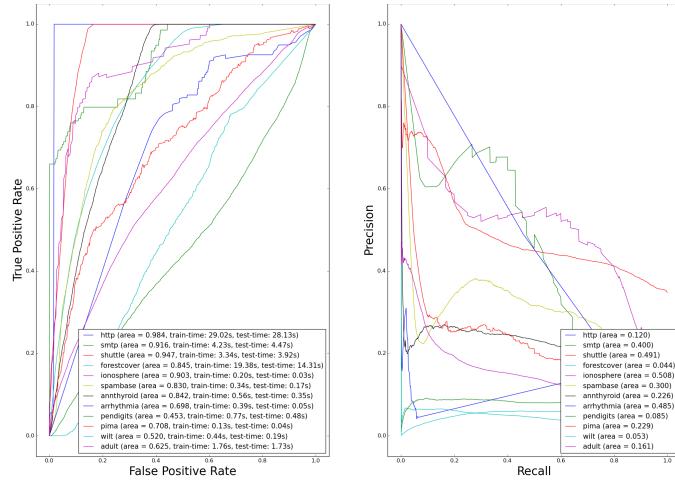


Figure D.9: ROC and PR curves for IFOREST (novelty detection framework)

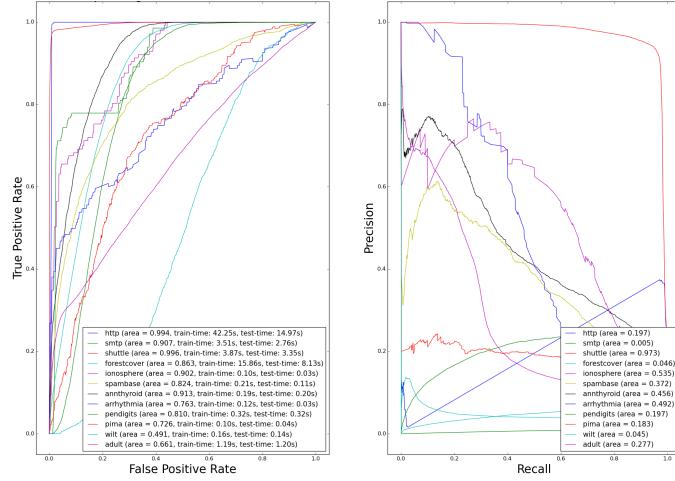


Figure D.10: ROC and PR curves for IFOREST (outlier detection framework)

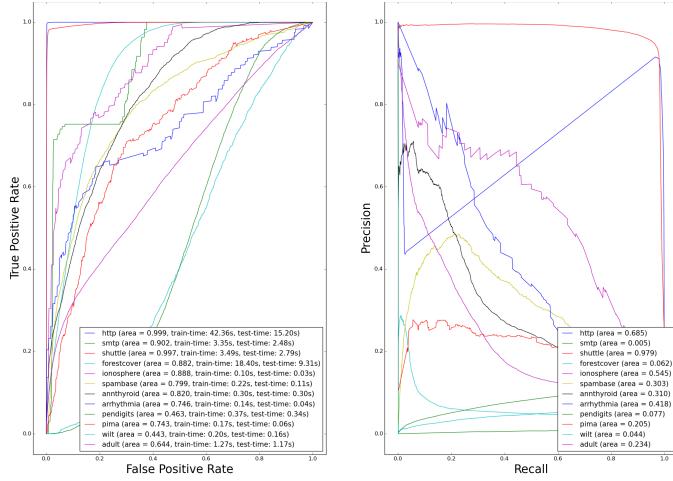


Figure D.11: ROC and PR curves for OCRCFSAMPLING (novelty detection framework)

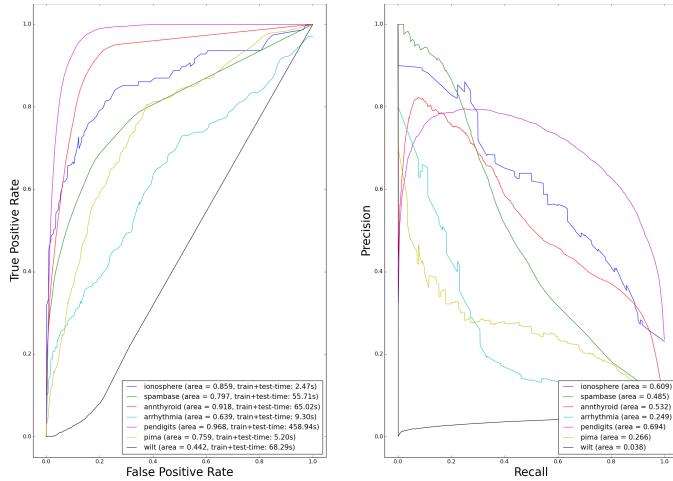


Figure D.12: ROC and PR curves for OCRCFSAMPLING (outlier detection framework)

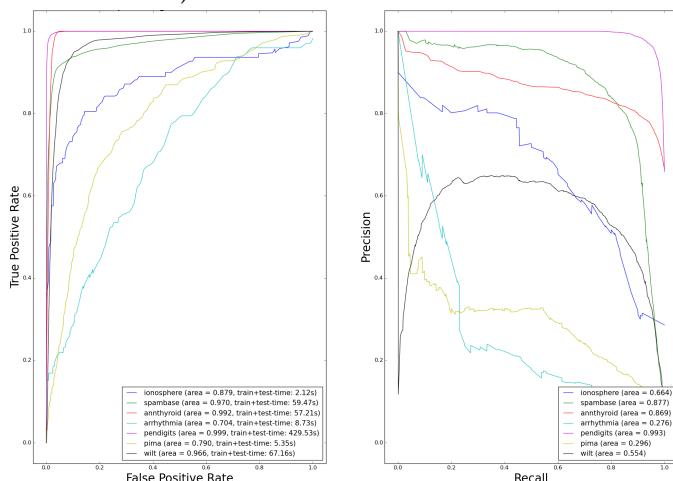


Figure D.13: ROC and PR curves for OCSM (novelty detection framework)

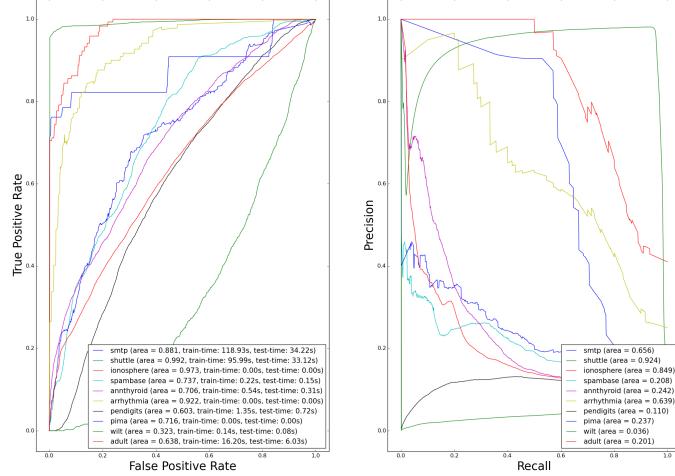


Figure D.14: ROC and PR curves for OCSM (outlier detection framework)

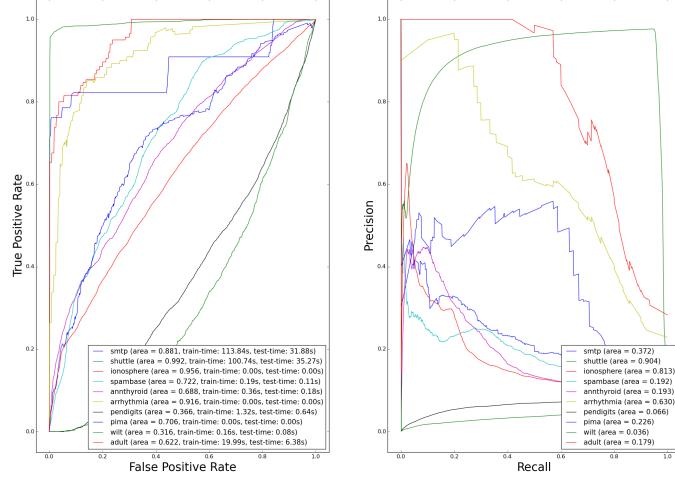


Figure D.15: ROC and PR curves for LOF (novelty detection framework)

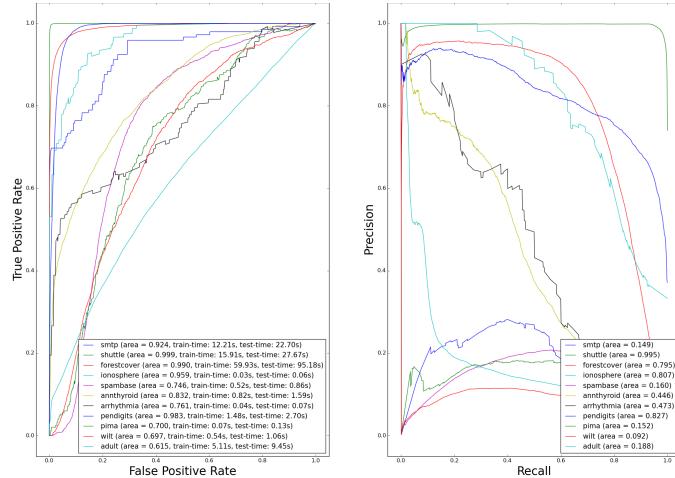


Figure D.16: ROC and PR curves for LOF (outlier detection framework)

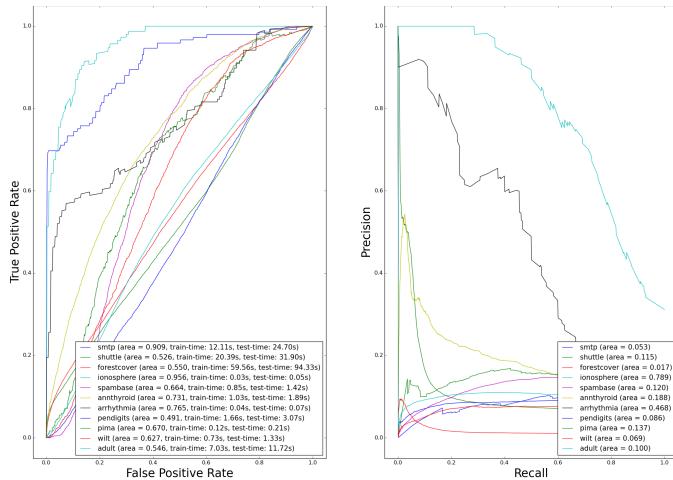


Figure D.17: ROC and PR curves for Orca (novelty detection framework)

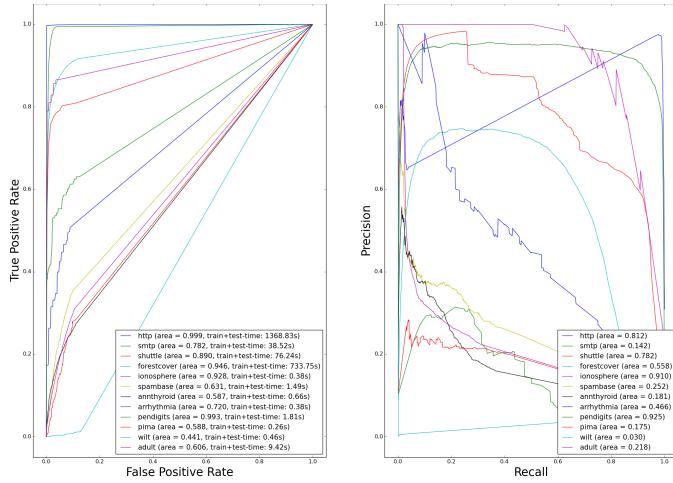


Figure D.18: ROC and PR curves for Orca (outlier detection framework)

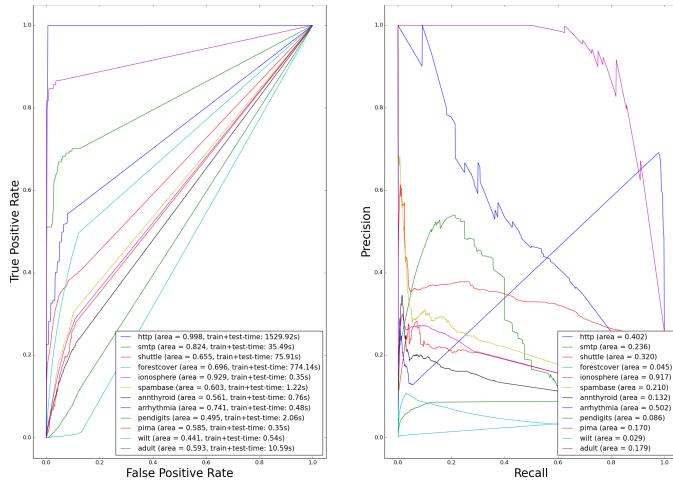


Figure D.19: ROC and PR curves for LSAD (novelty detection framework)

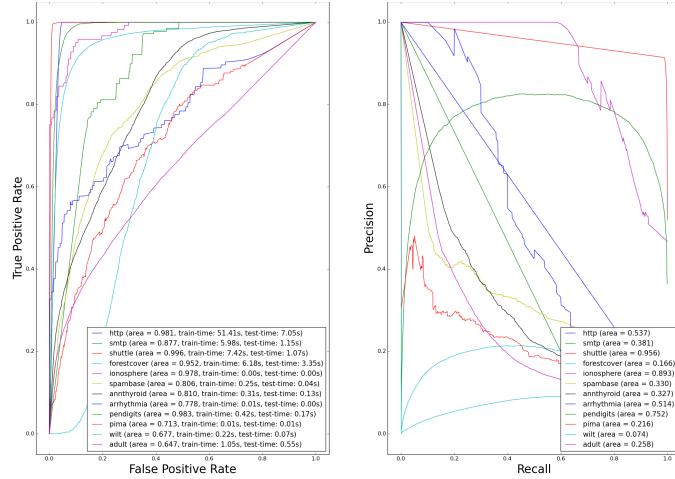


Figure D.20: ROC and PR curves for LSAD (outlier detection framework)

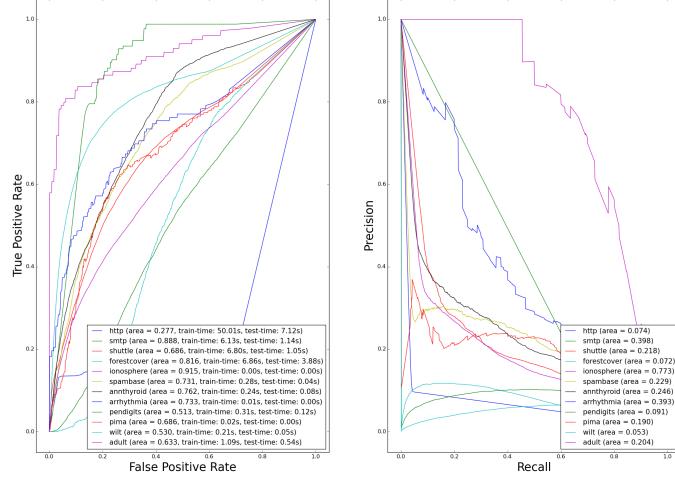


Figure D.21: ROC and PR curves for RFC (novelty detection framework)

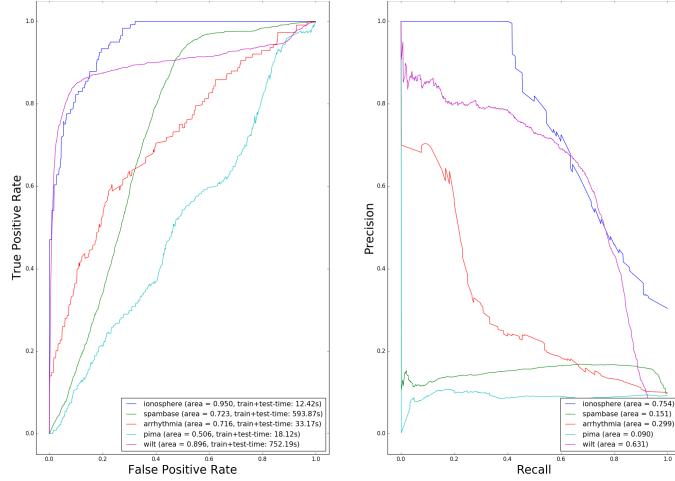


Figure D.22: ROC and PR curves for RFC (outlier detection framework)

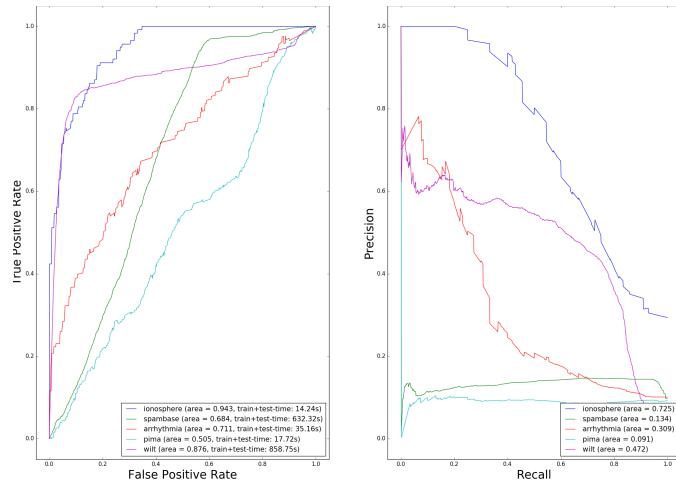


Table D.4: Results for the outlier detection setting

Dataset	ONECLASSRF		IFOREST		OCRFsAMPLING		OCSM		LOF		Orca		LSAD		RFC	
	ROC	PR	ROC	PR	ROC	PR	ROC	PR								
AUC																
adult	0.625	0.161	0.644	0.234	N.A.	N.A.	0.622	0.179	0.546	0.100	0.593	0.179	0.633	0.204	N.A.	N.A.
annthyroid	0.842	0.226	0.820	0.310	0.992	0.869	0.688	0.193	0.731	0.188	0.561	0.132	0.762	0.246	N.A.	N.A.
arrhythmia	0.698	0.485	0.746	0.418	0.704	0.276	0.916	0.630	0.765	0.468	0.741	0.502	0.733	0.393	0.711	0.309
forestcover	0.845	0.044	0.882	0.062	N.A.	N.A.	N.A.	N.A.	0.550	0.017	0.696	0.045	0.816	0.072	N.A.	N.A.
http	0.984	0.120	0.999	0.685	N.A.	N.A.	N.A.	N.A.	N.A.	0.998	0.402	0.277	0.074	N.A.	N.A.	
ionosphere	0.903	0.508	0.888	0.545	0.879	0.664	0.956	0.813	0.956	0.789	0.929	0.917	0.915	0.773	0.943	0.725
pendigits	0.453	0.085	0.463	0.077	0.999	0.993	0.366	0.066	0.491	0.086	0.495	0.086	0.513	0.091	N.A.	N.A.
pima	0.708	0.229	0.743	0.205	0.790	0.296	0.706	0.226	0.670	0.137	0.585	0.170	0.686	0.190	0.505	0.091
shuttle	0.947	0.491	0.997	0.979	N.A.	N.A.	0.992	0.904	0.526	0.115	0.655	0.320	0.686	0.218	N.A.	N.A.
smtp	0.916	0.400	0.902	0.005	N.A.	N.A.	0.881	0.372	0.909	0.053	0.824	0.236	0.888	0.398	N.A.	N.A.
spambase	0.830	0.300	0.799	0.303	0.970	0.877	0.722	0.192	0.664	0.120	0.603	0.210	0.731	0.229	0.684	0.134
wilt	0.520	0.053	0.443	0.044	0.966	0.554	0.316	0.036	0.627	0.069	0.441	0.029	0.530	0.053	0.876	0.472
average	0.773	0.259	0.777	0.322	0.900	0.647	0.717	0.361	0.676	0.195	0.677	0.269	0.681	0.245	0.744	0.346
cum. train time	61s		70s		N.A.		N.A.		N.A.		2432s		72s		N.A.	



REFERENCES ON OPERATOR-VALUED KERNELS

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. “Tensorflow: Large-scale machine learning on heterogeneous distributed systems.” In: *arXiv preprint arXiv:1603.04467* (2016) (cit. on p. 157).
- [3] R. Ahlswede and A. Winter. “Strong converse for identification via quantum channels.” In: *IEEE Trans. Inform. Theory* 48(3) (2002), pp. 569–679 (cit. on p. 172).
- [4] E. M. Alfsen. “A simplified constructive proof of the existence and uniqueness of Haar measure.” In: *Mathematica Scandinavica* 12.1 (1964), pp. 106–116 (cit. on p. 29).
- [5] 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 (cit. on pp. 4, 45, 48, 97, 139, 144).
- [6] P.-O. Amblard and H. Kadri. “Operator-valued kernel recursive least squares algorithm.” In: *Signal Processing Conference (EUSIPCO), 2015 23rd European*. IEEE. 2015, pp. 2376–2380 (cit. on p. 49).
- [8] A. Argyriou, C. A. Micchelli, and M. Pontil. “When Is There a Representer Theorem? Vector vs Matrix Regularizers.” In: *Journal of Machine Learning Research* 10 (2009), pp. 2507–2529 (cit. on p. 97).
- [9] N. Aronszajn. “Theory of reproducing kernels.” In: *Transactions of the American mathematical society* (1950), pp. 337–404 (cit. on pp. 4, 7, 11, 12, 36).
- [10] J. Audiffren and H. Kadri. “Online learning with operator-valued kernels.” In: *European symposium on artificial neural networks (ESANN)*. 2015 (cit. on pp. 16, 49, 119, 142, 162).
- [11] I. Avramidi. “Notes on Hilbert Spaces.” In: (2000) (cit. on p. 28).
- [12] F. Bach. *On the Equivalence between Quadrature Rules and Random Features*. HAL-report-/hal-01118276. 2015 (cit. on pp. 4, 18).
- [13] F. Bach. “On the equivalence between quadrature rules and random features.” In: *arXiv preprint arXiv:1502.06800* (2015) (cit. on p. 135).

- [14] 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 (cit. on p. 45).
- [15] L. Baldassarre, L. Rosasco, A. Barla, and A. Verri. “Multi-output learning via spectral filtering.” In: *Machine Learning* 87.3 (2012), pp. 259–301 (cit. on pp. 46–48).
- [17] P. L. Bartlett and S. Mendelson. “Rademacher and Gaussian complexities: Risk bounds and structural results.” In: *Journal of Machine Learning Research* 3.Nov (2002), pp. 463–482 (cit. on pp. 18, 127–129).
- [19] M. Belkin, P. Niyogi, and V. Sindhwani. “Manifold regularization: A geometric framework for learning from labeled and unlabeled examples.” In: *Journal of machine learning research* 7.Nov (2006), pp. 2399–2434 (cit. on p. 184).
- [20] A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Springer, 2003 (cit. on p. 11).
- [24] B. E. Boser, I. M. Guyon, and V. N. Vapnik. “A training algorithm for optimal margin classifiers.” In: *Proceedings of the fifth annual workshop on Computational learning theory*. ACM, 1992, pp. 144–152 (cit. on p. 11).
- [25] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities*. Oxford Press, 2013 (cit. on pp. 83, 85).
- [26] O. Bousquet and A. Elisseeff. “Stability and generalization.” In: *Journal of Machine Learning Research* 2.Mar (2002), pp. 499–526 (cit. on pp. 130, 131).
- [27] G. E. Box. “Robustness in the strategy of scientific model building.” In: *Robustness in statistics* 1 (1979), pp. 201–236 (cit. on p. 7).
- [31] U. Brefeld, T. Gärtner, T. Scheffer, and S. Wrobel. “Efficient co-regularised least squares regression.” In: *Proceedings of the 23rd international conference on Machine learning*. ACM, 2006, pp. 137–144 (cit. on p. 184).
- [34] C. Brouard, F. d’Alché-Buc, and M. Szafranski. “Semi-supervised Penalized Output Kernel Regression for Link Prediction.” In: *Proc. of the 28th Int. Conf. on Machine Learning*. 2011 (cit. on pp. 4, 48, 98, 184).
- [35] C. Brouard, F. d’Alché-Buc, and M. Szafranski. “Input Output Kernel Regression.” In: *to appear in JMLR* (2016) (cit. on pp. 48, 97, 162, 184).

- [36] C. Brouard, H. Shen, K. Dührkop, F. d'Alché Buc, S. Böcker, and J. Rousu. "Fast metabolite identification with input output Kernel regression." In: *Bioinformatics* 32.12 (2016), pp. i28–i36 (cit. on p. 49).
- [37] R. H. Byrd, P. Lu, J. Nocedal, and C. Zhu. "A limited memory algorithm for bound constrained optimization." In: *SIAM Journal on Scientific Computing* 16.5 (1995), pp. 1190–1208 (cit. on pp. 121, 140).
- [39] A. Caponnetto, C. A. Micchelli, M., and Y. Ying. "Universal MultiTask Kernels." In: *Journal of Machine Learning Research* 9 (2008), pp. 1615–1646 (cit. on pp. 45, 47, 48, 140).
- [40] 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 (cit. on pp. 36–39, 42, 43, 48, 56).
- [41] 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 (cit. on pp. 4, 5, 36, 37, 41–43, 45, 48, 53, 55, 57, 69, 139, 141, 151, 154).
- [43] C. Ciliberto, Y. Mroueh, T. Poggio, and L. Rosasco. "Convex Learning of Multiple Tasks and their Structure." In: *Proc. of the 32nd International Conference on Machine Learning*. 2015 (cit. on p. 49).
- [46] J. B. Conway. *A course in functional analysis*. Vol. 96. Springer Science & Business Media, 2013 (cit. on pp. 29, 106).
- [47] A. Cornuéjols and L. Miclet. *Apprentissage artificiel: concepts et algorithmes*. Editions Eyrolles, 2011 (cit. on p. 9).
- [49] I. I. Cotaescu. "Elements of Linear Algebra. Lecture Notes." In: *arXiv preprint arXiv:1602.03006* (2016) (cit. on p. 22).
- [50] F. Cucker and S. Smale. "On the mathematical foundation of learning." In: *AMERICAN MATHEMATICAL SOCIETY* 39.1 (2001), pp. 1–49 (cit. on pp. 87, 170).
- [51] B. Dai, B. Xie, N. He, Y. Liang, A. Raj, M.-F. F. Balcan, and L. Song. "Scalable kernel methods via doubly stochastic gradients." In: *Advances in Neural Information Processing Systems*. 2014, pp. 3041–3049 (cit. on pp. 4, 6, 20, 137, 142).
- [54] F. Dinuzzo, C. Ong, P. Gehler, and G. Pillonetto. "Learning Output Kernels with Block Coordinate Descent." In: *Proc. of the 28th Int. Conf. on Machine Learning*. 2011 (cit. on pp. 45, 48, 49).

- [55] P. Drineas and M. W. Mahoney. “On the Nyström method for approximating a Gram matrix for improved kernel-based learning.” In: *journal of machine learning research* 6.Dec (2005), pp. 2153–2175 (cit. on pp. 4, 18).
- [56] R. M. Dudley. “The sizes of compact subsets of Hilbert space and continuity of Gaussian processes.” In: *Journal of Functional Analysis* 1.3 (1967), pp. 290–330 (cit. on p. 85).
- [59] T. Evgeniou, C. A. Micchelli, and M. Pontil. “Learning Multiple Tasks with kernel methods.” In: *JMLR* 6 (2005), pp. 615–637 (cit. on pp. 45, 140).
- [60] P. L. Falb. “On a theorem of Bochner.” In: *Publications Mathématiques de l'IHÉS* 36 (1969), pp. 59–67 (cit. on p. 26).
- [61] X. Y. Felix, A. T. Suresh, K. M. Choromanski, D. N. Holtmann-Rice, and S. Kumar. “Orthogonal random features.” In: *Advances in Neural Information Processing Systems*. 2016, pp. 1975–1983 (cit. on pp. 4, 19).
- [62] O. Fercoq and R. Peter. “Accelerated, parallel, and proximal coordinate descent.” In: *SIAM Journal on Optimization* 25.4 (2015), pp. 1997–2023 (cit. on p. 140).
- [63] O. Fercoq and P. Richtárik. “Smooth minimization of nonsmooth functions with parallel coordinate descent methods.” In: *arXiv preprint arXiv:1309.5885* (2013) (cit. on p. 140).
- [65] G. B. Folland. *A course in abstract harmonic analysis*. CRC press, 1994 (cit. on pp. 5, 16, 21, 29–35).
- [68] E. Fuselier. “Refined Error Estimates for Matrix-Valued Radial Basis Functions.” PhD thesis. Texas A&M University, 2006 (cit. on p. 46).
- [74] L. Górniewicz. *Topological fixed point theory of multivalued mappings*. Vol. 495. Springer, 1999 (cit. on p. 96).
- [75] G. Guennebaud, B. Jacob, et al. *Eigen v3*. 2010. URL: <http://eigen.tuxfamily.org> (cit. on p. 114).
- [76] M. Ha Quang, S. H. Kang, and T. M. Le. “Image and video colorization using vector-valued reproducing kernel Hilbert spaces.” In: *Journal of Mathematical Imaging and Vision* 37.1 (2010), pp. 49–65 (cit. on p. 49).
- [77] R. Hamid, Y. Xiao, A. Gittens, and D. DeCoste. “Compact Random Feature Maps.” In: *ICML*. 2014, pp. 19–27 (cit. on p. 19).
- [80] T. Jaakkola, M. Diekhans, and D. Haussler. “Using the Fisher kernel method to detect remote protein homologies.” In: *ISMB*. Vol. 99. 1999, pp. 149–158 (cit. on p. 233).
- [81] E. Jones, T. Oliphant, and P. Peterson. “[SciPy]: open source scientific tools for [Python].” In: (2014) (cit. on p. 110).

- [83] H. Kadri, M. Ghavamzadeh, and P. Preux. “A generalized kernel approach to structured output learning.” In: *Proc. of the 30th International Conference on Machine Learning*. 2013 (cit. on p. 48).
- [84] H. Kadri, E. Duflos, P. Preux, S. Canu, and M. Davy. “Non-linear functional regression: a functional RKHS approach.” In: *JMLR Proc. of International Conference on Artificial Intelligence and Statistics*. Vol. 9. 2010 (cit. on pp. 4, 49).
- [85] H. Kadri, A. Rakotomamonjy, P. Preux, and F. R. Bach. “Multiple operator-valued kernel learning.” In: *Advances in NIPS*. 2012, pp. 2429–2437 (cit. on p. 49).
- [86] H. Kadri, E. Duflos, P. Preux, S. Canu, A. Rakotomamonjy, and J. Audiffren. “Operator-valued kernels for learning from functional response data.” In: *Journal of Machine Learning Research* 16 (2015), pp. 1–54 (cit. on pp. 47–49, 97, 98, 130, 131, 135, 151).
- [87] P. Kar and H. Karnick. “Random Feature Maps for Dot Product Kernels.” In: *AISTATS*. Vol. 22. 2012, pp. 583–591 (cit. on p. 19).
- [89] G. Kimeldorf and G. Wahba. “Some results on Tchebychefian spline functions.” In: *Journal of Mathematical Analysis and Applications* 33.1 (1971), pp. 82 –95 (cit. on p. 11).
- [90] J. Kivinen, A. J. Smola, and R. C. Williamson. “Online learning with kernels.” In: *IEEE transactions on signal processing* 52.8 (2004), pp. 2165–2176 (cit. on pp. 16, 142).
- [91] R. Koenker and G. Bassett Jr. “Regression quantiles.” In: *Econometrica: journal of the Econometric Society* (1978), pp. 33–50 (cit. on p. 150).
- [92] V. Koltchinskii et al. “A remark on low rank matrix recovery and noncommutative Bernstein type inequalities.” In: *From Probability to Statistics and Back: High-Dimensional Models and Processes—A Festschrift in Honor of Jon A. Wellner*. Institute of Mathematical Statistics, 2013, pp. 213–226 (cit. on pp. 5, 81, 86, 88, 170, 172, 174).
- [93] A. J. Kurdila and M. Zabarankin. *Convex functional analysis*. 2006 (cit. on pp. 5, 21, 22, 96, 99, 106).
- [94] Q. V. Le, T. Sarlós, and A. J. Smola. “Fastfood - Computing Hilbert Space Expansions in loglinear time.” In: *Proc. of ICML 2013, Atlanta, USA, 16-21 June 2013*. 2013, pp. 244–252 (cit. on pp. 4, 19, 20).
- [95] Y. LeCun, Y. Bengio, et al. “Convolutional networks for images, speech, and time series.” In: *The handbook of brain theory and neural networks* 3361.10 (1995), p. 1995 (cit. on pp. 20, 166).

- [96] M. Ledoux. *The concentration of measure phenomenon*. 89. American Mathematical Soc., 2005 (cit. on p. 83).
- [97] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media, 2013 (cit. on p. 86).
- [98] 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 (cit. on pp. 16, 19, 31, 32, 63, 77).
- [100] N. Lim, Y. Senbabaoglu, G. Michailidis, and F. d’Alché-Buc. “OKVAR-Boost: a novel boosting algorithm to infer nonlinear dynamics and interactions in gene regulatory networks.” In: *Bioinformatics* 29.11 (2013), pp. 1416–1423 (cit. on pp. 138, 146).
- [101] 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 (cit. on pp. 6, 45, 137–140, 146).
- [102] N. Lim, Y. Şenbabaoglu, G. Michailidis, and F. d’Alché Buc. “OKVAR-Boost: a novel boosting algorithm to infer nonlinear dynamics and interactions in gene regulatory networks.” In: *Bioinformatics* 29.11 (2013), pp. 1416–1423 (cit. on pp. 47, 49).
- [103] 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 (cit. on pp. 47, 49, 162).
- [106] Y. Liu, A. Niculescu-Mizil, A. C. Lozano, and Y. Lu. “Learning temporal causal graphs for relational time-series analysis.” In: *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*. 2010, pp. 687–694 (cit. on p. 145).
- [108] Z. Lu, A. May, K. Liu, A. B. Garakani, D. Guo, A. Bellet, L. Fan, M. Collins, B. Kingsbury, M. Picheny, et al. “How to scale up kernel methods to be as good as deep neural nets.” In: *arXiv preprint arXiv:1411.4000* (2014) (cit. on pp. 4, 20, 165).
- [109] 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 (cit. on pp. 46–48).

- [110] L. Mackey, M. I. Jordan, R. Chen, B. Farrel, and J. Tropp. "Matrix Concentration Inequalities via the method of exchangeable pairs." In: *The Annals of Probability* 42:3 (2014), pp. 906–945 (cit. on pp. 170, 172).
- [112] A. Maurer. "A vector-contraction inequality for Rademacher complexities." In: *International Conference on Algorithmic Learning Theory*. Springer. 2016, pp. 3–17 (cit. on pp. 5, 48, 125, 127–129, 132, 135).
- [113] C. A. Micchelli and M. A. Pontil. "On Learning Vector-Valued Functions." In: *Neural Computation* 17 (2005), pp. 177–204 (cit. on pp. 4, 5, 36, 37, 45, 48, 97, 98, 139, 140).
- [114] C. A. Micchelli and M. Pontil. "Kernels for Multi-task Learning." In: *NIPS*. Vol. 86. 2004, p. 89 (cit. on p. 48).
- [115] M. Micheli and J. Glaunes. *Matrix-valued kernels for shape deformation analysis*. Tech. rep. Arxiv report, 2013 (cit. on p. 46).
- [116] H. Q. Minh, L. Bazzani, and V. Murino. "A unifying framework for vector-valued manifold regularization and multi-view learning." In: *Proc. of the 30th International Conference on Machine Learning*. 2013 (cit. on p. 97).
- [117] H. Q. Minh and V. Sindhwani. "Vector-valued Manifold Regularization." In: *Proc. of the 28th International Conference on Machine Learning*. 2011 (cit. on p. 184).
- [118] H. Q. Minh. "Operator-Valued Bochner Theorem, Fourier Feature Maps for Operator-Valued Kernels, and Vector-Valued Learning." In: *arXiv preprint arXiv:1608.05639* (2016) (cit. on pp. 60, 89, 171, 180).
- [119] H. Q. Minh, L. Bazzani, and V. Murino. "A unifying framework for vector-valued manifold regularization and multi-view learning." In: *ICML* (2). 2013, pp. 100–108 (cit. on pp. 48, 183, 184).
- [120] H. Q. Minh, L. Bazzani, and V. Murino. "A unifying framework in vector-valued reproducing kernel Hilbert spaces for manifold regularization and co-regularized multi-view learning." In: *Journal of Machine Learning Research* 17.25 (2016), pp. 1–72 (cit. on pp. 97, 98, 183–185, 194).
- [121] S. Minsker. "On some extensions of Bernstein's inequality for self-adjoint operators." In: *arXiv preprint arXiv:1112.5448* (2011) (cit. on pp. 5, 81, 86, 90, 173, 175).
- [122] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. The MIT press, 2012 (cit. on pp. 9, 12).
- [123] Y. Mroueh, T. Poggio, L. Rosasco, and J.-j. Slotine. "Multiclass learning with simplex coding." In: *Advances in NIPS*. 2012, pp. 2789–2797 (cit. on pp. 48, 119, 192).

- [124] Y. Mukuta and T. Harada. “Kernel Approximation via Empirical Orthogonal Decomposition for Unsupervised Feature Learning.” In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016, pp. 5222–5230 (cit. on p. 20).
- [125] A. Naor. “On the Banach-space-valued Azuma inequality and small-set isoperimetry of Alon–Roichman graphs.” In: *Combinatorics, Probability and Computing* 21.04 (2012), pp. 623–634 (cit. on p. 83).
- [126] K.-H. Neeb. “Operator-valued positive definite kernels on tubes.” In: *Monatshefte für Mathematik* 126.2 (1998), pp. 125–160 (cit. on pp. 54, 55).
- [127] T. E. Oliphant. *A guide to NumPy*. Vol. 1. Trelgol Publishing USA, 2006 (cit. on p. 110).
- [128] R. I. Oliveira. “Concentration of the adjacency matrix and of the Laplacian in random graphs with independent edges.” In: *arXiv preprint arXiv:0911.0600* (2009) (cit. on p. 172).
- [130] N. Parikh, S. Boyd, et al. “Proximal algorithms.” In: *Foundations and Trends® in Optimization* 1.3 (2014), pp. 127–239 (cit. on p. 140).
- [132] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. “Scikit-learn: Machine learning in Python.” In: *Journal of Machine Learning Research* 12.Oct (2011), pp. 2825–2830 (cit. on pp. 6, 110, 149, 160, 233).
- [133] G. Pedrick. *Theory of reproducing kernels for Hilbert spaces of vector-valued functions*. Tech. rep. University of Kansas, Department of Mathematics, 1957 (cit. on pp. 36, 139).
- [134] T. Penzl. “Numerical solution of generalized Lyapunov equations.” In: *Advances in Computational Mathematics* 8.1 (1998), pp. 33–48 (cit. on p. 112).
- [135] N. Pham and R. Pagh. “Fast and scalable polynomial kernels via explicit feature maps.” In: *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM. 2013, pp. 239–247 (cit. on p. 19).
- [136] I. Pinelis. “Optimum bounds for the distributions of martingales in Banach spaces.” In: *The Annals of Probability* (1994), pp. 1679–1706 (cit. on pp. 83, 86, 133).
- [139] A. Rahimi and B. Recht. “Random Features for Large-Scale Kernel Machines.” In: *NIPS 2007*. 2007, pp. 1177–1184 (cit. on pp. 4, 5, 7, 16–19, 76, 81, 84, 85, 87, 89, 90, 138, 141, 164, 170, 180).

- [140] A. Rahimi and B. Recht. "Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning." In: *Advances in neural information processing systems*. 2009, pp. 1313–1320 (cit. on pp. 5, 19, 125, 132, 133, 135).
- [141] P. Richtárik and M. Takáč. "Parallel coordinate descent methods for big data optimization." In: *Mathematical Programming* 156.1-2 (2016), pp. 433–484 (cit. on p. 140).
- [142] L. Rosasco, M. Belkin, and E. D. Vito. "On learning with integral operators." In: *Journal of Machine Learning Research* 11.Feb (2010), pp. 905–934 (cit. on pp. 18, 98, 112).
- [143] D Rosenberg, V. Sindhwani, P Bartlett, and P. Niyogi. "A kernel for semi-supervised learning with multi-view point cloud regularization." In: *IEEE Signal Processing Magazine* 26.5 (2009), pp. 145–150 (cit. on p. 184).
- [144] A. Rudi, R. Camoriano, and L. Rosasco. "Generalization properties of learning with random features." In: *arXiv preprint arXiv:1602.04474* (2016) (cit. on pp. 4, 18, 19, 135, 165).
- [145] A. L. Samuel. "Some studies in machine learning using the game of checkers." In: *IBM Journal of research and development* 3.3 (1959), pp. 210–229 (cit. on p. vii).
- [146] M. Sangnier, O. Fercoq, and F. d'Alché Buc. "Joint quantile regression in vector-valued RKHSs." In: *Neural Information Processing Systems* (2016) (cit. on pp. 48, 49, 97, 150, 155, 157, 162).
- [152] E. Senkene and A. Tempel'man. "Hilbert Spaces of operator-valued functions." In: *Lithuanian Mathematical Journal* 13.4 (1973), pp. 665–670 (cit. on pp. 37, 139).
- [154] J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge press, 2004 (cit. on p. 11).
- [158] V. Sindhwani, H. Q. Minh, and A. 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 (cit. on pp. 45, 48, 49, 112).
- [159] V. Sindhwani and D. S. Rosenberg. "An RKHS for multi-view learning and manifold co-regularization." In: *Proceedings of the 25th international conference on Machine learning*. ACM. 2008, pp. 976–983 (cit. on p. 184).
- [160] G. L. Sleijpen, P. Sonneveld, and M. B. Van Gijzen. "Bi-CGSTAB as an induced dimension reduction method." In: *Applied Numerical Mathematics* 60.11 (2010), pp. 1100–1114 (cit. on p. 112).

- [161] S. Smale and D.-X. Zhou. “Learning theory estimates via integral operators and their approximations.” In: *Constructive approximation* 26.2 (2007), pp. 153–172 (cit. on p. 83).
- [162] A. J. Smola, B. Schölkopf, and K.-R. Müller. “The connection between regularization operators and support vector kernels.” In: *Neural networks* 11.4 (1998), pp. 637–649 (cit. on p. 79).
- [163] P. Sonneveld and M. B. van Gijzen. “IDR (s): A family of simple and fast algorithms for solving large nonsymmetric systems of linear equations.” In: *SIAM Journal on Scientific Computing* 31.2 (2008), pp. 1035–1062 (cit. on p. 142).
- [164] B. Sriperumbudur and Z. Szabo. “Optimal Rates for Random Fourier Features.” In: *Advances in NIPS 28*. Ed. by C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett. 2015, pp. 1144–1152 (cit. on pp. 4, 17, 84, 85).
- [166] S. Sun. “Multi-view Laplacian support vector machines.” In: *International Conference on Advanced Data Mining and Applications*. Springer. 2011, pp. 209–222 (cit. on p. 184).
- [167] D. J. Sutherland and J. G. Schneider. “On the Error of Random Fourier Features.” In: *Proc. of UAI 2015, July 12-16, 2015, Amsterdam, The Netherlands*. 2015, pp. 862–871 (cit. on pp. 4, 5, 17, 18, 81, 84, 85, 89–91, 138, 170–172, 180).
- [171] A. Tewari and P. L. Bartlett. “Learning theory.” In: *Academic Press Library in Signal Processing* 1 (), pp. 775–816 (cit. on p. 130).
- [172] T. Tieleman and G Hinton. “Lecture 6.5-RMSProp, COURSE-ERA: Neural networks for machine learning.” In: *University of Toronto, Tech. Rep* (2012) (cit. on p. 157).
- [173] S. Tornier. *Haar measures*. https://people.math.ethz.ch/~torniers/download/2014/haar_measures.pdf. 2014 (cit. on p. 29).
- [174] J. A Tropp. “User-friendly tail bounds for sums of random matrices.” In: *Foundations of computational mathematics* 12.4 (2012), pp. 389–434 (cit. on pp. 170, 172).
- [175] J. A. Tropp et al. “An introduction to matrix concentration inequalities.” In: *Foundations and Trends® in Machine Learning* 8.1-2 (2015), pp. 1–230 (cit. on pp. 5, 81, 86, 173).
- [176] A. W Van Der Vaart and J. A Wellner. “Weak Convergence.” In: *Weak Convergence and Empirical Processes*. Springer, 1996, pp. 16–28 (cit. on p. 86).
- [177] V. N. Vapnik. *Statistical learning theory*. Vol. 1. Wiley New York, 1998 (cit. on pp. 4, 7, 126).

- [178] V. Vapnik. "Principles of risk minimization for learning theory." In: *Advances in neural information processing systems*. 1992, pp. 831–838 (cit. on pp. 9, 130).
- [179] E. Vazquez and E. Walter. "Multi-output support vector regression." In: *13th IFAC Symposium on System Identification*. 2003, pp. 1820–1825 (cit. on p. 47).
- [180] J.-P. Vert. "Regularization of Kernel Methods by Decreasing the Bandwidth of the Gaussian Kernel." In: () (cit. on p. 79).
- [181] U. Von Luxburg. "A tutorial on spectral clustering." In: *Statistics and computing* 17.4 (2007), pp. 395–416 (cit. on pp. 189, 192).
- [182] G. Wahba. *Spline model for observational data*. Philadelphia, Society for Industrial and Applied Mathematics, 1990 (cit. on pp. 5, 13, 95, 98).
- [183] N. Wahlström, M. Kok, T. Schön, and F. Gustafsson. "Modeling magnetic fields using Gaussian processes." In: *in Proc. of the 38th ICASSP*. 2013 (cit. on pp. 46, 48).
- [184] S. v. d. Walt, S. C. Colbert, and G. Varoquaux. "The NumPy array: a structure for efficient numerical computation." In: *Computing in Science & Engineering* 13.2 (2011), pp. 22–30 (cit. on p. 111).
- [185] C. K. I. Williams and M. Seeger. "Using the Nyström Method to Speed Up Kernel Machines." In: *Advances in Neural Information Processing Systems 13*. Ed. by T. K. Leen, T. G. Dietterich, and V. Tresp. MIT Press, 2001, pp. 682–688 (cit. on pp. 4, 7, 17).
- [186] B. Xie, Y. Liang, and L. Song. "Scale up nonlinear component analysis with doubly stochastic gradients." In: *Advances in Neural Information Processing Systems*. 2015, pp. 2341–2349 (cit. on p. 20).
- [188] J. Yang, V. Sindhwani, H. Avron, and M. Mahoney. "Quasi-Monte Carlo feature maps for shift-invariant kernels." In: *Proceedings of The 31st International Conference on Machine Learning (ICML-14)*. 2014, pp. 485–493 (cit. on p. 19).
- [189] J. Yang, V. Sindhwani, Q. Fan, H. Avron, and M. W. Mahoney. "Random laplace feature maps for semigroup kernels on histograms." In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2014, pp. 971–978 (cit. on p. 19).
- [190] T. Yang, Y.-F. Li, M. Mahdavi, R. Jin, and Z. Zhou. "Nyström Method vs Random Fourier Features: A Theoretical and Empirical Comparison." In: *NIPS 25*. Ed. by F. Pereira, C. Burges, L. Bottou, and K. Weinberger. 2012, pp. 476–484 (cit. on pp. 18, 19, 79).

- [191] Z. Yang, A. G. Wilson, A. J. Smola, and L. Song. “A la Carte - Learning Fast Kernels.” In: *Proc. of AISTATS 2015, San Diego, California, USA, 2015*. 2015 (cit. on pp. 4, 20, 165).
- [192] Z. Yang, M. Moczulski, M. Denil, N. de Freitas, A. Smola, L. Song, and Z. Wang. “Deep fried convnets.” In: *Proceedings of the IEEE International Conference on Computer Vision*. 2015, pp. 1476–1483 (cit. on pp. 4, 20, 166).
- [193] H. Zhang, Y. Xu, and Q. Zhang. “Refinement of Operator-valued Reproducing Kernels.” In: *Journal of Machine Learning Research* 13 (2012), pp. 91–136 (cit. on pp. 54, 55).
- [194] D.-X. Zhou. “The covering number in learning theory.” In: *Journal of Complexity* 18.3 (2002), pp. 739–767 (cit. on p. 130).

PUBLICATIONS

- [29] R. Brault, M. Heinonen, and F. d’Alché Buc. “Random Fourier Features For Operator-Valued Kernels.” In: *Proceedings of The 8th Asian Conference on Machine Learning*. 2016, pp. 110–125 (cit. on pp. 5, 53, 54, 81, 86, 112, 120, 162, 171, 173).

WORKSHOPS

- [30] R. Brault, N. Lim, and F. d’Alché Buc. “Scaling up Vector Autoregressive Models With Operator-Valued Random Fourier Features.” In: *Proceedings of AALTD 2016: Second ECML/PKDD International Workshop on Advanced Analytics and Learning on Temporal Data*. 2016, p. 3 (cit. on pp. 6, 137, 141).

COMMUNICATION ACTS

- [28] R. Brault and F. d’Alché Buc. “Borne sur l’approximation de noyaux à valeurs opérateurs à l’aide de transformées de Fourier.” In: *Société Française des Statistiques*. SFDS. 2016 (cit. on pp. 5, 81).

REFERENCES ON EXTREMS

- [2] C. Aggarwal and P. Yu. “Outlier detection for high dimensional data.” In: *ACM Sigmod Record*. 2001 (cit. on p. 214).
- [7] Y. Amit and D. Geman. “Shape quantization and recognition with randomized trees.” In: *Neural comp.* (1997) (cit. on p. 223).
- [16] V. Barnett and T. Lewis. *Outliers in statistical data*. Wiley New York, 1994 (cit. on p. 214).
- [18] S. Bay and M. Schwabacher. “Mining distance-based outliers in near linear time with randomization and a simple pruning rule.” In: *KDD*. 2003 (cit. on p. 224).
- [21] G. Biau, L. Devroye, and G. Lugosi. “Consistency of random forests and other averaging classifiers.” In: *Journal of Machine Learning Research* 9.Sep (2008), pp. 2015–2033 (cit. on p. 214).
- [22] G. Biau and E. Scornet. “A random forest guided tour.” In: *Test* 25.2 (2016), pp. 197–227 (cit. on p. 214).
- [23] G. Blanchard, C. Schäfer, and Y. Rozenholc. “Oracle bounds and exact algorithm for dyadic classification trees.” In: *International Conference on Computational Learning Theory*. Springer. 2004, pp. 378–392 (cit. on p. 215).
- [32] L. Breiman. “Random Forests.” English. In: *Machine Learning* (2001). ISSN: 0885-6125 (cit. on pp. 214, 223, 231).
- [33] M. Breunig, H. Kriegel, R. Ng, and J. Sander. “LOF: identifying density-based local outliers.” In: *ACM Sigmod Rec.* 2000 (cit. on pp. 214, 224).
- [38] N. A. Campbell. “Robust procedures in multivariate analysis I: Robust covariance estimation.” In: *Applied statistics* (1980), pp. 231–237 (cit. on p. 160).
- [42] V. Chandola, A. Banerjee, and V. Kumar. “Anomaly detection: A survey.” In: *ACM Comput. Surv.* (2009) (cit. on p. 214).
- [44] S. Clémençon and S. Robbiano. “Anomaly Ranking as Supervised Bipartite Ranking.” In: *ICML*. 2014 (cit. on pp. 215, 216).
- [45] S. Clémençon and N. Vayatis. “Tree-based ranking methods.” In: *IEEE Transactions on Information Theory* 55.9 (2009), pp. 4316–4336 (cit. on p. 215).
- [48] C. Cortes and V. Vapnik. “Support-Vector Networks.” In: *Machine Learning* (1995) (cit. on pp. 157, 214).

- [52] C. Désir, S. Bernard, C. Petitjean, and L. Heutte. “One Class Random Forests.” In: *Pattern Recogn.* (2013) (cit. on pp. 214, 215, 219, 224).
- [53] R. Díaz-Uriarte and S. De Andres. “Gene selection and classification of microarray data using random forest.” In: *BMC bioinformatics* (2006) (cit. on p. 214).
- [57] J. H. Einmahl and D. M. Mason. “Generalized quantile processes.” In: *The Annals of Statistics* (1992), pp. 1062–1078 (cit. on p. 215).
- [58] E. Eskin. “Anomaly Detection over Noisy Data using Learned Probability Distributions.” In: *ICML*. 2000 (cit. on p. 214).
- [64] P. Flach. “The geometry of ROC space: understanding ml metrics through ROC isometrics.” In: *ICML*. 2003 (cit. on p. 219).
- [66] Y. Freund, R. Schapire, et al. “Experiments with a new boosting algorithm.” In: *ICML*. 1996 (cit. on p. 214).
- [69] R. Genuer, J.-M. Poggi, and C. Tuleau-Malot. “Variable selection using random forests.” In: *Pattern Recog. Letters* (2010) (cit. on p. 214).
- [70] R. Genuer, J.-M. Poggi, and C. Tuleau. “Random Forests: some methodological insights.” In: *arXiv:0811.3619* (2008) (cit. on p. 214).
- [71] P. Geurts, D. Ernst, and L. Wehenkel. “Extremely randomized trees.” In: *Machine learning* (2006) (cit. on pp. 214, 217, 222).
- [72] C. Gini. “Variabilità e mutabilità.” In: *Memorie di metodologia statistica* (1912) (cit. on p. 217).
- [78] T. Ho. “The random subspace method for constructing decision forests.” In: *TPAMI* (1998) (cit. on p. 223).
- [79] V. Hodge and J. Austin. “A survey of outlier detection methodologies.” In: *Artif. Intel. Review* (2004) (cit. on p. 214).
- [82] KDDCup. “The third international knowledge discovery and data mining tools competition dataset.” In: (1999) (cit. on p. 233).
- [99] M. Lichman. *UCI Machine Learning Repository*. 2013. URL: <http://archive.ics.uci.edu/ml> (cit. on pp. 224, 233).
- [104] F. Liu, K. Ting, and Z.-H. Zhou. “Isolation-Based Anomaly Detection.” In: *ACM Trans. Knowl. Discov. Data* (Mar. 2012). ISSN: 1556-4681. doi: [10.1145/2133360.2133363](https://doi.org/10.1145/2133360.2133363). URL: <http://doi.acm.org/10.1145/2133360.2133363> (cit. on p. 232).
- [105] F. Liu, K. Ting, and Z. Zhou. “Isolation Forest.” In: *ICDM*. 2008 (cit. on pp. 160, 214, 222, 224, 230, 233).
- [107] G. Louppe. “Understanding random forests: From theory to practice.” In: *arXiv:1407.7502* (2014) (cit. on p. 214).

- [111] M. Markou and S. Singh. "Novelty detection: a review part 1: statistical approaches." In: *Signal proc.* (2003) (cit. on p. 214).
- [129] P. Panov and S. Džeroski. *Combining bagging and random subspaces to create better ensembles*. Springer, 2007 (cit. on p. 223).
- [131] A. Patcha and J. Park. "An overview of anomaly detection techniques: Existing solutions and latest technological trends." In: *Computer Networks* (2007) (cit. on p. 214).
- [137] W. Polonik. "Minimum volume sets and generalized quantile processes." In: *Stochastic Processes and their Applications* (1997) (cit. on p. 215).
- [138] J. Quinn and M. Sugiyama. "A least-squares approach to anomaly detection in static and sequential data." In: *Pattern Recognition Letters* (2014) (cit. on pp. 214, 224).
- [147] B. Schölkopf, J. Platt, J. Shawe-Taylor, A. Smola, and R. Williamson. "Estimating the support of a high-dimensional distribution." In: *Neural computation* (2001) (cit. on pp. 157, 160, 214, 224).
- [148] E. Schubert, R. Wojdanowski, A. Zimek, and H.-P. Kriegel. "On Evaluation of Outlier Rankings and Outlier Scores." In: *SDM*. 2012 (cit. on p. 234).
- [149] E. Scornet, G. Biau, J.-P. Vert, et al. "Consistency of random forests." In: *The Annals of Statistics* 43.4 (2015), pp. 1716–1741 (cit. on p. 216).
- [150] C. Scott and G. Blanchard. "Novelty detection: Unlabeled data definitely help." In: *AISTATS*. 2009, pp. 464–471 (cit. on p. 221).
- [151] C. Scott and R. Nowak. "Learning minimum volume sets." In: *JMLR* (2006) (cit. on pp. 214, 215).
- [153] C. E. Shannon. "A mathematical theory of communication." In: *ACM SIGMOBILE MC2R* (2001) (cit. on p. 217).
- [155] J. Shawe-Taylor and N. Cristianini. *Kernel methods for pattern analysis*. Cambridge univ. press, 2004 (cit. on pp. 157, 214).
- [156] T. Shi and S. Horvath. "Unsupervised learning with random forest predictors." In: *J. Comp. Graph. Stat.* (2012) (cit. on pp. 214, 215, 219, 224).
- [157] M. Shyu, S. Chen, K. Sarinnapakorn, and L. Chang. *A novel anomaly detection scheme based on principal component classifier*. Tech. rep. DTIC Document, 2003 (cit. on p. 214).
- [165] M. Sugiyama. "Superfast-trainable multi-class probabilistic classifier by least-squares posterior fitting." In: *IEICE Transactions on Information and Systems* (2010) (cit. on p. 214).

- [168] V. Svetnik, A. Liaw, C. Tong, J. C. Culberson, R. Sheridan, and B. Feuston. “Random forest: a classification and regression tool for compound classification and QSAR modeling.” In: *J. Chem. Inf. Model.* (2003) (cit. on p. [214](#)).
- [169] M. Tavallaei, E. Bagheri, W. Lu, and A. Ghorbani. “A detailed analysis of the KDD CUP 99 data set.” In: *IEEE CISDA*. 2009 (cit. on p. [233](#)).
- [170] D. Tax and R. Duin. “Uniform object generation for optimizing one-class classifiers.” In: *JMLR* (2002) (cit. on p. [215](#)).
- [187] K. Yamanishi, J. Takeuchi, G. Williams, and P. Milne. “On-line unsupervised outlier detection using finite mixtures with discounting learning algorithms.” In: *KDD*. 2000 (cit. on p. [233](#)).

PUBLICATIONS UNDER REVIEW

- [73] N. Goix, R. Brault, N. Drougard, and M. Chiapino. “One Class Splitting Criteria for Random Forests.” In: *arXiv preprint arXiv:1611.01971* (2016) (cit. on p. [213](#)).



DECLARATION

I hereby declare that this manuscript entitled "*Data Are Not Reals! Large-scale Operator-Valued Kernel Regression*" is a presentation of my original research work done during my thesis, carried out at Université d'Évry-Val-d'Essonne, Évry, Télécom ParisTech and Université Paris-Saclay, Paris, for the degree of Doctor of Philosophy in Computer Science of Université Paris-Saclay.

The interpretations put forth are based on my reading and understanding of the original texts and they are not published anywhere in the form of books, monographs or articles. Wherever contributions of others are involved, every effort is made to indicate this clearly, with due reference to the literature, and acknowledgement of collaborative research and discussions. For the present thesis, which I am submitting to the University, no degree or diploma or distinction has been conferred on me before, either in this or in any other University. The present work was done under the guidance of Professor Florence d'Alché-Buc, at Télécom ParisTech, Université Paris-Saclay, Paris.

Place: (Mr. Romain Raymond Brault)

Date: Research student

This is to certify that the work incorporated in the manuscript "*Data Are Not Reals! Large-scale Operator-Valued Kernel Regression*" submitted by Romain Raymond Brault was carried out by the candidate under my guidance. Such materials as has been obtained from other sources have been duly acknowledged in the thesis. In my capacity as supervisor of the candidates thesis, I certify that the above statements of this page are true to the best of my knowledge.

Place: (Prof. Florence d'Alché-Buc)

Date: Research supervisor

46, Rue Barrault, 75013 — Paris, France, October 6, 2017.

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 L^AT_EX and LyX at

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

Université Paris-Saclay

ED STIC — 580

Université Paris Sud, Bâtiment 650 Ada Lovelace, 91405 Orsay Cedex,
France.

LARGE-SCALE OPERATOR-VALUED KERNEL REGRESSION

Keywords: Operator-Valued Kernels, Large Scale Learning, Random Fourier Features

Abstract: Many problems in Machine Learning can be cast into vector-valued functions approximation. *Operator-Valued Kernels* and vector-valued Reproducing Kernel Hilbert Spaces provide a theoretical and practical framework to address that issue, extending nicely the well-known setting of scalar-valued kernels. However large scale applications are usually not affordable with these tools that require an important computational power along with a large memory capacity. In this thesis, we propose and study scalable methods to perform regression with *Operator-Valued Kernels*. To achieve this goal, we extend Random Fourier Features, an approximation technique originally introduced for scalar-valued kernels, to *Operator-Valued Kernels*. The idea is to take advantage of an approximated operator-valued feature map in order to come up with a linear model in a finite-dimensional space.

This thesis is structured as follows. First we develop a general framework devoted to the approximation of shift-invariant Mercer kernels on Locally Compact Abelian groups and study their properties along with the complexity of the algorithms based on them. Second we show theoretical guarantees by bounding the error due to the approximation, with high probability. Third, we study various applications of Operator Random Fourier Features (ORFF) to different tasks of Machine learning such as multi-class classification, multi-task learning, time serie modeling, functional regression and anomaly detection. We also compare the proposed framework with other state of the art methods. Fourth, we conclude by drawing short-term and mid-term perspectives of this work.

RÉGRESSION À NOYAUX À VALEURS OPÉRATEURS POUR GRANDS ENSEMBLES DE DONNÉES

Mots clefs: Noyaux à Valeurs Opérateurs, Passage à l'échelle, Random Fourier Features

Résumé: De nombreuses problématiques d'apprentissage artificiel peuvent être modélisées grâce à des fonctions à valeurs vectorielles. Les noyaux à valeurs opérateurs et leur espace de Hilbert à noyaux reproduisant à valeurs vectorielles associés donnent un cadre théorique et pratique pour apprendre de telles fonctions, étendant la littérature existante des noyaux scalaires. Cependant, lorsque les données sont nombreuses, ces méthodes sont peu utilisables, ne passant pas à l'échelle, car elle nécessite une quantité de mémoire évoluant quadratiquement et un temps de calcul évoluant cubiquement vis à vis du nombre de données, dans leur implémentation la plus naïve. Afin de faire passer les noyaux à valeurs opérateurs à l'échelle, nous étendons une technique d'approximation stochastique introduite dans le cadre des noyaux scalaires. L'idée est de tirer parti d'une fonction de redescription caractérisant le noyau à valeurs opérateurs, dont les fonctions associées vivent dans un espace de

dimension infinie, afin d'obtenir un problème d'optimisation linéaire de dimension finie.

Dans cette thèse nous développons dans un premier temps un cadre général afin de permettre l'approximation de noyaux de Mercer définis sur des groupes commutatifs localement compacts et étudions leurs propriétés ainsi que la complexité des algorithmes en découlant. Dans un second temps nous montrons des garanties théoriques en bornant l'erreur commise par l'approximation, avec grande probabilité. Enfin, nous mettons en évidence plusieurs applications des Représentations Opérateurs Aléatoires de Fourier (ORFF) telles que la classification multiple, l'apprentissage multi-tâche, la modélisation de séries temporelles, la régression fonctionnelle et la détection d'anomalies. Nous comparons également ce cadre avec d'autres méthodes de la littérature et concluons par des perspectives à moyen et long terme.

