

Mathematical Pictures at a Data Science Exhibition

In the past few decades, heuristic methods adopted by big tech companies have complemented existing scientific disciplines to form the new field of Data Science. This text provides deep and comprehensive coverage of the mathematical theory supporting the field. Composed of 27 lecture-length chapters with exercises, it embarks the readers on an engaging itinerary through key subjects in data science, including machine learning, optimal recovery, compressive sensing (also known as compressed sensing), optimization, and neural networks. While standard material is covered, the book also includes distinctive presentations of topics such as reproducing kernel Hilbert spaces, spectral clustering, optimal recovery, compressive sensing, group testing, and applications of semidefinite programming. Students and data scientists with less mathematical background will appreciate the appendices that supply more details on some of the abstract concepts.

SIMON FOUCART is Professor of Mathematics at Texas A&M University, where he was named Presidential Impact Fellow in 2019. He has previously written, together with Holger Rauhut, the influential book *A Mathematical Introduction to Compressive Sensing* (2013).

Cambridge University Press
978-1-316-51888-5 — Mathematical Pictures at a Data Science Exhibition
Simon Foucart
Frontmatter
[More Information](#)

Cambridge University Press
978-1-316-51888-5 — Mathematical Pictures at a Data Science Exhibition
Simon Foucart
Frontmatter
[More Information](#)

Mathematical Pictures
at a
Data Science Exhibition

SIMON FOUCART
Texas A&M University



CAMBRIDGE
UNIVERSITY PRESS

Cambridge University Press
978-1-316-51888-5 — Mathematical Pictures at a Data Science Exhibition
Simon Foucart
Frontmatter
[More Information](#)

CAMBRIDGE
UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom
One Liberty Plaza, 20th Floor, New York, NY 10006, USA
477 Williamstown Road, Port Melbourne, VIC 3207, Australia
314–321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre,
New Delhi – 110025, India
103 Penang Road, #05–06/07, Visioncrest Commercial, Singapore 238467

Cambridge University Press is part of the University of Cambridge.
It furthers the University's mission by disseminating knowledge in the pursuit of
education, learning, and research at the highest international levels of excellence.

www.cambridge.org
Information on this title: www.cambridge.org/9781316518885
DOI: 10.1017/9781009003933

© Simon Foucart 2022

This publication is in copyright. Subject to statutory exception
and to the provisions of relevant collective licensing agreements,
no reproduction of any part may take place without the written
permission of Cambridge University Press.

First published 2022

A catalogue record for this publication is available from the British Library.

ISBN 978-1-316-51888-5 Hardback
ISBN 978-1-009-00185-4 Paperback

Cambridge University Press has no responsibility for the persistence or accuracy of
URLs for external or third-party internet websites referred to in this publication
and does not guarantee that any content on such websites is, or will remain,
accurate or appropriate.

Cambridge University Press
978-1-316-51888-5 — Mathematical Pictures at a Data Science Exhibition
Simon Foucart
Frontmatter
[More Information](#)

Pour Jeanne, à nouveau,
and now also for Émile and Léopold

Cambridge University Press
978-1-316-51888-5 — Mathematical Pictures at a Data Science Exhibition
Simon Foucart
Frontmatter
[More Information](#)

Contents

<i>Preface</i>	page xiii
<i>Notation</i>	xvii
PART ONE MACHINE LEARNING	1
<i>Executive Summary</i>	3
1 Rudiments of Statistical Learning Theory	4
1.1 True and Empirical Risks	4
1.2 PAC-Learnability	6
1.3 Validation	8
Exercises	9
2 Vapnik–Chervonenkis Dimension	10
2.1 Definitions	10
2.2 Examples	11
2.3 Sauer Lemma	13
Exercises	15
3 Learnability for Binary Classification	16
3.1 Uniform Convergence Property	16
3.2 Finite VC-Dimension Implies PAC-Learnability	17
3.3 No-Free-Lunch Theorem	20
Exercises	22
4 Support Vector Machines	23
4.1 Linear Separability	23
4.2 Hard and Soft SVM	25
4.3 Kernel Trick	27
Exercises	29

viii	<i>Contents</i>	
5	Reproducing Kernel Hilbert Spaces	31
5.1	Abstract Definition	31
5.2	Moore–Aronszajn Theorem	33
5.3	Mercer Theorem	37
	Exercises	39
6	Regression and Regularization	41
6.1	Empirical Risk Minimization	41
6.2	Regularization	43
6.3	Classification via Regression	44
	Exercises	45
7	Clustering	47
7.1	Single-Linkage Clustering	48
7.2	Center-Based Clustering	49
7.3	Spectral Clustering	51
	Exercises	54
8	Dimension Reduction	56
8.1	Principal Component Analysis	56
8.2	Johnson–Lindenstrauss Lemma	58
8.3	Locally Linear Embedding	61
	Exercises	64
	PART TWO OPTIMAL RECOVERY	65
	<i>Executive Summary</i>	67
9	Foundational Results of Optimal Recovery	68
9.1	Models, Errors, and Optimality	69
9.2	Linearity of Optimal Recovery Maps	71
9.3	An Extremal Property of Splines	73
	Exercises	75
10	Approximability Models	76
10.1	The Model Set	76
10.2	Optimality in a Hilbert Setting	78
10.3	Optimality for Linear Functionals	82
	Exercises	85
11	Ideal Selection of Observation Schemes	86
11.1	Hilbert Setting	86
11.2	Integration of Lipschitz Functions	88
11.3	Adaptivity Does Not Help Much	91
	Exercises	92

<i>Contents</i>		ix
12	Curse of Dimensionality	94
	12.1 Notions of Tractability	94
	12.2 Integration of Trigonometric Polynomials	96
	12.3 Integration in Weighted Sobolev Spaces	98
	Exercises	100
13	Quasi-Monte Carlo Integration	102
	13.1 Variation and Discrepancy	103
	13.2 Koksma–Hlawka Inequality	105
	13.3 Low-Discrepancy Sets	107
	Exercises	111
	PART THREE COMPRESSIVE SENSING	113
	<i>Executive Summary</i>	115
14	Sparse Recovery from Linear Observations	116
	14.1 ℓ_0 -Minimization	117
	14.2 ℓ_1 -Minimization	118
	14.3 ℓ_1 -Restricted Isometry Property	119
	Exercises	122
15	The Complexity of Sparse Recovery	123
	15.1 Limitations Imposed by Stability and Robustness	123
	15.2 Gelfand Width of the ℓ_1 -Ball	127
	15.3 Irrelevance of ℓ_2 -Stability	130
	Exercises	130
16	Low-Rank Recovery from Linear Observations	132
	16.1 Nuclear Norm Minimization	132
	16.2 ℓ_1 -Rank Restricted Isometry Property	134
	16.3 Semidefinite Programming Formulation	136
	Exercises	137
17	Sparse Recovery from One-Bit Observations	139
	17.1 Estimating the Direction via Hard Thresholding	139
	17.2 Estimating the Direction via Linear Programming	141
	17.3 Estimating Both the Direction and the Magnitude	145
	Exercises	147
18	Group Testing	149
	18.1 Properties of the Test Matrix	149
	18.2 Satisfying the Separability Condition	151
	18.3 Recovery via a Linear Feasibility Program	154
	Exercises	155

x	<i>Contents</i>	
	PART FOUR OPTIMIZATION	157
	<i>Executive Summary</i>	159
19	Basic Convex Optimization	160
	19.1 Gradient Descent for Unconstrained Convex Programs	160
	19.2 Rates of Convergence for Steepest Descent	162
	19.3 Stochastic Gradient Descent	165
	Exercises	168
20	Snippets of Linear Programming	169
	20.1 Maximizers of a Convex Function	169
	20.2 The Simplex Algorithm	171
	20.3 Illustrative Linear Programs	173
	Exercises	175
21	Duality Theory and Practice	177
	21.1 Duality in Linear Programming	178
	21.2 Examples of Robust Optimization	181
	21.3 Duality in Conic Programming	183
	Exercises	185
22	Semidefinite Programming in Action	187
	22.1 Schur Complement	188
	22.2 Sum-of-Squares Technique	189
	22.3 Method of Moments	191
	Exercises	193
23	Instances of Nonconvex Optimization	195
	23.1 Quadratically Constrained Quadratic Programs	195
	23.2 Dynamic Programming	198
	23.3 Projected Gradient Descent	200
	Exercises	203
	PART FIVE NEURAL NETWORKS	205
	<i>Executive Summary</i>	207
24	First Encounter with ReLU Networks	208
	24.1 Some Terminology	209
	24.2 Shallow ReLU Networks and CPwL Functions	210
	24.3 Deep ReLU Networks and CPwL Functions	213
	Exercises	215

<i>Contents</i>		xi
25	Expressiveness of Shallow Networks	216
	25.1 Activation Functions and Universal Approximation	216
	25.2 Approximation Rate with ReLU: Upper Bound	219
	25.3 Approximation Rate with ReLU: Lower Bound	221
	Exercises	224
26	Various Advantages of Depth	226
	26.1 Omnipotent Activation Functions	226
	26.2 Compact Supports	229
	26.3 Approximation Power	230
	Exercises	237
27	Tidbits on Neural Network Training	239
	27.1 Backpropagation	239
	27.2 Overparametrized Empirical-Risk Landscapes	241
	27.3 Convolutional Neural Networks	245
	Exercises	246
	APPENDICES	247
<i>Appendix A</i>	High-Dimensional Geometry	249
	A.1 Volumes	250
	A.2 Covering and Packing Numbers	252
	A.3 Projected Cross-Polytope	255
	Exercises	257
<i>Appendix B</i>	Probability Theory	259
	B.1 Tails and Moment Generating Functions	259
	B.2 Concentration Inequalities	262
	B.3 Restricted Isometry Properties	267
	Exercises	272
<i>Appendix C</i>	Functional Analysis	274
	C.1 Completeness	274
	C.2 Convexity	277
	C.3 Extreme Points	281
	Exercises	283
<i>Appendix D</i>	Matrix Analysis	285
	D.1 Eigenvalues of Self-Adjoint Matrices	285
	D.2 Singular Values	288
	D.3 Matrix Norms	291
	Exercises	295

Appendix E	Approximation Theory	297
E.1	Classic Uniform Approximation Theorems	297
E.2	Riesz–Fejér and Carathéodory–Toeplitz Theorems	304
E.3	Kolmogorov Superposition Theorem	307
	Exercises	310
	<i>References</i>	311
	<i>Index</i>	315

Preface

Traditional scientific disciplines have lately been complemented by heuristics adopted in big tech companies to form the new field of Data Science. Now making its way into university curricula, this loosely defined field immediately brings computer science and statistics to mind. But mathematics, too, plays a central role by laying foundations and developing new theories. This book focuses on the subfield of Mathematical Data Science. Since its content is also loosely defined, a selection of topics was made to provide summaries only of Machine Learning, Optimal Recovery, Compressive Sensing (also known as Compressed Sensing), Optimization, and Neural Networks.

Audience: This book is intended for mathematicians who wish to know bits and pieces about Data Science. Ideally, it will convince them that there is some elegant theory behind this trendy field. Although the book may also be valuable for genuine data scientists in search of mathematical sophistication, it should primarily serve as a resource for a graduate course on Data Science given in a department of mathematics. In brief, the most important word of the title is the first one.

Theme: The common thread throughout this book is the processing of data given in the form

$$y_i = f(x^{(i)}), \quad i \in [1 : m],$$

toward the ultimate goal of learning/approximating/recovering¹ the unknown function f . In PART ONE (Machine Learning), one mostly thinks of the $x^{(i)}$ as independent realizations of a random variable and one targets results valid in expectation or with high probability. In PART TWO (Optimal Recovery), one thinks instead of the $x^{(i)}$ as fixed entities and one targets results valid

¹ The favored terminology seems to depend on one's inclination and training.

with certainty in a worst-case setting, given some side information about f . In PART THREE (Compressive Sensing), this side information conveys e.g. that f is a linear function depending on few variables. In this framework, it is actually possible to recover f exactly. A shared concern in these first three parts is complexity (sample or information complexity), i.e., the minimal number m of data that makes the learning/approximation/recovery task possible. The task almost invariably requires solving a minimization program, so PART FOUR (Optimization) reviews the necessary material. Finally, PART FIVE (Neural Networks) studies tools for the approximation of f that have recently proved very effective in Deep Learning.

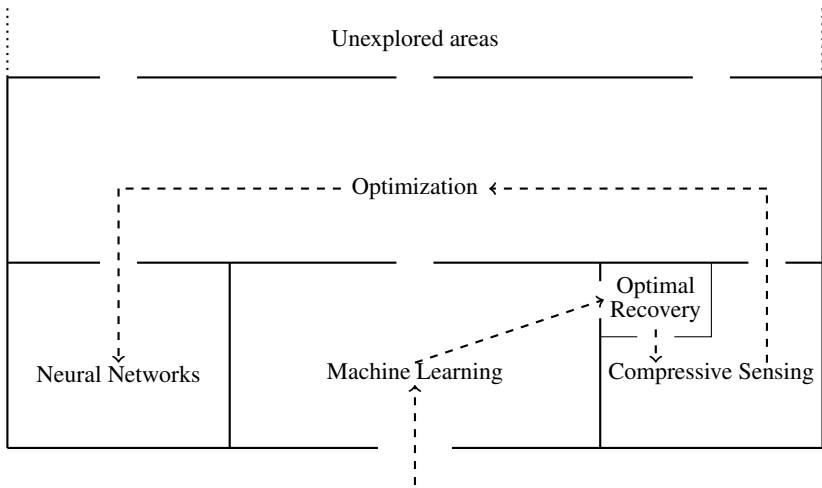


Figure 0.1 Map of the exhibition.

Content: The field of Data Science is too vast to be covered in a single book. As a matter of fact, each of the five parts picked out here is by itself worth a whole book, if not more. The executive summaries for each of the five parts suggest further readings that go into further detail. The selection of topics was merely dictated by my personal taste and interests. The route through the selected topics is metaphorically similar to an exhibition's itinerary; see Figure 0.1. Indeed, we enter through a hall (of Machine Learning) where the brightest lights lead us; continuing our path, we stumble upon a vestibule (of Optimal Recovery) filled with charming but neglected works, before arriving at a chamber (of Compressive Sensing) that we are particularly fond of; pausing

for a while, we realize that the previous rooms were all connected to the large court (of Optimization) that we visit next; finally, as time runs out, we decide to come back later to see some unexplored areas, but we cannot avoid a stop at the most fashionable parlor (of Neural Networks).

Unexplored Areas: Many important topics are left out of this biased overview of mathematical Data Science. They include data assimilation (Law et al., 2015), data streams (Muthukrishnan, 2005), uncertainty quantification (Smith, 2013), reinforcement learning (Sutton and Barto, 2018), and topological data analysis (Dey and Wang, 2022), among others.

Novelty: Some of the topics covered here can already be found in book form in other places. This is particularly true for PART ONE, whose novelty lies mostly in the presentation. Other topics are unlikely to appear elsewhere. For instance, PART TWO is rather uncommon—the content of Chapter 10 there is found only in research articles. PART THREE is original as its presentation relies fully on a modified version of the standard restricted isometry property. This property plays the central role in the exposition of One-Bit Compressive Sensing offered in Chapter 17, which follows the survey article (Foucart, 2017). Most of PART FOUR is rather standard, except Chapter 22 where semidefinite programming techniques are applied to Optimal Recovery. The ingredients of PART FIVE are currently scattered around the literature. Finally, a sizable appendix is included in order to make the text almost self-contained, so that outside references are not required in the main text (with the exception of a few footnotes). It can serve as a toolkit for mathematical scientists who lack a formal training in high-dimensional geometry, probability theory, functional analysis, matrix analysis, or approximation theory. The results recalled in the appendix are of course not new, but some proofs may be innovative (e.g. the von Neumann trace inequality, the Birkhoff theorem). Some other results may not be very familiar (e.g. the Korovkin theorem, the Kolmogorov theorem).

Computational illustrations: Arguably, the field of Data Science would be inconsequential without computations. Although this book focuses on theory, most of its chapters are accompanied by unpretentious implementations, both in MATLAB and in Python. They can be found at

github.com/foucart/Mathematical_Pictures_at_a_Data_Science_Exhibition

Acknowledgment: This book originated from the lecture notes I wrote for a graduate course entitled Topics in Mathematical Data Science and delivered

at Texas A&M University in Fall 2019 and in Fall 2020. Its completion was eased by a course development grant from the Texas A&M Institute of Data Science. The first bricks were actually laid while I was visiting the Institute for Foundations of Data Science at the University of Wisconsin–Madison during a sabbatical semester in Spring 2019. I am indebted to both these institutes for their support. I am also grateful to be associated with various grants from the NSF (DMS-1622134, DMS-1664803, CCF-1934904, DMS-2053172) and from the ONR (N00014-20-1-2787). Finally, I wish to thank a few colleagues for their feedback during the book’s development, namely Radu Balan, Albert Cohen, Rémi Gribonval, Mark Iwen, Philipp Petersen, Sebastien Roch, Jan Vybíral, and Stephan Wojtowytsch.

Notation

Commonly Used Notation

$[i : j]$	the set $\{i, i + 1, \dots, j\}$ of integers from i to j
\mathbb{N}	the set $\{0, 1, 2, \dots\}$ of natural numbers, including 0
\mathbb{N}^*	the set $\{1, 2, 3, \dots\}$ of natural numbers, excluding 0
\mathbb{Z}	the set of integers
\mathbb{Q}	the set of rational numbers
\mathbb{R}	the set of real numbers
\mathbb{C}	the set of complex numbers
i	the imaginary unit $\sqrt{-1}$
$\mathcal{A}, \mathcal{S}, \mathcal{X}$	generic sets
\mathcal{S}^c	the complement of \mathcal{S} (relative to $\mathcal{X} \supseteq \mathcal{S}$, i.e., $\mathcal{X} \setminus \mathcal{S}$)
$\mathcal{S} \Delta \mathcal{S}'$	the symmetric difference $(\mathcal{S} \cup \mathcal{S}') \setminus (\mathcal{S} \cap \mathcal{S}')$ of \mathcal{S} and \mathcal{S}'
$\mathbb{1}_{\text{event}}$	the number equal to 1 if event is true and to 0 otherwise
$\mathbb{1}_{\mathcal{S}}$	the indicator function of a set \mathcal{S} (so that $\mathbb{1}_{\mathcal{S}}(x) = \mathbb{1}_{\{x \in \mathcal{S}\}}$); it can also represent a vector in $\{0, 1\}^n$ when $\mathcal{S} \subseteq [1 : n]$
$ \mathcal{S} $	the cardinality of a finite set \mathcal{S}
F, X	generic vector spaces
$\text{span}(\mathcal{S})$	the linear subspace spanned by a set $\mathcal{S} \subseteq F$
$\text{conv}(\mathcal{S})$	the convex hull of a set $\mathcal{S} \subseteq F$
$\text{Ex}(\mathcal{S})$	the set of extreme points of a set $\mathcal{S} \subseteq F$
$\text{cl}(\mathcal{X})$	the closure of a set $\mathcal{S} \subseteq F$
$\text{vol}(\mathcal{S})$	the volume of a set $\mathcal{S} \subseteq F$
H	a Hilbert space
$\langle x, x' \rangle$	the inner product between two vectors $x, x' \in H$
\mathcal{S}^\perp	the linear space orthogonal to the set $\mathcal{S} \subseteq H$
$P_{\mathcal{V}}$	the orthogonal projector onto the linear subspace \mathcal{V}
T^*	the adjoint of a linear operator T defined on H

xviii	<i>Notation</i>
$\ x\ _F$	the norm of a vector $x \in F$
$\text{dist}_F(x, S)$	the distance from $x \in F$ to a subset S of F
$B(c, r)$	the ball centered at $c \in F$ with radius $r \geq 0$
B_F	the unit ball $B(0, 1)$ of a normed space F
F^*	the dual space of a normed space F
ℓ_p^n	the space \mathbb{R}^n or \mathbb{C}^n normed with $\ x\ _p = \left[\sum_{i=1}^n x_i ^p \right]^{1/p}$
B_p^n	the unit ball of the space ℓ_p^n
(e_1, \dots, e_n)	the canonical basis for \mathbb{R}^n or \mathbb{C}^n
$[x_1; \dots; x_n]$	a column vector with entries x_1, x_2, \dots, x_n
$F(X, \mathcal{Y})$	the space of functions from a set X to a set \mathcal{Y}
$C(X)$	the space of continuous functions from X to \mathbb{R}
$C^k(X)$	the space of k -times continuously differentiable functions
$L_p(X)$	the space of functions with integrable p th power
$W_p^k(X)$	the Sobolev space of functions with k th derivative in $L_p(X)$
\mathcal{K}_{Lip}	the set of functions f with Lipschitz constant $ f _{\text{Lip}} \leq 1$
δ_x	the evaluation functional at a point $x \in X$
μ, ν	generic Borel measures
$\mathcal{M}(X)$	the set of Borel measures on X
$\mathcal{M}_+(X)$	the set of nonnegative Borel measures on X
\mathcal{P}_n	the space of algebraic polynomials of degree $\leq n$
\mathcal{T}_n	the space of trigonometric polynomials of degree $\leq n$
A, B, X	generic matrices
$A_{i,j}$	the entry of a matrix A on the i th row and j th column
A^*	the adjoint of a matrix A , defined by $A_{i,j}^* = \overline{A_{j,i}}$
A^\top	the transpose of a matrix A , defined by $A_{i,j}^\top = A_{j,i}$
$A^{-\top}$	the matrix $(A^{-1})^\top = (A^\top)^{-1}$
$\text{diag}[x_1; \dots; x_n]$	a diagonal matrix with diagonal entries x_1, x_2, \dots, x_n
$\lambda_j(A)$	the j th eigenvalue of A (in nonincreasing order)
$\sigma_j(A)$	the j th singular value of A (in nonincreasing order)
$A \geq 0$	means that the matrix A is positive semidefinite
$\langle A, B \rangle_F$	the Frobenius inner product between $A, B \in \mathbb{R}^{m \times n}$
$\ A\ _F$	the Frobenius norm of $A \in \mathbb{R}^{m \times n}$
$\ A\ _{2 \rightarrow 2}$	the operator norm of $A \in \mathbb{R}^{m \times n}$, i.e., $\sigma_1(A)$
$\ A\ _*$	the nuclear norm $A \in \mathbb{R}^{m \times n}$, i.e., $\sum_{i=1}^m \sigma_i(A)$
$\ker(A)$	the null space of (a matrix or linear map) A
$\text{ran}(A)$	the range of (a matrix or linear map) A
$u * v$	the (discrete or continuous) convolution product of u and v
$\log_2(x)$	the logarithm in base 2 of $x \in (0, +\infty)$
$\ln(x)$	the natural logarithm (in base e) of $x \in (0, +\infty)$
$\exp(x)$	the exponential of $x \in \mathbb{R}$

$\lfloor x \rfloor$	the floor of $x \in \mathbb{R}$, i.e., the integer satisfying $x - 1 < \lfloor x \rfloor \leq x$
$\lceil x \rceil$	the ceiling of $x \in \mathbb{R}$, i.e., the integer satisfying $x \leq \lceil x \rceil < x + 1$
$\mathbb{P}[\mathcal{E}]$	the probability of an event \mathcal{E}
$\mathbb{E}[Z]$	the expectation of a random variable Z
$\mathbb{V}[Z]$	the variance of a random variable Z
$\mathcal{N}(0, \sigma^2)$	the normal distribution with mean zero and variance σ^2
g	a standard gaussian random variable, i.e., $g \sim \mathcal{N}(0, 1)$

Machine-Learning-Specific Notation

\mathcal{X}	the set where the instances $x^{(i)}$ (aka datapoints) live
\mathcal{Y}	the set where the targets y_i (aka observations) live, often $\mathcal{Y} = \mathbb{R}$, $\mathcal{Y} = \{0, 1\}$, or $\mathcal{Y} = \{-1, +1\}$
\mathcal{H}	a hypothesis class, i.e., a subset of $F(\mathcal{X}, \mathcal{Y})$
Loss	a function from $\mathcal{Y} \times \mathcal{Y}$ into $[0, +\infty)$ such that $L(y, y) = 0$
$\text{Risk}_f(h)$	the risk of a predictor $h \in \mathcal{H}$ given $f \in F(\mathcal{X}, \mathcal{Y})$
S	an element of $(\mathcal{X} \times \mathcal{Y})^m$ representing a sample
$\widehat{\text{Risk}}_S(h)$	the empirical risk of $h \in \mathcal{H}$ relative to S
ε_{app}	the approximation error
ε_{est}	the estimation error
$m_{\mathcal{H}}(\varepsilon, \delta)$	the sample complexity
$\text{vc}(\mathcal{H})$	the Vapnik–Chervonenkis dimension of $\mathcal{H} \subseteq F(\mathcal{X}, \{0, 1\})$
K	a kernel, i.e., a symmetric function defined on $\mathcal{X} \times \mathcal{X}$

Optimal-Recovery-Specific Notation

\mathcal{K}	the model set
Q	a quantity of interest, typically a linear map
λ_i	the i th observation functional, often equal to $\delta_{x^{(i)}}$
Λ	the observation map
Δ	a recovery map
$\text{Err}_{\mathcal{K}, Q}(\Lambda, \Delta)$	the worst-case error of Δ (for Q over \mathcal{K} given Λ)
$\text{Err}_{\mathcal{K}, Q}^*(\Lambda)$	the intrinsic error (for Q over \mathcal{K} given Λ)
$\text{Err}_{\mathcal{K}, Q}^0(\Lambda)$	the null error (for Q over \mathcal{K} given Λ)
$\text{Err}_{\mathcal{K}, Q}^*(m)$	the m th minimal intrinsic error (for Q over \mathcal{K})
$m_{\mathcal{K}, Q}(\varepsilon, d)$	the information complexity
$\text{Var}_{HK}(f)$	the variation of f (in the sense of Hardy and Krause)
$\text{Disc}^*(\mathfrak{X})$	the star discrepancy of a finite set \mathfrak{X}

Compressive-Sensing-Specific Notation

Σ_s^N	the set of s -sparse vectors in \mathbb{R}^N or \mathbb{C}^N
$\text{supp}(x)$	the support of a vector x in \mathbb{R}^N or \mathbb{C}^N
S	an index set, i.e., a subset of $[1 : N]$
x_S	the vector x whose entries outside of S are zeroed out
H_s	the hard thresholding operator with parameter s
A	an $m \times N$ observation matrix
A_S	the submatrix of A with columns indexed by S^c removed
\mathcal{A}	an observation map (defined on a space of matrices)
μ	the coherence of a matrix
δ	ℓ_1 -restricted isometry constant of a matrix or linear map
α	lower ℓ_1 -restricted isometry constant of a matrix or linear map
β	upper ℓ_1 -restricted isometry constant of a matrix or linear map
Δ	a recovery map
$d^m(\mathcal{K}, F)$	the m th Gelfand width of a set \mathcal{K} in a space F
$\chi(v)$	the binary vector with i th entry given by $\mathbb{1}_{\{v_i > 0\}}$

Optimization-Specific Notation

$(x^t)_{t \geq 0}$	a sequence of vectors produced by some iterative algorithm
∇f	the gradient of a multivariate function f
L	the Lagrangian of a minimization program
λ, ν	the dual optimization variables, aka Lagrange multipliers
C^*	the dual cone of a set C
$\text{Toep}_\infty(u)$	the infinite symmetric Toeplitz matrix built from $(u_n)_{n \geq 0}$
$\text{Toep}_{N+1}(u)$	the finite symmetric Toeplitz matrix built from $(u_n)_{n=0}^N$

Neural-Networks-Specific notation

ϕ	a generic activation function
ReLU	the rectified linear unit defined by $\text{ReLU}(x) = \max\{x, 0\}$
n_ℓ	the width of the ℓ th layer
$x^{[\ell]}$	the state vector at the ℓ th layer
$W^{[\ell]}$	the weight matrix (in $\mathbb{R}^{n_\ell \times n_{\ell-1}}$) producing the ℓ th layer
$b^{[\ell]}$	the bias vector (in \mathbb{R}^{n_ℓ}) producing the ℓ th layer
\mathcal{N}_ϕ	the linear space of functions generated by shallow networks
\mathcal{N}_ϕ^n	the set of functions generated by width- n shallow networks
$\mathcal{N}_\phi^{n,L}$	the set of functions generated by width- n , depth- L networks