



SAPIENZA
UNIVERSITÀ DI ROMA

SAPIENZA UNIVERSITY OF ROME
DEPARTMENT OF INFORMATION ENGINEERING,
ELECTRONICS AND TELECOMMUNICATIONS

PHD THESIS

**DISTRIBUTED SUPERVISED LEARNING
USING NEURAL NETWORKS**

DISSERTATION SUBMITTED IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE OF DOCTOR OF PHILOSOPHY IN
INFORMATION AND COMMUNICATION ENGINEERING
XXVIII CYCLE

Supervisor
Prof. Aurelio Uncini

Candidate
Simone Scardapane

Rome, Italy
April 2016

Abstract

Distributed learning is the problem of inferring a function in the case where training data is distributed among multiple geographically separated sources. Particularly, the focus is on designing learning strategies with low computational requirements, in which communication is restricted only to neighboring agents, with no reliance on a centralized authority. In this thesis, we analyze multiple distributed protocols for a large number of neural network architectures. The first part of the thesis is devoted to a definition of the problem, followed by an extensive overview of the state-of-the-art. Next, we introduce different strategies for a relatively simple class of single layer neural networks, where a linear output layer is preceded by a nonlinear layer, whose weights are stochastically assigned in the beginning of the learning process. We consider both batch and sequential learning, with horizontally and vertically partitioned data. In the third part, we consider instead the more complex problem of semi-supervised distributed learning, where each agent is provided with an additional set of unlabeled training samples. We propose two different algorithms based on diffusion processes for linear support vector machines and kernel ridge regression. Subsequently, the fourth part extends the discussion to learning with time-varying data (e.g. time-series) using recurrent neural networks. We consider two different families of networks, namely echo state networks (extending the algorithms introduced in the second part), and spline adaptive filters. Overall, the algorithms presented throughout the thesis cover a wide range of possible practical applications, and lead the way to numerous future extensions, which are briefly summarized in the conclusive chapter.

Table of Contents

Abstract	i
List of Acronyms	vi
List of Figures	x
List of Tables	xi
List of Algorithms	xii
1 Introduction	1
Structure of the Thesis	3
Research Contributions	5
Notation	5
I Background Material	7
2 Centralized Supervised Learning	8
2.1 General definitions	8
2.2 ANN models for SL	11
2.2.1 Linear neuron	11
2.2.2 Fixed nonlinear projection	12
2.2.3 Kernel methods	13
2.2.4 Multiple adaptable hidden layers	15
3 Distributed Learning: Formulation and State-of-the-art	17
3.1 Formulation of the Problem	18
3.2 Categorization of DL algorithms	19
3.3 Relation to other research fields	20
3.4 State-of-the-art	22
3.4.1 Distributed linear regression	22
3.4.2 Diffusion filtering and adaptation	23
3.4.3 Distributed sparse linear regression	25
3.4.4 Distributed linear models with a fixed nonlinear projection layer	26

3.4.5	Kernel filtering on sensor networks	26
3.4.6	Distributed support vector machines	28
3.4.7	Distributed multilayer perceptrons	30
II	Distributed Training Algorithms for RVFL Networks	31
4	Distributed Learning for RVFL Networks	32
4.1	Basic concepts of RVFL networks	33
4.1.1	An historical perspective on RVFL networks	34
4.1.2	On the effectiveness of random-weights ANNs	34
4.2	Distributed training strategies for RVFL networks	35
4.2.1	Consensus-based distributed training	35
4.2.2	ADMM-based distributed training	36
4.3	Experimental Setup	39
4.3.1	Description of the Datasets	39
4.3.2	Algorithms and Software Implementation	40
4.4	Results and Discussion	42
4.4.1	Accuracy and Training Times	42
4.4.2	Effect of Network Topology	44
4.4.3	Early Stopping for ADMM	46
4.4.4	Experiment on Large-Scale Data	46
5	Extending Distributed RVFL Networks to a Sequential Scenario	49
5.1	Derivation of the algorithm	49
5.2	Experiments on Distributed Music Classification	51
5.2.1	The Distributed Music Classification Problem	51
5.2.2	Experiment Setup	52
5.2.3	Results and Discussion	54
5.3	Comparison of DAC strategies	56
5.3.1	Description of the strategies	57
5.3.2	Experimental Results	59
6	Distributed RVFL Networks with Vertically Partitioned Data	61
6.1	Derivation of the algorithm	61
6.2	Experimental setup	64
6.3	Results and discussion	65
III	Distributed Semi-Supervised Learning	68
7	Decentralized Semi-supervised Learning via Privacy-Preserving Matrix Completion	69

7.1	Introduction	69
7.2	Preliminaries	72
7.2.1	Semi-supervised learning	72
7.2.2	(Euclidean) matrix completion	74
7.2.3	Privacy-preserving similarity computation	75
7.3	Distributed Laplacian Estimation	77
7.3.1	Formulation of the problem	77
7.3.2	Decentralized block estimation	78
7.3.3	Diffusion gradient descent	80
7.4	Distributed Semi-supervised Manifold Regularization	81
7.5	Experimental results	84
7.5.1	Experiments setup	84
7.5.2	Distributed Laplacian estimation	84
7.5.3	Distributed semi-supervised manifold regularization	86
7.5.4	Privacy preservation	88
8	Distributed Semi-Supervised Support Vector Machines	91
8.1	Introduction	91
8.2	Semi-Supervised Support Vector Machines	93
8.3	Distributed learning for S^3VM	94
8.3.1	Formulation of the problem	95
8.3.2	Solution 1: Distributed gradient descent	95
8.3.3	Solution 2: In-network successive convex approximation	97
8.4	Experimental Results	99
8.4.1	Experimental Setup	99
8.4.2	Results and discussion	101
IV	Distributed Learning from Time-Varying Data	107
9	Distributed Training for Echo State Networks	108
9.1	Introduction	108
9.2	A primer on ESNs	109
9.3	Distributed training for ESNs	110
9.4	Experimental Setup	112
9.4.1	Description of the Datasets	112
9.4.2	Description of the Algorithms	114
9.4.3	ESN Architecture	115
9.5	Experimental Results	116
9.6	Extension to ESNs with Sparse Readouts	119
9.6.1	Comparisons in the centralized case	120
9.6.2	Comparisons in the distributed case	120

10 Diffusion Spline Filtering	125
10.1 Introduction	125
10.2 Spline Adaptive Filter	126
10.3 Diffusion SAF	129
10.4 Experimental Setup	132
10.5 Experimental Results	133
10.5.1 Experiment 1 - Small Network ($L = 10$)	133
10.5.2 Experiment 2 - Large Network ($L = 30$)	135
10.5.3 Experiment 3 - Strong nonlinearity ($L = 15$)	136
V Conclusions and Future Works	139
11 Conclusions and Future Works	140
Appendices	144
A Elements of Graph Theory	145
A.1 Algebraic graph theory	145
A.2 Decentralized average consensus	146
B Software Libraries	148
B.1 Lynx MATLAB Toolbox (Chapters 4-6)	148
B.2 Additional software implementations	149
B.2.1 Distributed LapKRR (Chapter 7)	149
B.2.2 Distributed S ³ VM (Chapter 8)	149
B.2.3 Distributed ESN (Chapter 9)	150
B.2.4 Diffusion Spline Filtering (Chapter 10)	150
Acknowledgments	151
References	165

List of Acronyms

AMC	Automatic Music Classification
ANN	Artificial Neural Network
ATC	Adapt-Then-Combine
BP	Back-Propagation
BRLS	Blockwise Recursive Least Square (see also RLS)
CTA	Combine-Then-Adapt (see also ATC)
DA	Diffusion Adaptation
DAC	Decentralized Average Consensus
DAI	Distributed Artificial Intelligence
DF	Diffusion Filtering
DGD	Distributed Gradient Descent (see also DA)
DL	Distributed Learning
DSO	Distributed Sum Optimization
EDM	Euclidean Distance Matrix
ESN	Echo State Network
ESP	Echo State Property
FL	Functional Link
GD/SGD	(Stochastic) Gradient Descent
HP	Horizontally Partitioned (see also VP)
KAF	Kernel Adaptive Filtering
KRR	Kernel Ridge Regression
LASSO	Least Angle Shrinkage and Selection Operator
LIP	Linear-In-the-Parameters
LMS	Least Mean Square
LRR	Linear Ridge Regression (see also KRR)

MEB	Minimum Enclosing Ball
MFCC	Mel Frequency Cepstral Coefficient
MIR	Music Information Retrieval
ML	Machine Learning
MLP	Multilayer Perceptron
MR	Manifold Regularization
MSD	Mean-Squared Deviation
MSE	Mean-Squared Error
NEXT	In-Network Nonconvex Optimization
NRMSE	Normalized Root Mean-Squared Error (see also MSE)
P2P	Peer-to-Peer
PSD	Positive Semi-Definite
QP	Quadratic Programming
RBF	Radial Basis Function
RKHS	Reproducing Kernel Hilbert Space
RLS	Recursive Least Square
RNN	Recurrent Neural Network
RVFL	Random Vector Functional-Link (see also FL)
SAF	Spline Adaptive Filter
SL	Supervised Learning
SSL	Semi-Supervised Learning
SV	Support Vector (see also SVM)
S³VM	Semi-Supervised Support Vector Machine
SVM	Support Vector Machine
VP	Vertically Partitioned (see also HP)
WLS	Weighted Least Square
WSN	Wireless Sensor Network

List of Figures

1.1	Schematic organization of the algorithms presented in the thesis.	3
2.1	Architecture of an ANN with one fixed hidden layer and a linear output layer.	12
2.2	Architecture of an MLP with T hidden layers and a linear output layer with a single output neuron.	15
3.1	Schematic depiction of DL in a network of 4 agents.	18
3.2	Example of a diffusion step for the first node in the 4-nodes network of Fig. 3.1.	24
3.3	Example of cascade SVM in a network with 6 nodes.	29
4.1	Example of network (with 8 nodes) considered in the experimental sections throughout the thesis.	42
4.2	Average error and standard deviation of CONS-RVFL and ADMM-RVFL on four datasets, when varying the number of nodes in the network from 5 to 50.	43
4.3	Average training time for CONS-RVFL and ADMM-RVFL on a single node.	43
4.4	Consensus iterations needed to reach convergence in CONS-RVFL when varying the network topology.	45
4.5	Relative decrease in error of ADMM-RVFL with respect to L-RVFL, when using an early stopping procedure at different iterations.	46
4.6	Average misclassification error and training time of CONS-RVFL and ADMM-RVFL on the CIFAR-10 dataset, when varying the nodes of the network from 2 to 12.	47
5.1	Evolution of the testing error for the sequential S-CONS-RVFL after every iteration.	56
5.2	Training time required by the sequential S-CONS-RVFL, for varying sizes of the network, from 2 to 14 by steps of 2.	57
5.3	Number of consensus iterations required to reach convergence in the S-CONS-RVFL, when varying the number of nodes in the network from 2 to 14.	57

5.4	Evolution of the DAC iterations required by four different strategies, when processing successive amounts of training batches.	59
5.5	Evolution of the relative network disagreement for four different DAC strategies as the number of DAC iterations increases.	60
6.1	Schematic description of the proposed algorithm for training an RVFL with vertically partitioned data.	62
6.2	Evolution of the error for VP-ADMM-RVFL and ENS-RVFL when varying the size of the network from $L = 4$ to $L = 12$	67
7.1	Depiction of distributed SSL over a network of agents.	70
7.2	Average EDM completion error of the two EDM completion strategies on the considered datasets.	85
7.3	Average EDM completion time required by the two EDM completion strategies on the considered datasets.	87
7.4	Average classification error of the privacy-preserving transformations on the considered datasets when varying the ratio m/d	89
8.1	The hinge loss approximation is shown in blue for varying values of $f(\mathbf{x}_i)$, while in dashed red we show the approximation given by $\exp\{-5f(\mathbf{x}_i)^2\}$	94
8.2	Convergence behavior of DG-VS3VM and NEXT-VS3VM, compared to C-VS3VM.	103
8.3	Box plots for the classification accuracy of the centralized and distributed S^3 VM algorithms.	104
8.4	Training time and test error of GD-VS3VM and NEXT-VS3VM when varying the number of nodes in the network from $L = 5$ to $L = 40$	105
9.1	Schematic depiction of a ESN with multiple outputs.	110
9.2	Evolution of the testing error for ADMM-ESN, for networks going from 5 agents to 25 agents.	117
9.3	Evolution of the training time for ADMM-ESN, for networks going from 5 agents to 25 agents.	118
9.4	Evolution of test error, training time and sparsity when testing L1-ESN.	121
9.5	Evolution of test error, training time and sparsity when testing ADMM-L1-ESN.	124
9.6	Evolution of the (primal) residual of ADMM-L1-ESN for $L = 5$ and $L = 15$	124
10.1	Example of spline interpolation scheme.	127
10.2	Schematic depiction of SAF interpolation performed over a network of agents.	130

10.3	Nonlinear distortion applied to the output signal in experiments 1 and 2 for testing D-SAF.	132
10.4	Dataset setup for the first experiment of D-SAF.	134
10.5	Average MSE evolution for experiment 1 of D-SAF.	135
10.6	MSD evolution for experiment 1 of D-SAF.	136
10.7	Final estimation of the nonlinear model in experiment 1 of D-SAF. .	137
10.8	Average MSE evolution for experiment 2 of D-SAF.	137
10.9	Final estimation of the nonlinear model in experiment 2 of D-SAF. .	138
10.10	Average MSE evolution for experiment 3 of D-SAF.	138
10.11	Final estimation of the nonlinear model in experiment 3 of D-SAF. .	138

List of Tables

1.1	Schematic overview of the research contributions related to the thesis	5
4.1	General description of the datasets for testing CONS-RVFL and ADMM-RVFL	39
4.2	Optimal parameters found by the grid-search procedure for CONS-RVFL and ADMM-RVFL	41
5.1	General description of the datasets for testing the sequential S-CONS-RVFL algorithm.	52
5.2	Optimal parameters found by the grid-search procedure for S-CONS-RVFL	54
5.3	Final misclassification error and training time for the sequential S-CONS-RVFL algorithm, together with one standard deviation. . .	55
6.1	Misclassification error and training time for VP-ADMM-RVFL. . . .	65
7.1	Description of the datasets used for testing Distr-LapKRR.	84
7.2	Values for the parameters used in the simulations of Distr-LapKRR.	86
7.3	Average values for classification error and computational time, together with one standard deviation, for Distr-LapKRR and comparisons.	88
8.1	Description of the datasets used for testing the distributed S^3 VM. .	100
8.2	Optimal values of the parameters used in the experiments for the distributed S^3 VM.	101
8.3	Average value for classification error and computational time for the centralized SVMs.	102
9.1	Optimal parameters found by the grid-search procedure for testing ADMM-ESN.	116
9.2	Final misclassification error and training time for C-ESN, provided as a reference, together with one standard deviation.	116
9.3	The results of Fig. 9.4, shown in tabular form, together with one standard deviation.	122

List of Algorithms

4.1	CONS-RVFL: Consensus-based training for RVFL networks (k th node).	36
4.2	ADMM-RVFL: ADMM-based training for RVFL networks (k th node).	39
5.1	S-CONS-RVFL: Extension of CONS-RVFL to the sequential setting (k th node).	51
6.1	VP-ADMM-RVFL: Extension of ADMM-RVFL to vertically partitioned data (k th node).	64
7.1	Distr-LapKRR: Pseudocode of the proposed distributed SSL algorithm (k th node).	83
8.1	Distributed ∇S^3VM using a distributed gradient descent procedure.	96
8.2	Distributed ∇S^3VM using the In-Network Convex Optimization framework.	99
9.1	ADMM-ESN: Local training algorithm for ADMM-based ESN (k th node).	112
10.1	SAF: Summary of the SAF algorithm with first-order updates. . . .	129
10.2	D-SAF: Summary of the D-SAF algorithm (CTA version).	132

Introduction

SUPERVISED learning (SL) is the task of automatically inferring a mathematical function, starting from a finite set of examples [67]. Together with unsupervised learning and reinforcement learning, it is one of the three main subfields of machine learning (ML). Its roots as a scientific discipline can be traced back to the introduction of the first fully SL algorithms, namely the perceptron rule around 1960 [141], and the k -nearest neighbors (k -NN) in 1967 [39].¹ The perceptron, in particular, became the basis for a wider family of models, which are known today as artificial neural networks (ANNs). ANNs model the unknown desired relation using the interconnection of several building blocks, denoted as artificial neurons, which are loosely inspired to the biological neuron. Over the last decades, hundreds of variants of ANNs, and associated learning algorithms, have been proposed. Their development was sparked by a few fundamental innovations, including the Widrow-Hoff algorithm in 1960 [194], the popularization of the back-propagation (BP) rule in 1986 [142] (and its later extension for dynamical systems [191]), the support vector machine (SVM) in 1992 [17], and additional recent developments on ‘deep’ architectures from 2006 onwards [158].

As a fundamentally *data-driven* technology, SL has been changed greatly by the impact of the so-called ‘big data’ revolution [197]. Big data is a general terminology, which is used to refer to any application where data cannot be processed using ‘conventional’ means. As such, big data is not defined axiomatically, but only through its possible characteristics. These include, among others, its volume and speed of arrival. Each of these aspects has influenced SL theory and algorithms [197], although in many cases solutions were developed prior to the emergence of the big data paradigm itself. As an example, handling large volumes of data is known in the SL community as the large-scale learning problem [18]. This has brought forth multiple developments in parallel solutions for training SL models [60], particularly with the use of *commodity computing* frameworks such as MapReduce [37]. Similarly, learning with continuously arriving streaming data is

¹While this is a generally accepted convention, one may easily choose earlier works to denote a starting point, such as the work by R. A. Fisher on linear discriminant analysis in 1936 [54].

at the center of the subfield of online SL [206].

In this thesis, we focus on another characteristic of several real-world big data applications, namely, their *distributed* nature [197]. In fact, an ever-increasing number of authors is starting to consider this last aspect as a defining property of big data in many real-world scenarios, which complements the more standard characteristics (e.g. volume, velocity, etc.). As an example, Wu *et al.* [197] state that “*autonomous data sources with distributed and decentralized controls are a main characteristic of Big Data applications*”. In particular, we consider the case where training data is distributed among a network of interconnected agents, a setting denoted as distributed learning (DL). If we assume that the agents can communicate with one (or more) coordinating nodes, then it is possible to apply a number of parallel SL algorithm, such as those described above. In this thesis, however, we focus on a more general setting, in which nodes can communicate exclusively with a set of neighboring agents, but none of them is allowed to coordinate in any way the training process. This is a rather general formalization, which subsumes multiple applicative domains, including learning on wireless sensor networks (WSNs) [9, 129], peer-to-peer (P2P) networks [42], robotic swarms, smart grids, distributed databases [87], and several others.

More specifically, we develop distributed SL algorithms for multiple classes of ANN models. We consider first the standard SL setting, and then multiple extensions of it, including online learning [206], semi-supervised learning (SSL) [31], and learning with time-varying signals [184]. Due to the generality of our setting, we assume that computational constraints may be present at each agent, such as a sensor in a WSN. Thus, we focus mostly on relatively simple classes of ANNs, where both training and prediction can be performed with possibly low computational capabilities. In particular, in the first and last part of this thesis, we will be concerned with two-layered ANNs, where the weights of the first layer are stochastically assigned from a predefined probability distribution. These include random vector functional-link (RVFLs) [77, 121], and echo state networks (ESNs) [104]. Despite this simplification, these models are capable of high accuracies in most real-world settings. Additional motivations and an historical perspective on this are provided in the respective chapters. The rest of the thesis deals with linear SVMs, kernel ridge regression, and single neurons with flexible activation functions.

Another important point to note is that the algorithms developed here do not require the exchange of examples between the nodes, but only of a finite subset of parameters of the ANN itself (and a few auxiliary variables).² Thus, they are

²With the exception of Chapter 7, where nodes are allowed to compute a small subset of similarities between their training samples. As is shown in the chapter, privacy can still be kept with the use of privacy-preserving protocols for computing Euclidean distances [95].

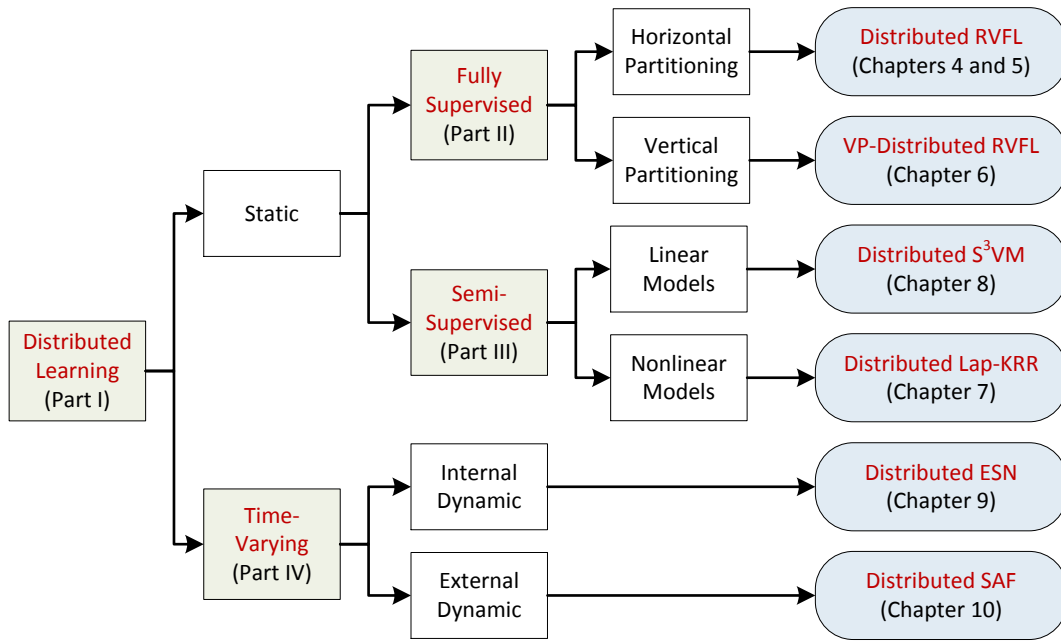


Figure 1.1: Schematic organization of the algorithms presented in the thesis.

able to scale easily to large, and possibly time-varying, networks, while keeping a fixed communication overhead. Constraining the exchange of data points is a reasonable assumption in big data scenarios, where datasets are generally large. However, it might be desirable even in other contexts, e.g. whenever privacy concerns are present [186]. A prototypical example in this case is that of distributed *medical* databases, where sensible information on each patient must go through strict controls on its diffusion [181].

Structure of the Thesis

A schematic categorization of the algorithms presented in the thesis is given in Fig. 1.1. A group corresponding to a specific part of the thesis is shown with a green rectangle, while the algorithms (together with their corresponding chapters) are given with light blue rectangles.

Part I is devoted to introducing the required background material.

Chapter 2 describes the basic tools of SL theory in the centralized case. It starts by stating formally the SL problem, and then moves on to introduce the ANN models to be used successively.

Chapter 3 provides a formal definition of the distributed (multi-agent) SL problem. Additionally, we provide an in-depth overview of previous works dealing with distributed SL with ANN models. The overview combines

works coming from multiple research fields, and tries to give an unified discussion by taking a model-based approach.

Part II introduces algorithms for training RVFL networks in the DL setting.

In **Chapter 4**, we develop two fully distributed training algorithms for them. Strength and weaknesses of both approaches are analyzed and compared to the pre-existing literature.

Chapter 5 extends one algorithm presented in the previous chapter to the sequential setting, where new data is arriving continuously at every node. Additionally, we present an application to the problem of distributed music classification, and we analyze how different strategies for computing a distributed average can influence the convergence time of the algorithm.

Chapter 6 presents a second extension of Chapter 4, to the situation where each example is partitioned across multiple nodes. Technically, this is known as ‘vertical partitioning’ in the data mining literature.

In **Part III** we consider distributed SL with the presence of additional *unlabeled* data at every node, thus extending the standard theory of SSL [31]. This part focuses on kernel models. To the best of our knowledge, these are the first algorithms for general purpose distributed SSL.

In **Chapter 7** we provide a distributed protocol for a kernel-based algorithm belonging to the manifold regularization (MR) framework [11]. To this end, we also derive a novel algorithm for decentralized Euclidean distance matrix (EDM) completion, inspired to the theory of diffusion adaptation (DA) [145].

In **Chapter 8** we propose two distributed algorithms for a family of semi-supervised linear SVMs, derived from the transductive literature. The first algorithm is again inspired to the DA theory, while the second builds on more recent developments in the field of distributed non-convex optimization.

Part IV considers the more general setting of DL in a *time-varying* scenario.

In **Chapter 9**, we exploit a well-known recurrent extension of the RVFL network, called ESN [104]. We leverage on this to provide an extension of one algorithm presented in Chapter 4. The algorithm is then tested on four large-scale prediction problems. We also present an extension for training ESNs with a sparse output layer.

Then, **Chapter 10** considers learning from time-varying signals with the use of particular neurons with flexible nonlinear activation functions, called

spline adaptive filters (SAF) [155]. Again, the theory of DA is used to derive a fully distributed training protocol for SAFs, with local interactions between neighboring nodes. It requires only a small, fixed overhead with respect to a linear counterpart.

Finally, **Chapter 11** summarizes the main contributions of this thesis, along with the possible further developments. The thesis is complemented by two appendices. In **Appendix A**, we provide a general overview of algebraic graph theory (which is used to model networks of agents), and of the decentralized average consensus (DAC) protocol [119, 199]. DAC is a flexible routine to compute global averages over a network, which is used extensively throughout the thesis, including Chapters 4-6 and Chapter 9. **Appendix B** details the open-source software implementations which can be used to replicate the algorithms presented here.

Research contributions

Part of this thesis is adapted from material published (or currently under review) on several journals and conferences. Table 1.1 shows a global overview, while an introductory footnote on each chapter provides more information whenever required.

Table 1.1: Schematic overview of the research contributions related to the thesis

Part II	Chapter 4	Published on <i>Information Sciences</i> [154]
	Chapter 5	Presented at the <i>2015 International Joint Conference on Neural Networks</i> [148]; one section is published as a book chapter in [53]
	Chapter 6	Presented at the <i>2015 INNS Conference on Big Data</i> [150]
Part III	Chapter 7	Conditionally accepted at <i>IEEE Transactions on Neural Networks and Learning Systems</i>
	Chapter 8	Published in <i>Neural Networks</i> [151]; final section is in final editorial review at <i>IEEE Computational Intelligence Magazine</i>
Part IV	Chapter 9	Published in <i>Neural Networks</i> [153]
	Chapter 10	Submitted for presentation at the <i>2016 European Signal Processing Conference</i>

Notation

Throughout the thesis, vectors are denoted by boldface lowercase letters, e.g. \mathbf{a} , while matrices are denoted by boldface uppercase letters, e.g. \mathbf{A} . All vectors are assumed to be column vectors, with \mathbf{a}^T denoting the transpose of \mathbf{a} . The notation A_{ij}

denotes the (i, j) th entry of matrix \mathbf{A} , and similarly for vectors. $\|\mathbf{a}\|_p$ is used for the L_p -norm of a generic vector \mathbf{a} . For $p = 2$ this is the standard Euclidean norm, while for $p = 1$ we have $\|\mathbf{a}\|_1 = \sum_i a_i$. The notation $a[n]$ is used to denote dependence with respect to a time-instant, both for time-varying signals (in which case n refers to a time-instant) and for elements in an iterative procedure (in which case n is the iteration's index). The spectral radius of a generic matrix \mathbf{A} is $\rho(\mathbf{A}) = \max_i \{|\lambda_i(\mathbf{A})|\}$, where $\lambda_i(\mathbf{A})$ is the i th eigenvalue of \mathbf{A} . Finally, we use $\mathbf{A} \geq 0$ to denote a positive semi-definite (PSD) matrix, i.e. a matrix for which $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ for any vector \mathbf{x} of suitable dimensionality.

Part I

Background Material

Centralized Supervised Learning

Contents

2.1	General definitions	8
2.2	ANN models for SL	11
2.2.1	Linear neuron	11
2.2.2	Fixed nonlinear projection	12
2.2.3	Kernel methods	13
2.2.4	Multiple adaptable hidden layers	15

THIS chapter is devoted to the exposition of the basic concepts of SL in the centralized (single-agent) case. It starts with the formalization of the SL problem, using standard tools from regularization theory, in Section 2.1. Next, we introduce the ANN models (and associated learning algorithms) that are used successively, going from the simplest one (i.e., a linear regression) to a more complex multilayer perceptron (MLP). The exposition focuses on a few fundamental concepts, without going into the details on consistency, stability, and so on. The interested reader is referred to any introductory book on the subject for a fuller treatment, e.g. [67].

2.1 General definitions

SL is concerned on automatically extracting a mathematical relation between an *input space* \mathcal{X} , and an *output space* \mathcal{Y} . Throughout the thesis, we assume that the input is a d -dimensional vector of real numbers, i.e. $\mathcal{X} \subseteq \mathbb{R}^d$. The input \mathbf{x} is also called *example* or *pattern*, while a single element x_i of \mathbf{x} is called a *feature*. For ease of notation, we also assume that the output is a single scalar number, such that $\mathcal{Y} \subseteq \mathbb{R}$. However, everything that follows can be extended straightforwardly to the case of a multi-dimensional output vector. It is worth noting here that many representations can be transformed to a vector of real numbers through suitable pre-processing procedures, including categorical variables, complex inputs, texts, sequences, and

so on. Hence, restricting ourselves to this case is a reasonable assumption in most real-world applications. Possible choices for the output space are discussed at the end of this section.

Generally speaking, in a stationary environment, it is assumed that the relation between \mathcal{X} and \mathcal{Y} can be described in its entirety with a joint probability distribution $p(\mathbf{x} \in \mathcal{X}, y \in \mathcal{Y})$. This probabilistic point of view takes into account the fact that the entries in the input \mathbf{x} may not identify univocally a single output y , that noise may be present in the measurements, and so on. The only information we are given is in the form of a *training* dataset of N samples of the relation:

Definition 1 (Dataset)

A dataset D of size N is a collection of N samples of the unknown relation, in the form $D = \{\mathbf{x}_i, y_i\}_{i=1}^N$. The set of all datasets of size N is denoted as $\mathcal{D}(N)$.

Informally, the task of SL is to infer from D a function $f(\cdot)$ such that $f(\mathbf{x}) \approx y$ for any unseen pair (\mathbf{x}, y) sampled from $p(\mathbf{x}, y)$.¹ This process is denoted as *training*. To make this definition more formal, let us assume that the unknown function belongs to a functional space \mathcal{H} . We refer to each element $f(\cdot) \in \mathcal{H}$ as an *hypothesis* or, more commonly, as a *model*. Consequently, we call \mathcal{H} the hypothesis (or model) space. Additionally, a non-negative *loss* function $l(y, f(\mathbf{x})) : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ is used to determine the error incurred in estimating $f(\mathbf{x})$ instead of the true y for any possible pair (\mathbf{x}, y) . Using these elements, we are ready to define the (ideal) SL problem.

Definition 2 (Ideal SL Problem)

Given an hypothesis space \mathcal{H} and a loss function $l(\cdot, \cdot)$, the ideal solution to the SL problem is the function $f(\cdot)$ minimizing the following expected risk functional:

$$I_{exp}[f] = \int l(y, f(\mathbf{x})) dp(\mathbf{x}, y), \quad f \in \mathcal{H}. \quad (2.1)$$

The function f^* minimizing Eq. (2.1) is called Bayes estimator, while $I[f^*]$ is called Bayes risk. Since the probability distribution is unknown, Eq. (2.1) can be

¹The emphasis on *predicting* an output instead of explaining the underlying process distinguishes the ML field from a large part of previous statistics literature, see e.g. [166].

approximated using a generic dataset D :

$$I_{\text{emp}}[f] = \sum_{i=1}^N l(y_i, f(\mathbf{x}_i)) . \quad (2.2)$$

Eq. (2.2) is known as the empirical risk functional. It is relatively easy to show that minimizing Eq. (2.2) instead of Eq. (2.1) may lead to a risk of *overfitting*, i.e., a function which is not able to generalize efficiently to unseen data. A common solution is to include in the optimization process an additional “regularizing” term, imposing reasonable assumptions on the unknown function, such as smoothness, sparsity, and so on. This gives rise to the regularized SL problem.

Definition 3 (Regularized SL problem)

Given a dataset $D \in \mathcal{D}(N)$, an hypothesis space \mathcal{H} , a loss function $l(\cdot, \cdot)$, a regularization functional $\phi[f] : \mathcal{H} \rightarrow \mathbb{R}$, and a scalar coefficient $\lambda > 0$, the regularized SL problem is defined as the minimization of the following functional:

$$I_{\text{reg}}[f] = \sum_{i=1}^N l(y_i, f(\mathbf{x}_i)) + \lambda \phi[f] . \quad (2.3)$$

Problem in Eq. (2.3) can be justified, and analyzed, from a wide variety of viewpoints, including the theory of linear inverse problems, statistical learning theory and Bayes’ theory [41, 51]. Throughout the rest of this thesis, we will be concerned with solving it for different choices of its elements. In particular, we consider models belonging to the class of ANNs. These are briefly summarized in the rest of this chapter, going from the simplest one, a linear neuron trained via least-square regression, to the more complex MLP trained using SGD and the BP rule.

Before this, however, it is necessary to spend a few words on the possible choices for the output space \mathcal{Y} . We distinguish two different cases. In a *regression* task, the output can take any real value in a proper subset of \mathbb{R} . Conversely, in a *binary classification* task, the output can take only two values, which are customarily denoted as -1 and $+1$. More generally, in multi-class classification, the output can assume any value in the set $\{1, 2, \dots, M\}$, where M is the total number of classes. This problem can be addressed by a proper transformation of the output (if the model allows for a multi-dimensional output), or by converting it to a set of binary classification problems, using well-known strategies [140].

2.2 ANN models for SL

2.2.1 Linear neuron

The simplest ANN model for SL is given by the linear neuron, which performs a linear combination of its input vector:

$$f(\mathbf{x}) = \boldsymbol{\beta}^T \mathbf{x} + b, \quad (2.4)$$

where $\boldsymbol{\beta} \in \mathbb{R}^d$ and $b \in \mathbb{R}$. For ease of notation, in the following we drop the bias term b , since it can always be included by considering an additional constant unitary input. A standard choice in this case is minimizing the squared loss $l(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$, subject to an L_2 regularization term on the weights. This gives rise to the well-known linear ridge regression (LRR) problem.

Definition 4 (Linear ridge regression)

Given a dataset $D \in \mathcal{D}(N)$, the LRR problem is defined as the solution to the following optimization problem:

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\beta}\|_2^2 \right\}, \quad (2.5)$$

where $\mathbf{X} = [\mathbf{x}_1^T \dots \mathbf{x}_N^T]^T$ and $\mathbf{y} = [y_1 \dots y_N]^T$.

Assuming that $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})$ is invertible (which is always true for sufficiently large λ), where \mathbf{I} is the identity matrix of suitable dimensionality, the solution of the LRR problem can be expressed in closed form as:

$$\boldsymbol{\beta}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}. \quad (2.6)$$

Work on unregularized LRR dates as far back as Gauss and Legendre [172], and it poses a cornerstone on which most of this thesis is built. It is interesting to note that the effect of the regularization term amounts in adding a fixed scalar value on the diagonal of $\mathbf{X}^T \mathbf{X}$, which is a common heuristic in linear algebra to ensure both the existence of an inverse matrix, and stability in its computation.

Three additional points are worth mentioning here, as they will be used in subsequent chapters. First of all, whenever $N < d$, it is possible to reformulate Eq. (2.6) in order to obtain a computationally cheaper expression:

$$\boldsymbol{\beta}^* = \mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{y}. \quad (2.7)$$

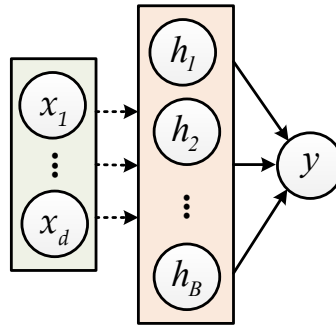


Figure 2.1: Architecture of an ANN with one fixed hidden layer and a linear output layer. Fixed connections are shown with a dashed line, while trainable connections are shown with a solid line.

Secondly, it is possible to modify the standard LRR problem in order to obtain a *sparse* solution, meaning that only a subset of the entries of the optimal weight vector β^* are non-zero. This is achieved by substituting the L_2 norm in Eq. (2.5) with the L_1 norm $\|\beta\|_1$, which provides a convex approximation to the L_0 norm. The resulting algorithm is known as the least absolute shrinkage and selection operator (LASSO) problem [178]. It provides an efficient feature selection strategy, as well as being central to multiple developments in sparse signal processing, including compressed sensing [24]. While the optimization problem of LASSO cannot be solved in closed form anymore, efficient algorithms are available for its solution.

A third aspect that we briefly consider is the use of the LRR problem in binary classification tasks. In this case, in the testing phase, the obtained linear model in Eq. (2.4) is generally binarized using a predefined threshold, making it similar to the original perceptron [141]. The squared loss acts as a convex proxy (or, more technically, as a surrogate loss) of the more accurate misclassification error. Other choices for binary classification might be more accurate, including the hinge loss commonly used in SVMs (introduced in Section 2.2.3), or the logistic loss [67].

2.2.2 Fixed nonlinear projection

Linear models, as described in the previous section, have been widely investigated in the literature due to their simplicity, particularly in terms of training efficiency. Clearly, their usefulness is limited to cases where the assumption of linearity in the underlying process is reasonable. One possibility of maintaining the general theory of linear models, while at the same time obtaining nonlinear modeling capability, is to add an additional *fixed* layer of nonlinearities in front of the linear neuron. This is shown schematically in Fig. 2.1, where fixed and adaptable connections are shown with dashed and solid lines, respectively. In the context of binary classification, the usefulness of such transformations is known since the seminal work of Cover [40].

Mathematically, we consider a model of the form:

$$f(\mathbf{x}) = \sum_{i=1}^B \beta_i h_i(\mathbf{x}) = \boldsymbol{\beta}^T \mathbf{h}(\mathbf{x}) , \quad (2.8)$$

where $\boldsymbol{\beta} \in \mathbb{R}^B$ and we defined $\mathbf{h}(\mathbf{x}) = [h_1(\mathbf{x}), \dots, h_B(\mathbf{x})]^T$. Clearly, Eq. (2.8) is equivalent to a linear model over the transformed vector $\mathbf{h}(\mathbf{x})$, hence it can be trained by considering the linear methods described in the previous section. Due to their characteristics, these models are widespread in SL, including functional link (FL) networks [122], kernel methods (introduced in the next section), radial basis function networks (once the centers are chosen) [124], wavelet expansions, and others [62]. One particular class of FL networks, namely RVFL networks, is introduced in Chapter 4 and further analyzed in Chapters 5 and 6.

2.2.3 Kernel methods

Although the methods considered in the previous section possess good nonlinear modeling capabilities, they may require an extremely large hidden layer (i.e., large B), possibly even infinite. An alternative approach, based on the idea of kernel functions, has been popularized by the introduction of the SVM [17].² The starting observation is that, for a wide range of feature mappings $\mathbf{h}(\cdot)$, there exists a function $\mathcal{K}(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, such that:

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathbf{h}^T(\mathbf{x})\mathbf{h}(\mathbf{x}') \quad \forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d . \quad (2.9)$$

The function $\mathcal{K}(\cdot, \cdot)$ is called a *kernel* function, while Eq. (2.9) is known informally as the kernel trick. It allows transforming any dot product in the transformed space to a function evaluation over the original space. To understand its importance, we first need to introduce a particular class of model spaces.

Definition 5 (Reproducing Kernel Hilbert Space)

A Reproducing Kernel Hilbert Space (RKHS) \mathcal{H} defined over X is an Hilbert space of functions such that any evaluation functional defined as:

$$\mathcal{F}_{\mathbf{x}}[f] = f(\mathbf{x}) \quad \forall f \in \mathcal{H} , \quad (2.10)$$

is linear and bounded.

²The notion of kernel itself was known long before the introduction of the SVM, particularly in statistics and functional analysis, see for example [72, Section 2.3.3].

It can be shown that any RKHS has an associated kernel function. More importantly, solving a regularized SL problem over an RKHS has a fundamental property.

Theorem 1 (Representer's Theorem)

Consider the regularized SL problem in Eq. (2.3). Suppose that \mathcal{H} is an RKHS, and $\phi[f] = \Phi(\|f\|_{\mathcal{H}})$, where $\|f\|_{\mathcal{H}}$ is the norm in the RKHS, and $\Phi(\cdot)$ is a monotonically increasing function. Then, any $f^* \in \mathcal{H}$ minimizing it admits a representation of the form:

$$f^*(\mathbf{x}) = \sum_{i=1}^N \alpha_i \mathcal{K}(\mathbf{x}, \mathbf{x}_i), \quad (2.11)$$

where $\alpha_i, i = 1, \dots, N \in \mathbb{R}$.

Proof 1

See [161]. □

The representer's theorem shows that an optimization problem over a possibly infinite dimensional RKHS is equivalent to an optimization problem over the finite dimensional set of linear coefficients $\alpha_i, i = 1, \dots, N$. SL methods working on RKHSs are known as kernel methods, and we conclude this section by introducing two of them. First, by employing the standard squared loss as error function, and $\Phi(\|f\|_{\mathcal{H}}) = \|f\|_{\mathcal{H}}^2$, we obtain a kernel extension of LRR, which we denote as KRR. Similarly to LRR, the coefficients α of the kernel expansion can be computed in closed form as [51]:

$$\alpha^* = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}, \quad (2.12)$$

where $K_{ij} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$, and \mathbf{K} is called the kernel matrix. KRR is used in Chapter 7 to derive a distributed algorithm for SSL.

In the binary case, an alternative algorithm is given by the SVM, which considers the same squared norm, but substitutes the squared loss with the hinge loss $l(y, f(\mathbf{x})) = \max(0, 1 - yf(\mathbf{x}))$ [171]. In this case, the optimization problem does not allow for a closed-form solution anymore, since it results in a quadratic programming (QP) problem. However, the resulting optimal weight vector is sparse, and the patterns \mathbf{x}_i corresponding to its non-zero elements are called support vectors (SV). Similar formulations can be obtained for regression, such as

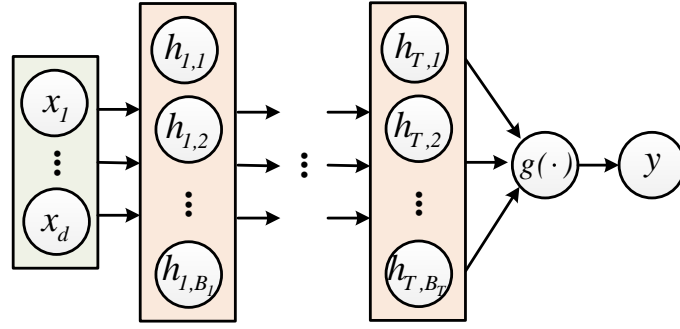


Figure 2.2: Architecture of an MLP with T hidden layers and a linear output layer with a single output neuron. All connections are adaptable.

the ν -SVM and the ε -SVM [171]. Due to the sparseness property of SVs, SVMs have been used extensively in the distributed scenario, as detailed more in depth in the next chapter.

2.2.4 Multiple adaptable hidden layers

The ANN models discussed in Sections 2.2.2 and 2.2.3 are built on a *single* layer of nonlinearities, followed by an adaptable linear layer. While this is enough in many practical situations (and indeed these methods generally possess universal approximation capabilities), more complex real-world applications may require the presence of multiple layers of adaptable nonlinearities, e.g. in the case of classification of multimedia signals [158]. An ANN with these characteristics is called an MLP, and it is shown schematically in Fig. 2.2. In this case, the input vector \mathbf{x} is propagated through T hidden layers. The activation of the i th neuron in the j th layer, with $i \in 1, \dots, B_j$, and $j \in 1, \dots, T$, is given by:

$$h_{i,j}(\mathbf{x}) = \sum_{t=1}^{B_{j-1}} w_{t,j,i} h_{t,j-1}(\mathbf{x}), \quad (2.13)$$

where $h_{i,j}(\cdot)$ is the scalar activation function of the neuron, and we define axiomatically B_0 as $B_0 = d$ and $h_{t,0}$ as $h_{t,0} = x_i$ ($t = 1, \dots, B_0$). In the one-dimensional output case, the output of the MLP is then given by:

$$y = g\left(\sum_{t=1}^{B_T} w_{t,T+1,1} h_{t,T}(\mathbf{x})\right). \quad (2.14)$$

Generally speaking, adapting the full set of weights $\{w_{t,j,i}\}$ results in a non-convex optimization problem, differently from the previous, simpler architectures [67]. This is commonly solved with the use of stochastic gradient descent (SGD), or Quasi-Newton optimization methods, where the error at the output layer can be analytically computed, while it is computed recursively (by back-propagating the

outer error [142]) for the hidden layers.

As we stated in Chapter 1, due to the generality of our distributed setting, in this thesis we focus on the simpler methods described previously, as they provide cheaper algorithms for training and prediction. However, we mention some works on DL for MLPs in the next chapter. Additionally, extending the algorithms presented subsequently to MLPs is a natural future research line, as we discuss in Chapter 11.

Distributed Learning: Formulation and State-of-the-art

Contents

3.1	Formulation of the Problem	18
3.2	Categorization of DL algorithms	19
3.3	Relation to other research fields	20
3.4	State-of-the-art	22
3.4.1	Distributed linear regression	22
3.4.2	Diffusion filtering and adaptation	23
3.4.3	Distributed sparse linear regression	25
3.4.4	Distributed linear models with a fixed nonlinear projection layer	26
3.4.5	Kernel filtering on sensor networks	26
3.4.6	Distributed support vector machines	28
3.4.7	Distributed multilayer perceptrons	30

THIS chapter is devoted to an analysis of the problem of DL using ANN models. We provide a categorization of DL algorithms, in terms of required network topology, communication capabilities, and data exchange. After this, the biggest part of the chapter is devoted to an overview of previous work on DL using ANN models. For readability, the exposition follows the same structure as the previous chapter, i.e. it moves from the simplest ANN model, corresponding to a linear regression, to the more complex MLP. For each model, we describe relative strengths and weaknesses when applied in a distributed scenario. These comments will also serve as motivating remarks for the algorithms introduced in the next chapters. The review aggregates works coming from multiple interdisciplinary fields, including signal processing, machine learning, distributed databases, and several others. When possible, we group works coming from the same research field, in order to provide coherent pointers to the respective literature.

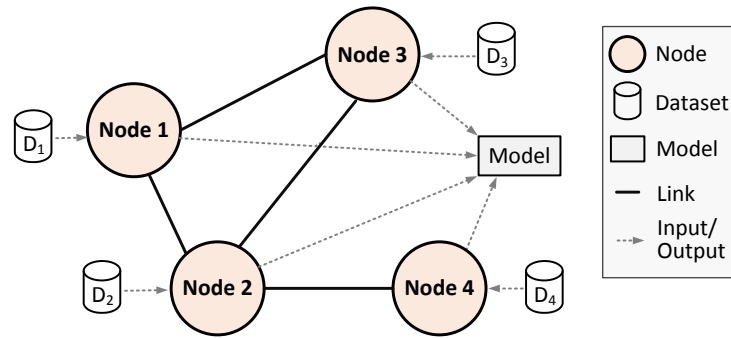


Figure 3.1: DL in a network of agents: training data is distributed throughout the nodes, and all of them must converge to the optimal parameters of a single model. For readability, we assume undirected connections between agents.

3.1 Formulation of the Problem

In the previous chapter, it was assumed that the training dataset S is available on a centralized location for processing. In many contexts, however, this assumption does not hold. As a motivating example for the following, consider the case of a distributed music classification task on a P2P network. Each peer in the network has access to a personal set of labeled songs, e.g., every user has categorized a certain number of its own songs with respect to a predefined set of musical genres. Clearly, solving efficiently this task requires leveraging over *all* local datasets, since we can safely assume that no single dataset alone is sufficient for obtaining adequate performance. Practically, this means that the peers in the network must implement a suitable training protocol for converging to an optimal solution to this DL task. Other examples of DL abound, and a few of them will be mentioned successively.

More formally, we consider the setting described schematically in Fig. 3.1. We have L agents (or nodes), each of which has access to a local training dataset $D_k \in \mathcal{D}(N_k)$, such that $\bigcup_{k=1}^L D_k = D$ and $\sum_{k=1}^L N_k = N$.¹ The connectivity of the agents can be described entirely by a matrix $\mathbf{C} \in \mathbb{R}^{L \times L}$, as detailed in Appendix A. Given these elements, we are now ready to provide a formal definition of the DL problem.

Definition 6 (Distributed learning)

Given L datasets $D_k \in \mathcal{D}(N_k)$, $k = 1, \dots, L$ distributed over a network, an hypothesis space \mathcal{H} , a loss function $l(\cdot, \cdot)$, a regularization functional $\phi[f] : \mathcal{H} \rightarrow \mathbb{R}$, and a scalar coefficient $\lambda > 0$, the distributed learning problem is defined as the minimization

¹In the data mining literature, this is known as ‘horizontal partitioning’ [128]. Chapter 6 considers the complementary case of ‘vertical partitioning’, where the features of every pattern \mathbf{x} are distributed throughout the network.

of the following (joint) functional:

$$I_{dist}[f] = \sum_{k=1}^L \sum_{(\mathbf{x}_i, y_i) \in S_k} l(y_i, f(\mathbf{x}_i)) + \lambda \phi[f]. \quad (3.1)$$

We distinguish between *batch* DL algorithms and *sequential* DL algorithms. In the latter case, each dataset D_k is assumed to be observed in a set of successive batches $D_{k,1}, \dots, D_{k,T}$, such that $D_k = \bigcup_{i=1}^T D_{k,i}$. In the extreme case where each batch is composed of a single element, the resulting formulation is closely linked to the distributed adaptive filtering problem [145]. New batches may arrive synchronously or asynchronously at every agent, as detailed next. The objective in this case is to produce a sequence of estimates $f_{k,1}, \dots, f_{k,T}$ converging as rapidly as possible to the global solution of Eq. (3.1) computed over the overall dataset.

3.2 Categorization of DL algorithms

Despite the generality of the DL setting, existing algorithms can be categorized with respect to a few broad characteristics, which are briefly summarized next.

Coordination Generally speaking, no node is allowed to coordinate specific aspects of the training process, and we assume that there is no shared memory capability. This lack of centralization is in fact the major difference with respect to prior work on parallel SL [60].² Still, some DL algorithms may require the presence of a given subset of dynamically chosen nodes aggregating results from their local neighborhood, such as the clusterheads in a WSN or the super-peers in a P2P network [5].

Connectivity The minimum assumption in DL is that the overall network is connected, i.e. each node can be reached from any other node in a finite number of steps. DL algorithms differentiate themselves on whether they require specific additional properties on the connectivity (e.g. undirected connections). Additionally, some algorithms may assume that the connectivity graph is time-varying.

Communication Distributed training protocols can be categorized based on the communication infrastructure that is required. In particular, messages

²Clearly, there are also important overlaps between parallel SL algorithms and the DL problem considered here, such as the Cascade SVM detailed in Section 3.4.6.

can be exchanged via one-hop or multi-hop connectivity. In multi-hop communication (e.g. ip-based protocols), messages can be routed from any node to any other node, while in single-hop communication, nodes can exchange messages only with their neighbors. At the extreme, each node is allowed to communicate with a single other node at every time slot, as in gossip algorithms [19]. It is easy to understand that multi-hop protocols are not able to efficiently scale to large networks, while they make the design of the algorithm simpler. Similarly, multi-hop communication may not be feasible in particularly unstructured scenarios (e.g. ad-hoc WSNs). This distinction is blurred in some contexts, as it is possible to design broadcast protocols starting from one-hop communication.

Privacy In our context, a privacy violation refers to the need of exchanging local training patterns to other nodes in the network. Algorithms that are designed to preserve privacy are important for two main aspects. First, datasets are generally large, particularly in big data scenarios, and their communication can easily become the main bottleneck in a practical implementation. Secondly, in some context privacy has to be preserved due to the sensitivity of the data, especially in medical applications [181].

Primitives Algorithms can be categorized according to the specific mathematical primitives that are requested on the network. Some algorithms do not require operations in addition to the one-hop exchange. Others may require the possibility of computing vector-sums over the network, Hamiltonian cycles, or even more complex operations. These primitives can then be implemented differently depending on the specific technology of the network, e.g. a sum implemented via a DAC protocol in a WSN [119].

Synchronization Lastly, the algorithms differentiate themselves on whether synchronization among the different agents is required, e.g. in case of successive optimization steps. Most of the literature makes this assumption, as the resulting protocols are easier to analyze and implement. However, designing asynchronous strategies can lead to enormous speed-ups in terms of computational costs and training time.

3.3 Relation to other research fields

Before continuing on to the state-of-the-art, we spend a few words on the relationships between the DL problem and other research fields.

First of all, the problem in Eq. (3.1) is strictly related to a well-known problem in the distributed optimization field, known as distributed sum optimization (DSO). We briefly introduce it here, as its implications will be used extensively in the

subsequent sections. Suppose that the k th agent must minimize a generic function $J_k(\mathbf{w})$ parameterized by the vector \mathbf{w} . In SL, the function can represent a specific form of the loss functional in Eq. (2.3), minimized over the local dataset S_k , and where the vector \mathbf{w} embodies the parameters of the learning model $h \in \mathcal{H}$. DSO is the problem of minimizing the global joint cost function given by:

$$J(\mathbf{w}) = \sum_{k=1}^L J_k(\mathbf{w}). \quad (3.2)$$

Note the relation between Eq. (3.2) and Eq. (3.1). For a single-agent minimizing a differentiable cost function, the most representative algorithm for minimizing it is given by the gradient descent (GD) procedure. In this case, denote by $\mathbf{w}_k[n]$ the estimate of the single node k at the n th time instant. GD computes the minimum of $J_k(\mathbf{w})$ by iteratively updating the estimate as:

$$\mathbf{w}_k[n+1] = \mathbf{w}_k[n] - \eta_k \nabla_{\mathbf{w}} J(\mathbf{w}_k[n]), \quad (3.3)$$

where η_k is the local step-size at time k , whose sequence should be sufficiently small in order to guarantee convergence to the global optimum. Much work on DSO is sparked by the additivity property of the gradient update in the previous equation. In particular, a GD step for the joint cost function in Eq. (3.2) can be computed by summing the gradient contributions from each local node.

Starting from this observation and the seminal work of Tsitsiklis *et al.* [183], a large number of approaches for DSO have been developed. These include [82, 116] for convex unconstrained problems, [20, 48, 174] for convex constrained problems, and [16] for the extension to non-convex problems. Beside GD, representative approaches include subgradient descents [116], dual averaging [48], ADMM [20], and others. Many of these algorithms can be (and have been) applied seamlessly to the setting of DL. For simplicity, in the following we mention only works that have been directly applied or conceived for the DL setting.

In signal processing, instead, the problem of distributed *parametric* inference has a long history [129], and it was revived recently thanks to the interest in large, unstructured WSN networks. Novel approaches in this context are discussed in Section 3.4.2 (linear distributed filtering) and Section 3.4.5 (distributed kernel-based filtering).

More in general, distributed AI (DAI) and distributed problem solving have always been two major themes in the AI community, particularly due to the diffusion of parallel and concurrent programming paradigms [28]. From a philosophical perspective, this is due also to the realization that “*a system may be so complicated and contain so much knowledge that it is better to break it down into different cooperative entities in order to obtain more efficiency*” [28]. Recently, DAI has received renowned

attention in the context of multi-agent systems theory.

Finally, distributed learning has also received attention from the data mining [123] and P2P fields [42]. However, before 2004, almost no work was done in this context using ANN models.

3.4 State-of-the-art

3.4.1 Distributed linear regression

We start our analysis of the state-of-the-art in DL from the LRR algorithm introduced in Section 2.2.1. It is easy to show that this training algorithm is extremely suitable for a distributed implementation. In fact, denote as \mathbf{X}_k and \mathbf{y}_k the input matrix and output vector computed with respect to the k th local dataset. Eq. (2.6) can be rewritten as:

$$\boldsymbol{\beta}^* = \left(\sum_{k=1}^L (\mathbf{X}_k^T \mathbf{X}_k) + \lambda \mathbf{I} \right)^{-1} \sum_{k=1}^L (\mathbf{X}_k^T \mathbf{y}_k) . \quad (3.4)$$

Thus, distributed LRR can be implemented straightforwardly with two sums over the network, the first one on the $d \times d$ matrices $\mathbf{X}_k^T \mathbf{X}_k$, the second one on the d -dimensional vectors $\mathbf{X}_k^T \mathbf{y}_k$. Generally speaking, sums can be considered as primitive on most networks, even the most unstructured ones, e.g. with the use of the DAC protocol introduced in Appendix A.2.

Due to this, the basic idea underlying Eq. (3.4) has been discussed multiple times in the literature. In the following, we consider three representative examples. Karr *et al.* [85] were among the first to exploit it, with the additional use of a secure summation protocol for ensuring data privacy. In the same paper, secure summation is also used to compute diagnostic statistics, in order to confirm the validity of the linear model.

A similar idea is derived in Xiao *et al.* [200, 201], where it is applied to a generalization of LRR denoted as ‘weighted least-square’ (WLS). In WLS, we assume that each output is corrupted by Gaussian noise with mean zero and covariance matrix $\boldsymbol{\Sigma}$. In the centralized case, in the absence of regularization, the solution of the WLS problem is then given by:

$$\boldsymbol{\beta}_{\text{WLS}}^* = (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} . \quad (3.5)$$

In [200], this is solved for a single example at every node using two DAC steps as detailed above. In [201], this is extended to the case of multiple examples arriving asynchronously. In this case, the WLS solution is obtained by interleaving temporal updates with respect to the newly arrived data, with spatial updates corresponding to a single DAC iteration.

Similar concepts are also explored in Bhaduri and Kargupta [14]. As in the

previous case, new data is arriving continuously at every node. Differently than before, however, the LRR solution is recomputed with a global sum only when the error over the training set exceeds a predefined threshold, to reduce communication.

The LRR problem has also been solved in a distributed fashion with the use of distributed optimization techniques, including the subgradient algorithm [136] and the ADMM procedure [20]. These techniques can also be adapted to handle different loss functions for the linear neuron, including the Huber loss [84] and the logistic loss [82].

An alternative approach is followed in [109] for the case where the overall output vector \mathbf{y} is globally known. In this case, each agent projects its local matrix \mathbf{X}_k to a lower-dimensional space with the use of random projections. Next, these projections are broadcasted to the rest of the network. By concatenating the resulting matrices, each agent can independently solve the global LRR problem with bounded error.

If we allow the nodes to exchange data points, Balcan *et al.* [8] and Daumé III *et al.* [43] independently derive bounds on the number of patterns that must be exchanged between the agents for obtaining a desired level of accuracy in the context of probably approximately correct (PAC) theory. As an example, [8, Section 7.3] shows that, if data is separated by a margin γ , the model is trained using the perceptron rule, and the nodes communicate in a round robin fashion, learning the model requires $\mathcal{O}\left(\frac{1}{\gamma^2}\right)$ rounds. Making additional assumptions on the distribution of the data allows to reduce this bound [8].

3.4.2 Diffusion filtering and adaptation

Next, we consider the problem of solving the LRR problem in Eq. (2.5) with continuously arriving data. Additionally, we suppose that the nodes have stringent computational requirements, so that solving multiple times Eq. (2.6) is infeasible. This setting is closely linked to the problem of adaptive filtering in signal processing [184], where the input typically represents a buffer of the last observed samples of an unknown linear process. Two widespread solutions in the centralized case are the least mean square (LMS) algorithm, which is strictly related to the GD minimization procedure, and the recursive least square (RLS) algorithm, which recursively computes Eq. (2.6)³ [184].

In the context of distributed filtering, these algorithms were initially extended using incremental gradient updates [100, 146], where information on the update steps is propagated on a Hamiltonian cycle over the network. This includes incremental LMS [100] and incremental RLS [146]. These methods, however, have a major drawback, in that computing such a cycle is an NP-hard problem. Due to this, an

³The RLS is formally introduced in Chapter 5.

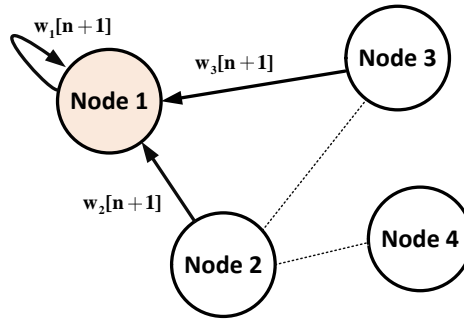


Figure 3.2: Example of a diffusion step for the first node in the 4-nodes network of Fig. 3.1. Links that are deactivated are shown with dashed lines. Note that node 4 is not directly connected to node 1, thus its estimate will only reach it in an indirect way through node 2.

alternative formulation, denoted as diffusion filtering (DF), was popularized in Lopes and Sayed [101] for the LMS and in Cattivelli *et al.* [25] for the RLS. In a DF, local updates are interpolated with ‘diffusion’ steps, where the estimates are locally weighted with information coming from the neighbors. Multiple extensions over this basic scheme have been introduced in the following years, including DF with adaptive combiners [175], total least-square [6], sparse models [47], asynchronous networks [205], and so on.

The popularity of the DF field has led to its application to the wider problem of DSO, under the name of diffusion adaptation (DA) or distributed gradient descent (DGD) [145]. Since DA will be used extensively in Chapters 7 and 10, we briefly detail it here. DA works by interleaving local gradient descents as in Eq. (3.3) with averaging steps given by:

$$\mathbf{w}_k[n+1] = C_{kk}\mathbf{w}_k[n+1] + \sum_{t \in \mathcal{N}_k} C_{kt}\mathbf{w}_t[n+1], \quad (3.6)$$

where the weights C_{kt} have the same meaning as the connectivity matrix of the DAC protocol (see Appendix A.2). In fact, the previous equation can be understood as a single DAC step. An example of a diffusion step is shown in Fig. 3.2. In particular, this strategy is known as adapt-then-combine (ATC), while an equivalent combine-then-adapt (CTA) formulation can be obtained by interchanging the two updates. Different choices of $J_k(\mathbf{w})$ give rise to different algorithms, including the diffusion LMS and RLS mentioned before. For a recent exposition on the theory of DA, its convergence properties and applications to stochastic optimization, see the monograph by Sayed [145], where the author also mentions the application to a diffusion logistic regression in Section V-C.

Before concluding this section, we mention that distributed linear filters without the use of DF theory have also been proposed in the literature. As an example, Schizas *et al.* [156] present a distributed LMS, where ADMM is used to enforce

consensus on a set of ‘bridging’ sensors. A similar formulation for the RLS is derived in Mateos *et al.* [108]. An alternative RLS algorithm, which bypasses the need for bridge sensors, is analyzed instead in Mateos and Giannakis [107].

3.4.3 Distributed sparse linear regression

Distributed training of a sparse linear method has also been investigated extensively in the literature. Mateos *et al.* [106] reformulate the problem of LASSO in a separable form, and then solve it by enforcing consensus constraints with the use of the ADMM procedure. They present three different versions, which differ in the amount of computational resources required by the single node. Particularly, in the orthonormal design case, it is shown that the local update step can be computed by an elementary thresholding operation. Mota *et al.* [114] solve in a similar way a closely related problem, denoted as basis pursuit. In [106, Section V], the authors discuss also a distributed cross-validation procedure for selecting an optimal λ in a decentralized fashion.

An alternative formulation is presented in Chen and Sayed [35], where the L_1 norm is approximated with the twice-differentiable regularization term given by:

$$\|\beta\|_1 \approx \sum_{i=1}^d \sqrt{\beta_i^2 + \varepsilon^2}, \quad (3.7)$$

where ε is a small number. The problem is solved with DA (see Section 3.4.2).

A third approach, based on the method of iterative thresholding, is instead presented in Ravazzi *et al.* [137], for both the LASSO problem and the optimally sparse LRR problem with an L_0 regularization term. Results are similar to [107], but the algorithm requires significantly less computations at every node.

Much work has been done also in the case of sequential distributed LASSO problems. Liu *et al.* [97] extend the standard diffusion LMS with the inclusion of L_0 and L_1 penalties, showing significant improvements with respect to the standard formulation when the underlying vector is sparse. A similar formulation is derived in Di Lorenzo and Sayed [47], with two important differences. First, they consider two different sets of combination coefficients, allowing for a faster rate of convergence. Secondly, they consider an adaptive procedure for selecting an optimal λ coefficient.

In the case of sparse RLS, Liu *et al.* [98] present an algorithm framed on the principle of maximum likelihood, with the use of expectation maximization and thresholding operators. An alternative, more demanding formulation, is presented in Barbarossa *et al.* [9, Section IV-A4], where the optimization problem is solved with the use of the ADMM procedure.

3.4.4 Distributed linear models with a fixed nonlinear projection layer

We now consider DL with nonlinear ANN models, starting from the linear neuron with a fixed nonlinear projection layer introduced in Section 2.2.2. We already remarked that, in the centralized case, this family of models offers a good compromise between speed of training and nonlinear modeling capabilities. In the distributed scenario, however, their use has been limited to a few cases, which are briefly summarized next. This remark, in fact, offers a substantial motivation for the algorithms introduced in the next chapters.

Hershberger and Kargupta [70] analyze a vertically partitioned scenario, where the input is projected using a set of wavelet basis functions. Particularly, they consider dilated and translated instances of the scaling function, denoted as “box” functions. The coefficients of the wavelet representation are then transmitted to a fusion center, which is in charge of computing the global LRR solution.

Sun *et al.* [173] perform DL in a P2P network with an ensemble of extreme learning machine (ELM) networks. In ELM, the parameters of the nonlinear functions in the hidden layer are stochastically assigned at the beginning of the learning procedure (see Section 4.1.1). In [173], an ensemble of ELM functions are handled by a set of ‘super-peers’, using an efficient data structure in order to minimize data exchange over the network. An alternative approach for training an ELM network is presented in Samet and Miri [143], both for horizontally and vertically partitioned data, which makes use of secure protocols for computing vector products and the SVD decomposition. A third approach is presented in Huang and Li [76], where the output layer is trained with the use of diffusion LMS and diffusion RLS (see Section 3.4.2).

3.4.5 Kernel filtering on sensor networks

We now begin our analysis of distributed kernel methods, starting with the KRR algorithm described in Section 2.2.3. On first glance, this algorithm is not particularly suited for a distributed implementation, as the kernel model in Eq. (2.11) requires knowledge of all the local datasets. This is particularly daunting for the implementation of incremental gradient (and subgradient) methods, as is already discussed in Predd *et al.* [129]: *“In consequence, all the data will ultimately propagate to all the sensors, since exchanging [the examples] is necessary to compute [the gradient] and hence to share [the model] (assuming that the sensors are preprogrammed with the kernel).”*. Despite this apparent limitation, much work has been done in this context, particularly for non-parametric inference on WSNs.

Possibly the first investigation in this sense was done in Simić [169]. In a WSN, in many cases, we can assume that the input \mathbf{x} represents the geographical coordinates of the sensor itself, e.g. in the case of sensors measuring a specific field. In this

case, if we use a translational kernel, i.e. a kernel that depends only on Euclidean distances, we have that $\mathcal{K}(\mathbf{x}_1, \mathbf{x}_2) = 0$ for any two sensors which are sufficiently far away. Thus, the resulting kernel matrix \mathbf{K} is sparse. In [169], each node solves its local KRR model, sending its optimal coefficients to a fusion center, which combines them by taking into consideration the previous observation. A similar procedure without the need for a fusion center is presented in Guestrin *et al.* [66], where the problem is solved with a distributed Gaussian elimination procedure. The approach in [66] has strong convergence guarantees and can be used even in cases where the matrix \mathbf{K} is not sparse, albeit losing most of its attractiveness in terms of communication efficiency.

A second approach is proposed for a more general case in Rabbat and Nowak [133], and applied to the KRR algorithm in [129]. The method is based on incrementally passing the subgradient updates over the network, thus it requires the presence of a Hamiltonian cycle. Additionally, as we discussed before, the data must be propagated throughout the network. Hence, this approach is feasible only in specific cases, e.g. when the RKHS admits a lower dimensional parameterization.

A third approach is investigated in Predd *et al.* [131]. The overall DL setting is represented as a bipartite graph. Each node in the first part corresponds to an agent, while each node in the second part corresponds to an example. An edge between the two parts means that a node has access to a given pattern. A relaxed version of the overall optimization problem is solved, by imposing that the agents reach a consensus only on the patterns that are shared among them. Due to this, it is possible to avoid sending the complete datasets, while communication is restricted to a set of state messages. Some extensions, particularly to asynchronous updates, are discussed in Pérez-Cruz and Kulkarni [127]. All the three approaches introduced up to now are summarized in [129].

In the centralized case, solving the KRR optimization problem in a sequential setting has received considerable attention in the field of kernel adaptive filtering (KAF) [96], giving rise to multiple kernel-based extensions of the linear adaptive filters discussed in Section 3.4.2. The fact that the resulting model grows linearly with the number of processed patterns is also one of the main drawbacks of KAFs, where it is known as the ‘growth’ problem. Much work has been done to curtail it [96], and in a limited part it has been extended to the decentralized case. In particular, Honeine *et al.* [73] propose a criterion for discarding examples based on a previously introduced concept of ‘coherence’. Given a set of patterns $\mathbf{x}_1, \dots, \mathbf{x}_m$, the coherence with respect to a new pattern \mathbf{x} is defined as:

$$\min_{\gamma_1, \dots, \gamma_m} \left\| \mathcal{K}(\mathbf{x}, \cdot) - \sum_{i=1}^m \gamma_i \mathcal{K}(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}}. \quad (3.8)$$

In [73], the pattern \mathbf{x} is discarded if the coherence is greater than a certain threshold.

The authors discuss efficient implementations of this idea. In [74], similar ideas are developed for *removing* elements that have already been processed.

Finally, we mention the distributed algorithm presented in Chen *et al.* [36], where the previous setting is extended by considering non-negativity constraints on the model's coefficients. This is particularly important in applications imposing non-negativity constraints on the parameters to estimate.

3.4.6 Distributed support vector machines

As we discussed in Section 2.2.3, an alternative widespread kernel method is the SVM. Intuitively, this algorithm is preferable to KRR for a distributed implementation, due to the sparseness property of the resulting output vector. In fact, the SVs embed all the information which is required from a classification point of view, providing for a theoretically efficient way of compressing information to be sent throughout the network. However, this idea is hindered by a practical problem; namely, the SVs of a reduced dataset may not correspond to the SVs of the entire dataset. More formally, denote by $SV(S)$ the set of SVs obtained by solving the QP problem with dataset S . Given two partitions S_1, S_2 such that $S_1 \cup S_2 = S$, we have:

$$SV(S_1) \cup SV(S_2) \neq SV(S). \quad (3.9)$$

Nonetheless, for a proper subdivision of the dataset, we may expect that the two terms in the previous equation may still share a good amount of SVs. Initial work on distributed SVM was fueled by an algorithm exploiting this idea, the cascade SVM [64], originally developed for parallelizing the solution to the global QP problem. In a cascade SVM, the network is organized in a set of successive layers. Nodes in the first layer receive parts of the input dataset, and propagate forward their SVs. Nodes in the next layers receive the set of SVs from their ancestors, merge them, and solve again the QP problem, up to a final node which outputs a final set of SVs. As stated by the authors, "*Often a single pass through this Cascade produces satisfactory accuracy, but if the global optimum has to be reached, the result of the last layer is fed back into the first layer*" [64]. This is shown schematically in Fig. 3.3.

The first work we are aware of to explore training an SVM in a fully distributed setting, without constraining the network topology as in the cascade SVM, is the PhD thesis by Pedersen [126], which explores informally multiple data exchange protocols in a distributed Java implementation. Another early work was the Distributed Semiparametric Support Vector Machine (DSSVM) presented in [115]. In the DSSVM, every local node selects a number of centroids from the training data. These centroids are then shared with the other nodes, and their corresponding weights are updated locally based on an iterated reweighted least squares (IRWLS) procedure. Privacy can be preserved by adding noise to the elements of the training

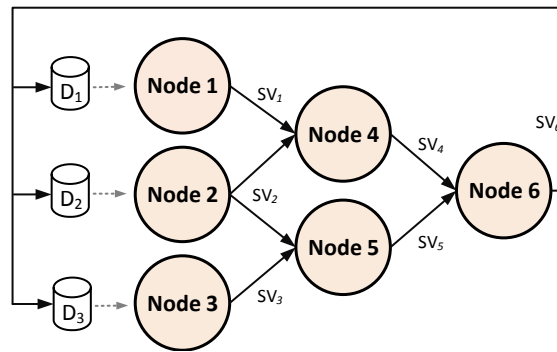


Figure 3.3: Example of cascade SVM in a network with 6 nodes. For readability, SV_k denotes the output SVs of the k th node.

data when selecting the centroids. The DSSVM may be suboptimal depending on the particular choice of centroids. Moreover, it requires incremental passing of the SVs, or centroids, between the nodes, which in turn requires the computation of a Hamiltonian cycle between them. An alternative Distributed Parallel SVM (DPSVM) is presented in [103]. Differently from the DSSVM, the DPSVM does not depend on the particular choice of centroids, and it is guaranteed to reach the global optimal solution of the centralized SVM in a finite number of steps. Moreover, it considers general strongly connected networks, with only exchanges of SVs between neighboring nodes. Still, the need of exchanging the set of SVs reduces the capability of the algorithm to scale to very large networks. A third approach is presented in [58], where the problem is recast as multiple convex subproblems at every node, and solved with the use of the ADMM procedure. Most of the attempts described up to now are summarized in the overview by Wang and Zhou [188].

To conclude this section, we shortly present a set of alternative formulations for distributed SVM implementations that were proposed in the last years. Flouri *et al.* [55] consider the possibility of exchanging SVs cyclically, through a set of ‘clusterheads’ distributed over the network. In [56, 57], the algorithm is refined in order to consider only one-hop communication. Additionally, the authors propose the idea of exchanging only the SVs which are on the borders of the convex hulls for each class, in order to accelerate the convergence speed.

Ang *et al.* [4] combines the idea of the cascade SVM and the Reduced SVM (similar to the Semiparametric SVM in [115]) in order to reduce the communication cost on a P2P network. In [5], the cost is further reduced by considering a bagging procedure at each node.

Hensel and Duta [68] investigate a gradient descent procedure, where the overall gradient update is computed with a ‘Push-Sum’ protocol, a special gossip procedure for computing approximate sums over a network, allowing the algorithm to scale linearly with respect to the size of the agents’ network. Differently, in Wang *et al.* [187], gradient descent is used locally to update the local models, which are

then fused at the prediction phase with a DAC procedure.

Lodi *et al.* [99] explicitly consider the problem of multi-class SVM, exploiting the equivalence between a multi-class extension of SVM and the Minimum Enclosing Ball (MEB) problem (see [99, Section 3.1] for the definition of MEB). Each node computes its own MEB, and forwards the results to a central processor, which is in charge of computing the global solution. The authors mention that this last step can be parallelized by reformulating the algorithms in the Cascade SVM [64] or the DPSVM [103], substituting the SVs with the solution to the MEB problems.

Finally, several authors have considered the use of distributed optimization routines for solving the distributed linear SVM problem (i.e. with a kernel corresponding to dot products in the input space, see Chapter 8). Among these, we may cite the random projection algorithm [89], dual coordinate ascent [81], and the box-constrained QP [88].

3.4.7 Distributed multilayer perceptrons

We conclude this chapter with a brief overview on distributed MLPs. Remember from Section 2.2.4 that MLPs are generally trained with SGD. As we stated previously in Section 3.3, gradient descent is relatively easy to implement in a distributed fashion, due to the additivity of the gradient. In fact, the literature on online learning and prediction has considered multiple distributed implementations of SGD, starting from the work of Zinkevich *et al.* [206, 207], including variants with asynchronous updates [139], and without the need for a parameters' server [45].

However, these ideas have been rarely applied to the distributed training of MLPs, except in a handful of cases. As examples, Georgopoulos and Hasler [61] train it by summing the gradient updates directly with a DAC procedure. Similarly, Schlitter [157] and Samet and Miri [143] investigate the use of secure summation protocols for ensuring privacy during the updates. This scarcity of results has a strong motivation. In fact, relatively large MLPs may possess millions of free parameters, resulting in millions of gradient computations to be exchanged throughout each node, making it impractical (e.g. [164]). This problem has started to being addressed in different ways. The first is model parallelism, where the MLPs itself is split over multiple machines, such as in DistBelief [44]. The other is quantization, where the gradient updates are heavily quantized in order to reduce the communication cost [163]. Additionally, the problem involved in training the MLP is non-convex, making it more complex (both theoretically and practically) to apply the aforementioned ideas.

An alternative approach is to construct an *ensemble* of MLPs, one for each node, in order to avoid the gradient exchanges. Lazarevic and Obradovic [87] were among the first to consider this idea, with a distributed version of the standard AdaBoost algorithm. A similar approach is presented in Zhang and Zhong [204].

Part II

Distributed Training Algorithms for RVFL Networks

Distributed Learning for RVFL Networks

Contents

4.1	Basic concepts of RVFL networks	33
4.1.1	An historical perspective on RVFL networks	34
4.1.2	On the effectiveness of random-weights ANNs	34
4.2	Distributed training strategies for RVFL networks	35
4.2.1	Consensus-based distributed training	35
4.2.2	ADMM-based distributed training	36
4.3	Experimental Setup	39
4.3.1	Description of the Datasets	39
4.3.2	Algorithms and Software Implementation	40
4.4	Results and Discussion	42
4.4.1	Accuracy and Training Times	42
4.4.2	Effect of Network Topology	44
4.4.3	Early Stopping for ADMM	46
4.4.4	Experiment on Large-Scale Data	46

THIS chapter introduces two distributed algorithms for RVFL networks, which are a special case of the fixed hidden layer ANN models presented in Section 2.2.2. As we said in Section 2.2.2, the use of these models is widespread in the centralized case, due to their good trade-off of algorithmic simplicity and nonlinear modeling capabilities. At the same time, as detailed in Section 3.4.4, their use in the DL setting has been relatively limited, which is the main motivation for this chapter. After introducing the RVFL network, we describe the two distributed strategies, based on the DAC protocol and the ADMM optimization algorithm. Next, we evaluate them on multiple real-world scenarios.

The content of this chapter, except Sections 4.1.1 and 4.1.2, is adapted from the material published in [154].

4.1 Basic concepts of RVFL networks

RVFL networks are a particular class of ANN models with a fixed hidden layer, as depicted in Section 2.2.2. Mathematically, their most common variation is given by [121]:

$$f(\mathbf{x}) = \sum_{m=1}^B \beta_m h_m(\mathbf{x}; \mathbf{w}_m) = \boldsymbol{\beta}^T \mathbf{h}(\mathbf{x}; \mathbf{w}_1, \dots, \mathbf{w}_B), \quad (4.1)$$

where the m th transformation is parametrized by the vector \mathbf{w}_m .¹ The parameters $\mathbf{w}_1, \dots, \mathbf{w}_B$ are chosen in the beginning of the learning process, in particular, they are extracted randomly from a predefined probability distribution. Conceptually, this is similar to the well-known technique of random projections [105], which is a common procedure in statistics for dimensionality reduction. Differently from it, however, in RVFL networks the stochastic transformation of the input vector is not required to preserve distances and, more importantly, can increase the dimensionality. In the following, dependence of the hidden functions with respect to the stochastic parameters is omitted for readability. If we define the hidden matrix $\mathbf{H} \in \mathbb{R}^{N \times B}$ as:

$$\mathbf{H} = \begin{pmatrix} h_1(\mathbf{x}_1) & \cdots & h_B(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ h_1(\mathbf{x}_N) & \cdots & h_B(\mathbf{x}_N) \end{pmatrix}, \quad (4.2)$$

it is straightforward to show that training of an RVFL network can be implemented efficiently with the use of the LRR algorithm described in Eq. (2.5) and Eq. (2.6), by substituting the input matrix \mathbf{X} with the hidden matrix \mathbf{H} . The resulting output weights are given by:

$$\boldsymbol{\beta}^* = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^T \mathbf{y}. \quad (4.3)$$

Throughout this chapter (and subsequent ones), we will use sigmoid activation functions given by:

$$h(\mathbf{x}; \mathbf{w}, b) = \frac{1}{1 + \exp\{- (\mathbf{w}^T \mathbf{x} + b)\}}. \quad (4.4)$$

The derivation in this section extends trivially also to the situation of $M > 1$ outputs. In this case, $\boldsymbol{\beta}$ becomes a $B \times M$ matrix and the output vector \mathbf{y} becomes an $N \times M$ matrix, where the i th row corresponds to the M -dimensional output \mathbf{y}_i^T of the training set. Additionally, we replace the L_2 -norm on vectors in Eq. (2.5) with a suitable matrix norm.

¹The original derivation in [121] had additional connections from the input layer to the output layer, however, this is a trivial extension with respect to our formulation.

4.1.1 An historical perspective on RVFL networks

This kind of random-weights ANNs have a long history in the field of SL. The original perceptron, indeed, considered a fixed layer of binary projections, which was loosely inspired to the biological vision [141]. In 1992, Schmidt *et al.* [160] investigated a model equivalent to the RVFL network, however, this was not “presented as an alternative learning method”, but “only to analyse the functional behavior of the networks with respect to learning”. The RVFL network itself was presented by Pao and his coworkers [121] as a variant of the more general functional-link network [122]. In [77], it was shown to be an universal approximator for smooth functions, provided that the weights were extracted in a proper range and B was large enough. In particular, the rate of convergence to zero of the approximation error is $O(\frac{C}{\sqrt{B}})$, with the constant C independent of B . Further analyses of approximation with random bases were obtained successively in [63, 134, 135]. Recently, similar models were popularized under the name extreme learning machine (ELM) [75], raising multiple controversies due to the lack of proper acknowledgment of previous material, particularly RVFL networks [189]. Historically, the RVFL network is also connected to the radial basis function (RBF) network investigated by Broomhead and Lowe in 1988 [102], to the statistical test for neglected nonlinearities presented by White [192], and to the later QuickNet family of networks by the same author [193]. Another similar algorithm has been proposed in 2013 as the ‘no-prop’ algorithm [196].²

4.1.2 On the effectiveness of random-weights ANNs

Despite the stochastic assignment of weights in the first layer, RVFL networks are known to provide excellent performance in many real-world scenario, giving a good trade-off between accuracy and training simplicity. This was shown clearly in a 2014 analysis by Fernandez-Delgado *et al.* [52]. Over 179 classifiers, a kernel-based variation of RVFL networks with RBF functions was shown to be among the top-three performing algorithms over 121 different datasets. In the words of B. Widrow [195]: “we [...] have independently discovered that it is not necessary to train the hidden layers of a multi-layer neural network. Training the output layer will be sufficient for many applications.”. Clearly, randomly selecting bases is at most a naive approach, which can easily be outperformed by proper adaptation of the hidden layer. Worse, in some cases this choice can introduce a large variance in the results, as stated by Principe and Chen [132]: “[random-weights models] still suffer from design choices, translated in free parameters, which are difficult to set optimally with the current mathematical framework, so practically they involve many trials and cross validation to

²Which, ironically, has been criticized for its similarity to the ELM network [195].

find a good projection space, on top of the selection of the number of hidden [processing elements] and the nonlinear functions.”. Although we use RVFL networks for their efficiency in the DL setting, these limitations should be kept in mind.

4.2 Distributed training strategies for RVFL networks

Let us now consider the problem of training an RVFL network in the DL setting. By combining the DL problem in Eq. (3.1) with the RVFL least-square optimization criterion, the global optimization problem of the distributed RVFL can be stated as:

$$\beta^* = \arg \min_{\beta \in \mathbb{R}^B} \frac{1}{2} \left(\sum_{k=1}^L \|\mathbf{H}_k \beta - \mathbf{y}_k\|_2^2 \right) + \frac{\lambda}{2} \|\beta\|_2^2, \quad (4.5)$$

where \mathbf{H}_k and \mathbf{y}_k are the hidden matrix and output vector computed over the local dataset S_k . Remember from Section 3.4.1 that the optimal weight vector in this case can be expressed as:

$$\begin{aligned} \beta^* &= \left(\sum_{k=1}^L (\mathbf{H}_k^T \mathbf{H}_k) + \lambda \mathbf{I} \right)^{-1} \sum_{k=1}^L (\mathbf{H}_k^T \mathbf{y}_k) = \\ &= (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^T \mathbf{y}. \end{aligned} \quad (4.6)$$

This can be implemented in a fully distributed fashion by executing two sequential DAC steps³: the first on the matrices $\mathbf{H}_k^T \mathbf{H}_k$, and the second on $\mathbf{H}_k^T \mathbf{y}_k$. However, since the matrices $\mathbf{H}_k^T \mathbf{H}_k$ have size $B \times B$, this approach is feasible only for small hidden expansions, i.e. small B . Otherwise, the free exchange of these matrices over the network can become a computational bottleneck or, worse, be infeasible. For this reason, we do not consider this idea further in this chapter, and we focus on computationally cheaper algorithms which are able to scale better with large hidden layers. Two strategies to this end are introduced next.

4.2.1 Consensus-based distributed training

The first strategy that we investigate for training an RVFL network in a fully decentralized way is simple, yet it results in a highly efficient training algorithm. It is composed of three steps:

1. **Initialization:** Parameters $\mathbf{w}_1, \dots, \mathbf{w}_B$ of the activation functions are agreed between nodes. For example, one node can draw these parameters from a uniform distribution and broadcast them to the rest of the network. This can be achieved in a decentralized way using a basic leader election strategy

³The DAC protocol is introduced in Appendix A.2.

[7]. Alternatively, they can be generated during the design of the distributed system (i.e. hardcoded in the network's design), so that they are already available when the system is actually started.

2. **Local training:** Each node solves its local training problem, considering only its own training dataset S_k . Solution is given by Eq. (4.3), obtaining a local set of output weights β_k^* , $k = 1 \dots L$.
3. **Averaging:** Local parameters vectors are averaged using a DAC strategy. After running DAC, the final weight vector at every node is given by:

$$\beta_{\text{CONS}}^* = \frac{1}{L} \sum_{k=1}^L \beta_k^*. \quad (4.7)$$

Despite its simplicity, consensus-based RVFL (denoted as CONS-RVFL) results in an interesting algorithm. It is easy to implement, even on low-cost hardware [119]; it requires low training times (i.e., local training and a short set of consensus iterations); moreover, our results show that it achieves a very low error, in many cases comparable to that of the centralized problem. From a theoretical standpoint, this algorithm can be seen as an ensemble of multiple linear predictors defined over the feature space induced by the mapping $\mathbf{h}(\cdot)$, i.e. it is similar to a bagged ensemble of linear predictors [21]. The overall algorithm is summarized in Algorithm 4.1.

Algorithm 4.1 CONS-RVFL: Consensus-based training for RVFL networks (k th node).

Inputs: Training set S_k , number of nodes L (global), regularization factor λ (global)

Output: Optimal weight vector β_k^*

- 1: Select parameters $\mathbf{w}_1, \dots, \mathbf{w}_B$, in agreement with the other $L - 1$ nodes.
 - 2: Compute \mathbf{H}_k and \mathbf{y}_k from S_k .
 - 3: Compute β_k^* via Eq. (4.3).
 - 4: $\beta^* \leftarrow \text{DAC}(\beta_1^*, \dots, \beta_L^*)$. ▷ Run in parallel, see Appendix A.
 - 5: **return** β^*
-

4.2.2 ADMM-based distributed training

Another strategy for training in a decentralized way a RVFL network is to optimize directly the global problem in Eq. (4.5) in a distributed fashion. Although potentially more demanding in computational time, this would ensure convergence to the global optimum. We can obtain a fully decentralized solution to problem in Eq.

(4.5) using the well-known ADMM. Most of the following derivation will follow [20, Section 8.2].

Derivation of the training algorithm

First, we reformulate the problem in the so-called ‘global consensus’ form, by introducing local variables β_k for every node, and forcing them to be equal at convergence. Hence, we rephrase the optimization problem as:

$$\begin{aligned} \beta_{\text{ADMM}}^* = \underset{\mathbf{z}, \beta_1, \dots, \beta_L \in \mathbb{R}^B}{\text{minimize}} \quad & \frac{1}{2} \left(\sum_{k=1}^L \|\mathbf{H}_k \beta_k - \mathbf{y}_k\|_2^2 \right) + \frac{\lambda}{2} \|\mathbf{z}\|_2^2 \\ \text{subject to} \quad & \beta_k = \mathbf{z}, \quad k = 1 \dots L. \end{aligned} \quad (4.8)$$

Then, we construct the augmented Lagrangian:

$$\begin{aligned} \mathcal{L} = \frac{1}{2} \left(\sum_{k=1}^L \|\mathbf{H}_k \beta_k - \mathbf{y}_k\|_2^2 \right) + \frac{\lambda}{2} \|\mathbf{z}\|_2^2 + \\ + \sum_{k=1}^L \mathbf{t}_k^T (\beta_k - \mathbf{z}) + \frac{\gamma}{2} \sum_{k=1}^L \|\beta_k - \mathbf{z}\|_2^2, \end{aligned} \quad (4.9)$$

where $\mathcal{L} = \mathcal{L}(\mathbf{z}, \beta_1, \dots, \beta_L, \mathbf{t}_1, \dots, \mathbf{t}_L)$, the vectors $\mathbf{t}_k, k = 1 \dots L$, are the Lagrange multipliers, $\gamma > 0$ is a penalty parameter, and the last term is introduced to ensure differentiability and convergence [20]. ADMM solves problems of this form using an iterative procedure, where at each step we optimize separately for β_k, \mathbf{z} , and we update the Lagrangian multipliers using a steepest-descent approach:

$$\beta_k[n+1] = \arg \min_{\beta_k \in \mathbb{R}^B} \mathcal{L}(\mathbf{z}[n], \beta_1, \dots, \beta_L, \mathbf{t}_1[n], \dots, \mathbf{t}_L[n]), \quad (4.10)$$

$$\mathbf{z}[n+1] = \arg \min_{\mathbf{z} \in \mathbb{R}^B} \mathcal{L}(\mathbf{z}, \beta_1[n+1], \dots, \beta_L[n+1], \mathbf{t}_1[n], \dots, \mathbf{t}_L[n]), \quad (4.11)$$

$$\mathbf{t}_k[n+1] = \mathbf{t}_k[n] + \gamma (\beta_k[n+1] - \mathbf{z}[n+1]). \quad (4.12)$$

In our case, the updates for $\beta_k[n+1]$ and $\mathbf{z}[n+1]$ can be computed in a closed form:

$$\beta_k[n+1] = (\mathbf{H}_k^T \mathbf{H}_k + \gamma \mathbf{I})^{-1} (\mathbf{H}_k^T \mathbf{y}_k - \mathbf{t}_k[n] + \gamma \mathbf{z}[n]), \quad (4.13)$$

$$\mathbf{z}[n+1] = \frac{\gamma \hat{\boldsymbol{\beta}} + \mathbf{t}}{\lambda/L + \gamma}, \quad (4.14)$$

where we introduced the averages $\hat{\boldsymbol{\beta}} = \frac{1}{L} \sum_{k=1}^L \beta_k[n+1]$ and $\mathbf{t} = \frac{1}{L} \sum_{k=1}^L \mathbf{t}_k[n]$. These averages can be computed in a decentralized fashion using a DAC step. We refer to [20] for a proof of the asymptotic convergence of ADMM.

Remark 1

In cases where, on a node, $N_k \ll B$, we can exploit the matrix inversion lemma to obtain a more convenient matrix inversion step [106]:

$$\left(\mathbf{H}_k^T \mathbf{H}_k + \gamma \mathbf{I}\right)^{-1} = \gamma^{-1} \left[\mathbf{I} - \mathbf{H}_k^T \left(\gamma \mathbf{I} + \mathbf{H}_k \mathbf{H}_k^T\right)^{-1} \mathbf{H}_k \right]. \quad (4.15)$$

Moreover, with respect to the training complexity, we note that the matrix inversion and the term $\mathbf{H}_k^T \mathbf{y}_k$ in Eq. (4.13) can be precomputed at the beginning and stored into memory. More advanced speedups can also be obtained with the use of Cholesky decompositions. Hence, time complexity is mostly related to the DAC step required in Eq. (4.14). Roughly speaking, if we allow ADMM to run for T iterations (see next subsection), the ADMM-based strategy is approximately T times slower than the consensus-based one.

Stopping criterion

Convergence of the algorithm at the k th node can be tracked by computing the ‘primal residual’ $\mathbf{r}_k[n]$ and ‘dual residual’ $\mathbf{s}[n]$, which are defined as:

$$\mathbf{r}_k[n] = \boldsymbol{\beta}_k[n] - \mathbf{z}[n], \quad (4.16)$$

$$\mathbf{s}[n] = -\gamma (\mathbf{z}[n] - \mathbf{z}[n-1]). \quad (4.17)$$

A possible stopping criterion is that both residuals should be less (in norm) than two thresholds:

$$\|\mathbf{r}_k[n]\|_2 < \epsilon_{\text{primal}}, \quad (4.18)$$

$$\|\mathbf{s}[n]\|_2 < \epsilon_{\text{dual}}. \quad (4.19)$$

A way of choosing the thresholds is given by [20]:

$$\epsilon_{\text{primal}} = \sqrt{L} \epsilon_{\text{abs}} + \epsilon_{\text{rel}} \max \left\{ \|\boldsymbol{\beta}_k[n]\|_2, \|\mathbf{z}[n]\|_2 \right\}, \quad (4.20)$$

$$\epsilon_{\text{dual}} = \sqrt{L} \epsilon_{\text{abs}} + \epsilon_{\text{rel}} \|\mathbf{t}_k[n]\|_2, \quad (4.21)$$

where ϵ_{abs} and ϵ_{rel} are user-specified absolute and relative tolerances, respectively. Alternatively, as in the previous case, the algorithm can be stopped after a maximum number of iterations is reached. The pseudocode for the overall algorithm, denoted as ADMM-RVFL, at a single node is given in Algorithm 4.2.

Algorithm 4.2 ADMM-RVFL: ADMM-based training for RVFL networks (k th node).

Inputs: Training set S_k , number of nodes L (global), regularization factors λ, γ (global), maximum number of iterations T (global)

Output: Optimal vector β_k^*

- 1: Select parameters $\mathbf{w}_1, \dots, \mathbf{w}_B$, in agreement with the other $L - 1$ nodes.
 - 2: Compute \mathbf{H}_k and \mathbf{y}_k from S_k .
 - 3: Initialize $\mathbf{t}_k[0] = \mathbf{0}, \mathbf{z}[0] = \mathbf{0}$.
 - 4: **for** n from 0 to T **do**
 - 5: Compute $\beta_k[n + 1]$ according to Eq. (4.13).
 - 6: Compute averages $\hat{\beta}$ and \mathbf{t} using DAC.
 - 7: Compute $\mathbf{z}[n + 1]$ according to Eq. (4.14).
 - 8: Update $\mathbf{t}_k[n]$ according to Eq. (4.12).
 - 9: Check termination with residuals.
 - 10: **end for**
 - 11: **return** $\mathbf{z}[n]$
-

4.3 Experimental Setup

4.3.1 Description of the Datasets

We tested our algorithms on four publicly available datasets, whose characteristics are summarized in Table 4.1.

Table 4.1: General description of the datasets

Dataset name	Features	Instances	Desired output	Task type
G50C	50	550	Gaussian of origin	Classification (2 classes)
Garageband	44	1856	Genre recognition	Classification (9 classes)
Skills	18	3338	User's level	Regression
Sylva	216	14394	Forest Type	Classification (2 classes)

We have chosen them to represent different applicative domains of our algorithms, and to provide enough diversity in terms of size, number of features, and imbalance of the classes:

- *Garageband* is a music classification problem [111], where the task is to discern among 9 different genres. As we stated in the previous chapter, in the distributed case we can assume that the songs are present over different computers, and we can use our strategies as a way of leveraging over the entire dataset without a centralized controller.

- *Skills* is a regression dataset taken from the UCI repository [177]. The task is to assess the skill level of a video game user, based on a set of recordings of its actions in the video game itself. This is useful for letting the game adapt to the user's characteristics. In this case, data is distributed by definition throughout the different players in the network. By employing our strategy several computers, each playing their own version of the game, can learn to adapt better by exploiting collective data.
- *Sylva* is a binary classification task for distinguishing classes of trees (Ponderosa pine vs. everything else).⁴ It is an interesting dataset since it has a large imbalance between the positive and negative examples (approximately 15 : 1), and a large subset of the features are not informative from a classification point of view. In the distributed case, we can imagine that data is collected by different sensors.
- *G50C*, differently from the others, is an artificial dataset [110], whose main interest is given by the fact that the optimal (Bayes) error rate is designed to be equal exactly to 5%.

In all cases, input variables are normalized between 0 and 1, and missing values are replaced with the average computed over the rest of the dataset. Multi-class classification is handled with the standard M bit encoding for the output, associating to an input \mathbf{x}_i a single output vector \mathbf{y}_i of M bits, where if its elements are $y_{ij} = 1$ and $y_{ik} = 0, k \neq j$, then the corresponding pattern is of class j . We can retrieve the actual class from the M -dimensional RVFL output as:

$$\text{Class of } \mathbf{x} = \arg \max_{j=1 \dots M} f_j(\mathbf{x}), \quad (4.22)$$

where $f_j(\mathbf{x})$ is the j th element of the M -dimensional output $f(\mathbf{x})$. For all the models, testing accuracy and training times is computed by executing a 5-fold cross-validation over the available data. This 5-fold procedure is then repeated 15 times by varying the topology of the agents and the initial weights of the RVFL net. Final misclassification error and training time is then collected for all the $15 \times 5 = 75$ repetitions, and the average values and standard deviations are computed.

4.3.2 Algorithms and Software Implementation

We compare the following algorithms:

- **Centralized RVFL (C-RVFL):** this is a RVFL trained with all the available training data. It is equivalent to a centralized node collecting all the data, and

⁴http://www.causality.inf.ethz.ch/al_data/SYLVA.html

it can be used as a baseline for the other approaches.

- **Local RVFL (L-RVFL)**: in this case, training data is distributed evenly across the nodes. Every node trains a standard RVFL with its own local dataset, but no communication is performed. Testing error is averaged throughout the nodes.
- **Consensus-based RVFL (CONS-RVFL)**: as before, data is evenly distributed in the network, and the consensus strategy explained in Section 4.2.1 is executed. We set a maximum of 300 iterations and $\delta = 10^{-3}$.
- **ADMM-based RVFL (ADMM-RVFL)**: similar to before, but we employ the ADMM-based strategy described in Section 4.2.2. In this case, we set a maximum of 300 iterations, $\epsilon_{\text{rel}} = \epsilon_{\text{abs}} = 10^{-3}$ and $\gamma = 1$.

In all cases, we use sigmoid hidden functions given by Eq. (4.4), where parameters \mathbf{w} and b in (4.4) are extracted randomly from a uniform distribution over the interval $[-1, +1]$. To compute the optimal number of hidden nodes and the regularization parameter λ , we execute an inner 3-fold cross-validation on the training data only for C-RVFL. In particular, we search the uniform interval $\{50, 100, 150, \dots, 1000\}$ for the number of hidden nodes, and the exponential interval $2^j, j \in \{-10, -9, \dots, 9, 10\}$ for λ . The step size of 50 in the hidden nodes interval was found to provide a good compromise between final accuracy and the computational cost of the grid-search procedure. These parameters are then shared with the three remaining models. We experimented with a separate fine-tuning for each model, but no improvement in performance was found. Optimal parameters averaged over the runs are shown in Table 4.2.

Table 4.2: Optimal parameters found by the grid-search procedure

Dataset	Hidden nodes	λ
G50C	500	2^3
Garageband	200	2^{-3}
Skills	400	2^{-2}
Sylva	450	2^{-5}

We have implemented CONS-RVFL and ADMM-RVFL in the open-source Lynx MATLAB toolbox (see Appendix B). Throughout this thesis, we are not concerned with the analysis of communication overhead over a realistic channel; hence, we employ a serial version of the code where the network is simulated artificially. However, in the aforementioned toolbox we also provide a fully parallel version, able to work on a cluster architecture, in order to test the accuracy of the system in a more realistic setting.

4.4 Results and Discussion

4.4.1 Accuracy and Training Times

The first set of experiments is to show that both algorithms that we propose are able to approximate very closely the centralized solution, irrespective of the number of nodes in the network. The topology of the network in these experiments is constructed according to the so-called ‘Erdős–Rényi model’ [117], i.e., once we have selected a number L of nodes, we randomly construct an adjacency matrix such that every edge has a probability p of appearing, with p specified a-priori. For the moment, we set $p = 0.2$; an example of such a network for $L = 8$ is shown in Fig. 4.1.

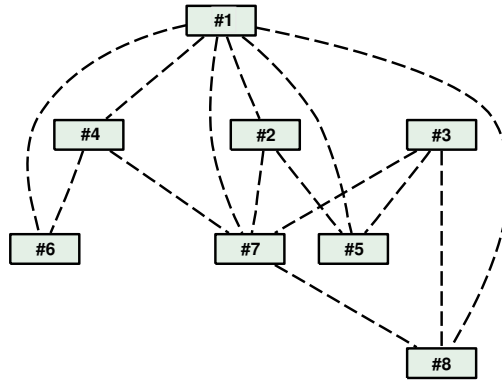


Figure 4.1: Example of network used in the experiments with 8 nodes. Connectivity is generated at random, with a 20% probability for each link of being present.

To test the accuracy of the algorithms, we vary L from 5 to 50 by steps of 5. Results are presented in Fig. 4.2 (a)-(d). For the three classification datasets, we show the averaged misclassification error. For the Skills dataset, instead, we show the Normalized Root Mean-Squared Error (NRMSE), defined for a test set \mathcal{T} as:

$$\text{NRMSE}(\mathcal{T}) = \sqrt{\frac{\sum_{(x_i, y_i) \in \mathcal{T}} [f(x_i) - y_i]^2}{|\mathcal{T}| \hat{\sigma}_y}}, \quad (4.23)$$

where $|\mathcal{T}|$ denotes the cardinality of the set \mathcal{T} and $\hat{\sigma}_y$ is an empirical estimate of the variance of the output samples $y_i, i = 1, \dots, |\mathcal{T}|$. For every fold, the misclassification error of L-RVFL is obtained by averaging the error over the L different nodes. While this is a common practice, it can introduce a small bias with respect to the other curves. Nonetheless, we stress that it does not influence the following discussion, which mostly focuses on the comparison of the other three algorithms.

The first thing to observe is that L-RVFL has a steady decrease in performance in all situations, ranging from a small decrease in the Skills dataset with 5 nodes, to more than 20% of classification accuracy when considering networks of 50 agents in the G50C and Garageband datasets. Despite being obvious, because of the decrease

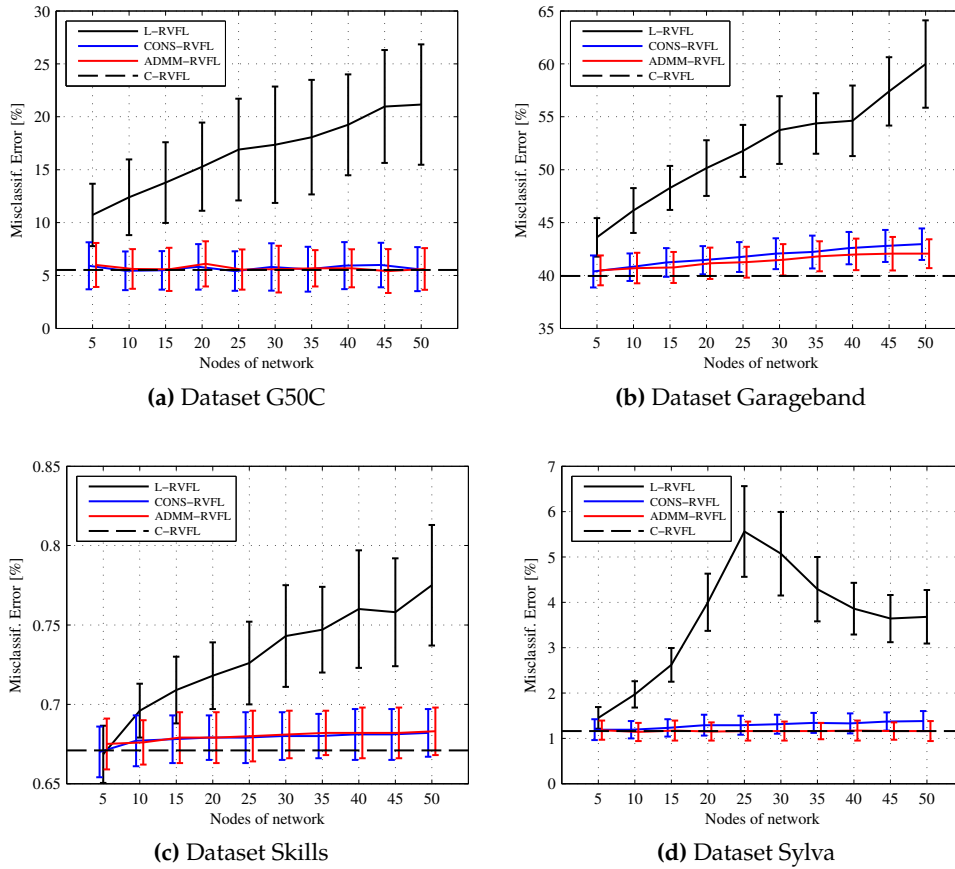


Figure 4.2: Average error and standard deviation of the models on the four datasets, when varying the number of nodes in the network from 5 to 50. For G50C, Garageband, and Sylva we show the misclassification error, while for Skills we show the NRMSE. Lines for CONS-RVFL and ADMM-RVFL are slightly separated for better readability. Vertical bars represent the standard deviation from the average result.

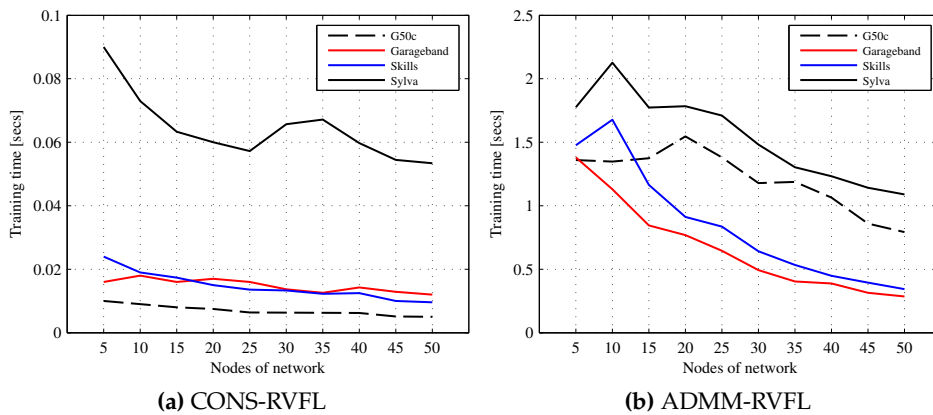


Figure 4.3: Average training time for (a) CONS-RVFL and (b) ADMM-RVFL on a single node.

of available data at each node, it is an experimental confirmation of the importance of leveraging over *all* possible data in terms of accuracy. It is also interesting to note that the gap between L-RVFL and C-RVFL does not always increase monotonically with respect to the size of the network, as shown by Fig. 4.2-(d). A possible explanation of this fact is that, by keeping fixed the λ parameter, the effect of the regularization factor in Eq. (4.5) is proportionally higher when decreasing the amount of training data.

The second important aspect is that CONS-RVFL and ADMM-RVFL are *both* able to match very closely the performance of C-RVFL, irrespective of the network's size. In particular, they have the same performance on the G50C and Skills datasets, whilst a small gap is present in the Garageband and Sylva cases, although it is not significant.

Next, let us analyze the training times of the distributed algorithm, shown in Fig. 4.3a for CONS-RVFL and Fig. 4.3b for ADMM-RVFL, respectively. In particular, we show the average training time spent at a single node. Generally speaking, CONS-RVFL is approximately one order of magnitude faster than ADMM-RVFL, which requires multiple iterations of consensus. In both cases, the average training time spent at a single node is monotonically decreasing with respect to the overall number of nodes. Hence, the computational time of the matrix inversion is predominant compared to the overhead introduced by the DAC and ADMM procedures.

4.4.2 Effect of Network Topology

Now that we have ascertained the convergence properties of both algorithms, we analyze an interesting aspect: how does the topology of the network influences the convergence time? Clearly, as long as the network stays connected, the accuracy is not influenced. However, the time required for the consensus to achieve convergence is dependent on how the nodes are interconnected. At the extreme, in a fully connected network, two iterations are always sufficient to achieve convergence at any desired level of accuracy. More in general, the time will be roughly proportional to the average distance between any two nodes. To test this, we compute the iterations needed to reach consensus for several topologies of networks composed of 50 nodes:

- **Random network:** this is the network constructed according to the Erdős–Rényi model described in the previous subsection. We experiment with $p = 0.2$ and $p = 0.5$, and denote the corresponding graphs as $R(0.2)$ and $R(0.5)$ respectively.
- **Linear network:** in this network the nodes are ordered, and each node in the sequence is connected to its most immediate K successors, with K specified a-priori, except the last $K - 1$ nodes, which are connected only to the remaining

ones. We experiment with $K = 1$ and $K = 4$, and denote the networks as $K(1)$ and $K(4)$ respectively.

- **Small world:** this is a network constructed according to the well-known ‘Watts-Strogatz’ mechanism [190]. First, a cyclic topology is constructed, i.e., nodes are ordered in a circular sequence, and every node is connected to K nodes to its left and K to its right. Then, every link is ‘rewired’ with probability set by a parameter $\alpha \in [0, 1]$. In our case, we have $K = 6$ and $\alpha = 0.15$, and denote the resulting topology as SW .
- **Scale-free:** this is another topology that tries to reflect realistic networks, in this case exhibiting a power law with respect to the degree distribution. We construct it according to the ‘Barabási-Albert’ model of preferential attachment [2], and denote the resulting topology as SF .

Results are presented in Fig. 4.4 (a)-(d).

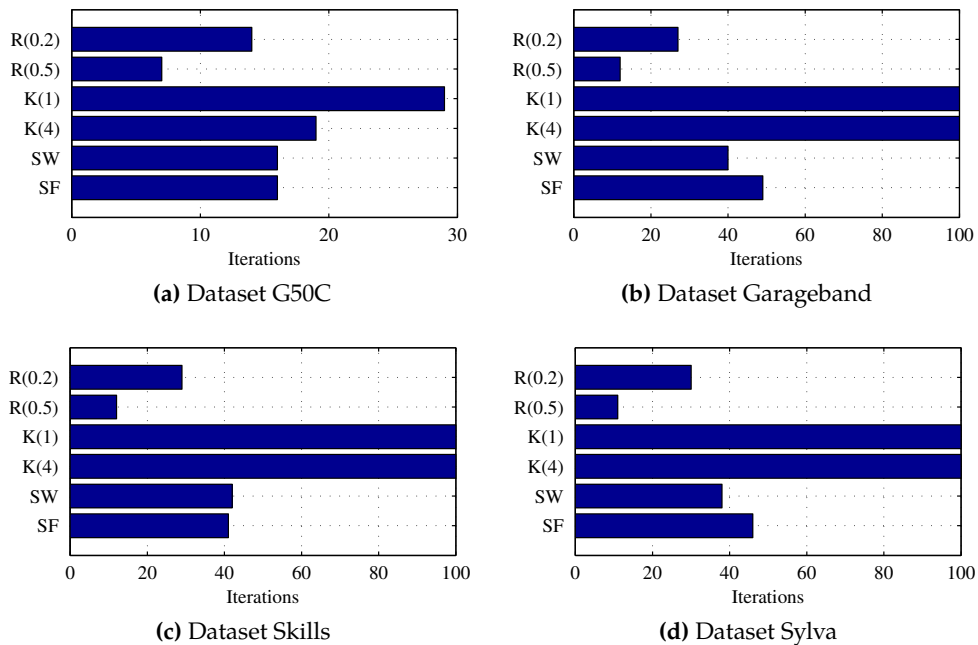


Figure 4.4: Consensus iterations needed to reach convergence when varying the network topology. For the explanation of the topologies see Sect. 4.4.2. The number of iterations is truncated at 100 for better readability.

We see that the algorithm has very similar results on all four datasets. In particular, as we expected, consensus is extremely slow in reaching agreement when considering linear topologies, where information takes several iterations to reach one end of the graph from the other. At the other extreme, it takes a very limited number of iterations in the case of a highly connected graph, as in the case of $R(0.5)$. In between, we can see that consensus is extremely robust to a change

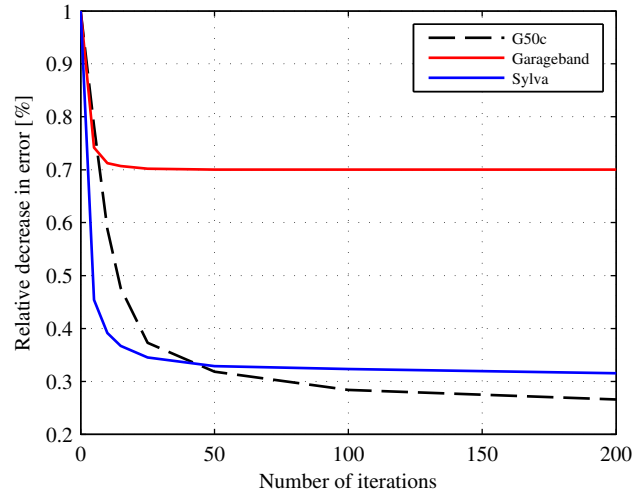


Figure 4.5: Relative decrease in error of ADMM-RVFL with respect to L-RVFL, when using an early stopping procedure at different iterations.

in topology and its performance is not affected when considering small-world or scale-free graphs.

4.4.3 Early Stopping for ADMM

Next, we explore a peculiar difference between CONS-RVFL and ADMM-RVFL. In the case of CONS-RVFL, no agreement is reached between the different nodes until the consensus procedure is completed. Differently from it, an intermediate solution is available at every iteration in ADMM-RVFL, given by the vector $\mathbf{z}[n]$. This allows for the use of an early stopping procedure, i.e., the possibility of stopping the optimization process before actual convergence, by fixing in advance a predefined (small) number of iterations. In fact, several experimental findings support the idea that ADMM can achieve a reasonable degree of accuracy in the initial stages of optimization [20]. To test this, we experiment early stopping for ADMM-RVFL at $\{5, 10, 15, 25, 50, 100, 200\}$ iterations for the three classification datasets. In Fig. 4.5 we plot the relative decrease in performance with respect to L-RVFL. We can see that 10-15 iterations are generally enough to reach a good performance, while the remaining iterations are proportionally less useful. As a concrete example, misclassification error of ADMM-RVFL for G50C is 16.55% after only 5 iterations, 12.44% after 10, 7.89% after 25, while the remaining 175 iterations are used to decrease the error only by an additional 2 percentage points.

4.4.4 Experiment on Large-Scale Data

As a final experimental validation, we analyze the behavior of CONS-RVFL and ADMM-RVFL on a realistic large-scale dataset, the well-known CIFAR-10 image

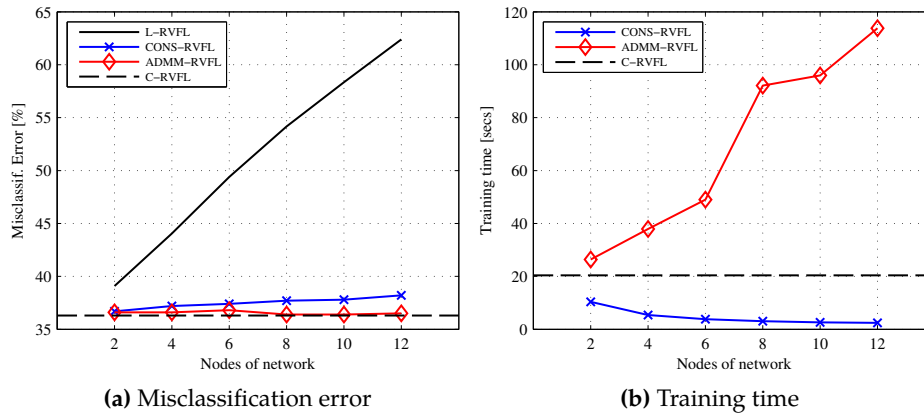


Figure 4.6: Average misclassification error and training time of CONS-RVFL and ADMM-RVFL on the CIFAR-10 dataset, when varying the nodes of the network from 2 to 12.

classification database [86]. It is composed of 50000 images labeled in 10 different classes, along with a standard testing set of additional 10000 images. Each image is composed of exactly 32×32 pixels, and each pixel is further represented by 3 integer values in the interval $[1, 255]$, one for each color channel in the RGB color space. Classes are equally distributed between the training patterns, i.e., every class is represented by exactly 5000 images. Since we are mostly interested into the relative difference in performance between the algorithms, and not in achieving the lowest possible classification error, we preprocess the images using the relatively simple procedure detailed in [38]. In particular, we extract 1600 significant patches from the original images, and represent each image using their similarity with respect to each of the patches. We refer to [38] for more details on the overall workflow. In this experiment, we use $B = 3000$ and $\lambda = 10$, a setting which was found to work consistently on all situations. Moreover, we use the $R(0.2)$ graph explained before, but we experiment with lower number of nodes in the network, which we vary from 2 to 12 by steps of 2. All the other parameters are set as in the previous experiments. Although the test set is fixed in this case, we repeat each experiment 15 times to average out the effect of randomness in the RVFL and connectivity initializations.

The average misclassification error of the four models is shown in Fig. 4.6a. In this case, the effect of splitting data is extremely pronounced, and the average misclassification error of L-RVFL goes from 39% with 2 nodes, up to 62.4% with 12 nodes. Both CONS-RVFL and ADMM-RVFL are able to track very efficiently the centralized solution, although there is a small gap in performance between the two (of approximately 1%), when distributing over more than 8 nodes. This is more than counter-balanced, however, by considering the advantage of CONS-RVFL with respect to the required training time. To show this, we present the average training time (averaged over the nodes) in Fig. 4.6b. In this case, due to the large

expansion block, time required to perform the multiple consensus iterations in ADMM-RVFL prevails over the rest, and the average training time tends to increase when increasing the size of the network. This is not true of CONS-RVFL, however, which obtains an extremely low training time with respect to C-RVFL, up to an order of magnitude for sufficiently large networks. Hence, we can say that CONS-RVFL can also be an efficient way of computing an approximate solution to a standard RVFL, with good accuracy, by distributing the computation over multiple machines.

Extending Distributed RVFL Networks to a Sequential Scenario

Contents

5.1	Derivation of the algorithm	49
5.2	Experiments on Distributed Music Classification	51
5.2.1	The Distributed Music Classification Problem	51
5.2.2	Experiment Setup	52
5.2.3	Results and Discussion	54
5.3	Comparison of DAC strategies	56
5.3.1	Description of the strategies	57
5.3.2	Experimental Results	59

THE algorithms presented in the previous chapter have been designed for working in a batch setting. In this chapter, we extend CONS-RVFL to the distributed training of RVFL networks in the case where data is arriving sequentially at every node. Particularly, we combine the DAC-based strategies with local updates based on the blockwise RLS (BRLS) training algorithm. Next, we present a case study for distributed music classification in Section 5.2. Finally, we compare the impact of using advanced choices for the connectivity matrix \mathbf{C} of the DAC protocol in Section 5.3.

5.1 Derivation of the algorithm

Remember from Section 3.2 that in a sequential setting, the local dataset S_k is not processed as a whole, but it is presented in a series of batches (or chunks)

The content of this chapter is adapted from the material published in [148] and [53].

$S_{k,1}, \dots, S_{k,T}$ such that:

$$\bigcup_{i=1}^T S_{k,T} = S_k \quad k = 1, \dots, L. \quad (5.1)$$

This encompasses situations where training data arrives in a streaming fashion, or the case where the dataset S_k is too large for the matrix inversion in Eq. (4.3) to be practical. In this section, we assume that new batches arrive synchronously at every node. In the single-agent case, an RVFL network can be trained efficiently in the sequential setting by the use of the BRLS algorithm [184]. Denote by $\beta[n]$ the estimate of its optimal weight vector after having observed the first n chunks, and by \mathbf{H}_{n+1} and \mathbf{y}_{n+1} the matrices collecting the hidden nodes values and outputs of the $(n+1)$ th chunk S_{n+1} . BRLS recursively computes Eq. (4.3) by the following two-step update:

$$\mathbf{P}[n+1] = \mathbf{P}[n] - \mathbf{P}[n]\mathbf{H}_{n+1}^T\mathbf{M}_{n+1}^{-1}\mathbf{H}_{n+1}\mathbf{P}[n], \quad (5.2)$$

$$\beta[n+1] = \beta[n] + \mathbf{P}[n+1]\mathbf{H}_{n+1}^T(\mathbf{y}_{n+1} - \mathbf{H}_{n+1}\beta[n]), \quad (5.3)$$

where we have defined:

$$\mathbf{M}_{n+1} = \mathbf{I} + \mathbf{H}_{n+1}\mathbf{P}[n]\mathbf{H}_{n+1}^T. \quad (5.4)$$

The matrix \mathbf{P} in Eq. (5.3) and Eq. (5.4) can be initialized as $\mathbf{P}[0] = \lambda^{-1}\mathbf{I}$, while the weights $\beta[0]$ as the zero vector. For a derivation of the algorithm, based on the Sherman-Morrison formula, and an analysis of its convergence properties we refer the interested reader to [184]. The BRLS gives rise to a straightforward extension of the DAC-based training algorithm presented in Section 4.2.1 for the DL setting, consisting in interleaving local update steps with global averaging over the output weight vector. Practically, we consider the following algorithm:

1. **Initialization:** the nodes agree on parameters $\mathbf{w}_1, \dots, \mathbf{w}_B$ in Eq. (4.1). The same considerations made in Section 4.2.1 apply here. Moreover, all the nodes initialize their own local estimate of the \mathbf{P} matrix in Eq. (5.2) and Eq. (5.4) as $\mathbf{P}_k[0] = \lambda^{-1}\mathbf{I}$, and their estimate of the output weight vector as $\beta_k[0] = \mathbf{0}$.
2. At every iteration $n+1$, each node k receives a new batch $S_{k,n+1}$. The following steps are performed:
 - (a) **Local update:** every node computes (locally) its estimate $\beta_k[n+1]$ using Eqs. (5.2)-(5.3) and local data $S_{k,n+1}$.
 - (b) **Global average:** the nodes agree on a single parameter vector by averaging their local estimates with a DAC protocol. The final weight vector at

iteration $n + 1$ is then given by:

$$\boldsymbol{\beta}[n + 1] = \frac{1}{L} \sum_{k=1}^L \boldsymbol{\beta}_k[n + 1]. \quad (5.5)$$

The overall algorithm, denoted as S-CONS-RVFL, is summarized in Algorithm 5.1.

Algorithm 5.1 S-CONS-RVFL: Extension of CONS-RVFL to the sequential setting (k th node).

Inputs: Number of nodes L (global), regularization factor λ (global)

Output: Optimal weight vector $\boldsymbol{\beta}_k^*$

- 1: Select parameters $\mathbf{w}_1, \dots, \mathbf{w}_B$, in agreement with the other $L - 1$ nodes.
 - 2: $\mathbf{P}_k[0] = \lambda^{-1} \mathbf{I}$.
 - 3: $\boldsymbol{\beta}_k[0] = \mathbf{0}$.
 - 4: **for** $n = 1, \dots, T$ **do**
 - 5: Receive batch $S_{k,n}$.
 - 6: Update $\boldsymbol{\beta}_k[n + 1]$ using Eqs. (5.2)-(5.3).
 - 7: $\boldsymbol{\beta}_k[n + 1] \leftarrow \text{DAC}(\boldsymbol{\beta}_k[n + 1], \dots, \boldsymbol{\beta}_L[n + 1])$. ▶ Run in parallel, see Appendix A.
 - 8: **end for**
 - 9: **return** $\boldsymbol{\beta}_k[T]$
-

5.2 Experiments on Distributed Music Classification

5.2.1 The Distributed Music Classification Problem

As an experimental setting, we consider the problem of distributed automatic music classification (AMC). AMC is the task of automatically assigning a song to one (or more) classes, depending on its audio content [147]. It is a fundamental task in many music information retrieval (MIR) systems, whose broader scope is to efficiently retrieve songs from a vast database depending on the user's requirements [59]. Examples of labels that can be assigned to a song include its musical genre, artist [50], induced mood [59] and leading instrument. Classically, the interest in music classification is two-fold. First, being able to correctly assess the aforementioned characteristics can increase the efficiency of a generic MIR system (see survey [59] and references therein). Secondly, due to its properties, music classification can be considered as a fundamental benchmark for supervised learning algorithms [147]: apart from the intrinsic partial subjectivity of assigning labels, datasets tend to be relatively large, and a wide variety of features can be used to describe each song. These features can also be supplemented by meta-informations and social tags.

More formally, we suppose that the input $\mathbf{x} \in \mathbb{R}^d$ to the model is given by a suitable d -dimensional representation of a song. Examples of features that can be used in this sense include temporal features such as the zero-crossing count, compact statistics in the frequency and cepstral domain [59], higher-order descriptors (e.g. timbre [50]), meta-information on the track (e.g., author), and social tags extracted from the web. The output is instead given by one of M predefined classes, where each class represents a particular categorization of the song, such as its musical genre. In the distributed AMC setting, these songs are distributed over a network, as is common in distributed AMC on peer-to-peer (P2P) systems, and over wireless sensor networks [138].

5.2.2 Experiment Setup

We use four freely available AMC benchmarks. A schematic description of their characteristics is given in Table 5.1.

Table 5.1: General description of the datasets for testing the sequential S-CONSRVFL algorithm.

Dataset name	Features	Instances	Task	Classes	Reference
Garageband	49	1856	Genre recognition	9	[111]
Latin Music Database (LMD)	30	3160	Genre recognition	10	[168]
Artist20	30	1413	Artist recognition	20	[50]
YearPredictionMSD	90	200000	Decade identification	2	[12]

Below we provide more information on each of them.

- *Garageband* [111] is a genre classification dataset, considering 1856 songs and 9 different genres (alternative, blues, electronic, folkcountry, funksoulrnb, jazz, pop, raphiphop and rock). The input is given by 49 features extracted according to the procedure detailed in [111]. It is the same as the one used in the previous chapter.
- *LMD* is another genre classification task, of higher difficulty [168]. In this case, we have 3160 different songs categorized in 10 Latin American genres (tango, bolero, batchata, salsa, merengue, ax, forr, sertaneja, gacha and pagode). The input is a 30-dimensional feature vector, extracted from the middle 30 seconds of every song using the Marsyas software.¹ Features are computed both in the frequency domain (e.g. the spectral centroid) and in the cepstral domain, i.e. Mel Frequency Cepstral Coefficients (MFCC).

¹<http://marsyas.info/>

- *Artist20* is an artist recognition task comprising 1413 songs distributed between 20 different artists [50]. The 30-dimensional input vector comprises both MFCC and chroma features (see [50] for additional details).
- *YearPredictionMSD* is a year recognition task derived from the subset of the million song dataset [12] available on the UCI machine learning repository.² It is a dataset of 500000 songs categorized by year. In our experiment, we consider a simplified version comprising only the initial 200000 songs, and the following binary classification output: a song is of class (a) if it was written previously than 2000, and of class (b) otherwise. This is a meaningful task due to the unbalance of the original dataset with respect to the decade 2001 – 2010.

In all cases, input features were normalized between -1 and $+1$ before the experiments. Testing accuracy is computed over a 10-fold cross-validation of the data, and every experiment is repeated 50 times to average out randomness effects due to the initialization of the parameters. Additionally, to increase the dataset size, we artificially replicate twice the training data for all datasets, excluding YearDatasetMSD.

We consider networks of 8 nodes, whose topology is constructed according to the ‘Erdős–Rényi model’ (see Section 4.3.2). In particular, every pair of nodes in the network has a 20% probability of being connected, with the only constraint that the overall network is connected. Training data is distributed evenly across the nodes, and chunks are constructed such that every batch is composed of approximately 20 examples (100 for the YearPredictionMSD dataset). We compare the following algorithms:

- **Sequential CONS-RVFL (S-CONS-RVFL)**: this is trained according to the consensus-based sequential algorithm. For the DAC procedure, we set the maximum number of iterations to 300, and $\delta = 10^{-4}$.
- **Centralized RVFL (C-RVFL)**: this is a RVFL trained by first collecting all the local chunks and aggregating them in a single batch. It can be considered as an upper bound on the performance of S-CONS-RVFL.
- **Local RVFL (L-RVFL)**: in this case, nodes update their estimate using their local batch, but no communication is performed. Final misclassification error is averaged across the nodes. This can be considered as a worst-case baseline for the performance of any distributed algorithm for RVFL networks.

In all cases, we use sigmoid hidden functions given by Eq. (4.4). Optimal parameters for C-RVFL are found by executing an inner 3-fold cross-validation on the training

²<https://archive.ics.uci.edu/ml/>

data. In particular, we search the uniform interval $\{50, 100, 150, \dots, 1000\}$ for the number of hidden nodes, and the exponential interval 2^j , $j \in \{-10, -9, \dots, 9, 10\}$ for λ . These parameters are then shared with L-RVFL and S-CONS-RVFL. Resulting parameters from the grid search procedure are listed in Table 5.2.

Table 5.2: Optimal parameters found by the grid-search procedure.

Dataset	Hidden nodes	λ
Garageband	300	2^{-3}
LMD	400	2^{-2}
Artist20	200	2^{-4}
YearPredictionMSD	300	1

We note that C-RVFL can be considered as a benchmark for audio classification using shallow neural networks. In fact, in [147] it is shown that it outperforms a standard MLP trained using SGD.

5.2.3 Results and Discussion

We start our discussion of the results by analyzing the final misclassification error and training time for the three models, reported in Table 5.3. Results of the proposed algorithm, S-CONS-RVFL, are highlighted in bold.

Whenever we consider medium-sized datasets, the performance of L-RVFL is strictly worse than the performance of C-RVFL (similarly to the previous chapter), ranging from an additional 5% misclassification error for Garageband and LMD, up to an additional 10% for Artist20. The most important fact highlighted in Table 5.3, however, is that S-CONS-RVFL is able to efficiently match the performance of C-RVFL in all situations, except for a small decrease in the LMD dataset. From a computational perspective, this performance is achieved with a very small overhead in terms of training time with respect to L-RVFL in all cases (as evidenced by the fourth column in Table 5.3).

In a sequential setting, the evolution of the testing error after every batch is equally as important as the final accuracy obtained. We report it in Fig. 5.1(a)-(d) for the four datasets.

Performance of C-RVFL, L-RVFL and S-CONS-RVFL are shown with dashed black, solid red and solid blue lines respectively. Moreover, performance of L-RVFL is averaged across the nodes. Once again, we see that S-CONS-RVFL is able to track very efficiently the accuracy obtained by C-RVFL. The performance is practically equivalent in the Garageband and YearPredictionMSD datasets (Fig. 5.1(a) and Fig.

Table 5.3: Final misclassification error and training time for the three models, together with standard deviation. The proposed algorithm is highlighted in bold. Training time for S-CONS-RVFL and L-RVFL is averaged over the nodes.

Dataset	Algorithm	Error	Time [secs]
Garageband	C-RVFL	0.40 ± 0.02	0.24 ± 0.09
	L-RVFL	0.45 ± 0.03	0.13 ± 0.03
	S-CONS-RVFL	0.40 ± 0.02	0.15 ± 0.04
LMD	C-RVFL	0.25 ± 0.02	0.70 ± 0.17
	L-RVFL	0.31 ± 0.03	0.46 ± 0.08
	S-CONS-RVFL	0.26 ± 0.02	0.49 ± 0.10
Artist20	C-RVFL	0.37 ± 0.04	0.13 ± 0.07
	L-RVFL	0.47 ± 0.04	0.06 ± 0.01
	S-CONS-RVFL	0.37 ± 0.04	0.09 ± 0.02
YearPredictionMSD	C-RVFL	0.27 ± 0.01	8.66 ± 0.93
	L-RVFL	0.27 ± 0.01	2.35 ± 0.48
	S-CONS-RVFL	0.27 ± 0.01	2.46 ± 0.62

5.1(d)), while convergence speed is slightly slower in the LMD and Artist20 case (Fig. 5.1(b) and Fig. 5.1(c)), although by a small amount. This gap depends on the fact that S-CONS-RVFL remains an approximation of C-RVFL. In particular, in the current version of S-CONS-RVFL no information is exchanged with respect to the state matrices $\mathbf{P}_k[n]$, which would be infeasible for large B (see also the similar observation for the batch case in Section 4.2).

Next, we investigate the behavior of S-CONS-RVFL when varying the size of the network. In fact, due to its parallel nature, we expect that, the higher the number of nodes, the lower the training time (apart from communication bottlenecks, depending on the real channel of the network). The following experiments show that the increase in time required by the DAC procedure for bigger networks is more than compensated by the gain in time obtained by processing a lower number of samples per node. To this end, we consider the training time required by CONS-RVFL when varying the number of nodes of the network from 2 to 14 by steps of 2, keeping the same topology model as before. Results of this experiment are presented in Fig. 5.2(a) for datasets Garageband and Artist20, and in Fig. 5.2(b) for datasets LMD and YearPredictionMSD.

The decrease in training time is extremely pronounced for Garageband, with a five-fold decrease going from 2 to 14 nodes, and for YearPredictionMSD, with a seven-fold decrease. This result is especially important, showing that S-CONS-RVFL can be efficiently used in large-scale situations. It is also consistent with the analysis of the batch CONS-RVFL in the previous chapter. Similarly, the number of

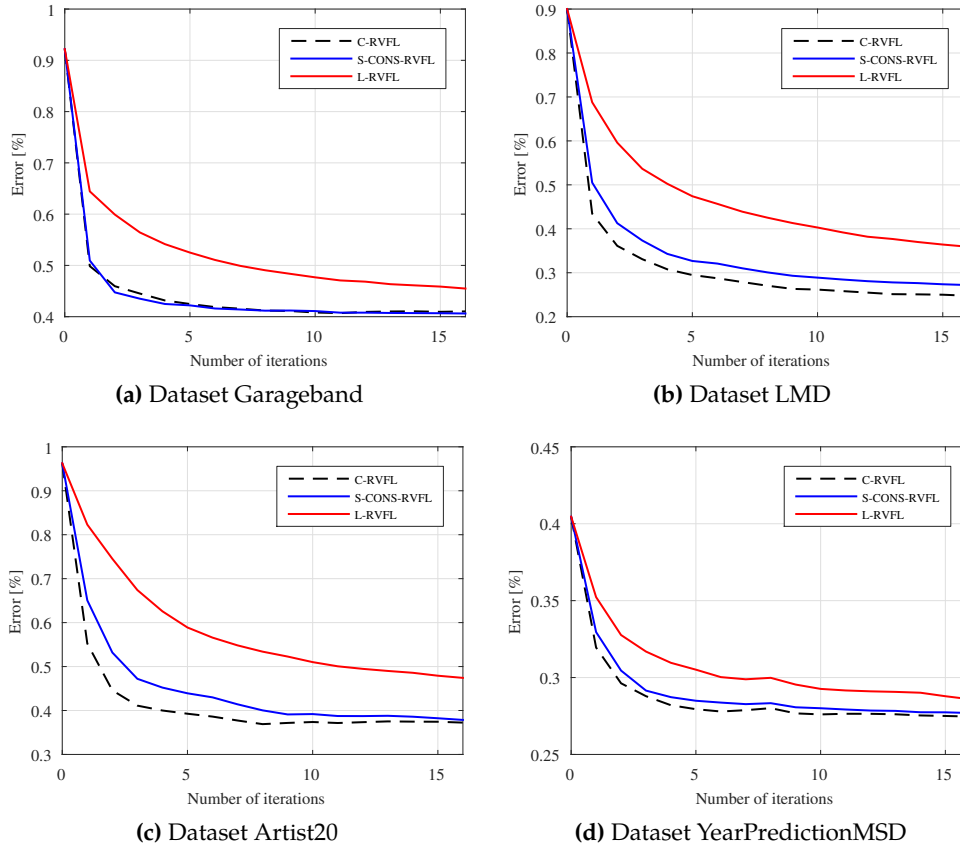


Figure 5.1: Evolution of the testing error after every iteration. Performance of L-RVFL is averaged across the nodes.

consensus iterations needed to reach the desired accuracy is shown in Fig. 5.3.

Although the required number of iterations grows approximately linearly with respect to the size of the network, a low number of iterations is generally enough to reach convergence to a very good accuracy. In fact, no experiment in this section required more than 35 iterations in total. Additionally, the consensus procedure is extremely robust to a change in the network topology, as shown in the previous chapter. The same considerations apply here.

5.3 Comparison of DAC strategies

Up to this point, we have considered only the ‘max-degree’ strategy for the DAC protocol, detailed in Appendix A. However, it is known that the performance of the DAC protocol, and by consequence the performance of any distributed training algorithm based on its application, can improve significantly with proper choices of the mixing parameters [198]. At the same time, a thorough investigation of multiple strategies for choosing the weights is missing in the literature. In this section, we compare four of them in the context of S-CONS-RVFL, ranging from

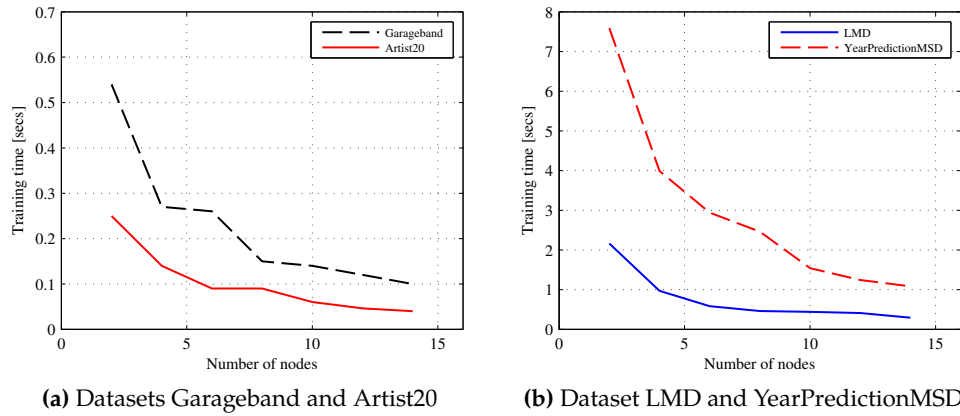


Figure 5.2: Training time required by the sequential S-CONS-RVFL, for varying sizes of the network, from 2 to 14 by steps of 2.

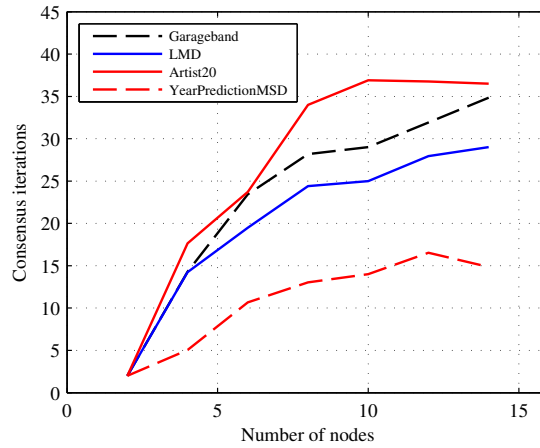


Figure 5.3: Number of consensus iterations required to reach convergence in the S-CONS-RVFL, when varying the number of nodes in the network from 2 to 14.

choosing a fixed value for every coefficient, to more complex choices satisfying strong optimality conditions. Our experimental results show that the performance of the DAC protocol, and by consequence the performance of any distributed training algorithm based on its application, can improve significantly with proper choices of the mixing parameters.

5.3.1 Description of the strategies

Different strategies for the DAC protocol corresponds to different choices of the weights matrix. Clearly, the choice of a particular weight matrix depends on the available information at every node about the network topology, and on their specific computational requirements. Apart from the max-degree strategy, we consider three additional strategies, which are briefly detailed next.

Metropolis-Hastings

The Metropolis-Hastings weights matrix is defined as:

$$C_{ij} = \begin{cases} \frac{1}{\max\{d_i, d_j\}+1} & j \in \mathcal{N}_i \\ 1 - \sum_{j \in \mathcal{N}_i} \frac{1}{\max\{d_i, d_j\}+1} & i = j \\ 0 & \text{otherwise} \end{cases}. \quad (5.6)$$

Differently from the max-degree strategy, the Metropolis-Hastings strategy does not require the knowledge of global information (the maximum degree) about the network topology, but requires that each node knows the degrees of all its neighbors.

Minimum Asymptotic

The third matrix strategy considered here corresponds to the optimal strategy introduced in [198], wherein the weights matrix is constructed to minimize the asymptotic convergence factor $\rho(\mathbf{C} - \frac{\mathbf{1}\mathbf{1}^T}{L})$, where $\rho(\cdot)$ denotes the spectral radius operator. This is achieved by solving the constrained optimization problem:

$$\begin{aligned} & \text{minimize} && \rho(\mathbf{C} - \frac{\mathbf{1}\mathbf{1}^T}{L}) \\ & \text{subject to} && \mathbf{C} \in \mathcal{C}, \quad \mathbf{1}^T \mathbf{C} = \mathbf{1}^T, \quad \mathbf{C} \mathbf{1} = \mathbf{1} \end{aligned} \quad (5.7)$$

where \mathcal{C} is the set of possible weight matrices. Problem (5.7) is non-convex, but it can be shown to be equivalent to a semidefinite programming (SDP) problem [198], solvable using efficient ad-hoc algorithms.

Laplacian Heuristic

The fourth and last matrix considered here is an heuristic approach based on constant edge weights matrix [198]:

$$\mathbf{C} = \mathbf{I} - \alpha \mathbf{L}, \quad (5.8)$$

where $\alpha \in \mathbb{R}$ is a user-defined parameter, and \mathbf{L} is the Laplacian matrix associated to the network (see Appendix A). For weights matrices in the form of (5.8), the asymptotic convergence factor satisfies:

$$\begin{aligned} \rho(\mathbf{C} - \frac{\mathbf{1}\mathbf{1}^T}{L}) &= \max\{\lambda_2(\mathbf{C}), -\lambda_n(\mathbf{C})\} \\ &= \max\{1 - \alpha \lambda_{n-1}(\mathbf{L}), \alpha \lambda_1(\mathbf{L}) - 1\}, \end{aligned} \quad (5.9)$$

where $\lambda_i(\mathbf{C})$ denotes the i -th eigenvalue associated to \mathbf{C} . The value of α that minimizes (5.9) is given by:

$$\alpha^* = \frac{2}{\lambda_1(\mathbf{L}) + \lambda_{L-1}(\mathbf{L})}. \quad (5.10)$$

5.3.2 Experimental Results

We compare the performance of the 4 different strategies illustrated in the previous section in terms of number of iterations required to converge to the average and speed of convergence. In order to avoid that a particular network topology compromises the statistical significance of the experiments, we perform 25 rounds of simulation. In each round, we generate a random topology for an 8-nodes network, according to the Erdős-Rényi model with $p = 0.5$. We consider 2 datasets: G50C (detailed in Section 4.3.1); and CCPP, a regression dataset with 4 features and 9568 examples taken from the UCI repository.³ At each round, datasets are subdivided in batches following the procedure detailed in the previous experimental section. Since in real applications the value of the average is not available to the nodes, in order to evaluate the number of iterations, we consider that all the nodes reached consensus when $\|\beta_k(t) - \beta_k(t-1)\|^2 \leq 10^{-6}$ for any possible value of k . In Fig. 5.4 we show the average number of iterations required by the DAC protocol, averaged over the rounds.

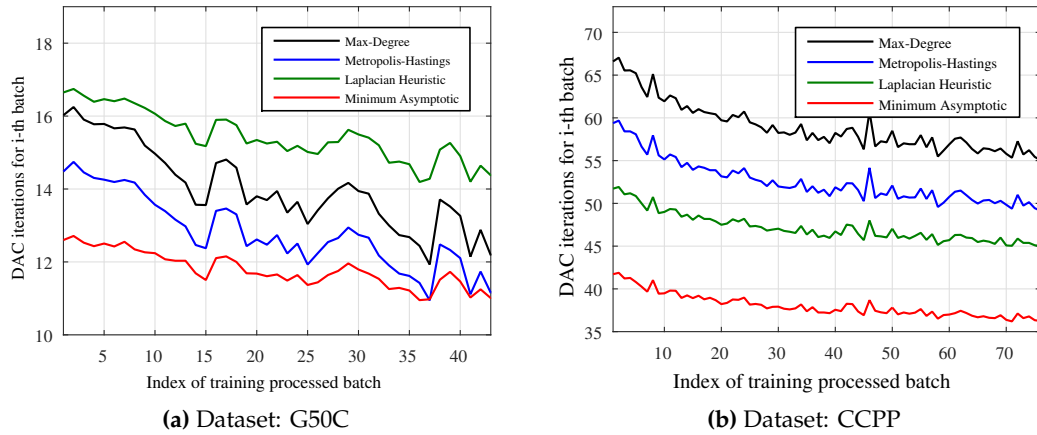


Figure 5.4: Evolution of the DAC iterations required by the considered strategies to converge to the average, when processing successive amounts of training batches.

The x -axis in Fig. 5.4 shows the index of the processed batch. As expected, the number of DAC iterations shows a decreasing trend as the number of processed training batches grows, since the nodes are slowly converging to a single RVFL model. The main result in Fig. 5.4, however, is that a suitable choice of the mixing

³<https://archive.ics.uci.edu/ml/datasets/Combined+Cycle+Power+Plant>

strategy can significantly improve the convergence time (and hence the training time) required by the algorithm. In particular, the optimal strategy defined by Eq. (5.7) achieves the best performance, with a reduction of the required number of iterations up to 35% and 28% when compared with max-degree and Metropolis-Hasting strategies respectively. On the other side, the strategy based on constant edge matrix in Eq. (5.8) shows different behaviors for the 2 datasets, probably due to the heuristic nature of this strategy.

The second experiment, whose results are shown in Fig. 5.5, is to show the speed of convergence for the considered strategies. This is made by evaluating the trend of the relative network disagreement:

$$RND(t) = \frac{1}{N} \sum_{i=1}^N \frac{\|\beta_i(t) - \hat{\beta}\|^2}{\|\beta_i(0) - \hat{\beta}\|^2}, \quad (5.11)$$

as the number of DAC iterations increases. The value of $\hat{\beta}$ in Eq. (5.11) is the true

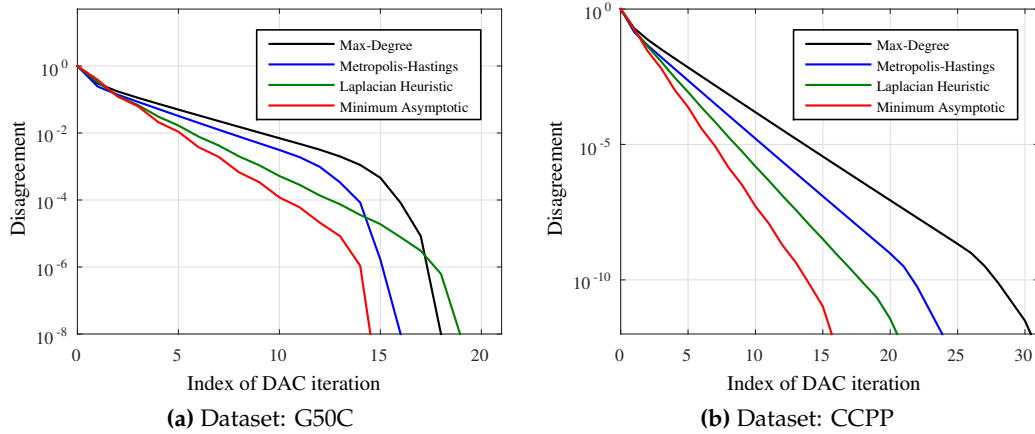


Figure 5.5: Evolution of the relative network disagreement for the considered strategies as the number of DAC iterations increases. The y -axis is shown with a logarithmic scale.

average. The y -axis in Fig. 5.5 is shown with a logarithmic scale. Results show that the “optimal” strategy has the fastest speed of convergence, as expected, while it is interesting to notice how, when compared to max-degree and Metropolis-Hastings weights, the heuristic strategy achieves a rapid decrease in disagreement in the initial iterations, while its speed become slower in the end (this is noticeable in Fig. 5.5a). This may help to explain the lower performance of this strategy in Fig. 5.4a.

Overall, this set of experimental results show how an appropriate choice of the weights matrix can lead to considerable improvements both in the number of iterations required by the protocol to converge to the average, and in the speed of convergence. In particular, when compared to other strategies, an “optimal” choice of the weights matrix can save up to 30% in time.

Distributed RVFL Networks with Vertically Partitioned Data

Contents

6.1	Derivation of the algorithm	61
6.2	Experimental setup	64
6.3	Results and discussion	65

THIS chapter presents an extension of the ADMM-RVFL algorithm presented in Section 4.2.2 to the case of vertically partitioned (VP) data. In the VP scenario, the features of every pattern are partitioned over the nodes. A prototypical example of this is found in the field of distributed databases [87], where several organizations possess only a partial view on the overall dataset (e.g., global health records distributed over multiple medical databases). In the centralized case, this is also known as the problem of learning from heterogeneous sources, and it is typically solved with the use of ensemble procedures [91]. However, as we show in our experimental results, in the VP setting naive ensembles over a network tend to achieve highly sub-optimal results, with respect to a fully centralized solution.

6.1 Derivation of the algorithm

We suppose that the k th agent has access to a subset \mathbf{x}_k of features, such that:

$$\mathbf{x} = [\mathbf{x}_1 | \dots | \mathbf{x}_L]$$

The main problem for the distributed training of an RVFL network in this setting is that the computation of any functional link in Eq. (4.1) requires knowledge of the full sample. However, as we stated in the previous chapters, we would like to avoid

The content of this chapter is adapted from the material published in [150].

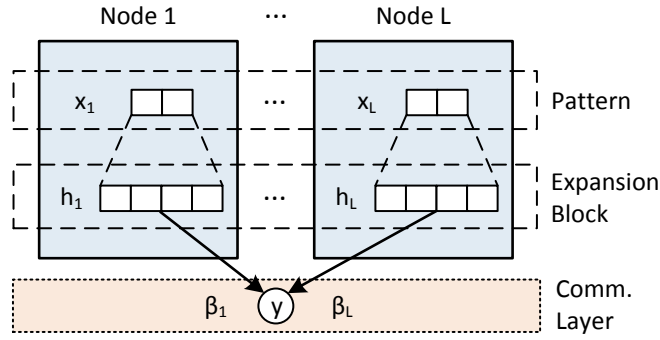


Figure 6.1: Schematic description of the proposed algorithm. Each node has access to a subset of the global pattern. This local feature vector is projected to a local expansion block, and the overall output is computed by a linear combination of the local expansions, through a suitable communication layer.

exchange of data patterns, due to both size and privacy concerns. To this end, we approximate model in Eq. (4.1) by considering local expansion blocks:

$$f(\mathbf{x}) = \sum_{k=1}^L (\boldsymbol{\beta}_k)^T \mathbf{h}_k(\mathbf{x}_k). \quad (6.1)$$

In this way, each term $\mathbf{h}_k(\mathbf{x}_k)$ can be computed locally. Input vectors and expansion blocks may have different lengths at every node, depending on the application and on the local computational requirements. This is shown pictorially in Fig. 6.1.

The overall optimization problem becomes:

$$\arg \min_{\boldsymbol{\beta}} \frac{1}{2} \left\| \sum_{k=1}^L \mathbf{H}_k \boldsymbol{\beta}_k - \mathbf{y} \right\|_2^2 + \frac{\lambda}{2} \sum_{k=1}^L \|\boldsymbol{\beta}_k\|_2^2, \quad (6.2)$$

where \mathbf{H}_k denotes the hidden matrix computed at the k -th node, such that $\mathbf{H} = [\mathbf{H}_1 | \dots | \mathbf{H}_L]$. The ADMM optimization algorithm can be adapted to this setting, as shown in [20, Section 8.3]. To this end, we consider the equivalent optimization problem:

$$\begin{aligned} \underset{\boldsymbol{\beta}}{\text{minimize}} \quad & \frac{1}{2} \left\| \sum_{k=1}^L \mathbf{z}_k - \mathbf{y} \right\|_2^2 + \frac{\lambda}{2} \sum_{k=1}^L \|\boldsymbol{\beta}_k\|_2^2 \\ \text{subject to} \quad & \mathbf{H}_k \boldsymbol{\beta}_k - \mathbf{z}_k = \mathbf{0}, \quad k = 1 \dots L. \end{aligned} \quad (6.3)$$

where we introduced local variables $\mathbf{z}_k = \mathbf{H}_k \boldsymbol{\beta}_k$. The augmented Lagrangian of this

problem is given by:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\beta}_k, \mathbf{z}_k, \mathbf{t}_k) &= \frac{1}{2} \left\| \sum_{k=1}^L \mathbf{z}_k - \mathbf{y} \right\|_2^2 + \frac{\lambda}{2} \sum_{k=1}^L \|\boldsymbol{\beta}_k\|_2^2 + \\ &+ \sum_{k=1}^L \mathbf{t}_k^T (\mathbf{H}_k \boldsymbol{\beta}_k - \mathbf{z}_k) + \frac{\rho}{2} \sum_{k=1}^L \|\mathbf{H}_k \boldsymbol{\beta}_k - \mathbf{z}_k\|_2^2, \end{aligned} \quad (6.4)$$

where \mathbf{t}_k are the Lagrange multipliers, $\rho \in \mathbb{R}^+$ is a regularization factor, and the last term is added to ensure convergence. The solution to problem (6.2) can be computed by iterating the updates in Eqs. (4.10)-(4.12). Following the derivation in [20, Section 8.3], and computing the gradient terms, the final updates can be expressed as:

$$\boldsymbol{\beta}_k[n+1] = \left(\frac{\lambda}{\rho} \mathbf{I} + \mathbf{H}_k^T \mathbf{H}_k \right)^{-1} \mathbf{H}_k^T (\mathbf{H}_k \boldsymbol{\beta}_k[n] + \bar{\mathbf{z}}[n] - \overline{\mathbf{H}\boldsymbol{\beta}}[n] - \mathbf{t}[n]), \quad (6.5)$$

$$\bar{\mathbf{z}}[n+1] = \frac{1}{L + \rho} (\mathbf{y} + \overline{\mathbf{H}\boldsymbol{\beta}}[n+1] + \mathbf{t}[n]), \quad (6.6)$$

$$\mathbf{t}[n+1] = \mathbf{t}[n] + \overline{\mathbf{H}\boldsymbol{\beta}}[n+1] - \bar{\mathbf{z}}[n+1], \quad (6.7)$$

where we defined the averages $\overline{\mathbf{H}\boldsymbol{\beta}}[n] = \frac{1}{L} \sum_{k=1}^L \mathbf{H}_k \boldsymbol{\beta}_k[n]$, and $\bar{\mathbf{z}}[n] = \sum_{k=1}^L \mathbf{z}_k[n]$. Additionally, the variables \mathbf{t}_k can be shown to be equal between every node [20], so we removed the subscript. Convergence of the algorithm can be tracked locally by computing the residual:

$$\mathbf{r}_k[n] = \mathbf{H}_k^T \boldsymbol{\beta}_k[n] - \mathbf{z}_k[n]. \quad (6.8)$$

It can be shown that, for the iterations defined by Eqs. (6.5)-(6.7), $\|\mathbf{r}_k[n]\|_2 \rightarrow 0$ as $n \rightarrow +\infty$, with the solution converging asymptotically to the solution of problem in Eq. (6.2). The overall algorithm, denoted as VP-ADMM-RVFL, is summarized in Algorithm 6.1.

After training, every node has access to its own local mapping $\mathbf{h}_k(\cdot)$, and to its subset of coefficients $\boldsymbol{\beta}_k$. Differently from the horizontally partitioned (HP) scenario, when the agents require a new prediction, the overall output defined by Eq. (6.1) has to be computed in a decentralized fashion. Once again, this part will depend on the actual communication layer available to the agents. As an example, it is possible to run the DAC protocol over the values $(\boldsymbol{\beta}_k)^T \mathbf{h}_k(\mathbf{x}_k)$, such that every node obtain a suitable approximation of $\frac{1}{L} f(\mathbf{x})$. For smaller networks, it is possible to compute an Hamiltonian cycle between the nodes [129]. Once the cycle is known to the agents, they can compute Eq. (6.1) by forward propagating the partial sums up to the final node of the cycle, and then back-propagating the result. Clearly, many other choices are possible, depending on the network.

Algorithm 6.1 VP-ADMM-RVFL: Extension of ADMM-RVFL to vertically partitioned data (k th node).

Inputs: Training set S_k , number of nodes L (global), regularization factors λ, γ (global), maximum number of iterations T (global)

Output: Optimal weight vector β^*

- 1: Select parameters $\mathbf{w}_1, \dots, \mathbf{w}_B$, in agreement with the other $L - 1$ nodes.
 - 2: Compute \mathbf{H}_k and \mathbf{y}_k from S_k .
 - 3: Initialize $\mathbf{t}[0] = \mathbf{0}$, $\bar{\mathbf{z}}[0] = \mathbf{0}$.
 - 4: **for** n from 0 to T **do**
 - 5: Compute $\beta_k[n + 1]$ according to Eq. (6.5).
 - 6: Compute averages $\bar{\mathbf{H}}\beta[n]$ and $\bar{\mathbf{z}}[n]$ with DAC.
 - 7: Compute $\bar{\mathbf{z}}[n + 1]$ according to Eq. (6.6).
 - 8: Update $\mathbf{t}[n]$ according to Eq. (6.7).
 - 9: Check termination with residuals.
 - 10: **end for**
 - 11: **return** $\bar{\mathbf{z}}[n + 1]$
-

6.2 Experimental setup

In this section, we present an experimental validation of the proposed algorithm on three classification tasks: Garageband, G50C and Sylva (detailed in Sections 4.3.1 and 5.2.2). Optimal parameters for the RVFL network are taken from the corresponding sections. In our first set of experiments, we consider networks of 8 agents, whose connectivity is randomly generated such that every pair of nodes has a 60% probability of being connected, with the only global requirement that the overall network is connected. The input features are equally partitioned through the nodes, i.e., every node has access to roughly $d/8$ features, where d is the dimensionality of the dataset. We compare the following algorithms:

Centralized RVFL (C-RVFL): this corresponds to the case where a fusion center is available, collecting all local datasets and solving directly Eq. (4.3). Settings for this model are the optimal ones.

Local RVFL (L-RVFL): this is a naive implementation, where each node trains a local model with its own dataset, and no communication is performed. Accuracy of the models is then averaged throughout the L nodes. As a general settings, we employ the same regularization coefficient for every node as C-RVFL, and $B_k = \lceil B/8 \rceil$ expansions in every agent.

Ensemble RVFL (ENS-RVFL): this corresponds to a basic distributed ensemble. As for L-RVFL, during the training phase every node trains a local model

with its own dataset. In the testing phase, the nodes agree on a single class prediction by taking a majority vote over their local predictions. Parameters are the same as for L-RVFL.

Distributed RVFL (VP-ADMM-RVFL): this is trained using the distributed protocol introduced in the previous section. Settings are the same as L-RVFL, while for the ADMM we set $\rho = 0.1$ and a maximum number of 200 iterations.

To compute the misclassification rate, we perform a 3-fold cross-validation on the overall dataset, and repeat the procedure 15 times.

6.3 Results and discussion

Results of the experiments are presented in Table 6.1.

Table 6.1: Misclassification error and training time for the four algorithms. Results are averaged over the 8 different nodes of the network. Standard deviation is provided between brackets.

Dataset	Algorithm	Misclassification error [%]	Training time [secs.]
Garageband	C-RVFL	41.32 (± 1.24)	0.03 (± 0.02)
	L-RVFL	82.79 (± 3.82)	0.01 (± 0.01)
	ENS-RVFL	61.01 (± 1.97)	0.01 (± 0.01)
	VP-ADMM-RVFL	41.34 (± 1.34)	2.35 (± 0.58)
Sylva	C-RVFL	1.18 (± 0.13)	0.44 (± 0.06)
	L-RVFL	49.80 (± 36.35)	0.05 (± 0.02)
	ENS-RVFL	6.04 (± 0.12)	0.06 (± 0.02)
	VP-ADMM-RVFL	1.22 (± 0.15)	1.94 (± 0.40)
G50C	C-RVFL	5.80 (± 1.19)	0.05 (± 0.02)
	L-RVFL	49.51 (± 6.37)	0.01 (± 0.01)
	ENS-RVFL	10.98 (± 2.32)	0.01 (± 0.01)
	VP-ADMM-RVFL	5.80 (± 1.37)	0.38 (± 0.16)

It can be seen that, despite we approximate the global expansion block of C-RVFL using L distinct local expansions, this has minimal or no impact on the global solution. In fact, VP-ADMM-RVFL is able to achieve performance comparable to C-RVFL in all three datasets, while the ensemble approach is performing relatively poorly: it has a 20%, 5% and 5% increase in error respectively in each dataset. This shows that the relatively common approach of averaging over the local models may be highly sub-optimal in practical situations.

As a reference, in Table 6.1 we also provide the average training time spent at every node. However, we note that in our experiments the network was simulated

in a serial architecture, removing all communication costs. Clearly, a practical analysis of this point would require knowledge of the communication layer, which goes beyond the scope of the thesis. Still, we can see from the fourth column of Table 6.1 that the proposed algorithm requires an acceptable computational time for performing the 200 iterations, since the matrix inversions in Eq. (6.5) can be pre-computed at the beginning of the training process. Additionally, we add that the training time of VP-ADMM-RVFL can be greatly reduced in practice by the implementation of an efficient stopping criterion as in the previous chapter.

Finally, we show the evolution of the misclassification error for VP-ADMM-RVFL and ENS-RVFL when varying the size of the network from $L = 4$ to $L = 12$. Results of this experiment are given in Fig. 6.2 (a)-(c). Settings are kept fixed with respect to the previous experiment, while the features are equally partitioned as before (hence, for smaller networks each node has access to a larger subset of features). Performance of C-RVFL is given as a comparison with a dashed black line. As expected, we see that, although the behavior of ENS-RVFL strongly depends on the number of nodes in the network, VP-ADMM-RVFL is resilient to such change, always approximating very well the centralized performance. It is also interesting to note that the behavior of ENS-RVFL is not always monotonically increasing, as is shown in Fig. 6.2-(c), possibly due to its ensembling characteristics and to the artificial nature of the G50C dataset.

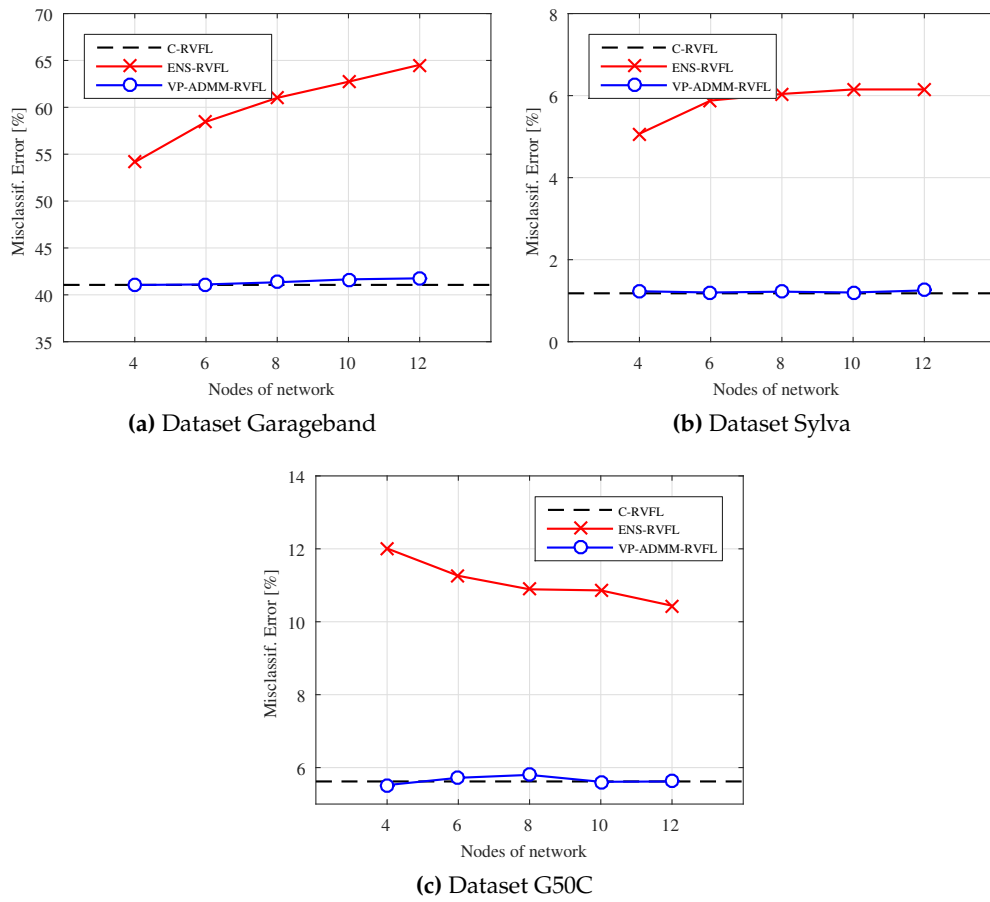


Figure 6.2: Evolution of the error for VP-ADMM-RVFL and ENS-RVFL when varying the size of the network from $L = 4$ to $L = 12$.

Part III

Distributed Semi-Supervised Learning

Decentralized Semi-supervised Learning via Privacy-Preserving Matrix Completion

Contents

7.1	Introduction	69
7.2	Preliminaries	72
7.2.1	Semi-supervised learning	72
7.2.2	(Euclidean) matrix completion	74
7.2.3	Privacy-preserving similarity computation	75
7.3	Distributed Laplacian Estimation	77
7.3.1	Formulation of the problem	77
7.3.2	Decentralized block estimation	78
7.3.3	Diffusion gradient descent	80
7.4	Distributed Semi-supervised Manifold Regularization	81
7.5	Experimental results	84
7.5.1	Experiments setup	84
7.5.2	Distributed Laplacian estimation	84
7.5.3	Distributed semi-supervised manifold regularization	86
7.5.4	Privacy preservation	88

7.1 Introduction

As we saw in the previous chapters, many centralized SL algorithms have been extended successfully to the distributed setting. However, many crucial sub-areas of machine learning remain to be extended to the fully distributed scenario. Among these, the DL setting could benefit strongly from the availability

The content of this chapter has been (conditionally) accepted for publication at IEEE Transactions on Neural Networks and Learning Systems.

of distributed protocols for semi-supervised learning (SSL) [31]. In SSL, it is assumed that the labeled training data is supplemented by some additional unlabeled data, which has to be suitably exploited in order to improve the test accuracy. State-of-the-art research on SSL is concerned on the single-agent (centralized) case, e.g. with the use of manifold regularization (MR) [11, 110], transductive learning [30], and several others. To the best of our knowledge, the case of SSL over multiple agents has been addressed only in very specific settings, such as localization over WSNs [34], while no algorithm is available for the general case. However, we argue that such an algorithm would be well suited for a wide range of applications. As an example, consider the case of medical diagnosis, with labeled and unlabeled data distributed over multiple clinical databases. Other examples include distributed text classification over peer-to-peer networks, distributed music classification (which we considered in Chapter 5), and so on. In all of them, labeled data at every agent is costly to obtain, while unlabeled data is plentiful. The overall setting is summarized in Fig. 7.1, where each agent in a network receives two training datasets, one composed of labeled patterns and one composed of unlabeled patterns.

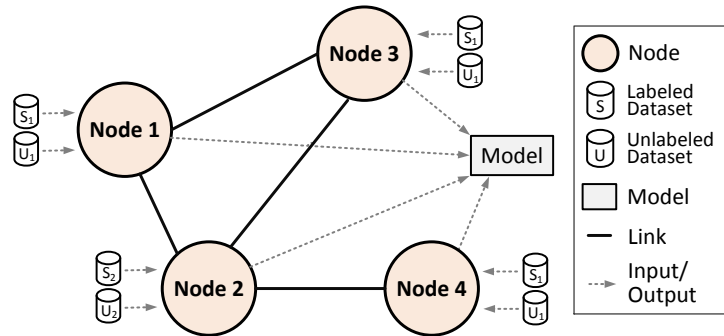


Figure 7.1: Depiction of SSL over a network of agents. Each agent receives a labeled training dataset, together with an unlabeled one. The task is for all the nodes to converge to a single model, by exploiting all their local datasets.

In this chapter, we propose the first fully distributed algorithm for SSL over networks, satisfying the above requirements. In particular, we extend an algorithm belonging to the MR family, namely laplacian kernel ridge regression (LapKRR) [11]. MR algorithms, originated in the seminal works of [10] and [11], are based on the assumption that data often lie in a low-dimensional manifold \mathcal{M} embedded in the higher-dimensional input space. When the structure of the manifold is unknown, it can be approximated well by a weighted graph where the vertexes are represented by the data points and the weights of the edges represent a measure of similarity between the points. In the MR framework, the classification function is obtained by solving an extension of the classical regularized optimization problem, with an additional regularization term, which incorporates information about the function's smoothness on the manifold.

The algorithm presented in this chapter starts from the observation that, in the MR optimization problem, information is mostly encoded in a matrix \mathbf{D} of pairwise distances between patterns. In fact, both the additional regularization term, and the kernel matrix (for any translation-invariant kernel function) can be computed using the information about the distance between points. In the distributed setting, each agent can compute this matrix relatively only to its own training data, while information about the distance between points belonging to different agents are unknown. Obtaining this information would allow a very simple protocol for solving the overall optimization problem. As a consequence, we subdivide the training algorithm in two steps: a distributed protocol for computing \mathbf{D} , followed by a distributed strategy for solving the optimization problem.

For the former step, in the initial phase of the algorithm, we allow a small exchange of data patterns between agents. In this phase, privacy can be preserved with the inclusion of any privacy-preserving protocol for the computation of distances [186]. For completeness, we describe the strategies that are used in our experiments in Section 7.2.3. As a second step, we recover the rest of the global distance matrix \mathbf{D} by building on previous works on Euclidean distance matrix (EDM) completion [23, 112]. To this end, we consider two strategies. The first one is a simple modification of the state-of-the-art algorithm presented in [92, 94], which is based on a column-wise partitioning of \mathbf{D} over the agents. In this chapter, we modify it to take into account the specific nature of Euclidean distance matrices, by the incorporation of non-negativity and symmetry constraints. As a second strategy, we propose a novel algorithm for EDM completion, which is inspired to the framework of diffusion adaptation (DA) (see Section 3.4.2). The algorithm works by interleaving gradient descent steps with local interpolation of a suitable low-rank factorization of \mathbf{D} . While the first algorithm has a lower computational cost, we found that this comes at the cost of a worse performance, particularly when the sampling set of the matrix to complete is small. On the opposite, our algorithm exploits the particular structure of EDMs, at the cost of a possibly greater computational demanding. We discuss in more detail the advantages and disadvantages of the two approaches in Section 7.3 and in the experimental section.

As we stated before, once the matrix \mathbf{D} is known, solving the rest of the optimization problem is trivial. In this chapter we focus on the LapKRR algorithm, and we show that its distributed version can be solved using a single operation of sum over the network. Our experimental results show that, in most cases, the performance of the novel diffusion adaptation-based algorithm for distributed EDM completion overcome those of the state-of-the-art column-wise partitioning strategy. Secondly, experiments show that the distributed LapKRR is competitive with a centralized LapKRR model trained on the overall dataset.

The rest of the chapter is structured as follows: in Section 7.2 we introduce the

theoretical tools upon which our algorithm is based. In particular, we detail the problem of SSL in the framework of MR in Section 7.2.1, some notions of EDM completion in Section 7.2.2, and two strategies for privacy-preserving similarity computation in Section 7.2.3. In Section 7.3 we propose our algorithm to complete an EDM in a decentralized fashion. Then, Section 7.4 details the proposed framework for distributed LapKRR. In Section 7.5 we present the results for both the distributed EDM completion and distributed LapKRR.

7.2 Preliminaries

In this section, we introduce some concepts that are used in the development of our algorithm. We start by describing the basic setting of SSL in Section 7.2.1. Then, we introduce the matrix completion problem and its application to the EDMs in Section 7.2.2. As the last point, in Section 7.2.3 we report some results on privacy-preserving similarity computation.

7.2.1 Semi-supervised learning

In the SSL setting, we are provided with a set of l input/output labeled data $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_l, y_l)\}$ and an additional set of u unlabeled data $U = \{\mathbf{x}_{l+1}, \dots, \mathbf{x}_{l+u}\}$ [31]. As before, in the following inputs are assumed to be d -dimensional real vectors $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d$, while outputs are assumed to be scalars $y \in \mathcal{Y} \subseteq \mathbb{R}$. The discussion can be extended straightforwardly to the case of a multi-dimensional output. In this chapter, we consider one particular class of SSL algorithms belonging to the family of MR [11]. Practically, MR learning algorithms are based on three assumptions.

- Smoothness assumption: if two points $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ are close in the intrinsic geometry of their marginal distribution, then their conditional distributions $p(y | \mathbf{x}_1)$ and $p(y | \mathbf{x}_2)$ are similar.
- Cluster assumption: the decision boundary should lie in a low-density region of the input space \mathcal{X} .
- Manifold assumption: the marginal distribution $p(\mathbf{x})$ is supported on a low-dimensional manifold \mathcal{M} embedded in \mathcal{X} .

We now define the SLL problem formally.

Definition 7 (SSL problem with manifold regularization)

Let $\mathcal{H}_{\mathcal{K}}$ be a Reproducing Kernel Hilbert Space defined by the kernel function $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ with norm $\|f\|_{\mathcal{K}}^2$, the approximation function for the SSL problem is

estimated by solving:

$$f^* = \underset{f \in \mathcal{H}_K}{\operatorname{argmin}} \sum_{i=1}^l l(y_i, f(\mathbf{x}_i)) + \gamma_A \|f\|_K^2 + \gamma_I \|f\|_I^2, \quad (7.1)$$

where $l(\cdot, \cdot)$ is a suitable loss function, $\|f\|_I^2$ is a penalty term that penalizes the structure of f with respect to the manifold and $\gamma_A, \gamma_I \geq 0$ are the regularization parameters.

Usually, the structure of the manifold \mathcal{M} is unknown and it must be estimated from both labeled and unlabeled data. In particular, we can define an adjacency matrix $\mathbf{W} \in \mathbb{R}^{l+u \times l+u}$, where each entry W_{ij} is a measure of similarity between patterns \mathbf{x}_i and \mathbf{x}_j (see [11] for possible ways of constructing this matrix). Using this, the regularization term $\|f\|_I^2$ can be rewritten as [11]:

$$\|f\|_I^2 = f^T \mathbf{L} f, \quad (7.2)$$

where $\mathbf{L} \in \mathbb{R}^{l+u \times l+u}$ is the data adjacency graph Laplacian (see Appendix A). Practically, the overall manifold \mathcal{M} is approximated with an adjacency graph, which can be computed from both labeled and unlabeled data. In order to obtain better performances, usually a normalized Laplacian $\hat{\mathbf{L}} = \mathbf{G}^{-1/2} \mathbf{L} \mathbf{G}^{-1/2}$, or an iterated version $\hat{\mathbf{L}}^q$, $q \geq 0$, is used [11]. An extension of the classical Representer Theorem proves that the function f^* is in the form of:

$$f^*(\mathbf{x}) = \sum_{i=1}^N \alpha_i \mathcal{K}(\mathbf{x}, \mathbf{x}_i), \quad (7.3)$$

where $N = l + u$ and α_i are weight parameters. As we stated in the introduction, for simplicity in this chapter we focus on a particular algorithm belonging to this framework, denoted as LapKRR. This is obtained by substituting Eq. (7.3) into problem (7.1) and setting a squared loss function:

$$l(y_i, f(\mathbf{x}_i)) = \|y_i - f(\mathbf{x}_i)\|_2^2. \quad (7.4)$$

Considering the dual optimization problem, by the optimality conditions the final parameters vector $\boldsymbol{\alpha}^* = [\alpha_1, \dots, \alpha_N]^T$ is easily obtained as:

$$\boldsymbol{\alpha}^* = (\mathbf{J}\mathbf{K} + \gamma_A \mathbf{I} + \gamma_I \mathbf{L}\mathbf{K})^{-1} \hat{\mathbf{y}}, \quad (7.5)$$

where $\hat{\mathbf{y}}$ is an N -dimensional vector with components:

$$\hat{y}_i = \begin{cases} y_i & \text{if } i \in \{1, \dots, l\} \\ 0 & \text{if } i \in \{l+1, \dots, l+u\} \end{cases}, \quad (7.6)$$

\mathbf{J} is an $N \times N$ diagonal matrix with elements:

$$J_{ii} = \begin{cases} 1 & \text{if } i \in \{1, \dots, l\} \\ 0 & \text{if } i \in \{l+1, \dots, l+u\} \end{cases}, \quad (7.7)$$

and finally \mathbf{K} is the $N \times N$ kernel matrix defined by $\{K_{ij} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)\}$.

7.2.2 (Euclidean) matrix completion

The second notion that will be used in the proposed algorithm is the EDM completion problem [3]. A matrix completion problem is defined as the problem of recovering the missing entries of a matrix only from a set of known entries [23]. This problem has many practical applications, i.e. sensors localization, covariance estimation and customer recommendations, and it was largely investigated in the literature.

In this chapter, we focus on completion of the square matrix $\mathbf{D} \in \mathbb{R}^{N \times N}$ containing the pairwise distances among the training patterns, i.e.:

$$D_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \quad \forall i, j = 1, \dots, N. \quad (7.8)$$

\mathbf{D} is called an Euclidean Distance Matrix (EDM). Clearly, Eq. (7.8) implies that \mathbf{D} is symmetric and $D_{ii} = 0$ for all the elements on the main diagonal. It is possible to show that the rank r of \mathbf{D} is upper bounded by $d + 2$, meaning that \mathbf{D} is low-rank whenever $d \ll N$, which is common in all practical applications.

In the following, we suppose to have observed only a subset of entries of \mathbf{D} , in the form of a matrix $\hat{\mathbf{D}}$. More formally, there exists a matrix with binary entries $\mathbf{\Omega} \in [0, 1]^{N \times N}$ such that:

$$\hat{\mathbf{D}} = \begin{cases} \hat{D}_{ij} = D_{ij} & \text{if } \Omega_{ij} = 1 \\ \hat{D}_{ij} = 0 & \text{otherwise} \end{cases}. \quad (7.9)$$

We wish to recover the original matrix \mathbf{D} from $\hat{\mathbf{D}}$, i.e. we want to solve the following optimization problem:

$$\min_{\mathbf{D} \in \text{EDM}(N)} \|\mathbf{\Omega} \circ (\hat{\mathbf{D}} - \mathbf{D})\|_F^2, \quad (7.10)$$

where \circ denotes the Hadamard product between two matrices, $\text{EDM}(N)$ is the set of all EDMs of size N , and $\|\mathbf{A}\|_F$ is the Frobenius norm of matrix \mathbf{A} . It is possible to reformulate problem in Eq. (7.10) as a semidefinite problem by considering the

Schoenberg mapping between EDMs and positive semidefinite matrices [3]:

$$\begin{aligned} \min_{\mathbf{D}} \quad & \left\| \boldsymbol{\Omega} \circ [\hat{\mathbf{D}} - \kappa(\mathbf{D})] \right\|_{\text{F}}^2, \\ \text{s. t.} \quad & \mathbf{D} \geq 0 \end{aligned} \quad (7.11)$$

where $\mathbf{D} \geq 0$ means that \mathbf{D} is positive semidefinite and:

$$\kappa(\mathbf{D}) = \text{diag}(\mathbf{D})\mathbf{1}^{\text{T}} + \mathbf{1}\text{diag}(\mathbf{D})^{\text{T}} - 2\mathbf{D}, \quad (7.12)$$

such that $\text{diag}(\mathbf{D})$ extracts the main diagonal of \mathbf{D} as a column vector. This observation motivated most of the initial research on EDM completion [3]. Recently, an alternative formulation was proposed in [112], which exploits the fact that every positive semidefinite matrix \mathbf{D} with rank r admits a factorization $\{\mathbf{D} = \mathbf{V}\mathbf{V}^{\text{T}}\}$, where $\mathbf{V} \in \mathbb{R}_*^{N \times r} = \{\mathbf{V} \in \mathbb{R}^{N \times r} : \det(\mathbf{V}^{\text{T}}\mathbf{V}) \neq 0\}$. Using this factorization and assuming we know the rank of \mathbf{D} , problem (7.11) can be reformulated as:

$$\min_{\mathbf{V}\mathbf{V}^{\text{T}} \in S_+(r, N)} \left\| \boldsymbol{\Omega} \circ [\hat{\mathbf{D}} - \kappa(\mathbf{V}\mathbf{V}^{\text{T}})] \right\|_{\text{F}}^2, \quad (7.13)$$

where we have:

$$S_+(r, N) = \{\mathbf{U} \in \mathbb{R}^{N \times N} : \mathbf{U} = \mathbf{U}^{\text{T}} \geq 0, \text{rank}(\mathbf{U}) = r\}. \quad (7.14)$$

7.2.3 Privacy-preserving similarity computation

As we stated in the Introduction, a fundamental step in the algorithm presented in this chapter is a distributed computation of similarity between two training patterns, i.e. a distributed computation of a particular entry of \mathbf{D} . If these patterns cannot be exchanged over the network, e.g. for privacy reasons, there is the need of implementing suitable protocols for privacy-preserving similarity computation. To show the applicability of the proposed approach, in our experimental simulations we make use of two state-of-the-art solutions to this problem. For completeness, we detail them here briefly.

More formally, the problem can be stated as follows. Given two training patterns $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^d$, belonging to different agents, we want to compute $\mathbf{x}_i^{\text{T}}\mathbf{x}_j$, without revealing the two patterns. Clearly, computing the inner product allows the computation of several other distance metrics, including the standard L_2 Euclidean norm. The first strategy that we investigate here is the random projection-based technique developed in [95]. Suppose that both agents agree on a projection matrix $\mathbf{R} \in \mathbb{R}^{m \times d}$, with $m < d$, such that each entry R_{ij} is independent and chosen from a normal distribution with mean zero and variance σ^2 . We have the following lemma:

Lemma 1

Given two input patterns $\mathbf{x}_i, \mathbf{x}_j$, and the respective projections:

$$\mathbf{u}_i = \frac{1}{\sqrt{m\sigma}} \mathbf{R}\mathbf{x}_i, \text{ and } \mathbf{u}_j = \frac{1}{\sqrt{m\sigma}} \mathbf{R}\mathbf{x}_j, \quad (7.15)$$

we have that:

$$\mathbb{E} \{ \mathbf{u}_i^T \mathbf{u}_j \} = \mathbf{x}_i^T \mathbf{x}_j. \quad (7.16)$$

Proof 2

See [95, Lemma 5.2]. □

In light of Lemma 1, exchanging the projected patterns instead of the original ones allows preserving, on average, the inner product. A thorough investigation on the privacy-preservation guarantees of this protocol can be found in [95]. Additionally, we can observe that this protocol provides a reduction on the communication requirements of the application, since it effectively reduces the dimensionality of the patterns to be exchanged by a factor m/d .

The second protocol that we investigate in our experimental section is a more general (nonlinear) transformation introduced in [13]. It is given by:

$$\mathbf{v} = \mathbf{b} + \mathbf{Q} \tanh(\mathbf{a} + \mathbf{C}\mathbf{x}), \quad (7.17)$$

for a generic input pattern \mathbf{x} , where $\mathbf{b} \in \mathbb{R}^m$, $\mathbf{Q} \in \mathbb{R}^{m \times t}$, $\mathbf{a} \in \mathbb{R}^t$, $\mathbf{C} \in \mathbb{R}^{t \times d}$ are matrices whose entries are drawn from normal distributions with mean zero and possibly different variances. As in the previous method, it is possible to show that the inner product is approximately preserved, provided that the input patterns are not “outliers” in a specific sense. See [13] for more details and an analysis of the privacy-preservation capabilities of this scheme. Again, choosing t and m allows to balance between a more accurate reconstruction and a reduction on the input dimensionality.

The field of privacy-preserving similarity computation, and more in general privacy-preserving data mining, is vast and with more methods introduced each year. Although we have chosen these two protocols due to their wide diffusion and

simplicity, we stress that our algorithm does not depend specifically on any of them. We refer to [186] and references therein for more general investigations on this field.

7.3 Distributed Laplacian Estimation

In this section, we start by formulating a problem of distributed estimation of \mathbf{L} in Section 7.3.1. Then, we focus on two algorithms for its solution. The first is a modification of a state-of-the-art algorithm, described in Section 7.3.2, while the second is a fully novel protocol which is based on the ideas of ‘diffusion adaptation’ [145] introduced in Section 7.3.3.

7.3.1 Formulation of the problem

In the distributed Laplacian estimation problem, we suppose that both the labeled data and the unlabeled data are distributed through a network of L interconnected agents, as shown in Fig. 7.1 and described in Appendix A. Without loss of generality, we assume that data is organized as follows: the k th agent is provided with N_k patterns, such that $N = \sum_{k=1}^L N_k$. For each agent, the first l_k patterns are labeled: $S_k = \{(\mathbf{x}_{k,1}, y_{k,1}), \dots, (\mathbf{x}_{k,l_k}, y_{k,l_k})\}$, while the last u_k are unlabeled: $U_k = \{\mathbf{x}_{k,l_k+1}, \dots, \mathbf{x}_{k,l_k+u_k}\}$. The local data sets are non-overlapping, so we have $S = \cup_{k=1}^L S_k$ and $U = \cup_{k=1}^L U_k$.

Let $\mathbf{L}_k \in \mathbb{R}^{N_k \times N_k}$, $k = 1 \dots L$, be the Laplacian matrices computed by each agent using its own data; we are interested in estimating in a totally decentralized fashion the Laplacian matrix \mathbf{L} calculated with respect to all the N patterns. The local Laplacian matrices can be always expressed, rearranging the rows and the columns, as block matrices on the main diagonal of \mathbf{L} :

$$\mathbf{L} = \begin{bmatrix} \mathbf{L}_1 & ? & ? \\ ? & \ddots & ? \\ ? & ? & \mathbf{L}_L \end{bmatrix} \quad (7.18)$$

The same structure of (7.18) applies also to matrices \mathbf{D} and \mathbf{K} , with \mathbf{D}_k and \mathbf{K}_k representing the distance matrix and kernel matrix computed over the local dataset. This particular structure implies that the sampling set is not random, and makes non-trivial the problem of completing \mathbf{L} solely from the knowledge of the local matrices. At the opposite, the idea of exchanging the entire local datasets between nodes is unfeasible because of the amount of data to share. Instead of completing in a distributed manner the global Laplacian matrix, in this chapter we consider the alternative approach of computing the global EDM \mathbf{D} first, and then using it to calculate the Laplacian. This approach has two advantages:

- We can exploit the structure of EDMs to design efficient algorithms.

- From the global EDM we can compute, in addition to the Laplacian, the kernel matrix \mathbf{K} for all kernel functions \mathcal{K} based on Euclidean distance (e.g. the Gaussian kernel).

Based on these considerations, we propose a framework for the distributed estimation of \mathbf{L} , which consists in five steps:

1. Patterns exchange: every agent exchanges a fraction p of the available input data (both labeled and unlabeled) with its neighbors. This step is necessary so that the agents can increase the number of known entries in their local matrices. In order to maximize the diffusion of the data within the network, this step is iterated $n_{\max}^{(1)}$ times; at every iteration an increasing percentage of shared data is constituted by pattern received by the neighbors in previous iterations. A simple strategy to do this consists, at the iteration n , to choose $\frac{n_{\max}-n+1}{n_{\max}}p$ patterns from the local dataset, and $\frac{n-1}{n_{\max}}p$ patterns received in the previous $n-1$ iterations. In order to preserve privacy, this step can include one of the privacy-preserving strategies showed in Section 7.2.3.
2. Local EDM computation: each agent computes, using its original dataset and the data received from its neighbors, an incomplete approximation $\hat{\mathbf{D}}_k \in \mathbb{R}^{N \times N}$ of the real EDM matrix \mathbf{D} .
3. Entries exchange: the agents exchange a sample of their local EDMs $\hat{\mathbf{D}}_k$ with their neighbors. Again, this step is iterated $n_{\max}^{(2)}$ times using the same rule of step 1.
4. Distributed EDM completion: the agents complete the estimate $\tilde{\mathbf{D}}$ of the global EDM using one of the distributed algorithms presented in the following sections.
5. Global Laplacian estimation: using $\tilde{\mathbf{D}}$ the agents compute the global Laplacian estimate $\tilde{\mathbf{L}}$ and the kernel matrix estimate $\tilde{\mathbf{K}}$.

7.3.2 Decentralized block estimation

As stated in the Introduction, the first algorithm that we take into account for the decentralized completion of \mathbf{D} is a modified version of the algorithm named *D-LMaFit* [92, 94]. To the best of our knowledge, this is the only existing algorithm for distributed matrix completion available in the literature.

Let $\hat{\mathbf{D}}$ be the incomplete global EDM matrix and denote with \mathcal{I} the set of indexes corresponding to its known entries. In a centralized setting, without taking into account the structure of distance matrices, and assuming that the rank r is known,

$\tilde{\mathbf{D}}$ can be completed by solving the problem:

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{B}, \tilde{\mathbf{D}}} \quad & \|\mathbf{A}\mathbf{B} - \tilde{\mathbf{D}}\|_{\text{F}}^2 \\ \text{s. t.} \quad & \tilde{D}_{ij} = \hat{D}_{ij}, \forall (i, j) \in \mathcal{I} \end{aligned} \quad (7.19)$$

where $\mathbf{A} \in \mathbb{R}^{N \times r}$, $\mathbf{B} \in \mathbb{R}^{r \times N}$ represent a suitable low-rank factorization of $\tilde{\mathbf{D}}$.

In extending problem (7.19) to a decentralized setting, the algorithm presented in [94] considers a column-wise partitioning of $\hat{\mathbf{D}}$ over the agents. For simplicity of notation, we suppose here that this partitioning is such that the k th agent stores only the columns corresponding to its local dataset. Thus, the block partitioning has the form $\hat{\mathbf{D}} = [\hat{\mathbf{D}}_1, \dots, \hat{\mathbf{D}}_L]$, where $\hat{\mathbf{D}}_k \in \mathbb{R}^{N \times N_k}$ is the block of the matrix held by the k th agent, and \mathcal{I}_k is the set of indexes of known entries of $\hat{\mathbf{D}}_k$. The same block partition applies also to matrices $\mathbf{B} = [\mathbf{B}_1, \dots, \mathbf{B}_L]$, with $\mathbf{B}_k \in \mathbb{R}^{r \times N_k}$, and $\tilde{\mathbf{D}} = [\tilde{\mathbf{D}}_1, \dots, \tilde{\mathbf{D}}_L]$, with $\tilde{\mathbf{D}}_k \in \mathbb{R}^{N \times N_k}$. The matrix \mathbf{A} cannot be partitioned, but each agent stores a local copy \mathbf{A}_k to use in computations. The *D-LMaFit* algorithm consists in an alternation of matrix factorizations and inexact average consensus, formalized in the following steps:

1. Initialization: For each agent, the matrices $\mathbf{A}_k [0]$ and $\mathbf{B}_k [0]$ are initialized as random matrices of appropriate dimensions. Matrix $\tilde{\mathbf{D}}_k [0]$ is initialized as $\tilde{\mathbf{D}}_k [0] = \hat{\mathbf{D}}_k$.
2. Update of \mathbf{A} : At time n , the k th agent updates its local copy of the matrix \mathbf{A} . If $n = 0$, the updating rule is:

$$\mathbf{A}_k [1] = \sum_{i=1}^L C_{ki} \mathbf{A}_i [0] - \alpha \left(\mathbf{A}_k [0] - \tilde{\mathbf{D}}_k [0] \mathbf{B}_k^{\text{T}} [0] \right), \quad (7.20)$$

where α is a suitable positive step-size. If $n > 0$, the updating rule is given by:

$$\begin{aligned} \mathbf{A}_k [n + 1] = \mathbf{A}_k [n] - \sum_{i=1}^L \left(C_{ki} \mathbf{A}_i [n] - \tilde{C}_{ki} \mathbf{A}_i [n - 1] \right) - \\ \alpha \left(\mathbf{A}_k [n] - \mathbf{A}_k [n - 1] - \tilde{\mathbf{D}}_k [n] \mathbf{B}_k^{\text{T}} [n] + \tilde{\mathbf{D}}_k [n - 1] \mathbf{B}_k^{\text{T}} [n - 1] \right). \end{aligned} \quad (7.21)$$

In Eq. (7.21), $\tilde{\mathbf{C}}$ is a mixing matrix that satisfies some properties [92]. A suitable choice is $\tilde{\mathbf{C}} = (1/2)(\mathbf{I} + \mathbf{C})$.

3. Update of \mathbf{B} and $\tilde{\mathbf{D}}$: At the n th iteration, agent k updates matrices \mathbf{B}_k and $\tilde{\mathbf{D}}_k$

according to:

$$\mathbf{B}_k[n+1] = \mathbf{A}_k^\dagger[n+1] \mathbf{A}_k^T[n+1] \tilde{\mathbf{D}}_k[n] \quad (7.22)$$

$$\begin{aligned} \tilde{\mathbf{D}}_k[n+1] &= \mathbf{A}_k[n+1] \mathbf{B}_k[n+1] + \\ &P_{\mathcal{I}_k} \left(\hat{\mathbf{D}}_k - \mathbf{A}_k[n+1] \mathbf{B}_k[n+1] \right) \end{aligned} \quad (7.23)$$

where $\mathbf{A}_k^\dagger[n+1]$ is the Moore-Penrose inverse of $\mathbf{A}_k[n+1]$, and $P_{\mathcal{I}}(\mathbf{M}) : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^{n \times m}$ is a projection operator defined by:

$$P_{\mathcal{I}}(\mathbf{M})_{ij} = \begin{cases} M_{ij} & \text{if } (i, j) \in \mathcal{I} \\ 0 & \text{otherwise} \end{cases}. \quad (7.24)$$

The algorithm stops when the maximum number of iterations n_{\max}^{EDM} is reached.

As we stated, *D-LMaFit* is not specifically designed for EDM completion. Consequently, it has some important limitations in our context. In particular, the resulting matrix $\tilde{\mathbf{D}}$ can have negative entries and could be non-symmetric; moreover, it is distributed across the nodes and so, if an agent wants access to the complete matrix, it has to collect the local matrices $\tilde{\mathbf{D}}$ through all the network. In order to at least satisfy the constraint that $\tilde{\mathbf{D}}$ be an appropriate EDM, we introduce the following modifications into the original algorithm:

- The updating equation for $\tilde{\mathbf{D}}_k$ is modified by setting to 0 all the negative entries. This projection operator is a standard approach in non-negative matrix factorization to enforce non-negativity constraints [93].
- When all the agents gathered the complete matrix $\tilde{\mathbf{D}}$, this is symmetrized as $\check{\mathbf{D}} = \frac{\tilde{\mathbf{D}} + \tilde{\mathbf{D}}^T}{2}$.

7.3.3 Diffusion gradient descent

The second algorithm for distributed EDM completion proposed in this chapter exploits the low-rank factorization $\mathbf{D} = \kappa(\mathbf{V}\mathbf{V}^T)$ showed in Section 7.2.2. In particular, we consider the general framework of DA (see Section 3.3). To begin with, we can observe that the objective function in Eq. (7.13) can be approximated locally by:

$$J_k(\mathbf{V}) = \left\| \boldsymbol{\Omega}_k \circ \left[\hat{\mathbf{D}}_k - \kappa(\mathbf{V}\mathbf{V}^T) \right] \right\|_{\text{F}}^2 \quad k = 1, \dots, L, \quad (7.25)$$

where $\boldsymbol{\Omega}_k$ is the local auxiliary matrix associated with $\hat{\mathbf{D}}_k$. Hence, we can exploit a DA algorithm to minimize the joint cost function given by $\check{J}(\mathbf{V}) = \sum_{k=1}^L J_k(\mathbf{V})$. The DGD for the distributed completion of an EDM is defined by an alternation of updating and diffusion equations in the form of:

1. Initialization: All the agents initialize the local matrices \mathbf{V}_k as random $N \times r$ matrices.
2. Update of \mathbf{V} : At time n , the k th agent updates the local matrix \mathbf{V}_k using a gradient descent step with respect to its local cost function:

$$\tilde{\mathbf{V}}_k[n+1] = \mathbf{V}_k[n] - \eta_k[n] \nabla_{\mathbf{V}_k} J_k(\mathbf{V}). \quad (7.26)$$

where $\eta_k[n]$ is a positive step-size. It is straightforward to show that the gradient of the cost function is given by:

$$\begin{aligned} \nabla_{\mathbf{V}_k} J_k(\mathbf{V}) = & \kappa^* \left\{ \mathbf{\Omega}_k \circ \right. \\ & \left. \circ \left(\kappa \left(\mathbf{V}_k[n] \mathbf{V}_k^T[n] \right) - \hat{\mathbf{D}}_k \right) \right\} \mathbf{V}_k[n], \end{aligned} \quad (7.27)$$

where $\kappa^*(\mathbf{A}) = 2 [\text{diag}(\mathbf{A}\mathbf{1}) - \mathbf{A}]$ is the adjoint operator of κ .

3. Diffusion: The updated matrices are combined according to the mixing weights \mathbf{C} :

$$\mathbf{V}_k[n+1] = \sum_{i=1}^L C_{ki} \tilde{\mathbf{V}}_i[n+1]. \quad (7.28)$$

Compared with the state-of-the-art decentralized block algorithm presented in the previous section, the diffusion-based approach has two main advantages. First, it is able to take into account naturally the properties of EDM matrices. Secondly, at every step each node has a complete estimate of the overall matrix, instead of a single column-wise block. Thus, there is no need of gathering the overall matrix at the end of the optimization process.

7.4 Distributed Semi-supervised Manifold Regularization

In this section, we consider the more general distributed SSL setting, as illustrated in Fig. 7.1. We suppose that the agents in the network have performed a distributed matrix completion step, using either the algorithm in Section 7.3.2 or the one in Section 7.3.3, so that the estimates $\tilde{\mathbf{D}}$, $\tilde{\mathbf{L}}$ and $\tilde{\mathbf{K}}$ are globally known. For the k th agent, we denote with $\hat{\mathbf{y}}_k$ the N_k dimensional vector with elements:

$$\hat{y}_{k,i} = \begin{cases} y_{k,i} & \text{if } i \in \{1, \dots, l_k\} \\ 0 & \text{if } i \in \{l_k + 1, \dots, l_k + u_k\} \end{cases}, \quad (7.29)$$

and $\hat{\mathbf{J}}_k$ the $N_k \times N$ matrix defined by $\hat{\mathbf{J}}_k = [\bar{\mathbf{0}}_k \quad \Lambda_k \quad \underline{\mathbf{0}}_k]$, where Λ_k is a $N_k \times N_k$ diagonal matrix with elements:

$$\Lambda_{k,ii} = \begin{cases} 1 & \text{if } i \in \{1, \dots, l_k\} \\ 0 & \text{if } i \in \{l_k + 1, \dots, l_k + u_k\} \end{cases}, \quad (7.30)$$

$\bar{\mathbf{0}}_k$ is a $N_k \times \sum_{j < k} N_j$ null matrix and $\underline{\mathbf{0}}_k$ is a $N_k \times \sum_{j > k} N_j$ null matrix. Using this notation, the optimization problem of LapKRR can be reformulated in distributed form as:

$$\min_{\boldsymbol{\alpha}} \sum_{k=1}^L \|\hat{\mathbf{y}}_k - \hat{\mathbf{J}}_k \tilde{\mathbf{K}} \boldsymbol{\alpha}\|_2^2 + \gamma_A \boldsymbol{\alpha}^T \tilde{\mathbf{K}} \boldsymbol{\alpha} + \gamma_I \boldsymbol{\alpha}^T \tilde{\mathbf{L}} \tilde{\mathbf{K}} \boldsymbol{\alpha}. \quad (7.31)$$

Denoting with $\hat{\mathbf{J}}_{\text{tot}} = \sum_{k=1}^L \hat{\mathbf{J}}_k^T \hat{\mathbf{J}}_k$ and $\hat{\mathbf{y}}_{\text{tot}} = \sum_{k=1}^L \hat{\mathbf{J}}_k^T \hat{\mathbf{y}}_k$, we can derive the expression for the optimal weights vector $\boldsymbol{\alpha}^*$:

$$\boldsymbol{\alpha}^* = (\hat{\mathbf{J}}_{\text{tot}} \tilde{\mathbf{K}} + \gamma_A \mathbf{I} + \gamma_I \tilde{\mathbf{L}} \tilde{\mathbf{K}})^{-1} \hat{\mathbf{y}}_{\text{tot}}. \quad (7.32)$$

The particular structure of $\boldsymbol{\alpha}^*$ implies that the distributed solution can be decomposed as $\boldsymbol{\alpha}^* = \sum_{k=1}^L \boldsymbol{\alpha}_k^*$, where:

$$\boldsymbol{\alpha}_k^* = (\hat{\mathbf{J}}_{\text{tot}} \tilde{\mathbf{K}} + \gamma_A \mathbf{I} + \gamma_I \tilde{\mathbf{L}} \tilde{\mathbf{K}})^{-1} \hat{\mathbf{J}}_k^T \hat{\mathbf{y}}_k. \quad (7.33)$$

To compute the local solution $\boldsymbol{\alpha}_k^*$, the k th agent requires only the knowledge of matrix $\hat{\mathbf{J}}_{\text{tot}}$, which can be computed with a distributed sum over the network using the DAC protocol. Clearly, the sum can be obtained by post-multiplying the final estimate by L . Overall, the distributed LapKRR algorithm can be summarized in five main steps:

1. Distributed Laplacian estimation: this step corresponds to the process illustrated in Sec. 7.3. It includes the patterns exchange (with the inclusion of a privacy-preserving strategy, if needed) and the points exchange procedures, the distributed EDM completion, and the computation of $\tilde{\mathbf{L}}$ and $\tilde{\mathbf{K}}$.
2. Global sum of $\hat{\mathbf{J}}_{\text{tot}}$: in this step the local matrices $\hat{\mathbf{J}}_k^T \hat{\mathbf{J}}_k$ are summed up using the DAC protocol.
3. Local training: using the matrix $\hat{\mathbf{J}}_{\text{tot}}$ computed in the previous step, each agent calculates its local solution, given by:

$$\boldsymbol{\alpha}_k^* = (\hat{\mathbf{J}}_{\text{tot}} \tilde{\mathbf{K}} + \gamma_A \mathbf{I} + \gamma_I \tilde{\mathbf{L}} \tilde{\mathbf{K}})^{-1} \hat{\mathbf{J}}_k^T \hat{\mathbf{y}}_k. \quad (7.34)$$

4. Global sum of $\boldsymbol{\alpha}^*$: in this step, using the DAC protocol, the local vectors $\boldsymbol{\alpha}_k^*$ are summed up to compute the global weight vector.

5. Output estimation: when a new unlabeled pattern \mathbf{x} is available to the network, each agent can initialize a partial output as:

$$f_k(\mathbf{x}) = \sum_{i=1}^{N_k} \mathcal{K}(\mathbf{x}, \mathbf{x}_{k,i}) \beta_{k,i}^*, \quad (7.35)$$

where β_k^* is a N_k -dimensional vector containing the entries of \mathbf{a}^* corresponding to the patterns belonging to the k th agent. The global output is then computed as:

$$f(\mathbf{x}) = \sum_{k=1}^L f_k(\mathbf{x}), \quad (7.36)$$

which can be obtained efficiently with the use of the DAC protocol.

A pseudocode of the algorithm, from the point of view of a single agent, is provided in Algorithm 7.1.

Algorithm 7.1 Distr-LapKRR: Pseudocode of the proposed distributed SSL algorithm (k th node).

Inputs: Labeled S_k and unlabeled U_k training data, number of nodes L (global), regularization parameters γ_A, γ_I (global)

Output: Optimal vector \mathbf{a}_k^*

- 1: **for** $n = 1$ to n_{\max}^1 **do**
 - 2: Select a set of input patterns and share them with the neighbors \mathcal{N}_k , using a privacy-preserving transformation if needed.
 - 3: Receive patterns from the neighbors.
 - 4: **end for**
 - 5: Compute the incomplete EDM matrix $\hat{\mathbf{D}}_k$.
 - 6: **for** $n = 1$ to n_{\max}^2 **do**
 - 7: Select a set of entries from $\hat{\mathbf{D}}_k$ and share them with the neighbors.
 - 8: Receive entries from the neighbors.
 - 9: Update $\hat{\mathbf{D}}_k$ with the entries received.
 - 10: **end for**
 - 11: Complete the matrix $\tilde{\mathbf{D}}$ using the algorithm presented in Sec. 7.3.2 or in Sec. 7.3.3.
 - 12: Compute the Laplacian matrix $\tilde{\mathbf{L}}$ and the kernel matrix $\tilde{\mathbf{K}}$ using $\tilde{\mathbf{D}}$.
 - 13: Compute the sum $\hat{\mathbf{J}}_{\text{tot}}$ over the network using the DAC protocol.
 - 14: **return** \mathbf{a}_k^* according to Eq. (7.33).
-

7.5 Experimental results

7.5.1 Experiments setup

We tested the performance of our proposed algorithm over five publicly available datasets. In order to get comparable results with state-of-the-art SSL algorithms, the datasets were chosen among a variety of benchmarks for SSL. A schematic overview of their characteristics is given in Tab. 7.1. For further information about the datasets, we refer to [11] for 2Moons, to [31] for BCI, and to [110] for the rest of the datasets. The COIL dataset is used in two different versions, one with 2 classes (COIL2) and a harder version with 20 classes (COIL20). In all the cases, input variables are normalized between -1 and 1 before the experiments.

Table 7.1: Description of the datasets used for testing Distr-LapKRR.

Name	Features	Size	N. Classes	TR	TST	U
2Moons	2	400	2	14	200	186
BCI	117	400	2	14	100	286
G50C	50	550	2	50	186	314
COIL20	1024	1440	20	40	400	1000
COIL2	1024	1440	2	40	400	1000

In our experimental setup we considered a 7-nodes network, whose topology is kept fixed for all the experiments. The topology is generated such that each pair of agents is connected with a probability c . In particular, in our implementation we set $c = 0.5$, while we choose the weights matrix \mathbf{C} using the ‘max-degree’ strategy. This choice ensures both convergence of the DAC protocol [198] and it satisfies the requirements of the DA framework [145]. All the experiments are repeated 25 times, to average possible outliers results due to the randomness in the processes of exchange and in the initialization of the matrices in the EDM completion algorithms. At every run, data are randomly shuffled and then partitioned in a labeled training set TR, a test set TST, and an unlabeled set U, whose cardinalities are reported in Tab. 7.1. Both the labeled and unlabeled training sets are then partitioned evenly across the nodes. All the experiments are performed using MATLAB R2014a on an Intel i7-3820 @3.6 GHz and 32 GB of memory.

7.5.2 Distributed Laplacian estimation

In this section we compare the performance of the two strategies for distributed EDM completion illustrated in Section 7.3. We analyze the matrix completion error, together with the overall computational time for the two strategies. Given an

estimate $\tilde{\mathbf{D}}$ of \mathbf{D} , we define the matrix completion error as:

$$E(\tilde{\mathbf{D}}) = \frac{\|\tilde{\mathbf{D}} - \mathbf{D}\|_F}{\|\mathbf{D}\|_F}. \quad (7.37)$$

The first set of experiments consists in comparing the completion error and the time required by the two algorithms, for different sizes of the sampling set of \mathbf{D} . In our context, the size of the sampling set depends only on the amount of data that are exchanged before the algorithm runs. To this end, we consider the completion error when varying the number of iterations for both the patterns exchange and the entries exchange steps, while keeping fixed the exchange fraction p . In particular, for all the datasets we varied the maximum number of iterations $n_{\max}^{(1)}$ and $n_{\max}^{(2)}$ from 0 to 150, by steps of 10. Results of this experiment are presented in Fig. 7.2. The solid red and the solid blue lines show the performance of Decentralized Block Estimation and DGD, respectively. Since the value of the completion error only depends on the input \mathbf{x} , the results for datasets COIL20 and COIL2 are reported together. The values for the patterns exchange fraction p_1 and the entries exchange

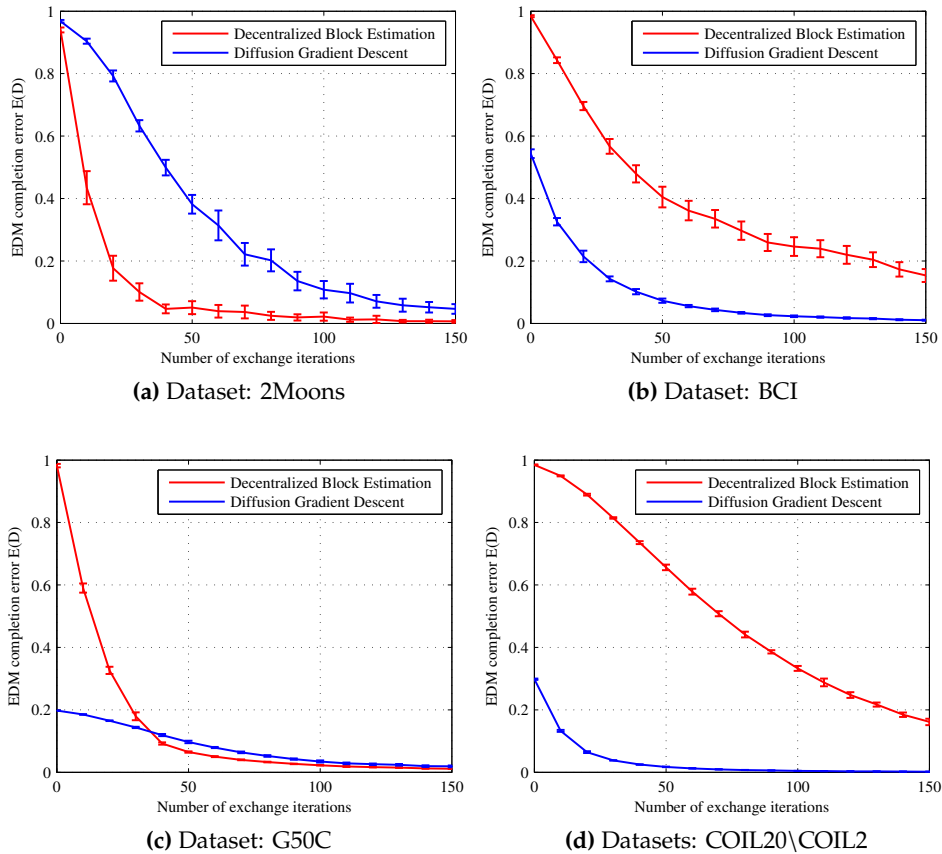


Figure 7.2: Average EDM completion error of the two strategies on the considered datasets, when varying the number of iterations for the patterns and entries exchange protocols. The vertical bars represent the standard deviation from the average.

Table 7.2: Values for the parameters used in the simulations. The values in the first group are used in the distributed protocols and in the DGD algorithm (p1 and p2 are in percentages). Those in the second group are used to build the Laplacian and kernel matrices. In the third group are reported the parameters used in the privacy-preserving transformations.

Dataset	p1	$n_{\max}^{(1)}$	p2	$n_{\max}^{(2)}$	η	γ_A	γ_I	nn	σ_K	q	t	σ_a	σ_b	σ_Q	σ_C
2Moons	3.5	100	3.5	100	10^{-3}	2^{-5}	4	6	0.03	1	–	–	–	–	–
BCI	2.5	100	2.5	100	10^{-6}	10^{-6}	1	5	1	2	10^4	0	0	1	10^{-6}
G50C	2	150	2.5	150	10^{-6}	10^{-6}	10^{-2}	50	17.5	5	$2e^4$	0	1	1	$1.1e^{-6}$
COIL	2	150	2.5	150	10^{-7}	10^{-6}	1	2	0.6	1	10^3	0	0	1	10^{-6}

fraction p2 are chosen to balance the communication overhead and the size of the sampling set of \mathbf{D} . For both the algorithms, we set the maximum number of iterations n_{\max}^{EDM} to 1500, and we used a fixed step-size strategy. In particular for the Decentralized Block Estimation we set $\alpha = 0.4$, as suggested in [92], while for the Diffused Gradient Descent, the optimal values for η are chosen singularly for each dataset by searching in the interval 10^j , $j \in \{-10, \dots, -3\}$. These parameters, together with the values for p1 and p2, are reported in Tab. 7.2, and are used in all the experiments.

We see that, with the solely exception of the 2Moons dataset (see Fig. 7.2a), the novel Diffused Gradient Descent algorithm achieves better performance when compared to the Decentralized Block Estimation, in particular when few information is exchanged before the completion process. For all the datasets, as the number of the exchange iterations increases, the diffusion strategy is able to converge rapidly to the real EDM \mathbf{D} , while the performance is poorer for the block partitioning strategy, resulting for datasets BCI and COIL in a completion error of 19% even for high quantity of information exchanged (see Fig. 7.2b and Fig. 7.2d).

When considering the time required by the two algorithms, which is shown in Fig. 7.3, we observe that the block partition strategy requires for datasets 2Moons and G50C less than half the time required by the diffused strategy, while, as the number of the features increases, the diffusion strategy tends to be less computational expensive. In fact, the time required by both strategies is nearly the same for the dataset BCI, while for COIL the diffusion strategy is 1.2 times faster. We remark that the Decentralized Block Estimation requires an additional step for all the agents to gather the columns-wise blocks through the network, which has not been taken into account in calculating the computational time.

7.5.3 Distributed semi-supervised manifold regularization

The second experiment analyzes the performance of the distributed algorithm when compared to a centralized learning strategy and to a local learning strategy.

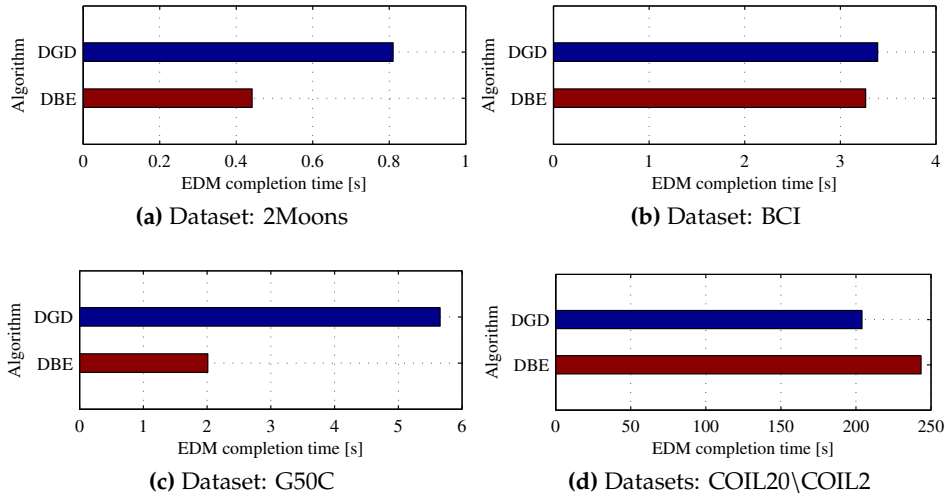


Figure 7.3: Average EDM completion time required by the two strategies on the considered datasets. DGD and DBE are the acronyms for Decentralized Block Estimation and DGD respectively.

We compare the following algorithms:

- **Centr-LapKRR**: this is the algorithm depicted in Sec. 7.2.1. It is equivalent to a single agent collecting all the training data.
- **Local-LapKRR**: in the local setting, the training set is distributed across the agents and every agent trains a LapKRR on its own dataset, without any communication with other agents. The error is averaged throughout the nodes.
- **Distr-LapKRR**: as before, the training set is distributed within the network, but the agents converge to a centralized solution using the strategy detailed in Sec. 7.4. In this experiment, the EDM completion is achieved by the DGD algorithm.

For all the algorithms, we build the Laplacian and the kernel matrices according to the method detailed in [110], using the parameters reported in Tab. 7.2. In particular the parameters for datasets G50C and COIL come from [110], while those for 2Moons and BCI come from [11] and [31], respectively. Lower values for the exchange iterations in datasets 2Moons and BCI are chosen to balance the higher values for the exchange fractions. The classification error and the computational time for the three models over the five datasets are reported in Table 7.3. Results of the proposed algorithm, Distr-LapKRR, are highlighted in bold.

We can see that Distr-LapKRR is generally able to match the same performance of the Centr-LapKRR, both in mean and variance, except for a small decrease in the G50C dataset. Clearly, the performance of Local-LapKRR is noticeably worse than

Table 7.3: Average values for classification error and computational time, together with standard deviation, for the three algorithms. Results for the proposed algorithm are highlighted in bold.

Dataset	Algorithm	Error [%]	Time [s]
2Moons	Centr-LapKRR	0.005 ± 0.001	0.006 ± 0.015
	Distr-LapKRR	0.01 ± 0.03	0.875 ± 0.030
	Local-LapKRR	0.41 ± 0.28	0.000 ± 0.000
BCI	Centr-LapKRR	0.49 ± 0.04	0.021 ± 0.012
	Distr-LapKRR	0.49 ± 0.05	3.396 ± 0.028
	Local-LapKRR	0.54 ± 0.14	0.001 ± 0.000
G50C	Centr-LapKRR	0.07 ± 0.02	0.101 ± 0.017
	Distr-LapKRR	0.12 ± 0.10	5.764 ± 0.066
	Local-LapKRR	0.45 ± 0.06	0.001 ± 0.000
COIL20	Centr-LapKRR	0.13 ± 0.02	1.565 ± 0.019
	Distr-LapKRR	0.13 ± 0.02	195.933 ± 2.176
	Local-LapKRR	0.78 ± 0.07	0.056 ± 0.001
COIL2	Centr-LapKRR	0.10 ± 0.03	1.556 ± 0.028
	Distr-LapKRR	0.10 ± 0.03	191.478 ± 0.864
	Local-LapKRR	0.43 ± 0.12	0.055 ± 0.000

the other two algorithms, because the local models are built on considerably smaller training sets. The computational time required by the distributed algorithm is given by the sum of the time required by both the exchange protocols, the distributed Laplacian estimation, the DAC protocol, and the matrix inversion in (7.33). When comparing the results with the values for EDM completion time obtained in the previous experiment, we notice that the order of magnitude of the time required by Distr-LapKRR is given by the time necessary to complete the distance matrix.

7.5.4 Privacy preservation

As a final experiment, we include in our algorithm the two privacy-preserving strategies presented in Sec. 7.2.3. In particular, we analyze the evolution of the classification error when varying the ratio m/d from 0.1 to 0.95, i.e. when varying the dimensionality m of the transformed patterns. In this experiment we do not consider the 2Moons dataset, because of its limited number of features. Since the value of σ in the linear random projection has no influence on the error of the transformed patterns, we set $\sigma = 1$ for all the datasets. As for the nonlinear transformation, the values for the parameters are searched inside a grid and then optimized locally. Possible values for t are searched in 10^i , $i = \{1, \dots, 5\}$, while values for the variances are searched in 10^j , $j = \{-6, \dots, 6\}$. The optimal values for

the datasets are reported in the third group of Table 7.2.

Results of the experiment are presented in Fig. 7.4. The classification error for the linear random projection and nonlinear transformation are shown with solid red and dashed blue lines, respectively. In addition, the mean value for Distr-LapKRR (together with its confidence interval) is reported as a baseline, shown with a dashed black line.

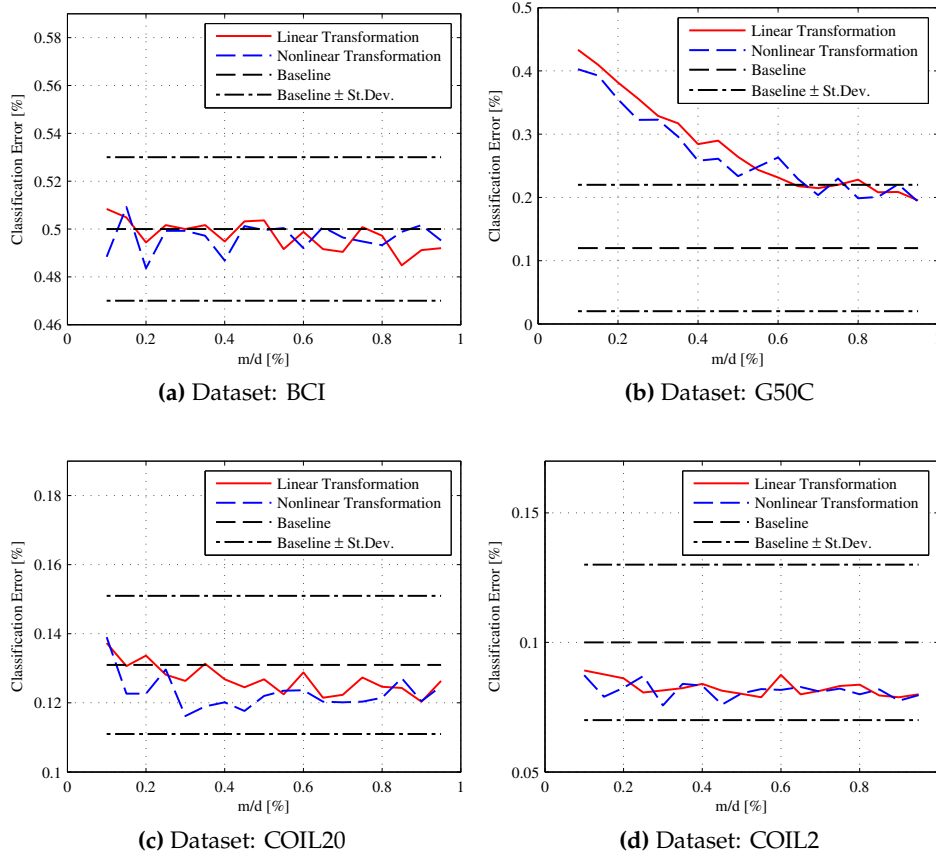


Figure 7.4: Average classification error of the privacy-preserving transformations on the considered datasets when varying the ratio m/d .

By observing the results, we can see that when compared to Distr-LapKRR, the privacy-preserving strategies show different behaviors depending on the dataset. In particular, for dataset BCI, the error is nearly the same of Distr-LapKRR, while it is slightly lower for COIL2 and COIL20, and somewhat higher for G50C, where it shows a decreasing trend. For all the datasets, we see that the error achieved using the privacy-preserving strategies remains inside the limits of Distr-LapKRR error’s confidence interval, denoting how the variability introduced does not have significant influence on the algorithm’s performance.

We notice that in most cases, we can obtain a comparable or even better performance with respect to the privacy-free algorithm, with significantly fewer

features, leading to a reduction of the information exchanged and therefore of the overall computational requirements. For all the datasets, both the transformations present a non-smooth trend, caused by the heuristic nature of these methods. Moreover, the error is very similar between the strategies, suggesting that the use of a nonlinear transformation, potentially safer than a linear one, does not influence the performance.

Distributed Semi-Supervised Support Vector Machines

Contents

8.1	Introduction	91
8.2	Semi-Supervised Support Vector Machines	93
8.3	Distributed learning for S^3VM	94
8.3.1	Formulation of the problem	95
8.3.2	Solution 1: Distributed gradient descent	95
8.3.3	Solution 2: In-network successive convex approximation	97
8.4	Experimental Results	99
8.4.1	Experimental Setup	99
8.4.2	Results and discussion	101

8.1 Introduction

IN the previous chapter, we have explored the problem of training a semi-supervised Laplacian KRR using a distributed computation of the underlying kernel matrix. However, despite its good performance, the resulting algorithm requires a large amount of computational and/or communication resources, which might not be available on specific devices or communication channels. To this end, in this chapter we propose two simpler algorithms for a different family of semi-supervised SVM, denoted as S^3VM . The S^3VM has attracted a large amount of attention over the last decades [30]. It is based on the idea of minimizing the training error and maximizing the margin over both labeled and unlabeled data, whose labels are included as additional variables in the optimization problem. Since its first practical implementation in [83], numerous researchers have proposed

The content of this chapter is adapted from the material published in [153].

alternative solutions for solving the resulting mixed integer optimization problem, including branch and bound algorithms [32], convex relaxations, convex-concave procedures [30], and others. It has been applied in a wide variety of practical problems, such as text inference [83], and it has given birth to numerous other algorithms, including semi-supervised least-square SVMs [1], and semi-supervised random vector functional-link networks [152].

In order to simplify our derivation, in this chapter we focus on the *linear* S³VM formulation, whose decision boundary corresponds to an hyperplane in the input space. Due to this, the algorithms presented in this chapter can be implemented even on agents with stringent requirements in terms of power, such as sensors in a WSN. At the same time, it is known that limiting ourselves to a linear decision boundary can be reasonable, as the linear S³VM can perform well in a large range of settings, due to the scarcity of labeled data [30].

Specifically, starting from the smooth approximation to the original S³VM presented in [33], we show that the distributed training problem can be formulated as the joint minimization of a sum of non-convex cost functions. This is a complex problem, which has been investigated only very recently in the distributed optimization literature [16, 46]. In our case, we build on two different solutions. The first one is based on the idea of diffusion gradient descent (DGD), similarly to the previous chapter. Nevertheless, since it is a gradient-based algorithm exploiting only first order information of the objective function, it generally suffers of slow practical convergence speed, especially in the case of non-convex and large-scale optimization problems. Recently, it was showed in [46, 162] that exploiting the structure of nonconvex functions by replacing their linearization (i.e., their gradient) with a “better” approximant can enhance practical convergence speed. Thus, we propose a distributed algorithm based on the recently proposed In-Network Successive Convex Approximation (NEXT) framework [46]. The method hinges on successive convex approximation techniques while leveraging dynamic consensus as a mechanism to distribute the computation among the agents as well as diffuse the needed information over the network. Both algorithms are proved convergent to a stationary point of the optimization problem. Moreover, as shown in our experimental results, the NEXT exhibits a faster practical convergence speed with respect to DGD, which is paid by a larger computation cost per iteration.

The rest of the chapter is structured as follows. In Section 8.2 we introduce the S³VM model together with the approximation presented in [33]. In Section 8.3, we first formulate the distributed training problem for S³VMs, and subsequently we derive our two proposed solutions. Finally, Section 8.4 details an extensive set of experimental results.

8.2 Semi-Supervised Support Vector Machines

Let us consider the standard SSL problem, where we are interested in learning a binary classifier starting from L labeled samples $(\mathbf{x}_i, y_i)_{i=1}^L$ and U unlabeled samples $(\mathbf{x}_i)_{i=1}^U$. As before, each input is a d -dimensional real vector $\mathbf{x}_i \in \mathbb{R}^d$, while each output can only take one of two possible values $y_i \in \{-1, +1\}$. The linear S³VM optimization problem can be formulated as [33]:

$$\min_{\mathbf{w}, b, \hat{\mathbf{y}}} \frac{C_1}{2L} \sum_{i=1}^L l(y_i, f(\mathbf{x}_i)) + \frac{C_2}{2U} \sum_{i=1}^U l(\hat{y}_i, f(\mathbf{x}_i)) + \frac{1}{2} \|\mathbf{w}\|_2^2, \quad (8.1)$$

where $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, $\hat{\mathbf{y}} \in \{-1, +1\}^U$ is a vector of unknown labels, $l(\cdot, \cdot)$ is a proper loss function and $C_1, C_2 > 0$ are coefficients weighting the relative importance of labeled and unlabeled samples. The main difference with respect to the standard SVM formulation is the inclusion of the unknown labels $\hat{\mathbf{y}}$ as variables of the optimization problem. This makes Problem (8.1) a mixed integer optimization problem, whose exact solution can be computed only for relatively small datasets, e.g. using standard branch-and-bound algorithms. We note that, for $C_2 = 0$, we recover the standard SVM formulation. The most common choice for the loss function is the hinge loss, given by:

$$l(y, f(\mathbf{x})) = \max(0, 1 - y f(\mathbf{x}))^p, \quad (8.2)$$

where $p \in \mathbb{N}$. In this chapter, we use the choice $p = 2$, which leads to a smooth and convex function. Additionally, it is standard practice to introduce an additional constraint in the optimization problem, so that the resulting vector $\hat{\mathbf{y}}$ has a fixed proportion $r \in [0, 1]$ of positive labels:

$$\frac{1}{U} \sum_{i=1}^U \max(0, \hat{y}_i) = r. \quad (8.3)$$

This constraint helps achieve a balanced solution, especially when the ratio r reflects the true proportion of positive labels in the underlying dataset.

A common way of solving Problem (8.1) stems from the fact that, for a fixed \mathbf{w} and b , the optimal $\hat{\mathbf{y}}$ is given in closed form by

$$\hat{y}_i = \text{sign}(\mathbf{w}^T \mathbf{x}_i + b), \quad i = 1, \dots, U.$$

Exploiting this fact, it is possible to devise a *continuous* approximation of the cost function in (8.1) [30]. In particular, to obtain a smooth optimization problem solvable by standard first-order methods, [33] propose to replace the hinge loss over the unknown labels with the approximation given by $\exp\{-s f(\mathbf{x})^2\}$, $s > 0$. In

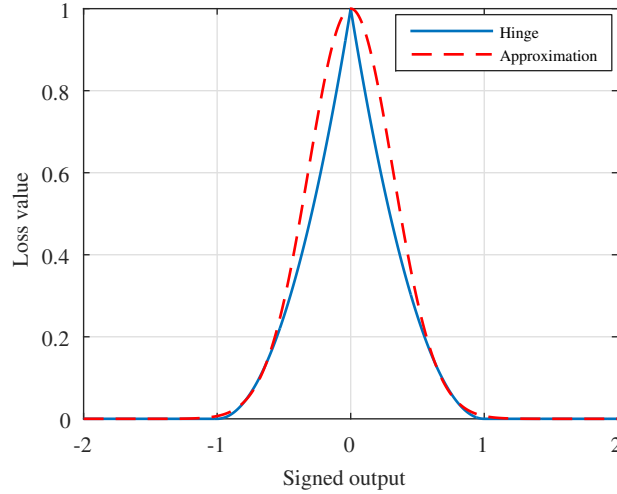


Figure 8.1: For a fixed choice of \mathbf{w} and b , $\max(0, 1 - \hat{y}_i f(\mathbf{x}_i))^2 = \max(0, 1 - |f(\mathbf{x}_i)|)^2$. This is shown in blue for varying values of $f(\mathbf{x}_i)$, while in dashed red we show the approximation given by $\exp\{-5f(\mathbf{x}_i)^2\}$.

the following, we choose in particular $s = 5$, as suggested by [30]. A visual example of the approximation is illustrated in Fig. 8.1. The resulting ∇S^3VM optimization problem writes as:

$$\min_{\mathbf{w}, b} \frac{C_1}{2L} \sum_{i=1}^L l(y_i, f(\mathbf{x}_i)) + \frac{C_2}{2U} \sum_{i=1}^U \exp\{-s f(\mathbf{x}_i)^2\} + \frac{1}{2} \|\mathbf{w}\|_2^2. \quad (8.4)$$

Problem (8.4) does not incorporate the constraint in (8.3) yet. A possible way to handle the balancing constraint in (8.3) is a relaxation that uses the following linear approximation [33]:

$$\frac{1}{U} \sum_{i=1}^U \mathbf{w}^T \mathbf{x}_i + b = 2r - 1, \quad (8.5)$$

which can easily be enforced for a fixed r by first translating the unlabeled points so that their mean is $\mathbf{0}$, and then fixing the offset b as $b = 2r - 1$. The resulting problem can then be solved using standard first-order procedures.

8.3 Distributed learning for S^3VM

In this section, we first formulate a distributed optimization problem for a ∇S^3VM over a network of agents in Section 8.3.1. Then, we present two alternative methods for solving the overall optimization problem in a fully decentralized fashion in Sections 8.3.2 and 8.3.3.

8.3.1 Formulation of the problem

For the rest of this chapter, we assume that labeled and unlabeled training samples are not available on a single processor. Instead, they are distributed over a network of N agents. In particular, as in the previous chapter, we assume that the k th node has access to L_k labeled samples, and U_k unlabeled ones, such that $\sum_{k=1}^N L_k = L$ and $\sum_{k=1}^N U_k = U$. Assumptions on the topology are similar to the rest of the thesis. The distributed ∇S^3VM problem can be cast as:

$$\min_{\mathbf{w}} \sum_{k=1}^N l_k(\mathbf{w}) + \sum_{k=1}^N g_k(\mathbf{w}) + r(\mathbf{w}), \quad (8.6)$$

where we have defined the following shorthands:

$$l_k(\mathbf{w}) = \frac{C_1}{2L} \sum_{i=1}^{L_k} l(y_{k,i}, f(\mathbf{x}_{k,i})), \quad (8.7)$$

$$g_k(\mathbf{w}) = \frac{C_2}{2U} \sum_{i=1}^{U_k} \exp\{-s f(\mathbf{x}_{k,i})^2\}, \quad (8.8)$$

$$r(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2. \quad (8.9)$$

In the previous equations, we use the double subscript (k, i) to denote the i th sample available at the k th node, and we assume that the bias b has been fixed *a-priori* using the strategy detailed in the previous section. In a distributed setting, this requires that each agent knows the mean of all unlabeled points given by $\frac{1}{U} \sum_{i=1}^U \mathbf{x}_i$. This can easily be achieved, before starting the training process, with a number of different in-network algorithms. For example, the agents can compute the average using a DAC procedure, push-sum protocols [68] in a P2P network, or a number of alternative techniques.

8.3.2 Solution 1: Distributed gradient descent

The first solution is based on the DGD procedure, which has already been used extensively in the previous chapter for the distributed EDM completion problem. The main problem is that all the previous art on DGD focused on the solution of convex versions of problem the DSO problem. In our case, the $g_k(\mathbf{w})$ are non-convex, and the analysis in the aforementioned papers cannot be used. However, convergence of a similar family of algorithms in the case of non-convex (smooth) cost functions has been recently studied in [16]. Customizing the DGD method in (3.3) to Problem (8.6), we obtain the following local update at each agent:

$$\psi_k = \mathbf{w}_k[n] - \alpha_k[n] \left(\nabla l_k(\mathbf{w}_k[n]) + \nabla g_k(\mathbf{w}_k[n]) + \frac{1}{N} \nabla r(\mathbf{w}_k[n]) \right). \quad (8.10)$$

Note that we have included a factor $\frac{1}{N}$ in (8.10) in order to be consistent with the formulation in (3.2). Defining the margin $m_{k,i} = y_{k,i}f(\mathbf{x}_{k,i})$, we can easily show that:

$$\nabla l_k(\mathbf{w}) = -\frac{C_1}{L} \sum_{i=1}^{L_k} \mathbb{I}(1 - m_{k,i}) \cdot m_{k,i} (1 - m_{k,i}) , \quad (8.11)$$

$$\nabla g_k(\mathbf{w}) = -s \frac{C_2}{U} \sum_{i=1}^{U_k} \exp \left\{ -s f(\mathbf{x}_{k,i})^2 \right\} f(\mathbf{x}_{k,i}) \mathbf{x}_{k,i} , \quad (8.12)$$

$$\nabla r(\mathbf{w}) = \mathbf{w} , \quad (8.13)$$

where $\mathbb{I}(\cdot)$ is the indicator function defined for a generic scalar $o \in \mathbb{R}$ as:

$$\mathbb{I}(o) = \begin{cases} 1 & \text{if } o \leq 0 \\ 0 & \text{otherwise} \end{cases} .$$

The overall algorithm is summarized in Algorithm 8.1. Its convergence properties are illustrated in following theorem.

Algorithm 8.1 Distributed VS³VM using a distributed gradient descent procedure.

Inputs: Regularization factors C_1, C_2 , maximum number of iterations T .

- 1: **Initialization:**
 - 2: $\mathbf{w}_k[0] = \mathbf{0}, k = 1, \dots, N$.
 - 3: **for** n from 0 to T **do**
 - 4: **for** k from 1 to N **do in parallel**
 - 5: Compute auxiliary variable ψ_k using (8.10).
 - 6: Combine estimates as $\mathbf{w}_k[n+1] = \sum_{t=1}^N C_{kt} \psi_t$.
 - 7: **end for**
 - 8: **end for**
-

Theorem 2

Let $\{\mathbf{w}_k[n]\}_{k=1}^N$ be the sequence generated by Algorithm 1, and let $\bar{\mathbf{w}}[n] = \frac{1}{N} \sum_{k=1}^N \mathbf{w}_k[n]$ be its average across the agents. Let us select the step-size sequence $\{\alpha[n]\}_n$ such that i) $\alpha[n] \in (0, 1]$, for all n , ii) $\sum_{n=0}^{\infty} \alpha[n] = \infty$; and iii) $\sum_{n=0}^{\infty} \alpha[n]^2 < \infty$. Then, if the sequence $\{\bar{\mathbf{w}}[n]\}_n$ is bounded ¹, (a) [convergence]: all its limit points are stationary solutions of problem (8.6); (b) [consensus]: all the sequences $\mathbf{w}_k[n]$ asymptotically agree, i.e. $\lim_{n \rightarrow +\infty} \|\mathbf{w}_k[n] - \bar{\mathbf{w}}[n]\|_2 = 0, k = 1, \dots, N$.

Proof 3

See [16].

□

8.3.3 Solution 2: In-network successive convex approximation

The DGD algorithm is extremely efficient to implement, however, as we discussed in the introduction, its convergence is often sub-optimal due to two main reasons. First, the update in (8.10) considers only first order information and does not take into account the fact that the local cost function has some hidden convexity (since it is composed by the sum of a convex term plus a non-convex term) that one can properly exploit. Second, each agent k obtains information on the cost functions $J_t(\cdot)$, $t \neq k$, only in a very indirect way through the averaging step. In this section, we use a recent framework for in-network non-convex optimization from [46], which exploits the structure of nonconvex functions by replacing their linearization (i.e., their gradient) with a “better” approximant, thus typically resulting in enhanced practical convergence speed. In this section we customize the NEXT algorithm from [46] to our case, and we refer to the original chapter for more details.

The main idea of NEXT is to parallelize the problem in (8.6) such that, at each agent, the original (global) non-convex cost function is replaced with a strongly convex surrogate that preserves the first order conditions, see [46]. To this aim, we associate to agent k the surrogate $F_k(\mathbf{w}; \mathbf{w}_k[n])$, which is obtained by: i) keeping unaltered the local convex function $l_k(\mathbf{w})$ and the regularization function $r(\mathbf{w})$; ii) linearizing the local non-convex cost $g_k(\mathbf{w})$ and all the other (non-convex and unknown) terms $f_l(\mathbf{w})$ and $g_l(\mathbf{w})$, $l \neq k$, around the current local iterate $\mathbf{w}_k[n]$. As a result, the surrogate at node k takes the form:

$$F_k(\mathbf{w}; \mathbf{w}_k[n]) = l_k(\mathbf{w}) + \tilde{g}_k(\mathbf{w}; \mathbf{w}_k[n]) + r(\mathbf{w}) + \boldsymbol{\pi}_k(\mathbf{w}_k[n])^T (\mathbf{w} - \mathbf{w}_k[n]), \quad (8.14)$$

where

$$\tilde{g}_k(\mathbf{w}; \mathbf{w}_k[n]) = g_k(\mathbf{w}_k[n]) + \nabla g_k^T(\mathbf{w}_k[n]) (\mathbf{w} - \mathbf{w}_k[n]), \quad (8.15)$$

and $\boldsymbol{\pi}_k(\mathbf{w}_k[n])$ is defined as:

$$\boldsymbol{\pi}_k(\mathbf{w}_k[n]) = \sum_{t \neq k} \nabla h_t(\mathbf{w}_k[n]), \quad (8.16)$$

with $\nabla h_k(\cdot) = \nabla l_k(\cdot) + \nabla g_k(\cdot)$. Clearly, the information in (8.16) related to the knowledge of the other cost functions is not available at node k . To cope with this

issue, the NEXT approach consists in replacing $\pi_k(\mathbf{w}_k[n])$ in (8.14) with a local estimate $\tilde{\pi}_k[n]$ that asymptotically converges to $\pi_k(\mathbf{w}_k[n])$, thus considering the local approximated surrogate $\tilde{F}(\mathbf{w}; \mathbf{w}_k[n], \tilde{\pi}_k[n])$ given by:

$$\begin{aligned} \tilde{F}_k(\mathbf{w}; \mathbf{w}_k[n], \tilde{\pi}_k[n]) &= l_k(\mathbf{w}) + \tilde{g}_k(\mathbf{w}; \mathbf{w}_k[n]) + r(\mathbf{w}) \\ &+ \tilde{\pi}_k[n]^T (\mathbf{w} - \mathbf{w}_k[n]) . \end{aligned} \quad (8.17)$$

In the first phase of the algorithm, each agent solves a convex optimization problem involving the surrogate function in (8.17), thus obtaining a new estimate $\tilde{\mathbf{w}}_k[n]$. Then, an auxiliary variable $\mathbf{z}_k[n]$ is computed as a convex combination of the current estimate $\mathbf{w}_k[n]$ and the new $\tilde{\mathbf{w}}_k[n]$, as:

$$\mathbf{z}_k[n] = \mathbf{w}_k[n] + \alpha[n] (\tilde{\mathbf{w}}_k[n] - \mathbf{w}_k[n]) . \quad (8.18)$$

where $\alpha[n]$ is a possibly time-varying step-size sequence. This concludes the optimization phase of NEXT. The consensus phase of NEXT consists of two main steps. First, to achieve asymptotic agreement among the estimates at different nodes, each agent updates its local estimate combining the auxiliary variables from the neighborhood, i.e., for all k ,

$$\mathbf{w}_k[n+1] = \sum_{t=1}^N C_{kt} \mathbf{z}_t[n] . \quad (8.19)$$

This is similar to the diffusion step of the DGD procedure. Second, the update of the local estimate $\tilde{\pi}_k[n]$ in (8.17) is computed in two steps: i) an auxiliary variable $\mathbf{v}_k[n]$ is updated through a dynamic consensus step as:

$$\mathbf{v}_k[n+1] = \sum_{t=1}^N C_{kt} \mathbf{v}_t[n] + \left(\nabla h_k(\mathbf{w}_k[n+1]) - \nabla h_k(\mathbf{w}_k[n]) \right) ; \quad (8.20)$$

ii) the variable $\tilde{\pi}_k[n]$ is updated as:

$$\tilde{\pi}_k[n+1] = N \mathbf{v}_k[n+1] - \nabla h_k(\mathbf{w}_k[n+1]) . \quad (8.21)$$

The steps of the NEXT algorithm for Problem (8.6) are described in Algorithm 8.2. Its convergence properties are described by a Theorem completely similar to Theorem 1, and the details on the proof can be found in [46].

Algorithm 8.2 Distributed VS³VM using the In-Network Convex Optimization framework.

Inputs: Regularization factors C_1, C_2 , maximum number of iterations T .

1: **Initialization:**

2: $\mathbf{w}_k[0] = \mathbf{0}$, $k = 1, \dots, N$.

3: $\mathbf{v}_k[0] = \nabla h_k(\mathbf{w}_k[0])$, $k = 1, \dots, N$.

4: $\tilde{\pi}_k[0] = (N - 1)\mathbf{v}_k[0]$, $k = 1, \dots, N$.

5: **for** n from 0 to T **do**

6: **for** k from 1 to N **do in parallel**

7: Solve the local optimization problem:

$$\tilde{\mathbf{w}}_k[n] = \arg \min \tilde{F}_k(\mathbf{w}; \mathbf{w}_k[n], \tilde{\pi}_k[n]).$$

8: Compute $\mathbf{z}_k[n]$ using (8.18).

9: **end for**

10: **for** k from 1 to N **do in parallel**

11: Perform consensus step in (8.19).

12: Update auxiliary variable using (8.20).

13: Set $\tilde{\pi}_k[n + 1]$ as (8.21)..

14: **end for**

15: **end for**

8.4 Experimental Results

8.4.1 Experimental Setup

We tested the proposed distributed algorithms on three semi-supervised learning benchmarks, whose overview is given in Tab. 8.1. For more details on the datasets see [31] and the previous chapter for the first two, and [111] and Chapter 5 for GARAGEBAND. For this one, the original dataset comprises 9 different musical genres. In order to obtain a binary classification task, we select the two most prominent ones, namely ‘rock’ and ‘pop’, and discard the rest of the dataset. For G50C and GARAGEBAND, input variables are normalized between -1 and 1 . The experimental results are computed over a 10-fold cross-validation, and all the experiments are repeated 15 times. For each repetition, the training folds are partitioned in one labeled and one unlabeled datasets, according to the proportions given in Tab. 8.1. Results are then averaged over the 150 repetitions.

We compare the following models:

- **LIN-SVM:** this is a fully supervised SVM with a linear kernel, trained only on the labeled data. The model is trained using the LIBSVM library [29].
- **RBF-SVM:** similar to before, but a RBF kernel is used instead. The parameter

Table 8.1: Description of the datasets. The fourth and fifth columns denote the size of the training and unlabeled datasets, respectively.

Name	Features	Instances	L	U	Ref.
G50C	50	550	40	455	[31]
PCMAC	7511	1940	40	1700	[31]
GARAGEBAND	44	790	40	670	[111]

for the kernel is set according to the internal heuristic of LIBSVM.

- **C-VS3VM**: this is a centralized ∇S^3VM trained on both the labeled and the unlabeled data using a gradient descent procedure.
- **DG-VS3VM**: in this case, training data (both labeled and unlabeled) is distributed evenly across the network, and the distributed model is trained using the diffusion gradient algorithm detailed in Section 8.3.2.
- **NEXT-VS3VM**: data is distributed over the network as before, but the model is trained through the use of the NEXT framework, as detailed in Section 8.3.3. The internal optimization problem in (8.17) is solved using a standard gradient descent procedure.

For C-VS3VM, DG-VS3VM and NEXT-VS3VM we set $s = 5$ and a maximum number of iterations $T = 500$. In order to obtain a fair comparison between the algorithms, we also introduce a stopping criterion, i.e. the algorithms terminate when the norm of the gradient of the global cost function in (8.6) at the current iteration is less than 10^{-5} . Clearly, this is only for comparison purposes, and a truly distributed implementation would require a more sophisticated mechanism, which however goes outside the scope of the present chapter. The same value for the threshold is set for the gradient descent algorithm used within the NEXT framework to optimize the local surrogate function in (8.17). In this case, we let the gradient run for a maximum of $T = 50$ iterations. We note that, in general, we do not need to solve the internal optimization problem to optimal accuracy, as convergence of NEXT is guaranteed as long as the the problems are solved with increasing accuracy for every iteration [46].

We searched the values of C_1 and C_2 by executing a 5-fold cross-validation in the interval $\{10^{-5}, 10^{-4}, \dots, 10^3\}$ using C-VS3VM as in [33]. The values of these parameters are then shared with DG-VS3VM and NEXT-VS3VM. For all the models, included NEXT's internal gradient descent algorithm, the step-size α is chosen using a decreasing strategy given by:

$$\alpha[n] = \frac{\alpha_0}{(n+1)^\delta}, \quad (8.22)$$

Table 8.2: Optimal values of the parameters used in the experiments. In the first group are reported the values of the regularization coefficients for the three models, averaged over the 150 repetitions. In the following groups are reported the values of the initial step-size and of the diminishing factor for C-VS3VM, DG-VS3VM and NEXT-VS3VM respectively.

Dataset	C_1	C_2	α_0^C	δ^C	α_0^{DG}	δ^{DG}	α_0^{NEXT}	δ^{NEXT}
G50C	1	1	1	0.55	1	0.55	0.6	0.8
PCMAC	100	100	0.1	0.55	1	0.9	0.5	0.8
GARAGEBAND	2	5	0.09	0.8	0.1	0.1	0.05	0.55

where $\alpha_0, \delta > 0$ are set by the user. In particular, this strategy satisfies the convergence conditions for both the DGD algorithm and NEXT. After preliminary tests, we selected for every model the values of α_0 and δ that guarantee the fastest convergence. The optimal values of the parameters are shown in Tab. 8.2.

The network topologies are generated according to the ‘Erdős-Rényi model’, such that every edge has a 25% probability of appearing. The only constraint is that the network is connected. The topologies are generated at the beginning of the experiments and kept fixed during all the repetitions. We choose the weight matrix C using the *Metropolis-Hastings* strategy as in previous chapters. This choice of the weight matrix satisfies the convergence conditions for both the distributed approaches.

8.4.2 Results and discussion

The first set of experiments consists in analyzing the performance of C-VS3VM, when compared to a linear SVM and RBF SVM trained only on the labeled data. While these results are well known in the semi-supervised literature, they allow us to quantitatively evaluate the performance of C-VS3VM, in order to provide a coherent benchmark for the successive comparisons. Results of this experiment are shown in Tab. 8.3.

We can see that, for all the datasets, C-VS3VM outperforms standard SVMs trained only on labeled data, with a reduction of the classification error ranging from 2.37% on GARAGEBAND to 15.22% on PCMAC. Clearly, the training time required by C-VS3VM is higher than the time required by a standard SVM, due to the larger number of training data, and to the use of the gradient descent algorithm. Another important aspect to be considered is that, with the only exception of G50C, the RBF-SVM fails in matching the performance of the linear model due to higher complexity of the model in relationship to the amount of training data.

Next, we investigate the convergence behavior of DG-VS3VM and NEXT-VS3VM, compared to the centralized implementation. In particular, we test the algorithm

Table 8.3: Average value for classification error and computational time for the centralized algorithms.

Dataset	Algorithm	Error [%]	Time [s]
G50C	LIN-SVM	13.79	0.0008
	RBF-SVM	13.36	0.0005
	C-VS3VM	6.36	0.024
PCMAC	LIN-SVM	21.32	0.0035
	RBF-SVM	36.68	0.0032
	C-VS3VM	6.10	35.12
GARAGEBAND	LIN-SVM	23.87	0.0010
	RBF-SVM	27.92	0.0007
	C-VS3VM	21.50	0.2872

on randomly generated networks of $L = 25$ nodes. Results are presented in Fig. 8.2. Particularly, panels on the left show the evolution of the global cost function in (8.6), while panels on the right show the evolution of the squared norm of the gradient. For readability, the graphs use a logarithmic scale on the y -axis, while on the left we only show the first 50 iterations of the optimization procedure. The results are similar for all three datasets, namely, NEXT-VS3VM is able to converge faster (up to one/two orders of magnitude) than DG-VS3VM, which can only exploit first order information on the local cost functions. Indeed, both NEXT-VS3VM and the centralized implementation are able to converge to a stationary point in a relatively small number of iterations, as shown by the panels on the left. The same can be seen from the gradient norm evolution, shown on the right panels, where the fast convergence of NEXT-VS3VM is even more pronounced. Similar insights can be obtained by the analysis of the box plots in Fig. 8.3, where we also compare with the results of LIN-SVM and RBF-SVM obtained previously.

As a final experiment, we investigate the scalability of the distributed algorithms, by analyzing the training time and the test error of DG-VS3VM and NEXT-VS3VM when varying the number of nodes in the network from $L = 5$ to $L = 40$ by steps of 5. Results of this experiment are shown in Fig. 8.4. The three panels on the left show the evolution of the classification error, while the three panels on the right show the evolution of the training time. Results of LIN-SVM, RBF-SVM and C-VS3VM are shown with dashed lines for comparison. It is possible to see that NEXT-VS3VM can track efficiently the centralized solution in all settings, regardless of the size of the network, while DG-VS3VM is not able to properly converge (in the required number of iterations) for larger networks on PCMAC. With respect to training time, results are more varied. Generally speaking, NEXT-VS3VM requires in average more training time than DG-VS3VM. However, for large datasets (PCMAC and GARAGEBAND) both algorithms are comparable in training time with the centralized solution and, more notably, their training time generally decreases for

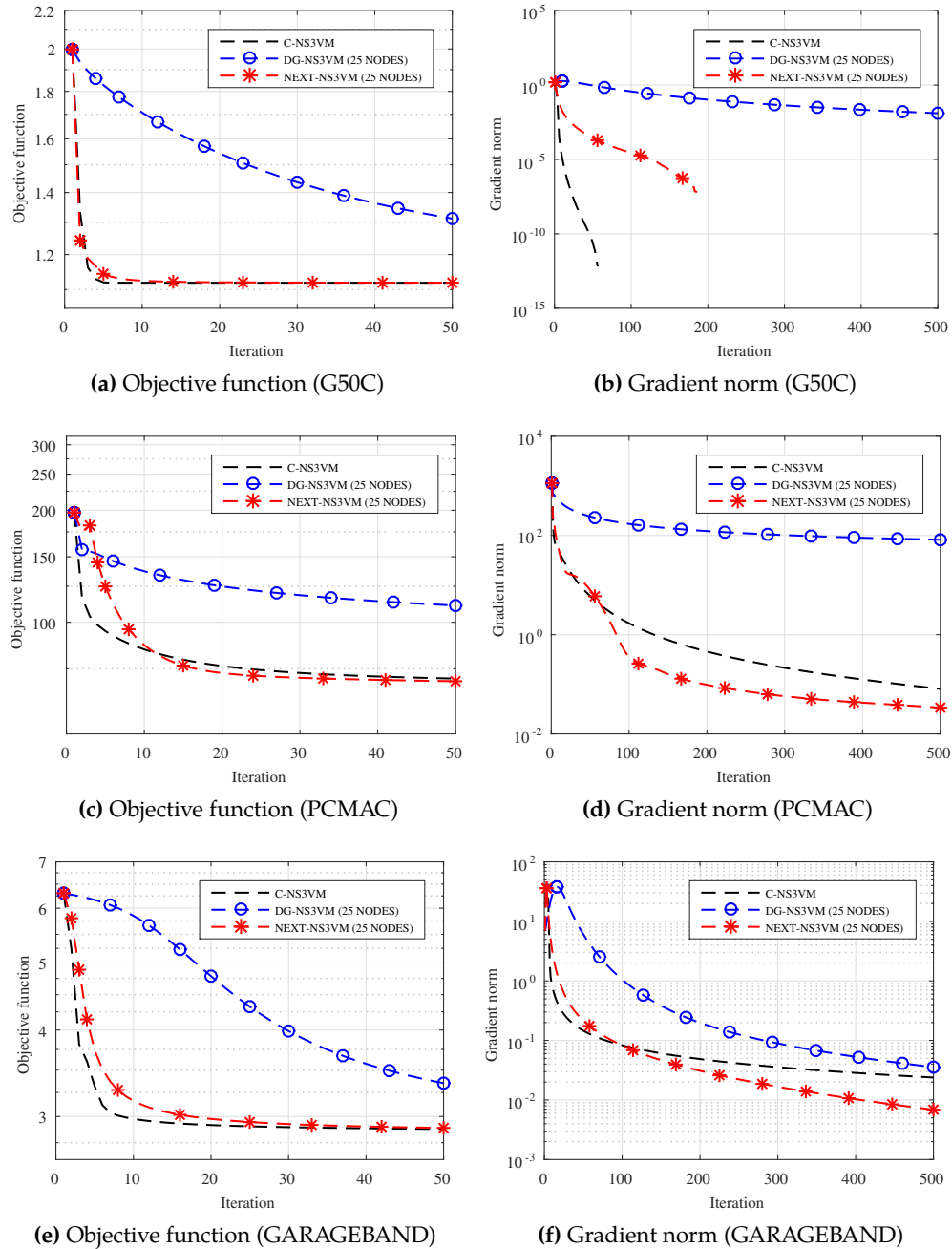


Figure 8.2: Convergence behavior of DG-VS3VM and NEXT-VS3VM, compared to C-VS3VM. The panels on the left show the evolution of the global cost function, while the panels on the right show the evolution of the squared norm of the gradient.

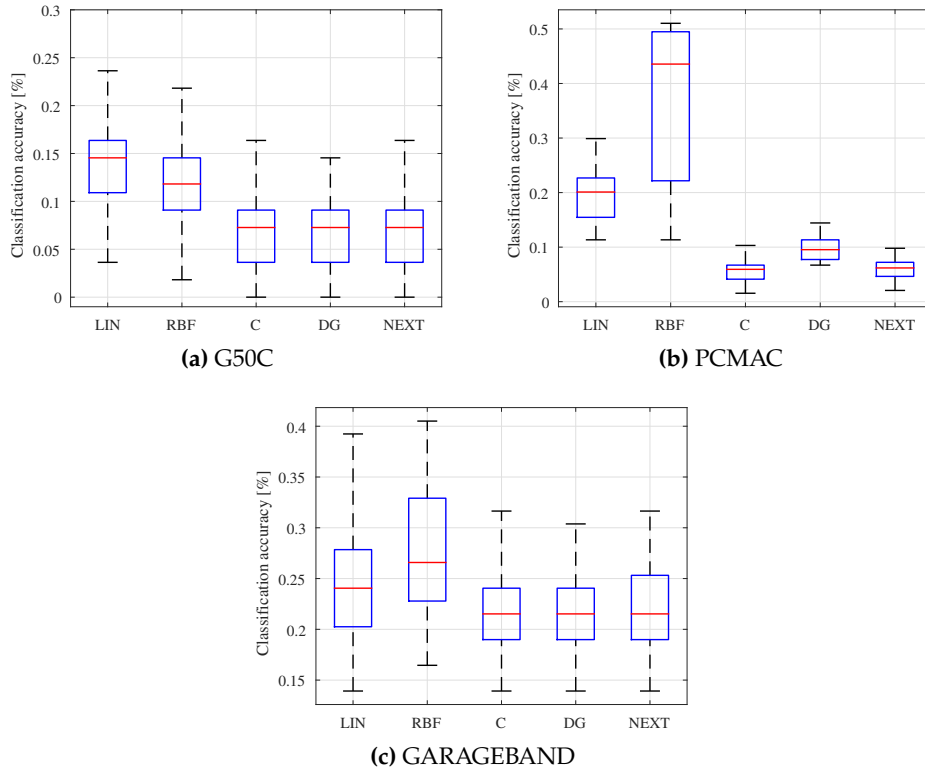


Figure 8.3: Box plots for the classification accuracy of the 5 algorithms, in the case $N = 25$. The central line is the median, the edges are the 25th and 75th percentiles, and the whiskers extend to the most extreme data points. For readability, the names of the algorithms have been abbreviated to LIN (LIN-SVM), RBF (RBF-SVM), C (C-VS3VM), DG (DG-VS3VM) and NEXT (NEXT-VS3VM).

bigger networks.

It is worth mentioning here that the results presented in this chapter strongly depend on our selection of the step-size sequences, and the specific surrogate function in (8.17). In the former case, it is known that the convergence speed of any gradient descent procedure can be accelerated by considering a proper adaptable step-size criterion. Along similar reasonings, the training time of NEXT-VS3VM can in principle be decreased by loosening the precision to which the internal surrogate function is optimized, due to the convergence properties of NEXT already mentioned above. Finally, we can also envision a different choice of surrogate function for NEXT-VS3VM, in order to achieve a different trade-off between training time and speed of convergence. As an example, we can replace the hinge loss $l_k(\mathbf{w})$ with its first-order linearization $\tilde{l}_k(\mathbf{w})$, similarly to (8.15). In this case, the resulting optimization problem would have a closed form solution, resulting in a faster training time per iteration (at the cost of more iterations required for convergence).

Overall, the experimental results suggest that both algorithms can be efficient tools for training a VS3VM in a distributed setting, wherein NEXT-VS3VM is able to converge extremely faster, at the expense of a larger training time. Thus, the choice

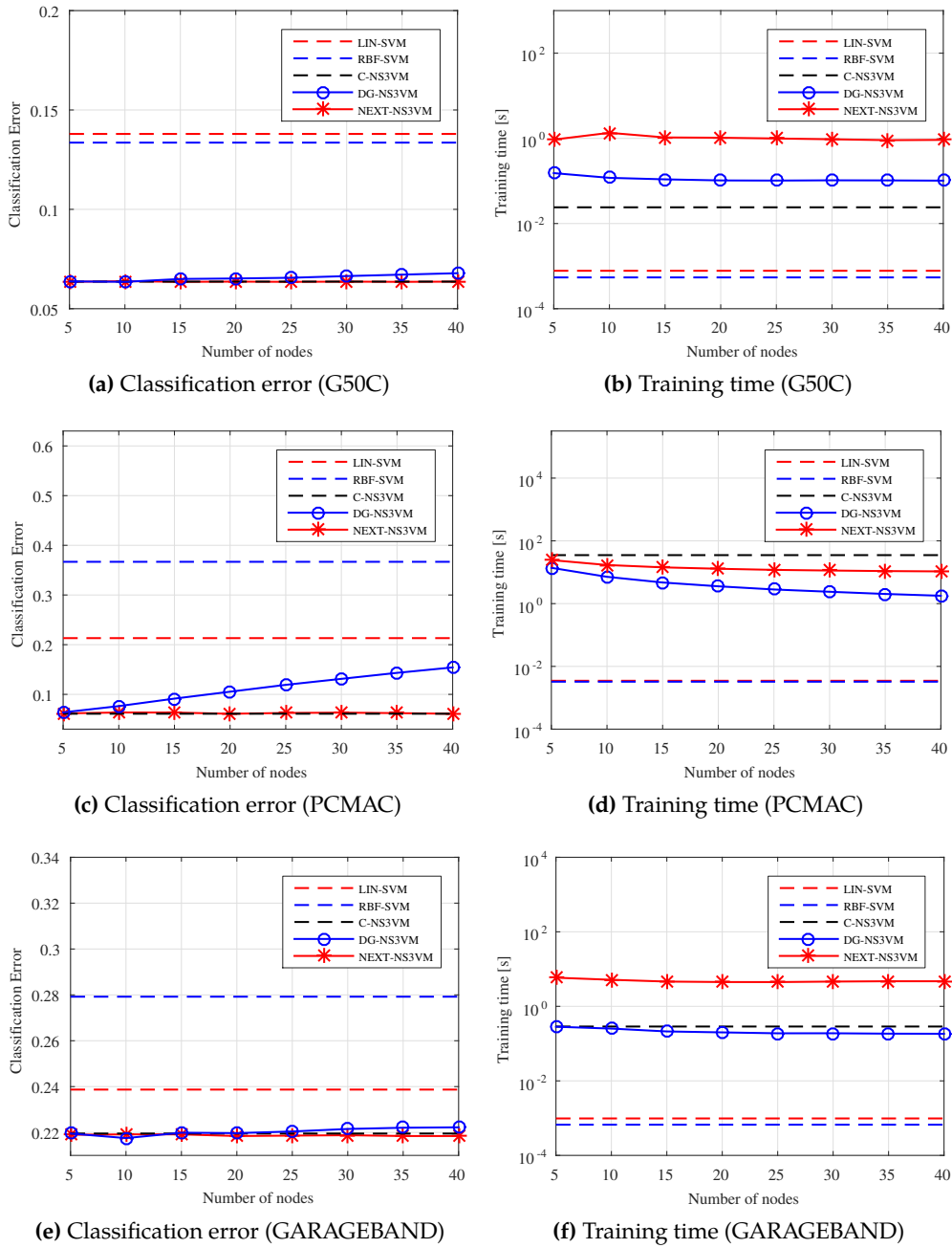


Figure 8.4: Training time and test error of GD-VS3VM and NEXT-VS3VM when varying the number of nodes in the network from $L = 5$ to $L = 40$. Results for LIN-SVM, RBF-SVM and C-VS3VM are shown with dashed lines for comparison.

of a specific algorithm will depend on the applicative domain, and on the amount of computational resources (and size of the training dataset) available to each agent.

Part IV

Distributed Learning from Time-Varying Data

Distributed Training for Echo State Networks

Contents

9.1	Introduction	108
9.2	A primer on ESNs	109
9.3	Distributed training for ESNs	110
9.4	Experimental Setup	112
9.4.1	Description of the Datasets	112
9.4.2	Description of the Algorithms	114
9.4.3	ESN Architecture	115
9.5	Experimental Results	116
9.6	Extension to ESNs with Sparse Readouts	119
9.6.1	Comparisons in the centralized case	120
9.6.2	Comparisons in the distributed case	120

9.1 Introduction

IN the previous part of this thesis, we considered *static* classification and regression tasks, where the order of presentation of the different examples does not matter. In many real world applications, however, the patterns exhibit a temporal dependence among them, as in time-series prediction. In this case, it is necessary to include some form of memory of the previously observed patterns in the ANN models. In this respect, there are two main possibilities. The first is to include an external memory, by feeding as input a buffer of the last K patterns, with K chosen *a priori*. Differently, it is possible to consider *recurrent* connections inside the ANN, which effectively create an internal memory of the previous state, making the ANN

The content of this chapter is adapted from the material published in [151], except Section 9.6, whose content is currently under final editorial review at IEEE Computational Intelligence Magazine.

a dynamic model. This last class of ANNs are called recurrent neural networks (RNNs).

In the DL setting, the former option has been investigated extensively, particularly using linear and kernel adaptive filters (see Section 3.4.2 and Section 3.4.5). The latter option, however, has received considerably less attention. In fact, despite numerous recent advances (e.g. [69]), RNN training remains a daunting task even in the centralized case, mostly due to the well-known problems of the exploding and vanishing gradients [125]. A decentralized training algorithm for RNNs, however, would be an invaluable tool in multiple large-scale real world applications, including time-series prediction on WSNs [130], and multimedia classification over P2P networks.

In this chapter we aim to bridge (partially) this gap, by proposing a distributed training algorithm for a recurrent extension of the RVFL, the ESN. ESNs were introduced by H. Jaeger [78] and together with liquid state machines and backpropagation-decorrelation, they form the family of RNNs known as reservoir computing [104]. The main idea of ESNs, similar to RVFLs, is to separate the recurrent part of the network (the so-called ‘reservoir’), from the non-recurrent part (the ‘readout’). The reservoir is typically fixed in advance, by randomly assigning its connections, and the learning problem is reduced to a standard linear regression over the weights of the readout. Due to this, ESNs do not required complex back-propagation algorithms over the recurrent portion of the network, thus avoiding the problems of the exploding and vanishing gradients. Over the last years, ESNs have been applied successfully to a wide range of domains, including chaotic time-series prediction [80, 90], load prediction [15], grammatical inference [179], and acoustic modeling [182], between others. While several researchers have investigated the possibility of spatially distributing the reservoir [118, 167, 185], to the best of our knowledge, no algorithm has been proposed to train an ESN in the DL setting.

The remaining of the chapter is formulated as follows. In Section 9.2 we introduce the basic concepts on ESNs and a least-square criterion for training them. Section 9.3 details a distributed algorithm for ESNs, extending the ADMM-RVFL presented in Chapter 4. After some experimental results, we also present an extension to consider ESNs with *sparse* readouts in Section 9.6.

9.2 A primer on ESNs

An ESN is a recurrent neural network which can be partitioned in three components, as shown in Fig. 9.1.

The N_i -dimensional input vector $\mathbf{x}[n] \in \mathbb{R}^{N_i}$ is fed to an N_r -dimensional reservoir,

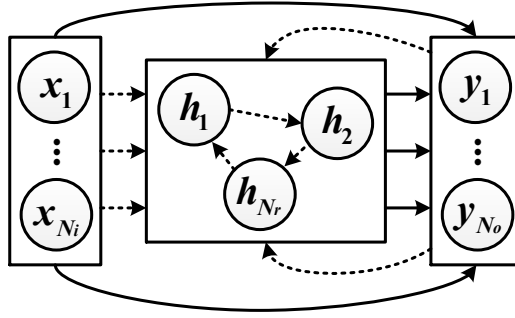


Figure 9.1: Schematic depiction of an ESN. Random connections are shown with dashed lines, while trainable connections are shown with solid lines.

whose internal state $\mathbf{h}[n-1] \in \mathbb{R}^{N_r}$ is updated according to the state equation:

$$\mathbf{h}[n] = f_{\text{res}}(\mathbf{W}_i^r \mathbf{x}[n] + \mathbf{W}_r^r \mathbf{h}[n-1] + \mathbf{W}_o^r \mathbf{y}[n-1]), \quad (9.1)$$

where $\mathbf{W}_i^r \in \mathbb{R}^{N_r \times N_i}$, $\mathbf{W}_r^r \in \mathbb{R}^{N_r \times N_r}$ and $\mathbf{W}_o^r \in \mathbb{R}^{N_r \times N_o}$ are randomly generated matrices, $f_{\text{res}}(\cdot)$ is a suitably defined nonlinear function, and $\mathbf{y}[n-1] \in \mathbb{R}^{N_o}$ is the previous N_o -dimensional output of the network. To increase stability, it is possible to add a small uniform noise term to the state update, before computing the nonlinear transformation $f_{\text{res}}(\cdot)$ [79]. Then, the current output is computed according to:

$$\mathbf{y}[n] = f_{\text{out}}(\mathbf{W}_i^o \mathbf{x}[n] + \mathbf{W}_r^o \mathbf{h}[n]), \quad (9.2)$$

where $\mathbf{W}_i^o \in \mathbb{R}^{N_o \times N_i}$, $\mathbf{W}_r^o \in \mathbb{R}^{N_o \times N_r}$ are adapted based on the training data, and $f_{\text{out}}(\cdot)$ is an invertible nonlinear function. For simplicity, in the rest of the chapter we will consider the case of one-dimensional output, i.e. $N_o = 1$, but everything we say extends straightforwardly to the case $N_o > 1$.

To be of use in any learning application, the reservoir must satisfy the so-called ‘echo state property’ (ESP) [104]. Informally, this means that the effect of a given input on the state of the reservoir must vanish in a finite number of time-instants. A widely used rule-of-thumb that works well in most situations is to rescale the matrix \mathbf{W}_r^r to have $\rho(\mathbf{W}_r^r) < 1$, where $\rho(\cdot)$ denotes the spectral radius operator. For simplicity, we adopt this heuristic strategy in this chapter, but we refer the interested reader to [203] for recent theoretical studies on this aspect. If the ESP is satisfied, an ESN with a suitably large N_r can approximate any nonlinear filter with bounded memory to any given level of accuracy [104].

9.3 Distributed training for ESNs

To train the ESN, suppose we are provided with a sequence of Q desired input-output pairs $(\mathbf{x}[1], d[1]) \dots, (\mathbf{x}[Q], d[Q])$. The sequence of inputs is fed to the

reservoir, giving a sequence of internal states $\mathbf{h}[1], \dots, \mathbf{h}[Q]$ (this is known as ‘warming’). During this phase, since the output of the ESN is not available for feedback, the desired output is used instead in Eq. (9.1) (so-called ‘teacher forcing’). Define the hidden matrix \mathbf{H} and output vector \mathbf{d} as:

$$\mathbf{H} = \begin{bmatrix} \mathbf{x}^T[1] \mathbf{h}^T[1] \\ \vdots \\ \mathbf{x}^T[1] \mathbf{h}^T[Q] \end{bmatrix} \quad (9.3)$$

$$\mathbf{d} = \begin{bmatrix} f_{\text{out}}^{-1}(d[1]) \\ \vdots \\ f_{\text{out}}^{-1}(d[Q]) \end{bmatrix} \quad (9.4)$$

The optimal output weight vector is then given by solving the following regularized least-square problem:

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^{N_i + N_r}} \frac{1}{2} \|\mathbf{H}\mathbf{w} - \mathbf{d}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2, \quad (9.5)$$

where $\mathbf{w} = [\mathbf{W}_i^o \ \mathbf{W}_r^o]^T$ and λ is the standard regularization factor.¹ Solution of problem (9.5) is a standard LRR problem as in Eq. (2.5), and can be obtained in closed form as:

$$\mathbf{w}^* = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^T \mathbf{d}. \quad (9.6)$$

Whenever $N_r + N_i > Q$, Eq. (9.6) can be computed more efficiently by rewriting it as:

$$\mathbf{w}^* = \mathbf{H}^T (\mathbf{H}\mathbf{H}^T + \lambda \mathbf{I})^{-1} \mathbf{d}. \quad (9.7)$$

More in general, we are provided with a training set S of multiple desired sequences. In this case, we can simply stack the resulting hidden matrices and output vectors, and solve Eq. (9.5). Additionally, we note that in practice we can remove the initial D elements (denoted as ‘wash-out’ elements) from each sequence when solving the least-square problem, with D specified *a-priori*, due to their transient state. In the DL setting, we suppose that the S sequences are distributed among the L agents. Clearly, since training results in a LRR problem, at this point we can directly apply any of the algorithms presented in Chapter 4. In particular, we choose to apply the ADMM algorithm due to its convergence properties. The resulting distributed protocol is summarized in Algorithm 9.1.

¹Since we consider one dimensional outputs, \mathbf{W}_i^o and \mathbf{W}_r^o are now row vectors, of dimensionality N_i and N_r respectively.

Algorithm 9.1 ADMM-ESN: Local training algorithm for ADMM-based ESN (k th node).

Inputs: Training set S_k (local), size of reservoir N_r (global), regularization factors λ, γ (global), maximum number of iterations T (global)

Output: Optimal output weight vector \mathbf{w}^*

- 1: Assign matrices $\mathbf{W}_i^r, \mathbf{W}_r^r$ and \mathbf{W}_o^r , in agreement with the other agents in the network.
 - 2: Gather the hidden matrix \mathbf{H}_k and teacher signal \mathbf{d}_k from S_k .
 - 3: Initialize $\mathbf{t}_k[0] = \mathbf{0}, \mathbf{z}[0] = \mathbf{0}$.
 - 4: **for** n from 0 to T **do**
 - 5: $\mathbf{w}_k[n+1] = (\mathbf{H}_k^T \mathbf{H}_k + \gamma \mathbf{I})^{-1} (\mathbf{H}_k^T \mathbf{d}_k - \mathbf{t}_k[n] + \gamma \mathbf{z}[n])$.
 - 6: Compute averages $\hat{\mathbf{w}}$ and $\hat{\mathbf{t}}$ by means of the DAC procedure (see Appendix A.2).
 - 7: $\mathbf{z}[n+1] = \frac{\gamma \hat{\mathbf{w}} + \hat{\mathbf{t}}}{\lambda/L + \gamma}$.
 - 8: $\mathbf{t}_k[n+1] = \mathbf{t}_k[n] + \gamma (\mathbf{w}_k[n+1] - \mathbf{z}[n+1])$.
 - 9: Check termination with residuals (see Section 4.2.2).
 - 10: **end for**
 - 11: **return** $\mathbf{z}[n]$
-

Remark

A large number of techniques have been developed to increase the generalization capability of ESNs without increasing its computational complexity [104]. Provided that the optimization problem in Eq. (9.5) remains unchanged, and the topology of the ESN is not modified during the learning process, many of them can be applied straightforwardly to the distributed training case with the algorithm presented in this chapter. Examples of techniques that can be used in this context include lateral inhibition [202] and random projections [22]. Conversely, techniques that cannot be straightforwardly applied include intrinsic plasticity [170] and reservoir's pruning [149].

9.4 Experimental Setup

In this section we describe our experimental setup. Simulations were performed on MATLAB R2013a, on a 64bit operative system, using an Intel® Core™ i5-3330 CPU with 3 GHZ and 16 GB of RAM.

9.4.1 Description of the Datasets

We validate the proposed ADMM-ESN on four standard artificial benchmarks applications, related to nonlinear system identification and chaotic time-series

prediction. These are tasks where ESNs are known to perform at least as good as the state of the art [104]. Additionally, they are common in distributed scenarios. To simulate a large-scale analysis, we consider datasets that are approximately 1–2 orders of magnitude larger than previous works. In particular, for every dataset we generate 50 sequences of 2000 elements each, starting from different initial conditions, summing up to 100.000 samples for every experiment. This is roughly the limit at which a centralized solution is amenable for comparison. Below we provide a brief description of the four datasets.

The NARMA-10 dataset (denoted by N10) is a nonlinear system identification task, where the input $x[n]$ to the system is white noise in the interval $[0, 0.5]$, while the output $d[n]$ is computed from the recurrence equation [79]:

$$d[n] = 0.1 + 0.3d[n-1] + 0.05d[n-1] \prod_{i=1}^{10} d[n-i] + 1.5x[n]x[n-9]. \quad (9.8)$$

The output is then squashed to the interval $[-1, +1]$ by the nonlinear transformation:

$$d[n] = \tanh(d[n] - \hat{d}), \quad (9.9)$$

where \hat{d} is the empirical mean computed from the overall output vector.

The second dataset is the extended polynomial (denoted by EXTPOLY) introduced in [22]. The input is given by white noise in the interval $[-1, +1]$, while the output is computed as:

$$d[n] = \sum_{i=0}^p \sum_{j=0}^{p-i} a_{ij} x^i[n] x^j[n-l], \quad (9.10)$$

where $p, l \in \mathbb{R}$ are user-defined parameters controlling the memory and nonlinearity of the polynomial, while the coefficients a_{ij} are randomly assigned from the same distribution as the input data. In our experiments, we use a mild level of memory and nonlinearity by setting $p = l = 7$. The output is normalized using Eq. (9.9).

The third dataset is the prediction of the well-known Mackey-Glass chaotic time-series (denoted as MKG). This is defined in continuous time by the differential equation:

$$\dot{x}[n] = \beta x[n] + \frac{\alpha x[n-\tau]}{1 + x^\gamma[n-\tau]}. \quad (9.11)$$

We use the common assignment $\alpha = 0.2, \beta = -0.1, \gamma = 10$, giving rise to a chaotic behavior for $\tau > 16.8$. In particular, in our experiments we set $\tau = 30$. Time-series (9.11) is integrated with a 4-th order Runge-Kutta method using a time step of 0.1, and then sampled every 10 time-instants. The task is a 10-step ahead prediction

task, i.e.:

$$d[n] = x[n + 10]. \quad (9.12)$$

The fourth dataset is another chaotic time-series prediction task, this time on the Lorenz attractor. This is a 3-dimensional time-series, defined in continuous time by the following set of differential equations:

$$\begin{cases} \dot{x}_1[n] &= \sigma (x_2[n] - x_1[n]) \\ \dot{x}_2[n] &= x_1[n] (\eta - x_3[n]) - x_2[n] , \\ \dot{x}_3[n] &= x_1[n]x_2[n] - \zeta x_3[n] \end{cases} \quad (9.13)$$

where the standard choice for chaotic behavior is $\sigma = 10$, $\eta = 28$ and $\zeta = 8/3$. The model in Eq. (9.13) is integrated using an ODE45 solver, and sampled every second. For this task, the input to the system is given by the vector $[x_1[n] \ x_2[n] \ x_3[n]]$, while the required output is a 1-step ahead prediction of the x_1 component, i.e.:

$$d[n] = x_1[n + 1]. \quad (9.14)$$

For all four datasets, we supplement the original input with an additional constant unitary input, as is standard practice in ESNs' implementations [104].

9.4.2 Description of the Algorithms

In our simulations we generate a network of agents, using a random topology model for the connectivity matrix, where each pair of nodes can be connected with 25% probability. The only global requirement is that the overall network is connected. We experiment with a number of nodes going from 5 to 25, by steps of 5. To estimate the testing error, we perform a 3-fold cross-validation on the 50 original sequences. For every fold, the training sequences are evenly distributed across the nodes, and the following three algorithms are compared:

Centralized ESN (C-ESN): This simulates the case where training data is collected on a centralized location, and the net is trained by directly solving problem (9.5).

Local ESN (L-ESN): In this case, each node trains a local ESN starting from its data, but no communication is performed. The testing error is then averaged throughout the L nodes.

ADMM-based ESN (ADMM-ESN): This is an ESN trained with the distributed protocol introduced in the previous section. We set $\rho = 0.01$, a maximum number of 400 iterations, and $\epsilon_{\text{abs}} = \epsilon_{\text{rel}} = 10^{-4}$.

All algorithms share the same ESN architecture, which is detailed in the following section. The 3-fold cross-validation procedure is repeated 15 times by varying the ESN initialization and the data partitioning, and the errors for every iteration and every fold are collected. To compute the error, we run the trained ESN on the test sequences, and gather the predicted outputs $\tilde{y}_1, \dots, \tilde{y}_K$, where K is the number of testing samples after removing the wash-out elements from the test sequences. Then, we compute the Normalized Root Mean-Squared Error (NRMSE), defined as:

$$\text{NRMSE} = \sqrt{\frac{\sum_{i=1}^K [\tilde{y}_i - d_i]^2}{|K| \hat{\sigma}_d}}, \quad (9.15)$$

where $\hat{\sigma}_d$ is an empirical estimate of the variance of the true output samples d_1, \dots, d_K .

9.4.3 ESN Architecture

As stated previously, all algorithms share the same ESN architecture. In this section we provide a brief overview on the selection of its parameters. First, we choose a default reservoir's size of $N_r = 300$, which was found to work well in all situations. Secondly, since the datasets are artificial and noiseless, we set a small regularization factor $\lambda = 10^{-3}$. Four other parameters are instead selected based on a grid search procedure. The validation error for the grid-search procedure is computed by performing a 3-fold cross-validation over 9 sequences, which are generated independently from the training and testing set. Each validation sequence has length 2000. In particular, we select the following parameters:

- The matrix \mathbf{W}_i^r , connecting the input to the reservoir, is initialized as a full matrix, with entries assigned from the uniform distribution $[-\alpha_i \ \alpha_i]$. The optimal parameter α_i is searched in the set $\{0.1, 0.3, \dots 0.9\}$.
- Similarly, the matrix \mathbf{W}_o^r , connecting the output to the reservoir, is initialized as a full matrix, with entries assigned from the uniform distribution $[-\alpha_f \ \alpha_f]$. The parameter α_f is searched in the set $\{0, 0.1, 0.3, \dots 0.9\}$. We allow $\alpha_f = 0$ for the case where no output feedback is needed.
- The internal reservoir matrix \mathbf{W}_r^r is initialized from the uniform distribution $[-1 \ 1]$. Then, on average 75% of its connections are set to 0, to encourage sparseness. Finally, the matrix is rescaled so as to have a desired spectral radius ρ , which is searched in the same interval as α_i .
- We use $\tanh(\cdot)$ nonlinearities in the reservoir, while a scaled identity $f(s) = \alpha_t s$ as the output function. The parameter α_t is searched in the same interval as α_i .

Additionally, we insert uniform noise in the state update of the reservoir, sampled uniformly in the interval $[0, 10^{-3}]$, and we discard $D = 100$ initial elements from each sequence.

9.5 Experimental Results

The final settings resulting from the grid-search procedure are listed in Table 9.1.

Table 9.1: Optimal parameters found by the grid-search procedure. For a description of the parameters, see Section 9.4.2.

Dataset	ρ	α_i	α_t	α_f	N_r	λ
N10	0.9	0.5	0.1	0.3		
EXTPOLY	0.7	0.5	0.1	0	300	2^{-3}
MKG	0.9	0.3	0.5	0		
LORENZ	0.1	0.9	0.1	0		

It can be seen that, except for the LORENZ dataset, there is a tendency towards selecting large values of ρ . Output feedback is needed only for the N10 dataset, while it is found unnecessary in the other three datasets. The optimal input scaling α_f is ranging in the interval $[0.5, 0.9]$, while the optimal teacher scaling α_t is small in the majority of cases.

The average NRMSE and training times for C-ESN are provided in Table 9.2 as a reference.

Table 9.2: Final misclassification error and training time for C-ESN, provided as a reference, together with one standard deviation.

Dataset	NRMSE	Time [secs]
N10	0.08 ± 0.01	9.26 ± 0.20
EXTPOLY	0.39 ± 0.01	8.96 ± 0.19
MKG	0.09 ± 0.01	9.26 ± 0.20
LORENZ	0.67 ± 0.01	9.47 ± 0.14

Clearly, these values do not depend on the size of the network, and they can be used as an upper baseline for the results of the distributed algorithms. Since we are considering the same amount of training data for each dataset, and the same reservoir's size, the training times in Table 9.2 are roughly similar, except for the LORENZ dataset, which has 4 inputs compared to the other three datasets (considering also the unitary input). As we stated earlier, performance of C-ESN are competitive with the state-of-the-art for all the four datasets. Moreover, we can see that it is extremely efficient to train, taking approximately 9 seconds in all cases.

To study the behavior of the decentralized procedures when training data is distributed, we plot the average error for the three algorithms, when varying the number of nodes in the network, in Fig. 9.2 (a)-(d). The average NRMSE of C-ESN is shown as dashed black line, while the errors of L-ESN and ADMM-ESN are shown with blue squares and red circles respectively.

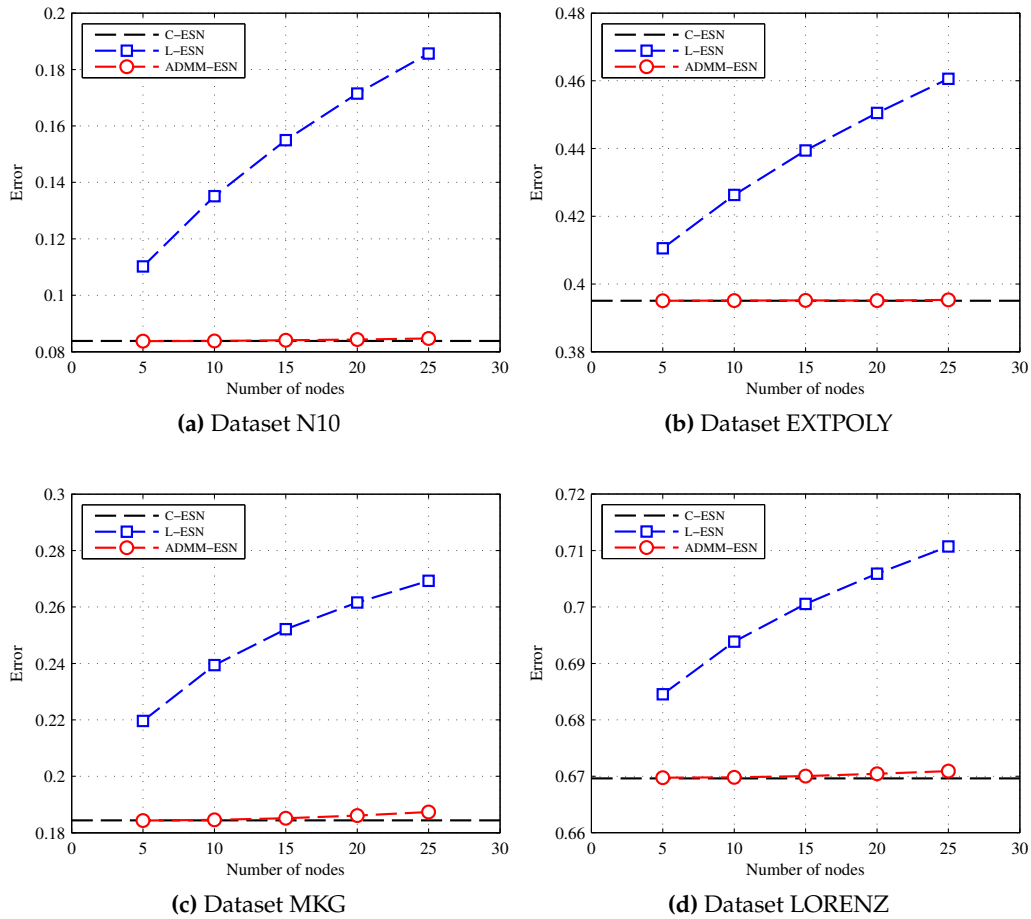


Figure 9.2: Evolution of the testing error, for networks going from 5 agents to 25 agents. Performance of L-ESN is averaged across the nodes.

Clearly, L-ESN is performing worse than C-ESN, due to its partial view on the training data. For small networks of 5 nodes, this gap may not be particularly pronounced. This goes from a 3% worse performance on the LORENZ dataset, up to a 37% decrease in performance for the N10 dataset (going from an NRMSE of 0.8 to an NRMSE of 0.11). The gap is instead substantial for large networks of up to 25 nodes. For example, the error of L-ESN is more than twice that of C-ESN for the N10 dataset, and its performance is 50% worse in the MKG dataset. Albeit these results are expected, they are evidence of the need of a decentralized training protocol for ESNs, able to take into account all the local datasets.

As is clear from Fig. 9.2, ADMM-ESN is able to perfectly track the performance of the centralized solution in all situations. A small gap in performance is present for the two predictions tasks when considering large networks. In particular, the performance of ADMM-ESN is roughly 1% worse than C-ESN for networks of 25 nodes in the datasets MKG and LORENZ. In theory, this gap can be reduced by considering additional iterations for the ADMM procedure, although this would be impractical in real world applications.

Training time requested by the three algorithms is shown in Fig. 9.3 (a)-(d). The training time for L-ESN and ADMM-ESN is averaged throughout the agents.

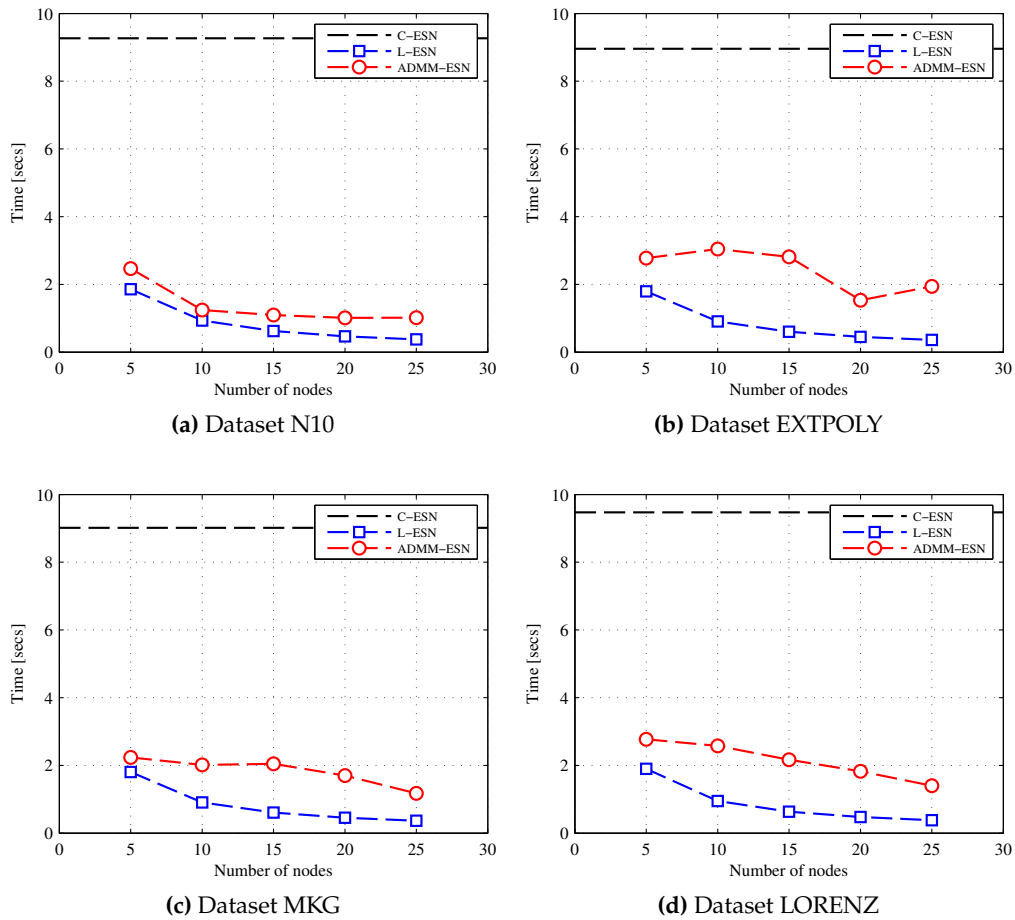


Figure 9.3: Evolution of the training time, for networks going from 5 agents to 25 agents. Time of L-ESN is averaged across the nodes.

Since the computational time of training an ESN is mostly related to the matrix inversion in Eq. (9.6), training time is monotonically decreasing in L-ESN with respect to the number of nodes in the network (the higher the number of agents, the lower the amount of data at every local node). Fig. 9.3 shows that the computational overhead requested by the ADMM procedure is limited. In the best case, the N10 dataset with 10 nodes, it required only 0.3 seconds more than L-ESN, as shown from

Fig. 9.3(a). In the worst setting, the EXTPLY dataset with 15 nodes, it required 2.2 seconds more, as shown from Fig. 9.3(b). In all settings, the time requested by ADMM-ESN is significantly lower compared to the training time of its centralized counterpart, showing its usefulness in large-scale applications.

9.6 Extension to ESNs with Sparse Readouts

Up to now, the chapter has focused on training an ESN with a ridge regression routine. Still, it is known that standard ridge regression may not be the most suitable training algorithm for ESNs. Specifically, a large number of authors have been concerned with training an ESN with a sparse readout, i.e. a readout where the majority of the connections are set to zero. In the centralized case, this has been initially explored in depth in [49]. The authors investigated different greedy methods to this end, including backward selection (where connections are removed one at a time based on an iterative procedure), random deletion, and others. Significant improvements are found, both in terms of generalization accuracy, and in terms of computational requirements. Moreover, having only a small amount of connections can lead to extremely efficient implementations [149], particularly on low-cost devices. Thus, having the possibility of training sparse readouts for an ESN in a decentralized case can be a valuable tool.

Since the readout is linear, sparsity can be enforced by including an additional L_1 regularization term to be minimized, resulting in the LASSO algorithm. For ESNs, this is derived for the first time in Ceperic and Baric [27]. In the distributed case under consideration, the ADMM can be used for solving the LASSO problem quite efficiently, with only a minor modification with respect to the ADMM-ESN [20]. In particular, it is enough to replace the update for $\mathbf{z}[n + 1]$ with:

$$\mathbf{z}[n + 1] = S_{\lambda/N\rho}(\hat{\mathbf{w}}[n + 1] + \mathbf{t}[n]), \quad (9.16)$$

where the soft-thresholding operator $S_\alpha(\cdot)$ is defined for a generic vector \mathbf{a} as:

$$S_\alpha(\mathbf{a}) = \left(1 - \frac{\alpha}{\|\mathbf{a}\|_2}\right)_+ \mathbf{a},$$

and $(\cdot)_+$ is defined element-wise as $(\cdot)_+ = \max(0, \cdot)$. In order to test the resulting sparse algorithm, we consider the MKG and N10 datasets with the same setup as before, but a lower number of elements (in total 2500 for training and 2000 for testing). Additionally, in order to have a slightly redundant reservoir, we select $N_r = 500$.

9.6.1 Comparisons in the centralized case

We begin our experimental evaluation by comparing the standard ESN and the ESN trained using the LASSO algorithm (denoted as L1-ESN) in the centralized case. This allows us to better investigate their behavior, and to choose an optimal regularization parameter λ . Particularly, we analyze test error, training time, and sparsity of the resulting L1-ESN when varying λ in 10^{-j} , with j going from -1 to -6 . The LASSO problems are solved using a freely available implementation of the iterated ridge regression algorithm by M. Schmidt [159].² The algorithm works by approximating the L_1 term with $\|w_i\|_1 \approx \frac{w_i^2}{\|w_i\|_1}$, and iteratively solving the resulting ridge regression problem. Results are presented in Fig. 9.4, where results for MG and N10 are shown in the left and right columns, respectively.

First of all, we can see clearly from Figs 9.4a and 9.4b that the regularization effect of the two algorithms is similar, a result in line with previous works [49]. Particularly, for large regularization factors, the estimates tend to provide an unsatisfactory test error, which however is relatively stable for sufficiently small coefficients. The tendency to select such a small factor is to be expected, due to the artificial nature of the datasets. A minimum in test error is reached for j around -5 for MG, and j around -4 for N10.

With respect to the training time, ridge regression is relatively stable to the amount of regularization, as the matrix to be inverted tends to be already well conditioned. Training time of LASSO is regular for MG, while it slightly increases for larger values of j in the N10 case, as shown in Fig. 9.4d. In all cases, however, it is comparable to that of ridge regression, with a small increase of 0.5 seconds in average.

The most important aspect, however, is evidenced in Figs 9.4e and 9.4f. Clearly, sparsity of the readout goes from almost 100% to 0% as the regularization factor decreases. At the point of best test accuracy, the resulting readout has an average sparsity of 70% for MG and 38% for N10. This, combined with the simultaneous possibility of pruning the resulting reservoir [49, 149], can lead to an extreme saving of computational resources requested at the single sensor during the prediction phase. In order to provide a simpler comparison of the results, we also display them in tabular form in Table 9.3.

9.6.2 Comparisons in the distributed case

We now consider the implementation of the distributed L1-ESN over a network of agents. More in detail, training observations are uniformly subdivided among the L agents in a predefined network, with L varying from 5 to 30 by steps of 5.

²<http://www.cs.ubc.ca/~schmidtm/Software/lasso.html>

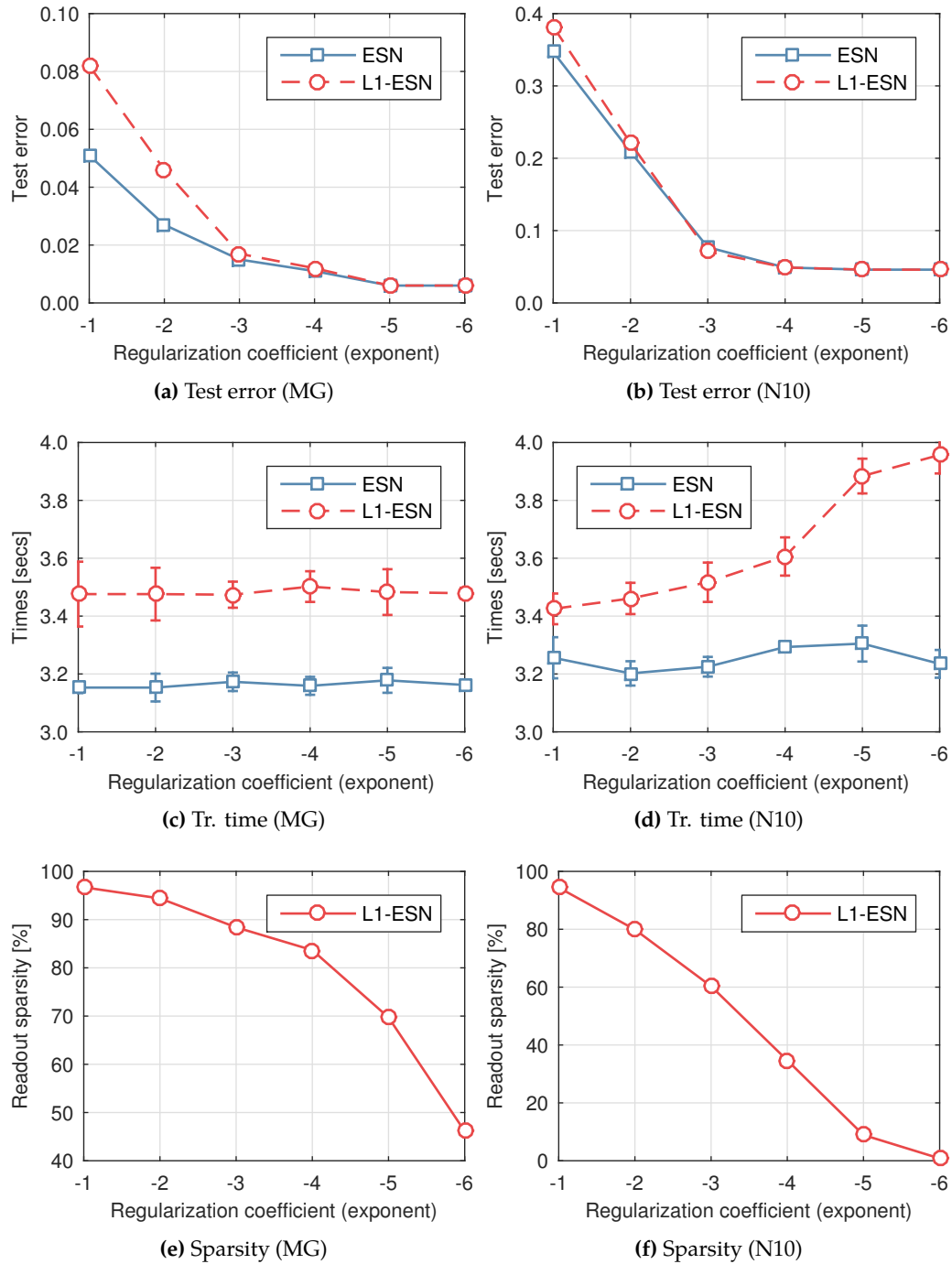


Figure 9.4: Evolution of (a-b) test error, (c-d) training time and (e-f) sparsity of the output vector when varying the regularization coefficient in 10^j , $j = -1, \dots, -6$. Results for the MG dataset are shown on the left column, while results for the N10 dataset are shown in the right column.

Table 9.3: The results of Fig. 9.4, shown in tabular form, together with one standard deviation.

Dataset	λ	Algorithm	Test error (NRMSE)	Tr. time [secs]	Sparsity [%]
MG	10^{-1}	ESN	0.051 ± 0.010	3.153 ± 0.019	0
		L1-ESN	0.082 ± 0.001	3.476 ± 0.112	0.967 ± 0.01
	10^{-2}	ESN	0.027 ± 0.006	3.153 ± 0.048	0
		L1-ESN	0.046 ± 0.001	3.476 ± 0.091	0.944 ± 0.06
	10^{-3}	ESN	0.015 ± 0.003	3.173 ± 0.032	0
		L1-ESN	0.017 ± 0.001	3.474 ± 0.045	0.884 ± 0.01
	10^{-4}	ESN	0.011 ± 0.001	3.159 ± 0.031	0
		L1-ESN	0.012 ± 0.001	3.502 ± 0.053	0.837 ± 0.04
	10^{-5}	ESN	0.006 ± 0.001	3.178 ± 0.043	0
		L1-ESN	0.006 ± 0.001	3.483 ± 0.079	0.697 ± 0.09
	10^{-6}	ESN	0.006 ± 0.001	3.162 ± 0.011	0
		L1-ESN	0.006 ± 0.001	3.976 ± 0.019	0.461 ± 0.03
N10	10^{-1}	ESN	0.347 ± 0.006	3.256 ± 0.071	0
		L1-ESN	0.382 ± 0.008	3.425 ± 0.053	0.944 ± 0.01
	10^{-2}	ESN	0.209 ± 0.004	3.202 ± 0.042	0
		L1-ESN	0.221 ± 0.007	3.461 ± 0.054	0.799 ± 0.04
	10^{-3}	ESN	0.077 ± 0.001	3.225 ± 0.034	0
		L1-ESN	0.071 ± 0.001	3.517 ± 0.068	0.603 ± 0.03
	10^{-4}	ESN	0.049 ± 0.001	3.293 ± 0.015	0
		L1-ESN	0.049 ± 0.001	3.606 ± 0.066	0.347 ± 0.02
	10^{-5}	ESN	0.046 ± 0.001	3.305 ± 0.062	0
		L1-ESN	0.046 ± 0.001	3.884 ± 0.060	0.089 ± 0.01
	10^{-6}	ESN	0.046 ± 0.001	3.235 ± 0.048	0
		L1-ESN	0.046 ± 0.001	3.957 ± 0.064	0.008 ± 0.01

For every run, the connectivity among the agents is generated randomly, such that each pair of agents has a 25% probability of being connected, with the only global requirement that the overall network is connected. The following three algorithms are compared:

1. **Centralized ESN (C-L1-ESN)**: this simulates the case where training data is collected on a centralized location, and the net is trained by directly solving the LRR problem. This is equivalent to the ESN analyzed in the previous section, and following the results obtained there, we set $\lambda = 10^{-5}$ for MG, and $\lambda = 10^{-4}$ for N10.
2. **Local ESN (L-L1-ESN)**: in this case, each agent trains an L1-ESN starting from its local measurements, but no communication is performed. The testing error is averaged throughout the L agents.
3. **ADMM-based ESN (ADMM-L1-ESN)**: this is trained with the algorithm

introduced previously. We select $\gamma = 0.01$ and a maximum number of 400 iterations. For the DAC protocol, we set a maximum number of 300 iterations. DAC also stops whenever the updates (in norm) at every agent are smaller than a predefined threshold $\delta = 10^{-8}$:

$$\|q_k[n+1; j] - q_k[n+1; j-1]\|_2^2 < \delta, \quad k \in \{1, 2, \dots, L\}. \quad (9.17)$$

Results of this set of experiments are presented in Fig. 9.5. Similarly to before, results from the two datasets are presented in the left and right columns, respectively. From Figs 9.5a and 9.5b we see that, although L-L1-ESN achieves degrading performance for bigger networks (due to the lower number of measurements per agent), ADMM-L1-ESN is able to effectively track the performance of the centralized counterpart, except for a small deviation in MG. Indeed, it is possible to reduce this gap by increasing the number of iterations; however, the performance gain is not balanced by the increase in computational cost.

With respect to the training time, it is possible to see from Figs 9.5c and 9.5d that the training time is relatively steady for larger networks in ADMM-L1-ESN, showing its feasibility in the context of large sensor networks. Moreover, the computational cost requested by the distributed procedure is low and, in the worst case, it requires no more than 1 second with respect to the cost of a centralized counterpart. Overall, we can see that our distributed protocol allows for an efficient implementation in terms of performance and training time, while at the same time guaranteeing a good level of sparsity of the resulting readout. This, in turn, is essential for many practical implementations where computational savings are necessary.

Some additional insights into the convergence behavior of ADMM-L1-ESN can also be obtained by analyzing the evolution of the so-called (primal) residual, given by [20]:

$$r[n+1] = \frac{1}{L} \sum_{k=1}^L \|\mathbf{w}[n+1] - \mathbf{z}[n+1]\|_2. \quad (9.18)$$

As can be seen from Fig. 9.6 (shown with a logarithmic y -axis), this rapidly converges towards 0, ensuring that the algorithm is able to reach a stationary solution in a relatively small number of iterations.

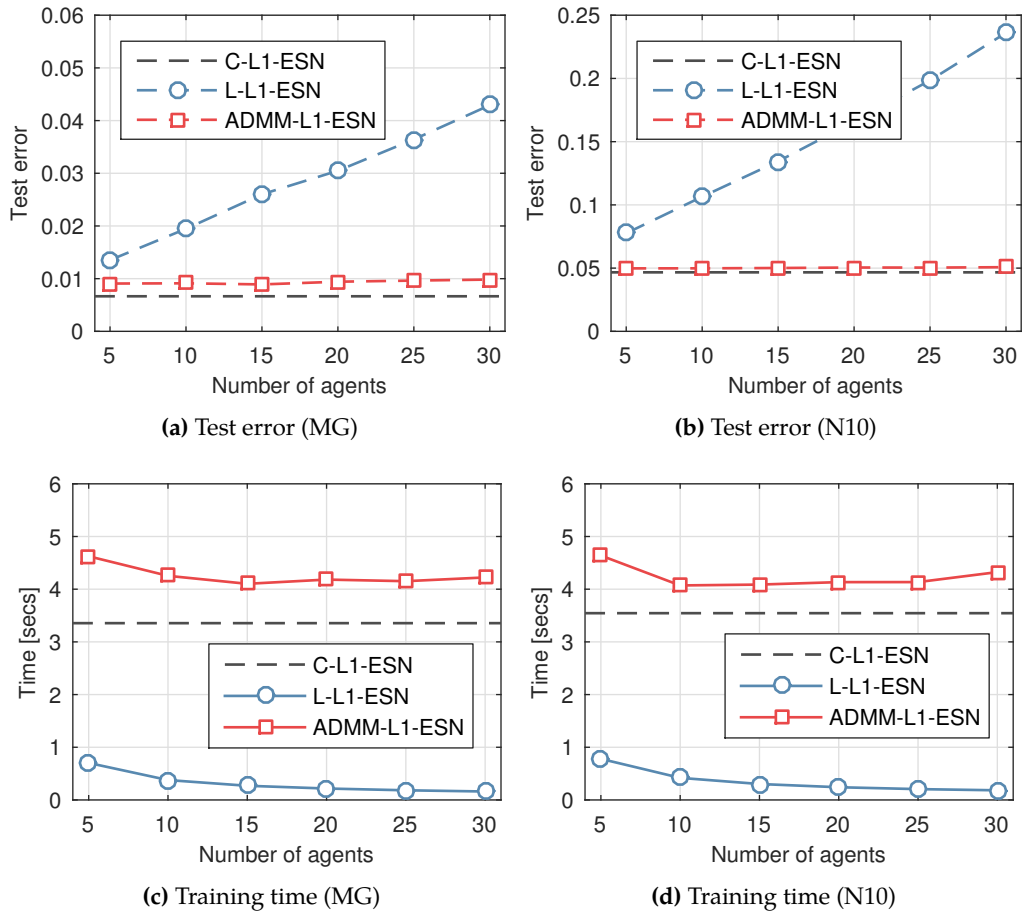


Figure 9.5: Evolution of (a-b) test error, (c-d) training time when varying the number of agents in the network from 5 to 30 by steps of 5.

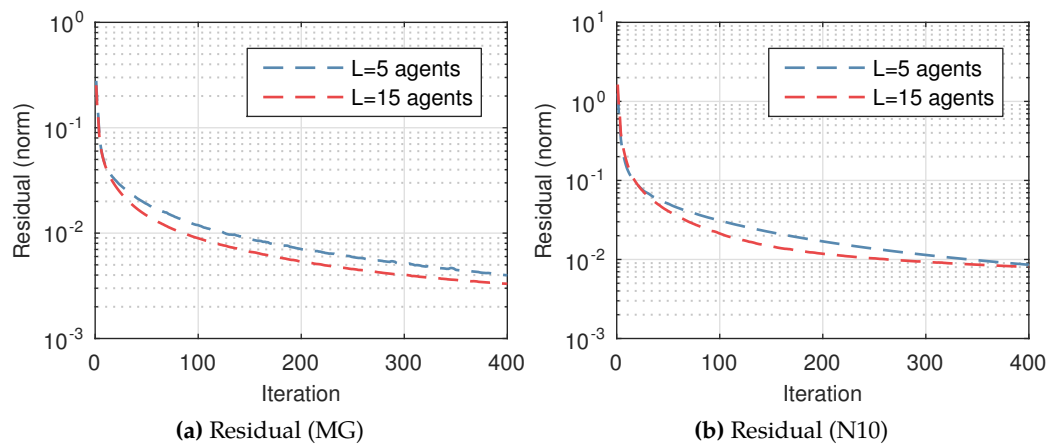


Figure 9.6: Evolution of the (primal) residual of ADMM-L1-ESN for $L = 5$ and $L = 15$.

Diffusion Spline Filtering

Contents

10.1	Introduction	125
10.2	Spline Adaptive Filter	126
10.3	Diffusion SAF	129
10.4	Experimental Setup	132
10.5	Experimental Results	133
10.5.1	Experiment 1 - Small Network ($L = 10$)	133
10.5.2	Experiment 2 - Large Network ($L = 30$)	135
10.5.3	Experiment 3 - Strong nonlinearity ($L = 15$)	136

10.1 Introduction

THIS chapter continues the investigation on distributed training algorithms for time-varying data. Particularly, we focus on models with external memory (i.e., with a buffer of the last input elements), adequate to devices with extremely low computation resources. Available approaches in this sense include the linear diffusion filters (see Section 3.4.2), and kernel-based distributed filters (see Section 3.4.5). However, the former applicability is limited to scenarios where the assumption of a linear model between the output and the observed variables is meaningful. Kernel methods, instead, are hindered by the fact that a kernel model depends by definition on the full observed dataset, as we analyzed extensively in Chapters 3 and 7.

In this chapter, we propose a novel nonlinear distributed filtering algorithm based on the recently proposed spline adaptive filter (SAF) [155]. Specifically, we focus on the Wiener SAF filter [155], where a linear filter is followed by an *adaptive*

A partial content of this chapter is currently under review for the 2016 European Signal Processing Conference (EUSIPCO).

nonlinear transformation, obtained with spline interpolation. They are attractive nonlinear filters for two main reasons. First, the nonlinear part is linear-in-the-parameters (LIP), allowing for the possibility of adapting both parts of the filter using standard linear filtering techniques. Secondly, while the spline can be defined by a potentially large number of parameters, only a small subset of them must be considered and adapted at each time step (4 in our experiments). Due to this, they allow to approximate non-trivial nonlinear functions with a small increase in computational complexity with respect to linear filters.

Based on the general theory of DA,¹ in this chapter we propose a diffused version of the SAF filter, denoted as D-SAF. In particular, we show that a cooperative behavior can be implemented by considering two subsequent diffusion operations, on the linear and non-linear components of the SAF respectively. Due to this, the D-SAF inherits the aforementioned characteristics of the centralized SAF, namely, it enables the agents to collectively converge to a non-linear function, with a small overhead with respect to a purely linear diffusion filter. In fact, D-LMS can be shown to be a special case of D-SAF, where adaptation is restricted to the linear part only. To demonstrate the merits of the proposed D-SAF, we perform an extensive set of experiments, considering medium and large-sized networks, coupled with mild and strong non-linearities. Simulations show that the D-SAF is able to efficiently learn the underlying model, and strongly outperform D-LMS and a purely non-cooperative SAF.

The rest of the chapter is organized as follows. Section 10.2 introduces the basic framework of spline interpolation and SAFs. Section 10.3 formulates the D-SAF algorithm. Subsequently, we details our experimental setup and results in Section 10.4 and Section 10.5 respectively.

10.2 Spline Adaptive Filter

Denote by $x[n]$ the input to the SAF filter at time n , and by $\mathbf{x}_n = [x[n], \dots, x[n - M + 1]]^T$ a buffer of the last M samples. As in the previous chapters, we assume to be dealing with real inputs. Additionally, we assume that an unknown Wiener model is generating the desired response as follows:

$$d[n] = f_0(\mathbf{w}_0^T \mathbf{x}_n) + v[n], \quad (10.1)$$

where $\mathbf{w}_0 \in \mathbb{R}^M$ are the linear coefficients, $f_0(\cdot)$ is a desired nonlinear function, which is supposed continuous and derivable, and $v[n] \sim \mathcal{N}(0, \sigma^2)$ is a Gaussian noise term. Similarly, a SAF computes the output in a two-step fashion. First, it

¹Described in Chapter 3.

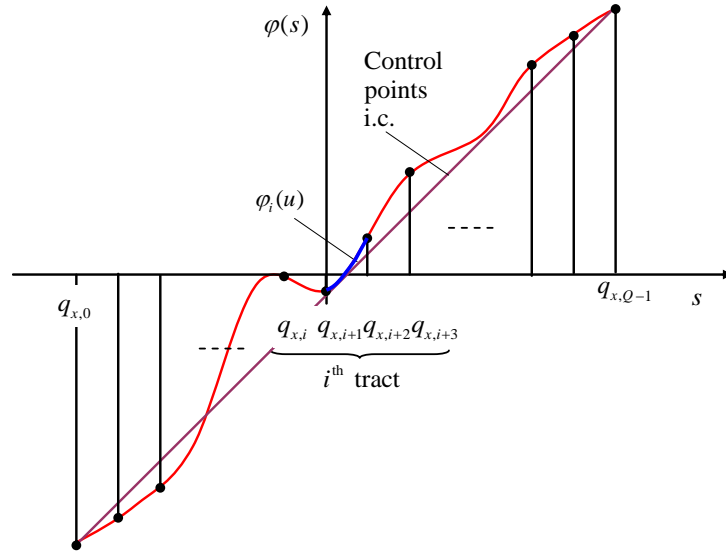


Figure 10.1: Example of spline interpolation scheme. We suppose that the control points are equispaced on the x -axis, and symmetrically spaced around the origin.

performs a linear filtering operation given by:

$$s[n] = \mathbf{w}_n^T \mathbf{x}_n. \quad (10.2)$$

Then, the final output is computed via spline interpolation over $s[n]$. A spline is a flexible polynomial defined by a set of Q control points (called *knots*), and denoted as $\mathbf{Q}_i = [q_{x,i} \ q_{y,i}]$. We suppose that the knots are uniformly distributed, i.e. $q_{x,i+1} = q_{x,i} + \Delta x$, for a fixed $\Delta x \in \mathbb{R}$. Without lack of generality, we also constrain the knots to be symmetrically spaced around the origin. This pair of assumptions are at the base of the SAF family of algorithms, and dates back to earlier work on spline neurons for multilayer perceptrons [65]. Practically, they allow for a simple derivation of the adaptation rule, while sacrificing only a small part of the flexibility of the spline interpolation framework. This is shown pictorially in Fig. 10.1.

Given the output of the linear filter $s[n]$, the spline is defined as an interpolating polynomial of order P , passing by the closest knot to $s[n]$ and its P successive knots. In particular, due to our earlier assumptions, the index of the closest knot can be computed as:

$$i = \left\lfloor \frac{s[n]}{\Delta x} \right\rfloor + \frac{Q-1}{2}. \quad (10.3)$$

Given this, we can define the normalized abscissa value between $q_{x,i}$ and $q_{x,i+1}$ as:

$$u = \frac{s[n]}{\Delta x} - \left\lfloor \frac{s[n]}{\Delta x} \right\rfloor. \quad (10.4)$$

From u we can compute the normalized reference vector $\mathbf{u} = [u^P \ u^{P-1} \ \dots \ u \ 1]^T$,

while from i we can extract the relevant control points $\mathbf{q}_{i,n} = [q_{y,i} \ q_{y,i+1} \ \dots \ q_{y,i+P}]^T$. We refer to the vector $\mathbf{q}_{i,n}$ as the i th *span*. The output of the filter is then given by:

$$y[n] = \varphi(s[n]) = \mathbf{u}^T \mathbf{B} \mathbf{q}_{i,n}, \quad (10.5)$$

where $\varphi(s[n])$ is the adaptable nonlinearity as shown in Fig. 10.1, and $\mathbf{B} \in \mathbb{R}^{(P+1) \times (P+1)}$ is called the spline basis matrix. In this chapter, we use the Catmul-Rom (CR) spline with $P = 3$, given by:

$$\mathbf{B} = \frac{1}{2} \begin{bmatrix} -1 & 3 & -3 & 1 \\ 2 & -5 & 4 & -1 \\ -1 & 0 & 1 & 0 \\ 0 & 2 & 0 & 0 \end{bmatrix}. \quad (10.6)$$

Several alternative choices are available, such as the B-spline matrix [155]. Different bases give rise to alternative interpolation schemes, e.g. a spline defined by a CR basis passes through all the control points, but its second derivative is not continuous, while the opposite is true for the B-spline basis. Note that both (10.2) and (10.5) are LIP, and can be adapted with the use of any standard linear filtering technique. Applying the chain rule, it is straightforward to compute the derivative of the SAF output with respect to the linear coefficients:

$$\begin{aligned} \frac{\partial \varphi(s[n])}{\partial \mathbf{w}_n} &= \frac{\partial \varphi(s[n])}{\partial u} \cdot \frac{\partial u}{\partial s[n]} \cdot \frac{\partial s[n]}{\partial \mathbf{w}_n} = \\ &= \dot{\mathbf{u}} \mathbf{B} \mathbf{q}_{i,n} \left(\frac{1}{\Delta x} \right) \mathbf{x}_n, \end{aligned} \quad (10.7)$$

where:

$$\dot{\mathbf{u}} = \frac{\partial \mathbf{u}}{\partial u} = [P u^{P-1} \ (P-1) u^{P-2} \ \dots \ 1 \ 0]^T. \quad (10.8)$$

Similarly, for the nonlinear part we obtain:

$$\frac{\partial \varphi(s[n])}{\partial \mathbf{q}_{i,n}} = \mathbf{B}^T \mathbf{u}. \quad (10.9)$$

We consider a first-order adaptation for both the linear and the nonlinear part of the SAF. Defining the error $e[n] = d[n] - y[n]$, we aim at minimizing the expected mean-squared error given by:

$$J(\mathbf{w}, \mathbf{q}) = \mathbb{E} \{ e[n]^2 \}, \quad (10.10)$$

where $\mathbf{q} = [q_{y,1}, \dots, q_{y,Q}]^T$. As is standard approach, we approximate (10.10) with

the instantaneous error given by:

$$\hat{J}(\mathbf{w}, \mathbf{q}) = e^2[n]. \quad (10.11)$$

Then, we apply two simultaneous steepest-descent steps to solve the overall optimization problem:

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu_w e[n] \varphi'(s[n]) \mathbf{x}_n, \quad (10.12)$$

$$\mathbf{q}_{i,n+1} = \mathbf{q}_{i,n} + \mu_q e[n] \mathbf{B}^T \mathbf{u}, \quad (10.13)$$

where we defined $\varphi'(s[n]) = \dot{\mathbf{u}} \mathbf{B} \mathbf{q}_{i,n} \left(\frac{1}{\Delta x} \right)$, and we use two possibly different step-sizes $\mu_w, \mu_q > 0$. For simplicity, we consider adaptation with constant step sizes. Additionally, note that in (10.13) we adapt only the coefficients related to the i th span, since it can easily be shown that $\frac{\partial \hat{J}(\mathbf{w}_n, \mathbf{q}_n)}{\partial \mathbf{q}_n}$ is 0 for all the coefficients outside the span. Convergence properties of this scheme are analyzed in a number of previous works [155]. The overall algorithm is summarized in Algorithm 10.1. A standard way to initialize the coefficients of the spline is to consider:

$$q_{x,i} = q_{y,i}, \quad i = 1, \dots, Q, \quad (10.14)$$

such that $\varphi(s[n]) = s[n]$. Using this initialization criterion, the LMS filter can be considered as a special case of the SAF, where adaptation is restricted to the linear part, i.e. $\mu_q = 0$.

Algorithm 10.1 SAF: Summary of the SAF algorithm with first-order updates.

- 1: Initialize $\mathbf{w}_{-1} = \delta[n], \mathbf{q}_0$
 - 2: **for** $n = 0, 1, \dots$ **do**
 - 3: $s[n] = \mathbf{w}_n^T \mathbf{x}_n$
 - 4: $u = s[n]/\Delta x - \lfloor s[n]/\Delta x \rfloor$
 - 5: $i = \lfloor s[n]/\Delta x \rfloor + (Q - 1)/2$
 - 6: $y[n] = \mathbf{u}^T \mathbf{B} \mathbf{q}_{i,n}$
 - 7: $e[n] = d[n] - y[n]$
 - 8: $\mathbf{w}_{n+1} = \mathbf{w}_n + \mu_w e[n] \varphi'(s[n]) \mathbf{x}_n$
 - 9: $\mathbf{q}_{i,n+1} = \mathbf{q}_{i,n} + \mu_q e[n] \mathbf{B}^T \mathbf{u}$
 - 10: **end for**
-

10.3 Diffusion SAF

Consider a network model as in the previous chapters. At a generic time instant n , each agent receives some input/output data denoted by $(\mathbf{x}_n^{(k)}, d^{(k)}[n])$, where

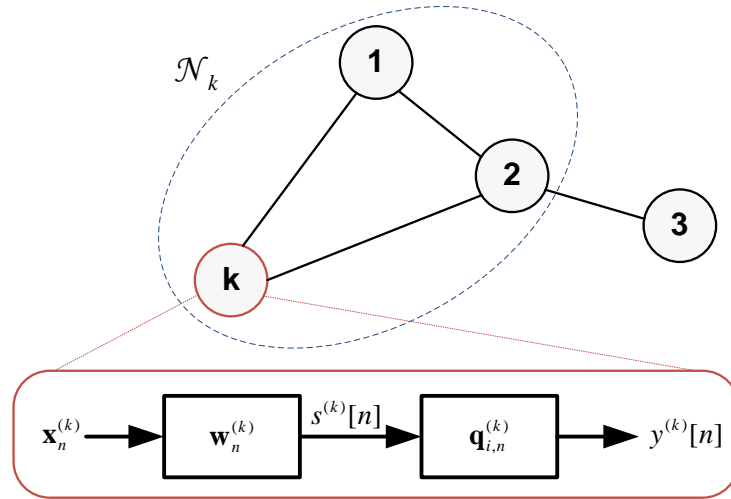


Figure 10.2: Schematic depiction of SAF interpolation performed over a network of agents. Each agent is connected to a neighborhood of other agents, and at every time instant it updates a local estimate of the optimal SAF model.

we introduce an additional superscript (k) for explicating the node dependence. We assume that streaming data at the local level is generated similarly to (10.1), according to:

$$d^{(k)}[n] = f_0(\mathbf{w}_0^T \mathbf{x}_n^{(k)}) + v^{(k)}[n]. \quad (10.15)$$

More in particular, we assume that \mathbf{w}_0 and $f_0(\cdot)$ are shared over the network, which is a reasonable assumption in many situations [26, 145]. Each node, however, receives input data with possibly different autocorrelation $\mathbf{R}_u^{(k)} = \mathbb{E}\{\mathbf{x}^{(k)T} \mathbf{x}^{(k)}\}$, and different additive noise terms $v^{(k)}[n] \sim \mathcal{N}(0, \sigma_k^2)$. Additionally, we assume that the nodes have agreed beforehand on a specific spline basis matrix \mathbf{B} , and on a set of initial control points \mathbf{q}_0 . Both quantities are common throughout the network. This is shown schematically in Fig. 10.2.

Given these assumptions, the network objective is to find the optimal SAF parameters (\mathbf{w}, \mathbf{q}) such that the following global cost function is minimized:

$$J_{\text{glob}}(\mathbf{w}, \mathbf{q}) = \sum_{k=1}^L J_{\text{loc}}^{(k)}(\mathbf{w}, \mathbf{q}) = \sum_{k=1}^L \mathbb{E}\{e^{(k)}[n]^2\}, \quad (10.16)$$

where each expectation is defined with respect to the local input statistics. Remember that the main idea of DA techniques is to interleave parallel adaptation steps with diffusion steps, where information on the current estimates are locally combined based on the mixing matrix \mathbf{C} (see for example [145, Section V-B]). Denote by $(\mathbf{w}_n^{(k)}, \mathbf{q}_n^{(k)})$ the SAF estimate of node k at time-instant n . In the diffusion SAF (D-SAF), each node starts by diffusing its own estimate of the linear part of the SAF

filter:

$$\boldsymbol{\psi}_n^{(k)} = \sum_{l \in \mathcal{N}_k} C_{kl} \mathbf{w}_n^{(k)}. \quad (10.17)$$

w-diffusion

Next, we can use the new weights $\boldsymbol{\psi}_n^{(k)}$ to compute the linear output of the filter as $s^{(k)}[n] = \boldsymbol{\psi}_n^{(k)T} \mathbf{x}_n^{(k)}$. From this, each node can identify its current span index i with (10.3) and (10.4). In the second phase, the nodes performs a second diffusion step over their span:

$$\boldsymbol{\xi}_{i,n}^{(k)} = \sum_{l \in \mathcal{N}_k} C_{kl} \mathbf{q}_{i,n}^{(k)}. \quad (10.18)$$

q-diffusion

Note that the q-diffusion step requires combination of the coefficients in the span $\mathbf{q}_{i,n}^{(k)}$, hence its complexity is independent on the number of control points in the spline, being defined only by the spline order P .

Once the nodes have diffused their information, they can proceed to a standard adaptation step as in the single-agent case. In particular, the spline output given the new span is obtained as:

$$\mathbf{y}^{(k)}[n] = \varphi_k(s^{(k)}[n]) = \mathbf{u}^T \mathbf{B} \boldsymbol{\xi}_{i,n}^{(k)}. \quad (10.19)$$

From this, the local error is given as $e^{(k)}[n] = d^{(k)}[n] - \mathbf{y}^{(k)}[n]$. The two gradient descent steps are then:

$$\mathbf{w}_{n+1}^{(k)} = \boldsymbol{\psi}_n^{(k)} + \mu_w^{(k)} e^{(k)}[n] \varphi'(s^{(k)}[n]) \mathbf{x}_n^{(k)}, \quad (10.20)$$

w-adapt

$$\mathbf{q}_{i,n+1}^{(k)} = \boldsymbol{\xi}_{i,n}^{(k)} + \mu_q^{(k)} e^{(k)}[n] \mathbf{B}^T \mathbf{u}. \quad (10.21)$$

q-adapt

where the two step sizes $\mu_w^{(k)}, \mu_q^{(k)}$ are possibly different across different agents. The overall algorithm is summarized in Algorithm 10.2. Note that in this chapter we consider a diffusion step prior to the adaptation step. In the DA literature, this is known as a combine-then-adapt (CTA) strategy [175]. This is true even if the two diffusion steps are not consecutive in Algorithm 10.2. In fact, Algorithm 10.2 is equivalent to the case where the full vector $\mathbf{q}_n^{(k)}$ is exchanged before selecting the proper span. Following similar reasonings, we can easily obtain an adapt-then-combine (ATC) strategy by inverting the two steps. Additionally, similarly to what we remarked in Section 10.2, we note that D-LMS [101] is a special case of the D-SAF, where each node initialize its nonlinearity with (10.14), and $\mu_q^{(k)} = 0, k = 1, \dots, L$.

Algorithm 10.2 D-SAF: Summary of the D-SAF algorithm (CTA version).

-
- 1: Initialize $\mathbf{w}_{-1}^{(k)} = \delta[n], \mathbf{q}_0^{(k)}$, for $k = 1, \dots, L$
 - 2: **for** $n = 0, 1, \dots$ **do**
 - 3: **for** $k = 1, \dots, L$ **do**
 - 4: $\boldsymbol{\psi}_n^{(k)} = \sum_{l \in \mathcal{N}_k} C_{kl} \mathbf{w}_n^{(k)}$
 - 5: $s^{(k)}[n] = \boldsymbol{\psi}_n^{(k)T} \mathbf{x}_n^{(k)}$
 - 6: $u = s^{(k)}[n]/\Delta x - \lfloor s^{(k)}[n]/\Delta x \rfloor$
 - 7: $i = \lfloor s^{(k)}[n]/\Delta x \rfloor + (Q - 1)/2$
 - 8: $\boldsymbol{\xi}_{i,n}^{(k)} = \sum_{l \in \mathcal{N}_k} C_{kl} \mathbf{q}_{i,n}^{(k)}$
 - 9: $y^{(k)}[n] = \mathbf{u}^T \mathbf{B} \boldsymbol{\xi}_{i,n}^{(k)}$
 - 10: $e^{(k)}[n] = d^{(k)}[n] - y^{(k)}[n]$
 - 11: $\mathbf{w}_{n+1}^{(k)} = \boldsymbol{\psi}_n^{(k)} + \mu_w^{(k)} e^{(k)}[n] \varphi'(s^{(k)}[n]) \mathbf{x}_n^{(k)}$
 - 12: $\mathbf{q}_{i,n+1}^{(k)} = \boldsymbol{\xi}_{i,n}^{(k)} + \mu_k^{(k)} e^{(k)}[n] \mathbf{B}^T \mathbf{u}$
 - 13: **end for**
 - 14: **end for**
-

10.4 Experimental Setup

To test the proposed D-SAF, we consider network topologies with L agents, whose connectivity is generated randomly, such that every pair of nodes has a 60% probability of being connected. To provide sufficient diversity, we experiment with a small network with $L = 10$ and a larger network with $L = 30$. Data is generated according to the Wiener model in (10.15), where the optimal weights \mathbf{w}_0 are extracted randomly from a normal distribution, while the nonlinearity $f_0(\cdot)$ for the initial experiments is depicted in Fig. 10.3. This represents a mild nonlinearity.

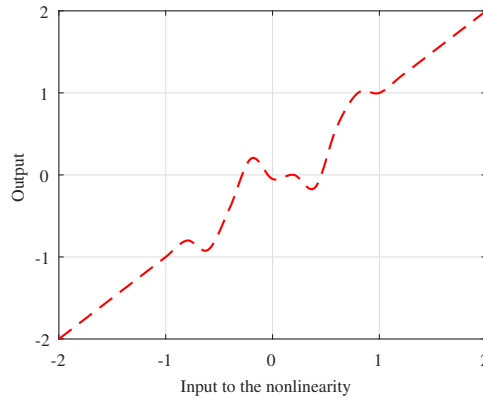


Figure 10.3: Nonlinear distortion applied to the output signal in experiments 1 and 2.

The input signal at each node is generated following the experiments in [155],

and it consists of 25000 samples generated according to:

$$x_k[n] = a_k x_k[n-1] + \sqrt{1 - a_k^2} \xi[n], \quad (10.22)$$

where the correlation coefficients a_k are assigned randomly at every node from an uniform probability distribution in $[0, 0.8]$, while $\xi[n]$ is a white Gaussian noise term with zero mean and unitary variance. The desired signal is then given by (10.15), where the noise variances $\sigma^{(k)}[n]$ at every node are assigned randomly in $[-10, -25]$ dB. The mixing coefficients are chosen according to the ‘metropolis’ strategy as in previous chapters. In all experiments, knots are equispaced in $[-2, +2]$ with $\Delta x = 0.2$.

We compare D-SAF with a non-cooperative SAF (denoted as NC-SAF), which corresponds in choosing a diagonal mixing matrix $\mathbf{C} = \mathbf{I}$. Similarly, we compare with the standard D-LMS [101], and a non-cooperative LMS, denoted as NC-LMS. To average out statistical effects, experiments are repeated 15 times, by keeping fixed the topology of the network and the optimal parameters of the system. Results are then averaged throughout the nodes.

10.5 Experimental Results

10.5.1 Experiment 1 - Small Network ($L = 10$)

In the first experiment, we consider a network with $L = 10$. Details on the signal generation are provided in Fig. 10.4. In particular, the local correlation coefficients are shown in Fig. 10.4a, and the amount of noise variance in Fig. 10.4b. The step-sizes are instead given in Fig. 10.4c. These settings allow a certain amount of variety on the network. As an example, input values at node 3 are highly correlated, while node 2 has the strongest amount of noise. Similarly, speed of adaptation (and consequently steady-state convergence) covers a large range of settings, as depicted in Fig. 10.4c. The first measure of error that we consider is the mean-squared error (MSE), defined in dB as:

$$\text{MSE}_k[n] = 10 \log (d_k[n] - y_k[n])^2 . \quad (10.23)$$

Results in term of MSE are given in Fig. 10.5, where the proposed algorithm is shown with a solid violet line. Here and in the following figures, the MSE is computed by averaging (10.23) over the different nodes.

As expected, due to the nonlinear distortion, LMS achieves a generally poor performance, with a steady-state MSE of -12 dB. Additionally, there is almost no improvement when considering D-LMS compared to NC-LMS. The SAF filters are instead able to approximate extremely well the desired system. The diffusion

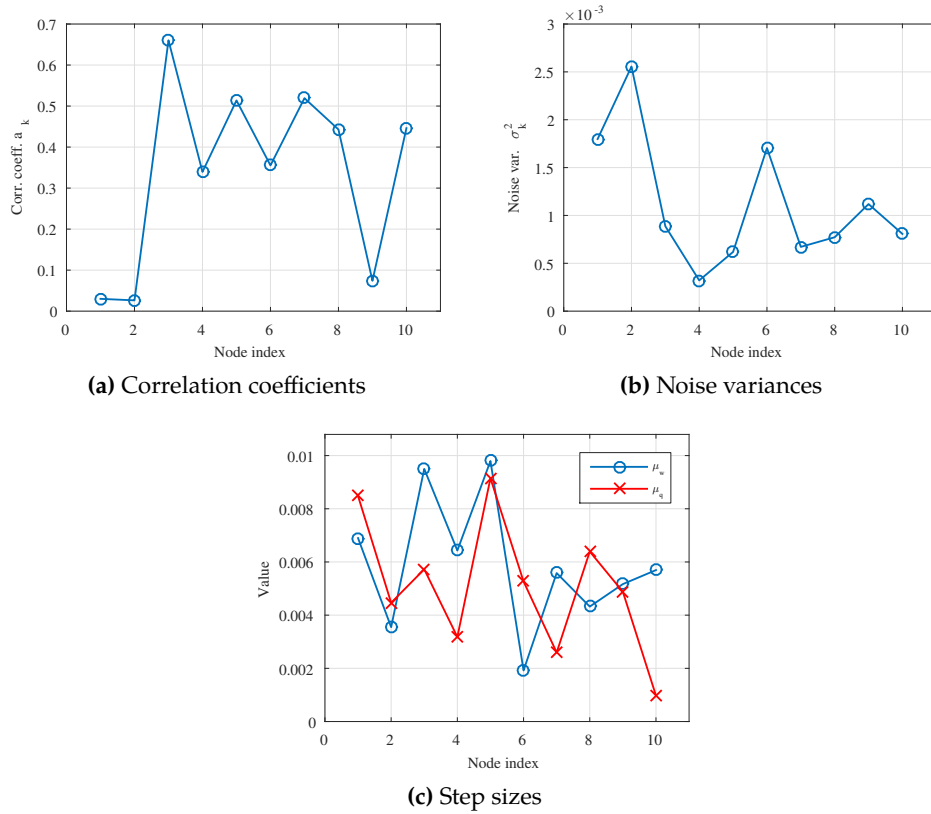


Figure 10.4: Dataset setup for experiment 1. (a) Correlation coefficients in (10.22); (b) Noise variances in (10.15); (c) Step sizes.

strategy, however, provides a significant improvement in convergence time with respect to the non-cooperative version, as is evident from Fig. 10.5. Further clarifications on the two algorithms can be obtained by considering the linear mean-squared deviation (MSD) given by:

$$\text{MSD}_k^l = 10 \log \left(\left\| \mathbf{w}_0 - \mathbf{w}_n^{(k)} \right\|_2 \right), \quad (10.24)$$

and the nonlinear MSD given by:

$$\text{MSD}_k^{\text{nl}} = 10 \log \left(\left\| \mathbf{q}_0 - \mathbf{q}_n^{(k)} \right\|_2 \right). \quad (10.25)$$

The overall behavior of the MSD is shown in Fig. 10.6. In particular, we show in Fig. 10.6a and Fig. 10.6b the global MSD of the network, which is obtained by averaging the local MSDs at every node.

It can be seen that the MSD achieved with a diffusion algorithm strongly outperform the average MSD obtained with a non-cooperative solution. Additionally, the gap in the linear and nonlinear case is similar, with a steady-state difference of 4 dB. The reason of this difference is shown in Fig. 10.6c and Fig. 10.6d, where we plot the MSD evolution for D-SAF and for three representative agents running

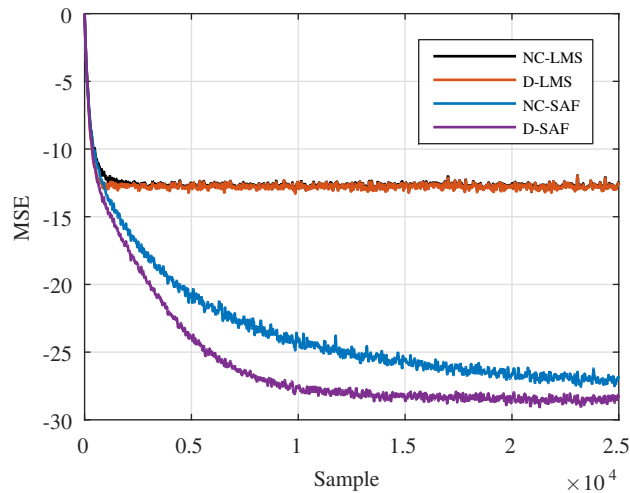


Figure 10.5: MSE evolution for experiment 1, averaged across the nodes.

NC-SAF. It can be seen that, due to the differences in configuration, some nodes have a much slower convergence than other, such as node 6 compared to node 1. However, these statistical variations are successfully averaged out by the diffusion algorithm, which is able to outperform even the fastest node in the network. This is shown visually in Fig. 10.7, where we show the resulting nonlinear models for three representative nodes running NC-SAF, and for the nodes running D-SAF.

10.5.2 Experiment 2 - Large Network ($L = 30$)

For the second set of experiments, we consider a larger network with $L = 30$ agents. Everything else is kept fixed as before, in particular, each pair of nodes in the network has a 60% probability of being connected, with the only requirement that the global network is connected. The local correlation coefficients in (10.22), noise variances in (10.15), and local step-sizes are extracted randomly at every node using the same settings as the previous experiment. In this case, this setting provides a larger range of configurations for the different nodes. The results in term of MSE evolution are shown in Fig. 10.8, where the proposed D-SAF is again shown with a violet line.

While the performance of NC-LMS and D-LMS are similar to those exhibited in the previous experiment, it is interesting to observe that, by increasing the amount of nodes in the network, the convergence of NC-SAF is slower in this case, to the point that the algorithm is not able to converge efficiently in the provided number of samples. D-SAF, instead, is robust to this increase in network's size, and it is able to reach almost complete convergence in less than 15000 samples. Clearly, this is expected from the behavior of the algorithm. The larger the network, the higher the amounts of neighbors each single agent has. Thus, the diffusion steps are able to fuse more information, providing a faster convergence, as also evidenced by

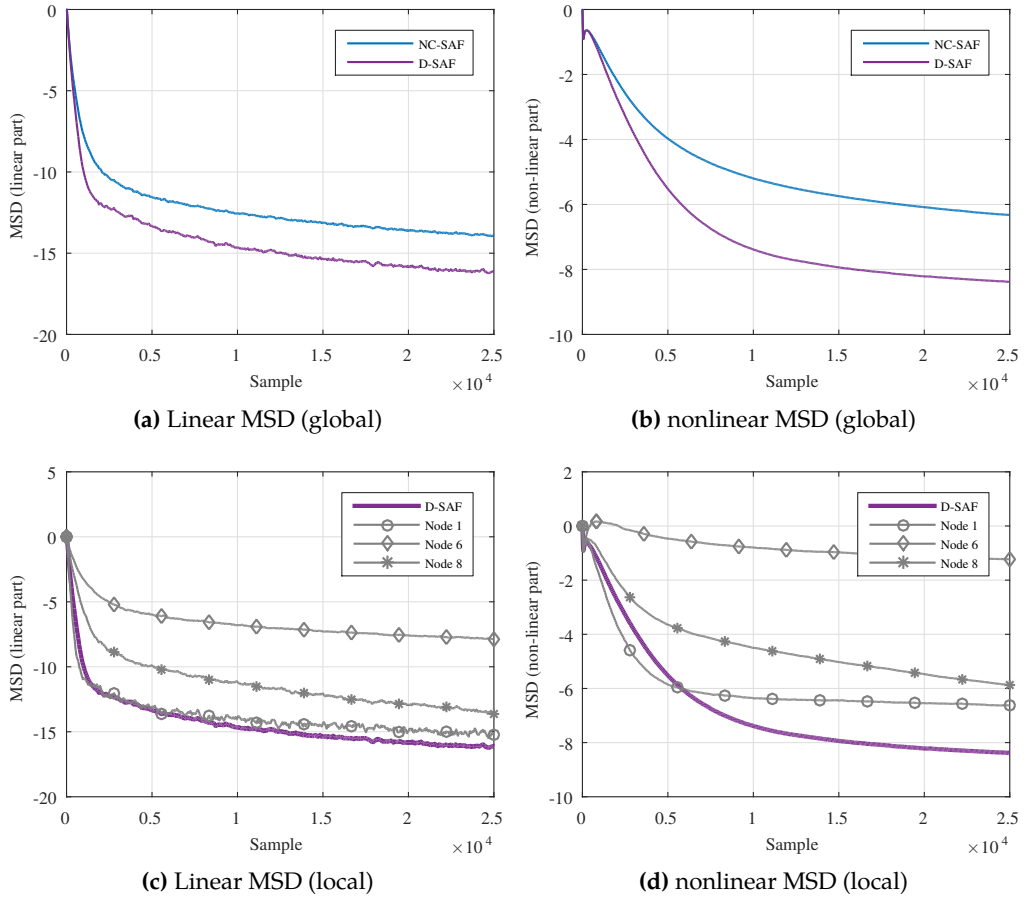


Figure 10.6: MSD evolution for experiment 1. (a-b) Global MSD evolution. (c-d) MSD evolution for D-SAF and 3 representative nodes running NC-SAF.

previous literature in DA [145]. Due to this, the algorithm is able to average out the performance of isolated nodes, where convergence is not achieved. This can be seen from Fig. 10.9, where we plot the splines obtained from 3 representative nodes running NC-SAF in Fig. 10.9a, and the spline resulting from D-SAF in Fig. 10.9b. In Fig. 10.9a, it is possible to see that some nodes achieve perfect convergence, while others would require a larger amount of samples. Even worse, some nodes are actually diverging from the optimal solution, due to their peculiar configuration. Despite this, D-SAF is converging globally to an optimal solution, as shown by the black line in Fig. 10.9b.

10.5.3 Experiment 3 - Strong nonlinearity ($L = 15$)

As a final validation, we consider an intermediate network with $L = 15$, but we change the output nonlinearity in (10.22) with the one showed in Fig. 10.11. This is a stronger nonlinearity, with two larger peaks. As before, the correlation coefficients in (10.22), the variances of the noise, and the local step-sizes are assigned randomly

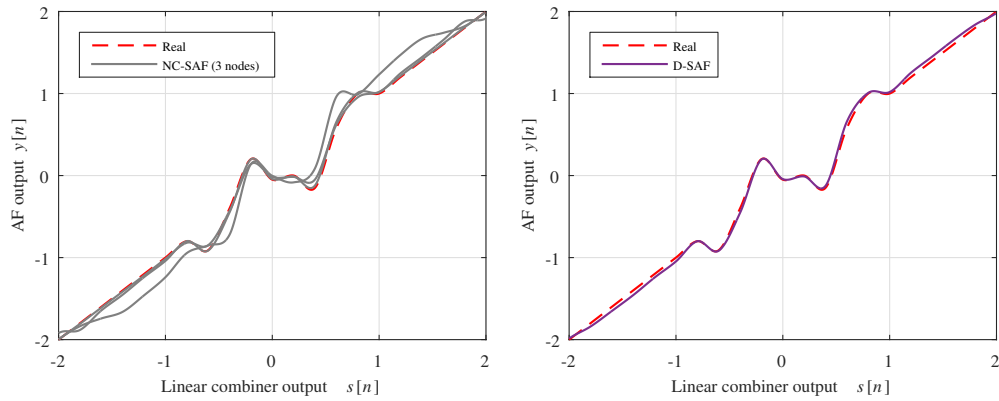


Figure 10.7: Final estimation of the nonlinear model in experiment 1. (a) Three representative nodes running NC-SAF. (b) Final spline of the nodes running D-SAF.

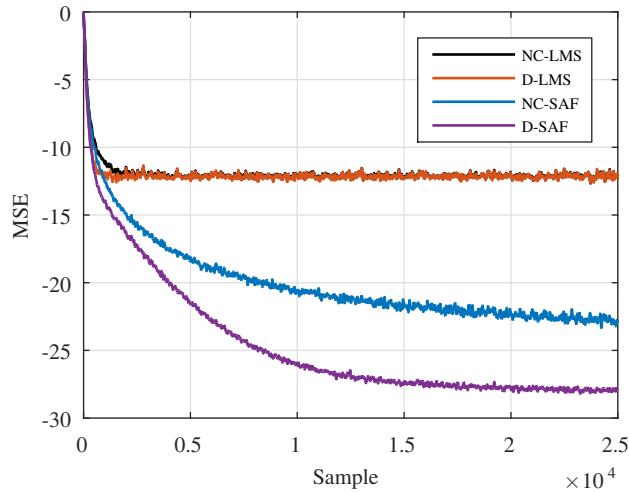


Figure 10.8: MSE evolution for experiment 2, averaged across the nodes.

at every node. Results of the experiment are shown in Fig. 10.10. Due to the increased nonlinearity, the two linear filters are performing poorly, with a steady-state MSE of -5 dB. Convergence is also slowed for NC-SAF, while there is now a gap of almost 10 dB between the final MSE of the cooperative and non-cooperative versions of SAF. It is particularly interesting to observe the final nonlinearities at every node. These are shown for three representative nodes running NC-SAF in Fig. 10.11a, and for D-SAF in Fig. 10.11b. Overall, a small portion of nodes running NC-SAF is achieving a satisfactory convergence, while several of them are only achieving a moderate convergence, or no convergence at all. Despite this, and the smaller size of the network with respect to the second experiment, D-SAF is still obtaining an almost complete convergence at a global level. Overall, this shows that the algorithm is robust to variations in the network's size, local configurations, and amount of nonlinearity.

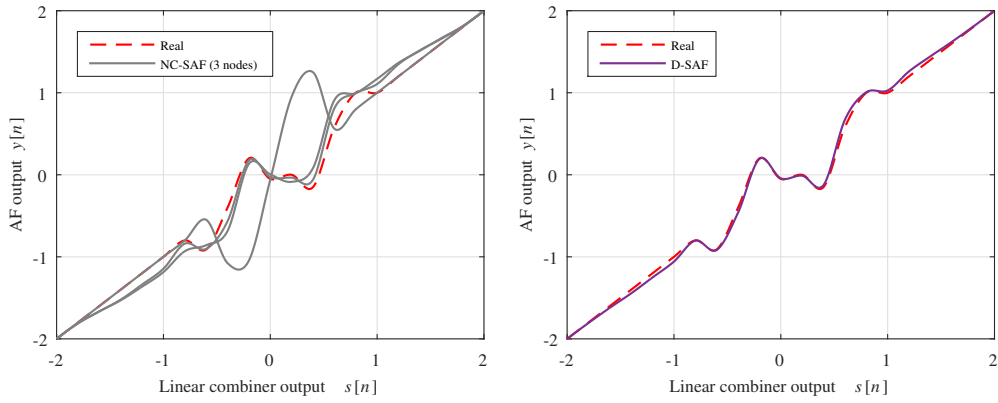


Figure 10.9: Final estimation of the nonlinear model in experiment 2. (a) Three representative nodes running NC-SAF. (b) Final spline of the nodes running D-SAF.

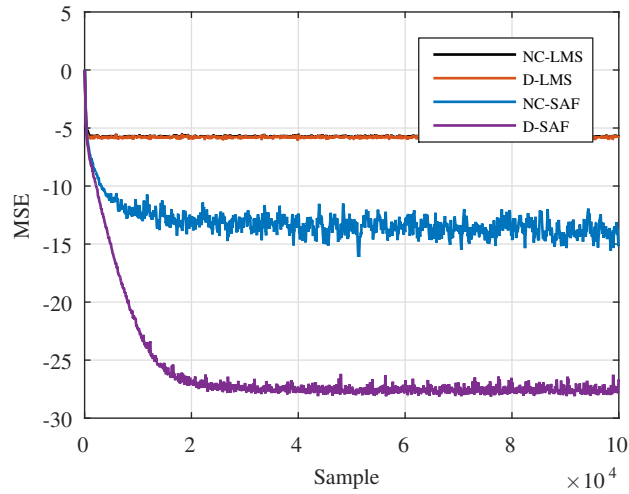


Figure 10.10: MSE evolution for experiment 3, averaged across the nodes.

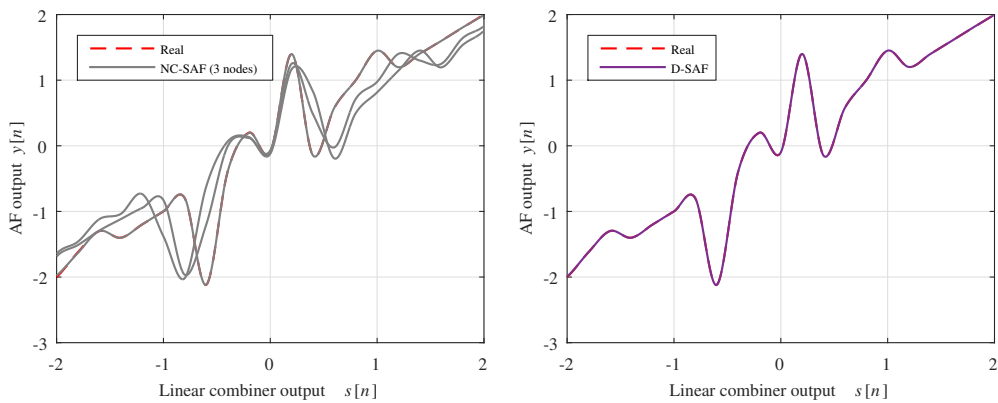


Figure 10.11: Final estimation of the nonlinear model in experiment 3. (a) Three representative nodes running NC-SAF. (b) Final spline of the nodes running D-SAF.

Part V

Conclusions and Future Works

Conclusions and Future Works

DISTRIBUTED learning has received considerable attention over the past years due to its broad real-world applications. It is common nowadays that data must be collected, stored locally and data exchange is not allowed for specific reasons, such as technological bottlenecks or privacy concerns. In such a circumstance, it is necessary and useful to build in a decentralized fashion an ANN model. Motivated by this, throughout this thesis we have put forth multiple algorithms to such end.

Initially, we have explored extensions to the DL setting of the well-known RVFL network (Chapters 4-6). As in the centralized case, distributed RVFL networks are able to provide strong nonlinear modeling capabilities, while at the same time allowing for a fast and simple set of training algorithms, which are fundamentally framed in the linear regression literature. Thus, they provide a good compromise between a linear model and more complex nonlinear ANNs, such as distributed SVMs [58].

The successive chapters have considered the more complex problem of distributed training in the presence of labeled and unlabeled data, thus extending the theory of SSL [31]. This is a relatively new problem in the literature, with a large set of possible real world applications. In this sense, the two distributed models explored in Chapters 7 and 8 are only initial explorations of a field which can potentially reveal much promise.

Finally, in the last part of the thesis we have been concerned with learning from time-varying data. Although this is a well-known setting, both of the algorithms that we presented are relatively novel. Indeed, Chapter 9 has introduced one of the first available algorithms for training recurrent networks, while the diffusion SAF in Chapter 10 can be seen as a general nonlinear extension of the much celebrated D-LMS [145].

Below we provide a set of possible future lines of research, which refer to specific portions of the thesis, along with the main content of each chapter.

- In **Chapters 4-6** we have detailed distributed algorithms for learning a RVFL network, in the case of batch and online learning, both for HP and VP partitioned data. In Chapter 5, in particular, we have focused on the applica-

tion to multiple distributed music classification tasks, including genre and artist recognition. These problems arise frequently in real-world scenarios, including P2P and mobile networks. Our experimental results show that the proposed algorithms can be efficiently applied in these situations, and compares favorably with a centralized solution in terms of accuracy and speed. Clearly, the algorithms can be successfully applied to distributed learning problems laying outside this specific applicative domain, particularly in real-world big data scenarios. Moreover, although in Chapter 5 we have focused on local updates based on the BRLS algorithm, nothing prevents the framework from being used with different rules, including efficient stochastic gradient descent updates. Similar considerations also apply for Chapter 4 and Chapter 6.

- In **Chapter 7** we have proposed a totally decentralized algorithm for SSL in the framework of MR. The core of our proposal is constituted by a distributed protocol designed to compute the Laplacian matrix. Our experimental results show that, also in this case, the proposed algorithm is able to match efficiently the performance of a centralized model built on the overall training set. Although we have focused on a particular algorithm belonging to MR, namely LapKRR, the framework is easily applicable to additional algorithms, including the laplacian SVM (LapSVM) [11], and others. Moreover, extensions beyond MR are possible, i.e. to all the methods that encode information in the form of a matrix of pairwise distances, such as spectral dimensionality reduction, spectral clustering, and so on. In the case of kernels that directly depend on the dot product between patterns (e.g. the polynomial one), particular care must be taken in designing appropriate privacy-preserving protocols for distributed margin computation [165], an aspect which is left to future investigations. Currently, the main limit of our algorithm is the computation time required by the distributed algorithm for completing the Laplacian matrix. This is due to a basic implementation of the two optimization algorithms. In this sense, in future works we intend to improve the distributed algorithm to achieve better computational performance. Examples of possible modifications include adaptive strategies for the choice of the step-size, as well as early stopping protocols.
- Next, in **Chapter 8** we have solved the problem of distributed SSL via another type of semi-supervised SVM, framed in the transductive literature. Particularly, we have leveraged over recent advances on distributed non-convex optimization, in order to provide two flexible mechanisms with a different balance in computational requirements and speed of convergence. A natural extension would be to consider different semi-supervised techniques

to be extended to the distributed setting, particularly among those developed for the S^3VM [30].

- In **Chapter 9** we have introduced a decentralized algorithm for training an ESN. Experimental results on multiple benchmarks, related to non-linear system identification and chaotic time-series prediction, demonstrated that it is able to efficiently track a purely centralized solution, while at the same time imposing a small computational overhead in terms of vector–matrix operations requested to the single node. This represents a first step towards the development of data-distributed strategies for general RNNs, which would provide invaluable tools in real world applications. Future lines of research involve considering different optimization procedures with respect to ADMM, or more flexible DAC procedures. More in general, it is possible to consider other distributed training criteria beyond ridge regression and LASSO (such as training via a support vector algorithm) to be implemented in a distributed fashion. Finally, ESN are known to perform worse for problems that require a long memory. In this case, it is necessary to devise distributed strategies for other classes of recurrent networks, such as LSTM architectures [71, 113].
- Finally, in **Chapter 10** we have investigated a distributed algorithm for adapting a particular class of nonlinear filters, called SAF, using the general framework of DA. The algorithm inherits the properties of SAFs in the centralized case, namely, it allows for a flexible nonlinear estimation of the underlying function, with a relatively small increase in computational complexity. In particular, the algorithm can be implemented with two diffusion steps, and two gradient descent steps, thus requiring in average only twice as much computations as the standard D-LMS. Our experimental results show that D-SAF is able to efficiently learn hard nonlinearities, with a definite increase in convergence time with respect to a non-cooperative implementation. In the respective chapter, we have focused on a first-order adaptation algorithm, with CTA combiners. In future works, we plan to extend the D-SAF algorithm to the case of second-order adaptation with Hessian information, ATC combiners, and asynchronous networks. Additionally, we plan to investigate diffusion protocols for more general architectures, including Hammerstein and IIR spline filters.

A few general considerations on the thesis are also worth mentioning here:

- **Fixed topology:** for simplicity, in this thesis we have supposed that the network of agents is fixed, and connectivity is known at the agent level. This is not a necessary condition (indeed, many practical applications might require

time-varying connectivities), and work along this sense is planned in the near future. Indeed, many of the tools employed throughout the thesis, e.g. ADMM, already possess extensions to this scenario, which can in principle be applied to the problems considered here.

- **Synchronization:** similar considerations apply for the problem of having synchronized updates, requiring a general mechanism for coordinating the agents. Investigations along this line can start by considering asynchronous versions of SGD [120], or DA [205].
- **Specific ML fields:** another important aspect is that, similarly to SSL, many subfields of ML remain to be extended to the distributed setting. As an example, there is limited literature for distributed implementation of *active learning* strategies [180], where agents are allowed to request a set of labels on items that they assume to be interesting. This can potentially reduce drastically the training time and the amount of communication overhead.
- **Multilayer networks:** we have investigated distributed methods only for ANN having at most one hidden layer of nonlinearities, which are known as ‘shallow’ in the current terminology [158]. Indeed, we saw in Section 3.4.7 that investigations on distributed deep neural networks have been limited. This is due to the large number of parameters to be exchanged, and to the resulting non-convex optimization problem. Both these problems require additional investigations in order to be properly addressed.
- **Additional distributed techniques:** finally, we expect that techniques originally developed for distributed signal processing and distributed AI might be applied to the problem of DL, resulting in beneficial effects in term of in-network communication and/or computational requirements. This is the case, for example, of message censoring [176], a set of techniques allowing each individual agent to decide whether to take a specific measurement and propagate it over the network.

Appendices

A

Elements of Graph Theory

A.1 Algebraic graph theory

Consider a graph composed by L nodes, whose connectivity is fixed and known in advance. Mathematically, this graph can be represented by the so-called adjacency matrix $\mathbf{A} \in \{0, 1\}^{L \times L}$, defined as:

$$A_{ij} = \begin{cases} 1 & \text{if node } i \text{ is connected to node } j \\ 0 & \text{otherwise} \end{cases} . \quad (\text{A.1})$$

The symbol \mathcal{N}_k denotes the exclusive neighborhood of node k , i.e. the set of nodes connected to k , with the exclusion of k itself. In this thesis, we are concerned with graphs which are undirected, meaning that \mathbf{A} is symmetric, and connected, meaning that each node can be reached from every other node in a finite sequence of steps. Additionally, we suppose that there are no self-loops. We can define the degree d_k of node k as the number of nodes which are connected to it:

$$d_k = |\mathcal{N}_k| = \sum_{l=1}^L A_{kl} . \quad (\text{A.2})$$

The degree d of the network is defined as the maximum degree of its composing nodes:

$$d = \max_{k=1, \dots, L} d_k . \quad (\text{A.3})$$

The degree matrix $\mathbf{D} \in \mathbb{N}^{L \times L}$ is then defined as:

$$\mathbf{D} = \text{diag} \{d_1, \dots, d_L\} , \quad (\text{A.4})$$

where $\text{diag} \{\cdot\}$ constructs a diagonal matrix from its arguments. Lastly, the Laplacian matrix $\mathbf{L} \in \mathbb{Z}^{L \times L}$ is defined as:

$$\mathbf{L} = \mathbf{D} - \mathbf{A} . \quad (\text{A.5})$$

From the previous definitions, we obtain:

$$L_{ij} = \begin{cases} d_i & \text{if } i = j \\ -1 & \text{if } i \in \mathcal{N}_j \\ 0 & \text{otherwise} \end{cases} . \quad (\text{A.6})$$

It is known that an analysis of the Laplacian matrix allows to derive multiple important properties of the underlying graph. As an example, it can be shown that $\lambda_0(\mathbf{L}) = 0$, while the second-smallest eigenvalue is directly related to the connectivity of the graph itself [117]. In Chapter 7 we make use of a variant of \mathbf{L} , called the normalized Laplacian matrix and defined as:

$$\mathbf{L}' = \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}} . \quad (\text{A.7})$$

It follows straightforwardly that:

$$L'_{ij} = \begin{cases} 1 & \text{if } i = j \\ -\frac{1}{\sqrt{d_i d_j}} & \text{if } i \in \mathcal{N}_j \\ 0 & \text{otherwise} \end{cases} . \quad (\text{A.8})$$

A.2 Decentralized average consensus

Suppose now that the nodes in the graph represent agents in a physical network. Additionally, each of them has access to a measurement vector $\mathbf{m}_k \in \mathbb{R}^S$. The task is for each of them to compute the global average given by:

$$\hat{\mathbf{m}} = \frac{1}{L} \sum_{k=1}^L \mathbf{m}_k . \quad (\text{A.9})$$

For generality, however, we allow every node to communicate only with its direct neighbors. With respect to the categorization of Section 3.2, this is denoted as one-hop communication. DAC is an iterative network protocol to compute the global average (or, equivalently, sum) with respect to the local measurement vectors, requiring only local communications between them [9, 119, 199]. Its simplicity makes it suitable for implementation even in the most basic networks, such as robot swarms. Each agent initializes its estimate of the global average as $\mathbf{m}_k[0] = \mathbf{m}_k$. Then, at a generic iteration n , the local DAC update is given by:

$$\mathbf{m}_k[n] = \sum_{j=1}^L C_{kj} \mathbf{m}_j[n-1] , \quad (\text{A.10})$$

where the weight C_{kj} is a real-valued scalar denoting the confidence that the k th node has with respect to the information coming from the j th node. By reorganizing these weights in a $L \times L$ connectivity matrix \mathbf{C} , and defining:

$$\mathbf{M}[n] = [\mathbf{m}_1[n] \dots \mathbf{m}_L[n]] , \quad (\text{A.11})$$

Eq. (A.10) can be rewritten more compactly as:

$$\mathbf{M}[n] = \mathbf{C}\mathbf{M}[n - 1]. \quad (\text{A.12})$$

If the weights of the connectivity matrix \mathbf{C} are chosen appropriately, this recursive procedure converges to the global average given by Eq. (A.9) [119]:

$$\lim_{n \rightarrow +\infty} \mathbf{m}_k[n] = \frac{1}{L} \sum_{k=1}^L \mathbf{m}_k[0] = \hat{\mathbf{m}}, \quad k = 1, 2, \dots, L. \quad (\text{A.13})$$

Practically, the procedure can be stopped after a certain predefined number of iterations is reached, or when the norm of the update is smaller than a certain user-defined threshold δ :

$$\|\mathbf{m}_k[n] - \mathbf{m}_k[n - 1]\|_2^2 < \delta, \quad k = 1, 2, \dots, L. \quad (\text{A.14})$$

In the case of undirected, connected networks, convergence is guaranteed provided that the connectivity matrix \mathbf{C} respects the following properties:

$$\mathbf{C} \cdot \mathbf{1} = \mathbf{1}, \quad (\text{A.15})$$

$$\rho \left(\mathbf{C} - \frac{\mathbf{1} \cdot \mathbf{1}^T}{L} \right) < 1. \quad (\text{A.16})$$

A simple way of ensuring this is given by choosing the so-called ‘max-degree’ weights [119]:

$$C_{kj} = \begin{cases} \frac{1}{d+1} & \text{if } k \text{ is connected to } j \\ 1 - \frac{d_k}{d+1} & \text{if } k = j \\ 0 & \text{otherwise} \end{cases}, \quad (\text{A.17})$$

In practice, many variations on this standard procedure can be implemented to increase the convergence rate, such as the ‘definite consensus’ [61], or the strategy introduced in [144]. In this thesis, Eq. (A.17) is used for choosing \mathbf{C} unless otherwise specified. Other strategies are explored in Section 5.3.

B

Software Libraries

In this appendix, we present the open-source software libraries developed during the course of the PhD. The libraries can be used to replicate most of the experiments and simulations presented in the previous chapters. All of them are implemented in the MATLAB environment.

B.1 Lynx MATLAB Toolbox (Chapters 4-6)

Lynx is a research-oriented MATLAB toolbox, to provide a simple environment for performing large-scale comparisons of SL algorithms.¹ Basically, the details of a comparison, in terms of algorithms, datasets, etc., can be specified in a human understandable configuration file, which is loaded at runtime by the toolbox. In this way, it is possible to abstract the elements of the simulation from the actual code, and to repeat easily previously defined experiments. An example of configuration file is provided below, where two different algorithms (an SVM and a RVFL network) are compared on a well-known UCI benchmark.

Listing B.1: Demo code for the Lynx Toolbox

```
1 % Add the SVM to the simulation
2 add_model('SVM', 'Support Vector Machine', @SupportVectorMachine);
3
4 % Add the RVFL network
5 add_model('RVFL', 'Random Vector Network', @RandomVectorFunctionallink);
6
7 % Add a dataset
8 add_dataset('Y', 'Yacht', 'uci_yacht');
```

Inside the toolbox, we implemented a set of utilities in order to simplify the development of distributed algorithms. Specifically, the toolbox allows to develop additional ‘features’, which are objects that perform specific actions during the course of a simulation. We implemented a `InitializeTopology()` feature, which takes care of partitioning the training data evenly across a network of agents, and

¹<https://github.com/ispamm/Lynx-Toolbox>

provides the algorithms with a set of network specific functions, such as DAC protocols. Below is an example of enabling this feature in a configuration file:

Listing B.2: Example of data partitioning

```
1 % Define a network topology (in this case, a randomly generated one with 20
  agents)
2 r = RandomTopology(20, 0.2);
3
4 % Initialize the feature (with a set of user-specified flags)
5 feat = InitializeTopology(r, 'disable_parallel', 'distribute_data', '
  disable_plot');
6
7 % Add the feature to the simulation
8 add_feature(feat);
```

Due to the way in which the toolbox is structured, it is possible to combine distributed and non-distributed algorithms in the same simulation, leaving to the software the task of choosing whether to partition or not the data, and to collect the results from the different agents in the former case. The configuration files for Chapter 4 and Chapter 6 are available on the author's website,² together with a set of additional examples of usage.

B.2 Additional software implementations

B.2.1 Distributed LapKRR (Chapter 7)

The code for this chapter is available on BitBucket.³ The network utilities (e.g. random graph generation) are adapted from the Lynx toolbox (see previous section). Each set of experiments can be repeated by running the corresponding script in the 'Scripts' folder. Specifically, there are three scripts, which can be used to replicate the experiments on EDM completion, distributed SSL, and privacy preservation, respectively.

B.2.2 Distributed S³VM (Chapter 8)

Similarly to the previous chapter, the code has been released on BitBucket,⁴ with some adaptations from the Lynx toolbox in terms of network utilities. With respect to the library for distributed LapKRR, the code has been designed in a more flexible fashion, as it allows to define a variable number of centralized and distributed algorithms to be compared in the `test_script.m` file:

²<http://ispac.diet.uniroma1.it/scardapane/software/code/>.

³<https://bitbucket.org/robertofierimonte/distributed-semisupervised-code/>

⁴<https://bitbucket.org/robertofierimonte/code-distributed-s3vm>

Listing B.3: Definition of semi-supervised distributed algorithms

```

1 % Define a centralized RBF and a centralized Nabla-S3VM
2 centralized_algorithms = {...
3     CentralizedSVM('RBF-SVM', 1, '2', ''), ...
4     NablaS3VM('NS3VM (GD)', 1, 1, 5)
5 };
6
7 % Define a distributed Nabla-S3VM
8 distributed_algorithms = {...
9     DiffusionNablaS3VM('D-NS3VM', 1, 1, 5)
10 };

```

New algorithms can be defined by extending the abstract class `LearningAlgorithm`.

B.2.3 Distributed ESN (Chapter 9)

This code is released on a different package on BitBucket,⁵ following the general ideas detailed above. Specifically, configuration and execution are divided in two different scripts, which can be easily customized. The ESN is implemented using a set of functions adapted from the Simple ESN toolbox developed by the research group of H. Jaeger.⁶

B.2.4 Diffusion Spline Filtering (Chapter 10)

This package has been developed in order to provide an effective testing ground for distributed filtering applications.⁷ As for the previous libraries, it is possible to declare dynamically new distributed filters to be tested, as long as they derive correctly from the base abstract class `DiffusionFilter`.

⁵<https://bitbucket.org/ispamm/distributed-esn>

⁶<http://organic.elis.ugent.be/node/129>

⁷<https://bitbucket.org/ispamm/diffusion-spline-filtering>

Acknowledgments

Throughout these years, I have had the pleasure of working with a large number of people, all of whom have taught me something. I am indebted to all of them, and this thesis, in its smallness, is dedicated to them. Needless to say, this thesis is also dedicated to those that have been the closest to me: family and loved ones.

To start with, I would like to express my gratitude to my supervisor, Prof. Aurelio Uncini, for its never-ending support. The same gratitude also extends to current and past members of the ISPAMM research group, including (in strict alphabetical order): Andrea Alemanno, Danilo Communiello, Francesca Ortolani, Prof. Raffaele Parisi, and Michele Scarpiniti. It has been a real pleasure working with all of you so far.

I am indebted to Prof. Massimo Panella for giving me the opportunity of spending a few months in this strange land which is Australia. My most sincere gratitude goes to my Australian supervisor, Prof. Dianhui Wang, for his warmth and inspiration. Another thanks goes to anyone, and particularly his students, who has welcomed me and cheered me during my stay there.

During my PhD program, I have had the possibility of collaborating with many researchers from my department. Thus, I would like to thank, in no particular order, Prof. Antonello Rizzi, Filippo Bianchi, Marta Bucciarelli, Prof. Fabiola Colone, Andrea Proietti, Paolo Di Lorenzo, Luca Liparulo and Rosa Altilio. Most of all, I thank Prof. Sergio Barbarossa for introducing me with enthusiasm to many topics on distributed optimization, that are used extensively here.

Thanks to the students that I have had the pleasure of partially supervising, including Roberto Fierimonte, Valentina Ciccarelli, Marco Biagi and Gabriele Medeot.

Thanks to Prof. Giandomenico Boffi and Prof. Carlo Cirotto for having organized the wonderful SEFIR schools in Perugia, and to Prof. Giovanni Iacovitti and Prof. Giulio Iannello for allowing me to participate there twice. Of course, I thank all those that I have met there, including in particular Prof. Flavio Keller.

Before starting my PhD program, I had a fantastic yearlong work experience. Many

of the colleagues I met there have remained my friends, and among them, I would like to particularly thank Fernando Nigro, Alessandra Piccolo and Ilaria Piccolo.

I thank Prof. Asim Roy from Arizona State University and Prof. Plamen Angelov from Lancaster University for giving me a great opportunity to see close-hand how an international conference is organized. Even if I could not participate in the end, I have had a great experience, and I thank my 'colleagues' Bill Howell, Teng Teck-Hou, and José Iglesias.

Thanks are also in order to Prof. Amir Hussain of Stirling University. As an honorary fellow of his research group, I hope in a long and fruitful collaboration in the future.

Thanks to Antonella Blasetti for giving me space to talk about machine learning in her fantastic events as head of the Lazio-Abruzzo Google Developer Group (LAB-GDG). A general thanks goes to all the members of the LAB-GDG.

A few final thanks to Prof. Stefano Squartini, for his insightful discussions, to Dr. Steven van Vaerenbergh for his advice and help on the Lynx MATLAB toolbox, and to the awesome staff at the 2014 IEEE International Conference on Acoustics, Speech, and Signal Processing.

References

- [1] M. M. Adankon, M. Cheriet, and A. Biem, "Semisupervised least squares support vector machine," *IEEE Transactions on Neural Networks*, vol. 20, no. 12, pp. 1858–1870, 2009.
- [2] R. Albert and A.-L. Barabási, "Statistical mechanics of complex networks," *Reviews of modern physics*, vol. 74, no. 1, p. 47, 2002.
- [3] A. Y. Alfakih, A. Khandani, and H. Wolkowicz, "Solving Euclidean distance matrix completion problems via semidefinite programming," *Computational optimization and applications*, vol. 12, no. 1-3, pp. 13–30, 1999.
- [4] H. H. Ang, V. Gopalkrishnan, S. C. H. Hoi, and W. K. Ng, "Cascade RSVM in peer-to-peer networks," in *Machine Learning and Knowledge Discovery in Databases*. Springer, 2008, pp. 55–70.
- [5] H. H. Ang, V. Gopalkrishnan, W. K. Ng, and S. C. H. Hoi, "Communication-efficient classification in P2P networks," in *Machine Learning and Knowledge Discovery in Databases*. Springer, 2009, pp. 83–98.
- [6] R. Arablouei, S. Werner, and K. Dogancay, "Diffusion-based distributed adaptive estimation utilizing gradient-descent total least-squares," in *Proceedings of the 2013 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP'13)*. IEEE, 2013, pp. 5308–5312.
- [7] B. Awerbuch, "Optimal distributed algorithms for minimum weight spanning tree, counting, leader election, and related problems," in *Proceedings of the nineteenth annual ACM symposium on Theory of computing*. ACM, 1987, pp. 230–240.
- [8] M. Balcan, A. Blum, S. Fine, and Y. Mansour, "Distributed Learning, Communication Complexity and Privacy," in *Proceedings of the 25th Annual Conference on Learning Theory, (COLT'12)*, 2012, pp. 26.1–26.22.
- [9] S. Barbarossa, S. Sardellitti, and P. Di Lorenzo, "Distributed detection and estimation in wireless sensor networks," in *Academic Press Library in Signal Processing, Vol. 2, Communications and Radar Signal Processing*, R. Chellappa and S. Theodoridis, Eds., 2014, pp. 329–408.
- [10] M. Belkin and P. Niyogi, "Semi-supervised learning on Riemannian manifolds," *Machine learning*, vol. 56, no. 1-3, pp. 209–239, 2004.
- [11] M. Belkin, P. Niyogi, and V. Sindhvani, "Manifold regularization: A geometric framework for learning from labeled and unlabeled examples," *Journal of Machine Learning Research*, vol. 7, pp. 2399–2434, 2006.
- [12] T. Bertin-Mahieux, D. P. W. Ellis, B. Whitman, and P. Lamere, "The million song dataset," in *Proceedings of the 12th International Society for Music Information Retrieval Conference (ISMIR'11)*. University of Miami, 2011, pp. 591–596.
- [13] K. Bhaduri, M. D. Stefanski, and A. N. Srivastava, "Privacy-preserving outlier detection through random nonlinear data distortion," *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, vol. 41, no. 1, pp. 260–272, 2011.

-
- [14] K. Bhaduri and H. Kargupta, "A scalable local algorithm for distributed multivariate regression," *Statistical Analysis and Data Mining*, vol. 1, no. 3, pp. 177–194, 2008.
- [15] F. M. Bianchi, S. Scardapane, A. Uncini, A. Rizzi, and A. Sadeghian, "Prediction of telephone calls load using Echo State Network with exogenous variables," *Neural Networks*, vol. 71, pp. 204–213, 2015.
- [16] P. Bianchi and J. Jakubowicz, "Convergence of a multi-agent projected stochastic gradient algorithm for non-convex optimization," *IEEE Transactions on Automatic Control*, vol. 58, no. 2, pp. 391–405, 2013.
- [17] 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 (COLT'92)*. ACM, 1992, pp. 144–152.
- [18] L. Bottou and O. Bousquet, "The Tradeoffs of Large Scale Learning," *Artificial Intelligence*, vol. 20, pp. 161–168, 2008.
- [19] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, "Randomized gossip algorithms," *IEEE Transactions on Information Theory*, vol. 52, no. 6, pp. 2508–2530, 2006.
- [20] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends® in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [21] L. Breiman, "Bagging predictors," *Machine learning*, vol. 24, no. 2, pp. 123–140, 1996.
- [22] J. B. Butcher, D. Verstraeten, B. Schrauwen, C. R. Day, and P. W. Haycock, "Reservoir computing and extreme learning machines for non-linear time-series data analysis," *Neural Networks*, vol. 38, pp. 76–89, 2013.
- [23] E. J. Candès and B. Recht, "Exact matrix completion via convex optimization," *Foundations of Computational mathematics*, vol. 9, no. 6, pp. 717–772, 2009.
- [24] E. Candes and M. Wakin, "An Introduction To Compressive Sampling," *IEEE Signal Processing Magazine*, vol. 25, no. 2, pp. 21–30, 2008.
- [25] F. S. Cattivelli, C. G. Lopes, and A. H. Sayed, "Diffusion recursive least-squares for distributed estimation over adaptive networks," *IEEE Transactions on Signal Processing*, vol. 56, no. 5, pp. 1865–1877, 2008.
- [26] F. S. Cattivelli and A. H. Sayed, "Diffusion LMS strategies for distributed estimation," *IEEE Transactions on Signal Processing*, vol. 58, no. 3, pp. 1035–1048, 2010.
- [27] V. Ceperic and A. Baric, "Reducing complexity of echo state networks with sparse linear regression algorithms," in *Proceedings of the 2014 UKSim-AMSS 16th International Conference on Computer Modelling and Simulation (UKSim'14)*, March 2014, pp. 26–31.
- [28] B. Chaib-Draa, B. Moulin, R. Mandiau, and P. Millot, "Trends in distributed artificial intelligence," *Artificial Intelligence Review*, vol. 6, no. 1, pp. 35–66, 1992.
- [29] C.-C. Chang and C.-J. Lin, "Libsvm: a library for support vector machines," *ACM Transactions on Intelligent Systems and Technology*, vol. 2, no. 3, p. 27, 2011.
- [30] O. Chapelle, V. Sindhwani, and S. S. Keerthi, "Optimization techniques for semi-supervised support vector machines," *Journal of Machine Learning Research*, vol. 9, pp. 203–233, 2008.

- [31] O. Chapelle, B. Schölkopf, and A. Zien, *Semi-Supervised Learning*. MIT Press, 2006.
- [32] O. Chapelle, V. Sindhwani, and S. S. Keerthi, "Branch and bound for semi-supervised support vector machines," in *Advances in Neural Information Processing Systems*, 2006, pp. 217–224.
- [33] O. Chapelle and A. Zien, "Semi-supervised classification by low density separation," in *Proceedings of the tenth international workshop on artificial intelligence and statistics*, vol. 1, 2005, pp. 57–64.
- [34] J. Chen, C. Wang, Y. Sun, and X. S. Shen, "Semi-supervised Laplacian regularized least squares algorithm for localization in wireless sensor networks," *Computer Networks*, vol. 55, no. 10, pp. 2481–2491, 2011.
- [35] J. Chen and A. H. Sayed, "Diffusion adaptation strategies for distributed optimization and learning over networks," *IEEE Transactions on Signal Processing*, vol. 60, no. 8, pp. 4289–4305, 2012.
- [36] J. Chen, C. Richard, P. Honeine, and J. C. M. Bermudez, "Non-negative distributed regression for data inference in wireless sensor networks," in *Proceedings of the 2010 Forty Fourth Asilomar Conference on Signals, Systems and Computers (ASILOMAR'10)*. IEEE, 2010, pp. 451–455.
- [37] C.-T. Chu, S. K. Kim, Y.-A. Lin, Y. Yu, G. Bradski, A. Y. Ng, and K. Olukotun, "Map-Reduce for Machine Learning on Multicore," in *Advances in Neural Information Processing Systems*, 2007, pp. 281–288.
- [38] A. Coates, A. Y. Ng, and H. Lee, "An analysis of single-layer networks in unsupervised feature learning," in *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics*, 2011, pp. 215–223.
- [39] T. M. Cover and P. Hart, "Nearest neighbor pattern classification," *IEEE Transactions on Information Theory*, vol. 13, no. 1, 1967.
- [40] T. M. Cover, "Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition," *IEEE Transactions on Electronic Computers*, no. 3, pp. 326–334, 1965.
- [41] F. Cucker and S. Smale, "On the mathematical foundations of learning," *Bulletin of the American Mathematical Society*, vol. 39, no. 1, pp. 1–49, 2002.
- [42] S. Datta, K. Bhaduri, C. Giannella, R. Wolff, and H. Kargupta, "Distributed data mining in peer-to-peer networks," *IEEE Internet Computing*, vol. 10, no. 4, pp. 18–26, 2006.
- [43] H. Daumé III, J. M. Phillips, A. Saha, and S. Venkatasubramanian, "Efficient protocols for distributed classification and optimization," in *Algorithmic Learning Theory*. Springer, 2012, pp. 154–168.
- [44] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, A. Senior, P. Tucker, K. Yang, Q. V. Le *et al.*, "Large scale distributed deep networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1223–1231.
- [45] O. Dekel, R. Gilad-Bachrach, O. Shamir, and L. Xiao, "Optimal distributed online prediction using mini-batches," *Journal of Machine Learning Research*, vol. 13, no. 1, pp. 165–202, 2012.
- [46] P. Di Lorenzo and G. Scutari, "NEXT: In-Network Nonconvex Optimization," *IEEE Transactions on Signal and Information Processing over Networks*, 2016, in press.

- [47] P. Di Lorenzo and A. H. Sayed, "Sparse distributed learning based on diffusion adaptation," *IEEE Transactions on Signal Processing*, vol. 61, no. 6, pp. 1419–1433, 2013.
- [48] J. C. Duchi, A. Agarwal, and M. J. Wainwright, "Dual averaging for distributed optimization: convergence analysis and network scaling," *IEEE Transactions on Automatic Control*, vol. 57, no. 3, pp. 592–606, 2012.
- [49] X. Dutoit, B. Schrauwen, J. Van Campenhout, D. Stroobandt, H. Van Brussel, and M. Nuttin, "Pruning and regularization in reservoir computing," *Neurocomputing*, vol. 72, no. 7, pp. 1534–1546, 2009.
- [50] D. P. W. Ellis, "Classifying music audio with timbral and chroma features," in *Proceedings of the 8th International Conference on Music Information Retrieval*. Austrian Computer Society, 2007, pp. 339–340.
- [51] T. Evgeniou, M. Pontil, and T. Poggio, "Regularization networks and support vector machines," *Advances in Computational Mathematics*, vol. 13, no. 1, pp. 1–50, 2000.
- [52] M. Fernández-Delgado, E. Cernadas, S. Barro, and D. Amorim, "Do we Need Hundreds of Classifiers to Solve Real World Classification Problems?" *Journal of Machine Learning Research*, vol. 15, pp. 3133–3181, 2014.
- [53] R. Fierimonte, S. Scardapane, M. Panella, and A. Uncini, "A Comparison of Consensus Strategies for Distributed Learning of Random Vector Functional-Link Networks," in *Advances in Neural Networks*. Springer, 2016.
- [54] R. A. Fisher, "The use of multiple measurements in taxonomic problems," *Annals of eugenics*, vol. 7, no. 2, pp. 179–188, 1936.
- [55] K. Flouri, B. Beferull-Lozano, and P. Tsakalides, "Training a SVM-based classifier in distributed sensor networks," in *Proceedings of 14th European Signal Processing Conference (EUSIPCO'06)*, vol. 2006, 2006, pp. 1–5.
- [56] —, "Distributed consensus algorithms for SVM training in wireless sensor networks," in *Proceedings of 16th European Signal Processing Conference (EUSIPCO'08)*, 2008, pp. 25–29.
- [57] —, "Optimal gossip algorithm for distributed consensus svm training in wireless sensor networks," in *Proceedings of the 16th International Conference on Digital Signal Processing (DSP'09)*. IEEE, 2009, pp. 1–6.
- [58] P. A. Forero, A. Cano, and G. B. Giannakis, "Consensus-based distributed support vector machines," *Journal of Machine Learning Research*, vol. 11, pp. 1663–1707, 2010.
- [59] Z. Fu, G. Lu, K. M. Ting, and D. Zhang, "A survey of audio-based music classification and annotation," *IEEE Transactions on Multimedia*, vol. 13, no. 2, pp. 303–319, 2011.
- [60] N. García-Pedrajas and A. Haro-García, "Scaling up data mining algorithms: review and taxonomy," *Progress in Artificial Intelligence*, vol. 1, no. 1, pp. 71–87, 2012.
- [61] L. Georgopoulos and M. Hasler, "Distributed machine learning in networks by consensus," *Neurocomputing*, vol. 124, pp. 2–12, 2014.
- [62] A. N. Gorban, "Approximation of continuous functions of several variables by an arbitrary nonlinear continuous function of one variable, linear functions, and their superpositions," *Applied mathematics letters*, vol. 11, no. 3, pp. 45–49, 1998.

- [63] A. N. Gorban, I. Yu, D. V. Prokhorov, and N. A. Jun, "Approximation with Random Bases: Pro et Contra," 2015, arXiv preprint arXiv:1506.04631.
- [64] H. P. Graf, E. Cosatto, L. Bottou, I. Dourdanovic, and V. Vapnik, "Parallel support vector machines: The cascade svm," in *Advances in Neural Information Processing Systems*, 2004, pp. 521–528.
- [65] S. Guarnieri, F. Piazza, and A. Uncini, "Multilayer feedforward networks with adaptive spline activation function," *IEEE Transactions on Neural Networks*, vol. 10, no. 3, pp. 672–683, 1999.
- [66] C. Guestrin, P. Bodik, R. Thibaux, M. Paskin, and S. Madden, "Distributed regression: an efficient framework for modeling sensor network data," in *Proceedings of the 3rd International Symposium on Information Processing in Sensor Networks (IPSN'04)*. IEEE, 2004, pp. 1–10.
- [67] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, 2nd ed. Springer, 2009.
- [68] C. Hensel and H. Dutta, "GADGET SVM: a gossip-based sub-gradient svm solver," in *Proceedings of the 2009 International Conference on Machine Learning (ICML'2009)*, 2009.
- [69] M. Hermans and B. Schrauwen, "Training and analysing deep recurrent neural networks," in *Advances in Neural Information Processing Systems*, 2013, pp. 190–198.
- [70] D. E. Hershberger and H. Kargupta, "Distributed multivariate regression using wavelet-based collective data mining," *Journal of Parallel and Distributed Computing*, vol. 61, no. 3, pp. 372–400, 2001.
- [71] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [72] T. Hofmann, B. Schölkopf, and A. J. Smola, "Kernel methods in machine learning," *Annals of Statistics*, vol. 36, no. 3, pp. 1171–1220, 2008.
- [73] P. Honeine, M. Essoloh, C. Richard, and H. Snoussi, "Distributed regression in sensor networks with a reduced-order kernel model," in *Proceedings of the 2008 IEEE Global Telecommunications Conference (GLOBECOM'08)*. IEEE, 2008, pp. 1–5.
- [74] P. Honeine, C. Richard, J. C. M. Bermudez, and H. Snoussi, "Distributed prediction of time series data with kernels and adaptive filtering techniques in sensor networks," in *Proceedings of the 42nd Asilomar Conference on Signals, Systems and Computers*. IEEE, 2008, pp. 246–250.
- [75] G.-B. Huang, Q.-Y. Zhu, and C.-K. Siew, "Extreme Learning Machine: Theory and Applications," *Neurocomputing*, vol. 70, no. 1-3, pp. 489–501, Dec. 2006.
- [76] S. Huang and C. Li, "Distributed Extreme Learning Machine for Nonlinear Learning over Network," *Entropy*, vol. 17, no. 2, pp. 818–840, 2015.
- [77] B. Igel'nik and Y.-H. Pao, "Stochastic choice of basis functions in adaptive function approximation and the functional-link net," *IEEE Transactions on Neural Networks*, vol. 6, no. 6, pp. 1320–1329, 1995.
- [78] H. Jaeger, "The echo state approach to analysing and training recurrent neural networks," Technical Report GMD Report 148, German National Research Center for Information Technology, Tech. Rep., 2001.

- [79] —, “Adaptive nonlinear system identification with echo state networks,” in *Advances in Neural Information Processing Systems*, 2002, pp. 593–600.
- [80] H. Jaeger and H. Haas, “Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication,” *Science*, vol. 304, no. 5667, pp. 78–80, 2004.
- [81] M. Jaggi, V. Smith, M. Takác, J. Terhorst, S. Krishnan, T. Hofmann, and M. I. Jordan, “Communication-efficient distributed dual coordinate ascent,” in *Advances in Neural Information Processing Systems*, 2014, pp. 3068–3076.
- [82] D. Jakovetic, J. Xavier, and J. M. Moura, “Fast distributed gradient methods,” *IEEE Transactions on Automatic Control*, vol. 59, no. 5, pp. 1131–1146, 2014.
- [83] T. Joachims, “Transductive Inference for Text Classification Using Support Vector Machines,” *Proceedings of the 1999 International Conference on Machine Learning (ICML’99)*, pp. 200–209, 1999.
- [84] B. Johansson, M. Rabi, and M. Johansson, “A simple peer-to-peer algorithm for distributed optimization in sensor networks,” in *Proceedings of the 2007 46th IEEE Conference on Decision and Control (CDC’07)*. IEEE, 2007, pp. 4705–4710.
- [85] A. F. Karr, X. Lin, A. P. Sanil, and J. P. Reiter, “Secure regression on distributed databases,” *Journal of Computational and Graphical Statistics*, vol. 14, no. 2, pp. 263–279, 2005.
- [86] A. Krizhevsky and G. Hinton, “Learning multiple layers of features from tiny images,” *Computer Science Department, University of Toronto, Tech. Rep.*, 2009.
- [87] A. Lazarevic and Z. Obradovic, “Boosting algorithms for parallel and distributed learning,” *Distributed and Parallel Databases*, vol. 11, no. 2, pp. 203–229, 2002.
- [88] C.-P. Lee and D. Roth, “Distributed Box-Constrained Quadratic Optimization for Dual Linear SVM,” in *Proceedings of the 32nd International Conference on Machine Learning (ICML’15)*. ICML, 2015.
- [89] S. Lee and A. Nedic, “Distributed random projection algorithm for convex optimization,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 7, no. 2, pp. 221–229, 2013.
- [90] D. Li, M. Han, and J. Wang, “Chaotic time series prediction based on a novel robust echo state network,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 23, no. 5, pp. 787–799, 2012.
- [91] W. Li, D. Wang, and T. Chai, “Flame image-based burning state recognition for sintering process of rotary kiln using heterogeneous features and fuzzy integral,” *IEEE Transactions on Industrial Informatics*, vol. 8, no. 4, pp. 780–790, 2012.
- [92] A. Y. Lin and Q. Ling, “Decentralized and privacy-preserving low-rank matrix completion,” *Journal of the Operations Research Society of China*, vol. 3, no. 2, pp. 189–205, 2015.
- [93] C.-B. Lin, “Projected gradient methods for nonnegative matrix factorization,” *Neural computation*, vol. 19, no. 10, pp. 2756–2779, 2007.
- [94] Q. Ling, Y. Xu, W. Yin, and Z. Wen, “Decentralized low-rank matrix completion,” in *Proceedings of the 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP’12)*. IEEE, 2012, pp. 2925–2928.

- [95] K. Liu, H. Kargupta, and J. Ryan, "Random projection-based multiplicative data perturbation for privacy preserving distributed data mining," *IEEE Transactions on Knowledge and Data Engineering*, vol. 18, no. 1, pp. 92–106, 2006.
- [96] W. Liu, J. C. Príncipe, and S. Haykin, *Kernel Adaptive Filtering: a Comprehensive Introduction*. Wiley Press, 2010.
- [97] Y. Liu, C. Li, and Z. Zhang, "Diffusion sparse least-mean squares over networks," *IEEE Transactions on Signal Processing*, vol. 60, no. 8, pp. 4480–4485, 2012.
- [98] Z. Liu, Y. Liu, and C. Li, "Distributed sparse recursive least-squares over networks," *IEEE Transactions on Signal Processing*, vol. 62, no. 6, pp. 1386–1395, 2014.
- [99] S. Lodi, R. Nanculef, and C. Sartori, "Single-pass distributed learning of multi-class svms using core-sets," *Proceedings of the 2010 SIAM International Conference on Data Mining (SDM'10)*, vol. 14, no. 27, p. 2, 2010.
- [100] C. G. Lopes and A. H. Sayed, "Incremental adaptive strategies over distributed networks," *IEEE Transactions on Signal Processing*, vol. 55, no. 8, pp. 4064–4077, 2007.
- [101] —, "Diffusion least-mean squares over adaptive networks: Formulation and performance analysis," *IEEE Transactions on Signal Processing*, vol. 56, no. 7, pp. 3122–3136, 2008.
- [102] D. Lowe, "Multi-variable functional interpolation and adaptive networks," *Complex Systems*, vol. 2, pp. 321–355.
- [103] Y. Lu, V. Roychowdhury, and L. Vandenberghe, "Distributed parallel support vector machines in strongly connected networks," *IEEE Transactions on Neural Networks*, vol. 19, no. 7, pp. 1167–1178, 2008.
- [104] M. Lukoševičius and H. Jaeger, "Reservoir computing approaches to recurrent neural network training," *Computer Science Review*, vol. 3, no. 3, pp. 127–149, 2009.
- [105] M. W. Mahoney, "Randomized algorithms for matrices and data," *Foundations and Trends® in Machine Learning*, vol. 3, no. 2, pp. 123–224, 2011.
- [106] G. Mateos, J. A. Bazerque, and G. B. Giannakis, "Distributed sparse linear regression," *IEEE Transactions on Signal Processing*, vol. 58, no. 10, pp. 5262–5276, 2010.
- [107] G. Mateos and G. B. Giannakis, "Distributed recursive least-squares: Stability and performance analysis," *IEEE Transactions on Signal Processing*, vol. 60, no. 7, pp. 3740–3754, 2012.
- [108] G. Mateos, I. D. Schizas, and G. B. Giannakis, "Distributed recursive least-squares for consensus-based in-network adaptive estimation," *IEEE Transactions on Signal Processing*, vol. 57, no. 11, pp. 4583–4588, 2009.
- [109] B. McWilliams, C. Heinze, N. Meinshausen, G. Krummenacher, and H. P. Vanchinathan, "LOCO: Distributing Ridge Regression with Random Projections," *arXiv preprint arXiv:1406.3469*, 2014.
- [110] S. Melacci and M. Belkin, "Laplacian support vector machines trained in the primal," *Journal of Machine Learning Research*, vol. 12, pp. 1149–1184, 2011.
- [111] I. Mierswa and K. Morik, "Automatic feature extraction for classifying audio data," *Machine learning*, vol. 58, no. 2-3, pp. 127–149, 2005.

- [112] B. Mishra, G. Meyer, and R. Sepulchre, "Low-rank optimization for distance matrix completion," in *Proceedings of the 2011 50th IEEE conference on Decision and control and European control conference (CDC-ECC'11)*. IEEE, 2011, pp. 4455–4460.
- [113] D. Monner and J. A. Reggia, "A generalized lstm-like training algorithm for second-order recurrent neural networks," *Neural Networks*, vol. 25, pp. 70–83, 2012.
- [114] J. F. Mota, J. M. Xavier, P. M. Aguiar, and M. Püschel, "Distributed basis pursuit," *IEEE Transactions on Signal Processing*, vol. 60, no. 4, pp. 1942–1956, 2012.
- [115] A. Navia-Vázquez, D. Gutierrez-Gonzalez, E. Parrado-Hernández, and J. J. Navarro-Abellan, "Distributed support vector machines," *IEEE Transactions on Neural Networks*, vol. 17, no. 4, pp. 1091–1097, 2006.
- [116] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [117] M. Newman, *Networks: an introduction*. Oxford University Press, 2010.
- [118] O. Obst, "Distributed fault detection in sensor networks using a recurrent neural network," *Neural Processing Letters*, vol. 40, no. 3, pp. 261–273, 2014.
- [119] R. Olfati-Saber, J. A. Fax, and R. M. Murray, "Consensus and cooperation in networked multi-agent systems," *Proceedings of the IEEE*, vol. 95, no. 1, pp. 215–233, 2007.
- [120] R. Ormándi, I. Hegedűs, and M. Jelasity, "Asynchronous peer-to-peer data mining with stochastic gradient descent," in *Euro-Par 2011 Parallel Processing*. Springer, 2011, pp. 528–540.
- [121] Y.-H. Pao, G.-H. Park, and D. J. Sobajic, "Learning and generalization characteristics of the random vector functional-link net," *Neurocomputing*, vol. 6, no. 2, pp. 163–180, 1994.
- [122] Y.-H. Pao and Y. Takefji, "Functional-link net computing," *IEEE Computer Journal*, vol. 25, no. 5, pp. 76–79, 1992.
- [123] B.-H. Park and H. Kargupta, "Distributed data mining: Algorithms, systems, and applications," in *The Handbook of Data Mining*, N. Ye, Ed. Lawrence Erlbaum Associates, Incorporated, 2002, pp. 341–358.
- [124] J. Park and I. W. Sandberg, "Universal approximation using radial-basis-function networks," *Neural computation*, vol. 3, no. 2, pp. 246–257, 1991.
- [125] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in *Proceedings of the 30th International Conference on Machine Learning (ICML'12)*, 2012.
- [126] R. U. Pedersen, "Using support vector machines for distributed machine learning," Ph.D. dissertation, Dept. of Computer Science, University of Copenhagen, 2005.
- [127] F. Perez-Cruz and S. R. Kulkarni, "Robust and low complexity distributed kernel least squares learning in sensor networks," *IEEE Signal Processing Letters*, vol. 17, no. 4, pp. 355–358, 2010.
- [128] D. Peteiro-Barral and B. Guijarro-Berdiñas, "A survey of methods for distributed machine learning," *Progress in Artificial Intelligence*, vol. 2, no. 1, pp. 1–11, 2013.
- [129] J. B. Predd, S. R. Kulkarni, and H. V. Poor, "Distributed kernel regression: An algorithm for training collaboratively," in *Proceedings of the 2006 IEEE Information Theory Workshop (ITW'06)*. IEEE, 2006, pp. 332–336.

- [130] —, “Distributed learning in wireless sensor networks,” *IEEE Signal Processing Magazine*, vol. 23, no. 4, pp. 56–69, 2006.
- [131] —, “A collaborative training algorithm for distributed learning,” *IEEE Transactions on Information Theory*, vol. 55, no. 4, pp. 1856–1871, 2009.
- [132] J. C. Principe and B. Chen, “Universal Approximation with Convex Optimization: Gimmick or Reality?” *IEEE Computational Intelligence Magazine*, vol. 10, no. 2, pp. 68–77, 2015.
- [133] M. G. Rabbat and R. D. Nowak, “Quantized incremental algorithms for distributed optimization,” *IEEE Journal on Selected Areas in Communications*, vol. 23, no. 4, pp. 798–808, 2005.
- [134] A. Rahimi and B. Recht, “Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning,” *Advances in Neural Information Processing Systems*, vol. 1, no. 1, pp. 1–8, 2009.
- [135] —, “Uniform approximation of functions with random bases,” in *Proceedings of the 46th Annual Allerton Conference on Communication, Control, and Computing*, 2008, pp. 555–561.
- [136] S. S. Ram, A. Nedic, and V. V. Veeravalli, “Incremental stochastic subgradient algorithms for convex optimization,” *SIAM Journal on Optimization*, vol. 20, no. 2, pp. 691–717, 2009.
- [137] C. Ravazzi, S. M. Fosson, and E. Magli, “Distributed Iterative Thresholding for 0/1-Regularized Linear Inverse Problems,” *IEEE Transactions on Information Theory*, vol. 61, no. 4, pp. 2081–2100, 2015.
- [138] S. Ravindran, D. Anderson, and M. Slaney, “Low-power audio classification for ubiquitous sensor networks,” in *Proceedings of the 2004 IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP'04)*, vol. 4. IEEE, 2004, pp. iv–337–iv–340.
- [139] B. Recht, C. Re, S. Wright, and F. Niu, “Hogwild: A lock-free approach to parallelizing stochastic gradient descent,” in *Advances in Neural Information Processing Systems*, 2011, pp. 693–701.
- [140] R. Rifkin and A. Klautau, “In defense of one-vs-all classification,” *Journal of Machine Learning Research*, vol. 5, pp. 101–141, 2004.
- [141] F. Rosenblatt, “The perceptron: a probabilistic model for information storage and organization in the brain,” *Psychological review*, vol. 65, no. 6, pp. 386–408, 1958.
- [142] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning Representations by Back-Propagating Errors,” *Nature*, vol. 323, pp. 533–536, 1986.
- [143] S. Samet and A. Miri, “Privacy-preserving back-propagation and extreme learning machine algorithms,” *Data & Knowledge Engineering*, vol. 79, pp. 40–61, 2012.
- [144] S. Sardellitti, M. Giona, and S. Barbarossa, “Fast distributed average consensus algorithms based on advection-diffusion processes,” *IEEE Transactions on Signal Processing*, vol. 58, no. 2, pp. 826–842, 2010.
- [145] A. H. Sayed, “Adaptive networks,” *Proceedings of the IEEE*, vol. 102, no. 4, pp. 460–497, 2014.
- [146] A. H. Sayed and C. G. Lopes, “Distributed recursive least-squares strategies over adaptive networks,” in *Proceedings of the Fortieth Asilomar Conference on Signals, Systems and Computers (ACSSC'06)*. IEEE, 2006, pp. 233–237.

- [147] S. Scardapane, D. Comminiello, M. Scarpiniti, and A. Uncini, "Music classification using extreme learning machines," in *Proceedings of the 2013 8th International Symposium on Image and Signal Processing and Analysis (ISPA'13)*. IEEE, 2013, pp. 377–381.
- [148] S. Scardapane, R. Fierimonte, D. Wang, M. Panella, and A. Uncini, "Distributed Music Classification Using Random Vector Functional-Link Nets," in *Proceedings of the 2015 International Joint Conference on Neural Networks (IJCNN'15)*. IEEE/INNS, 2015, pp. 1–8.
- [149] S. Scardapane, G. Nocco, D. Comminiello, M. Scarpiniti, and A. Uncini, "An effective criterion for pruning reservoir's connections in echo state networks," in *Proceedings of the 2014 International Joint Conference on Neural Networks (IJCNN'14)*. INNS/IEEE, 2014, pp. 1205–1212.
- [150] S. Scardapane, M. Panella, D. Comminiello, and A. Uncini, "Learning from Distributed Data Sources using Random Vector Functional-Link Networks," in *Procedia Computer Science*, 2015, vol. 53, pp. 468–477.
- [151] S. Scardapane, D. Wang, and M. Panella, "A Decentralized Training Algorithm for Echo State Networks in Distributed Big Data Applications," *Neural Networks*, vol. 78, pp. 65–74, 2016.
- [152] S. Scardapane, D. Comminiello, M. Scarpiniti, and A. Uncini, "A semi-supervised random vector functional-link network based on the transductive framework," *Information Sciences*, 2015, in press.
- [153] S. Scardapane, R. Fierimonte, P. Di Lorenzo, M. Panella, and A. Uncini, "Distributed semi-supervised support vector machines," *Neural Networks*, vol. 80, pp. 43–52, 2016.
- [154] S. Scardapane, D. Wang, M. Panella, and A. Uncini, "Distributed learning for Random Vector Functional-Link networks," *Information Sciences*, vol. 301, pp. 271–284, 2015.
- [155] M. Scarpiniti, D. Comminiello, R. Parisi, and A. Uncini, "Nonlinear spline adaptive filtering," *Signal Processing*, vol. 93, no. 4, pp. 772 – 783, 2013.
- [156] I. D. Schizas, G. Mateos, and G. B. Giannakis, "Distributed LMS for consensus-based in-network adaptive processing," *IEEE Transactions on Signal Processing*, vol. 57, no. 6, pp. 2365–2382, 2009.
- [157] N. Schmitter, "A protocol for privacy preserving neural network learning on horizontally partitioned data," in *Privacy in Statistical Databases (PSD)*, 2008.
- [158] J. Schmidhuber, "Deep Learning in Neural Networks: An Overview," *Neural Networks*, vol. 61, pp. 85–117, 2015.
- [159] M. Schmidt, "Least squares optimization with l1-norm regularization," *CS542B Project Report*, 2005.
- [160] W. F. Schmidt, M. A. Kraaijveld, and R. P. W. Duin, "Feedforward neural networks with random weights," in *Proceedings of the 11th IAPR International Conference on Pattern Recognition (ICPR'92)*. IEEE, 1992, pp. 1–4.
- [161] B. Schölkopf, R. Herbrich, and A. J. Smola, "A generalized representer theorem," in *Computational learning theory*. Springer, 2001, pp. 416–426.
- [162] G. Scutari, F. Facchinei, P. Song, D. P. Palomar, and J.-S. Pang, "Decomposition by partial linearization: Parallel optimization of multi-agent systems," *IEEE Transactions on Signal Processing*, vol. 62, no. 3, pp. 641–656, 2014.

- [163] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu, "1-Bit Stochastic Gradient Descent and its Application to Data-Parallel Distributed Training of Speech DNNs," in *Proceedings of the 15th Annual Conference of the International Speech Communication Association (INTERSPEECH'15)*, 2014.
- [164] —, "On parallelizability of stochastic gradient descent for speech dnns," in *Proceedings of the 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP'14)*. IEEE, 2014, pp. 235–239.
- [165] Q. Shi, C. Shen, R. Hill, and A. Hengel, "Is margin preserved after random projection?" in *Proceedings of the 29th International Conference on Machine Learning (ICML'12)*. ACM, 2012, pp. 591–598.
- [166] G. Shmueli, "To Explain or to Predict?" *Statistical Science*, vol. 25, no. 3, pp. 289–310, 2010.
- [167] D. Shutin and G. Kubin, "Echo state wireless sensor networks," in *2008 IEEE Workshop on Machine Learning for Signal Processing (MLSP'08)*, 2008, pp. 151–156.
- [168] C. N. Silla, C. A. A. Kaestner, and A. L. Koerich, "Automatic music genre classification using ensemble of classifiers," in *Proceedings of the 2007 IEEE International Conference on Systems, Man and Cybernetics (SMC'07)*. IEEE, 2007, pp. 1687–1692.
- [169] S. N. Simic, "A learning-theory approach to sensor networks," *IEEE Pervasive Computing*, vol. 2, no. 4, pp. 44–49, 2003.
- [170] J. J. Steil, "Online reservoir adaptation by intrinsic plasticity for backpropagation–decorrelation and echo state learning," *Neural Networks*, vol. 20, no. 3, pp. 353–364, 2007.
- [171] I. Steinwart and A. Christmann, *Support Vector Machines*, 2008.
- [172] S. M. Stigler, "Gauss and the invention of least squares," *The Annals of Statistics*, pp. 465–474, 1981.
- [173] Y. Sun, Y. Yuan, and G. Wang, "An OS-ELM based distributed ensemble classification framework in P2P networks," *Neurocomputing*, vol. 74, no. 16, pp. 2438–2443, 2011.
- [174] S. Sundhar Ram, A. Nedić, and V. V. Veeravalli, "A new class of distributed optimization algorithms: Application to regression of distributed data," *Optimization Methods and Software*, vol. 27, no. 1, pp. 71–88, 2012.
- [175] N. Takahashi, I. Yamada, and A. H. Sayed, "Diffusion least-mean squares with adaptive combiners: Formulation and performance analysis," *IEEE Transactions on Signal Processing*, vol. 58, no. 9, pp. 4795–4810, 2010.
- [176] W. P. Tay, J. N. Tsitsiklis, and M. Z. Win, "Asymptotic performance of a censoring sensor network," *IEEE Transactions on Information Theory*, vol. 53, no. 11, pp. 4191–4209, 2007.
- [177] J. J. Thompson, M. R. Blair, L. Chen, and A. J. Henrey, "Video game telemetry as a critical tool in the study of complex skill learning," *PloS one*, vol. 8, no. 9, p. e75129, 2013.
- [178] R. Tibshirani, "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [179] M. H. Tong, A. D. Bickett, E. M. Christiansen, and G. W. Cottrell, "Learning grammatical structure with echo state networks," *Neural Networks*, vol. 20, no. 3, pp. 424–432, 2007.

- [180] S. Tong and D. Koller, "Support vector machine active learning with applications to text classification," *Journal of Machine Learning Research*, vol. 2, pp. 45–66, 2002.
- [181] M. Torii, K. Waghlikar, and H. Liu, "Using machine learning for concept extraction on clinical documents from multiple data sources," *Journal of the American Medical Informatics Association*, vol. 18, no. 5, pp. 580–587, 2011.
- [182] F. Triefenbach, A. Jalalvand, K. Demuynck, and J.-P. Martens, "Acoustic modeling with hierarchical reservoirs," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 11, pp. 2439–2450, Nov 2013.
- [183] J. N. Tsitsiklis, D. P. Bertsekas, and M. Athans, "Distributed asynchronous deterministic and stochastic gradient optimization algorithms," *IEEE Transactions on Automatic Control*, vol. 31, no. 9, pp. 803–812, 1986.
- [184] A. Uncini, *Fundamentals of Adaptive Signal Processing*. Springer, 2015.
- [185] K. Vandoorne, J. Dambre, D. Verstraeten, B. Schrauwen, and P. Bienstman, "Parallel reservoir computing using optical amplifiers," *IEEE Transactions on Neural Networks*, vol. 22, no. 9, pp. 1469–1481, 2011.
- [186] V. S. Verykios, E. Bertino, I. N. Fovino, L. P. Provenza, Y. Saygin, and Y. Theodoridis, "State-of-the-art in privacy preserving data mining," *ACM Sigmod Record*, vol. 33, no. 1, pp. 50–57, 2004.
- [187] D. Wang, J. Zheng, Y. Zhou, and J. Li, "A scalable support vector machine for distributed classification in ad hoc sensor networks," *Neurocomputing*, vol. 74, no. 1, pp. 394–400, 2010.
- [188] D. Wang and Y. Zhou, "Distributed support vector machines: An overview," in *Proceedings of the 2012 24th Chinese Control and Decision Conference (CCDC'12)*. IEEE, 2012, pp. 3897–3901.
- [189] L. P. Wang and C. R. Wan, "Comments on The Extreme Learning Machine," *IEEE Transactions on Neural Networks*, vol. 19, no. 8, pp. 1494–1495, 2008.
- [190] D. J. Watts and S. H. Strogatz, "Collective dynamics of 'small-world' networks," *Nature*, vol. 393, no. 6684, pp. 440–442, 1998.
- [191] P. J. Werbos, "Backpropagation through time: What it does and how to do it," *Proceedings of the IEEE*, vol. 78, no. 10, pp. 1550–1560, 1990.
- [192] H. White, "An additional hidden unit test for neglected nonlinearity in multilayer feedforward networks," in *Proceedings of the 1989 International Joint Conference on Neural Networks (IJCNN'89)*. IEEE, 1989, pp. 451–455.
- [193] —, "Approximate nonlinear forecasting methods," *Handbook of economic forecasting*, vol. 1, pp. 459–512, 2006.
- [194] B. Widrow and M. A. Lehr, "30 years of adaptive neural networks: perceptron, Madaline, and backpropagation," *Proceedings of the IEEE*, vol. 78, no. 9, pp. 1415–1442, 1990.
- [195] B. Widrow, "Reply to the Comments on the No-Prop algorithm," *Neural Networks*, vol. 48, p. 204, 2013.
- [196] B. Widrow, A. Greenblatt, Y. Kim, and D. Park, "The no-prop algorithm: A new learning algorithm for multilayer neural networks," *Neural Networks*, vol. 37, pp. 182–188, 2013.

-
- [197] X. Wu, X. Zhu, G.-Q. Wu, and W. Ding, "Data mining with big data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 1, pp. 97–107, 2014.
- [198] L. Xiao and S. Boyd, "Fast linear iterations for distributed averaging," *Systems & Control Letters*, vol. 53, no. 1, pp. 65–78, 2004.
- [199] L. Xiao, S. Boyd, and S. J. Kim, "Distributed average consensus with least-mean-square deviation," *Journal of Parallel and Distributed Computing*, vol. 67, no. 1, pp. 33–46, 2007.
- [200] L. Xiao, S. Boyd, and S. Lall, "A scheme for robust distributed sensor fusion based on average consensus," in *Proceedings of the Fourth International Symposium on Information Processing in Sensor Networks (IPSN'05)*. IEEE, 2005, pp. 63–70.
- [201] —, "A space-time diffusion scheme for peer-to-peer least-squares estimation," in *Proceedings of the 5th international conference on Information processing in sensor networks (IPSN'06)*. ACM, 2006, pp. 168–176.
- [202] Y. Xue, L. Yang, and S. Haykin, "Decoupled echo state networks with lateral inhibition," *Neural Networks*, vol. 20, no. 3, pp. 365–376, 2007.
- [203] I. B. Yildiz, H. Jaeger, and S. J. Kiebel, "Re-visiting the echo state property," *Neural Networks*, vol. 35, pp. 1–9, 2012.
- [204] Y. Zhang and S. Zhong, "A privacy-preserving algorithm for distributed training of neural network ensembles," *Neural Computing and Applications*, vol. 22, no. 1, pp. 269–282, 2013.
- [205] X. Zhao and A. H. Sayed, "Asynchronous adaptation and learning over networks—Part I: Modeling and stability analysis," *IEEE Transactions on Signal Processing*, vol. 63, no. 4, pp. 811–826, 2015.
- [206] M. Zinkevich, J. Langford, and A. J. Smola, "Slow learners are fast," in *Advances in Neural Information Processing Systems*, 2009, pp. 2331–2339.
- [207] M. Zinkevich, M. Weimer, A. Smola, and L. J. Li, "Parallelized stochastic gradient descent," in *Advances in Neural Information Processing Systems*, 2010, pp. 2595–2603.