

Genuense Athenaeum

Cristian Jesús Vega Cereño

**Implicit Regularization: Insights from Iterative
Regularization and Overparameterization**

PhD thesis

DIMA - Dipartimento di eccellenza 2023-27
May 18, 2024



University of Genova

PhD Program in Mathematics and Applications
Curriculum Mathematical Methods for Data Analysis

Implicit Regularization: Insights from Iterative Regularization and Overparameterization

by

Cristian Jesús Vega Cereño

Thesis submitted for the degree of Doctor of Philosophy (Cycle XXXVI)

May 18, 2024

Silvia Villa
Lorenzo Rosasco
Stefano Vigni

Supervisor
Supervisor
Head of the PhD program



MaLGA Center, DIMA - Dipartimento di eccellenza 2023-27, Università di Genova, Via
Dodecaneso 35, Genoa, Italy

May 18, 2024

Cristian Jesús Vega Cereño

All Rights Reserved ©

For my grandmother Nelly Concha.

Abstract

The goal of many data-driven problems is to achieve a good prediction by estimating a quantity of interest based on a finite set of (possibly noisy) measurements, exploiting training data, and some property of the model that may be known or not a-priori. The most common methods to reach this objective are explicit and implicit regularization. The first technique consists of minimizing the sum of a loss function plus a regularizer, which is explicitly added to the objective function and entails some a priori knowledge or some desired property of the solutions that we want to select. The second technique, implicit regularization, which is the main focus of this thesis, consists of minimizing a regularizer subject to the constraints established by the loss function minimizers. In this thesis, we propose two different approaches to solving the implicit problem.

In the first part of the thesis, we follow a more traditional approach. We propose and study two new iterative regularization methods for inverse problems that are based on a primal-dual algorithm, where the bias and the loss is fixed. Our analysis, in the noise-free case, provides convergence rates for the Lagrangian and the feasibility gap. In the noisy case, it provides stability bounds and early stopping rules with theoretical guarantees. The main novelty of our work is the exploitation of some a priori knowledge about the solution set: we show that the linear equations determined by the data can be used more than once along the iterations. We discuss various approaches to reusing linear equations that are at the same time consistent with our assumptions and flexible in their implementation. Finally, we illustrate our theoretical findings with numerical simulations for robust sparse recovery and image reconstruction. We confirm the efficiency of the proposed regularization approaches by comparing the results with state-of-the-art methods.

In the second part of this thesis, inspired by the recent success of re- and over-parameterization trained with gradient descent in machine learning, we flip our perspective by fixing the loss and the algorithm (gradient flow) and reparameterize the linear model. Then, we aim to find the implicit bias introduced by the chosen optimization method and reparameterization. But there is still an open question of how to find systematically what the inductive bias hidden behind the model for a particular optimization scheme is. The goal of this thesis is to take a step in this direction by studying a unified framework encompassing various reparametrizations presented in the state of the art, called time-warped mirror flow. However, a theoretical analysis of the existence and convergence of the trajectory is missing in the state of the art. Here, we fill this gap by providing a comprehensive study. First, we prove the existence and uniqueness of the solution. Next, we establish the convergence of both the trajectory and the corresponding values of the loss function. Finally, for any convex loss function, we prove that the trajectory converges towards a minimizer of the loss function, and we provide an implicit bias. For a specific case of loss functions, including least squares, the implicit bias can be made explicit. Furthermore, we explore the flexibility of our formulation by applying the previous results to different examples related to weight normalization. Finally, we give a criterion to determine, for a

given function that depends only on the norm, a suitable weight normalization parameterization.

Keywords. Primal-dual splitting algorithms, Iterative regularization, Early stopping, Landweber method, Stability and convergence analysis, Overparameterization, Implicit Regularization, Time-warping mirror flow, Fully connected normalized linear networks, Weight normalization.

AMS Mathematics Subject Classification (2020): 34A55, 90C25, 65K10, 49M29.

Acknowledgements

First of all, my heartfelt appreciation goes to my advisor, Silvia Villa, for her continuous support, patience, and exceptional guidance as both my thesis director and a human being. Similarly, I express my gratitude to Lorenzo Rosasco, the co-director of this thesis, to whom I'm grateful for his valuable comments on this work and his extensive knowledge. Their invaluable guidance and support were pivotal throughout the development of this work. Their insightful direction and suggestions helped to shape the findings in this thesis. Their guidance in my writing, communication, and other essential skills has been important to my growth as a mathematician.

I extend special thanks to Cesare Molinari, whose interaction was a profoundly enriching experience. Continuing collaboration with him would be both an honor and a pleasure. His comments made significant enhancements to some of the results I had initially obtained. Similarly, I extend my gratitude to Vasilis Apidopoulos, who became my mentor during this doctoral journey. I also express my gratitude to the two reviewers, Juan Peypouquet and Nicolas Flammarion, for having gladly consented to read this thesis.

The development of this doctorate would not have been possible without the financial support of Marie Skłodowska-Curie Grant No. 861137 of the European Union's Horizon 2020 research and innovation program. I also want to thank the mathematics department at the University of Genova, especially MaLGA. Likewise, I thank my former university (Técnica Federico Santa María), especially Luis Briceño and Julio Deride, who sparked my interest in optimization.

The time I spent in Genoa (and Paris) was made all the more enjoyable thanks to Luis, Angelos, Eliana, Elena, and Fernando. I also must express my profound appreciation for the enormous professional support and wonderful company during these three years from my colleagues Cheik Traore, Jonathan Chirinos, Mouna Gharbi, and Marco Rando.

I especially thank my fiancée, Valeska Campos, for being the biggest source of unconditional support during these years of study, research, and writing of this thesis. I extend special thanks to my friends Omar and Jeremias for making this doctorate time a jolly one. I also must thank my family for their constant help and support during these three years. I would like to emphasize my grandmother, Nelly Concha, who was, is, and until the end of her life will be one of the most important people and loved ones of all time, but regrettably, she could not accompany me in this achievement. You were a big support in my life, and I appreciate all your unconditional love. I want you to know that this achievement and all others that come are due to you.

Finally, I take this opportunity to express my gratitude to each and everyone who directly or indirectly supported me until the completion of my degree, as this achievement would not have been possible without the guidance and blessings of many. Thank you.

Table of Contents

1	General introduction	1
1.1	Motivation	1
1.2	Main contributions	3
1.3	Outline	4
2	Notation and Preliminaries	5
2.1	Basic notions	5
2.2	Matrices	7
2.3	Convex optimization	8
2.4	Differential Calculus and Ordinary differential equations	13
3	Regularization techniques	17
3.1	Inverse problems	17
3.2	Tikhonov regularization, flows, and algorithms	19
3.2.1	$R = 0$: Continuous case	19
3.2.2	$R = 0$: Discrete case	20
3.2.3	Case $R \neq 0$	21
3.3	Iterative regularization	22
4	Fast iterative regularization by reusing data	25
4.1	Introduction	26
4.2	Main problem and algorithm	27
4.2.1	Primal-Dual Splittings with a priori Information	29
4.2.2	Equivalence between Primal-dual and Dual-primal algorithms.	30
4.2.3	Assumptions	31
4.3	Main results	32
4.4	Implementation details	44
4.5	Numerical results	47
4.5.1	ℓ^1 -norm regularization	48
4.5.2	Total variation	51
4.6	Conclusion and Future Work	53
5	Implicit regularization and reparameterization	55
5.1	Reparameterization	55
5.1.1	Optimization by reparameterization	56
5.2	Case $G(\cdot) = 1$: Reparameterizing gradient flow as mirror flow.	56
5.3	Case $G(\cdot) \neq 1$: Time warping Mirror Flow	59
6	Learning from data via overparameterization	61
6.1	Introduction	62
6.1.1	Related work	63
6.2	Global existence and asymptotic analysis	64
6.2.1	Well-posedness	65

6.2.2	Minimization properties and implicit bias	68
6.2.3	Application to fully connected linear networks	71
6.3	Reparameterizing Mirror Descent for radial functions as Projected Gradient Descent	74
6.3.1	Weight normalization of a fully connected network	81
6.3.2	Fully connected normalized linear network of depth 2	83
6.4	Conclusion and Future Work	86
A	Appendix for Chapter 6	89
A.1	Examples	89
	Bibliography	93

List of Figures

4.1	Graphical representation of early stopping. Note that the reconstruction error decreases and then increases, since the iterates first approach the exact solution and then converges to the noisy solution.	49
4.2	Early stopping with respect to the feasibility. Note that their behavior with respect to k is similar to that in Figure 4.1.	50
4.3	Reconstruction error of Tikhonov Method with different penalties.	50
4.4	Qualitative comparison of the 4 proposed methods.	54

List of Tables

- 4.1 Proposed algorithms for iterative regularization. 29
- 4.2 Run-time and number of iterations of each method until it reaches the best reconstruction error. We compare the proposed algorithms with Tikhonov regularization (Tik), Douglas-Rachford (DR), and iterative regularization (PD). 49
- 4.3 General form of the algorithms. 52
- 4.4 Quantitative comparison of the algorithms in terms of Structural similarity (SSIM), peak signal-to-noise ratio (PSNR), Mean square error (MSE), time, and iterations to reach the early stopping. 53

CHAPTER 1

General introduction

1.1 Motivation

Many data-driven problems involve estimating an input-output relation based on a finite set of (possibly noisy) measurements, i.e., $f(x_i) = y_i$. Assuming that f is linear, expressed as $f(x) = h \cdot x$, the problem simplifies to determine h from the following linear equation:

$$Xh = y;$$

where X represents a matrix with rows corresponding to input data points, and y can be seen as the vector of measurements of some unknown h that we want to recover.

In general, the solution might not exist or might not be unique, or might not depend continuously on the data y . To address the issue of existence, a loss or data-fitting function L , often the least squares, is utilized. To solve non-uniqueness and select a particular solution, a regularization term R is introduced, also called the bias of the solution. This regularization enforces prior knowledge or a desired property of the solution. A standard approach to recover h is to assume that it is a minimizer of the following constrained optimization problem:

$$\min_{h \in \mathbb{R}^p} fR(h) : \quad \text{argmin } Lg: \quad (1.1.1)$$

Typically, in (1.1.1), the regularizer R and the loss function L are fixed. Subsequently, a suitable algorithm that leverages the properties of both the loss and the regularizer is designed to efficiently find a solution. In this thesis, we propose two different approaches to solve the implicit problem.

We also assume that L is the least squares and the data is noisy, meaning that there exists y such that:

$$\|y - Xh\| \leq \epsilon;$$

where ϵ is the level of noise. A common method to solve the above problem and avoid instabilities is Tikhonov regularization [52, 128] which consists in minimizing the sum of a penalized loss term plus a regularizer. A trade-off parameter is then introduced to balance the fidelity term and the regularizer. In practice, this implies that the optimization problem has to be solved many times for different values of the parameter. Then, the best performing iterate is chosen according to some a priori data driven criterion and considered the regularized solution.

In the first part of our thesis, we solve (1.1.1) using a classical approach, to deal the case where R is a convex regularizer, which is neither smooth nor strongly convex. We explore an alternative, efficient approach known as early stopping [12, 20, 27, 30, 82, 112, 141, 146]. This method runs an iterative algorithm solving the exact problem but on the inexact data and early stopping to prevent convergence to the noisy solution. In this setting, the number of iterations plays the role of the regularization parameter. Compared to Tikhonov regularization, this procedure is very efficient since only one optimization problem is solved, and not even until convergence. Recently, in [40, 84] an algorithm was introduced that combined primal-dual methods with early stopping. On the other side, in [24], an extra activation step is added, improving the feasibility of the iterates, obtaining empirical speed up. Then, the idea is to combine these two approaches to design an efficient algorithm to solve (1.1.1).

We introduce and analyze two new iterative regularization methods based on a primal-dual algorithm with activations. Our analysis yields in the noise-free scenario ($\sigma = 0$) convergence rates for the Lagrangian and the feasibility gap of (P) , and stability bounds along with early stopping criteria with theoretical guarantees in the noisy scenario. The main novelty of our work is the exploitation of some a priori knowledge about the solution set: we show that the linear equations determined by the data can be used more than once along the iterations. We propose various strategies for reusing linear equations that align with our theoretical framework while remaining flexible for practical implementation. Finally, we illustrate our theoretical findings with numerical simulations for robust sparse recovery and image reconstruction. We confirm the efficiency of the proposed regularization approaches by comparing the results with state-of-the-art methods.

In the second part of this thesis, we flip this perspective by fixing the loss and the algorithm, which, to simplify the analysis, will be gradient flow. We also reparameterize the linear model as follows:

$$\mathcal{V} = q(\cdot)$$

which is called re- or over-parameterization. Then, we want to answer the following question:

For a given loss, algorithm, and reparameterization, what implicit bias is introduced by the chosen optimization method and reparameterization?

In simpler terms, does there exist a regularization function R such that the output of gradient flow applied to the loss function L for a specific reparameterization q is the solution of problem (1.1.1)?

This question was partially answered in [4, 77], which found that gradient flow on the reparameterization and mirror flow on \mathcal{V} are equivalent, for a suitable mirror map that only depends on reparameterization. While these articles cover a wide range of reparameterizations, their assumptions are too restrictive to include certain interesting reparameterizations, such as weight normalization and the multiplication of a matrix by a vector.

To address this limitation, we study a unified framework proposed in [7], called time-warp mirror flow. This consists of the vanilla mirror flow with a scalar preconditioning that allows us to encompass many reparameterizations presented in the state of the art. However, the theoretical analysis of the existence and convergence of the trajectory is missing. Then, we pose the following question:

Can we expect well-posedness, convergence to a stationary point, convergence rates, and implicit bias of the sequence (t) generated by time-warped mirror flow?

In this thesis, we provide a set of assumptions to prove the well-posedness of the trajectory of the time-warping mirror flow. Our analysis includes convergence of both the trajectory and the corresponding values of the loss function. In the case when the loss function is the composition of a linear operator with a strictly convex function, we provide an explicit expression of the implicit bias. Finally, we illustrate our theoretical findings by applying the obtained result to weight normalization and matrix vector overparameterization.

1.2 Main contributions

We briefly mention the main contributions of our work.

- **In Chapter 4**, we design and analyze two new iterative regularization methods for convex regularizers, which are not necessarily smooth nor strongly convex. The new iterative regularization methods are based on primal-dual algorithms [40, 46, 133] combined with the idea of reusing the linear equations determined by the data at every iteration [24]. The first method that we propose is a primal-dual algorithm with additional activations of the linear equations. The second method is a dual-primal algorithm, where a subset containing the dual solutions is activated at each step.

Our analysis, in the noisy case, provides stability bounds and early stopping rules with theoretical guarantees. In the noise-free case, we provide convergence rates for the Lagrangian and the feasibility gap of the problem (P) .

We propose different variants of our algorithm, using different extra activation steps, including a gradient descent step over the least square with a fixed or adaptive step size. We also compare their numerical performance with state-of-the-art methods, obtaining a considerable improvement in run-time.

- **In Chapter 6**, we provide an implicit bias for reparameterizations such that applying gradient flow over \mathcal{P} is equivalent to apply time-warped mirror flow over \mathcal{P} . The analysis consists of several steps. First, we establish conditions for the well-posedness of the time warped mirror flow. Second, we demonstrate that for any convex function, the sequence (t) generated by time-warping mirror flow converges to a stationary point that minimizes the loss function while avoiding the extra stationary points that the reparameterization produces. For the specific case of a strictly convex function composed with a linear operator, an explicit expression for the implicit bias is provided.

We apply our results to matrix vector parameterization and to various weight normalization reparameterizations, generalizing many results in the state of the art. Furthermore, we give a criterion to determine, for a given function that depends only on the norm, a suitable weight normalization parameterization. Finally, we explore the flexibility of our formulation by applying the previous results to different examples related to weight normalization.

Journal publications

- C. Vega, C. Molinari, S. Villa, and L. Rosasco, “Fast iterative regularization by reusing data”, Journal of Inverse and Ill-posed Problems, 2023.

Articles in preparation

- C. Vega, C. Molinari, S. Villa, and L. Rosasco, “Learning from data via overparameterization”, In preparation, 2024.

1.3 Outline

This section offers a concise overview of the thesis, which is organized into six chapters and supplemented by one appendix. Chapters 4 and 6 comprise the core articles developed during the doctoral research, one of them published in a journal [132]. Chapters 3 and 5 provide technical introductions that precede the discussions in Chapters 4 and 6, respectively.

Chapter 2 lays out the notations and mathematical framework essential for this thesis, organized into four sections. The first section is dedicated to Hilbert spaces, establishing the fundamental concepts used throughout the thesis. The second section delves into matrices, presenting the results necessary for our analysis. The third section explores convex optimization concepts, which are mainly, but not exclusively, applied in Chapter 4. In the last section, we introduce essential notions of calculus, ordinary differential equations, and the algorithms used in this thesis necessary for the discussion in Chapter 5.

Chapter 3 presents implicit and iterative regularization approaches. Initially, it outlines the four fundamental components of implicit regularization: the model, the loss, the bias, and the algorithm. Subsequently, the chapter provides a review of iterative algorithms.

Chapter 4 is organized as follows. In Section 4.2 we present the main problem and propose two algorithms to solve it numerically. In Section 4.3 we derive stability and feasibility gap bounds and related early stopping rules. In Section 4.4 we verify the performance of the algorithm on two numerical applications: robust sparse recovery problem and image reconstruction by total variation. Finally, we provide some conclusions.

Chapter 5 formulates our problem as an implicit regularization problem. Next, from the point of view of optimization, we explain how overparameterization induces some implicit bias and we give some examples. Finally, we present the main results in the state of the art for vanilla and time-warped mirror flow.

Chapter 6 is organized as follows. In Section 6.2, we derive the convergence of our algorithm to a stationary point and characterize this solution. In Section 6.3, we derive conditions to cast gradient flow on weight normalization reparameterization, which decomposes a vector into its direction and its norm, as a mirror flow over a radial function. Finally, we provide some conclusions and open problems for future research directions.

Additionally, the thesis includes an appendix: which offers the supplementary material of Chapter 6, omitted from the main text due to its extension. This section is mainly focused on examples of overparameterization.

CHAPTER 2

Notation and Preliminaries

We begin by introducing some notation used throughout this thesis.

2.1 Basic notions

Let H, G be two real vector spaces. Although in most of this thesis $H = \mathbb{R}^p$ and $G = \mathbb{R}^d$, which are finite dimensional, this section is written for general Hilbert spaces. For further results, the reader is referred to [15].

Definition 2.1.1. Now, we present some examples of norms in \mathbb{R}^p , which will be used in this thesis.

- **General ℓ_q -norm:** Let $x \in \mathbb{R}^p$. For $q \geq 1$, the ℓ_q -norm is defined as:

$$\|x\|_q = \left(\sum_{i=1}^p |x_i|^q \right)^{\frac{1}{q}}$$

For $q = 1, 2$; and ∞ ; we have that

$$\|x\|_1 = \sum_{i=1}^p |x_i|; \quad \|x\|_2 = \left(\sum_{i=1}^p |x_i|^2 \right)^{\frac{1}{2}} \quad \text{and} \quad \|x\|_\infty = \max_{1 \leq i \leq p} |x_i|$$

- **$\ell_{1,2}$ -norm:** Let A be a square matrix in $\mathbb{R}^{p \times p}$. The $\ell_{1,2}$ -norm is defined as:

$$\|A\|_{1,2} = \left(\sum_{i=1}^p \sum_{j=1}^p |a_{ij}|^2 \right)^{\frac{1}{2}}$$

where from now on a_{ij} is i -th component of the j -th column of A . This norm will be used in the numerical examples.

Remark 2.1.2. The following two fundamental inequalities are continuously employed throughout this thesis. For every $(x, z) \in H \times H$ and $\alpha > 0$, we have that:

$$\langle x, z \rangle \leq \|x\| \|z\| \leq \frac{\alpha \|x\|^2}{2} + \frac{\|z\|^2}{2\alpha}$$

The first and second inequalities are known as the Cauchy-Schwarz and Young inequalities, respectively.

The open ball of radius $r > 0$ centered at $x \in H$ is the set of all points in H whose distance from x is strictly less than r . This can be formally written as:

$$B(x; r) = \{z \in H \mid \|z - x\| < r\}$$

Analogously, the closed ball of radius $r > 0$ centered at $x \in H$ is the set of all points in H whose distance from x is less than or equal to r and is denoted by $\bar{B}(x; r)$. This can be formally written as:

$$\bar{B}(x; r) = \{z \in H \mid \|z - x\| \leq r\}$$

Definition 2.1.3. Given $\alpha \in [0, 1]$; an operator $T: H \rightarrow H$ is

- α -averaged non-expansive, if for all $(x; z) \in H \times H$;

$$\|Tx - Tz\| \leq \alpha \|x - z\| + (1 - \alpha) \|x - Tz\|$$

- Firmly nonexpansive if it is $\frac{1}{2}$ average.
- Non-expansive if, for all $(x; z) \in H \times H$;

$$\|Tx - Tz\| \leq \|x - z\|$$

- Quasi-non-expansive if

$$\|Tx - z\| \leq \|x - z\|$$

for all $x \in H$ and all $z \in \text{Fix } T$, where $\text{Fix } T = \{x \in H \mid Tx = x\}$ is the set of fixed points of T .

The next lemma allows us to bound a positive sequence (see [111, Lemma A.1] for further references).

Lemma 2.1.4. Assume that the sequence $(f_k)_{k \in \mathbb{N}}$ is a non-negative and satisfies the recursion

$$f_{N+1} = f_N + \sum_{k=1}^N f_k \tag{2.1.1}$$

for all $N \in \mathbb{N}$, where $(f_k)_{k \in \mathbb{N}}$ is an increasing sequence, $f_0 = 0$, and $f_k \geq 0$ for all $k \in \mathbb{N}$. Then for all $N \in \mathbb{N}$

$$f_N \leq \frac{1}{2} \sum_{k=1}^N f_k + \frac{1}{2} \sum_{k=1}^N f_k \tag{2.1.2}$$

The following lemma allows us to prove convergence of a given sequence.

Lemma 2.1.5 (Opial's Lemma). Let C be a nonempty subset of a finite dimensional Hilbert space H and let $(x_k)_{k \in \mathbb{N}}$ be a sequence in H : Assume

- For every $x \in C$ there exists $\lim_{k \rightarrow \infty} \|x_k - x\| = 0$; and
- Every limit point of $(x_k)_{k \in \mathbb{N}}$ belongs to C .

Then $(x_k)_{k \in \mathbb{N}}$ converges as $k \rightarrow \infty$ to some $x \in C$:

2.2 Matrices

This section introduces various fundamental definitions and results in matrix theory and linear algebra.

We use lowercase and uppercase letters, for matrices and vectors, respectively. For any positive integer p , the set of all integers from 1 to p is denoted by $[p]$. The Hadamard product of two vectors x and z in \mathbb{R}^p is represented as $x \odot z$, where, for every i in $[p]$, $(x \odot z)_i = x_i z_i$. For a positive scalar L , the vector $\mathbf{1}^L := \mathbf{1} \odot L$. The vector $\mathbf{1}$ in \mathbb{R}^p has all components equal to 1. The diagonal matrix $\text{Diag}(x) \in \mathbb{R}^{p \times p}$, for a vector $x \in \mathbb{R}^p$, has the components of x on its main diagonal and zeros elsewhere.

The operator norm for a matrix $X \in \mathbb{R}^{d \times p}$ is defined by

$$\|X\| = \sup_{\substack{z \in \mathbb{R}^p \\ \|z\| = 1}} \|Xz\|$$

The Frobenius norm of X , represented as $\|X\|_F$, is defined as

$$\|X\|_F^2 := \sum_{i=1}^d \|x_i\|^2;$$

where x_i is the i -th row of X . This norm extends the Euclidean norm to matrices.

The adjoint (or transpose) of a matrix $X \in \mathbb{R}^{d \times p}$, denoted by X^\top , is the unique matrix satisfying $\langle Xz, y \rangle = \langle z, X^\top y \rangle$, for every $z \in \mathbb{R}^p$ and $y \in \mathbb{R}^d$. Two vectors $(x; z) \in \mathbb{R}^p \times \mathbb{R}^p$ are orthogonal if and only if $\langle x; z \rangle = 0$. For a subset $V \subseteq \mathbb{R}^p$, its orthogonal complement is defined as:

$$V^\perp = \{y \in \mathbb{R}^p \mid \langle y, z \rangle = 0 \text{ for all } z \in V\}$$

The range of $X \in \mathbb{R}^{d \times p}$ is denoted by $\text{ran}(X)$ and defined as $\{Xz \mid z \in \mathbb{R}^p\}$. Similarly, the kernel of X is denoted by $\text{ker}(X)$ and defined as $\{z \in \mathbb{R}^p \mid Xz = 0\}$. In the finite-dimensional setting, it follows that both $\text{ker}(X)$ and $\text{ran}(X)$ are closed subspaces, and we have the relationship $\text{ran}(X) = \text{ker}(X)^\perp$: An orthogonal matrix, or orthonormal matrix, is a real square matrix $X \in \mathbb{R}^{p \times p}$ whose columns and rows are orthonormal vectors.

A matrix $M \in \mathbb{R}^{p \times p}$ is positive definite if $\langle Mx, x \rangle > 0$ for every non-zero vector $x \in \mathbb{R}^p$. Similarly, a matrix is positive semidefinite if $\langle Mx, x \rangle \geq 0$ for every vector $x \in \mathbb{R}^p$. The set of all $p \times p$ positive semidefinite and positive definite matrices are represented by S_+^p and S_{++}^p , respectively.

Definition 2.2.1. Given a matrix $X \in \mathbb{R}^{p \times d}$, its Singular Value Decomposition (SVD) is a factorization of the form $X = U \Sigma V^\top$, where:

- $U \in \mathbb{R}^{p \times p}$ and $V \in \mathbb{R}^{d \times d}$ are orthogonal matrices.
- $\Sigma \in \mathbb{R}^{p \times d}$ is a diagonal matrix with non-negative real numbers on the diagonal.

These diagonal entries of Σ are known as the singular values of X , and the columns of U and V are called the left-singular vectors and right-singular vectors of X , respectively.

Definition 2.2.2. Given a matrix $X \in \mathbb{R}^{p \times d}$, if the singular values of the X matrix are denoted by σ_i for every $i \in [1; p]$, then the nuclear-norm $\|X\|_1$ is defined by:

$$\|X\|_1 = \sum_{i=1}^{\min\{p,d\}} \sigma_i$$

The next formula provides an efficient method for calculating the inverse of a matrix plus a rank-one matrix, specifically for $(M + z z^T)$ where M is an invertible $p \times p$ matrix and $z \in \mathbb{R}^p$.

Lemma 2.2.3. If M is an invertible $p \times p$ matrix and z and β are two p -dimensional vectors such that $\beta^T M^{-1} z \neq -1$, the Sherman–Morrison formula gives

$$(M + z z^T)^{-1} = M^{-1} - \frac{M^{-1} z z^T M^{-1}}{1 + z^T M^{-1} z}. \quad (2.2.1)$$

The following definition will be used in the numerical examples.

Definition 2.2.4. The discrete gradient operator $D: \mathbb{R}^{p \times p} \rightarrow (\mathbb{R}^2)^{p^2}$, for every matrix $u \in \mathbb{R}^{p \times p}$ is defined by:

$$(Du)_{ij} = ((D_x u)_{ij}; (D_y u)_{ij});$$

where

$$(D_y u)_{ij} = \begin{cases} u_{i+1;j} - u_{i;j} & \text{if } 1 \leq i < N; \\ 0 & \text{if } i = N; \end{cases}$$

$$(D_x u)_{ij} = \begin{cases} u_{i;j+1} - u_{i;j} & \text{if } 1 \leq j < N; \\ 0 & \text{if } j = N; \end{cases}$$

2.3 Convex optimization

This section introduces various fundamental definitions and results in convex optimization. For further results on convex analysis and operator theory, the reader is referred to [15, 107].

Definition 2.3.1. Let $f: H \rightarrow [1; +\infty]$. The domain of f is

$$\text{dom}(f) = \{x \in H \mid f(x) < +\infty\};$$

the graph of f is

$$\text{gra}(f) = \{(x; g) \in H \times \mathbb{R} \mid f(x) = g\};$$

the epigraph of f is

$$\text{epi}(f) = \{(x; g) \in H \times \mathbb{R} \mid f(x) \leq g\};$$

The function f is proper if $1 \in \text{epi}(f)$ and $\text{dom}(f) \neq \emptyset$.

Definition 2.3.2. The argmin of a proper function $f: H \rightarrow [1; +\infty]$ is

$$\text{argmin}(f) := \{x \in H \mid f(x) = \inf_{x \in H} f(x)\} \text{ for all } x \in H;$$

Definition 2.3.3. A subset $C \subseteq H$ is convex if and only if

$$\lambda x_1 + (1 - \lambda) x_2 \in C$$

for every $x_1, x_2 \in C$ and $\lambda \in (0; 1)$.

Definition 2.3.4. A function $f: H \rightarrow \mathbb{R} \cup \{+\infty\}$ is convex if

$$f(\alpha x_1 + (1 - \alpha)x_2) \leq \alpha f(x_1) + (1 - \alpha)f(x_2)$$

for every $x_1, x_2 \in \text{dom}(f)$ and $\alpha \in (0, 1)$.

Definition 2.3.5. A function $f: H \rightarrow \mathbb{R} \cup \{+\infty\}$ is lower semicontinuous (l.s.c.) if and only if $\text{epi}(f)$ is closed in $H \times \mathbb{R}$.

We denote by $\mathcal{C}_0(H)$ the set of convex, lower semicontinuous, and proper functions on H .

Definition 2.3.6. Let $f: H \rightarrow \mathbb{R} \cup \{+\infty\}$ be a proper and convex function. A point $z \in H$ is a subgradient of f at x if, for all $z \in H$;

$$f(z) \geq f(x) + \langle z, z - x \rangle$$

The set of all subgradients of f at x is the subdifferential of f at x and is denoted by $\partial f(x)$. If $\partial f(x) \neq \emptyset$; we say that f is subdifferentiable at x . The domain of the subdifferential is

$$\text{dom}(\partial f) = \{x \in H \mid \partial f(x) \neq \emptyset\}$$

Note that, by definition,

$$\text{dom}(\partial f) \subseteq \text{dom}(f)$$

Proposition 2.3.7. Let $f: H \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. If $x_1 \in \partial f(x_1)$ and $x_2 \in \partial f(x_2)$, then $\langle x_1 - x_2, x_1 - x_2 \rangle \leq 0$:

Example 2.3.8. Let $f: \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$; $f(x) = |x|$. Then, $\partial f(x) = \{1\}$ if $x > 0$, $\partial f(x) = \{-1\}$ if $x < 0$, and $\partial f(0) = [-1, 1]$:

Example 2.3.9. Let $f: H \rightarrow \mathbb{R} \cup \{+\infty\}$; $f(x) = \frac{1}{k} \|x\|$. Then, $\partial f(x) = \frac{1}{k} \frac{x}{\|x\|}$ if $x \neq 0$, and $\partial f(0) = \overline{B}(0, 1)$, where $\|\cdot\|$ is the norm induced by the scalar product.

Example 2.3.10. Let C a subset of H . Define the indicator function of C as

$$c(x) = \begin{cases} 0; & \text{if } x \in C; \\ +\infty; & \text{if } x \notin C; \end{cases}$$

If $C \subseteq H$ is a nonempty, convex, and closed subset, then $c \in \mathcal{C}_0(H)$. Moreover, $\partial c = N_C$, where, for every $x \in H$;

$$N_C(x) = \begin{cases} \{z \in H \mid \langle z, x - y \rangle \leq 0 \text{ for all } y \in C\} & \text{if } x \in C; \\ \emptyset & \text{if } x \notin C; \end{cases}$$

The following definition allows us to define a Hilbert space by the product of a finite number of Hilbert spaces.

Definition 2.3.11. The Hilbert direct sum of a family of real Hilbert spaces $(H_i)_{i \in [d]}$ is the real Hilbert space:

$$\bigoplus_{i \in [d]} H_i = \left\{ (x_i)_{i \in [d]} \mid x_i \in H_i \right\}$$

For every $i \in [d]$; let $f_i: H_i \rightarrow \mathbb{R} \cup \{+\infty\}$. Then,

$$\bigoplus_{i \in [d]} f_i: \bigoplus_{i \in [d]} H_i \rightarrow \mathbb{R} \cup \{+\infty\}; \quad (x_i)_{i \in [d]} \mapsto \sum_{i=1}^d f_i(x_i)$$

Proposition 2.3.12. Consider a finite family of functions $f_i, g_i, i \in [d]$ such that, for every $i \in [d]$, $f_i: H_i \rightarrow \mathbb{R} \cup \{+\infty\}$ is proper and convex. Then,

$$\bigoplus_{i \in [d]} f_i = \bigoplus_{i \in [d]} f_i;$$

The following lemma allows us to characterize the minimizer of a function.

Theorem 2.3.13. Let $f: H \rightarrow \mathbb{R} \cup \{+\infty\}$ be a proper and convex function. Then

$$\operatorname{argmin} f = \{x \in \mathbb{R}^p \mid 0 \in \partial f(x)\};$$

For a given set C , we denote by \bar{C} the smallest closed subset containing C . The relative interior of C is

$$\operatorname{ri}(C) = \{x \in C \mid \exists \epsilon > 0, \operatorname{span}(C - x) \subseteq \operatorname{span}(C - x)\};$$

where

$$\mathbb{R}_{++}C = \{y \mid \exists \lambda > 0, y \in \lambda C\}$$

and $\operatorname{span}(C)$ is the smallest linear subspace of H containing C . From now on, every time we use ri we assume that we are in a finite dimensional space.

Proposition 2.3.14. Let $f: H \rightarrow \mathbb{R} \cup \{+\infty\}$, and $g: G \rightarrow \mathbb{R} \cup \{+\infty\}$ be two proper, l.s.c., and convex functions. Let $X: H \rightarrow G$ a nonzero, bounded, and linear operator. Suppose that $\operatorname{dom} g \setminus X \operatorname{dom} f \neq \emptyset$. Then, $\partial(f + X^*g) = \partial(f + g \circ X)$: Additionally, if we suppose that $0 \in \operatorname{ri}(\operatorname{dom} g \circ X \operatorname{dom} f)$; then $\partial(f + X^*g) \circ X = \partial(f + g \circ X)$:

Definition 2.3.15. Let $f: H \rightarrow \mathbb{R} \cup \{+\infty\}$ be in $\mathcal{C}_0(H)$. The conjugate (or Fenchel conjugate, or Legendre transform, or Legendre–Fenchel transform) of f is

$$f^*: H^* \rightarrow \mathbb{R} \cup \{+\infty\}, f^*(u) := \sup_{x \in H} \langle u, x \rangle - f(x);$$

Proposition 2.3.16. Consider a finite family of functions $f_i, g_i, i \in [d]$ such that, for every $i \in [d]$, $f_i: H_i \rightarrow \mathbb{R} \cup \{+\infty\}$ is proper and convex. Then,

$$\bigoplus_{i \in [d]} f_i^* = \left(\bigoplus_{i \in [d]} f_i\right)^*;$$

Example 2.3.17. Let C be a nonempty, closed, and convex set. Then,

$$(C^*)^* = C; \quad f^*(u) = \sup_{x \in C} \langle u, x \rangle;$$

Note that if $C = \{y \mid g(y) \leq 1\}$; then $f^*(u) = \sup_{y \in H} \langle u, y \rangle - g(y)$:

Example 2.3.18. Set $B_1 = \{x \in \mathbb{R}^n \mid \max_{i \in [n]} |x_i| \leq 1\}$. Then, $(B_1^*)^* = B_1$.

Proposition 2.3.19. Let $f \in \mathcal{C}_0(H)$. Then $(f^*)^* = f$, where the inverse of f^* , is defined through its graph

$$\operatorname{gra} (f^*)^* = \{(u, x) \mid x \in H, u \in H^*, \langle u, x \rangle = f(x)\}$$

Definition 2.3.20. Let $f: H \rightarrow \mathbb{R} \cup \{+\infty\}$ and $g: G \rightarrow \mathbb{R} \cup \{+\infty\}$ be two proper, l.s.c., and convex functions. Let $X: H \rightarrow G$ a nonzero, bounded, and linear operator. The primal problem is given by:

$$\inf_{x \in H} f(x) + g(Xx)$$

Its optimal value is denoted by p^* and the set of primal solutions is S . The dual problem is given by:

$$\sup_{u \in H} f^*(X^*u) + g^*(u)$$

Its optimal value is denoted by d^* and the set of dual solutions is S^* .

Theorem 2.3.21. Let $f: H \rightarrow \mathbb{R} \cup \{+\infty\}$ and $g: G \rightarrow \mathbb{R} \cup \{+\infty\}$ be two proper, l.s.c., and convex functions. Let $X: H \rightarrow G$ a nonzero, bounded, and linear operator. Then, $p^* = d^*$ if and only if $0 \in \text{ri}(\text{dom } g - X \text{dom } f)$;

then $p^* + d^* = 0$:

The next theorem characterizes the primal-dual solutions.

Theorem 2.3.22. Let $x \in H$ and $u \in G$. Then the following are equivalent:

1. $X^*u \in \partial f(x)$ and $u \in \partial g(Xx)$;
2. $f(x) + g(Xx) + f^*(X^*u) + g^*(u) = 0$; and
3. $x \in S$ and $u \in S^*$ and $p^* + d^* = 0$;

Definition 2.3.23. The proximal operator of $f \in \Gamma_0(H)$ is defined by

$$\text{prox}_f: H \rightarrow H, \quad \forall x \in H, \quad \text{prox}_f(x) = \arg \min_{z \in H} f(z) + \frac{1}{2} \|x - z\|^2$$

Proposition 2.3.24. Let $f \in \Gamma_0(H)$; and let x and p be in H : Then $p = \text{prox}_f(x)$ if and only if $x \in \partial f(p) + p$:

For every self-adjoint positive definite matrix K , we define the proximity operator of f relative to the metric induced by $\langle x, y \rangle_K := \langle Kx, y \rangle$ as $\text{prox}_f^K = (\text{Id} + \partial f)^{-1}_K$. If $K = \text{Id}$ for some real number $\alpha > 0$, it is customary to write prox_f rather than prox_f^K .

Proposition 2.3.25. Let $f \in \Gamma_0(H)$ and let $\alpha > 0$. Then $\text{prox}_f^\alpha(x) = \text{prox}_f(x/\alpha)$:

Proposition 2.3.26. Consider a finite family of functions $f_i \in \Gamma_0(H_i)$ such that, for every $i \in [d]$, $f_i: H_i \rightarrow \mathbb{R} \cup \{+\infty\}$ is in $\Gamma_0(H_i)$. Set $H = \prod_{i \in [d]} H_i$ and $f = \sum_{i \in [d]} f_i$. Then, $\text{prox}_f(x) = \prod_{i \in [d]} \text{prox}_{f_i}(x_i)$.

Example 2.3.27. Let $\alpha > 0$ and $x \in \mathbb{R}^p$. Then,

$$\text{prox}_{\frac{1}{2\alpha} \|\cdot\|^2}(x) = (\text{soft}(x/\alpha))_{i \in [p]}$$

where

$$\text{soft}(x) = \begin{cases} \alpha & \text{if } x < -\alpha; \\ x & \text{if } x \in [-\alpha, \alpha]; \\ -\alpha & \text{if } x > \alpha; \end{cases}$$

and

$$\text{prox}_{\frac{1}{2\alpha} \|\cdot\|^2}(x) = \frac{1}{\alpha} \frac{x}{\sqrt{1 + \|x\|^2/\alpha^2}}$$

Definition 2.3.28. The projection operator onto a non-empty closed convex set $C \subset H$ is defined by

$$P_C: H \rightarrow H; \quad \forall z \in H \quad \arg\min_{z \in C} \|z - z\|^2$$

Let $f = \{z \in H \mid \|z - z\|^2 = \min_{z \in C} \|z - z\|^2\}$, then $\text{prox}_C = P_C$. The distance to C is defined as $\text{dist}^2(z; C) = \min_{z \in C} \|z - z\|^2 = \|z - P_C z\|^2$.

Example 2.3.29. Let x be a non-zero vector in H , let $y \in \mathbb{R}$, and set

$$C = \{z \in H \mid \langle z, x \rangle = y\}$$

Then

$$P_C z = z + \frac{y - \langle z, x \rangle}{\|x\|^2} x$$

Example 2.3.30. Let x be a non-zero vector in H , let $(y_1, y_2) \in \mathbb{R}^2$, such that $y_2 \leq y_1$, and set $C = \{z \in H \mid y_1 \leq \langle z, x \rangle \leq y_2\}$. Then

$$P_C z = \begin{cases} z + \frac{y_1 - \langle z, x \rangle}{\|x\|^2} x; & \text{if } \langle z, x \rangle < y_1; \\ z; & \text{if } y_1 \leq \langle z, x \rangle \leq y_2; \\ z + \frac{y_2 - \langle z, x \rangle}{\|x\|^2} x; & \text{if } \langle z, x \rangle > y_2; \end{cases}$$

Theorem 2.3.31. Let C be a nonempty closed and convex subset of H . Then, for every $z \in H$; $P_C z \in C$, and for every $z \in C$,

$$\langle z - P_C z, z - P_C z \rangle = 0$$

Additionally, if C is a closed affine subspace, then P_C is an affine operator, and for every $z \in C$, $\langle z - P_C z, z - P_C z \rangle = 0$.

In the following two definitions and two propositions, we assume that H and G are finite dimensional.

Definition 2.3.32. Let $X: H \rightarrow G$ a nonzero, bounded, and linear operator and let $y \in G$. Then z is a least square solution to the linear system $Xz = y$ if

$$z \in \arg\min_{z \in H} \|Xz - y\|^2$$

Proposition 2.3.33. Let $X: H \rightarrow G$ a nonzero, bounded, and linear operator and let $y \in G$. Then the equation $Xz = y$ has at least one least squares solution. Moreover, for every $z \in H$, the following are equivalent:

1. z is the least square solution.
2. $Xz = P_{\text{ran}(X)} y$.
3. $X^> X z = X^> y$.

Definition 2.3.34. Let $X: H \rightarrow G$ a nonzero, bounded, and linear operator. Let $C_y = \{z \in H \mid X^> X z = X^> y\}$: The generalized (or Moore-Penrose) inverse of X is $X^y: G \rightarrow H$: $y \in G \mapsto P_{C_y}(0)$.

Proposition 2.3.35. Let $X: H \rightarrow G$ a nonzero, bounded, and linear operator, such that $\text{ran}(X)$ is closed. Then X^y is a bounded linear operator, $P_{\text{ran}(X)} = X X^y$, and $P_{\ker(X)} = \text{Id} - X^> (X^>)^y$.

2.4 Differential Calculus and Ordinary differential equations

In this section, we recall some basic facts of Differential Calculus and Ordinary differential equations.

Definition 2.4.1. Let $f: H \rightarrow \mathbb{R} [f+1 g$. The directional derivative of $f: H \rightarrow \mathbb{R} [f+1 g$ at $x \in \text{int}(\text{dom}(f))$ in the direction $d \in H$ is:

$$f^0(x; d) = \lim_{t \rightarrow 0^+} \frac{f(x + td) - f(x)}{t};$$

whenever the limit exists. The function f is Gâteaux differentiable at the point x if $f^0(x; d)$ exists for all $d \in H$ and the function $d \mapsto f^0(x; d)$ is linear and continuous. In this situation, the Gâteaux derivative (or gradient) of f at x is $\nabla f(x) = f^0(x; \cdot)$.

As before, we can define a directional derivative of ∇f at $x \in \text{int}(\text{dom}(f))$ in the direction $d \in H$ as:

$$(\nabla f)^0(x; d) = \lim_{t \rightarrow 0^+} \frac{\nabla f(x + td) - \nabla f(x)}{t}$$

whenever the limit exists. The function f is twice Gâteaux differentiable at the point x if $(\nabla f)^0(x; d)$ exists for all $d \in H$ and the function $d \mapsto (\nabla f)^0(x; d)$ is linear and continuous. In this situation, the second Gâteaux derivative (or hessian) of f at x is denoted by $\nabla^2 f(x) = (\nabla f)^0(x; \cdot)$.

Proposition 2.4.2. Let $f: H \rightarrow \mathbb{R} [f+1 g$ be proper and convex function, and let $x \in \text{dom}(f)$. If the function f is Gâteaux differentiable at the point x , then $\nabla f(x) = \nabla f(x)$.

Lemma 2.4.3 (Descent Lemma). If $f: H \rightarrow \mathbb{R} [f+1 g$ is Gâteaux-differentiable and ∇f is Lipschitz-continuous with constant L , then:

$$f(z) \leq f(x) + \langle \nabla f(x), z - x \rangle + \frac{L}{2} \|z - x\|^2;$$

for every $(z; x) \in \text{dom}(f) \times \text{int}(\text{dom}(f))$:

Definition 2.4.4. Let $f: H \rightarrow \mathbb{R} [f+1 g$ be strictly convex, proper, and Gâteaux differentiable on $\text{int}(\text{dom} f)$. The Bregman divergence D_f associated with f is defined as:

$$D_f: H \times H \rightarrow \mathbb{R} [f+1 g$$

$$(z; x) \mapsto \begin{cases} f(z) - f(x) - \langle \nabla f(x), z - x \rangle & \text{if } z \in \text{int}(\text{dom}(f)); \\ +\infty & \text{otherwise.} \end{cases} \tag{2.4.1}$$

Note that, for every $(z; x) \in \text{dom}(f) \times \text{int}(\text{dom}(f))$; the Bregman distance $D_f(z; x)$ is positive and $D_f(z; x) = 0$ if and only if $z = x$.

The next Lemma provides a sufficient condition to determine if a vector field is a gradient.

Lemma 2.4.5. Let $f = (f_1; \dots; f_p): U \rightarrow \mathbb{R}^p \rightarrow \mathbb{R}^p$ let be a continuously differentiable vector field on an open set U . If f is a gradient on U , then the partials derivatives of the components of f are related by the equations:

$$\frac{\partial f_i}{\partial x_j}(x) = \frac{\partial f_j}{\partial x_i}(x);$$

for all $i; j \in [p]$ and $x \in U$.

Definition 2.4.6. Let $U \subseteq \mathbb{R}^p$ be a non-empty open set. Let a function $f: U \rightarrow \mathbb{R} [f+1]g$ such that at each point of the set U all partial derivatives are continuous. Then we say that f is of the class C^1 on U . The set of all these functions is denoted by $C^1(U)$.

Analogously, a function $f: U \rightarrow \mathbb{R} [f+1]g$ is said to be of the class C^2 on U if it has continuous second partial derivatives at each point in the set U . The set of all such functions is denoted by $C^2(U)$.

Definition 2.4.7. Let $f: \mathbb{R}^p \rightarrow \mathbb{R} [f+1]g$ in $C^1(\mathbb{R}^p)$ such that $U := \text{int}(\text{dom}(f)) \neq \emptyset$. We say that f is essentially smooth if it satisfies the following two conditions

1. f is differentiable on U ;
2. For every sequence $\{x_k\}_{k=1}^{+\infty} \subset U$ converging to a point in $\partial U := U \cap \text{int}(U)$, we get that

$$\lim_{k \rightarrow +\infty} \frac{\|f(x_k) - f(x) - Df(x)(x_k - x)\|}{\|x_k - x\|} = 0$$

In addition, f is a Legendre type (or simply Legendre) function if it is essentially smooth and strictly convex on U .

For this part, we follow [11]. Now, we introduce the notion of a solution of an ordinary differential equation.

Definition 2.4.8. Let $U \subseteq \mathbb{R} \times \mathbb{R}^p$ be an open set. Let $f: U \rightarrow \mathbb{R}$ be a continuous function. Consider the (ordinary) differential equation:

$$\dot{x}(t) = f(t, x(t)) \quad (2.4.2)$$

A function $x: (t_-, t_+) \rightarrow \mathbb{R}^p$ of class C^1 (with $t_- > -\infty$ and $t_+ < +\infty$) is said to be a solution of (2.4.2) if:

1. $(t, x(t)) \in U$, for every $t \in (t_-, t_+)$;
2. $\dot{x}(t) = f(t, x(t))$, for every $t \in (t_-, t_+)$.

The ordinary differential equation (2.4.2) is said to be autonomous if f does not depend on t and $U \subseteq \mathbb{R}^p$, which is the setting that we consider from now on. Next we introduce the notion of initial value problem.

Definition 2.4.9. Given $(t_0, x_0) \in (t_-, t_+) \times U$, the initial value problem

$$\dot{x}(t) = f(x(t)); \quad x(t_0) = x_0 \quad (2.4.3)$$

consists of finding an interval (t_-, t_+) containing t_0 and a solution $x: (t_-, t_+) \rightarrow U$ of (2.4.2) such that $x(t_0) = x_0$. The condition $x(t_0) = x_0$ is called the initial condition of problem (2.4.3).

Proposition 2.4.10. Let $U \subseteq \mathbb{R}^p$ be an open set. Let $f: U \rightarrow \mathbb{R}$ be a continuous function. Given $(t_0, x_0) \in (t_-, t_+) \times U$, a continuous function $x: (t_-, t_+) \rightarrow \mathbb{R}^p$ in a open interval (t_-, t_+) containing t_0 is a solution of the initial value problem (2.4.3) if and only if

$$x(t) = x_0 + \int_{t_0}^t f(x(s)) ds;$$

for every $t \in (t_-, t_+)$:

Proposition 2.4.11. Let $\gamma : [t_0; +\infty) \rightarrow \mathbb{R}^p$ with $\gamma(t_0) = \gamma_0$ be a continuously differentiable function, whose derivative satisfies

$$\gamma'(t) = u(t) \gamma(t);$$

for an integrable function $u : [t_0; +\infty) \rightarrow (\mathbb{R}^{p \times p}; 0]$. Then we have

$$\gamma(t) = \gamma_0 \exp \int_{t_0}^t u(s) ds;$$

Definition 2.4.12. Let $U \subset \mathbb{R}^p$ be an open set. A function $f : U \rightarrow \mathbb{R}^p$ is said to be locally Lipschitz if for every compact set $K \subset U$ there exists $L > 0$ such that

$$\forall x, z \in K \quad \|f(x) - f(z)\| \leq L \|x - z\|;$$

Theorem 2.4.13. Let $U \subset \mathbb{R}^p$ be an open set. Let $f : U \rightarrow \mathbb{R}^p$. If f is continuous and locally Lipschitz function in U , then for every $(t_0; \gamma_0) \in \mathbb{R} \times U$ there exists an open interval containing t_0 such that the solution of the initial value problem (2.4.3) is unique.

Definition 2.4.14. Let $U \subset \mathbb{R}^p$ be an open set. Let $f : U \rightarrow \mathbb{R}^p$ a continuous and locally Lipschitz in an open set $U \subset \mathbb{R}^p$. The maximal interval of a solution $\gamma : I \rightarrow \mathbb{R}^p$ of the equation $\gamma'(t) = f(\gamma(t))$ is the largest open interval $(t_-; t_+)$ where there exists a solution.

Definition 2.4.15. If $\gamma = \gamma(t)$ is a solution of (2.4.2) with maximal interval $(t_-; t_+)$, then the set $\gamma : (t_-; t_+) \rightarrow U$ is called trajectory of the equation. This solution is said to be global if its maximal interval is \mathbb{R} .

Proposition 2.4.16. Any solution of equation (2.4.2) whose trajectory is contained in a compact subset of U is global.

Definition 2.4.17. A point $\gamma_1 \in U$ with $f(\gamma_1) = 0$ is called stationary point of (2.4.2).

Theorem 2.4.18 (Cauchy, Lipschitz, Picard). Let $f : \mathbb{R}^p \rightarrow \mathbb{R}^p$ be a Lipschitz continuous function in \mathbb{R}^p . Then for every $\gamma_0 \in \mathbb{R}^p$ there exists a unique solution $\gamma : [0; +\infty) \rightarrow \mathbb{R}^p$ in $C^1(\mathbb{R}^p)$ of the initial value problem (2.4.3).

The following definitions are collected from [77]. From now on, the map

$$\gamma : \mathbb{R}^k \rightarrow \mathbb{R}^p \quad \gamma = q(\gamma);$$

is called reparameterization.

Definition 2.4.19. The Jacobian of $q = (q_1; \dots; q_p)$ in \mathbb{R}^k is denoted by $J_q(\gamma)$ and is defined as

$$J_q := \begin{pmatrix} \frac{\partial q_1}{\partial \gamma_1} & \dots & \frac{\partial q_1}{\partial \gamma_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial q_p}{\partial \gamma_1} & \dots & \frac{\partial q_p}{\partial \gamma_k} \end{pmatrix} := \begin{pmatrix} 1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1 \end{pmatrix};$$

Definition 2.4.20. A regular parameterization $q : \mathbb{R}^k \rightarrow \mathbb{R}^p$ is a C^1 parameterization such that $J_q(\gamma)$ is of rank p for all $\gamma \in \mathbb{R}^k$.

Definition 2.4.21. A C^2 parameterization $q : \mathbb{R}^k \rightarrow \mathbb{R}^p$ is commuting if and only if for any $i, j \in [p]$, we have that,

$$r^2 q_j(\gamma) \cdot r q_i(\gamma) - r^2 q_i(\gamma) \cdot r q_j(\gamma) = 0;$$

for all $\gamma \in \mathbb{R}^k$.

Definition 2.4.22. For any C^1 function $f: \mathbb{R}^k \rightarrow \mathbb{R} [f+1 g$, we denote by ${}^t_f(x_0) = (t)$; where (t) is the solution at time t (when it exists) of

$$\begin{cases} \dot{(t)} = f((t)); & t > 0; \\ (0) = x_0 \in \mathbb{R}^k; \end{cases} \quad (2.4.4)$$

We say ${}^t_f(x_0)$ is well-defined at time t when the above differential equation has a solution at time t .

Definition 2.4.23. Given a C^2 parameterization $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$, for any $x_0 \in \mathbb{R}^p$ and $t \in \mathbb{R}^p$, we define

$$(x_0; t) := \begin{matrix} t_1 \\ \vdots \\ t_p \end{matrix} \begin{matrix} q_1(x_0) \\ \vdots \\ q_p(x_0) \end{matrix};$$

when it is well-defined, i.e., the corresponding differential equations have a solution. For any $x_0 \in \mathbb{R}^p$, we define the domain of $(x_0; \cdot)$ as:

$$U(x_0) = \{t \in \mathbb{R}^p \mid (x_0; t) \text{ is well defined}\};$$

Definition 2.4.24. For any parameterization $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$ in C^2 and for any function $L: \mathbb{R}^p \rightarrow \mathbb{R} [f+1 g$ in C^1 , given any starting point $x_0 \in \mathbb{R}^p$, we define the reachable set $(x_0; q)$ as

$$(x_0; q) = \{L(q(x_0)) \mid t > 0\};$$

For further results, the reader is referred to [3, 11, 15, 77].

CHAPTER 3

Regularization techniques

In this Chapter, we recall key ideas at the basis of our study, in particular with respect to the design regularization techniques.

3.1 Inverse problems

Many applied problems require estimating a function of interest from input/output data relation,

$$(x_i, y_i)_{i=1}^d \in \mathbb{R}^p \times \mathbb{R} \quad \forall \quad f: \mathbb{R}^p \rightarrow \mathbb{R}$$

In this chapter, we are interested in the case when f is a linear function, expressed as

$$f: \mathbb{R}^p \rightarrow \mathbb{R}; \quad x \mapsto \sum_{j=1}^p h_j x_j$$

for some $h \in \mathbb{R}^p$: Then, the problem simplifies to determining h from the following linear equation

$$Xh = y;$$

where X represents a matrix with rows corresponding to input data points, and y can be seen as the vector of measurements of the vector h that we want to recover.

Observe that, if $y \notin \text{ran}(X)$, the previous problem does not have a solution. For this reason, it is necessary to introduce a loss function $L: \mathbb{R}^p \rightarrow \mathbb{R} [f+1]g$ such that it ensures the existence of at least one solution, for example, least squares. In general, this function is defined via a function $\ell: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, through the following equation:

$$L(h) = \sum_{i=1}^d \ell(hx_i, y_i)$$

On the other hand, the solution of the linear system may be not unique and a selection criterion is needed to choose a suitable solution. To deal with non-uniqueness and select a particular solution, a regularization term $R: \mathbb{R}^p \rightarrow \mathbb{R} [f+1]g$ is introduced. This regularization includes the prior knowledge on the solution in the model, and enforces a bias on the solution. Summarizing, a standard approach to recover h is to assume that it is a minimizer of the following linearly constrained optimization problem:

$$\min_{h \in \mathbb{R}^p} fR(h) : \quad 2 S = \text{argmin} Lg; \quad (P)$$

The above ideas can be justified from different perspectives. For instance, compressed sensing [32, 34, 48, 115, 129], image processing [37, 39, 103, 105, 116, 117, 140], and various problems in machine learning [13, 50, 90, 114, 121, 140, 141].

Example 3.1.1. In this example, we show the most common losses used both in regression and classification. In classification tasks, loss functions are typically based on the concept of margin, i.e., $y^h - j x_i$, while in regression, they are often based on the difference $y - h - j x_i$. Some well-known examples include:

- **Least Squares:** It is defined as $L(\theta) = \frac{1}{d} \sum_{i=1}^d (h - j x_i - y_i)^2$. This loss function measures the average squared difference between the predicted values and the actual values.
- **Huber Loss:** A piece-wise function that combines squared error and absolute error, defined as $L(\theta) = \frac{1}{d} \sum_{i=1}^d \rho(h - j x_i - y_i)$, where ρ is defined as follow

$$\rho(\delta) = \begin{cases} \frac{1}{2} \delta^2 & \text{for } |\delta| \leq \frac{1}{2} \\ (|\delta| - \frac{1}{2}) & \text{otherwise.} \end{cases}$$

- **Exact Penalization:** It is defined as

$$L(\theta) = \begin{cases} 0 & \text{if } X = y; \\ + \gamma & \text{otherwise;} \end{cases}$$

This loss is the one we will use throughout this thesis.

- **Hinge Loss:** Typically used in Support Vector Machines (SVMs) for classification, it is defined as $L(\theta) = \frac{1}{d} \sum_{i=1}^d \max(0; 1 - y_i h - j x_i)$. Here, for every $i \in [d]$, $y_i \in \{-1; 1\}$.
- **Exponential Loss:** It is defined as $L(\theta) = \frac{1}{d} \sum_{i=1}^d e^{-y_i h - j x_i}$. Here, for every $i \in [d]$, $y_i \in \{-1; 1\}$.
- **Logistic Loss:** This is widely used in logistic regression for binary classification. It can be expressed as $L(\theta) = \frac{1}{d} \sum_{i=1}^d \log(1 + e^{-y_i h - j x_i})$, where, for every $i \in [d]$, $y_i \in \{-1; 1\}$.

Example 3.1.2. We illustrate some examples of regularizers presented in the state of the art [32, 34, 52, 105, 116, 117, 148].

- **Squared norm:** It is defined as $R(\theta) = k \|\theta\|_2^2$. This regularizer leads to the minimum norm solution that is the Moore-Penrose pseudoinverse $X^+ y$ [52, Theorem 2.5], when the loss is the least square.
- **ℓ_0 -norm:** It is defined as $R(\theta) = k \|\theta\|_0$, where $\|\theta\|_0$ is the function that counts the non-zero elements of the vector.
- **ℓ_1 -norm:** It is defined as $R(\theta) = k \|\theta\|_1$. This function is a convex surrogate for the ℓ_0 -norm, and promotes sparsity of the solution.
- **Nuclear norm:** It is defined as $R(\theta) = k \|\theta\|_*$, where $\theta \in \mathbb{R}^{p \times p}$ is a square matrix. This function is a convex surrogate for the rank of a matrix, enhancing row/column correlation [28]. This function will be used in the numerical examples in Section 4.5.
- **Group lasso norm:** It is defined as $R(\theta) = k \|\theta\|_J$, where $J = \{J_1; \dots; J_l\}$ is a partition of $[p]$.

- **Total variation:** It is defined as $R(x) = k \|Dx\|_{1,2}$; where the $k \in \mathbb{R}_{1,2}$ and D denote the $\|\cdot\|_{1,2}$ -norm and D is the discrete gradient operator, respectively (see Chapter 2). This regularizer is usually used in image reconstruction (deblurring + denoising) since it maintains sharp edges between different areas with constant color (see for example [23, 37, 39, 87, 103, 105, 116, 117, 140]).

This function will be used in the numerical examples in Section 4.5.

3.2 Tikhonov regularization, flows, and algorithms

A classical way to solve (P) is to relax the constraint, and use **Tikhonov regularization** [19, 52, 128]:

$$\min_{x \in \mathbb{R}^p} \frac{1}{2} L(x) + R(x); \quad (P')$$

Note that the loss is explicitly added to the objective function and is multiplied by a regularization parameter α . The minimizers of the above problem define a sequence $f_\alpha, g_\alpha > 0$, of possible solutions, called regularization path [51, 63]. Among these solutions, the best regularized solution is selected according to some criteria, such as the Morozov discrepancy principle in inverse problems [52] or cross-validation on left-out data in machine learning [57, 126]. Typically, this involves fixing a grid α within the interval $[\alpha_{\min}, \alpha_{\max}]$ and then solving (P') for each $\alpha \in \mathcal{G}$, where the best regularized solution is selected according to an appropriate data-driven criterion.

Example 3.2.1 (Ridge regression). The most classic example is when L is the least square, and R is the square norm. Then, problem (P') becomes

$$\min_{x \in \mathbb{R}^p} \frac{1}{2} \|Xx - y\|^2 + k \|x\|^2;$$

The above problem is also called ridge regression [57, 65, 128]. It is a classic fact that the sequence of solutions $f_\alpha := (Id + \alpha X^T X)^{-1} X^T y, g_\alpha > 0$, converges to a minimum norm solution [52, Theorem 5.2], when $\alpha \rightarrow 0$.

In the following section we present flows and algorithms that can be applied to solve (P'). From the above discussion, it is clear that optimization to solve (P') (and also (P)) plays a crucial role in the solution of machine learning and inverse problems.

3.2.1 $R = 0$: Continuous case

Note that when $R = 0$, the problems described in equations (P') and (P) simplify to solve:

$$\min_{x \in \mathbb{R}^p} L(x);$$

Additionally, if we assume that $L: \mathbb{R}^p \rightarrow \mathbb{R}$ is a convex and differentiable function such that ∇L is Lipschitz continuous on \mathbb{R}^p , a good approach to solving the previous problem is to consider the steepest descent trajectory. The gradient flow defines a trajectory for each point x_0 in the space, which can be described by the following ordinary differential equation:

$$\begin{aligned} \dot{x}(t) &= -\nabla L(x(t)); \quad t > 0 \\ x(0) &= x_0; \end{aligned} \quad (\text{GF})$$

It follows from Theorem 2.4.18 that the previous flow, for every $x_0 \in \mathbb{R}^p$, has unique solution.

An extension of gradient flow to non-Euclidean geometries is mirror flow [3], which can be modeled as the following ordinary differential equation:

$$r^2 \dot{F}(z(t)) - \dot{z}(t) = -r \nabla L(z(t)); \quad (\text{MF})$$

where $F: \text{dom}(F) \rightarrow \mathbb{R}$ is the mirror map and is strictly convex, Legendre and twice-differentiable in $U = \text{int}(\text{dom}(F))$, and the mapping $z \in U \mapsto [r^2 F(z)]$ is invertible with locally Lipschitz inverse. In the case when the mirror map is defined as the Euclidean norm, i.e., $F(z) = \frac{k \|z\|^2}{2}$, the equation describing the mirror flow becomes equivalent to the gradient flow.

Observe that, for every $z \in S$;

$$\frac{d(D_F(z; z(t)))}{dt} = h r \nabla L(z(t)) \cdot (z - z(t)) - L(z(t)) \leq 0; \quad (3.2.1)$$

where $L = \min_{z \in S} L(z)$. Since the Bregman distance is decreasing and bounded below we get that, for every $z \in S$; $D_F(z; z(t))$ converges. Moreover, computing the derivative of $L(z(t))$, we get that

$$\frac{d(L(z(t)) - L)}{dt} = \nabla L(z(t)) \cdot \dot{z}(t) = -r \nabla L(z(t)) \cdot \nabla L(z(t)) \leq 0;$$

Which implies that $(L(z(t)) - L)$ is decreasing with respect to t . Since $D_F(\cdot; \cdot)$ is non-negative, we can integrate (3.2.1) and obtain that

$$t(L(z(t)) - L) \leq \int_0^t (L(z(s)) - L) ds \leq D_F(z; z(0));$$

which implies convergence in value.

Note that, if $F(z) = \frac{k \|z\|^2}{2}$, we get that $k(z(t) - zk)$ converges for every $z \in S$; which is the first condition of Opial's Lemma. While the second condition, each limit point belongs to the solution set, is obtained by convergence in value. Then, by using Lemma 2.1.5, we conclude that $z(t)$ converges to some $z^* \in S$:

3.2.2 $R = 0$: Discrete case

If we discretize $z(t)$ with $t \in [k; k+1]$, using finite differences with a positive step-size $\Delta t = \frac{1}{L}$; where L is the Lipschitz constant of $r \nabla L$, then

$$z(k+1) - z(k) = -\frac{1}{L} r \nabla L(z(k));$$

Gradient descent [35]: If we take $r \nabla L(z(t)) = r \nabla L(z(k))$; then the algorithm has the following explicit update rule

$$z^{k+1} = z^k - \frac{1}{L} r \nabla L(z^k); \quad (3.2.2)$$

In the following theorem, we establish sufficient conditions for the convergence of gradient descent in \mathbb{R}^p .

Theorem 3.2.2. Let $L: \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex Gâteaux differentiable function, such that $S \in \arg \min L$. Suppose that rL is L -Lipschitz continuous on \mathbb{R}^p . Let $\{x^k\}_{k \in \mathbb{N}}$ be the sequence generated by (3.2.2) starting from some arbitrary $x^0 \in \mathbb{R}^p$ and with step-size $\alpha \in]0, \frac{2}{L}[$. Then,

$$\lim_{k \rightarrow +\infty} x^k = \arg \min L.$$

Example 3.2.3 (Gradient descent on the least squares). Let X be the matrix with rows corresponding the input data and let y be the vector of measurements. In this case the iterates of gradient descent are given by

$$x^{k+1} = x^k - \alpha (X^T X)^{-1} X^T y.$$

The previous method is also known as the Landweber method [52, 72]. In Chapter 4 we use a single iteration of this algorithm with different step-sizes to improve the feasibility of iterations.

If we assume that X is full rank and we continue assuming that L is the least square, we can use a well-known stochastic algorithm.

Randomized Kaczmarz: Let $x^0 \in \mathbb{R}^p$ and consider the following algorithm:

$$x^{k+1} = x^k + \alpha \frac{X_{j_t}^T (y_t - X_{j_t} x^k)}{\|X_{j_t}\|^2},$$

where $\{g_{k \in \mathbb{N}}\}$ is a sequence of independent $[d]$ -valued random variables. It is well-known and has been demonstrated that the sequence $\{x^k\}_{k \in \mathbb{N}}$ converges to a feasible point (see [41, 127]). For further information, refer to [36, 41, 62, 70, 73, 95–98, 127, 150]. In Chapter 4 we use a single iteration of this algorithm to improve the feasibility of iterations.

In the case when L is not necessary L -smooth and prox_L is easy to compute, we can use **proximal point algorithm** [81], which has the following implicit update rule

$$x^{k+1} = \text{prox}_L(x^k). \quad (3.2.3)$$

Note that this method can be applied, unlike gradient descent, to non-differentiable functions. However, since it is an implicit rule, in most cases the prox cannot be easily calculated. In the following theorem, we provide the conditions for the convergence of proximal point algorithm in \mathbb{R}^p .

Theorem 3.2.4. Let $L: \mathbb{R}^p \rightarrow \mathbb{R} \cup \{+\infty\}$ be a proper convex lower semi-continuous function such that $S \neq \emptyset$. Let $\{x^k\}_{k \in \mathbb{N}}$ be the sequence generated by (3.2.3) starting from some arbitrary $x^0 \in \mathbb{R}^p$ and with step-size $\alpha > 0$. Then,

$$\lim_{k \rightarrow +\infty} x^k = \arg \min L.$$

3.2.3 Case $R \neq 0$

If R and L are both L -smooth, then we can apply gradient descent, which has the following update rule

$$x^{k+1} = x^k - \alpha (L(x^k) - R(x^k)); \quad (3.2.4)$$

where $\alpha > 0$ is the regularization parameter. In Theorem 3.2.2 are provided conditions to ensure the convergence of (3.2.4). In the case when R is non-smooth and prox_R is easy to compute, we can use forward-backward algorithm, which is

$$x^{k+1} = \text{prox}_R(x^k - \alpha L(x^k)).$$

In the following theorem, we provide the conditions for the convergence of forward-backward in \mathbb{R}^p .

Theorem 3.2.5. Let $R: \mathbb{R}^p \rightarrow \mathbb{R} \cup \{+\infty\}$ and $L: \mathbb{R}^p \rightarrow \mathbb{R}$ be two functions in $\mathcal{C}_0(\mathbb{R}^p)$ such that $\operatorname{argmin}_{\mathbb{R}^p} (L(\cdot) + R(\cdot)) \neq \emptyset$. Suppose that L is Gâteaux differentiable such that rL is L -Lipschitz continuous on \mathbb{R}^p . Let $f^k g_{k \geq 2N}$ be the sequence generated by (3.2.2) starting from some arbitrary $x^0 \in \mathbb{R}^p$ and with step-size $\alpha \in]0, \frac{2}{L}[$. Then,

$$\lim_{k \rightarrow +\infty} f^k = \alpha \operatorname{argmin}_{\mathbb{R}^p} (L(\cdot) + R(\cdot));$$

For all of the methods presented above, selecting an appropriate α can be challenging and typically requires solving the optimization problem multiple times with different parameter values.

3.3 Iterative regularization

In this section, we explore a more efficient alternative known as iterative regularization. Unlike Tikhonov regularization, this technique involves running a single optimization procedure that stops before convergence. Typically, since a single problem is solved and the algorithm is not run even until convergence, iterative regularization is numerically more efficient compared to Tikhonov.

In practice, the exact data y is unknown, and only a noisy version is accessible. Given a noise level $\delta > 0$, we consider a worst-case scenario where the error is deterministic and the accessible data y^δ is such that

$$\|y^\delta - y\| \leq \delta;$$

which limits us to solve only a noisy version of (P).

Iterative regularization consists of finding an approximation of x^* by running an iterative algorithm that generates a regularizing sequence $f^k g_{k \geq 2N}$ and stopping when it is closed to the solution according to certain criteria. The algorithm must consider the properties of R ; L ; and X , and for any data $y^\delta \in \mathbb{R}^d$, the algorithm generates a sequence. To analyze the behavior of the sequence $f^k g_{k \geq 2N}$, it is necessary to define the auxiliary sequence $f^k g_{k \geq 2N}^\delta$, which is obtained by applying the same algorithm but using y^δ instead of y . We additionally assume that the sequence converges to x^* , the solution of (P). Then, we can decompose the error as

$$\|f^k g_{k \geq 2N}^\delta - x^*\| \leq \|f^k g_{k \geq 2N} - x^*\| + \|f^k g_{k \geq 2N}^\delta - f^k g_{k \geq 2N}\| \quad (3.3.1)$$

The first term corresponds to an optimization error, which decreases with respect to k . While the second error measures the stability of noise, which increases with respect to k , since the iterations converge to a noisy solution or diverges. An early stopping strategy takes in account the different behaviours of the two terms in equation (3.3.1) and stops the iterations when these terms are approximately equal. In this way, the number of iterations plays the role of the regularization parameter in the same way as α in Tikhonov regularization.

For the next two algorithms, we consider a simpler noisy version of (P),

$$\min_{x=y^\delta} R(x); \quad (3.3.2)$$

which is studied in the Chapter 4. The above problem is obtained by choosing L equal to least squares and taking y instead of y : Note that (3.3.2) may not be feasible. We also define the noise free problem as

$$\min_{X=y} R(\cdot) \tag{3.3.3}$$

We now present two algorithms: the first for strongly convex regularizers and the second for regularizers that are only convex. Additionally, both algorithms incorporate early stopping strategies.

Dual Gradient descent: Let $R: \mathbb{R}^p \rightarrow \mathbb{R} \cup \{+\infty\}$ be proper, lower semicontinuous, and μ -strongly convex function. Let $w^0 = 0$ and $\lambda = \frac{1}{kXk^2}$. Then, the iterations of dual gradient descent are:

$$\begin{aligned} w^{k+1} &= \text{prox}_R \left(\frac{X^T w^k}{\lambda} \right); \\ w^{k+1} &= w^k + \lambda X^{k+1} y; \\ p^{k+1} &= \frac{1}{k+1} \sum_{i=1}^{k+1} p_i. \end{aligned}$$

Now we present an early stopping result for the previous algorithm.

Theorem 3.3.1. [82, Theorem 4.1] Let $\lambda \in]0; 1]$. Assume that there exists $w \in \mathbb{R}^p$ such that $X^T w \in \text{ran}(R)$. Set $c_1 = 2kXk^{-1}$ and $c_2 = \frac{kXkw^k}{k^2}$, where w^k is a solution of the dual problem of (3.3.3). Then, for every $k \in \mathbb{N}$,

$$\|k p^k - w^k\| \leq c_1 \frac{1}{k} + \frac{c_2}{k} \tag{3.3.4}$$

In particular, choosing $\lambda = \frac{c}{c+1}$ for some $c > 0$, we derive

$$\|k p^k - w^k\| \leq c_1 (c^{\frac{1}{2}} + 1) + \frac{c_2}{c^{\frac{1}{2}}} \frac{1}{2}.$$

The previous result gives us a decomposition of the error into optimization and stability errors as in (3.3.1) in (3.3.4). In some cases the regularizer is neither smooth nor strongly convex. In this situation it is suitable to use the dual primal-dual algorithm.

Primal-dual: Let $R: \mathbb{R}^p \rightarrow \mathbb{R} \cup \{+\infty\}$ be a function in $\text{dom}(R^p)$. Let $w^0 = w^{-1} \in \mathbb{R}^d$, $w^0 \in \mathbb{R}^p$ and $\lambda > 0$ and $\mu > 0$ such that $kXk^2 < 1$. Then, the iterations of primal dual are:

$$\begin{aligned} w^{k+1} &= \text{prox}_R \left(\frac{X^T (2w^k - w^{k-1})}{\lambda} \right); \\ w^{k+1} &= w^k + \lambda X^{k+1} y; \\ p^{k+1} &= \frac{1}{k+1} \sum_{i=1}^{k+1} p_i. \end{aligned}$$

Note that the above algorithm is a specific case of the primal-dual algorithm proposed in [40], where the function in the primal is R , the function in the dual is the indicator $\text{f}_y g$, and the linear operator is X :

Theorem 3.3.2. [85, Proposition 7] Let $R: \mathbb{R}^p \rightarrow \mathbb{R} \cup \{+\infty\}$ be a proper, convex, and lower semicontinuous function. Assume that the problem (3.3.3) has at least one solution. Let $(x^*; w^*) \in \mathbb{R}^p \times \mathbb{R}^d$, where x^* is a solution of (3.3.3) and w^* is a solution of the dual problem of (3.3.3). Then, there exist strictly positive constants c_1 and c_2 such that, for every $k \in \mathbb{N}$,

$$\|x^k - x^*\| \leq c_1 \left(\frac{1}{k} + \frac{1}{k^2} \right) \quad (3.3.5)$$

and

$$R(x^k) + \sum_{j=1}^D w_j^* (x^k - y^j)^E - R(x^*) \leq \frac{c_2}{k} (1 + \frac{1}{k})^2 \quad (3.3.6)$$

Moreover, if we choose $k \geq \frac{1}{c}$ for some $c > 0$, then there exist constants c_3 and c_4 such that

$$\|x^k - x^*\| \leq c_3 \frac{1}{k^2}$$

and

$$R(x^k) + \sum_{j=1}^D w_j^* (x^k - y^j)^E - R(x^*) \leq c_4 \frac{1}{k}$$

Unlike (3.3.1), the previous result gives us a decomposition of the error into optimization and stability errors, but in this case applied to the feasibility gap in (3.3.5) and the duality gap in (3.3.6).

In Chapter 4, we introduced two algorithms inspired by the primal-dual approach with activation discussed in [24]. This method has shown significant numerical speed up in noise-free scenarios by improving the feasibility of the iterations. We combine this algorithm with the early stopping strategy presented in [84].

CHAPTER 4

Fast iterative regularization by reusing data

Abstract

Discrete inverse problems correspond to solving a system of equations in a stable way with respect to noise in the data. A typical approach to selecting a meaningful solution is to introduce a regularizer. While for most applications the regularizer is convex, in many cases it is neither smooth nor strongly convex. In this chapter, we propose and study two new iterative regularization methods, based on a primal-dual algorithm, to regularize inverse problems efficiently. Our analysis, in the noise-free case, provides convergence rates for the Lagrangian and the feasibility gap. In the noisy case, it provides stability bounds and early stopping rules with theoretical guarantees. The main novelty of our work is the exploitation of some a priori knowledge about the solution set: we show that the linear equations determined by the data can be used more than once along the iterations. We discuss various approaches to reusing linear equations that are at the same time consistent with our assumptions and flexible in their implementation. Finally, we illustrate our theoretical findings with numerical simulations for robust sparse recovery and image reconstruction. We confirm the efficiency of the proposed regularization approaches by comparing the results with state-of-the-art methods.

Keywords. Primal-dual splitting algorithms, Iterative regularization, Early stopping, Landweber method, Stability and convergence analysis.

AMS Mathematics Subject Classification (2020): 90C25, 65K10, 49M29.

4.1 Introduction

Many applied problems require the estimation of a quantity of interest from noisy linear measurements, for instance compressed sensing [32, 34, 48, 115, 129], image processing [37, 39, 103, 105, 116, 117, 140], matrix completion [28, 31, 33, 84], and various problems in machine learning [13, 50, 90, 114, 121, 140, 141]. In all these problems, we are interested in finding stable solutions to a system of equations where the accessible data is corrupted by noise. This is classically achieved by regularization. The most popular procedure in the literature is Tikhonov (or variational) regularization [52], which consists in minimizing the sum of a data fidelity term plus a regularizer, which is explicitly added to the objective function and entails some a priori knowledge or some desired property on the solutions that we want to select. A trade-off parameter is then introduced to balance the fidelity term and the regularizer. In practice, this implies that the optimization problem has to be solved many times for different values of the parameter. Finally, a parameter - and the correspondent solution - are chosen accordingly to the performance with respect to some criterion, such as the Morozov discrepancy principle in inverse problems [52] or cross-validation on left-out data in machine learning [57, 126].

A computationally efficient alternative to explicit regularization is iterative regularization, also known as implicit regularization [8, 21, 27, 52]. The minimization of the regularizer under noisy data constraints is considered, and a numerical algorithm to solve the optimization problem is chosen and early stopped, to avoid convergence to the noisy solution. In this setting, it is known that the number of iterations plays the role of the regularization parameter [52]. As for Tikhonov regularization, the best-performing iterate is chosen according to some a priori criterion and then considered as the regularized solution. Compared to Tikhonov regularization, this procedure is very efficient since only one optimization problem is solved, and not even until convergence.

The main novelty of this work is the design and analysis of two new iterative regularization methods for convex regularizers, which are neither necessarily smooth nor strongly convex. The new iterative regularization methods are based on primal-dual algorithms [40, 46, 133] combined with the idea of reusing the linear equations determined by the data at every iteration [24]. Primal-dual algorithms perform one minimization step on the primal variable followed by one on the dual and are well-suited for the large-scale setting, as only matrix-vector multiplications and the calculation of a proximity operator are required. The idea of exploiting redundant information was presented in [24] and turned out to be very effective in practice. The first method that we propose is a primal-dual algorithm (**PDA**) with additional activations of the linear equations: We propose different variants, depending on the extra activation steps. For instance, we are able to exploit the data constraints more than once at every iteration via gradient descent, with a fixed or adaptive step size. The second method is a dual-primal algorithm (**DPA**), where a subset containing the dual solutions is activated at each step. This subset is not affected by the noise in the data and is usually determined by a finite number of constraints.

These additional steps may seem artificial or inefficient. However, while maintaining easy implementation, our methods achieve better numerical performances and considerable speed-ups with respect to the vanilla primal-dual algorithm. We extend to the noisy case the techniques studied in [24, 25] for the exact case. The assumptions on the noise are the classical ones in inverse problems, see e.g. [27, 30, 82, 84], and the proposed results generalize the ones in [84], by including in the primal-dual procedure a diagonal preconditioning and an extra activation step. For the noisy case, we provide an early stopping

criterion to recover a stable approximation of an ideal solution, in the same spirit of [12, 20, 27, 30, 82, 112, 141, 146]. The early stopping rule is derived from theoretical stability bounds and feasibility gap rates for both algorithms, obtaining implicit regularization properties similar to those stated in [84] and [82]. Theoretical results are complemented by numerical experiments for robust sparse recovery and total variation, showing that state-of-the-art performances can be achieved with considerable computational speed-ups.

Related works. In this section, we briefly discuss the literature about variational and iterative regularization techniques. Tikhonov regularization has been introduced in [128]; see also [19, 52] and the references therein for an extensive treatment of the topic. The most famous iterative regularization method is the Landweber algorithm [52, 72], namely gradient descent on the least squares problem. Duality theory in optimization gives another interpretation that sheds light on the regularizing properties of this procedure. Indeed, consider the problem of minimizing the squared norm under linear constraints. Running gradient descent on its dual problem and mapping back to the primal variable, we obtain exactly the Landweber method. This provides another explanation of why the iterations of the Landweber algorithm converge to the minimal norm solution of the linear equation [82]. Stochastic gradient descent on the previous problem is the generalization of the Kaczmarz method [70, 78, 120, 127], which consists in applying cyclic or random projections onto single equations of the linear system. Accelerated and diagonal versions are also discussed in [52, 101] and [10, 71, 119], respectively. The regularization properties of other optimization algorithms for more general regularizers have also been studied. If strong convexity is assumed, mirror descent [16, 100] can also be interpreted as gradient descent on the dual problem, and its regularization properties (and those of its accelerated variant) have been studied in [82]. Diagonal approaches [9] with a regularization parameter that vanishes along the iterations have been studied in [54]; see [30] for an accelerated version. Another common approach relies on the linearized Bregman iteration [103, 140, 142, 143], which has found applications in compressed sensing [29, 104, 143] and image deblurring [29]. However, this method requires to solve non-trivial minimization problems at each iteration. For convex, but not strongly convex regularizers, the regularization properties of primal-dual algorithms have been investigated in [84]. Other optimization techniques are available to solve this kind of minimization problem (for instance, [86, 87] and [22, 83, 110]; see also [66, 123, 124]), but no iterative regularization properties have been studied so far for these algorithms.

4.2 Main problem and algorithm

Many applied problems require to estimate a quantity of interest $x \in \mathbb{R}^p$ based on linear measurements $y = Xx$, for some matrix $X \in \mathbb{R}^{d \times p}$. A standard approach to recover the desired solution is to assume that it is a minimizer of the following linearly constrained optimization problem:

$$\min_{x \in \mathbb{R}^p} f_R(x) : Xx = yg; \quad (P)$$

where $R \in \mathbb{R}^p$ encodes a priori information on the solution and is usually hand-crafted. Typical choices are: the squared norm [52]; the elastic net regularization [47, 69, 74, 82, 148, 149]; the ℓ^1 -norm [32, 34, 48, 129]; and the total variation [37, 105, 116, 117]. Note that, in the previous examples, the first two regularizers are strongly convex, while the second two are just convex and non-smooth.

If we use the indicator function of f_{yg} , the problem (P) can be written equivalently as

$$\min_{2R^p} R(\cdot) + f_{yg}(X(\cdot)):$$

We denote by \min the optimal value of (P) and by S the set of its minimizers. We assume that $S \neq \emptyset$. In order to build our regularization procedure, we consider the Lagrangian functional for problem (P) :

$$L: R^p \times R^d \rightarrow R \text{ [} f + 1 g \\ (\cdot; w) \mapsto R(\cdot) + hwj X(\cdot) - yi:$$

This approach allows us to split the contribution of the non-smooth term R and the one of the linear operator R , without requiring to compute the projection on the set

$$S := \{ \cdot \in R^p \mid X(\cdot) = y \}:$$

We define the set of saddle points of L by

$$Z = \{ (\cdot; w) \in R^p \times R^d \mid L(\cdot; v) \leq L(\cdot; w) \leq L(y; w) \leq \delta(y; v) \in R^p \times R^d \}:$$

The set Z is characterized by the first-order optimality condition:

$$Z = \{ (\cdot; w) \in R^p \times R^d \mid 0 \in \partial R(\cdot) + X^>w \text{ and } X(\cdot) = y \}:$$

In the following, we always assume that $Z \neq \emptyset$:

Remark 4.2.1 (Saddle points and primal-dual solutions). Observe that the objective function of (P) is the sum of two functions in $\text{co}(R^p)$ where one of the two is composed with a linear operator. This formulation is suitable to apply Fenchel-Rockafellar duality. Recalling that $f_{yg}(w) = hwj yi$, the dual problem of (P) is given by

$$\min_{w \in R^d} R(\cdot > w) + hwj yi: \tag{D}$$

We denote its optimal value by \max and its set of minimizers by S^* . Then, $\min = \max$, and equality holds if the qualification condition (see [15, Proposition 6.19] for special cases when it holds)

$$y \in \text{ri}(X(\text{dom } R)) \tag{4.2.1}$$

is satisfied [15, Proposition 19.21 (v)]. In addition, condition (4.2.1) implies that problem (D) has a solution. Then, under (4.2.1), since we assumed that $S \neq \emptyset$, we derive also that $Z \neq \emptyset$.

In practical situations, the exact data y is unknown and only a noisy version is accessible. Given a noise level $\epsilon > 0$, we consider a worst case scenario, where the error is deterministic and the accessible data y^ϵ is such that

$$\|y^\epsilon - y\| \leq \epsilon:$$

This is the classical model in inverse problems [52, 71]. The solution set of the inexact linear system $X(\cdot) = y^\epsilon$ is denoted by S^ϵ . Analogously, we denote by S^ϵ and $S^{\epsilon*}$ the sets of primal and dual solutions with noisy data, respectively. It is worth pointing out that, if $y^\epsilon \notin \text{ran}(X)$, then $S^\epsilon = \emptyset$; but our analysis and bounds still hold.

4.2.1 Primal-Dual Splittings with a priori Information

In this section, we propose an iterative regularization procedure to solve problem (P) , based on a primal-dual algorithm with preconditioning and arbitrary activations of a pre-defined set of operators. While the use of primal-dual algorithms [40] as iterative regularization methods is somewhat established [84], in this chapter we focus on the possibility of reusing the data constraints along the iterations. This idea was originally introduced in [24], where the authors studied the case in which the exact data is available, and consists in the activation of extra operators, which encode information about the solution set, to improve the feasibility of the updates. In our setting, we reuse data constraints, and we project, in series or in parallel, onto some equations given by the (noisy) linear constraints. But we will show that other interesting choices are possible, as projections onto the set of dual constraints.

More formally, for $i \in [m]$, we consider a finite number of operators $T_i: \mathbb{R}^p \rightarrow \mathbb{R}^p$ or $T_i: \mathbb{R}^d \rightarrow \mathbb{R}^d$, such that the set of noisy primal (or dual) solutions is contained in $\text{Fix } T_i$ for every $i \in [m]$. We refer to this as redundant a priori information. A list of operators suitable to our setting (and with a cheap practical implementation) can be found in Section 4.4.

The primal-dual algorithms with reuse of data which are given in Table 4.1 are a preconditioned and deterministic version of the one proposed in [24] applied to the case of linearly constrained minimization.

Primal-Dual splitting with activations	Dual-Primal splitting with activations
<p>Input: $(p^0; 0; w^0) \in \mathbb{R}^{2p} \times \mathbb{R}^d$.</p> <p>For $k = 1; \dots; N$:</p> $\begin{aligned} w^{k+1} &= w^k + (X^T p^k - y) \\ p^{k+1} &= \text{prox}_R(p^k, X^T w^{k+1}) \end{aligned}$ <p>Choose $i_{k+1} \in [m]$ and set (PDA)</p> $\begin{aligned} p^{k+1} &= T_{i_{k+1}} p^{k+1} \\ p^{k+1} &= p^{k+1} + \alpha_{k+1} (p^k - p^{k+1}); \end{aligned}$ <p>End</p>	<p>Input: $(v^0; w^0; 0) \in \mathbb{R}^{2d} \times \mathbb{R}^p$.</p> <p>For $k = 1; \dots; N$:</p> $\begin{aligned} v^{k+1} &= \text{prox}_R(v^k, X^T w^k) \\ w^{k+1} &= v^k + (X v^{k+1} - y) \end{aligned}$ <p>Choose $i_{k+1} \in [m]$ and set (DPA)</p> $\begin{aligned} v^{k+1} &= T_{i_{k+1}} w^{k+1} \\ v^{k+1} &= v^{k+1} + w^{k+1} - v^k; \end{aligned}$ <p>End</p>

Table 4.1: Proposed algorithms for iterative regularization.

We first focus on the Primal-Dual splitting. It is composed by four different steps, to be performed in series. The first step is the update of the dual variable, in which the residuals to the linear equation $X^T w = y$ are accumulated after preconditioning by the operator R . The second step is an implicit prox-step, with function R and norm $\|\cdot\|_R$, on the primal variable. The third one is the activation of the operator related to reusing data constraints, on the primal variable. Finally, the last step is an extrapolation again on the primal variable. Notice that, if no operator is activated, it corresponds simply to $p^{k+1} = 2 p^{k+1} - p^k$, which is the classical update in the primal-dual algorithm. Observe that, the Dual-Primal Splitting algorithm, except for permutation in the order of the steps, differs from the previous one because the activation of the operator is done not on the primal variable but on the dual one.

Remark 4.2.2. Observe that in the proof convergence and stability (Theorem 4.3.1 and Theorem 4.3.2) we will never use that \mathbb{R}^p belongs to a finite dimensional space. This is in line with previous research on the convergence guarantees of the plain methods in Hilbert and Banach spaces, as outlined in [46, 122, 133]. It follows that the primal-dual

algorithms above can be formulated exactly in the same way when the unknown vector x belongs to an infinite dimensional Hilbert space, and our analysis can be extended to that setting. Another possible extension of the algorithm, which we do not analyze explicitly in this work, is related with the stochastic version of primal-dual algorithm; see [1, 38, 61].

4.2.2 Equivalence between Primal-dual and Dual-primal algorithms.

The next lemma establishes that, if $T = \text{Id}$ and the initialization is the same, then there is an equivalence between the k -th primal variable generated by PDA and the ones generated by DPA.

Lemma 4.2.3. Let

$$(p_{PD}^0; p_{PD}^0; w_{PD}^0) \in \mathbb{R}^{2p} \times \mathbb{R}^d \text{ and } (v_{DP}^0; w_{DP}^0; p_{DP}^0) \in \mathbb{R}^{2d} \times \mathbb{R}^p$$

be the initialization PDA and DPA, respectively, in the case when $m = 1$ and $T = \text{Id}$. Suppose that

$$p_{PD}^0 = p_{PD}^0; \quad w_{DP}^0 = v_{DP}^0; \quad w_{PD}^0 = v_{DP}^0; \quad \text{and} \quad p_{PD}^1 = p_{DP}^1.$$

Then for every $k \in \mathbb{N}$, $p_{PD}^k = p_{DP}^k$.

Proof. Since $m = 1$ and $T = \text{Id}$ in both algorithms, for every $k \in \mathbb{N}$, we have that $p_{PD}^k = p_{PD}^k$ and $w_{DP}^k = v_{DP}^k$. On one hand, by the definition of PDA, we have that

$$\begin{aligned} w_{PD}^{k+1} &= w_{PD}^1 + \sum_{i=1}^k X p_{PD}^i - y \\ &= w_{PD}^1 + \sum_{i=1}^k X (p_{PD}^i - p_{PD}^{i-1}) + \sum_{i=1}^k X p_{PD}^i - y \\ &= w_{PD}^1 + X (p_{PD}^k - p_{PD}^0) + \sum_{i=1}^k X p_{PD}^i - y \\ &= w_{PD}^0 + (X p_{PD}^k - y) + \sum_{i=1}^k X p_{PD}^i - y; \end{aligned} \quad (4.2.2)$$

where the last equality is obtained since $p_{PD}^0 = p_{PD}^0$. Replacing (4.2.2) in the definition of p_{PD}^{k+1} , we obtain

$$p_{PD}^{k+1} = \text{prox}_R \left(p_{PD}^k - X \left(w_{PD}^0 + (X p_{PD}^k - y) + \sum_{i=1}^k X p_{PD}^i - y \right) \right);$$

On the other hand, by DPA we have that

$$w_{DP}^{k+1} = v_{DP}^{k+1} = v_{DP}^0 + \sum_{i=1}^k X p_{DP}^i - y;$$

and

$$v_{DP}^k = v_{DP}^k + w_{DP}^k \quad v_{DP}^{k+1} = v_{DP}^0 + (X p_{DP}^k - y) + \sum_{i=1}^k X p_{DP}^i - y; \quad (4.2.3)$$

Replacing (4.2.3) in DPA, for every $k > 1$, we can deduce that

$$w_{DP}^{k+1} = \text{prox}_R \left(w_{DP}^k - X^> v_{DP}^0 + (X T_k w_{DP}^k - y) + \sum_{i=1}^k X^i w_{DP}^i - y \right);$$

Since $w_{PD}^0 = v_{DP}^0$ and $w_{PD}^1 = w_{DP}^1$ the result follows by induction. □

Remark 4.2.4. An analysis similar to that in the proof of Lemma 4.2.3 shows that

$$w_{PD}^{k+1} = \text{prox}_R \left(w_{PD}^k - X T_k w_{PD}^k + y + \sum_{i=1}^k X^i w_{PD}^i - y \right);$$

which implies that the algorithm can be written in one step if we only care about the primal variable.

4.2.3 Assumptions

In the following, we list the assumptions that we require on the parameters and the operators involved in the algorithm.

Assumption 4.2.5. Consider the setting of PDA or DPA:

A1 The preconditioners $T \in \mathbb{R}^{p \times p}$ and $X \in \mathbb{R}^{d \times d}$ are two diagonal positive definite matrices such that

$$0 < \lambda := \min \left\{ \lambda_{\min}(T), \lambda_{\min}(X) \right\} > 0; \tag{4.2.4}$$

A2 For every $k \in \mathbb{N}$, $k \leq m$.

Consider the setting of PDA:

A3 The family of operators $f_{T_i g_{i \in [m]}}$ is from \mathbb{R}^p to \mathbb{R}^p and for every $i \in [m]$:

(a) Fix $T_i \in \mathbb{S}^p$;

(b) there exists $e_i \geq 0$ such that, for every $w \in \mathbb{R}^p$ and $s \in \mathbb{S}^p$,

$$\|T_i w - s\|_2 \leq \|w - s\|_2 + e_i \|w - s\|_2; \tag{4.2.5}$$

We set $e = \max_{i \in [m]} e_i$.

Now consider the setting of DPA:

A4 $f_{T_i g_{i \in [m]}}$ is a family of operators from \mathbb{R}^d to \mathbb{R}^d and for every $i \in [m]$:

(a) Fix $T_i \in \mathbb{S}^d$;

(b) for every $u \in \mathbb{R}^d$ and $s \in \mathbb{S}^d$,

$$\|T_i u - s\|_2 \leq \|u - s\|_2 + k \|u - s\|_2;$$

Remark 4.2.6 (Hypothesis about the operators). If Assumption A3a holds and $\beta = 0$, then Assumption A3b is implied by the quasi-nonexpansivity of T_i on S . This is a weaker condition than the one proposed in [24], where, due to the generality of the setting, β -averaged non-expansive operators were needed. A similar reasoning applies to Assumption A4.

4.3 Main results

In this section, we present and discuss the main results of the chapter. We derive convergence and stability properties of primal-dual and dual-primal splitting algorithms for linearly constrained optimization with a priori information.

First, we define the averaged iterates and the square weighted norm induced by \mathbb{P} and \mathbb{W} on $\mathbb{R}^p \times \mathbb{R}^d$, namely

$$\hat{z}^n := \frac{\sum_{k=1}^n z^k}{n} \quad \text{and} \quad V(z) := \frac{k}{2} \|z\|_{\mathbb{P}}^2 + \frac{k\mathbb{W}k}{2} \|z\|_{\mathbb{W}}^2;$$

where $z^k := (\rho^k; w^k)$ is the k -th iterate and $z := (\rho; w)$ is a primal-dual variable. We also recall the definition of the Lagrangian as $L(\rho; w) = R(\rho) + \langle w, X\rho - y \rangle$.

The first result establishes the stability properties of the algorithm PDA, both in terms of the Lagrangian and feasibility gap. We recall that here we use activation operators based on the noisy data and corresponding constraints in the primal space, namely the set C .

Theorem 4.3.1. Consider the setting of PDA under Assumptions A1, A2, and A3. Let $(\rho^0; w^0) \in \mathbb{R}^p \times \mathbb{R}^d$ be such that $\rho^0 = \rho^*$. Then, for every $z = (\rho; w) \in Z$ and for every $N \in \mathbb{N}$, we have

$$\begin{aligned} L(\hat{\rho}^N; w) - L(\rho^*; w^N) &\leq \frac{V(z^0 - z)}{N} + \frac{2Nk}{2} \|z^0 - z\|_{\mathbb{P}}^2 + k \|z^0 - z\|_{\mathbb{W}}^2 \\ &\quad + k \|z^0 - z\|_{\mathbb{W}} \frac{Ne}{2} + \frac{e}{2} \end{aligned} \quad (4.3.1)$$

and

$$\begin{aligned} k \|X\hat{\rho}^N - y\|^2 &\leq \frac{16Nk}{2} \|z^0 - z\|_{\mathbb{P}}^2 + 8k \|z^0 - z\|_{\mathbb{W}} \frac{2k}{3} \|V(z^0 - z)\|_{\mathbb{P}}^2 + 8k \|z^0 - z\|_{\mathbb{W}} \frac{k}{3} \|keN\|_{\mathbb{W}}^2 \\ &\quad + \frac{8k}{N} \|V(z^0 - z)\|_{\mathbb{P}}^2 + 2\|z^0 - z\|_{\mathbb{W}}^2 + \frac{4k}{N} \|ke\|_{\mathbb{W}}^2; \end{aligned} \quad (4.3.2)$$

where we recall that the constant β and e are defined in Assumptions A1 and A3, respectively.

Proof. From PDA, we deduce that:

$$\begin{aligned} \rho^{k+1} - \rho^* &\leq X^T w^{k+1} \in \mathcal{R}(C^{k+1}) \\ \rho^k - \rho^* + X\rho^k &= y \end{aligned} \quad (4.3.3)$$

Thus

$$\mathcal{R}(C^{k+1}) \cap \mathcal{D}(\rho^k - \rho^* + X\rho^k) \subseteq \mathcal{R}(C^{k+1}) \cap \mathcal{E}(\rho^k) \quad \text{for all } x \in \mathbb{R}^p; \quad (4.3.4)$$

and (4.3.4) yields

$$\begin{aligned}
 0 &= R^{(k+1)} - R^{(k)} + \frac{D}{2} (p^k - p^{k+1}) X^k w^{k+1} j w^{k+1} E \\
 &= R^{(k+1)} - R^{(k)} + \frac{kp^k - p^{k+1}}{2} + \frac{k - k^{k+1}}{2} \frac{k^2 - 1}{2} \\
 &\quad \frac{kp^k - p^{k+1}}{2} + \frac{D}{2} (p^k - p^{k+1}) X^k w^{k+1} j w^{k+1} E
 \end{aligned} \tag{4.3.5}$$

From (4.3.3) we get

$$\begin{aligned}
 0 &= \frac{D}{2} (w^k - w^{k+1}) + X p^k y j w^{k+1} E \\
 0 &= \frac{kw^{k+1} - w^k k^2 - 1}{2} + \frac{kw^{k+1} - w^k k^2 - 1}{2} \frac{kw^k - w^k k^2 - 1}{2} \\
 &\quad + y X p^k j w^{k+1} E
 \end{aligned} \tag{4.3.6}$$

Recall that

$$z := (; w) \in Z \subseteq \mathbb{R}^d; z^k := (; w^k) \text{ and } V(z) := \frac{k - k^2 - 1}{2} + \frac{kw^k k^2 - 1}{2}.$$

Summing (4.3.5) and (4.3.6), and by Assumption A3, we obtain

$$\begin{aligned}
 R^{(k+1)} - R^{(k)} &+ \frac{k - k^{k+1}}{2} \frac{p^k k^2 - 1}{2} + \frac{kw^{k+1} - w^k k^2 - 1}{2} \frac{kw^k - w^k k^2 - 1}{2} \\
 &+ V(z^{k+1} - z) - V(z^k - z) + X^{(k+1)} j w^{k+1} E \\
 &\quad + y X p^k j w^{k+1} E - \frac{e^2}{2} \leq 0.
 \end{aligned} \tag{4.3.7}$$

Now compute

$$\begin{aligned}
 &R^{(k+1)} - R^{(k)} + X^{(k+1)} j w^{k+1} E + y X p^k j w^{k+1} E \\
 &= L^{(k+1)}(; w) - L^{(k)}(; w^{k+1}) + X^{(k+1)} y j w^{k+1} E + X^{(k)} y j w^{k+1} E \\
 &\quad + X^{(k+1)} j w^{k+1} E + y X p^k j w^{k+1} E \\
 &= L^{(k+1)}(; w) - L^{(k)}(; w^{k+1}) + X^{(k+1)} j w^{k+1} E + h y j w^{k+1} E + X^{(k)} j w^{k+1} E \\
 &\quad + y j w^{k+1} E + X^{(k+1)} j w^{k+1} E + X^{(k)} j w^{k+1} E \\
 &\quad + y j w^{k+1} E + X p^k j w^{k+1} E \\
 &= L^{(k+1)}(; w) - L^{(k)}(; w^{k+1}) + y b j w^{k+1} E \\
 &\quad + X^{(k+1)} X p^k j w^{k+1} E \\
 &= L^{(k+1)}(; w) - L^{(k)}(; w^{k+1}) + \frac{k - k^2}{2} k w^{k+1} - w^k k^2 - 1 \\
 &\quad + X^{(k+1)} X p^k j w^{k+1} E
 \end{aligned} \tag{4.3.8}$$

From (4.3.8) and (4.3.7) we obtain

$$L(\cdot; w^{k+1}) - L(\cdot; w^k) + \frac{k^{k+1} p^k k^2}{2} + \frac{k w^{k+1} w^k k^2}{2} + V(z^{k+1}; z) - V(z^k; z) - k^{\frac{1}{2}} k w^{k+1} w^k + \frac{e^2}{2} \quad (4.3.9)$$

$$= \frac{D}{E} X(\cdot; p^k) j w^{k+1} w + \frac{D}{E} X(\cdot; p^{k-1}) j w^k w + \frac{D}{E} X(\cdot; p^k) j w^{k+1} w^k + \frac{D}{E} X(\cdot; p^{k-1}) j w^k w + \frac{D}{E} \frac{1}{2} X(\cdot; p^k) j \frac{1}{2} (w^{k+1} w^k) + \frac{D}{E} X(\cdot; p^k) j w^{k+1} w + \frac{D}{E} X(\cdot; p^{k-1}) j w^k w + k^{\frac{1}{2}} X(\cdot; p^k) j \frac{1}{2} k^2 \frac{k w^{k+1} w^k k^2}{2} + \frac{k^k p^{k-1} k^2}{2} \quad (4.3.10)$$

Then, recalling that $\frac{D}{E} = 1 - k^{\frac{1}{2}} X(\cdot; p^k) j \frac{1}{2} k^2$, we have the following estimate

$$L(\cdot; w^{k+1}) - L(\cdot; w^k) + \frac{k^{k+1} p^k k^2}{2} - \frac{k^k p^{k-1} k^2}{2} + \frac{k w^{k+1} w^k k^2}{2} + V(z^{k+1}; z) - V(z^k; z) - k^{\frac{1}{2}} k w^{k+1} w^k + \frac{D}{E} X(\cdot; p^k) j w^{k+1} w + \frac{D}{E} X(\cdot; p^{k-1}) j w^k w + \frac{e^2}{2}:$$

Summing from 1 to $N-1$, we obtain

$$\sum_{k=1}^{N-1} L(\cdot; w^{k+1}) - L(\cdot; w^k) + \frac{k^N p^{N-1} k^2}{2} - \frac{k^1 p^0 k^2}{2} + \frac{1}{2} \sum_{k=1}^{N-1} k w^{k+1} w^k k^2 + V(z^N; z) - V(z^1; z) - \sum_{k=1}^{N-1} X(\cdot; p^k) j w^{k+1} w + \frac{1}{2} \sum_{k=1}^{N-1} k X(\cdot; p^k) j w^k w + \frac{(N-1)e^2}{2} + \frac{1}{2} \sum_{k=1}^{N-1} X(\cdot; p^k) j \frac{1}{2} (w^{k+1} w^k) + \frac{1}{2} \sum_{k=1}^{N-1} X(\cdot; p^k) j w^{k+1} w + \frac{1}{2} \sum_{k=1}^{N-1} X(\cdot; p^{k-1}) j w^k w + \frac{(N-1)e^2}{2} + \frac{k^N p^{N-1} k^2}{2} + k^{\frac{1}{2}} X(\cdot; p^k) j \frac{1}{2} k^2 \frac{k w^N w^k}{2} \quad (4.3.11)$$

Now, by choosing $k = 0$ in (4.3.9) we get

$$\begin{aligned} L(-1; w) &= L(-; w^1) + \frac{k^{-1} p^0 k^2 - 1}{2} + \frac{1}{2} k w^1 - w^0 k^2 - 1 \\ &+ V(z^1 - z) - V(z^0 - z) + X(-1 - p^0) j w^1 - w \\ &= k^{-\frac{1}{2}} k w^1 - w k^{-1} + \frac{e^{-2}}{2}. \end{aligned} \quad (4.3.12)$$

Adding (4.3.11) and (4.3.12) we obtain

$$\begin{aligned} \sum_{k=0}^{N-1} L(-k+1; w) &= L(-; w^{k+1}) + \frac{1}{2} k w^N - w k^2 - 1 \\ &+ \sum_{k=1}^N \frac{1}{2} k w^k - w^{k-1} k^2 - 1 + \frac{k^{-N} - k^2 - 1}{2} \\ &= k^{-\frac{1}{2}} k \sum_{k=1}^N k w^k - w k^{-1} + V(z^0 - z) + \frac{N e^{-2}}{2} \end{aligned} \quad (4.3.13)$$

Next, by (4.3.10), we have the following estimate

$$\begin{aligned} &\frac{k^{-k+1} p^k k^2 - 1}{2} \sum_{j=0}^D X(-k - p^k - 1) j w^{k+1} - w^k \sum_{j=0}^E \\ &+ \frac{k w^{k+1} - w^k k^2 - 1}{2} + L(-k+1; w) - L(-; w^{k+1}) \\ &+ V(z^{k+1} - z) - V(z^k - z) \\ &= k^{-\frac{1}{2}} k \sum_{j=0}^D X(-k+1 - p^k) j w^{k+1} - w \sum_{j=0}^E \\ &+ \sum_{j=0}^D X(-k - p^k - 1) j w^k - w + \frac{e^{-2}}{2} \end{aligned}$$

Summing from 1 to $N-1$ we obtain

$$\begin{aligned} &\sum_{k=1}^{N-1} \frac{k^{-k+1} p^k k^2 - 1}{2} \sum_{j=0}^D X(-k - p^k - 1) j w^{k+1} - w^k \sum_{j=0}^E \\ &+ \sum_{k=1}^{N-1} L(-k+1; w) - L(-; w^{k+1}) + V(z^N - z) - V(z^1 - z) \\ &= k^{-\frac{1}{2}} k \sum_{k=1}^{N-1} k w^{k+1} - w k^{-1} + \frac{(N-1)e^{-2}}{2} \\ &\quad + \sum_{k=1}^{N-1} X(-N - p^N - 1) j w^N - w \\ &= k^{-\frac{1}{2}} k \sum_{k=1}^{N-1} k w^{k+1} - w k^{-1} + \frac{(N-1)e^{-2}}{2} \\ &\quad + \sum_{j=0}^D \frac{1}{2} X(-\frac{1}{2} - \frac{1}{2}(N - p^N - 1) j - \frac{1}{2}(w^N - w)) \sum_{k=1}^{N-1} \\ &= k^{-\frac{1}{2}} k \sum_{k=1}^{N-1} k w^{k+1} - w k^{-1} + \frac{(N-1)e^{-2}}{2} \\ &+ \frac{k^{-\frac{1}{2}} X(-\frac{1}{2} k^2 - N - p^N - 1) k^2 - 1}{2} + \frac{k w^N - w k^2 - 1}{2} \end{aligned} \quad (4.3.14)$$

Now, since $w^{k+1} - w^k = X p^k - y$ we derive that

$$\begin{aligned}
& \sum_{k=1}^N \frac{\|p^k\|^2}{2} \left(X^{(k)} p^{(k-1)} - y \right)^T \left(X^{(k)} p^{(k-1)} - y \right) + \frac{\|w^{k+1} - w^k\|^2}{2} \\
&= \sum_{k=1}^N \frac{\|p^k\|^2}{2} \left(X^{(k)} p^{(k-1)} - y \right)^T \left(X^{(k)} p^{(k-1)} - y \right) + \frac{\|w^{k+1} - w^k\|^2}{2} \\
&\quad + \frac{\|p^N\|^2}{2} \frac{\|p^0\|^2}{2} \\
&= \sum_{k=1}^N \frac{\|X^{(k)} p^{(k-1)} - y\|^2}{2} + \frac{\|X^{(k)} p^{(k-1)}\|^2}{2} + \frac{\|X p^k - y\|^2}{2} + \frac{\|X p^k - y\|^2}{2} \\
&\quad + \frac{\|p^N\|^2}{2} \frac{\|p^0\|^2}{2} \\
&= \sum_{k=1}^N \frac{\|X p^k - y\|^2}{2} + \frac{\|p^N\|^2}{2} \frac{\|p^0\|^2}{2} \\
&\quad + \sum_{k=1}^N \frac{\|p^k\|^2}{2} \frac{\|X^{(k)} p^{(k-1)}\|^2}{2} :
\end{aligned}$$

Furthermore, since

$$\|X^{(k)} p^{(k-1)}\|^2 > 0;$$

we obtain

$$\begin{aligned}
& \sum_{k=1}^N \frac{\|X p^k - y\|^2}{2} + \frac{\|p^N\|^2}{2} \frac{\|p^0\|^2}{2} \\
&\quad + \sum_{k=1}^N \frac{\|p^k\|^2}{2} \frac{\|X^{(k)} p^{(k-1)}\|^2}{2} \\
&= \sum_{k=1}^N \frac{\|X p^k - y\|^2}{2} + \frac{\|p^N\|^2}{2} \frac{\|p^0\|^2}{2} \\
&\quad + \frac{1}{2} \sum_{k=1}^N \|X^{(k)} p^{(k-1)}\|^2 \\
&= \sum_{k=1}^N \frac{\|X p^k - y\|^2}{2} + \frac{\|p^N\|^2}{2} \frac{\|p^0\|^2}{2} \\
&\quad + \frac{1}{2} \|X^{(N)} p^{(N-1)}\|^2 + \frac{1}{2} \|X^{(1)} p^{(0)}\|^2 \\
&\quad + \frac{1}{2} \sum_{k=1}^N \|X^{(k+1)} p^{(k)}\|^2 :
\end{aligned}$$

In turn, by the convexity of k^{-k^2} , we obtain

$$\begin{aligned}
& \sum_{k=1}^{N-1} \frac{k^{-\frac{1}{2}}(Xp^k - y)k^2}{2} + \frac{k^{-N} p^{N-1}k^2}{2} - \frac{k^{-1} p^0 k^2}{2} \\
& \quad - \frac{1}{2} k^{-\frac{1}{2}} X(N - p^{N-1})k^2 + \frac{1}{2} k^{-\frac{1}{2}} X(-1 - p^0)k^2 \\
& \quad + \frac{1}{2} \sum_{k=1}^{N-1} k^{-\frac{1}{2}} X(-k+1 - p^k)k^2 \\
& \quad - \frac{1}{4} \sum_{k=1}^{N-1} k^{-\frac{1}{2}} (X^{-k+1} - y)k^2 - \frac{k^{-1} p^0 k^2}{2} + \frac{1}{2} k^{-\frac{1}{2}} X(-1 - p^0)k^2 \\
& \quad + \frac{2 + k^{-\frac{1}{2}} X^{-\frac{1}{2}} k^2}{2} k^{-N} p^{N-1} k^2 \\
& \quad - \frac{1}{4} \sum_{k=2}^N k^{-\frac{1}{2}} (X^{-k} - y)k^2 - \frac{k^{-1} p^0 k^2}{2} + \frac{1}{2} k^{-\frac{1}{2}} X(-1 - p^0)k^2 \\
& \quad + \frac{2 + k^{-\frac{1}{2}} X^{-\frac{1}{2}} k^2}{2} k^{-N} p^{N-1} k^2 : \tag{4.3.15}
\end{aligned}$$

On the other hand, we get

$$\begin{aligned}
& k^{-\frac{1}{2}} (X^{-k} - y)k^2 - \frac{kX^{-k} - yk^2}{k^{-1}k} \\
& \quad - \frac{1}{k^{-1}k} \frac{kX^{-k} - yk^2}{2} - ky - yk^2 : \tag{4.3.16}
\end{aligned}$$

Combining (4.3.12), (4.3.14), (4.3.15), and (4.3.16) we have that

$$\begin{aligned}
& \sum_{k=0}^{N-1} L(-k+1; w) - L(-; w^{k+1}) + \frac{2}{2} k^{-N} p^{N-1} k^2 \\
& \quad - \sum_{k=1}^N \frac{1}{8k^{-1}k} kX^{-k+1} - bk^2 + \frac{k^{-N} k^2}{2} \\
& \quad - k^{-\frac{1}{2}} k^{-N} kW^k - w_{k-1} + V(z^0 - z) + \frac{Ne^{-2}}{2} + N \frac{1}{4k^{-1}k} \tag{4.3.17}
\end{aligned}$$

It remains to bound $\sum_{k=1}^N k^{-\frac{1}{2}} kW^k - w_{k-1}$. From (4.3.13) and since $(x; u)$ is a saddle-point of the Lagrangian we deduce that

$$kW^N - w_{k-1} \geq \frac{2k^{-\frac{1}{2}}}{k} \sum_{k=1}^N kW^k - w_{k-1} + \frac{2V(z^0 - z)}{2} + \frac{Ne^{-2}}{2} : \tag{4.3.18}$$

Applying [111, Lemma A.1] to Equation (4.3.18) with

$$k := \frac{2k^{-\frac{1}{2}}}{k} \quad \text{and} \quad S_N := \frac{2V(z^0 - z)}{2} + \frac{Ne^{-2}}{2}$$

we get

$$\begin{aligned} k w^N - w k & \leq \frac{N k^{\frac{1}{2}}}{2} + \frac{2V(z^0 - z)}{2} + \frac{N e^{-2}}{2} + \frac{N k^{\frac{1}{2}}}{2} \\ & \leq \frac{2N k^{\frac{1}{2}}}{2} + \frac{2V(z^0 - z)}{2} + \frac{N e^{-2}}{2}; \end{aligned} \quad (4.3.19)$$

Inserting (4.3.19) into (4.3.13) to obtain

$$\begin{aligned} \sum_{k=0}^{\infty} L(\cdot; w^{k+1}) - L(\cdot; w^{k+1}) & \leq \frac{2(N k^{\frac{1}{2}})^2}{2} + N k^{\frac{1}{2}} \frac{V(z^0 - z)}{2} \\ & \quad + N k^{\frac{1}{2}} \frac{N e^{-2}}{2} + V(z^0 - z) + \frac{N e^{-2}}{2}; \end{aligned}$$

By (4.3.17), we have

$$\begin{aligned} \sum_{k=1}^{\infty} k X^k - y k^2 & \leq \frac{16N^2 k^{\frac{1}{2}} k^{\frac{1}{2}}}{2} + 8N k^{\frac{1}{2}} \frac{2k^{\frac{1}{2}} V(z^0 - z)}{3} \\ & \quad + 8N^2 k^{\frac{1}{2}} \frac{k^{\frac{1}{2}} e N^{\frac{1}{2}}}{3} + \frac{8k^{\frac{1}{2}} V(z^0 - z)}{3} \\ & \quad + 2N^2 + \frac{4N k^{\frac{1}{2}} e^{-2}}{3} \end{aligned}$$

and both results are straightforward from Jensen's inequality. \square

Note that, the previous proof combines and extends the techniques developed in [24] and [84], based on the firm non-expansivity of the proximal point operator and discrete Bihari's lemma to deal with the error; see also [111].

In the next result, we establish upper bounds for the Lagrangian and feasibility gap analogous to those proposed in Theorem 4.3.1, but for the algorithm DPA. The main difference is that now the activation step is based on a priori information in the dual space \mathbb{R}^d , and not on S . This set is represented by the intersection of fixed point sets of a finite number of operators and encodes some knowledge about the dual solution.

Theorem 4.3.2. Consider the setting of DPA under Assumptions A1, A2, and A4. Let $(v^0; w^0; \lambda^0) \in \mathbb{R}^{2d} \times \mathbb{R}^p$ be such that $w^0 = v^0$. Then, for every $z = (\cdot; w) \in Z$ and for every $N \geq N$, we have that

$$L^{\wedge N}; w - L(\cdot; w^N) \leq \frac{V(z^0 - z)}{N} + 2k^{\frac{1}{2}} k^2 N^{-2} + k^{\frac{1}{2}} k^{\frac{1}{2}} 2V(z^0 - z)^{\frac{1}{2}}; \quad (4.3.20)$$

and

$$\begin{aligned} k X^{\wedge N} - y k^2 & \leq \frac{8k^{\frac{1}{2}} k^2 k^{\frac{1}{2}} N^{-2}}{3} + \frac{4k^{\frac{1}{2}} k^{\frac{1}{2}} k^{\frac{1}{2}} 2V(z^0 - z)^{\frac{1}{2}}}{3} \\ & \quad + \frac{4k^{\frac{1}{2}} k^{\frac{1}{2}} V(z^0 - z)}{3} + 2^{-2}; \end{aligned} \quad (4.3.21)$$

where we recall that the constant γ is defined in Assumptions A1.

Proof. It follows from [DPA](#) that

$$\begin{aligned} & \mathbb{1}(v^k - w^{k+1}) + X^{k+1} y = y \\ & \mathbb{1}(v^k - w^{k+1}) + X^{k+1} y = y \end{aligned} \tag{4.3.22}$$

Thus,

$$R^{(k+1)} + \mathbb{D} \mathbb{1}(v^k - w^{k+1}) + X^{k+1} y = y \tag{4.3.23}$$

and [\(4.3.23\)](#) yields

$$\begin{aligned} & 0 = R^{(k+1)} - R^{(k)} + \mathbb{D} \mathbb{1}(v^k - w^{k+1}) + X^{k+1} y - y \\ & = R^{(k+1)} - R^{(k)} + \frac{k^k - (k+1)^k}{2} + \frac{k^{k+1} - (k+1)^{k+1}}{2} \\ & \quad - \frac{k^k - (k+1)^k}{2} + \mathbb{D} \mathbb{1}(v^k - w^{k+1}) + X^{k+1} y - y \end{aligned} \tag{4.3.24}$$

From [\(4.3.22\)](#), it follows that

$$\begin{aligned} & 0 = \mathbb{D} \mathbb{1}(v^k - w^{k+1}) + X^{k+1} y - y \\ & 0 = \frac{k w^{k+1} - v^k k^2}{2} + \frac{k w^{k+1} - v^k k^2}{2} - \frac{k v^k - w k^2}{2} \\ & \quad + y - X^{k+1} y + w^{k+1} \end{aligned} \tag{4.3.25}$$

Recall that $z := (v; w) \in Z \subseteq \mathbb{R}^d$, $z^k := (v^k; w^k)$, and $V(z) := \frac{k^k - (k+1)^k}{2} + \frac{k w k^2 - v^k k^2}{2}$. Summing [\(4.3.24\)](#) and [\(4.3.25\)](#), we obtain

$$\begin{aligned} & R^{(k+1)} - R^{(k)} + \frac{k^{k+1} - (k+1)^{k+1}}{2} + \frac{k w^{k+1} - v^k k^2}{2} + V(z^{k+1} - z^k) \\ & \quad + V(z^k - z) + X^{k+1} y - y + w^{k+1} = 0 \end{aligned} \tag{4.3.26}$$

Now compute

$$\begin{aligned}
 & R(\begin{smallmatrix} k+1 \\ D \end{smallmatrix} ; w) L(\begin{smallmatrix} ; \\ E \end{smallmatrix} ; w^{k+1}) + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j v^k + y \begin{smallmatrix} E \\ D \end{smallmatrix} X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w^{k+1} \\
 &= L(\begin{smallmatrix} k+1 \\ D \end{smallmatrix} ; w) L(\begin{smallmatrix} ; \\ E \end{smallmatrix} ; w^{k+1}) X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w + X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w^{k+1} \\
 &\quad + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j v^k + y \begin{smallmatrix} E \\ D \end{smallmatrix} X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w^{k+1} \\
 &= L(\begin{smallmatrix} k+1 \\ D \end{smallmatrix} ; w) L(\begin{smallmatrix} ; \\ E \end{smallmatrix} ; w^{k+1}) X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w + h y j w i + X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w^{k+1} \\
 &\quad + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j v^k + y \begin{smallmatrix} E \\ D \end{smallmatrix} X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w^{k+1} \\
 &\quad + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j w^{k+1} + X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w \\
 &= L(\begin{smallmatrix} k+1 \\ D \end{smallmatrix} ; w) L(\begin{smallmatrix} ; \\ E \end{smallmatrix} ; w^{k+1}) + y \begin{smallmatrix} E \\ D \end{smallmatrix} y j w^{k+1} w \\
 &\quad + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j v^k w^{k+1} \tag{4.3.27}
 \end{aligned}$$

$$\begin{aligned}
 &= L(\begin{smallmatrix} k+1 \\ D \end{smallmatrix} ; w) L(\begin{smallmatrix} ; \\ E \end{smallmatrix} ; w^{k+1}) + y \begin{smallmatrix} E \\ D \end{smallmatrix} y j w^{k+1} w \\
 &\quad + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j v^k w^{k+1} + X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w^k v^{k-1} \\
 &= L(\begin{smallmatrix} k+1 \\ D \end{smallmatrix} ; w) L(\begin{smallmatrix} ; \\ E \end{smallmatrix} ; w^{k+1}) + y \begin{smallmatrix} E \\ D \end{smallmatrix} y j w^{k+1} w \\
 &\quad + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j v^k w^{k+1} + X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w^k v^{k-1} \\
 &\quad + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j w^k v^{k-1} \\
 &= L(\begin{smallmatrix} k+1 \\ D \end{smallmatrix} ; w) L(\begin{smallmatrix} ; \\ E \end{smallmatrix} ; w^{k+1}) + y \begin{smallmatrix} E \\ D \end{smallmatrix} y j w^{k+1} w \\
 &\quad + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j v^k w^{k+1} + X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j w^k v^{k-1} \\
 &\quad + \frac{1}{2} X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j \frac{1}{2} (w^k v^{k-1}) : \tag{4.3.28}
 \end{aligned}$$

From (4.3.28) and (4.3.26) we obtain

$$\begin{aligned}
 & L(\begin{smallmatrix} k+1 \\ D \end{smallmatrix} ; w) L(\begin{smallmatrix} ; \\ E \end{smallmatrix} ; w^{k+1}) + \frac{k^{k+1} k^{k^2-1}}{2} \\
 &+ \frac{k w^{k+1} v^{k^2-1}}{2} + V(z^{k+1}, z) V(z^k, z) \\
 &\quad + y \begin{smallmatrix} E \\ D \end{smallmatrix} y j w^{k+1} w X(\begin{smallmatrix} k+1 \\ E \end{smallmatrix}) j v^k w^{k+1} \\
 &\quad + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j v^k w^{k+1} + \frac{1}{2} X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j \frac{1}{2} (w^k v^{k-1}) \\
 &\quad + \frac{1}{2} k k w^{k+1} w^{k-1} X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j v^k w^{k+1} \\
 &\quad + X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j v^k w^{k+1} + \frac{k^{k+1} k^{k^2-1}}{2} \\
 &\quad + k \frac{1}{2} X(\begin{smallmatrix} k+1 \\ D \end{smallmatrix}) j \frac{1}{2} k^2 \frac{k w^k v^{k-1} k^2-1}{2}
 \end{aligned}$$

Therefore we have that

$$\begin{aligned}
L(z^{k+1}; w) &= L(z; w^{k+1}) + \frac{kw^{k+1} v^k k^2}{2} \\
&\quad + k \frac{1}{2} X \frac{1}{2} k^2 \frac{kw^k v^{k-1} k^2}{2} + V(z^{k+1} z) - V(z^k z) \\
&\quad + k \frac{1}{2} k w^{k+1} w k \frac{1}{2} X \frac{1}{2} k^2 \frac{kw^k v^{k-1} k^2}{2} + V(z^{k+1} z) - V(z^k z) \\
&\quad + X(z^k) j v^{k-1} w^k
\end{aligned} \tag{4.3.29}$$

Summing from 1 to $N-1$ we obtain

$$\begin{aligned}
\sum_{k=1}^{N-1} L(z^{k+1}; w) &= \sum_{k=1}^{N-1} L(z; w^{k+1}) + \frac{1}{2} \sum_{k=1}^{N-1} kw^{k+1} v^k k^2 \\
&\quad + V(z^N z) + k \frac{1}{2} X \frac{1}{2} k^2 \frac{kw^N v^{N-1} k^2}{2} \\
&\quad + k \frac{1}{2} k \sum_{k=1}^{N-1} kw^{k+1} w k \frac{1}{2} X \frac{1}{2} k^2 \frac{kw^k v^{k-1} k^2}{2} + V(z^1 z) \\
&\quad + X(z^1) j v^0 w^1 + V(z^1 z) \\
&\quad + k \frac{1}{2} k \sum_{k=1}^{N-1} kw^{k+1} w k + k \frac{1}{2} X \frac{1}{2} k^2 \frac{kw^N v^{N-1} k^2}{2} \\
&\quad + \frac{k^N}{2} + X(z^1) j v^0 w^1 + V(z^1 z)
\end{aligned} \tag{4.3.30}$$

Reordering (4.3.30) we obtain

$$\begin{aligned}
\sum_{k=1}^{N-1} L(z^{k+1}; w) &= \sum_{k=1}^{N-1} L(z; w^{k+1}) \\
&\quad + \frac{1}{2} \sum_{k=1}^{N-1} kw^{k+1} v^k k^2 + \frac{kw^N}{2} w k^2 \\
&\quad + k \frac{1}{2} k \sum_{k=1}^{N-1} kw^{k+1} w k + X(z^1) j v^0 w^1 + V(z^1 z):
\end{aligned} \tag{4.3.31}$$

On the other hand, from (4.3.26) and (4.3.27) we get

$$\begin{aligned}
L(z^1; w) &= L(z; w^1) + \frac{1}{2} kw^1 v^0 k^2 + kw^1 w k \\
&\quad + X(z^1) j v^0 w^1 \\
&\quad + V(z^0 z) - V(z^1 z)
\end{aligned} \tag{4.3.32}$$

Summing (4.3.31) and (4.3.32) yields

$$\begin{aligned}
\sum_{k=1}^N L(z^k; w) &= \sum_{k=1}^N L(z; w^k) \\
&\quad + \frac{1}{2} \sum_{k=1}^N kw^k v^{k-1} k^2 + \frac{kw^N}{2} w k^2 \\
&\quad + k \frac{1}{2} k \sum_{k=1}^N kw^k w k + V(z^0 z):
\end{aligned} \tag{4.3.33}$$

Moreover, since $w^{k+1} v^k = (X^{k+1} y)$ we have

$$\begin{aligned} k w^{k+1} v^k k^2 &= \frac{D}{(X^{k+1} y) j X^{k+1} y} E \\ &= \frac{k X^{k+1} y k^2}{k^{-1} k} \\ &= \frac{1}{k^{-1} k} \frac{k X^{k+1} y k^2}{2} k y y k^2 \end{aligned} \quad (4.3.34)$$

and from (4.3.33) and (4.3.34) we obtain

$$\begin{aligned} & \sum_{k=0}^{N-1} L(k+1; w) - L(0; w^{k+1}) \\ & + \frac{1}{4k^{-1}k} \sum_{k=1}^N k X^k y k^2 + \frac{k w^N w k^2}{2} \\ & - \sum_{k=1}^{N-1} k^{-\frac{1}{2}k} k w^k w k^{-1} + V(z^0 - z) + \frac{N^2}{2k^{-1}k} \end{aligned} \quad (4.3.35)$$

From (4.3.33) it follows that

$$k w^N w k^2 - \sum_{k=1}^{N-1} k^{-\frac{1}{2}k} k w^k w k^{-1} + 2V(z^0 - z); \quad (4.3.36)$$

Applying [111, Lemma A.1] to Equation (4.3.36) with $k := 2^{-k} k$ and $S_k := 2V(z^0 - z)$ to get

$$\begin{aligned} k w^k w k^{-1} & \leq N k^{-\frac{1}{2}k} + 2V(z^0 - z) + N k^{-\frac{1}{2}k} 2^{-\frac{1}{2}} \\ & \leq 2N k^{-\frac{1}{2}k} + 2V(z^0 - z)^{\frac{1}{2}} \end{aligned} \quad (4.3.37)$$

Insert the previous in Equation (4.3.33), to obtain

$$\begin{aligned} & \sum_{k=0}^{N-1} L(k+1; w) - L(0; w^{k+1}) \leq 2k^{-\frac{1}{2}k} N^2 + N k^{-\frac{1}{2}k} 2V(z^0 - z)^{\frac{1}{2}} \\ & \quad + V(z^0 - z) \end{aligned} \quad (4.3.38)$$

and by (4.3.35) and (4.3.37) we have

$$\begin{aligned} & \sum_{k=1}^N k A^k b k^2 \leq \frac{4k^{-1}k}{2k^{-\frac{1}{2}k} N^2 + N k^{-\frac{1}{2}k} 2V(z^0 - z)^{\frac{1}{2}}} \\ & \leq \frac{4k^{-1}k}{V(z^0 - z) + \frac{N^2}{2k^{-1}k}} \end{aligned} \quad (4.3.39)$$

and both results follows from Jensen's inequality. \square

First, we comment on the chosen optimality measures. As discussed in [84, 85, 111], the Lagrangian gap is equivalent to the Bregman distance of the iterates to the solution. If the penalty is strongly convex, the Bregman divergence is an upper bound of the squared norm of the difference between the reconstructed and the ideal solution, while if R is only convex, the Bregman divergence gives only limited information, and in general, it is a very

weak convergence measure. For instance, in the exact case, a vanishing Lagrangian gap does not imply that cluster points of the generated sequence are primal solutions. However, as can be derived from [85], a vanishing Lagrangian gap coupled with a vanishing feasibility gap implies that every cluster point of the primal sequence is a solution to the primal problem.

In both theorems, the established result ensures that the two optimality measures can be upper bounded with the sum of two terms. The first one, which can be interpreted as an optimization error, is of the order $O(N^{-1})$, and so it goes to zero as N tends to $+\infty$. Note that, in the exact case $\epsilon = 0$, only this term is present and both the Lagrangian and the feasibility gap are indeed vanishing, guaranteeing that every cluster point of the sequence is a primal solution. The second term, which can be interpreted as a stability control, collects all the errors due to the perturbation of the exact datum and also takes into account the presence of the activation operators T , when the data constraints are noisy. It is an increasing function of the number of iterations and the noise level ϵ .

Remark 4.3.3. Theorems 4.3.1 and 4.3.2 are an extension of [24, Theorem 1], where it is proved that the sequence generated by the algorithms converges to an element in Z when $\epsilon = 0$, but neither convergence rates nor stability bounds were given. In this work, we filled the gap for linearly constrained convex optimization problems. Moreover, in the noise-free case, our assumptions on the additional operators T are weaker than those proposed in [24], where ϵ -averagedness is required. For the noisy case, without the activation operators (and so with $\epsilon = 0$), our bounds are of the same order as in [84] in the number of iterations and noise level ϵ .

As mentioned above, in (4.3.1) and (4.3.2), when $\epsilon > 0$ and $N \rightarrow +\infty$, the upper bounds for the PDA iterates tend to infinity, and the iteration may not converge to the desired solution. The same comment can be made for the DPA iterates, based on (4.3.20) and (4.3.21). In both cases, to obtain a minimal reconstruction error, we need to impose a trade-off between convergence and stability. The next corollary introduces an early stopping criterion, depending only on the noise level and leading to stable reconstruction.

Corollary 4.3.4. (Early-stopping). Under the assumptions of Theorem 4.3.1 or Theorem 4.3.2, choose $N = c\epsilon^{-1}$ for some $c > 0$. Then, for every $z = (x; w) \in Z$, there exist constants c_1 , c_2 , and c_3 such that

$$\begin{aligned} L(x^N; w) - L(x; w) &\leq c_1 \\ \|X^N(x^N) - y\|^2 &\leq c_2 + c_3 \epsilon^2. \end{aligned}$$

The early stopping rule prescribed above is computationally efficient in the sense that the number of iterations is proportional to the inverse of the noise level. In particular, if the error ϵ is small, then more iterations are useful, while if ϵ is big, it is convenient to stop sooner. So, the number of iterations plays the role of a regularization parameter. Using the early stopping strategy proposed above, we can see that the error in the data transfers to the error in the solution with the same noise level, which is the best that one can expect for a general operator X .

Remark 4.3.5. Comparison with Tikhonov regularization. The reconstruction properties of the proposed algorithms are comparable to the ones obtained using Tikhonov regularization [18, 52], with the same dependence on the noise level. We underline that in [18, Theorem 5.1] only the Bregman divergence is considered and not the feasibility. In addition, iterative regularization is way more efficient from a computational point of view, as it requires the solution of only one optimization problem, while Tikhonov regularization amounts to solving a family of problems indexed by the regularization parameter. Let us also note that, when α is unknown, any principle used to determine a suitable α can be used to determine the stopping time.

4.4 Implementation details

In this section, we discuss some standard choices to construct non-expansive operators T that satisfy our assumptions and encode some redundant information on the solution set. We first present examples for PDA, and later for DPA.

To define the operators, we first recall how to compute the projection on the constraint determined by each datum. For every $j \in [d]$, we denote by x_j the j -th row of X and by P_j the projection onto the j -th linear equation; namely,

$$P_j: \mathbb{R}^p \rightarrow \mathbb{R}^p; \quad x \mapsto x - \frac{y_j - \langle x, x_j \rangle}{\|x_j\|^2} x_j.$$

Analogously, for every $j \in [d]$, we denote by P_j^y the projection operator as in the previous definition but with the noisy data y instead of y_j .

We proceed to define the four families of operators proposed in this chapter for PDA.

Definition 4.4.1. Consider the operator $T: \mathbb{R}^p \rightarrow \mathbb{R}^p$ is a

1. **Serial projection** if

$$T = P_{j_1} \circ \dots \circ P_{j_l};$$

where, for every $j \in [l]$, $j_j \in [d]$.

2. **Parallel projection** if

$$T = \sum_{j=1}^l \alpha_j P_{j_j} \tag{4.4.1}$$

where, for every $j \in [l]$, $j_j \in [d]$ and $(\alpha_j)_{j=1}^l$ are real numbers in $[0; 1]$, such that $\sum_{j=1}^l \alpha_j = 1$.

3. **Landweber operator** with parameter α if

$$T: \mathbb{R}^p \rightarrow \mathbb{R}^p; \quad x \mapsto X^\top (X - \alpha I)^{-1} X x \tag{4.4.2}$$

where $\alpha \in]0; \frac{2}{\|X\|^2}[$.

4. Landweber operator with adaptive step and parameter M if

$$T: \mathbb{R}^p \rightarrow \mathbb{R}^p; \quad T(x) = \begin{cases} X^T(X - y) & \text{if } \|X\| \leq \|X - y\| \\ \text{otherwise.} \end{cases} \quad (4.4.3)$$

where, for $M > 0$, $\| \cdot \| = \min_{kX \geq (X - y)^2} \frac{kX - y\|k^2}{kX\|k^2}; M$.

The next lemma states that the operators in Definition 4.4.1 satisfy Assumption A3.

Lemma 4.4.2. Let $T: \mathbb{R}^p \rightarrow \mathbb{R}^p$ be one of the operators given in Definition 4.4.1. Then Assumption A3 holds with the following

1. If T is a serial projection, then

$$e_T = \frac{1}{\max_{i=1, \dots, d} \|x_i\|}$$

2. If T is a parallel projection, then

$$e_T = \sum_{j=1}^d \frac{1}{\|x_j\|^2}$$

3. If T is the Landweber operator with parameter η , then

$$e_T = \frac{1}{2\eta \|X\|^2}$$

4. If T is the Landweber operator with adaptive step and parameter M , then $e_T = M$:

Proof. Let us first recall that

$$P_j = I + \frac{y_j}{\|x_j\|^2} x_j x_j^T$$

Note that the j -th equation of S and S^E are parallel. Then, for every $j \in [d]$ and $x_j \in S$, we get

$$\begin{aligned} \|kP_j - k\|^2 &= \|kP_j - k\|^2 + 2 \langle kP_j - k, x_j \rangle x_j^T + \|x_j\|^2 \\ &= \|kP_j - k\|^2 + 2 \langle kP_j - k, x_j \rangle \frac{y_j}{\|x_j\|^2} + \|x_j\|^2 \end{aligned} \quad (4.4.4)$$

Analogously, we have that

$$\|k - kP_j\|^2 = \|k - kP_j\|^2 + 2 \langle k - kP_j, x_j \rangle \frac{y_j}{\|x_j\|^2} + \|x_j\|^2 \quad (4.4.5)$$

It follows from (4.4.4) and (4.4.5) that

$$\|kP_j - k\|^2 + \|k - kP_j\|^2 = 2 \|k - kP_j\|^2 + 2 \langle k - kP_j, x_j \rangle \frac{y_j}{\|x_j\|^2} + 2 \langle kP_j - k, x_j \rangle \frac{y_j}{\|x_j\|^2} + 2 \|x_j\|^2$$

Hence,

$$\begin{aligned} \|kP_j - k\|^2 &= \|k - kP_j\|^2 + \frac{(y_j - y_j)^2}{\|x_j\|^2} \\ &= \|k - kP_j\|^2 + \frac{0}{\|x_j\|^2} \end{aligned}$$

1. Since $T = P_{j, \cdot} P_{\cdot, 1}$ it is clear that $S = \text{Fix } T$ and by induction we have that,

$$\|kT - k^2\|_k \leq \|k^2 - k\|_k + e^2;$$

where $e = \frac{1}{\max_{i=1, \dots, d} \|kx_i\|_k}$.

2. The proof follows from the convexity of $\|k - k^2\|_k$ which is obtained with $e = \frac{1}{\max_{i=1, \dots, d} \|kx_i\|_k}$.

3. Let $S \neq \emptyset$. By (4.4.2), we have

$$\begin{aligned} \|kT - k^2\|_k &= \|k^2 - k\|_k + \frac{1}{2} \|kX^2 - kX\|_k^2 \\ &+ \frac{1}{2} \|kX^2 - kX\|_k^2 \\ &= \|k^2 - k\|_k + \frac{1}{2} \|kX^2 - kX\|_k^2 + \frac{1}{2} \|kX^2 - kX\|_k^2 \\ &= \|k^2 - k\|_k + \frac{1}{2} \|kX^2 - kX\|_k^2 + \frac{1}{2} \|kX^2 - kX\|_k^2 \end{aligned}$$

Now using the Young inequality with parameter $\frac{1}{2} \|kX^2 - kX\|_k^2$, we have that

$$\|kT - k^2\|_k \leq \|k^2 - k\|_k + \frac{1}{2} \|kX^2 - kX\|_k^2 + \frac{1}{2} \|kX^2 - kX\|_k^2$$

It remains to prove that, if $S \neq \emptyset$, then $S = \text{Fix } T$, which is clear from (4.4.2).

4. Let $S \neq \emptyset$ and $x \in \mathbb{R}^p$. If $X^2 = Xy$, then (4.2.5) immediately holds. Otherwise, we have

$$\begin{aligned} \|kT - k^2\|_k &= \|k^2 - k\|_k + \frac{1}{2} \|kX^2 - kX\|_k^2 \\ &+ \frac{1}{2} \|kX^2 - kX\|_k^2 \\ &= \|k^2 - k\|_k + \frac{1}{2} \|kX^2 - kX\|_k^2 + \frac{1}{2} \|kX^2 - kX\|_k^2 \\ &= \|k^2 - k\|_k + \frac{1}{2} \|kX^2 - kX\|_k^2 + \frac{1}{2} \|kX^2 - kX\|_k^2 \\ &= \frac{1}{2} \|kX^2 - kX\|_k^2 + \frac{1}{2} \|kX^2 - kX\|_k^2 \end{aligned}$$

Now using the Young inequality with parameter $\frac{1}{2} \|kX^2 - kX\|_k^2$, we have that

$$\|kT - k^2\|_k \leq \|k^2 - k\|_k + \frac{1}{2} \|kX^2 - kX\|_k^2 + \frac{1}{2} \|kX^2 - kX\|_k^2$$

Finally, it is clear from (4.4.3) that, if $S \neq \emptyset$, then $S = \text{Fix } T$.

□

Remark 4.4.3. Relationship between Parallel projection and Landweber operator. A particular parallel projection is the one corresponding to $l = d$, $j = j$, and $\|j\|_j = \frac{\|kx_j\|_k^2}{\|kx\|_k^2}$. Then, (4.4.1) reduces to

$$T(\cdot) = \frac{1}{\|kx\|_k^2} X^2(X - y):$$

Observe that, since $kXk = kXk_F$, the previous is a special case of the Landweber operator with $\alpha = \frac{1}{kXk_F^2}$.

Remark 4.4.4. Steepest descent. Let $y \in \mathbb{R}^p$ be such that $Xy = b$. Then, from (4.4.3), we derive

$$\begin{aligned} kTy - yk^2 = k \sum_{j=1}^D y_j X_j - y \sum_{j=1}^E 2^{\alpha} kX_j - yk^2 \\ + \sum_{j=1}^E 2^{\alpha} kX_j - yk^2: \end{aligned} \quad (4.4.6)$$

If $\alpha = 0$, then the choice of α given in (4.4.3) minimizes the right-hand side of (4.4.6), if the minimizer is smaller than M . In this case, α is chosen in order to maximize the contractivity with respect to a fixed point of T . While we cannot repeat the same procedure for $\alpha > 0$, since we do not know y , we still keep the same choice. If $y \in \text{ran}(X)$, then $\sup_{y \in \mathbb{R}^p} kXy - yk^2 = kX(X - y)k^2 < +\infty$. However, in general, if $\alpha > 0$, this is not true and M is needed to ensure that α is bounded.

Remark 4.4.5. From a computational point of view, parallel projections and Landweber operators are more efficient than serial projections. In particular, note that the quantity $(X^k - y)$ needs to be computed anyway in the other steps of the algorithm.

While for the primal space the data constraints that we want to reuse are clearly given by the linear constraints, for the dual there is not always a natural choice. In the following we present an example related to the ℓ^1 norm. A similar implementation can be extended to the case of 1-homogenous penalty functions, for which the Fenchel conjugate is the indicator of a closed and convex subset of the dual space [15, Proposition 14.11 (ii)].

Example 4.4.6. Consider the noisy version of problem (P) with $R(\cdot) = k \cdot k_1$. Then the dual is given by

$$\min_{w \in \mathbb{R}^d} \sum_{j=1}^D y_j w_j : j(X^>w)_j \leq 1; \text{ for every } i \in [p] :$$

For every $i \in [p]$, set $D_i = \{w \in \mathbb{R}^d : j(X^>w)_j \leq 1\}$ and denote by T_i the projection over D_i . Note that this projection is easy to compute, see for example [15, Example 28.17], since it is the projection onto the intersection of two parallel half-hyperplane. Clearly Assumption A4 holds. Differently from the primal case, here we are projecting on exact constraints, which are independent of the noisy data y .

4.5 Numerical results

In this section, to test the efficiency of the proposed algorithms, we perform numerical experiments in two settings: sparse reconstruction with ℓ^1 -norm regularization and Image denoising and deblurring with total variation regularization. For the ℓ^1 -norm regularization, we compare our results with other regularization techniques. In the more complex problem of total variation we explore the properties of different variants of our procedure.

Code statement: All numerical examples are implemented in MATLAB[®] on a laptop. In the second experiment we also use the library Numerical tours [108]. The corresponding

code can be downloaded at <https://github.com/cristianvega1995/L1-TV-Experiments-of-Fast-iterative-regularization-by-reusing-data-constraints>

4.5.1 ℓ^1 -norm regularization

In this section, we apply the routines **PDA** and **DPA** when R is equal to the ℓ^1 -norm. We compare the results given by our method with two state-of-the-art regularization procedures: iterative regularization by vanilla primal-dual [84], and Tikhonov explicit regularization, solving each problem by using the forward-backward algorithm [45]. In addition, we compare to another classical optimization algorithm for the minimization of the sum of two non-differentiable functions, namely the Douglas-Rachford algorithm [26]. In the noise free case, this algorithm is very effective in terms of number of iterations, but at each iteration it requires the explicit projection on the feasible set. In the noisy case, a stability analysis of the latter is not available.

We use the four variants of the algorithm **PDA** corresponding to the different choices of the operators T in Definition 4.4.1 and the version of **DPA** described in Example 4.4.6. Unless otherwise stated, in all the experiments we use as preconditioners $\tilde{M} = \frac{0.99}{\|X\|_F} \text{Id}$, which both satisfy (4.2.4).

Let $d = 2260$, $p = 3000$, and let $X \in \mathbb{R}^{d \times p}$ be such that every entry of the matrix is an independent sample from $\mathcal{N}(0; 1)$, then normalized column by column. We set $y := X \tilde{w}$, where $\tilde{w} \in \mathbb{R}^p$ is a sparse vector with approximately 300 non-zero entries uniformly distributed in the interval $[0; 1]$. It follows from [53, Theorem 9.18] that \tilde{w} is the unique minimizer of the problem with probability bigger than 0.99. Let y^* be such that $y = y + \|y - y^*\|_2 w$ where the vector w is distributed, entry-wise, as uniformly on $[-0.2; 0.2]$. In this experiment, to test the reconstruction capabilities of our method, we use the exact datum to establish the best stopping time, i.e. the one minimizing $\|k - k^*\|_2$. The exact data solution is also used for the other regularization techniques. In a real practical situation, when both \tilde{w} and y^* are unknown, we would need to use parameter tuning techniques in order to select the optimal stopping time, but we do not address this aspect here.

We detail the used algorithms and their parameters below.

(Tik) **Tikhonov Regularization:** We consider a grid of penalty parameters

$$G = \left\{ \frac{1}{5}, \frac{1}{10}, \dots, \frac{1}{10^5} \right\} \cup \{d_k X^T y\} : d \in [5]; d \in [6]$$

and, for each value $\alpha \in G$, the optimization problem

$$\min_{\tilde{w} \in \mathbb{R}^p} \|k - k_1\|_2 + \frac{1}{2} \|k X - y\|_2^2 : \quad (4.5.1)$$

We solve each one of the previous problems with 300 iterations of forward-backward algorithm, unless the stopping criterion $\|k^{k+1} - k_k\|_2 \leq 10^{-3}$ is satisfied earlier. Moreover, to deal efficiently with the sequence of problems, we use warm restart [17]. We first solve problem (4.5.1) for the biggest value of α in G . Then, we initialize the algorithm for the next value of α , in decreasing order, with the solution reached for the previous one; and so on.

(DR) **Douglas Rachford:** see [26, Theorem 3.1].

(PD) **Primal-dual:** this corresponds to PDA with $m = 1$ and $T_1 = \text{Id}$.

- (PDS) **Primal-dual with serial projections:** at every iteration, we compute a serial projection using all the equations of the noisy system, where the order of the projections is given by a random shuffle.
- (PDP) **Primal-dual with parallel projections:** Set $m = 1$ and $T_1 = \frac{1}{kXk_F^2} X^>(X - y)$, see Remark 4.4.3.
- (PDL) **Primal-dual Landweber:** Set $m = 1$ and $T_1 = \frac{2}{kXk^2} X^>(X - y)$.
- (PDAL) **Primal-dual Landweber with adaptive step:** Set $m = 1$, and $T_1 = \frac{1}{kX^>(X - y)k^2} X^>(X - y)$, where $(\cdot) = \min \frac{kX - y k^2}{kX^>(X - y)k^2}; M$ for $M = 10^6$.
- (DPS) **Dual primal with serial projections:** at every iteration, we compute a serial projection over every inequality of $kX^>wk_1 = 1$, where the order is given by a random shuffle of the rows of $X^>$.

	Time [S]	Iteration	Reconstruction error
Tik	1.89	109	3.07
DR	3.08	5	5.01
PD	0.36	14	3.11
PDS	1.41	11	2.58
PDP	0.35	14	3.11
PDL	0.28	12	2.60
PDAL	0.27	11	2.56
DPS	0.54	17	2.83

Table 4.2: Run-time and number of iterations of each method until it reaches the best reconstruction error. We compare the proposed algorithms with Tikhonov regularization (Tik), Douglas-Rachford (DR), and iterative regularization (PD).

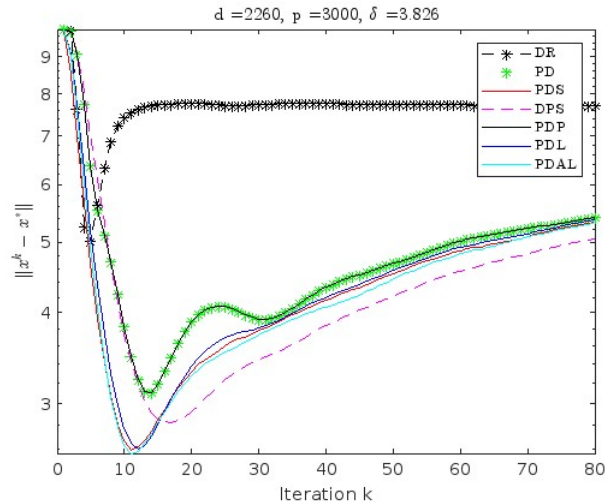


Figure 4.1: Graphical representation of early stopping. Note that the reconstruction error decreases and then increases, since the iterates first approach the exact solution and then converges to the noisy solution.

In Table 4.2, we reported also the number of iterations needed to achieve the best reconstruction error, but it is important to note that the iteration of each method has a different computational cost, so the run-time is a more appropriate comparison criterion.

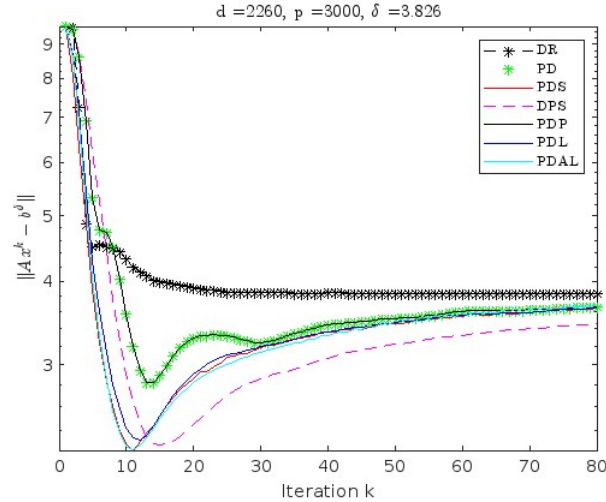


Figure 4.2: Early stopping with respect to the feasibility. Note that their behavior with respect to k is similar to that in Figure 4.1.

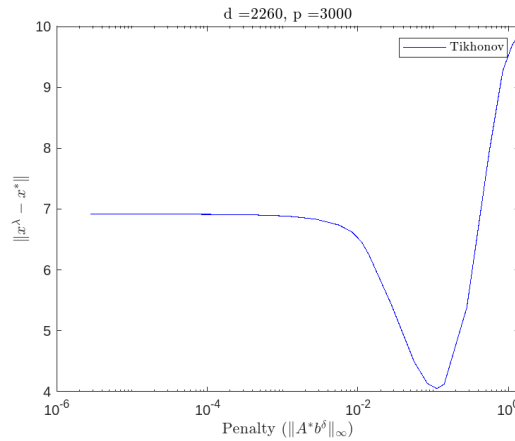


Figure 4.3: Reconstruction error of Tikhonov Method with different penalties.

Douglas-Rachford with early stopping is the regularization method performing worst on this example, both in terms of time and reconstruction error. This behavior may be explained by the fact that this algorithm converges fast (meaning in few iterations) convergence to the noisy solution, from which we infer that Douglas-Rachford is not a good algorithm for iterative regularization. Moreover, since we project on the noisy feasible set at every iteration, the resolution of a linear system is needed at every step. This also explains the cost of each iteration in terms of time. Note in addition that in our example y is in the range of X and so the noisy feasible set is non-empty. Tikhonov's regularization performs similarly in terms of time, but it requires many more (cheaper) iterations (see Figure 4.3). The achieved error is smaller than the one of DR, but bigger than the minimal one achieved by other methods.

Regarding our proposals, we observe that in Table 4.2 the proposed methods perform better than (PD). This supports the idea that reusing the constraints determined by the data is beneficial with respect to vanilla primal-dual. The benefit is not evident for (PDP), which achieves the worst reconstruction error, since kXk_F^2 is very big and so T_1 is very close to the identity. All other methods give better results in terms of reconstruction error. On the other hand, (PDS) is the slowest since it requires computing several projections

at each iteration in a serial manner. We also observe that (PDL) and (PDAL) have better performance improving 22.2% and 25.0% in reconstruction error and 16.4% and 17.7% in run-time.

Figure 4.1 empirically shows the existence of the trade-off between convergence and stability for all the algorithms, and therefore the advantage of early stopping. Similar results were obtained for the feasibility gap (see Figure 4.2).

4.5.2 Total variation

In this section, we perform several numerical experiments using the proposed algorithms for image denoising and deblurring. As done in the classical image denoising method introduced by Rudin, Osher, and Fantemi in [117], we rely on the total variation regularizer. See also [37, 39, 103, 105, 116, 117, 140]. We compare (PD) with (PDL) and (PDAL) algorithms, which were the algorithms performing the best in the previous application. In this section, we use two different preconditioners, which have been proved to be very efficient in practice [109].

Let $x \in \mathbb{R}^{N^2}$ represent an image with $N \times N$ pixels in $[0;1]$. We want to recover from a blurry and noisy measurement y , i.e. from

$$y = Kx + \epsilon;$$

where K is a linear bounded blurring operator and ϵ is a random noise vector. A standard approach is to assume that the original image is well approximated by the solution of the following constrained minimization problem:

$$\min_{u \in \mathbb{R}^{N^2}} \|Ku - y\|_{k_{1,2}} \quad (TV)$$

Here,

$$K \in \mathbb{R}^{N^2 \times N^2}; \quad p \times p \quad \sum_{i=1}^N \sum_{j=1}^N k_{p_{ij}};$$

and $D: \mathbb{R}^{N^2} \rightarrow (\mathbb{R}^2)^{N^2}$ is the discrete gradient operator for images, which is defined by

$$(Du)_{ij} = ((D_x u)_{ij}; (D_y u)_{ij})$$

with

$$(D_y u)_{ij} = \begin{cases} u_{i+1,j} - u_{i,j} & \text{if } 1 \leq i < N \\ 0 & \text{if } i = N \end{cases}$$

$$(D_x u)_{ij} = \begin{cases} u_{i,j+1} - u_{i,j} & \text{if } 1 \leq j < N \\ 0 & \text{if } j = N \end{cases}$$

In order to avoid the computation of the proximity operator of $KD \in \mathbb{R}^{2N^2}$, we introduce an auxiliary variable

$$v = Du \in Y := \mathbb{R}^{2N^2};$$

Since the value in each pixel must belong to $[0;1]$, we add the constraint $u \in X := [0;1]^{N^2}$. In this way, (TV) becomes

$$\min_{(u,v) \in X \times Y} (\|Ku - y\|_{k_{1,2}}; \|Du - v\|_{k_{1,2}}) \quad (TV)$$

Formulation and Algorithms

Problem (TV) is a special instance of (P), with

$$J: \mathbb{R}^{N^2} \rightarrow \mathbb{R}^2 \oplus \mathbb{R}^{N^2} \quad \forall x := (u; v) \quad \|x\|_{k_{1,2}} + \chi(u);$$

$$A = \begin{pmatrix} K & 0 \\ D & \text{Id} \end{pmatrix}; \quad y = \begin{pmatrix} y \\ 0 \end{pmatrix}; \quad \text{and } p = d = 3N^2;$$

Clearly, A is a linear non-zero operator, and $R \geq 0$ ($\mathbb{R}^{N^2} \oplus \mathbb{R}^2 \oplus \mathbb{R}^{N^2}$).

Primal-Dual for total variation	
Input:	$(p^0; p^1; v^0) \in \mathbb{R}^6 \oplus \mathbb{R}^{N^2} \oplus \mathbb{R}^{N^2}$ and $(q^0; q^1; z^0; w^0) \in \mathbb{R}^3 \oplus \mathbb{R}^{N^2} \oplus \mathbb{R}^{N^2}$.
For	$k = 1; \dots; L:$
	$\begin{aligned} v^{k+1} &= v^k + (K(p^k + p^{k-1}) - y) \\ w^{k+1} &= w^k - (q^k + z^k - q^{k-1}) + D(p^k + p^{k-1}) \\ p^{k+1} &= P_X(p^k - K(v^{k+1} + w^{k+1})) \\ z^{k+1} &= \text{prox}_{k_{1,2}}(q^k - D^* w^{k+1}) \\ p^{k+1} &= x^k - (x^k - K(x^k - y) + (Dx^k - z^k)) \\ q^{k+1} &= q^k - (x^k - Dx^k - z^k) \end{aligned} \tag{4.5.2}$
End	

Table 4.3: General form of the algorithms.

We compare the algorithms listed below. Note that all proposed algorithms are different instances of the general routine described in Table 4.3, and each of them corresponds to a different choice of (γ^k) :

1. PD, the vanilla primal-dual algorithm, corresponding to $(\gamma^k) = 0$;
2. PPD, the preconditioned primal-dual algorithm, obtained by $(\gamma^k) = 0$ and β^k as in [109, Lemma 2];
3. PDL, corresponding to $(\gamma^k) = 1/kXk^2$;
4. PDAL, corresponding to $(\gamma^k) = \wedge(\gamma^k)$ as (4.4.3).

Initializing by $p^0 = p^0 = 0$ and $q^0 = q^0 = z^0$, we recover the results of Theorem 4.3.1 and Corollary 4.3.4.

Remark 4.5.1. In order to implement the algorithm in 4.5.2, we first need to compute some operators.

1. It follows from [15, Proposition 24.11] and [15, Example 24.20] that

$$\text{prox}_{k_{1,2}}(v) = \text{prox}_{k_{1,2}}(v_i)_{i=1}^{N^2} = \begin{pmatrix} 1 \\ \max_{i=1}^{N^2} f_i \|v_i\| \end{pmatrix};$$

where $v_i \in \mathbb{R}^2$. The projection onto X can be computed as

$$P_X(u) = P_{[0,1]}(u_i)_{i=1}^{N^2};$$

where $P_{[0,1]}(u_i) = \min\{1, \max\{u_i, 0\}\}$:

2. It follows from [37] that

$$D^>p = \operatorname{div} p = \begin{cases} \begin{cases} (\rho_1)_{i,j} & (\rho_1)_{i-1,j} & \text{if } 1 < i < N \\ (\rho_1)_{i,j} & & \text{if } i = 1 \\ (\rho_1)_{i-1,j} & & \text{if } i = N \end{cases} \\ + \begin{cases} (\rho_2)_{i,j} & (\rho_2)_{i,j-1} & \text{if } 1 < j < N \\ (\rho_2)_{i,j} & & \text{if } j = 1 \\ (\rho_2)_{i,j-1} & & \text{if } j = N \end{cases} \end{cases}$$

Numerical results

Set $N = 256$, and let I be the image “boat” in the library Numerical tours [108]. We suppose that K is an operator assigning to every pixel the average of the pixels in a neighborhood of radius 8 and that $\eta \in U([0.025; 0.025])^{N^2}$. We use the original image as exact solution. For denoising and deblurring, we early stop the procedure at the iteration minimizing the mean square error (MSE), namely $k_X^k \times k^2 = N^2$, and we measure the time and the number of iterations needed to reach it. Another option for early stopping could be to consider the image with minimal structural similarity (SSIM). Numerically, in our experiments, this gives the same results. Additionally, we use the peak signal-to-noise ratio (PSNR) to compare the images. Note that the primal-dual algorithm with preconditioning is the method that needs less time and iterations among all procedures. Moreover, due to [40, Lemma 2], the condition (4.2.4) is automatically satisfied, while for the other methods we need to check it explicitly, which is computationally costly. However, (PPD) is the worst in terms of SSIM, PNSR, and MSE. We verify that all other algorithms have a superior performance in terms of reconstruction, with a small advantage for the Landweber with fixed and adaptive step-sizes, reducing the MSE of 94% with respect to the noisy image. In addition, compared to (PD), the algorithms (PDL) and (PDAL) require less iterations and time to satisfy the early stopping criterion. We believe that this is due to the fact that the extra Landweber operator improves the feasibility of the primal iterates. Visual assessment of the denoised and deblurred images are shown in Figure 4.4, which highlights the regularization properties achieved by the addition of the Landweber operator and confirms the previous conclusions.

	Iterations	Time	SSIM	PNSR	MSE
Noisy image	-	-	0.4468	21.4801	0.0071
PD	54	8.9773	0.8928	32.3614	0.0006
PD (precondition)	5	1.5515	0.8581	27.3753	0.0018
PDL	46	7.1846	0.9066	34.2174	0.0004
PDAL	31	5.4542	0.9112	34.3539	0.0004

Table 4.4: Quantitative comparison of the algorithms in terms of Structural similarity (SSIM), peak signal-to-noise ratio (PSNR), Mean square error (MSE), time, and iterations to reach the early stopping.

4.6 Conclusion and Future Work

In this chapter we studied two new iterative regularization methods for solving a linearly constrained minimization problem, based on an extra activation step reusing the data constraints. The analysis was carried out in the context of convex functions and worst-case deterministic noise. We proposed five instances of our algorithm and compared their numerical performance with state-of-the-art methods, and we observed considerable improvement in run-time.

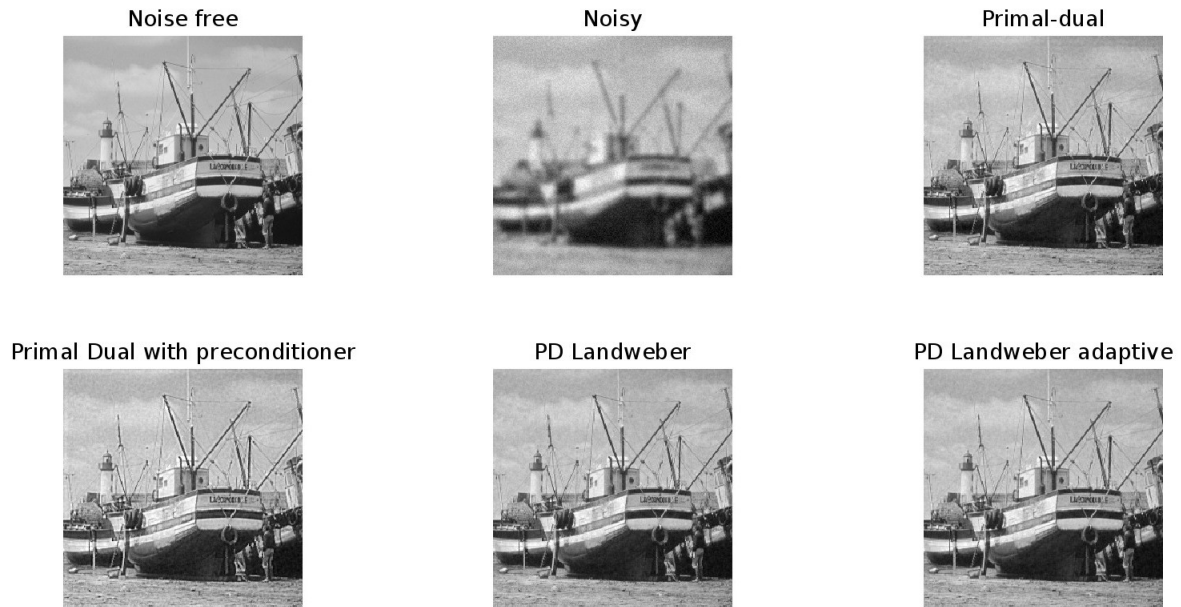


Figure 4.4: Qualitative comparison of the 4 proposed methods.

In the future, we would like to extend Theorem 4.3.1 to structured convex problems and other algorithms. Possible extensions are: 1) the study of problems including, in the objective function, a L -smooth term and a composite linear term; 2) the analysis of random updates in the dual variable (see [38]) and stochastic approximations for the gradient; 3) the theoretical study of the impact of different preconditioners; 4) the improvement of the convergence and stability rates for strongly convex objective functions.

CHAPTER 5

Implicit regularization and reparameterization

5.1 Reparameterization

The second part of this thesis is based on the recent success of overparameterization in machine learning trained with gradient descent. Instead of designing an algorithm with a fixed bias and loss, as we did in the first part of this thesis, we change our approach by reparametrizing the linear model by fixing the loss and the algorithm (which will be a gradient flow for simplicity). So, our goal is to find the implicit bias introduced by the selected optimization method and the reparameterization that has been introduced.

The map

$$\mathcal{R}^k \ni \theta = q(\cdot) \in \mathcal{R}^p$$

is the reparameterization. Typically, $k > p$, and in this view, the model becomes overparameterized. The key idea is that the sequence $\theta(t)$ obtained with gradient descent defines a sequence $\tilde{\theta}(t) = q(\theta(t))$, which corresponds to some suitable optimization procedure to find $\tilde{\theta}$ under a bias \mathcal{R} . Then, the reparameterization q defines a corresponding bias $\tilde{\mathcal{R}}$.

Then, for a given pair of models and corresponding reparameterization, the question is: what are the associated optimization procedures and corresponding biases?

In the following, we will tackle the above question, considering a simplified setting and recovering and extending a number of recent studies.

Further, we will consider a restricted set of models that are amenable to study; namely, we will assume that f is some kind of linear neural network. In particular, we will consider the following networks:

- **Two layers diagonal network:** $\theta = q(\cdot) = (\theta_1; \theta_2)$, where $\theta = (\theta_1; \theta_2) \in \mathcal{R}^p \times \mathcal{R}^p$; and $k = 2p$.
- **Deep diagonal networks:** $\theta = q(\cdot) = \theta^L$, where $\theta \in \mathcal{R}^p$, $L \geq 2$, and $k = p$.
- **Multi-neuron fully linear network of depth 2:** $\theta = q(\cdot) = Ww$, where W is a $p \times d$ matrix, $w \in \mathcal{R}^d$, $\theta = (W; w)$, and $k = (p + 1)d$.
- **Fully connected normalized linear network of depth 2:** $\theta = q(\cdot) = Ww$, where W is a $p \times d$ matrix, w is a unitary vector in \mathcal{R}^d , $\theta = (W; w)$, and $k = (p + 1)d$.

- **Weight normalized networks:** $\theta = h(\cdot)w$, where $h: \text{dom}(h) \subset \mathbb{R} \rightarrow \mathbb{R}_+$ is a given function, w is a unitary vector in \mathbb{R}^p , $\theta = (\cdot; w) \in \text{dom}(h) \subset \mathbb{R}^p$, and $k = p + 1$.

5.1.1 Optimization by reparameterization

Let $L: \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex functional, which can be interpreted as the data fit term. Let $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$, be a differentiable map, seen as the reparameterization. For some suitable initialization, consider the flow

$$\dot{\theta}(t) = -\gamma \nabla L(q(\theta(t))) \quad (5.1.1)$$

Instead of gradient descent, we consider the gradient flow, where we indicate for clarity the variable with respect to which the gradient is taken. The questions are:

1. If there exists a bias functional $R: \mathbb{R}^p \rightarrow \mathbb{R}$ such that $\theta(t) = q(\theta(t))$ converges in a suitable sense to some θ^* solving

$$\min_{\theta \in S} R(\theta); \quad S = \underset{\mathbb{R}^p}{\text{argmin}} L(\cdot) \quad (5.1.2)$$

2. If it is possible to show that the trajectory followed by $\theta(t)$ corresponds to some optimization approach to solve (5.1.2).
3. If it is possible to find an explicit expression of R in terms of q and the initialization.

Recent research has partially addressed these questions by establishing a link between gradient flow and mirror flow. The idea can be simply summarized as follows: Starting from the reparameterization $\theta = q(\cdot)$ and equation (5.1.1), a simple derivation of the composition shows that

$$\dot{\theta} = J_q(\cdot) \dot{\theta} = J_q(\cdot) \gamma \nabla L(q(\cdot)) = J_q(\cdot) J_q(\cdot)^T \gamma \nabla L(\cdot);$$

where we denote by J_q the Jacobian of q with respect to its variable. So, under the assumption that there exists a function $F: \text{dom}(F) \rightarrow \mathbb{R}$ and $G: \text{dom}(G) \rightarrow \mathbb{R}$ such that, for every $\theta \in \mathbb{R}^k$,

$$J_q(\cdot) J_q(\cdot)^T = G(q(\cdot)) \gamma^2 F(q(\cdot))^{-1}; \quad (5.1.3)$$

we get that the variable $\theta(t)$ follows the following generalized ("time-warping") mirror flow:

$$\dot{\theta}(t) = -\gamma G(\theta(t)) \gamma^2 F(\theta(t))^{-1} \nabla L(\theta(t)) \quad (5.1.4)$$

The previous computations can be understood as follows: Parameterizing $\theta = q(\cdot)$ and using gradient flow on θ is equivalent to applying "time-warping" mirror flow on θ , with some specific entropy F depending on the reparameterization q . The "time-warping" function can be interpreted as a non-linear preconditioner of the dynamical system. The flow on (5.1.4) has an implicit bias R towards a specific solution in the set of minimizers of the loss function. This bias, as we will see in the following, is related to F and G , and so to the reparameterization q .

5.2 Case $G(\cdot) = 1$: Reparameterizing gradient flow as mirror flow.

In the case when $G(\cdot) = 1$, the dynamical system (5.1.4) is indeed a mirror flow. In [4] the authors show that mirror flow [3, 16, 99] on the global variable is equivalent to

use gradient flow over the reparameterization, assuming the existence of a suitable mirror map. This map allows for characterizing the implicit bias of the optimization process for a given model. The equivalence between mirror flow on \mathbb{R}^k and gradient flow over \mathbb{R}^p was formalized in the following theorem:

Theorem 5.2.1. [4, Theorem 2] Let F be a strictly convex, continuously-differentiable function with domain in \mathbb{R}^p . Let $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$ be a reparameterization function, such that $k \leq p$, expressing parameters θ of F uniquely as $q(\phi)$ and $\text{ran}(q) = \text{dom}(F)$. Moreover, assume that

$$J_q(\phi) J_q(\phi)^T = r^2 F(\theta)^{-1};$$

for all $\theta = q(\phi)$. Then, the Mirror flow update on parameter ϕ for the convex function $F(\theta)$ and loss $L(\theta)$,

$$-\dot{\phi}(t) = r^{-2} F(\theta(t))^{-1} r L'(\theta(t));$$

coincides with the gradient flow update on parameters θ for the composite loss $L \circ q$,

$$-\dot{\theta}(t) = r L'(q(\theta(t)));$$

provided that $\theta(0) = q(\phi(0))$.

Remark 5.2.2. In the special case where $p = k$ and q is separable, meaning that the i -th component of q depends only on the corresponding component of ϕ , (5.1.3) can be reduced to the following ordinary differential equation:

$$r^2 F(q(\phi))^{-1} = \text{Diag} \{ \phi^i(\phi)^{-2} \}; \tag{5.2.1}$$

In the appendix, Example A.1.1 presents the explicit computations for the case where the reparameterization is separable and finding the mirror map is trivial, demonstrating that the mirror map can be straightforwardly determined by solving Equation (5.2.1).

In practice, for many models, finding F from (5.1.3) is not straightforward, and the resulting bias \mathcal{R} is unclear. The first technical step is to express the product of the Jacobians $J_q(\phi) J_q(\phi)^T$ in terms of the variable $\theta = q(\phi)$ (see Example A.1.2 in the appendix). The second is then to find a function $F(\theta)$ such that the inverse of its Hessian is equal to the product of the Jacobians (5.1.3).

While the previous theorem sheds light on certain aspects of how to find implicit bias, it does not tackle the existence of a mirror map. Moreover, even if such a map exists, uniqueness is not guaranteed, and finding the parameterization can be a challenging task. However, in [77] the previous results are generalized, providing conditions for the existence of the parameterization. Moreover, the convergence to a feasible point is also proved, whereas in previous works it was only assumed.

Before to present the existence result, we recall some definitions presented in Chapter 2

Definition 5.2.3. A regular parameterization $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$ is a C^1 parameterization such that $J_q(\phi)$ is of rank p for all $\phi \in \mathbb{R}^k$.

Definition 5.2.4. A C^2 parameterization $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$ is commuting if and only if for any $i, j \in [p]$, we have that,

$$r^2 q_j(\phi) r q_i(\phi) - r^2 q_i(\phi) r q_j(\phi) = 0;$$

for all $\phi \in \mathbb{R}^k$.

Definition 5.2.5. For any C^1 function $f: \mathbb{R}^k \rightarrow \mathbb{R} [f+1 g$, we denote by $\overset{t}{f}(\theta) = \overset{t}{f}(\theta)$; where $\overset{t}{f}(\theta)$ is the solution at time t (when it exists) of

$$\begin{cases} \dot{\overset{t}{f}}(\theta) = -\nabla f(\overset{t}{f}(\theta)); & t > 0; \\ \overset{t}{f}(0) = \theta \in \mathbb{R}^k: \end{cases}$$

We say $\overset{t}{f}(\theta)$ is well-defined at time t when the above differential equation has a solution at time t .

Definition 5.2.6. Given a C^2 parameterization $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$, for any $\theta \in \mathbb{R}^p$ and $t \in \mathbb{R}^p$, we define

$$\overset{t}{q}(\theta) := \begin{pmatrix} \overset{t_1}{q_1} \\ \vdots \\ \overset{t_p}{q_p} \end{pmatrix}(\theta);$$

when it is well-defined, i.e., the corresponding differential equations have a solution. For any $\theta \in \mathbb{R}^p$, we define the domain of $\overset{t}{q}(\theta)$ as:

$$U(\theta) = \{t \in \mathbb{R}^p \mid \overset{t}{q}(\theta) \text{ is well defined}\};$$

Definition 5.2.7. For any parameterization $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$ in C^2 and for any function $L: \mathbb{R}^p \rightarrow \mathbb{R} [f+1 g$ in C^1 , given any starting point $\theta \in \mathbb{R}^p$, we define the reachable set $\overset{t}{q}(\theta)$ as

$$\overset{t}{q}(\theta) = \{ \overset{t}{L} q(\theta) \mid t > 0 \};$$

Now, we present sufficient assumptions for the reparameterization, presented in [77], which ensure the existence of a mirror map and allow its implicit bias to be characterized.

Assumption 5.2.8. [77, Assumption 3.5.] Let $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$ be a parametrization. We assume that for any $\theta \in \mathbb{R}^p$ and $i \in [p]$, $\overset{t}{q_i}(\theta)$ is well-defined for $t \in T; T_+[$ such that either $\lim_{t \rightarrow T_+} k \overset{t}{q_i}(\theta) k = +1$ or $T_+ = +1$ and similarly for T^- . Also, we assume that for any $\theta \in \mathbb{R}^p$ and $i, j \in [p]$, we have that, $\overset{s}{q_i} \overset{t}{q_j}(\theta)$ is well-defined if and only if $\overset{t}{q_j} \overset{s}{q_i}(\theta)$ does, for every $t, s > 0$.

The following intermediate lemma states that the point reached by the gradient flow with any commuting parameterization is determined by the integral of the negative gradient of the loss along the trajectory.

Lemma 5.2.9. [77, Lemma 4.7.] Let $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$ be a commuting parametrization. For any initialization $\theta \in \mathbb{R}^k$, consider the gradient flow:

$$\dot{\overset{t}{f}}(\theta) = -\nabla L(q(\overset{t}{f}(\theta))) \quad \overset{t}{f}(0) = \theta;$$

Further define $\overset{t}{f}(\theta) = \int_0^t -\nabla L(q(\overset{s}{f}(\theta))) ds$: Suppose $\overset{t}{f}(\theta) \in U(\theta)$ for all $t \in [0; T)$ where $T \in \mathbb{R} [f+1 g$, then it holds that $\overset{t}{f}(\theta) = \overset{t}{q}(\overset{t}{f}(\theta))$ for all $t \in [0; T)$:

Now, we present the existence of a mirror map corresponding to a specified parameterization.

Theorem 5.2.10. [77, Lemma 4.8] Let $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$ be a commuting and regular parametrization satisfying Assumption 5.2.8. Then for any $\theta \in \mathbb{R}^k$, there exists a Legendre function $Q: \mathbb{R}^p \rightarrow \mathbb{R} [f+1 g$ such that $\nabla Q(\theta) = q(\overset{t}{q}(\theta))$ for all $\theta \in U(\theta)$. Moreover, let F be the convex conjugate of Q , then F is also a Legendre function and satisfies that $\text{int}(\text{dom}(F)) = \overset{t}{q}(\theta)$ and

$$\nabla^2 F(q(\overset{t}{q}(\theta))) = J_q(\overset{t}{q}(\theta)) J_q^T(\overset{t}{q}(\theta))^{-1};$$

for all $\theta \in U(\theta)$.

The next theorem characterizes the implicit bias for mirror flow when the loss function is the composition of a convex function with a linear operator. This implicit bias corresponds to the Bregman distance between the set of minimizers of the loss and the initialization, extending the results presented in [7, 139].

Theorem 5.2.11. [77, Lemma 4.16] Let $q: \mathbb{R}^k \rightarrow \mathbb{R}^p$ be a commuting and regular parametrization satisfying Assumption 5.2.8. Let $X \subseteq \mathbb{R}^d$, let $y \in \mathbb{R}^d$, let $\ell: \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable with a locally Lipschitz gradient. Assume that $L = \ell \circ X$ and $S = \{f \in \mathbb{R}^p \mid X^{-1}f = y\}$. Consider the gradient flow given by (5.1.1). Then, there exists function F given Theorem 5.2.10 and a solution $\gamma \in \mathbb{R}^p$ such that $\gamma = \lim_{t \rightarrow \infty} \gamma(t)$ and

$$\gamma \in \underset{z \in S}{\operatorname{argmin}} F(z) \quad \text{hr} F(\gamma) \mid z \in S = \underset{z \in S}{\operatorname{argmin}} D_F(z; \gamma)$$

Although [77] encompasses many of the existing parameterizations, there are also interesting reparameterizations for which a function F such that (5.1.3) holds does not exist if $G(\cdot) = 1$, making it necessary to consider $G \leq 1$. In this case, the dynamic cannot be expressed as vanilla mirror flow, but the flow can still be expressed as a mirror flow multiplied by a positive scalar function G , as it was studied in [7] for the case of a multi-neuron fully connected linear network.

5.3 Case $G(\cdot) \leq 1$: Time warping Mirror Flow

If $G(\cdot) \leq 1$, the dynamical system (5.1.4) is a time-warped mirror flow. In this case, the dynamics cannot be described by vanilla mirror flow; however, an entropy function can still be found, as was studied in [7] for multi-neuron fully connected linear networks. This flow can be understood as mirror descent with a non-linear time-warping, which gives degrees of freedom with respect to [77] on how to choose F . This map G allows for characterizing the implicit bias of the optimization process for a given model.

Consider, for instance, the following two leading examples along the this thesis.

Example 5.3.1 (Multi-neuron fully connected linear network of depth 2). In linear networks of depth 2, the reparameterization is given by $\gamma = q(\theta)$ with $\theta = (W; w) \in \mathbb{R}^{p \times m} \times \mathbb{R}^m$ and $q(\theta) = Ww$. The gradient flow on the reparameterization is then given by

$$\begin{aligned} \dot{W}(t) &= -\nabla_W L(W(t)w(t)) = -\nabla L(\gamma(t))w^\top(t); \\ \dot{w}(t) &= -\nabla_w L(W(t)w(t)) = -W^\top(t)\nabla L(\gamma(t)); \end{aligned} \quad (5.3.1)$$

where we suppose, for simplicity in the computations, that the initialization satisfies $W^\top(0)W(0) = w(0)w^\top(0)$. In this case, by Lemma 2.4.5, it is possible to prove that there is no function F such that (5.1.3) holds with $G(\cdot) = 1$. On the other hand, in the original variable θ , we obtain the following dynamical system:

$$\dot{\theta} = -G(\theta)\nabla L(\gamma(\theta)); \quad (5.3.2)$$

where

$$F(\theta) = (2/3)k \|\theta\|^{3/2} \quad \text{and} \quad G(\theta) = k \|\theta\|^{1/2}.$$

For the proofs of the statements above, see Section A.1.3 in the appendix. Note that, due to the presence of the function G , (5.3.2) is not a mirror flow but a generalization of it.

Example 5.3.2 (Standard weight normalization). Consider the reparameterization given by $w = q(\tilde{w})$ with $\tilde{w} = (w; \lambda) \in \mathbb{R}^p \times \mathbb{R}$ and $q(\tilde{w}) = \frac{w}{k\tilde{w}k}$. The gradient flow on the reparameterization is then given by

$$\begin{aligned} \dot{w}(t) &= -\text{r}_w L(w(t)) = \frac{\dot{w}(t)}{k\tilde{w}(t)k} \text{Id} - \frac{w(t)w^\top(t)}{k\tilde{w}(t)k^2} \text{r}_w L(w(t)) \\ \dot{\lambda}(t) &= -\text{r}_\lambda L(w(t)) = \frac{w(t)}{k\tilde{w}(t)k} j \text{r}_w L(w(t)) : \end{aligned}$$

with initialization $(w_0; \lambda_0) \in \mathbb{R}^p \times \mathbb{R}$ such that $k\tilde{w}_0k = 1$ and $\lambda_0 > 0$. In this case, by Lemma 2.4.5, it is possible to prove that there is no function F such that (5.1.3) holds with $G(\tilde{w}) = 1$. On the other hand, in the original variable w , we obtain the following dynamical system:

$$\text{r}^2 F(w(t)) - \dot{w}(t) = G(\tilde{w}(t)) \text{r}_w L(w(t));$$

where $G(\tilde{w}) = k \exp\left(\frac{k \cdot k^2}{2}\right)$ and F is a function that cannot be expressed explicitly. For the proofs of the statements above, see Example 6.3.6 with $L = 1$.

The above examples were presented in [7] and [89, 118, 139]. However, the convergence of the trajectory is assumed but not proven (see [7, 89, 139]). In the next chapter, we study the well-posedness, convergence of the iterates, and convergence in value for time-warped mirror flow applied to convex losses. Additionally, when the loss function consists of a strictly convex function composed with a linear operator, we give an explicit expression of its implicit bias.

CHAPTER 6

Learning from data via overparameterization

Abstract

The goal of machine learning is to achieve a good prediction exploiting training data and some a-priori information about the model. The most common methods to achieve the last objective are explicit and implicit regularization. In the first technique, a regularizer is explicitly introduced to find, among all the solutions, a good generalizing one. The second technique, i.e. implicit regularization, is based on the inductive bias intrinsically induced by the specific method used to optimize the parameters involved.

Recently, the success of learning is related to re- and over-parameterization, that are widely used - for instance - in neural networks applications and the optimization method used. However, there is still an open question of how to find systematically what is the inductive bias hidden behind the model for a particular optimization scheme. The goal of this chapter is taking a step in this direction, studying extensively many reparameterizations used in the state of the art and providing a common structure to analyze the problem in a unified way. We show that gradient descent on the empirical loss for many reparameterizations is equivalent, in the original problem, to a generalization of mirror descent. The mirror function depends on the reparameterization and introduces an inductive bias, which plays the role of the regularizer. Our theoretical results provide asymptotic behavior and convergence in the simplified setting of linear models.

Keywords. Overparameterization, Implicit Regularization, Time-warping Mirror Flow, Fully connected normalized linear networks, Weight normalization.

AMS Mathematics Subject Classification (2020): 34A55, 90C25, 65K10.

In this chapter, we study a unified framework observed in various reparameterizations presented in the state of the art, called time-warp mirror flow given in equation (6.2.1). We provide a complete analysis consisting of several steps. Firstly, in Section 6.2, we establish conditions for the well-posedness of the dynamical system. Then, we show that for any convex function, the sequence converges to a stationary point that minimizes the loss function and avoids the extra stationary points that are produced by the reparameterization. For the specific case of a function composed with a linear operator, an implicit bias is provided. Next, these results will be used in Section 6.3, where we look at weight normalization techniques for functions F and G that only depend on the norm of \cdot . These functions are also called radial functions. Furthermore, we provide a criterion for determining a suitable weight normalization parameterization for a given function that depends only on the norm. Finally, we explore the flexibility of our formulation by applying the previous results to different examples related to weight normalization. Finally, in Section 6.4, we conclude this chapter with some remarks and future works.

Compared to the previous chapter, the main difference is that, for the given reparameterization, we assume the existence of F and G . Unlike in Theorem 5.2.10, where conditions for the existence of such reparameterizations are provided and analyzed, our purpose here is not to study these conditions. Instead of focusing on reparameterization, we examine the convergence properties of the time-warping mirror flow and its implicit bias, as outlined in Theorem 5.2.11.

6.1 Introduction

Classic algorithm design in machine learning (ML) is based on fitting some chosen model to data while including some bias reflecting prior knowledge on the problem [64]. In fact, recent studies have shown that large overparameterized models can achieve excellent learning performance, even without enforcing any explicit bias [58, 59, 77, 88, 93, 102, 131, 135]. A possible explanation is that the bias is implicitly incorporated through the selection of the model and the optimization procedure used. However, uncovering such a bias is typically a complex challenge [5, 7]. An idea is to view the overparameterization as a reparameterization of some simpler, original model. This raises the question if the optimization of the over- or re-parameterized model has a clearer interpretation in the original model. A natural starting point for exploring the concept of overparameterization is to consider linear networks, which can be seen as overparameterizations of linear models. Indeed, recent results have started considering so-called diagonal networks and have shown that gradient flow for these models corresponds to a mirror flow for the original model [4, 77]. In this context, the bias is linked to a sparsity prior imposed on the linear model [43, 135, 147]. As observed in [7], a more general form of mirror flow, including a time-warping factor, allows for a much more general treatment. In this paper, we develop this latter line of work. First, we provide a general analysis of time-warped mirror flow. In particular, we study its convergence both in the value of the loss function and in the iterations. Additionally, we provide rates of convergence and characterize the implicit bias of the limit point. Second, we discuss in detail the case of one-hidden-layer linear networks and study the effect of different types of weight normalization. Our results show that for a given radial function, i.e., a function that only depends on the norm, we can find a suitable parameterization such that the implicit bias is given by the radial function. We note that weight normalization has been previously considered in [44, 89, 118, 139], but using the framework of time-warping mirror flow, we explore more general reparameterizations. To the best of our knowledge, the analysis of general linear one-hidden-layer neural networks is new.

6.1.1 Related work

In this subsection, we will briefly analyze the literature most related to this chapter.

Reparameterizing gradient flow as mirror flow: The first paper relating gradient flow on the reparametrization to mirror flow [3] on the original variable is [4]. The relation between the two flows holds under quite restrictive assumptions on the reparametrization, which ensure the existence of a mirror map. This map allows for characterizing the implicit bias of the optimization process for a given reparametrization. The existence of the mirror map and its analytic form given a reparametrization is not addressed in [4], neither the question of establishing whether different reparametrizations lead to the same implicit bias. A further step has been taken in [77]. This is the only paper where existence and convergence of a solution of the mirror flow is discussed. In addition, conditions on a parametrization are given for the existence of the mirror map. Although [77] encompasses many of the existing linear reparameterizations, there are still examples that cannot be cast under its setting, such as matrix times vector overparametrizations and weight normalization. In these two examples, the dynamic in the original variable cannot be expressed as a mirror flow. The paper [7] proposed a generalized version of mirror flow, that is the one analyzed in this paper, for the case of matrix times vector overparametrization. Their interpretation proposes to relate gradient descent on the reparametrization to mirror descent with a non-linear time-warping in the original variable. This gives more degrees of freedom with respect to [77] and allows to cover in principle a more general class of reparameterizations.

Diagonal linear neural network using least square: Deriving an Implicit bias for any general network is still not well understood. However, this problem has been addressed for simpler models. In [43, 135, 147], it is shown that for the simplest model, a diagonal linear neural network with depth L (where each component is raised to the power of L), vanilla gradient flow approximates the minimal ℓ_1 -norm solution. However, to achieve a good approximation, this method requires an initialization close to zero, which in practice makes the convergence slower. Similar results for stochastic gradient flow were obtained in [2, 106], where it is proved that slower convergence implies better bias. The main advantage obtained with respect to the deterministic approach is that, with the same initialization, the sequence generated by stochastic gradient flow achieves greater sparsity than the one obtained by vanilla gradient flow.

Implicit bias of least square: The bias of different first-order optimization algorithms and the effect of the step-size on the optimization process have been studied in [58, 138], [76], and [93], respectively. More general models have been addressed in [7, 49, 130], where the case of a fully connected linear neural network and a two-Layer Single (Leaky) ReLU neuron is studied. In [44], inspired by the weight normalization technique proposed in [118], which involves decoupling a vector into its norm and its direction, a diagonal linear network is combined with weight normalization, providing a robust implicit bias for large scale initialization. While in [139] is studied the inductive bias for the squared loss function using weight normalization in the discrete and continuous settings. Moreover, the authors also proved the convergence of the sequence generated by the algorithm to the minimal norm solution among all the solutions of least squares. In this chapter, we generalize this result by studying the effect of the choice of different reparameterizations of the norm, including classical and exponential weight normalization, among others. Furthermore, we give conditions to find a suitable weight normalization reparameterization for functions that depend on the norm, including the analysis of its convergence and the characterization of its implicit bias.

Implicit bias of logistic loss: Similar results have been obtained by applying gradient descent to the logistic loss function, which is commonly used in classification problems. In [67, 68, 92, 125], it was proved that using gradient descent, the linear classifier converges to the direction of L_2 maximum margin. This result was later extended for SGD in [94]. For the case of Diagonal linear neural networks, the same result was obtained in [88] when the initialization tends to infinity. The implicit bias of the weight normalization for the logistic loss was also studied in [89]. The inductive bias of a 2-layer infinitely wide ReLU neural network (i.e., a neural network with a large number of hidden units), it was studied in [42]. In [79] the implicit regularization of the gradient descent algorithm using the exponential loss function is studied. This work encompasses different types of parameterization of homogeneous neural networks, including fully connected and convolutional neural networks with ReLU or LeakyReLU activations, and proves the convergence to a minimal norm solution subject to margin constraints. The same technique was used in [59], for a full-width linear convolutional network applying gradient descent.

Deep Matrix Factorization: Classic theory also covers how to induce a bias using matrices. A common approach is to use matrix factorization, which is a powerful tool for simplifying its representation by finding the underlying structure of a matrix. A well-known method to induce low rank for semi-definite positive symmetric matrices is factorization, whereby a matrix is expressed as the product of a matrix by its transpose. This technique is useful to find a low-rank solution since, with the dimension of the parameterization, we can impose an additional rank constraint. In [60, 75] it was demonstrated that the factorization reaches the minimal nuclear norm for a sufficiently small initialization, assuming that the sensing matrices are commutative and semi-positive definite. The generalization for the multiplication of more than two matrices is in [5], which shows that gradient flow can act as a preconditioner that prefers the direction of the larger eigenvalues. The algorithmic bias of using mirror descent was studied in [134, 137], and the generalization for a tensor formulation was studied in [145]. The case where the sensing matrices satisfy the restricted isometry property was studied in [131, 136, 144]. A simplified version of a neural network is the linear neural network, i.e., when all the activation functions are the identity. This was studied in [6, 55, 56, 80].

6.2 Global existence and asymptotic analysis

In the setting of Section 5.1.1, given a differentiable convex loss function $L : \mathbb{R}^p \rightarrow \mathbb{R}$, $\nabla L(\cdot)$, we consider the trajectory $\theta(t)$ generated by the gradient flow dynamics on a reparameterization $\theta = q(\cdot)$. We analyze the case in which the gradient flow on the reparameterization may not be expressed as a mirror flow on the original variable θ ; but it can still be written as a dynamical system involving the inverse of the Hessian of F and a positive function G , as in (5.3.2):

$$\begin{cases} \dot{\theta}^2 F(\theta(t))^{-1} \dot{\theta} = G(\theta(t)) \nabla L(\theta(t)); & t > 0 \\ \theta(0) = \theta_0 \in \mathbb{R}^p; \end{cases} \quad (6.2.1)$$

The positive function G , already introduced in [7], can be interpreted as a non-linear “time-warping”. In the case when $G(\cdot) \equiv 1$, the dynamical system (6.2.1) is indeed a mirror flow.

We present a comprehensive framework to study the existence and uniqueness of the trajectory of (6.2.1) and its minimization properties, such as the decrease of the loss’

values along it and the corresponding convergence rates (see Theorems 6.2.4 and 6.2.8). Moreover, we show the convergence of the trajectory to $S = \operatorname{argmin} L$ and we explicitly compute its implicit bias towards a specific solution in S (see Theorem 6.2.8). As our analysis shows, with respect to the classical mirror flow setting, the presence of the function G , when supposed to be positive, does not affect the fact that the loss is a non-increasing function in time, neither the fact that the trajectory converges to the set of minimizers of the loss. On the other hand, it plays a role in the selection of the specific solution to which the trajectory is converging. The techniques used to study the mirror flows used in this chapter are presented in [3]. However, the main difference between their work and ours is that they focus on mirror flows in the context of constrained minimization problems, while in our work, the constraints are encoded within the loss function itself.

Differently from previous works [6, 7, 77], here we focus on the well-posedness and convergence over time of the dynamical system described by Equation (6.2.1). Instead of assuming the convergence of (t) to a minimizer of the loss function as an hypothesis, as in [6, 7], we prove this fact. Additionally, we derive the implicit bias induced by reparametrization for a general convex loss function L , while in previous work only the least squares loss function has been considered. For the least squares loss, we recover that the implicit bias R is the Bregman divergence of the entropy F as in [4, 77] (see Corollary 6.2.9).

6.2.1 Well-posedness

In this section, using techniques developed in [3], we show that the system (6.2.1) is well-posed by proving the existence of a unique solution that is defined in the entire time-interval $[0; +\infty)$. To do so, we first list the assumptions that we require on the loss function L (see A1), the initialization $(0) = \theta_0$ (see A3), the mirror map F and the time-warping function G (see A2, A4, and A5). These hypotheses are also useful to prove, in Theorem 6.2.8, the convergence of the global variable (t) to a specific minimizer of the loss function.

Assumption 6.2.1.

- A1 $L: \mathbb{R}^p \rightarrow \mathbb{R}$ is a convex and differentiable function such that rL is locally Lipschitz continuous on \mathbb{R}^p and $S = \operatorname{argmin}_{\mathbb{R}^p} L(\cdot) \neq \emptyset$.
- A2 There exists an open set $U \subset \mathbb{R}^p$ such that $S \setminus U \neq \emptyset$; the function $F: \operatorname{dom}(F) \rightarrow \mathbb{R}$ is strictly convex and twice-differentiable in U , and the mapping $\theta \in U \mapsto [r^2 F(\theta)]$ is invertible with locally Lipschitz inverse. The function $G: U \rightarrow \mathbb{R}_+$ is locally Lipschitz continuous. Moreover, one of the following two conditions is satisfied:
- $U := \mathbb{R}^p \setminus \{0\}$;
 - F is Legendre on $U = \operatorname{int}(\operatorname{dom}(f))$ (see Definition 2.4.7 in the Appendix).
- A3 The initialization satisfies $\theta_0 \in U$. If $U = \mathbb{R}^p \setminus \{0\}$, we additionally suppose that $L(\theta_0) < L(0)$.

A4 Set

$$r = \frac{L(0) - L(\theta_0)}{2krL(0)k} > 0:$$

There exists a constant $C > 0$ such that, for all $z \in B(0; r)$, $G(z) \leq C$.

A5 There exists $z \in S \setminus U$ such that the function $D_F(z; \cdot)$ is coercive (see (2.4.1) for the expression of the Bregman divergence D_F).

Remark 6.2.2. Assumption A3 implies in particular that $0 \notin S$.

Remark 6.2.3. Assumption 6.2.1 is key to establish the well-posedness and the convergence analysis of the flow given by (6.2.1). Assumption A1 requires convexity of the loss function L , which is important to derive global convergence results. Assumptions A1 and A2 are standard in the classical theory of local existence and uniqueness of solutions of ordinary differential equations [11, Subsection 1.2]. Similar assumptions have been considered in a related but different setting in [3, Section 3.2]. As we will see Assumptions A3 and A4 ensure a suitable initialization for the flow, which avoids convergence of the solution towards undesired stationary points of (6.2.1) which are not minimizers. Condition A3 is the generalization of to the one presented in [139, Lemma 2.3] for least squares problems. Assumption A5 is crucial to establish boundedness of the trajectory, which is an intermediate step for proving convergence. Indeed, Theorem 6.2.4 and Theorem 6.2.8 would still hold replacing Assumption A5 by boundedness of the trajectory. Observe that the coercivity hypothesis could be replaced for instance by the boundedness of the level set of the loss function.

The following theorem establishes the existence and uniqueness of the global solution of the dynamical system (6.2.1).

Theorem 6.2.4. Under assumptions A1-A5 the dynamical system (6.2.1) has a unique solution defined on $[0; +\infty)$. Moreover, for every $z \in S$, both the Bregman divergence $D_F(z; \cdot(t))$ and $L(\cdot(t))$ are non-increasing functions of $t \geq 0$.

Proof. Denote by T_M the maximal time for the existence of a solution of (6.2.1); namely,

$$T_M := \sup \{ T > 0 \mid \exists \text{ a } C^1 \text{ function } \gamma : [0; T) \rightarrow U \text{ solution of (6.2.1) } \}.$$

First we show that $T_M > 0$. Since $r^2 F$ is invertible on U , the dynamical system in (6.2.1) is equivalent to the following:

$$\begin{cases} \dot{\gamma}(t) = G(\gamma(t)) - r^2 F(\gamma(t))^{-1} r L(\gamma(t)); & 0 < t < T_M \\ \gamma(0) = z \in U; \end{cases} \quad (6.2.2)$$

Since $G(\cdot)$, $r^2 F(\cdot)^{-1}$, and $r L(\cdot)$ are locally Lipschitz on U , they are continuous and bounded on bounded subsets and their product is also locally Lipschitz. Since $z \in U$, it follows from [11, Theorem 1.18] that there exists a unique solution of the initial value problem (6.2.2), and so to (6.2.1), in some time-interval $[0; T_1)$ with $T_1 > 0$. Then, $T_M > 0$.

Next, we prove that $T_M = +\infty$; namely, that the solution of the dynamical system (6.2.1) is globally defined. Now we proceed to prove the claim that $T_M = +\infty$ separately for the two cases $U := \mathbb{R}^p \setminus \{0\}$ and F Legendre on U .

In order to prove Theorem 6.2.4, we establish an auxiliary result. It states that both the Bregman distance and the loss function, along the trajectory of the flow in $[0; T_M)$, do not increase.

Lemma 6.2.5. Let $z \in S$ and suppose that assumptions **A1** and **A2** hold. Then, both the Bregman divergence $D_F(z; \cdot(t))$ and the functional value $L(\cdot(t))$ are non-increasing functions for $t \in [0; T_M)$.

Proof. To prove that $t \mapsto D_F(z; \cdot(t))$ is decreasing, it is sufficient to note that, for every $t \in [0; T_M)$,

$$\begin{aligned} \frac{d}{dt} D_F(z; \cdot(t)) &= \frac{D}{D} r F(\cdot(t)) j \cdot(t) \quad \frac{E}{E} \frac{D}{D} r^2 F(\cdot(t)) - (t) j z \cdot(t) \\ &+ \frac{D}{D} r F(\cdot(t)) j \cdot(t) \\ &= \frac{D}{D} r^2 F(\cdot(t)) - (t) j z \cdot(t) \quad \frac{E}{E} \\ &= G(\cdot(t)) h r L(\cdot(t)) j z \cdot(t) i \\ &\quad G(\cdot(t)) (L(z) - L(\cdot(t))) \\ &\leq 0; \end{aligned} \tag{6.2.3}$$

where we used that \cdot is a solution of (6.2.1), the gradient inequality for the convex function L and the fact that, for every $\cdot \in U$, $L(\cdot) \leq L(z)$ and $G(\cdot) \geq 0$. On the other hand, for $t \in [0; T_M)$, $r^2 F(\cdot(t))$ is invertible (and the inverse is positive definite). Then,

$$\begin{aligned} \frac{d}{dt} L(\cdot(t)) &= \frac{D}{D} r L(\cdot(t)) j \cdot(t) \quad \frac{E}{E} \\ &= \frac{D}{D} G(\cdot(t)) r L(\cdot(t)) j r^2 F(\cdot(t))^{-1} r L(\cdot(t)) \quad \frac{E}{E} \\ &\leq 0; \end{aligned}$$

which completes the proof. \square

Case $U = \mathbb{R}^n \setminus \{0\}$: In this case, we start by providing a strictly-positive lower bound for the norm of \cdot along the trajectory, which implies that $\cdot \in ([0; T_M]) \cap U$, then, based on this, we prove that $T_M = +\infty$.

To finish this case, we need an intermediate lemma.

Lemma 6.2.6. Suppose that assumptions **A1**, **A2**, **A3** and **A4** hold with $U := \mathbb{R}^n \setminus \{0\}$. Then, for every $t \in [0; T_M)$,

$$\|\cdot(t)\| > r \quad \text{and} \quad G(\cdot(t)) \geq C \tag{6.2.4}$$

Proof. Let $\cdot \in U$ be such that $\|\cdot\| = r$. Then, the gradient inequality for the convex function L at the point 0, the Cauchy-Schwartz inequality, and Assumptions **A3** and **A4** imply

$$L(\cdot) \leq L(0) + h r L(0) j \cdot i \leq L(0) + r k r L(0) k = \frac{L(0) + L(\cdot)}{2} > L(\cdot): \tag{6.2.5}$$

Next, by contradiction, suppose that there exists a $\hat{t} \in [0; T_M)$ such that $\|\cdot(\hat{t})\| = r$. Then, by (6.2.5), we have that $L(\cdot) < L(\cdot(\hat{t}))$. This is a contradiction since Lemma 6.2.5 $t \mapsto L(\cdot(t))$ is non-increasing on $[0; T_M)$. Then $\|\cdot(t)\| > r$ for every $t \in [0; T_M)$. The second part of the statement in (6.2.4) is a consequence of the condition on G assumed in **A4**. \square

Lemma 6.2.6 yields $\|k(t)\| > r$ for every $t \in [0; T_M)$. Lemma 6.2.5 implies that $D_F(z; \cdot)$ is non-increasing for every $z \in S$ and by Assumption A5, $D_F(z; \cdot)$ is coercive for some $z \in S$. Then, $([0; T_M))$ is bounded. Therefore $([0; T_M))$ is contained in a compact subset of U . We then derive from [11, Proposition 1.57] that $T_M = +\infty$.

Case F is Legendre on U : We show that $T_M = +\infty$.

Let K be the closure of $([0; T_M))$. Assumption A5 implies that $([0; T_M))$ is bounded; therefore, K is compact. Recall that the trajectory is continuous and that $([0; T_M))$ is a subset of the open set U . If $K \subset U$, by applying [11, Proposition 1.57], we conclude that $T_M = +\infty$. If K is not contained in U , then there exists a point $z \in \bar{U} \cap U^c$ for which there exists a sequence $(t_j)_{j=1}^{+\infty} \subset [0; T_M)$ with $t_j \rightarrow T_M$ and $\|k(t_j)\| \rightarrow +\infty$ for $j \rightarrow +\infty$. Integrating in time equation (6.2.1), we get that, for every $j \in \mathbb{N}$,

$$\begin{aligned} \|k(t_j) - k(0)\| &= \int_0^{t_j} \|G(t) r L(t)\| dt \\ &\leq \int_0^{t_j} \|G(t)\| \|k(t)\| dt \\ &\leq T_M \max_{z \in K} \|G(z)\| \max_{z \in K} \|k(z)\|. \end{aligned} \quad (6.2.6)$$

The maxima are well-defined due to the fact that the functions G and $\|k\|$ are locally Lipschitz and so continuous, implying that they are bounded on the compact set K (see, for instance, [91, Theorem 27.4]). For $j \rightarrow +\infty$, we have that $\|k(t_j)\| \rightarrow +\infty$ and therefore, since F is a Legendre function, $\|k(t_j)\| \rightarrow +\infty$. Thus the left-hand side of (6.2.6) tends to infinity as well, and this implies that $T_M = +\infty$.

Uniqueness: Since the right hand side of (6.2.2) is locally Lipschitz at every point, local uniqueness holds for every initial state, and therefore globally. \square

Remark 6.2.7. To the best of our knowledge, in previous works, e.g. in [7], the existence and uniqueness of the solution of (6.2.1) is not proved. In contrast, we provide some sufficient conditions to ensure the well-posedness of the system. This is fundamental to guarantee that the implicit bias characterization is meaningful. Indeed, it is well known that the trajectory of a differential equation exists globally if and only if it is locally bounded. So, if the solution is defined on a bounded interval and not on $[0; +\infty)$, there is a finite time for which the norm of the trajectory diverges. In this case, the limit of the trajectory for that finite time cannot have any interesting properties, since clearly cannot converge to some specific minimizer of L .

6.2.2 Minimization properties and implicit bias

The next theorem illustrates the minimization properties of the solution of the dynamical system (6.2.1) for $t \rightarrow +\infty$. It states the convergence rates for the loss function and the convergence of the trajectory towards a specific minimizer of the loss, characterized by an implicit bias.

Theorem 6.2.8. Assume that A1-A5 hold, then $L(k(t)) \rightarrow L^* := \min L$ for $t \rightarrow +\infty$. In addition, $\|L(k(t)) - L^*\| \leq C e^{-\lambda t}$, and, for every $z \in S \setminus U$ and for every $t > 0$,

$$\|L(k(t)) - L^*\| \leq \frac{D_F(z; 0)}{Ct}.$$

Moreover, there exists $\gamma \in S \setminus U$ satisfying

$$\gamma \in \operatorname{argmin}_{z \in S} \int_0^{\gamma} F(z) \, dz = \int_0^{\gamma} F(z) \, dz + \int_0^{\gamma} G(t) \, dt = \int_0^{\gamma} L(t) \, dt; \quad (6.2.7)$$

such that

$$\lim_{t \rightarrow \gamma} D_F(z; t) = \gamma;$$

Proof. Let $z \in S$. By the Fundamental Theorem of Calculus, inequality (6.2.3) and Lemma 6.2.6, we have that, for all $t \geq 0$,

$$\begin{aligned} D_F(z; t) - D_F(z; 0) &= \int_0^t \frac{d}{ds} D_F(z; s) \, ds \\ &= \int_0^t G(s)(L(s) - L) \, ds \\ &\leq C \int_0^t (L(s) - L) \, ds \leq 0; \end{aligned} \quad (6.2.8)$$

Since $D_F(z; t)$ is non-negative, from the above inequality we derive that, for every $t \geq 0$,

$$\int_0^t (L(s) - L) \, ds \leq \frac{D_F(z; 0)}{C};$$

In particular, $t \in [L(t) - L]^{-1}([0; +\infty))$. Moreover, since $t \mapsto L(t)$ is non-increasing by Lemma 6.2.5, we have that $(L(t) - L) \leq \int_0^t (L(s) - L) \, ds$ and $D_F(z; 0) = C$ for every $t \geq 0$. From the latter we obtain the rate stated in the theorem.

Next we prove the convergence of the trajectory $(t)_{t \geq 0}$. Since $z \in S$ by Assumption A5, then $D_F(z; t)$ is non-negative and decreasing by Lemma 6.2.5, and therefore there exists $\lim_{t \rightarrow \infty} D_F(z; t) \in \mathbb{R}$. Since $D_F(z; \cdot)$ is coercive by assumption A5, we also have that (t) is bounded. Then there exist a sequence $(t_k)_{k \geq 0}$ and a point γ such that, for $k \rightarrow \infty$, $t_k \rightarrow \gamma$ and $(t_k) \rightarrow \gamma$. From continuity of L and by $[L(t) - L] \in L^1([0; +\infty))$, we deduce that

$$L(\gamma) = \lim_{k \rightarrow \infty} L(t_k) = \lim_{t \rightarrow \gamma} L(t) = L;$$

i.e., that $\gamma \in S$. We are going to show that indeed $\lim_{t \rightarrow \infty} (t) = \gamma$. Suppose that there exists another sequence $(t_n)_{n=1}^{\infty}$ such that $t_n \rightarrow \infty$ and $\lim_{n \rightarrow \infty} (t_n) = \beta$. With the same reasoning as above, we know that β belongs to S . Continuity of $D_F(\cdot; \cdot)$ and the existence of $\lim_{t \rightarrow \infty} D_F(\gamma; t)$ yields

$$\begin{aligned} D_F(\gamma; \beta) &= \lim_{n \rightarrow \infty} D_F(\gamma; (t_n)) = \lim_{t \rightarrow \infty} D_F(\gamma; t) \\ &= \lim_{k \rightarrow \infty} D_F(\gamma; (t_k)) = D_F(\gamma; \gamma) = 0; \end{aligned}$$

Since F is strictly convex, we get $\beta = \gamma$. Finally, we get that $\lim_{t \rightarrow \infty} (t) = \gamma$. To conclude, note that trivially

$$\gamma \in \operatorname{argmin}_{z \in S} D_F(z; \gamma);$$

where, by (6.2.8) and (6.2.3),

$$\begin{aligned} D_F(z; \gamma) &= \lim_{t \rightarrow \gamma} D_F(z; (t)) \\ &= \lim_{t \rightarrow \gamma} D_F(z; 0) + \int_0^{\gamma} G(s) \text{hr} L((s)) j z \quad (s) ds \\ &= F(z) - F(0) - \text{hr} F(0) j z \quad 0 + \int_0^{\gamma} G(s) \text{hr} L((s)) j z \quad (s) ds: \end{aligned}$$

Discarding the terms that are constant with respect to z , we obtain the claim. \square

The next corollary characterizes the implicit bias for time-warping mirror flow when the loss function is the composition of a convex function with a linear operator. This implicit bias corresponds to the Bregman distance between the set of minimizers of the loss and the initialization, extending the results presented in [7, 139].

Corollary 6.2.9. Suppose that assumptions A1-A5 hold. Let $X \subseteq \mathbb{R}^d$, let $y \in \mathbb{R}^d$, let $\cdot: \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable with a locally Lipschitz gradient. Assume that $L = \cdot \circ X$ and $S = \{z \in \mathbb{R}^p \mid Xz = y\}$. Then, there exists a solution $\gamma = \lim_{t \rightarrow \gamma} (t)$ satisfying

$$\gamma \in \underset{z \in S}{\operatorname{argmin}} F(z) - \text{hr} F(0) j z = \underset{z \in S}{\operatorname{argmin}} D_F(z; 0):$$

Proof. It follows from Theorem 6.2.8 that there exists $\gamma \in S \setminus U$ such that $\gamma = \lim_{t \rightarrow \gamma} (t)$ and

$$\gamma \in \underset{z \in S}{\operatorname{argmin}} F(z) - \text{hr} F(0) j z + \int_0^{\gamma} G(t) \text{hr} L((t)) j z \quad (t) dt :$$

To conclude the assertion, it is sufficient to prove that the integral is constant with respect to z . By the chain rule, for every $z \in S \setminus U$,

$$\begin{aligned} \int_0^{\gamma} G(t) \text{hr} L((t)) j z \quad (t) dt &= \int_0^{\gamma} G(t) \text{hr} \langle X^T \cdot \rangle (X(t)) j z \quad (t) dt \\ &= \int_0^{\gamma} G(t) \text{hr} \langle X(t) \rangle j y \quad X(t) dt \end{aligned}$$

which is negative, bounded below by $-D_F(\gamma; 0)$, and it is independent of $z \in S$. Consequently, from (6.2.7), we obtain that the limit of the trajectory satisfies

$$\gamma \in \underset{z \in S}{\operatorname{argmin}} F(z) - \text{hr} F(0) j z = \underset{z \in S}{\operatorname{argmin}} D_F(z; 0):$$

\square

Remark 6.2.10. (i) The point γ is the Bregman projection, with respect to the entropy F , of 0 on S . Since S is a convex set and F is a strictly-convex function, the minimizer is unique. Therefore,

$$\underset{z \in S}{\operatorname{argmin}} D_F(z; 0) = \gamma: \tag{6.2.9}$$

(ii) The limit $\gamma \in U$, the domain of the entropy F , which therefore can also implicitly enforce some constraints. For instance, the parametrization $z = X^{-1}(y - X\gamma)$ implicitly enforces the positivity of the solution: in this case the domain of F is the positive orthant, see the case $L = 2$ in equation (A.1.1).

(iii) Even if the term $\int_0^{R+1} G(\tau) h r L(\tau) j z(\tau) d\tau$ is not constant with respect z , it is upper bounded by the quantity $\int_0^{R+1} G(\tau) (L(\tau) - L(\tau_0)) d\tau < 0$, that measures the cumulative functional value decrease.

6.2.3 Application to fully connected linear networks

We consider the setting of Example 5.3.1, where $\theta = (W; w) \in \mathbb{R}^p \times \mathbb{R}^m$ and $\mathcal{L} = Ww$ for a loss function L satisfying Assumption A1. When L is the square loss, this setting has been considered in [7] and is a specific instance of a more general class of overparametrizations analyzed in [6].

Suppose that the initialization satisfies $W_0^T W_0 = w_0 w_0^T$: Then, for every $t > 0$, the trajectory $\theta(t)$ corresponding to the gradient flow trajectory in the reparameterization $(W(t); w(t))$ considered in (5.3.1) satisfies the differential equation

$$r^2 \frac{2}{3} k(t) k^{3=2} \dot{\theta}(t) = -k(t) k^{1=2} r L(\theta(t)); \tag{6.2.10}$$

namely (6.2.1) with $F(\theta) = (2=3)k k^{3=2}$ and $G(\theta) = k k^{1=2}$ (for more details, see Section A.1, Example A.1.3 in the appendix). The Hessian and its inverse are given by

$$r^2 F(\theta) = k k^{1=2} \text{Id} - \frac{2}{2k k^2} \theta \theta^T \quad \text{and} \quad r^2 F(\theta)^{-1} = k k^{1=2} \text{Id} + \frac{2}{k k^2} \theta \theta^T :$$

In $U = \mathbb{R}^p \times \mathbb{R}^m$, the function G is locally Lipschitz, and F is strictly convex, twice differentiable, and its Hessian is invertible, thus A2 is satisfied. If the initialization is such that $L(\theta_0) < L(0)$, then A4 holds with $C = \frac{1}{L(0) - L(\theta_0)}$. Combining this with the fact that $L(\theta_0) < L(0)$, as $U = \mathbb{R}^p \times \mathbb{R}^m$, it follows that $S \setminus U \neq \emptyset$. Consequently, A3 is also satisfied. To apply Theorem 6.2.4 and Theorem 6.2.8, we show that assumption A5 holds, namely that there exists $z \in S \setminus U$ such that $D_F(z; \cdot)$ is coercive. Indeed, the Bregman divergence $D_F(z; \cdot)$ is coercive for every $z \in \mathbb{R}^p$: from the Cauchy-Schwarz inequality,

$$\begin{aligned} D_F(z; \cdot) &= \frac{2}{3} k z k^{3=2} - \frac{2}{3} k k^{3=2} \langle k k^{1=2} h j z, \cdot \rangle \\ &= \frac{2}{3} k z k^{3=2} + \frac{1}{3} k k^{3=2} \langle k k^{1=2} h j z, \cdot \rangle \\ &= \frac{2}{3} k z k^{3=2} + \frac{1}{3} k k^{3=2} \langle k k^{1=2} k z k, \cdot \rangle \\ &= \frac{2}{3} k z k^{3=2} + k k^{1=3} \frac{1}{3} k k k z k : \end{aligned}$$

Then, when the norm of θ tends to infinity, $D_F(z; \cdot)$ also tends to infinity. Therefore, Theorem 6.2.4 yields the existence and uniqueness of the solution of the dynamical system (6.2.10). From Theorem 6.2.8 we get that the function $L(\theta(t))$ is non-increasing in $t \geq 0$ and that, for every $z \in S$ and for every $t > 0$,

$$L(\theta(t)) - L(\theta_0) \leq \frac{(2kr L(0)k)^{1=2} D_F(z; \theta_0)}{(L(0) - L(\theta_0))^{1=2} t}.$$

Moreover, $\lim_{t \rightarrow +\infty} \theta(t) = \theta_1$, where

$$\theta_1 \in \underset{z \in S}{\text{argmin}} \left\{ k z k^{3=2} - \frac{h_0 j z i}{k_0 k} + \int_0^{Z+1} \frac{2}{3} k(\tau) k^{3=2} h r L(\tau) j z(\tau) d\tau \right\}$$

Remark 6.2.11. In the specific case of the quadratic loss function $L(\cdot) = \|X\cdot - y\|^2$ for some $X \in \mathbb{R}^{d \times p}$ and $y \in \mathbb{R}^d$, Corollary 6.2.9, implies

$$\lim_{t \rightarrow \infty} \gamma(t) = \frac{1}{2} \operatorname{argmin}_{z \in S} \|z\|^{3-2} \frac{\langle z, y \rangle}{\|z\|}.$$

Therefore we recover the implicit bias presented in [7], and our result generalizes the latter to general losses. Moreover, in [7], the authors assume the convergence of the trajectory to a feasible point while we prove it.

Remark 6.2.12. Our results can be contrasted to those in [6]. In the latter paper, given N matrices $W_i \in \mathbb{R}^{p_i \times p_{i+1}}$ for every i with $p_1 = p$, $p_{N+1} = 1$, the authors consider the reparametrization $W = W_1 \cdots W_N$. For the special case $N = 2$, this setting coincides with ours. In [6, Claim 2], the flow (A.1.9) has been derived, and interpreted in a different way. Our results shed new light on their result, thanks to time warping mirror flow interpretation. Thanks to this, we can transform the impossibility result [6, Theorem 2], stating that the trajectory $\gamma(t)$ does not follow any gradient flow, into an existence result if the geometry is modified through a time warped mirror function. Moreover, the mirror flow interpretation gives the possibility of naturally deriving an implicit bias.

Remark 6.2.13. The paper [7] considered the case when $W^{\succ}(0)W(0) = w(0)w^{\succ}(0) = \operatorname{Id}$ for $\gamma \neq 0$. Observe that for suitable initializations, the assumptions of Theorem 6.2.8 still remain valid. In this case, we could compute the implicit bias, which is the generalization of the one proposed in [6][7, Theorem 2] for the least square, see [7, Appendix B]. The computations are very technical and are beyond the scope of this paper.

We can further characterize the limit point for fully connected linear networks of depth 2, in the linear setting of Corollary 6.2.9. The proof uses techniques similar to those in [139].

Corollary 6.2.14. In the setting of this section, suppose that, as in Corollary 6.2.9, $L(\cdot) = \|(X\cdot - y)\|^2$ and $S = \{z \in \mathbb{R}^p \mid Xz = y\}$, where $X \in \mathbb{R}^{d \times p}$ and $y \in \mathbb{R}^d$. Additionally, suppose that assumption A1 is satisfied, that the initialization satisfies $W^{\succ}(0)W(0) = w(0)w^{\succ}(0)$ and that $L(\gamma_0) < L(0)$. Let γ be the (unique) solution of

$$\min_{z \in S} \|z\|^2. \quad (6.2.11)$$

Let $P : \mathbb{R}^p \rightarrow \mathbb{R}^p$ be the projection operator onto the kernel of X . Then the trajectory $\gamma(t) = W(t)w(t)$ satisfies that $\lim_{t \rightarrow \infty} \gamma(t) = \gamma$, where

$$\gamma = \frac{S}{\|S\|} + P(\gamma_0) \frac{\|k_1\|}{\|k_0\|}; \quad \text{and}$$

$$\|k_1\| = \frac{\|kP(\gamma_0)\|^2 + \frac{P}{kP(\gamma_0)k^4 + 4k_0k^2k_1k^2}}{2k_0k}.$$

Proof. As shown in Section 6.2.3 all the assumptions of Theorem 6.2.4 and Theorem 6.2.8 are satisfied. Recall that $\ker(X) \subset \mathbb{R}^p$ is a closed subspace and that $\ker(X)^\perp = \overline{\operatorname{ran}(X^\succ)} = \operatorname{ran}(X^\succ)$ (see for instance [15, Fact 2.25 (iv)]). Let us define $P_\perp = \operatorname{Id} - P$, the projection onto $\ker(X)^\perp$. Then we can write γ as

$$\gamma = P(\gamma_1) + P_\perp(\gamma_1):$$

First, we compute $P(\gamma_1)$. Theorem 6.2.4 implies that equation (A.1.9),

$$-\dot{\gamma}(t) = \|k(\gamma(t))\| \operatorname{Id} + \frac{\langle \dot{\gamma}(t), \gamma(t) \rangle}{\|k(\gamma(t))\|^2} \gamma(t); \quad (6.2.12)$$

has a unique solution defined on $[0; +\infty)$ and $k(t)k > 0$ for every $t \geq 0$. Since $P(rL(\cdot)) = P(X^>r(X)) = 0$, if we apply the linear operator P to (6.2.12), we obtain that

$$\begin{aligned} \frac{d}{dt}P(k(t)) &= P(-\dot{k}(t)) = -k(t)kP(rL(\dot{k}(t))) - \frac{P(k(t))}{k(t)k}h(t)j r L(\dot{k}(t))i \\ &= -\frac{P(k(t))}{k(t)k}h(t)j r L(\dot{k}(t))i: \end{aligned} \quad (6.2.13)$$

On the other hand, (6.2.12) yields

$$\frac{d}{dt}k(t)k = \frac{\dot{k}(t)}{k(t)k}j r L(\dot{k}(t))i = 2h(t)j r L(\dot{k}(t))i: \quad (6.2.14)$$

Combining equations (6.2.13) and (6.2.14), we obtain the following differential equation:

$$\frac{d}{dt}(P(k(t))) = \frac{P(k(t))}{2k(t)k} \frac{d}{dt}k(t)k: \quad (6.2.15)$$

Equation (6.2.15) implies

$$\begin{aligned} \frac{d}{dt} \frac{P(k(t))}{k(t)k^{1-2}} &= \frac{1}{k(t)k^{1-2}} \frac{d}{dt}(P(k(t))) - \frac{P(k(t))}{2k(t)k^{3-2}} \frac{d}{dt}k(t)k \\ &= \frac{1}{k(t)k^{1-2}} \frac{d}{dt}(P(k(t))) - \frac{P(k(t))}{2k(t)k} \frac{d}{dt}k(t)k = 0: \end{aligned}$$

The previous result implies that the term $\frac{P(k(t))}{k(t)k^{1-2}}$ remains constant for every $t > 0$ and so it is equal to $\frac{P(k_0)}{k_0k^{1-2}}$. Consequently, we obtain the following expression:

$$P(k(t)) = P(k_0) \frac{k(t)k^{1-2}}{k_0k^{1-2}}:$$

Recall that, from Theorem 6.2.8, $\lim_{t \rightarrow +\infty} k(t) = k_1$. From the continuity of the projection and of the norm, taking the limit for $t \rightarrow +\infty$ in the previous equation, we get

$$P(k_1) = P(k_0) \frac{k_1k^{1-2}}{k_0k^{1-2}}:$$

Next we compute $P_? (k_1)$. Recall that $k_1 \in S$, where $S = \{f \in \mathbb{R}^p \mid jX = yf\}$. We have that

$$y = X k_1 = X(P(k_1) + P_?(k_1)) = X P_?(k_1):$$

Thus, $P_?(k_1) \in S$. Similarly, since $k_0 \in S$, we have that

$$y = X k_0 = X(P(k_0) + P_?(k_0)) = X P_?(k_0):$$

Then, $P_?(k_0) \in S$, meaning that it is a feasible point for the minimization problem (6.2.11). Moreover, since k_0 is a solution of problem (6.2.11), we have that

$$\|k_0\|^2 = \|k P_?(k_0)\|^2 = \|k P_?(k_1)\|^2 + \|k P(k_0)\|^2 = \|k_1\|^2:$$

Therefore, we conclude that $P(k_0) = 0$ and so that $k_0 \in \ker(X)$. Moreover, since $X k_0 = y = X k_1$, we obtain that $k_1 \in \ker X$. Then, since $k_1 \in \ker(X)$, we derive from

$$8p \in \ker(X) \quad hp \quad j \quad i = 0;$$

that $P_{\mathcal{S}}(z_1) = \dots$. Then,

$$z_1 = P(z_1) + P_{\mathcal{S}}(z_1) = P(z_0) \frac{\|z_1\|}{\|z_0\|} + \dots$$

Finally, we obtain the value of $\|z_1\|$ by solving the following second order equation:

$$\|z_1\|^2 = \|P_{\mathcal{S}}(z_1)\|^2 + \|P(z_1)\|^2 = \|z_0\|^2 + \|P(z_0)\|^2 \frac{\|z_1\|^2}{\|z_0\|^2}.$$

□

The previous result establishes that when we run gradient flow on $\mathcal{S} = (W; w)$, the corresponding trajectory is biased towards a specific element of \mathcal{S} , that is z_1 given by the sum of two terms: z_0 , the minimal-norm vector in \mathcal{S} ; and a rescaling of $P(z_0)$, the projection of the initialization onto $\ker(X)$.

Remark 6.2.15. We compare the previous result with the well-known one for vanilla gradient flow, namely

$$\begin{aligned} \dot{z} &= -r \nabla L(z(t)) \\ z(0) &= z_0 \end{aligned} \tag{6.2.16}$$

On the one hand, the system (6.2.16) is an instance of (6.2.1) with $F(z) = (1/2)\|z\|^2$ and $G(z) = 1$. Therefore, $D_F(z; z_0) = (1/2)\|z - z_0\|^2$, and we derive from (6.2.9) that $z(t) \rightarrow z_1$, where z_1 is the projection onto \mathcal{S} of z_0 . Following a similar reasoning to the one in the proof of Corollary 6.2.14, it is possible to show that

$$z_1 = z_0 + P(z_0).$$

We derive that the main difference in the implicit bias between running vanilla gradient flow on \mathcal{S} and gradient flow on the reparameterization $\mathcal{S} = (W; w)$ is given by the scaling factor on $P(z_0)$, namely $\frac{\|z_1\|}{\|z_0\|}$. In particular $z_1 = z_1$ if $z_0 \in \ker(X)$.

6.3 Reparameterizing Mirror Descent for radial functions as Projected Gradient Descent

In this section we consider reparameterizations in terms of polar coordinates, related to weight normalization [89, 118, 139], namely of the form

$$w = h(\theta) \frac{w}{\|w\|},$$

where $h: \mathbb{R} \rightarrow \mathbb{R}$ is a given function and $\theta \in \mathbb{R}$ and $w \in \mathbb{R}^p$ are the new variables. The associated gradient flow on the reparameterization is given by:

$$\dot{w}(t) = -r_w \nabla L(w(t)) = \frac{h'(\theta(t))}{\|w(t)\|} \text{Id} \frac{w(t)w^\top(t)}{\|w(t)\|^2} - r \nabla L(w(t)) \tag{6.3.1}$$

$$\dot{\theta}(t) = -r \nabla L(w(t)) = h''(\theta(t)) \frac{w(t)}{\|w(t)\|^3} \cdot r \nabla L(w(t)) \tag{6.3.2}$$

In the following we will show that this flow is equivalent to a generalized mirror one, for a mirror function of the form $F(z) = f(\|z\|)$ and a time warping of the form $g(\|z\|)$. To derive the main result of this section, we need an auxiliary Lemma, providing a reformulation of (6.2.1) for radial functions. From now on, unless explicitly stated otherwise, when we refer to assumption A2, we mean it with $U = \mathbb{R}^p \cap \text{ker}(X)$.

Lemma 6.3.1. Let assumptions **A1** and **A3** hold. Let $f: [0; +\infty) \rightarrow \mathbb{R}$ be a twice differentiable function on $(0; +\infty)$ such that $f'(s) > 0$, $f''(s) > 0$ for every $s > 0$ and $1=f''$ is locally Lipschitz on $(0; +\infty)$. Let $g: [0; +\infty) \rightarrow \mathbb{R}$ be locally Lipschitz on $(0; +\infty)$ and such that, defining r as in Assumption **A4**, there exists C such that $g(s) > C$ for every $s > r$. For every $\gamma \in \mathbb{R}^p$, define $F(\cdot) = f(k \cdot k)$ and $G(\cdot) = g(k \cdot k)$. Then assumptions **A2** and **A4** hold, and the dynamical system (6.2.1) corresponding to F and G can be written equivalently as

$$\begin{aligned} \dot{k}(t) &= -\frac{g(k(t)k'(t)k''(t))}{f'(k(t)k'(t))} \text{Id} - \frac{g'(k(t)k'(t))}{k(t)k'(t)^2} r L(k(t)) \\ \dot{g}(t) &= \frac{g(k(t)k'(t))}{f''(k(t)k'(t))} \frac{g'(k(t)k'(t))}{k(t)k'(t)^2} r L(k(t)); \quad t > 0 \\ (0) &= \gamma \in \mathbb{R}^p \text{ n } f'0g: \end{aligned} \tag{6.3.3}$$

Proof. The function F is strictly convex since it is the composition of the increasing strictly convex function f with the convex function $\|\cdot\|_k$. The gradient and the Hessian of $F(\cdot) = f(k \cdot k)$, for $\cdot \neq 0$ are

$$\begin{aligned} \nabla F(\cdot) &= f'(k \cdot k) \frac{\cdot}{k \cdot k} \quad \text{and} \\ \nabla^2 F(\cdot) &= f''(k \cdot k) \frac{\cdot \cdot}{k \cdot k^2} + \frac{f'(k \cdot k)}{k \cdot k} \text{Id} - \frac{f''(k \cdot k)}{k \cdot k^2} \cdot; \end{aligned} \tag{6.3.4}$$

and clearly F is twice differentiable on $\mathbb{R}^p \setminus \{0\}$. Using the Sherman-Morrison formula (2.2.1) with $a = f''(k \cdot k) \frac{\cdot \cdot}{k \cdot k^2}$, $b = \cdot$, and $M = \frac{f'(k \cdot k)}{k \cdot k} \text{Id}$, since $b^T M^{-1} a + 1 = \frac{k \cdot k f''(k \cdot k)}{f''(k \cdot k)} \neq 0$ we get that

$$\nabla^2 F(\cdot)^{-1} = \frac{1}{f''(k \cdot k)} \frac{\cdot \cdot}{k \cdot k^2} + \frac{k \cdot k}{f'(k \cdot k)} \text{Id} - \frac{\cdot \cdot}{k \cdot k^2} \cdot; \tag{6.3.5}$$

which is locally Lipschitz in $U \setminus \{0\}$ since $1=f''$ and f' are locally Lipschitz in $(0; +\infty)$. In addition, $S \setminus U \neq \emptyset$; since 0 is not a minimizer of L by Assumption **A3**. Local Lipschitzianity of g implies that G is locally Lipschitz. Therefore Assumption **A2** is satisfied. Assumption **A4** is also trivially satisfied, therefore Lemma 6.2.6 implies that the solution of the dynamical system is bounded away from zero; namely, that $k(t)k'(t) > r > 0$ for every $t \geq 0$. This implies that all the quantities in (6.3.4) and in (6.3.5) are well-defined along the trajectory. The statement then follows by plugging the explicit form of $\nabla^2 F(\cdot)^{-1}$ in equation (6.2.2), which is equivalent to (6.2.1). \square

Remark 6.3.2. By (6.3.3), we get that

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} k(t)k'(t)k''(t) &= \frac{d}{dt} \left(\frac{g(k(t)k'(t))}{f''(k(t)k'(t))} h(t) \right) \\ \frac{d}{dt} \frac{g(k(t)k'(t))}{f''(k(t)k'(t))} &= \frac{1}{k(t)k'(t)} \text{Id} - \frac{g'(k(t)k'(t))}{k(t)k'(t)^2} \cdot \end{aligned}$$

This computation shows that if the second term on the right-hand side of (6.3.3) is zero, then the norm is constant.

In the next theorem, we focus on reparametrizations of the form $\dot{w} = h(\cdot) \frac{w}{k \cdot w \cdot k}$. This theorem establishes a sufficient condition that allows us to derive from the gradient flow with respect to \dot{w} and w a time-warping mirror flow with respect to \dot{k} , as in (6.2.1).

Theorem 6.3.3. Let assumption **A1** hold. Consider the weight normalization parameterization $\tilde{w} = h(\cdot) \frac{w}{\|w\|}$ and the dynamics on $(\cdot; w)$ defined in (6.3.1)-(6.3.2) with initialization $(\cdot_0; w_0)$ such that $h(\cdot_0) > 0$ and $\|w_0\| = 1$ and assume that assumption **A3** holds for $\cdot_0 = h(\cdot_0) w_0 = \|w_0\| k$. Assume that there exists a function $f: [0; +\infty) \rightarrow [0; +\infty)$ twice differentiable on $(0; +\infty)$ such that $f'(s) > 0$ for every $s > 0$, $1=f''$ is locally Lipschitz on $(0; +\infty)$ and satisfying the following equality

$$(h'(\cdot))^2 = \frac{f''(h(\cdot))h(\cdot)}{f'(h(\cdot))}. \quad (6.3.6)$$

Suppose in addition that there exists $c > 0$ such that $f'(s) > c=s$ for every $s > r$ (r is defined in Assumption **A3**) and define $g: \mathbb{R}_+ \rightarrow \mathbb{R}$ as $g(s) = sf'(s)$. Then the assumptions **A2** and **A4** are satisfied, and for every $t > 0$, $\|w(t)\| = 1$, $h(\cdot(t)) > 0$, and $\tilde{w}(t)$ solves the following dynamical system, for $t \geq 0$

$$\begin{aligned} \dot{\tilde{w}}(t) &= -k(t)k^2 \text{Id} - \frac{h'(t) \tilde{w}(t)}{k(t)k^2} r L(\tilde{w}(t)) \\ \tilde{w}(0) &= \frac{f'(k(t)k)}{f''(k(t)k)} \frac{h'(t) \tilde{w}(t)}{k(t)k} r L(\tilde{w}(t)); t > 0 \end{aligned} \quad (6.3.7)$$

i.e., the dynamic of $\tilde{w}(t)$ satisfies (6.2.1), with $F(\cdot) = f(k \cdot k)$ and $G(\cdot) = g(k \cdot k)$.

Proof. We first show that $t \geq (0; +\infty) \ni \|w(t)\|$ is constant. Since w is a solution of the dynamical system in (6.3.1), we obtain

$$\frac{d}{dt} \|w(t)\|^2 = 2h w(t) \cdot \dot{w}(t) = 0;$$

therefore the norm of $w(t)$ is constant. Moreover, since $\|w(0)\| = 1$ by assumption, for every $t > 0$, $\|w(t)\| = 1$. We next derive an equation for $\tilde{w}(t)$ using (6.3.1)-(6.3.2). The derivative of the product yields

$$\begin{aligned} \dot{\tilde{w}}(t) &= h'(\tilde{w}(t)) w(t) + h(\tilde{w}(t)) \dot{w}(t) \\ &= \frac{(h'(\tilde{w}(t)))^2}{\|w(t)\|^2} \text{Id} - \frac{h(\tilde{w}(t)) \dot{w}(t)}{\|w(t)\|^2} r L(\tilde{w}(t)) \\ &= h'(\tilde{w}(t))^2 w(t) \dot{w}(t) r L(\tilde{w}(t)). \end{aligned} \quad (6.3.8)$$

We now prove that, for every $t > 0$, $h(\tilde{w}(t)) > 0$ and consequently $\|k(\tilde{w}(t))\| = h(\tilde{w}(t))$. It follows from (6.3.8) that

$$\begin{aligned} \frac{d}{dt} L(\tilde{w}(t)) &= \frac{D}{r L(\tilde{w}(t))} \dot{\tilde{w}}(t) \\ &= \frac{(h'(\tilde{w}(t)))^2}{\|w(t)\|^2} k r L(\tilde{w}(t)) k^2 + \frac{(h(\tilde{w}(t)))^2}{\|w(t)\|^2} r L(\tilde{w}(t)) \dot{w}(t) \cdot \frac{w(t)}{\|w(t)\|} \\ &= h'(\tilde{w}(t))^2 h r L(\tilde{w}(t)) \dot{w}(t) \cdot w(t) \\ &= h'(\tilde{w}(t))^2 h r L(\tilde{w}(t)) \dot{w}(t) \cdot w(t) \\ &= 0; \end{aligned}$$

Therefore $t \geq (0; +\infty) \ni L(\tilde{w}(t))$ is decreasing. Next, we provide a strictly-positive lower bound for the norm of \tilde{w} along the trajectory. Let $\tilde{w} \in \mathbb{R}^p \setminus \{0\}$ be such that $\|k(\tilde{w})\| = r =$

$\frac{L(0) - L(\theta)}{2krL(0)k}$. Then, convexity of L and the Cauchy-Schwarz inequality yield

$$\begin{aligned} L(\theta) &= L(0) + hrL(0)j - 0i \\ &= \frac{L(0) + L(\theta)}{2} \\ &> L(\theta); \end{aligned} \tag{6.3.9}$$

This implies on the one hand that $kh(\theta)w_0k = h(\theta) > r$ otherwise we would derive a contradiction from (6.3.9) choosing $\theta = 0$. On the other hand, from (6.3.9) we derive that $k(\dot{h}(t))k > r$ for every $t > 0$ since $L(\dot{h}(t))$ is decreasing. Then, for every $t > 0$, $kh(\dot{h}(t))k > r$, which implies that $h(\dot{h}(t)) = k(\dot{h}(t))k$, since the sign of $h(\dot{h}(t))$ remains positive. Replacing $w(t) = \frac{\dot{h}(t)}{k(\dot{h}(t))k}$, $kw(t)k = 1$, $k(\dot{h}(t))k = h(\dot{h}(t))$, and

$$(h^\theta(\theta))^2 = \frac{f^\theta(h(\theta))h(\theta)}{f^{\theta\theta}(h(\theta))};$$

in (6.3.8), we deduce that

$$\begin{aligned} -(\dot{h}) &= k(\dot{h}(t))k^2 \text{Id} - \frac{(\dot{h}(t)) \cdot (\dot{h}(t))}{k(\dot{h}(t))k^2} rL(\dot{h}(t)) \\ &= \frac{f^\theta(k(\dot{h}(t))k)k(\dot{h}(t))k - (\dot{h}(t)) \cdot (\dot{h}(t))}{f^{\theta\theta}(k(\dot{h}(t))k)} rL(\dot{h}(t)) \end{aligned}$$

Finally, Lemma 6.3.1 with $F(\theta) = f(k(\dot{h}(t))k)$ and $G(\theta) = g(k(\dot{h}(t))k) = k(\dot{h}(t))k f^\theta(k(\dot{h}(t))k)$ implies (6.3.7) and proves the statement. \square

Remark 6.3.4. 1. If $h(\theta) > 0$, from the definition of weight normalization reparameterization $\theta = h(\theta) \frac{w}{kwk}$, we have that $k(\dot{h}(t))k = h(\dot{h}(t))$. The previous theorem guarantees also that, if the initialization satisfies $h(\theta_0) > 0$, then $h(\dot{h}(t)) > 0$ for every $t > 0$ and, consequently, $k(\dot{h}(t))k = h(\dot{h}(t))$ for every $t > 0$.

2. The same dynamic for θ could be derived from other reparametrizations (for instance, see Example 5.3.1 and Remark 6.3.6(ii)). However, the proposed reparameterization only requires $p + 1$ parameters, which, in practice, makes this approach more efficient than other overparameterization schemes. Furthermore, this technique involves only one additional variable compared to vanilla gradient descent, and requires only matrix-vector multiplications and gradient calculations.
3. The previous theorem gives a sufficient condition on the functions h , f , and g to cast the gradient flow over θ and w into a time-warping mirror flow on \mathbb{R}^p , as in Theorems 6.2.4 and 6.2.8. This result allows us to establish the existence and uniqueness of the trajectory θ and convergence to an implicit bias.
4. If $(h^\theta(\theta))^2$ can be expressed in terms of $h(\theta)$ as $(h^\theta(\theta))^2 = \psi(h(\theta))$ for some function $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$, we can determine f^θ in (6.3.6) by solving the following differential equation:

$$\psi'(s) f^{\theta\theta}(s) - f^\theta(s)s = 0; \quad s > 0;$$

In the context of Corollary 6.2.9, where the loss is a composition of a convex function with a linear operator and the solution set is a linear subspace, the next corollary allows us to further characterize the limit point of the flow deriving from weight normalization. Similarly to Corollary 6.2.14, the limit point can be expressed as the sum of the minimal norm solution and a term in the kernel of the linear operator. We generalize the techniques used in [139, Theorem 2.6] to a broader class of weight-normalization reparameterizations.

Corollary 6.3.5. Under the same assumptions of Theorem 6.3.3, suppose that, $L(\cdot) = \langle X, \cdot \rangle$ and $S = \{f \in \mathbb{R}^p \mid Xf = y\}$, where $X \in \mathbb{R}^{d \times p}$ and $y \in \mathbb{R}^d$. Additionally, suppose that assumption A5 is satisfied. Let k^* be the (unique) solution of

$$\min_{z \in S} \|zk\|^2:$$

Moreover, let $P : \mathbb{R}^p \rightarrow \mathbb{R}^p$ be the projection operator onto the kernel of X . Then, for $t \in [0, +1)$, $k^*(t) = \operatorname{argmin}_{z \in S} f(k, z) = \frac{f'(k, 0)}{k} h_0; z_i$, and

$$k^*(t) = k^*(0) + P(k^*(0)) \frac{f'(k, 0)k^*(t) - k^*(0)}{k^*(0) f'(k, 0)k^*(t)}.$$

Proof. Let P be the projection operator onto $\ker(X)$ and let $P_\perp := \operatorname{Id} - P$, the projection onto $\ker(X)^\perp$. Then

$$k^*(t) = P(k^*(t)) + P_\perp(k^*(t)).$$

First, we compute $P(k^*(t))$. It follows from Theorem 6.3.3 that the differential equation

$$\begin{aligned} -\dot{k}^*(t) &= \|k^*(t)\|^2 \operatorname{Id} - \frac{\langle k^*(t), k^*(t) \rangle}{\|k^*(t)\|^2} r L(k^*(t)) \\ &= \frac{f''(k^*(t))k^*(t)k^*(t) - \langle k^*(t), k^*(t) \rangle r L(k^*(t))}{f''(k^*(t))k^*(t)k^*(t)} \end{aligned} \quad (6.3.10)$$

has a unique solution such that $k^*(0) = k_0$ defined on $[0, +1)$. Since $P(r L(k^*(t))) = P(X^T r(Xk^*(t))) = 0$, if we apply the linear operator P to both sides of the previous equation, we obtain that the differential equation

$$\begin{aligned} \frac{d}{dt} P(k^*(t)) &= P(-\dot{k}^*(t)) = -\|k^*(t)\|^2 P(r L(k^*(t))) \\ &\quad + P(k^*(t)) h_0; j r L(k^*(t)) i^{-1} \frac{f'(k^*(t))k^*(t)}{f''(k^*(t))k^*(t)k^*(t)} \\ &= -P(k^*(t)) h_0; j r L(k^*(t)) i^{-1} \frac{f'(k^*(t))k^*(t)}{f''(k^*(t))k^*(t)k^*(t)} \end{aligned} \quad (6.3.11)$$

Equation (6.3.10) implies

$$\frac{d}{dt} \|k^*(t)\|^2 = \frac{\langle k^*(t), -\dot{k}^*(t) \rangle}{\|k^*(t)\|^2} = \frac{f''(k^*(t))k^*(t)k^*(t) - \langle k^*(t), k^*(t) \rangle h_0; j r L(k^*(t)) i}{f''(k^*(t))k^*(t)k^*(t)} \quad (6.3.12)$$

Combining equations (6.3.11) and (6.3.12), we obtain the following differential equation:

$$\frac{d}{dt} P(k^*(t)) = P(k^*(t)) \frac{1}{\|k^*(t)\|^2} \frac{f''(k^*(t))k^*(t)k^*(t) - \langle k^*(t), k^*(t) \rangle h_0; j r L(k^*(t)) i}{f''(k^*(t))k^*(t)k^*(t)} \frac{d}{dt} \|k^*(t)\|^2 \quad (6.3.13)$$

By multiplying the equation (6.3.13) by the integrating factor $\frac{f''(k^*(t))k^*(t)k^*(t)}{\|k^*(t)\|^2}$ (which is strictly positive as proven in Lemma 6.2.6) and then rearranging the terms, we obtain the following result:

$$\begin{aligned} \frac{d}{dt} \frac{P(k^*(t)) f''(k^*(t))k^*(t)k^*(t)}{\|k^*(t)\|^2} &= \frac{f''(k^*(t))k^*(t)k^*(t)}{\|k^*(t)\|^2} \frac{d}{dt} (P(k^*(t))) \\ &\quad + P(k^*(t)) \frac{f''(k^*(t))k^*(t)k^*(t)}{\|k^*(t)\|^2} \frac{f''(k^*(t))k^*(t)k^*(t) - \langle k^*(t), k^*(t) \rangle h_0; j r L(k^*(t)) i}{f''(k^*(t))k^*(t)k^*(t)} \frac{d}{dt} \|k^*(t)\|^2 \\ &= \frac{f''(k^*(t))k^*(t)k^*(t)}{\|k^*(t)\|^2} \frac{d}{dt} (P(k^*(t))) \\ &\quad + \frac{f''(k^*(t))k^*(t)k^*(t)}{\|k^*(t)\|^2} P(k^*(t)) \frac{1}{\|k^*(t)\|^2} \frac{f''(k^*(t))k^*(t)k^*(t) - \langle k^*(t), k^*(t) \rangle h_0; j r L(k^*(t)) i}{f''(k^*(t))k^*(t)k^*(t)} \frac{d}{dt} \|k^*(t)\|^2 = 0. \end{aligned}$$

Thus the term $\frac{P(\tau) f^0(k(\tau))}{k(\tau)}$ remains constant for every $t > 0$ and it is equal to $\frac{P(0) f^0(k(0))}{k(0)}$. Consequently, we obtain the following expression:

$$P(\tau) = P(0) \frac{f^0(k(0))k(\tau)}{k(0)f^0(k(\tau))}.$$

Recall that, from Theorem 6.2.8, $\lim_{t \rightarrow \infty} k(t) = k^*$. From the continuity of the projection and the continuity of the norm, taking the limit for $t \rightarrow \infty$ in the previous equation, we get

$$P(k^*) = P(0) \frac{f^0(k(0))k^*}{k(0)f^0(k^*)}.$$

Next we proceed with the computation of the second term, namely $P_{\mathcal{C}}(k^*)$. Since $k^* \in S$,

$$y = X k^* = X(P(k^*) + P_{\mathcal{C}}(k^*)) = X P_{\mathcal{C}}(k^*);$$

Thus, $P_{\mathcal{C}}(k^*) \in S$. Similarly, since $k^* \in S$, $P(k^*) \in S$, meaning that it is a feasible point for the minimization problem (6.2.11). Moreover, since k^* is a solution of problem (6.2.11), we have that

$$\|k^*\|^2 - \|P(k^*)\|^2 = \|P_{\mathcal{C}}(k^*)\|^2 + \|P(k^*)\|^2 = \|k^*\|^2.$$

Therefore, we conclude that $P(k^*) = 0$ and so that $k^* \in \ker(X)$. Moreover, since $X k^* = y = X P_{\mathcal{C}}(k^*)$, we obtain that $P_{\mathcal{C}}(k^*) \in \ker X$. Then, since $k^* \in \ker(X)$,

$$\exists p \in \ker(X) \text{ s.t. } \|k^* - p\| = 0.$$

This, according to [15, Theorem 3.16], implies that $P_{\mathcal{C}}(k^*) = k^*$. Then,

$$k^* = P(k^*) + P_{\mathcal{C}}(k^*) = P(0) \frac{f^0(k(0))k^*}{k(0)f^0(k^*)} + k^*.$$

□

Note that, compared to the analogous result in Corollary 6.2.14, this is a fixed point equation, in the sense that we cannot compute the norm of k^* explicitly.

In the following two examples we explore the flexibility of our formulation. By choosing appropriately h and f , we obtain the two well-known reparameterizations corresponding to polynomial and exponential weight normalization [89, 118, 139].

Example 6.3.6. Polynomial Weight normalization. Let Assumption A1 be satisfied, and consider the dynamics defined in (6.3.1)-(6.3.2) corresponding to

$$h(s) = L - s;$$

with an initialization such that $k(0) = 1$ and $h(0) > 0$, so that Assumption A3 holds. Let $f: [0; +\infty)$ be such that

$$f(s) = \exp\left(-\frac{s^2}{2L}\right) \text{ for every } s > 0;$$

so that equation (6.3.6) is satisfied. Note that f is twice differentiable on $(0; +\infty)$, $f(s) > 0$, $f'(s) < 0$ for all $s > 0$, and f is locally Lipschitz. Define $g(s) = s f(s)$ for every $s \geq 0$ and let r be as in Assumption A4. Since g is increasing, $g(s) \geq g(r)$ for every $s \geq r$. Then Theorem 6.3.3 and Lemma 6.3.1 imply that Assumptions A2 and A5 are satisfied

with $U = \mathbb{R}^p \cap \text{ran}(X)$ and the differential equation in the original variable corresponding to the gradient flow on the variables $(\cdot; w)$ is

$$\begin{aligned} -\dot{w}(t) &= -k(t)k^2 \text{Id} - \frac{\dot{w}(t) \langle w(t), \cdot \rangle}{k(t)k^2} r_L(\cdot)(t) \\ L^2 \frac{\dot{w}(t) \langle w(t), \cdot \rangle}{k(t)k^2} &\in r_L(\cdot)(t): \quad t > 0 \end{aligned}$$

This equation has a unique solution $w(t)$ defined on $[0; +\infty)$ for the initial condition $w_0 = h(w_0)w_0$. Note that, from Lemma 6.2.6, along the trajectory $w(t)$, the scalar function G is bounded below by a strictly positive constant. This allows us to avoid the undesired stationary point zero. To derive the implicit bias, we need to apply Theorem 6.2.8. To do so, it remains to show that assumption A5 holds. Let $z \in S \setminus U$. We next prove that $D_F(z; \cdot)$ is coercive. It follows from convexity of f that for every $s > 0$

$$f(s) + f'(s)(k - s) \leq f(k):$$

Dividing the previous inequality by $k - s$ yields

$$f'(s) = \lim_{k \downarrow s} \frac{f(s) + f'(s)(k - s)}{k - s} = \liminf_{k \downarrow s} \frac{f(k) - f(s)}{k - s} = \liminf_{k \downarrow s} \frac{F(k) - F(s)}{k - s}.$$

Since $s > 0$ is arbitrary and $f'(s) \rightarrow +\infty$ as $s \rightarrow 0^+$, we deduce that F is supercoercive. It follows from [14, Lemma 7.3 (viii)] that the $D_F(z; \cdot)$ is coercive for every $z \in U$, which implies assumption A5. The implicit bias is given in Theorem 6.2.8 and specialized in Corollary 6.3.5. For the weight normalization proposed by [118], which corresponds to the choice $L = 1$, in the setting of Corollary 6.3.5, we obtain

$$w_1 = w_0 + P(w_0) \frac{k_1 - k_0}{k_0 k_1} e^{-\frac{k_0 k_1}{2}}; \quad (6.3.14)$$

where we recall that P is the projection operator on the kernel of X . Note that when $w_0 \in \text{ran}(X) = \ker(X)^\perp$, then w_1 is equal to the minimal norm solution w^* . Our results guarantee that the trajectory exists for $t \in [0; +\infty)$, which is something usually assumed, and not proved. The implicit bias for this example has been studied in [139] and [89] for the least squares and the logistic loss, respectively. In [139], the authors analyze the convergence to a stationary point (zero loss or zero norm) in the continuous and discrete settings, and they derive the expression in (6.3.14). The squared weight normalization corresponds to the case $L = 2$. For this choice, we can explicitly compute the function f . Indeed, let

$$h(w) = \|w\|^2; \quad f(s) = 4e^{-\frac{s}{4}}; \quad \text{and} \quad g(s) = se^{-\frac{s}{4}}.$$

In addition, in the setting of Corollary 6.3.5, we get that

$$w_1 = w_0 + P(w_0) \frac{k_1 - k_0}{k_0 k_1} e^{-\frac{k_0 k_1}{4}}.$$

Example 6.3.7. Exponential weight normalization. Under Assumption A1, consider the dynamic defined in (6.3.1)-(6.3.2) with an initialization $(w_0; w_0)$ such that $\|kw_0\| = 1$ and $L(w_0) < L(0)$, where $w_0 = h(w_0)w_0$. Let, for some $\alpha > 0$,

$$h(w) = e^{-\alpha \|w\|^2}; \quad f(s) = \frac{s^{2+\alpha}}{2+\alpha}; \quad \text{and} \quad g(s) = s^{2+\alpha}.$$

On the interval $(0; +\infty)$, we have $f^0 > 0$, $f^{(0)} > 0$ and $1 = f^{(0)}$ locally Lipschitz. In addition, equation (6.3.6) is satisfied, $g(s) = sf^0(s)$ for all $s \geq 0$ and $g(s) > r^{2+1}$ for every $s > r$. Then it follows from Theorem 6.3.3 that $\tilde{w}(t) = h(\tilde{w}(t)) \frac{w(t)}{k w(t) k}$, where the trajectories of $w(t)$ and $\tilde{w}(t)$ follow the dynamical system (6.3.1)-(6.3.2), solves the differential equation

$$\dot{\tilde{w}}(t) = -k(\tilde{w}(t))k^2 \text{Id} - \frac{1}{2} \tilde{w}(t) \langle \tilde{w}(t), r L(\tilde{w}(t)) \rangle; \quad t > 0$$

To apply Theorem 6.2.8, it remains to prove assumption A5, that requires the existence of $Z \subset S \setminus U$ such that $D_F(Z; \cdot)$ is coercive. Indeed,

$$\begin{aligned} D_F(Z; \cdot) &= \frac{kZk^{2+1}}{2+1} - \frac{k k^{2+1}}{2+1} - k k^{2+1} h_j z_i \\ &= \frac{kZk^{2+1}}{2+1} + \frac{k k^{2+1}}{2+1} - \frac{2+1}{k k^2} h_j z_i \\ &= \frac{kZk^{2+1}}{2+1} + \frac{k k^{2+1}}{2+1} - \frac{2+1}{k k} kZk \end{aligned}$$

and so $D_F(Z; \cdot)$ is coercive with respect to \tilde{w} for every $Z \subset U$. Theorem 6.2.8 then yields $\lim_{t \rightarrow +\infty} \tilde{w}(t) = \tilde{w}_1$, where

$$\tilde{w}_1 \in \underset{Z \subset S}{\text{argmin}} \frac{kZk^{2+1}}{2+1} - k k^{2+1} h_j z_i + \int_0^{Z+1} k(\tilde{w}(t))k^{2+1} h r L(\tilde{w}(t)) j z_i(\tilde{w}(t)) dt:$$

Moreover, in the setting of Corollary 6.3.5, we obtain that

$$\tilde{w}_1 = \tilde{w}_0 + P(\tilde{w}_0) \frac{k \tilde{w}_0 k^{2+1}}{k \tilde{w}_1 k^{2+1}};$$

Note that, when $\tilde{w}_0 = \frac{1}{2}$, the mirror function F is the same function obtained in Example 5.3.1, showing that different reparametrizations may have the same implicit bias.

6.3.1 Weight normalization of a fully connected network

In this subsection, we combine the two main examples presented in this chapter: fully connected linear network and weight normalization. The reparameterization involves decomposing the vector \tilde{w} into the product of three terms: a scalar that represents the norm of \tilde{w} , and the product of a matrix by a vector, product that represents the direction of \tilde{w} with unitary norm.

Consider the over-parameterization $\tilde{w} = h(\tilde{w}) \frac{Ww}{k W w k}$, where $h: \mathbb{R} \rightarrow \mathbb{R}_+$ is a given function, while $\tilde{w} \in \mathbb{R}$, $W \in \mathbb{R}^{p \times m}$ and $w \in \mathbb{R}^m$ are the new variables. The associated gradient flow on the reparameterization is given by the following dynamical system:

$$\dot{w}(t) = -r_w L(\tilde{w}(t)) = \frac{h(\tilde{w}(t))}{k W(t) w(t) k} W^{\top}(t) \text{Id} - \frac{\tilde{w}(t) \langle \tilde{w}(t), r L(\tilde{w}(t)) \rangle}{k(\tilde{w}(t))k^2} r L(\tilde{w}(t)) \quad (6.3.15)$$

$$\dot{W}(t) = -r_W L(\tilde{w}(t)) = \frac{h(\tilde{w}(t))}{k W(t) w(t) k} \text{Id} - \frac{\tilde{w}(t) \langle \tilde{w}(t), r L(\tilde{w}(t)) \rangle}{k(\tilde{w}(t))k^2} r L(\tilde{w}(t)) w^{\top}(t) \quad (6.3.16)$$

$$\dot{\tilde{w}}(t) = -r L(\tilde{w}(t)) = h^0(\tilde{w}(t)) \frac{\tilde{w}(t)}{k(\tilde{w}(t))k} j r L(\tilde{w}(t)) : \quad (6.3.17)$$

We consider, for simplicity of computations, an initialization $(w_0; W_0; w_0)$ such that $\|w_0\| = 1$ and $w_0^> = W_0^> w_0$. In the next theorem, we derive a time warping mirror flow on w as in (6.2.1) from the gradient flow defined in (6.3.15)-(6.3.17) over w , W , and λ , respectively. We then apply Theorems 6.2.4 and 6.2.8, providing the existence and uniqueness of solution for the dynamical system and convergence in value and in trajectory. In the case where the loss is a composition of a convex function and a linear operator, the limit point can be characterized as the Bregman projection on the solution set for a given mirror map, which depends on the choice of the reparameterization.

Theorem 6.3.8. Consider the overparameterized model $y = h(\lambda) \frac{w w^>}{k w w^> k}$ and the dynamic defined in (6.3.15)-(6.3.17) with initialization such that $\|w_0\| = 1$ and $w_0^> = W_0^> w_0$. Assume that assumptions A1 and A3 holds for $\lambda_0 = h^{-1}(\lambda_0) \frac{W_0 w_0}{k W_0 w_0 k}$. Assume that there exists a function $f: [0; +\infty) \rightarrow [0; +\infty)$ twice differentiable on $(0; +\infty)$ such that $f''(s) > 0$ for every $s > 0$, $1 = f''(s)$ is locally Lipschitz on $(0; +\infty)$ and satisfying the following equality

$$(h'(\lambda))^2 = \frac{f'(h(\lambda))h(\lambda)}{f''(h(\lambda))}. \quad (6.3.18)$$

Suppose in addition that there exists $c > 0$ such that $f'(s) > c = s$ for every $s > r$ (r is defined in Assumption A3) and define $g: \mathbb{R}_+ \rightarrow \mathbb{R}$ as $g(s) = s f'(s)$. Then the assumptions A2 and A4 are satisfied, and for every $t > 0$, $\|w(t)\| = 1$, $h(\lambda(t)) > 0$, and $\lambda(t)$ solves the following dynamical system, for $t \geq 0$

$$r^2 (f'(k w(t) k))' - \dot{\lambda} = g(k w(t) k) r - L(\lambda(t)); \quad (6.3.19)$$

i.e., the dynamic of $\lambda(t)$ satisfies (6.2.1), with $F(\lambda) = f'(k w(t) k)$ and $G(k w(t) k) = g(k w(t) k)$.

Proof. If we multiply (6.3.15) by $w^>(t)$ on the right and (6.3.16) by $W^>(t)$ on the left, we obtain that

$$W^>(t) W(t) = w(t) w^>(t); \quad (6.3.20)$$

Adding the transpose of (6.3.20) yields

$$W^>(t) W(t) + W^>(t) W(t) = w(t) w^>(t) + w(t) w^>(t):$$

Since $W^>(0) W(0) = w^>(0) w^>(0)$, we obtain

$$W^>(t) W(t) = w(t) w^>(t); \quad (6.3.21)$$

Multiplying (6.3.21) by $w^>(t)$ on the left and by $w(t)$ on the right, we have that $\|w(t)\|^4 = \|W(t) w(t)\|^2$, for every t . Analogously, multiplying by $W(t)$ the equation (6.3.21) on the left and by $W^>(t)$ on the right we deduce that

$$(W(t) W^>(t))^2 = \frac{w(t) w^>(t)}{k w(t) k^2} \|W(t) w(t)\|^2; \quad (6.3.22)$$

Thus, $\|w(t)\|^2 = \|W(t) w(t)\|$. In addition, since $W(t) W^>(t)$ is positive semidefinite and symmetric, it is the unique square root of the right hand side of (6.3.22) and therefore

$$W(t) W^>(t) = \frac{w(t) w^>(t)}{k w(t) k^2} \|W(t) w(t)\|;$$

since the right hand side of (6.3.22) is a rank one matrix. Moreover, note that

$$\begin{aligned}
 W(t)w(t) &= W(t)w(t) + W(t)\underline{w}(t) \\
 &= \frac{h(t)}{k\|W(t)w(t)\|} \|W(t)w(t)\|^2 \text{Id} + W(t)W^\top(t) \text{Id} - \frac{h(t)}{k\|w(t)\|^2} r_L(t) \\
 &= h(t) \text{Id} + \frac{h(t)}{k\|w(t)\|^2} \text{Id} - \frac{h(t)}{k\|w(t)\|^2} r_L(t) \\
 &= h(t) \text{Id} - \frac{h(t)}{k\|w(t)\|^2} r_L(t)
 \end{aligned}$$

and

$$\begin{aligned}
 \frac{d}{dt} \frac{\|W(t)w(t)\|^2}{2} &= \frac{d}{dt} \left(W(t)w(t) \cdot W(t)w(t) + W(t)\underline{w}(t) \cdot W(t)\underline{w}(t) \right) \\
 &= \frac{h(t)}{k\|w(t)\|} \left(W(t)w(t) + W(t)\underline{w}(t) \right) \\
 &= h(t) \text{Id} - \frac{h(t)}{k\|w(t)\|^2} r_L(t) = 0;
 \end{aligned}$$

which implies that $\|W(t)w(t)\| = 1$ for every t . Then,

$$\begin{aligned}
 -\dot{h}(t) &= -\dot{h}(t) \frac{h(t)}{k\|w(t)\|^2} \|W(t)w(t)\|^2 + h(t) \left(W(t)w(t) + W(t)\underline{w}(t) \right) \\
 &= h(t) \frac{h(t)}{k\|w(t)\|^2} r_L(t) - k\|w(t)\|^2 \text{Id} - \frac{h(t)}{k\|w(t)\|^2} r_L(t);
 \end{aligned}$$

Equation (6.3.19) then follows from (6.3.18) and Lemma 6.3.1. \square

Remark 6.3.9. Choosing h , f , and g as in Examples 6.3.6 and 6.3.7, the assumptions of Theorem 6.3.8 are satisfied. Consequently, (t) satisfies the properties derived in Theorem 6.2.4, Theorem 6.2.8 and Corollary 6.2.9.

6.3.2 Fully connected normalized linear network of depth 2

In this subsection, we study a new reparameterization of (t) . The reparameterization decomposes the vector (t) into the product of a matrix times a unitary vector. Then, we apply vanilla gradient flow to update the vector and the matrix. The main difference with respect to the previous examples is that the unitary vector is multiplied by a matrix which encode the norm of (t) , rather than a scalar. We name this parametrization connected normalized linear network of depth 2. More formally, we set $(t) = W \frac{w}{\|w\|}$ where $W \in \mathbb{R}^{p \times m}$ and $w \in \mathbb{R}^m$. We define the gradient flows of W and w as:

$$\dot{W}(t) = -r_{W,L} \left(W(t) \frac{w(t)}{\|w(t)\|} \right) = -r_L(t) \frac{w^\top(t)}{\|w(t)\|} \tag{6.3.23}$$

$$\begin{aligned}
 \dot{w}(t) &= -r_{w,L} \left(W(t) \frac{w(t)}{\|w(t)\|} \right) \\
 &= \frac{1}{\|w(t)\|} \text{Id} - \frac{w(t)w^\top(t)}{\|w(t)\|^2} - W^\top(t) r_L(t); \tag{6.3.24}
 \end{aligned}$$

The next theorem shows that the corresponding trajectory in (t) is a preconditioned gradient flow as in (6.2.1). It turns out that, for some specific initialization, the corresponding trajectory (t) is just vanilla gradient flow. We then apply Theorems 6.2.4 and 6.2.8 to get existence and uniqueness of the trajectory of the flow, optimization properties and

convergence to a specific solution asymptotically for t approaching infinity.

Theorem 6.3.10. Consider the dynamic defined in (6.3.23)-(6.3.24) initialized with $(W_0; w_0)$ such that $\|w_0\| = 1$, set $\phi_0 = W_0 w_0 = \|w_0\|$ and suppose that $L(\phi_0) < L(0)$. Then:

1. The variable $\phi(t)$ is governed by the following differential equation:

$$\dot{\phi}(t) = -(\text{Id} + S(t)) \phi(t); \tag{6.3.25}$$

where $S(t) := W(t)W^>(t) - \phi(t) \phi^>(t) \phi^{-2}(t)$.

2. If $S_0 = W_0 W_0^> - \phi_0 \phi_0^> = 0$, $\phi(t)$ is governed by the following differential equation:

$$\dot{\phi}(t) = -\phi(t) L(\phi(t));$$

Proof. We provide the proof of Theorem 6.3.10, which is divided in several steps.

1. We first prove that $\|w(t)\|$ is constant. Multiplying (6.3.24) by $w^>(t)$ we obtain

$$\frac{d}{dt} \|w(t)\|^2 = 2 \langle w(t), \dot{w}(t) \rangle = 0;$$

Then, if $\|w(0)\| = 1$, we get that $\|w(t)\| = 1$. Equations (6.3.23)-(6.3.24) and the derivative of the product's rule yield

$$\begin{aligned} \dot{\phi}(t) &= W(t) \frac{w(t)}{\|w(t)\|} + W(t) \frac{w(t)}{\|w(t)\|} \frac{w(t)w^>(t)w(t)}{\|w(t)\|^3} \\ &= W(t) \frac{w(t)}{\|w(t)\|} + W(t) \frac{w(t)}{\|w(t)\|} \\ &= -\phi(t) L(\phi(t)) \frac{w^>(t)w(t)}{\|w(t)\|^2} \\ &\quad + \frac{W(t)}{\|w(t)\|^2} \text{Id} \frac{w(t)w^>(t)}{\|w(t)\|^2} - W^>(t) \phi(t) L(\phi(t)) \\ &= -\text{Id} + W(t)W^>(t) - \phi(t) \phi^>(t) - \phi(t) L(\phi(t)) \\ &= (\text{Id} + S(t)) \phi(t); \end{aligned}$$

which establishes equation (6.3.25). Moreover, note that $S(t) = W(t)W^>(t) - \phi(t) \phi^>(t) \phi^{-2}(t) = 0$. Indeed,

$$\begin{aligned} \langle y, y \rangle &= \langle W W^> y, y \rangle = \langle y, W W^> y \rangle = \langle y, W \frac{w}{\|w\|} \frac{w^>}{\|w\|} W^> y \rangle \\ &= \|W^> y\|^2 \frac{w^>}{\|w\|} W^> y \\ &= \|W^> y\|^2 \frac{w^>}{\|w\|} \|W^> y\|^2 \\ &= 0; \end{aligned}$$

2. The proof is divided in three parts.

(a) **Derivation of the evolution of $S(t) = W(t)W^\top(t)$ $\dot{S}(t)$:**

$$\dot{S}(t) = \dot{W}(t)W^\top(t) - \dot{W}(t)^\top W(t) + W(t)\dot{W}^\top(t) - \dot{W}(t)^\top W(t): \quad (6.3.26)$$

It follows from (6.3.23) that

$$\dot{W}(t)W^\top(t) = -rL(\dot{W}(t))^\top W(t): \quad (6.3.27)$$

Multiplying (6.3.25) on the right by $W^\top(t)$, (6.3.27) implies

$$\begin{aligned} -\dot{W}(t)^\top W(t) &= \dot{W}(t)W^\top(t) \\ &= W(t)\dot{W}^\top(t) - rL(\dot{W}(t))^\top W(t) \\ &= W(t)\dot{W}^\top(t) - S(t)rL(\dot{W}(t))^\top W(t): \end{aligned} \quad (6.3.28)$$

Analogously

$$\begin{aligned} \dot{W}(t)^\top W(t) &= -rL(\dot{W}(t))^\top (I + S(t)) \\ &= -W(t)\dot{W}^\top(t) - rL(\dot{W}(t))^\top S(t): \end{aligned} \quad (6.3.29)$$

It follows from (6.3.26), (6.3.27), (6.3.28) and (6.3.29) that

$$\dot{S}(t) = S(t)rL(\dot{W}(t))^\top W(t) + \dot{W}(t)^\top rL(\dot{W}(t))S(t): \quad (6.3.30)$$

(b) **Show that $S(t) = 0$ in a neighborhood of zero:** Let us denote by $\|k\|_F$ the frobenius norm of a matrix. Then $\frac{1}{2}\|k\|_F^2 = \text{Tr}(S(t)S(t)) = \text{Tr}(S(t)S(t))$. Equation (6.3.30) and the linearity and cyclical property of the trace imply

$$\begin{aligned} \frac{1}{2}\|k\|_F^2 &= \text{Tr}(S(t)rL(\dot{W}(t))^\top W(t)S(t) \\ &\quad + \text{Tr}(\dot{W}(t)^\top rL(\dot{W}(t))S^2(t)) \\ &= \text{Tr}(rL(\dot{W}(t))^\top W(t)S^2(t)) \\ &= \text{Tr}(\dot{W}(t)^\top rL(\dot{W}(t))S^2(t)) \\ &= \text{Tr}(rL(\dot{W}(t))^\top W(t) + \dot{W}(t)^\top rL(\dot{W}(t))S^2(t)) \\ &= \|rL(\dot{W}(t))^\top W(t)\|_F + \|\dot{W}(t)^\top rL(\dot{W}(t))S^2(t)\|_F \\ &\leq 2\|rL(\dot{W}(t))\|_F \|W(t)\|_F + \|\dot{W}(t)^\top rL(\dot{W}(t))\|_F \|S^2(t)\|_F \\ &\leq 2\|rL(\dot{W}(t))\|_F \|k\|_F + \|\dot{W}(t)^\top rL(\dot{W}(t))\|_F \|k\|_F^2; \end{aligned} \quad (6.3.31)$$

where in the last inequality we used that, for every $u \in \mathbb{R}^p$ and $v \in \mathbb{R}^p$, $\|uv^\top\|_F = \|u\|_F \|v\|_F$.

We next prove that the term $\|rL(\dot{W}(t))\|_F \|k\|_F$ is bounded, to apply Grönwall's inequality. We proceed computing the derivative of the norm along the trajectory:

$$\begin{aligned} \frac{d}{dt} \|k\|_F^2 &= 2\text{Tr}(\text{Id} + S(t))rL(\dot{W}(t))^\top \dot{W}(t) \\ &\quad + 2\text{Tr}(S(t)rL(\dot{W}(t))^\top \dot{W}(t)) \\ &\quad + 2\text{Tr}(L(z)^\top L(\dot{W}(t))): \end{aligned}$$

If $L(W_0 w_0 = k w_0 k) = L$, then $r L(0) = 0$. So the norm is constant and the trajectory remains bounded. Otherwise, $L(z) - L(0) < 0$ and, from $S_0 = 0$ we derive

$$\frac{d}{dt} \|k(t) - z k^2(0)\| = L(z) - L(0) < 0:$$

Since the norm and $\|k(t) - z k^2\|$ are continuous, there exists an interval $[0; t_1]$ with nonempty interior such that $\|k(t) - z k^2\|$ is decreasing and bounded by $\|k(0) - z k^2\|$. This implies that $\|k(t) - z k^2\| : t \in [0; t_1]$ is bounded. Similarly, since $r L$ is locally Lipschitz, $\|r L(\|k(t) - z k^2\|)\|$ is bounded on $[0; t_1]$ (see, for instance, [91, Theorem 27.4]). Define, for every $t \in [0; t_1]$:

$$G(t) = e^{-2 \int_0^t \|r L(\|k(s) - z k^2(s)\|)\| ds} > 0,$$

then, for every $t \in [0; t_1]$, (6.3.31) implies that

$$\frac{d}{dt} (\|k_S(t)\|_{k_F}^2 G(t)) = G(t) (\|k_S(t)\|_{k_F}^2 - 2 \|r L(\|k(t) - z k^2(t)\|)\| \|k_S(t)\|_{k_F}^2) = 0$$

and therefore $\|k_S(t)\|_{k_F}^2 = \|k_S(0)\|_{k_F}^2 e^{2 \int_0^t \|r L(\|k(s) - z k^2(s)\|)\| ds} = 0$.

(c) **Extend the proof, for every $t > 0$:** For every $t \in [0; t_1]$, we obtain

$$\|k(t) - z k^2\| = L(z) - L(\|k(t) - z k^2\|) = 0:$$

Now define $t^* = \sup\{t \in [0; T] \mid S(t) = 0\}$. The previous step implies $t^* = t_1$. Suppose by contradiction that $t^* < +\infty$, i.e., S is null only in a finite interval. Since $S(t)$ is continuous $S(t^*) = 0$ and we can proceed analogously to the case when $S(0) = 0$ and prove that there exist a $t_2 > t^*$ such that, for every $t \in [t^*; t_2]$, $\|k_S(t)\|_{k_F}^2 = \|k_S(t^*)\|_{k_F}^2 e^{2 \int_{t^*}^t \|r L(\|k(s) - z k^2(s)\|)\| ds} = 0$. This leads to a contradiction with the maximality of t^* and we conclude that $t^* = +\infty$ and $S(t) = 0$, for every $t \geq 0$. It follows from (6.3.25) that

$$\frac{d}{dt} \|k(t) - z k^2\| = -\|r L(\|k(t) - z k^2\|)\|;$$

and 2 follows. □

Remark 6.3.11. Consider the dynamics defined in (6.3.23)-(6.3.24) and suppose that A1 holds. In addition, consider an initialization such that $\|k(0)\| = 1$, $W(0)W^>(0)W(0)w(0)w^>(0)W^>(0) = 0$ and $L(0) < L(0)$. Moreover, by Theorem 6.3.10, Assumption A2-A5 is satisfied with $F(\cdot) = \frac{k \cdot k^2}{2}$ and $G(\cdot) = 1$. Note that, in this case, the assumption A2 holds with $U = \mathbb{R}^p$. Consequently, $\|k(t)\|$ can be characterized as described in Theorem 6.2.4 and Theorem 6.2.8.

6.4 Conclusion and Future Work

In this chapter, we studied a unified framework encompassing several reparameterization schemes from the existing literature. We provided assumptions that guarantee the well-posedness of the dynamical system and convergence of the trajectory. This convergence is both in functional value and in the variable itself, that tends to a minimizer of the loss function avoiding the additional stationary point caused by the reparameterization. In the case where the function is composed by a linear operator, an explicit expression of the implicit bias is given. This comprehensive approach allows us to gain insights into the behavior of the time warped mirror flow, providing a more comprehensive understanding

of its properties.

We also provide a criterion to determine, for a certain function that depends only on the norm, a suitable weight normalization parameterization. Finally, we apply the obtained results to a general multi-neuron fully connected linear network of depth 2, different schemes of weight normalization, and two different extensions of these two previous models.

In the future, we would like to extend this work to models that are more closely related to neural networks and that introduce non-smooth activation functions. Possible extensions are: 1) Extend the results to more general models, such as neural networks; 2) Study new ways to train the overparameterization, for example Stochastic Gradient Descent or second-order methods; 3) Extend this approach to differential inclusions; 4) Study the problem in the matricial case.

APPENDIX A

Appendix for Chapter 6

A.1 Examples

The following three examples have been adapted from [43, Theorem 1], [6, Theorem 1], and [7, Theorem 4.1], respectively. The main objective is to illustrate how to use equation (5.1.3) and its limitations. First, in example A.1.1, we study the case when the reparameterization is separable and find the mirror map is trivial (calculate an integral twice). Second, in example A.1.2, we give a scenario where the derivation of the mirror flow is not straightforward and requires a technical effort. Lastly, in Example 5.3.1, we show that the derivation of Mirror Flow is not possible at all, but it is still possible to find a mirror map F to express the dynamics of (t) as mirror flow multiplied by a positive function G , emphasizing the limitations of equation (5.1.3) justifying our choice of model.

Example A.1.1. Consider a L -layer's overparameterized vector factorization, which can be mathematically expressed as $q(\theta) = \prod_{n=1}^L \mathbb{R}^p$ with $L \geq 2$. We define the flow over as:

$$\dot{q}(t) = -r \nabla L(q(t)) = -L' L'(q(t)) \quad (A.1.1)$$

where the last equality comes from the chain rule.

Note that q satisfies (5.1.3) with F equal to:

$$F: \begin{cases} \prod_{n=1}^L \mathbb{R} & \text{If } L = 2 \\ \prod_{n=1}^L \frac{1}{4} (n \log(n) - n) & \text{If } L > 2; \end{cases} \quad (A.1.1)$$

Indeed, if we compute $J_q(\theta) J_q(\theta)^T$, we get:

$$J_q(\theta) J_q(\theta)^T = q'(\theta)^2 = L^2 (2L - 2) = L^2 (2 - \frac{2}{L}) = r^2 F(\theta)^{-1};$$

Then, (t) follows the dynamic:

$$\dot{q}(t) = -r^2 F(q(t))^{-1} r \nabla L(q(t));$$

The previous flow coincides with the one derived in [43, 135]. Moreover, if $L(\theta) = kX^T yk^2$, for $X \in \mathbb{R}^{d \times p}$ and $y \in \mathbb{R}^d$, then by [77, Theorem 4.17], we get:

$$(t) \rightarrow \theta_1 \text{ and } X \theta_1 = y \quad (\text{Convergence and Feasibility}) \quad (A.1.2)$$

$$r F(\theta_1) = r F(\theta_0) \in \text{ran}(X^T) \quad (\text{Stationarity point}) \quad (A.1.3)$$

which are the KKT conditions of

$$\min_{XZ=y} D_F(z; 0): \tag{A.1.4}$$

For a deeper discussion about the KKT conditions we refer the reader to [113, Section 28].

Example A.1.2. Consider the Two layers vector-vector parameterization; i.e., $\theta = q(\gamma)$ with $\gamma = (\gamma_1; \gamma_2) \in \mathbb{R}^k \times \mathbb{R}^k$ and $q(\gamma) = \gamma_1 \gamma_2$, which is indeed an overparameterization because there are more learnable parameters than variables to predict. Then, the gradient flow on the reparameterization is given by

$$\begin{aligned} \dot{\gamma}_1(t) &= -r \gamma_1 L(\gamma_1(t) \gamma_2(t)) \\ \dot{\gamma}_2(t) &= -r \gamma_2 L(\gamma_1(t) \gamma_2(t)) \end{aligned} \quad \text{with } \gamma_1^2(0) \gamma_2^2(0) = \gamma^2;$$

where $\gamma^2 \in \mathbb{R}^p$. Let us compute the Jacobian of q

$$J_q(\gamma) = [\text{Diag}(\gamma_2) \text{Diag}(\gamma_1)] \text{ and } J_{q(\gamma)}(\gamma) J_{q(\gamma)}^T(\gamma) = \text{Diag}(\gamma_1^2 + \gamma_2^2);$$

Which is not straightforward to express the previous in terms of θ . However, if we multiply the equation of $\dot{\gamma}_1(t)$ component-wise by $\gamma_1(t)$ and we do the same for $\dot{\gamma}_2(t)$ by $\gamma_2(t)$, we get:

$$\dot{\gamma}_1(t) \gamma_1(t) = -r L(\gamma_1(t) \gamma_2(t)) \gamma_1(t) \quad \dot{\gamma}_2(t) \gamma_2(t) = -r L(\gamma_1(t) \gamma_2(t)) \gamma_2(t);$$

From the previous equality and the initialization $\gamma_1^2(0) \gamma_2^2(0) = \gamma^2$, we deduce that, for every $t \geq 0$,

$$\gamma_1^2(t) = \gamma_2^2(t) + \gamma^2; \tag{A.1.5}$$

Multiplying the latter component-wise by $\gamma_2^2(t)$ and recalling that $\gamma_2^2(t) = \gamma_1(t) \gamma_2(t)$, we obtain the following biquadratic equation:

$$\gamma_2^4(t) + \gamma^2 \gamma_2^2(t) - \gamma^2(t) = 0;$$

Since $\gamma_2^2(t) \geq 0$, the previous equation has only one positive solution and from (A.1.5), we have that:

$$\gamma_2^2(t) = \frac{\gamma^2 + \sqrt{\gamma^2(\gamma^2 + 4 \gamma_2^2(t))}}{2} \quad \text{and} \quad \gamma_1^2(t) = \frac{\gamma^2 + \sqrt{\gamma^2(\gamma^2 + 4 \gamma_2^2(t))}}{2};$$

Consider the following entropy function:

$$F(\gamma) := \frac{1}{4} \sum_{n=1}^k 2^n \log \left(2^n + \sqrt{\gamma^2(\gamma^2 + 4 \frac{2^n}{n})} \right);$$

It follows from the definition of the derivative of the product and the computation of the gradient that

$$\begin{aligned} \dot{F}(\gamma) &= \dot{\gamma}_1(t) \gamma_2(t) + \gamma_1(t) \dot{\gamma}_2(t) = \frac{\gamma_1^2 + \gamma_2^2}{2} r L(\gamma_1(t) \gamma_2(t)) \\ &= \frac{\gamma^2}{2} r L(\gamma_1(t) \gamma_2(t)) \\ &= r^2 F(\gamma) - \gamma^2 r L(\gamma_1(t) \gamma_2(t)); \end{aligned}$$

Finally by [77, Theorem 4.17], we have convergence to feasible stationary point (equations (A.1.2) and (A.1.3)), and we can characterize the implicit bias (A.1.4).

Example A.1.3. Consider a multi-neuron fully connected linear network of depth 2; i.e., $\dot{w} = q(w)$ with $\dot{W} = (W; w) \in \mathbb{R}^{p \times m} \times \mathbb{R}^m$ and $q(\cdot) = Ww$. The gradient flow on the reparameterization is then given by

$$\begin{cases} \dot{W}(t) &= -r_W L(W(t)w(t)) = -r L(W(t))w^>(t) \\ \dot{w}(t) &= -r_w L(W(t)w(t)) = -W^>(t)r L(w(t)) \end{cases} \quad (\text{A.1.6})$$

with $W^>(0)W(0) = w(0)w^>(0)$.

Observe that $J_q(W; w)J_q^>(W; w) = kwk^2 \text{Id} + WW^>$ and therefore it is not clear how to express it in terms of \cdot . Proceeding in the same way as in [5], if we multiply the expression of $\dot{W}(t)$ on the left by $W^>(t)$ and the one of $\dot{w}(t)$ on the right by $w^>(t)$, we get that Who derived this expression first

$$W^>(t)\dot{W}(t) = -W^>(t)r L(w(t))w^>(t) = -\dot{w}(t)w^>(t):$$

Adding to the previous equality the one obtained by transposing it, we have that,

$$\begin{aligned} \frac{d}{dt} W^>(t)W(t) &= W^>(t)\dot{W}(t) + \dot{W}^>(t)W(t) \\ &= -\dot{w}(t)w^>(t) + w(t)\dot{w}^>(t) \\ &= -\frac{d}{dt} w(t)w^>(t) : \end{aligned}$$

Then, from the latter equation and the initialization $W^>(0)W(0) = w(0)w^>(0)$, we deduce that, for every $t \geq 0$,

$$W^>(t)W(t) = w(t)w^>(t): \quad (\text{A.1.7})$$

Multiplying by $w(t)$ on the right, by $w^>(t)$ on the left and recalling that $\dot{w}(t) = W(t)w(t)$, we have that $k \dot{w}(t)k^2 = kw(t)k^4$ and so $k \dot{w}(t)k = kw(t)k^2$:

Analogously, multiplying (A.1.7) by $W(t)$ on the left and by $W^>(t)$ on the right, we deduce that $(W(t)W^>(t))^2 = \dot{w}(t)w^>(t)$. Since both $W(t)W^>(t)$ and $\dot{w}(t)w^>(t)$ are symmetric positive semi-definite, equality $(W(t)W^>(t))^2 = \dot{w}(t)w^>(t)$, implies

$$W(t)W^>(t) = \frac{\dot{w}(t)w^>(t)}{k \dot{w}(t)k}:$$

Equation (A.1.6) yields

$$-\dot{w}(t) = W(t)w(t) + W(t)\dot{w}(t) = -r L(w(t))kw(t)k^2 - W(t)W^>(t)r L(w(t)): \quad (\text{A.1.8})$$

Plugging in (A.1.8) the expressions of $kw(t)k^2$ and $W(t)W^>(t)$ in terms of $\dot{w}(t)$, we get

$$-\dot{w}(t) = -k \dot{w}(t)k \text{Id} + \frac{\dot{w}(t)w^>(t)}{k \dot{w}(t)k} -r L(w(t)): \quad (\text{A.1.9})$$

According to Lemma 2.4.5, there does not exist an entropy function F such that the previous dynamical system can be written as a mirror flow, therefore the flow cannot be formulated as vanilla mirror flow and the results in [77, Theorem 4.6] cannot be applied. However, as derived in [7]

$$r^2 \frac{2}{3} k^3 = 2 = k k^{1=2} \text{Id} \frac{w^>}{2k k^2} : \quad (\text{A.1.10})$$

Then, multiplying (A.1.9) on the left by (A.1.10), we obtain the following time-warped dynamical system:

$$r^{-2} \frac{2}{3} k(t) k^{3=2} \dot{x}(t) = k(t) k^{1=2} r^{-1} L(x(t)); \quad (\text{A.1.11})$$

which is a special case of (6.2.1) with $F; G : \mathbb{R}^p \rightarrow \mathbb{R}$ given by $F(x) = \frac{2}{3} k k^{3=2}$ and $G(x) = k k^{1=2}$. which comprises a time-warping factor.

Bibliography

- [1] A. Alacaoglu, O. Fercoq, and V. Cevher, “On the convergence of stochastic primal-dual hybrid gradient,” *SIAM Journal on Optimization*, vol. 32, no. 2, pp. 1288–1318, 2022.
- [2] A. Ali, E. Dobriban, and R. Tibshirani, “The implicit regularization of stochastic gradient flow for least squares,” in *International conference on machine learning*, PMLR, 2020, pp. 233–244.
- [3] F. Alvarez, J. Bolte, and O. Brahic, “Hessian riemannian gradient flows in convex programming,” *SIAM journal on control and optimization*, vol. 43, no. 2, pp. 477–501, 2004.
- [4] E. Amid and M. K. Warmuth, “Reparameterizing mirror descent as gradient descent,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 8430–8439, 2020.
- [5] S. Arora, N. Cohen, W. Hu, and Y. Luo, “Implicit regularization in deep matrix factorization,” in *NeurIPS*, 2019, pp. 7413–7424.
- [6] S. Arora, N. Cohen, and E. Hazan, “On the optimization of deep networks: Implicit acceleration by overparameterization,” in *International Conference on Machine Learning*, PMLR, 2018, pp. 244–253.
- [7] S. Azulay et al., “On the implicit bias of initialization shape: Beyond infinitesimal mirror descent,” in *International Conference on Machine Learning*, PMLR, 2021, pp. 468–477.
- [8] M. Bachmayr and M. Burger, “Iterative total variation schemes for nonlinear inverse problems,” *Inverse Problems*, vol. 25, no. 10, pp. 105004, 26, 2009, ISSN: 0266-5611. DOI: [10.1088/0266-5611/25/10/105004](https://doi.org/10.1088/0266-5611/25/10/105004).
- [9] M. A. Bahraoui and B. Lemaire, “Convergence of diagonally stationary sequences in convex optimization,” in 1-2, vol. 2, *Set convergence in nonlinear analysis and optimization*, Springer, 1994, pp. 49–61. DOI: [10.1007/BF01027092](https://doi.org/10.1007/BF01027092).
- [10] A. B. Bakushinsky and M. Y. Kokurin, *Iterative methods for approximate solution of inverse problems*. Springer Science & Business Media, 2005, vol. 577.

- [11] L. Barreira and C. Valls, Ordinary differential equations: Qualitative theory. American Mathematical Soc., 2012, vol. 137.
- [12] P. L. Bartlett and M. Traskin, “Adaboost is consistent,” the Journal of machine Learning research, vol. 8, pp. 2347–2368, 2007, ISSN: 1532-4435.
- [13] F. Bauer, S. Pereverzev, and L. Rosasco, “On regularization algorithms in learning theory,” Journal of complexity, vol. 23, no. 1, pp. 52–72, 2007.
- [14] H. H. Bauschke, J. M. Borwein, and P. L. Combettes, “Essential smoothness, essential strict convexity, and legendre functions in banach spaces,” Communications in Contemporary Mathematics, vol. 3, no. 04, pp. 615–647, 2001.
- [15] H. H. Bauschke and P. L. Combettes, Convex analysis and monotone operator theory in Hilbert spaces (CMS Books in Mathematics/Ouvrages de Mathématiques de la SMC), Second. Springer, Cham, 2017, pp. xix+619, With a foreword by Hédÿ Attouch, ISBN: 978-3-319-48311-5. DOI: [10.1007/978-3-319-48311-5](https://doi.org/10.1007/978-3-319-48311-5).
- [16] A. Beck and M. Teboulle, “Mirror descent and nonlinear projected subgradient methods for convex optimization,” Oper. Res. Lett., vol. 31, no. 3, pp. 167–175, 2003, ISSN: 0167-6377. DOI: [10.1016/S0167-6377\(02\)00231-6](https://doi.org/10.1016/S0167-6377(02)00231-6).
- [17] S. Becker, J. Bobin, and E. J. Candès, “NESTA: A fast and accurate first-order method for sparse recovery,” SIAM J. Imaging Sci., vol. 4, no. 1, pp. 1–39, 2011. DOI: [10.1137/090756855](https://doi.org/10.1137/090756855).
- [18] M. Benning and M. Burger, “Error estimates for general fidelities,” Electronic Transactions on Numerical Analysis, vol. 38, no. 44-68, p. 77, 2011.
- [19] M. Benning and M. Burger, “Modern regularization methods for inverse problems,” Acta Numer., vol. 27, pp. 1–111, 2018, ISSN: 0962-4929. DOI: [10.1017/S0962492918000016](https://doi.org/10.1017/S0962492918000016).
- [20] G. Blanchard and N. Krämer, “Optimal learning rates for kernel conjugate gradient regression,” in Proceedings of the 23rd International Conference on Neural Information Processing Systems-Volume 1, vol. 23, 2010, pp. 226–234.
- [21] R. I. Boç and T. Hein, “Iterative regularization with a general penalty term—theory and application to L1 and TV regularization,” Inverse Problems, vol. 28, no. 10, p. 104010, 2012, ISSN: 0266-5611. DOI: [10.1088/0266-5611/28/10/104010](https://doi.org/10.1088/0266-5611/28/10/104010).
- [22] L. M. Briceno-Arias, “Forward-douglas-rachford splitting and forward-partial inverse method for solving monotone inclusions,” Optimization, vol. 64, no. 5, pp. 1239–1261, 2015.
- [23] L. M. Briceno-Arias, P. L. Combettes, J.-C. Pesquet, and N. Pustelnik, “Proximal algorithms for multicomponent image recovery problems,” Journal of Mathematical Imaging and Vision, vol. 41, no. 1-2, pp. 3–22, 2011.

- [24] L. Briceño-Arias, J. Deride, and C. Vega, “Random activations in primal-dual splittings for monotone inclusions with a priori information,” *Journal of Optimization Theory and Applications*, pp. 1–26, 2021.
- [25] L. Briceño-Arias and S. López Rivera, “A projected primal-dual method for solving constrained monotone inclusions,” *J. Optim. Theory Appl.*, vol. 180, no. 3, pp. 907–924, 2019, ISSN: 0022-3239. DOI: [10.1007/s10957-018-1430-2](https://doi.org/10.1007/s10957-018-1430-2).
- [26] L. M. Briceño-Arias, “A Douglas-Rachford splitting method for solving equilibrium problems,” *Nonlinear Anal.*, vol. 75, no. 16, pp. 6053–6059, 2012, ISSN: 0362-546X. DOI: [10.1016/j.na.2012.06.014](https://doi.org/10.1016/j.na.2012.06.014). [Online]. Available: <https://doi.org/10.1016/j.na.2012.06.014>.
- [27] M. Burger, E. Resmerita, and L. He, “Error estimation for Bregman iterations and inverse scale space methods in image restoration,” *Computing*, vol. 81, no. 2-3, pp. 109–135, 2007, ISSN: 0010-485X. DOI: [10.1007/s00607-007-0245-z](https://doi.org/10.1007/s00607-007-0245-z).
- [28] J.-F. Cai, E. J. Candès, and Z. Shen, “A singular value thresholding algorithm for matrix completion,” *SIAM J. Optim.*, vol. 20, no. 4, pp. 1956–1982, 2010, ISSN: 1052-6234. DOI: [10.1137/080738970](https://doi.org/10.1137/080738970).
- [29] J.-F. Cai, S. Osher, and Z. Shen, “Linearized Bregman iterations for frame-based image deblurring,” *SIAM J. Imaging Sci.*, vol. 2, no. 1, pp. 226–252, 2009. DOI: [10.1137/080733371](https://doi.org/10.1137/080733371).
- [30] L. Calatroni, G. Garrigos, L. Rosasco, and S. Villa, “Accelerated iterative regularization via dual diagonal descent,” *SIAM J. Optim.*, vol. 31, no. 1, pp. 754–784, 2021, ISSN: 1052-6234. DOI: [10.1137/19M1308888](https://doi.org/10.1137/19M1308888).
- [31] E. Candès, “Matrix completion with noise,” *Proc. IEEE*, vol. 98, no. 6, pp. 925–936, 2010.
- [32] E. J. Candès and T. Tao, “Near-optimal signal recovery from random projections: Universal encoding strategies?” *IEEE Trans. Inform. Theory*, vol. 52, no. 12, pp. 5406–5425, 2006, ISSN: 0018-9448. DOI: [10.1109/TIT.2006.885507](https://doi.org/10.1109/TIT.2006.885507).
- [33] E. J. Candès and B. Recht, “Exact matrix completion via convex optimization,” *Found. Comput. Math.*, vol. 9, no. 6, pp. 717–772, 2009, ISSN: 1615-3375. DOI: [10.1007/s10208-009-9045-5](https://doi.org/10.1007/s10208-009-9045-5).
- [34] E. J. Candès, J. Romberg, and T. Tao, “Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information,” *IEEE Trans. Inform. Theory*, vol. 52, no. 2, pp. 489–509, 2006, ISSN: 0018-9448. DOI: [10.1109/TIT.2005.862083](https://doi.org/10.1109/TIT.2005.862083).
- [35] A. Cauchy et al., “Méthode générale pour la résolution des systemes d’équations simultanées,” *Comp. Rend. Sci. Paris*, vol. 25, no. 1847, pp. 536–538, 1847.

- [36] Y. Censor, P. P. Eggermont, and D. Gordon, “Strong underrelaxation in kaczmarz’s method for inconsistent systems,” *Numerische Mathematik*, vol. 41, pp. 83–92, 1983.
- [37] A. Chambolle, “An algorithm for total variation minimization and applications,” in 1-2, vol. 20, Special issue on mathematics and image analysis, 2004, pp. 89–97. DOI: [10.1023/B:JMI.V.0000011320.81911.38](https://doi.org/10.1023/B:JMI.V.0000011320.81911.38).
- [38] A. Chambolle, M. J. Ehrhardt, P. Richtárik, and C.-B. Schönlieb, “Stochastic primal-dual hybrid gradient algorithm with arbitrary sampling and imaging applications,” *SIAM J. Optim.*, vol. 28, no. 4, pp. 2783–2808, 2018, ISSN: 1052-6234. DOI: [10.1137/17M1134834](https://doi.org/10.1137/17M1134834).
- [39] A. Chambolle and P.-L. Lions, “Image recovery via total variation minimization and related problems,” *Numer. Math.*, vol. 76, no. 2, pp. 167–188, 1997, ISSN: 0029-599X. DOI: [10.1007/s002110050258](https://doi.org/10.1007/s002110050258).
- [40] A. Chambolle and T. Pock, “A first-order primal-dual algorithm for convex problems with applications to imaging,” *J. Math. Imaging Vision*, vol. 40, no. 1, pp. 120–145, 2011, ISSN: 0924-9907. DOI: [10.1007/s10851-010-0251-1](https://doi.org/10.1007/s10851-010-0251-1).
- [41] X. Chen and A. M. Powell, “Almost sure convergence of the kaczmarz algorithm with random measurements,” *Journal of Fourier Analysis and Applications*, vol. 18, no. 6, pp. 1195–1214, 2012.
- [42] L. Chizat and F. Bach, “Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss,” in *Conference on Learning Theory*, 2020, pp. 1305–1338.
- [43] H.-H. Chou, J. Maly, and H. Rauhut, “More is less: Inducing sparsity via overparameterization,” *Information and Inference: A Journal of the IMA*, vol. 12, no. 3, iaad012, 2023.
- [44] H.-H. Chou, H. Rauhut, and R. Ward, “Robust implicit regularization via weight normalization,” *arXiv preprint arXiv:2305.05448*, 2023.
- [45] P. L. Combettes and V. R. Wajs, “Signal recovery by proximal forward-backward splitting,” *Multiscale modeling & simulation*, vol. 4, no. 4, pp. 1168–1200, 2005.
- [46] L. Condat, “A primal-dual splitting method for convex optimization involving Lipschitzian, proximable and linear composite terms,” *J. Optim. Theory Appl.*, vol. 158, no. 2, pp. 460–479, 2013, ISSN: 0022-3239. DOI: [10.1007/s10957-012-0245-9](https://doi.org/10.1007/s10957-012-0245-9).
- [47] C. De Mol, E. De Vito, and L. Rosasco, “Elastic-net regularization in learning theory,” *Journal of Complexity*, vol. 25, no. 2, pp. 201–230, 2009.
- [48] D. L. Donoho, “Compressed sensing,” *IEEE Trans. Inform. Theory*, vol. 52, no. 4, pp. 1289–1306, 2006, ISSN: 0018-9448. DOI: [10.1109/TIT.2006.871582](https://doi.org/10.1109/TIT.2006.871582).

- [49] S. S. Du, W. Hu, and J. D. Lee, “Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced,” *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [50] J. Duchi and Y. Singer, “Efficient online and batch learning using forward backward splitting,” *J. Mach. Learn. Res.*, vol. 10, pp. 2899–2934, 2009, ISSN: 1532-4435.
- [51] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani, “Least angle regression,” *The Annals of statistics*, vol. 32, no. 2, pp. 407–499, 2004.
- [52] H. W. Engl, M. Hanke, and A. Neubauer, *Regularization of inverse problems*. Springer Science & Business Media, 1996, vol. 375.
- [53] S. Foucart, H. Rauhut, S. Foucart, and H. Rauhut, *An invitation to compressive sensing*. Springer, 2013.
- [54] G. Garrigos, L. Rosasco, and S. Villa, “Iterative regularization via dual diagonal descent,” *J. Math. Imaging Vision*, vol. 60, no. 2, pp. 189–215, 2018, ISSN: 0924-9907. DOI: [10.1007/s10851-017-0754-0](https://doi.org/10.1007/s10851-017-0754-0).
- [55] G. Gidel, F. Bach, and S. Lacoste-Julien, “Implicit regularization of discrete gradient dynamics in linear neural networks,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [56] D. Gissin, S. Shalev-Shwartz, and A. Daniely, “The implicit bias of depth: How incremental learning drives generalization,” in *International Conference on Learning Representations*, 2019.
- [57] G. H. Golub, M. Heath, and G. Wahba, “Generalized cross-validation as a method for choosing a good ridge parameter,” *Technometrics*, vol. 21, no. 2, pp. 215–223, 1979.
- [58] S. Gunasekar, J. Lee, D. Soudry, and N. Srebro, “Characterizing implicit bias in terms of optimization geometry,” in *International Conference on Machine Learning*, PMLR, 2018, pp. 1832–1841.
- [59] S. Gunasekar, J. D. Lee, D. Soudry, and N. Srebro, “Implicit bias of gradient descent on linear convolutional networks,” *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [60] S. Gunasekar, B. E. Woodworth, S. Bhojanapalli, B. Neyshabur, and N. Srebro, “Implicit regularization in matrix factorization,” *Advances in Neural Information Processing Systems*, vol. 30, pp. 6151–6159, 2017.
- [61] E. B. Gutiérrez, C. Delplancke, and M. J. Ehrhardt, “Convergence properties of a randomized primal-dual algorithm with applications to parallel mri,” in *International Conference on Scale Space and Variational Methods in Computer Vision*, Springer, New York, 2021, pp. 254–266.

- [62] M. Hanke and W. Niethammer, "On the acceleration of kaczmarz's method for inconsistent linear systems," *Linear Algebra and its Applications*, vol. 130, pp. 83–98, 1990.
- [63] T. Hastie, S. Rosset, R. Tibshirani, and J. Zhu, "The entire regularization path for the support vector machine," *Journal of Machine Learning Research*, vol. 5, no. Oct, pp. 1391–1415, 2004.
- [64] T. Hastie, R. Tibshirani, and J. Friedman, "The elements of statistical learning. springer series in statistics," New York, NY, USA, 2001.
- [65] A. E. Hoerl and R. W. Kennard, "Ridge regression: Biased estimation for non-orthogonal problems," *Technometrics*, vol. 12, no. 1, pp. 55–67, 1970.
- [66] M. Jaggi, "Revisiting frank-wolfe: Projection-free sparse convex optimization," in *International conference on machine learning*, PMLR, 2013, pp. 427–435.
- [67] Z. Ji and M. Telgarsky, "Gradient descent aligns the layers of deep linear networks," in *International Conference on Learning Representations*, 2018.
- [68] Z. Ji and M. Telgarsky, "The implicit bias of gradient descent on nonseparable data," in *Conference on Learning Theory*, PMLR, 2019, pp. 1772–1798.
- [69] B. Jin, D. A. Lorenz, and S. Schiffler, "Elastic-net regularization: Error estimates and active set methods," *Inverse Problems*, vol. 25, no. 11, p. 115 022, 2009.
- [70] S. Kaczmarz, "English translation: S. kaczmarz, approximate solution of systems of linear equations," *Int. J. Contr.*, vol. 57, pp. 355–357, 1937.
- [71] B. Kaltenbacher, A. Neubauer, and O. Scherzer, *Iterative regularization methods for nonlinear ill-posed problems (Radon Series on Computational and Applied Mathematics)*. Walter de Gruyter GmbH & Co. KG, Berlin, 2008, vol. 6, pp. viii+194, ISBN: 978-3-11-020420-9. DOI: [10.1515/9783110208276](https://doi.org/10.1515/9783110208276).
- [72] L. Landweber, "An iteration formula for Fredholm integral equations of the first kind," *Amer. J. Math.*, vol. 73, pp. 615–624, 1951, ISSN: 0002-9327. DOI: [10.2307/2372313](https://doi.org/10.2307/2372313).
- [73] D. Leventhal and A. S. Lewis, "Randomized methods for linear constraints: Convergence rates and conditioning," *Mathematics of Operations Research*, vol. 35, no. 3, pp. 641–654, 2010.
- [74] H. Li, N. Chen, and L. Li, "Error analysis for matrix elastic-net regularization algorithms," *IEEE transactions on neural networks and learning systems*, vol. 23, no. 5, pp. 737–748, 2012.
- [75] Y. Li, T. Ma, and H. Zhang, "Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations," in *Conference On Learning Theory*, PMLR, 2018, pp. 2–47.

- [76] Z. Li, T. Wang, and S. Arora, “What happens after sgd reaches zero loss?—a mathematical framework,” in International Conference on Learning Representations, 2021.
- [77] Z. Li, T. Wang, J. D. Lee, and S. Arora, “Implicit bias of gradient descent on reparametrized models: On equivalence to mirror descent,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 34 626–34 640, 2022.
- [78] D. A. Lorenz, “Convergence rates and source conditions for tikhonov regularization with sparsity constraints,” *Journal of Inverse and Ill-Posed Problems*, vol. 16, no. 5, pp. 463–478, 2008.
- [79] K. Lyu and J. Li, “Gradient descent maximizes the margin of homogeneous neural networks,” in International Conference on Learning Representations, 2019.
- [80] J. Maly, “Robust sensing of low-rank matrices with non-orthogonal sparse decomposition,” *Applied and Computational Harmonic Analysis*, 2023.
- [81] B. Martinet, “Regularisation d’inequations variationelles par approximations successives,” *Revue Francaise d’informatique et de Recherche operationelle*, vol. 4, pp. 154–159, 1970.
- [82] S. Matet, L. Rosasco, S. Villa, and B. L. Vu, “Don’t relax: Early stopping for convex regularization,” *arXiv preprint arXiv:1707.05422*, 2017.
- [83] C. Molinari, J. Liang, and J. Fadili, “Convergence rates of forward–douglas–rachford splitting method,” *Journal of Optimization Theory and Applications*, vol. 182, pp. 606–639, 2019.
- [84] C. Molinari, M. Massias, L. Rosasco, and S. Villa, “Iterative regularization for convex regularizers,” in International Conference on Artificial Intelligence and Statistics, PMLR, 2021, pp. 1684–1692.
- [85] C. Molinari, M. Massias, L. Rosasco, and S. Villa, “Iterative regularization for low complexity regularizers,” *Arxiv 2202.00420*, pp. 1684–1692, 2022.
- [86] C. Molinari and J. Peypouquet, “Lagrangian penalization scheme with parallel forward–backward splitting,” *Journal of Optimization Theory and Applications*, vol. 177, pp. 413–447, 2018.
- [87] C. Molinari, J. Peypouquet, and F. Roldan, “Alternating forward–backward splitting for linearly constrained optimization problems,” *Optimization Letters*, vol. 14, pp. 1071–1088, 2020.
- [88] E. Moroshko, B. E. Woodworth, S. Gunasekar, J. D. Lee, N. Srebro, and D. Soudry, “Implicit bias in deep linear classification: Initialization scale vs training accuracy,” *Advances in neural information processing systems*, vol. 33, pp. 22 182–22 193, 2020.

- [89] D. Morwani and H. G. Ramaswamy, “Inductive bias of gradient descent for weight normalized smooth homogeneous neural nets,” in *International Conference on Algorithmic Learning Theory*, PMLR, 2022, pp. 827–880.
- [90] E. Moulines and F. Bach, “Non-asymptotic analysis of stochastic approximation algorithms for machine learning,” *Advances in neural information processing systems*, vol. 24, pp. 451–459, 2011.
- [91] J. Munkres, “Topology james munkres second edition,”
- [92] M. S. Nacson, J. Lee, S. Gunasekar, P. H. P. Savarese, N. Srebro, and D. Soudry, “Convergence of gradient descent on separable data,” in *The 22nd International Conference on Artificial Intelligence and Statistics*, PMLR, 2019, pp. 3420–3428.
- [93] M. S. Nacson, K. Ravichandran, N. Srebro, and D. Soudry, “Implicit bias of the step size in linear diagonal neural networks,” in *International Conference on Machine Learning*, PMLR, 2022, pp. 16 270–16 295.
- [94] M. S. Nacson, N. Srebro, and D. Soudry, “Stochastic gradient descent on separable data: Exact convergence with a fixed learning rate,” in *The 22nd International Conference on Artificial Intelligence and Statistics*, PMLR, 2019, pp. 3051–3059.
- [95] I. Necoara, “Faster randomized block kaczmarz algorithms,” *SIAM Journal on Matrix Analysis and Applications*, vol. 40, no. 4, pp. 1425–1452, 2019.
- [96] D. Needell, “Randomized kaczmarz solver for noisy linear systems,” *BIT Numerical Mathematics*, vol. 50, pp. 395–403, 2010.
- [97] D. Needell, R. Ward, and N. Srebro, “Stochastic gradient descent, weighted sampling, and the randomized kaczmarz algorithm,” *Advances in neural information processing systems*, vol. 27, 2014.
- [98] D. Needell, R. Zhao, and A. Zouzias, “Randomized block kaczmarz method with projection for solving least squares,” *Linear Algebra and its Applications*, vol. 484, pp. 322–343, 2015.
- [99] A. S. Nemirovski and D. B. Yudin, *Problem complexity and method efficiency in optimization* (A Wiley-Interscience Publication). New York: John Wiley & Sons Inc., 1983, pp. xv+388.
- [100] A. S. Nemirovsky and D. B. a. Yudin, *Problem complexity and method efficiency in optimization* (Wiley-Interscience Series in Discrete Mathematics). John Wiley & Sons, Inc., New York, 1983, pp. xv+388, Translated from the Russian and with a preface by E. R. Dawson, ISBN: 0-471-10345-4.
- [101] A. Neubauer, “On Nesterov acceleration for Landweber iteration of linear ill-posed problems,” *J. Inverse Ill-Posed Probl.*, vol. 25, no. 3, pp. 381–390, 2017, ISSN: 0928-0219. DOI: [10.1515/jip-2016-0060](https://doi.org/10.1515/jip-2016-0060).

- [102] B. Neyshabur, R. Tomioka, and N. Srebro, “In search of the real inductive bias: On the role of implicit regularization in deep learning,” arXiv preprint arXiv:1412.6614, 2014.
- [103] S. Osher, M. Burger, D. Goldfarb, J. Xu, and W. Yin, “An iterative regularization method for total variation-based image restoration,” *Multiscale Model. Simul.*, vol. 4, no. 2, pp. 460–489, 2005, ISSN: 1540-3459. DOI: [10.1137/040605412](https://doi.org/10.1137/040605412).
- [104] S. Osher, Y. Mao, B. Dong, and W. Yin, “Fast linearized Bregman iteration for compressive sensing and sparse denoising,” *Commun. Math. Sci.*, vol. 8, no. 1, pp. 93–111, 2010, ISSN: 1539-6746.
- [105] S. Osher and L. I. Rudin, “Feature-oriented image enhancement using shock filters,” *SIAM Journal on numerical analysis*, vol. 27, no. 4, pp. 919–940, 1990.
- [106] S. Pesme, L. Pillaud-Vivien, and N. Flammarion, “Implicit bias of sgd for diagonal linear networks: A provable benefit of stochasticity,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 29 218–29 230, 2021.
- [107] J. Peypouquet, *Convex optimization in normed spaces (SpringerBriefs in Optimization)*. Springer, Cham, 2015, pp. xiv+124, Theory, methods and examples, With a foreword by Hedy Attouch, ISBN: 978-3-319-13710-0. DOI: [10.1007/978-3-319-13710-0](https://doi.org/10.1007/978-3-319-13710-0).
- [108] G. Peyré, “The numerical tours of signal processing-advanced computational signal and image processing,” *IEEE Computing in Science and Engineering*, vol. 13, no. 4, pp. 94–97, 2011.
- [109] T. Pock and A. Chambolle, “Diagonal preconditioning for first order primal-dual algorithms in convex optimization,” in *2011 International Conference on Computer Vision*, 2011, pp. 1762–1769.
- [110] H. Raguet, J. Fadili, and G. Peyré, “A generalized forward-backward splitting,” *SIAM J. Imaging Sci.*, vol. 6, no. 3, pp. 1199–1226, 2013.
- [111] J. Rasch and A. Chambolle, “Inexact first-order primal-dual algorithms,” *Comput. Optim. Appl.*, vol. 76, no. 2, pp. 381–430, 2020, ISSN: 0926-6003. DOI: [10.1007/s10589-020-00186-y](https://doi.org/10.1007/s10589-020-00186-y).
- [112] G. Raskutti, M. J. Wainwright, and B. Yu, “Early stopping and non-parametric regression: An optimal data-dependent stopping rule,” *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 335–366, 2014.
- [113] R. T. Rockafellar, *Convex analysis (Princeton Landmarks in Mathematics)*. Princeton University Press, Princeton, NJ, 1997, pp. xviii+451, Reprint of the 1970 original, Princeton Paperbacks, ISBN: 0-691-01586-4.
- [114] L. Rosasco and S. Villa, “Learning with incremental iterative regularization,” *NeurIPS*, pp. 1630–1638, 2015.

- [115] M. Rudelson and R. Vershynin, “Geometric approach to error-correcting codes and reconstruction of signals,” *International mathematics research notices*, vol. 2005, no. 64, pp. 4019–4041, 2005, ISSN: 1073-7928. DOI: [10.1155/IMRN.2005.4019](https://doi.org/10.1155/IMRN.2005.4019). [Online]. Available: <https://doi.org/10.1155/IMRN.2005.4019>.
- [116] L. I. Rudin and S. Osher, “Total variation based image restoration with free local constraints,” in *Proceedings of 1st International Conference on Image Processing*, IEEE, vol. 1, 1994, pp. 31–35.
- [117] L. I. Rudin, S. Osher, and E. Fatemi, “Nonlinear total variation based noise removal algorithms,” in 1-4, vol. 60, *Experimental mathematics: computational issues in nonlinear science* (Los Alamos, NM, 1991), 1992, pp. 259–268. DOI: [10.1016/0167-2789\(92\)90242-F](https://doi.org/10.1016/0167-2789(92)90242-F).
- [118] T. Salimans and D. P. Kingma, “Weight normalization: A simple reparameterization to accelerate training of deep neural networks,” *Advances in neural information processing systems*, vol. 29, 2016.
- [119] O. Scherzer, “A modified Landweber iteration for solving parameter estimation problems,” *Appl. Math. Optim.*, vol. 38, no. 1, pp. 45–68, 1998, ISSN: 0095-4616. DOI: [10.1007/s002459900081](https://doi.org/10.1007/s002459900081).
- [120] F. Schöpfer and D. A. Lorenz, “Linear convergence of the randomized sparse kaczmarz method,” *Mathematical Programming*, vol. 173, pp. 509–536, 2019.
- [121] S. Shalev-Shwartz and S. Ben-David, *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [122] A. Silveti-Falls, C. Molinari, and J. Fadili, “A stochastic bregman primal-dual splitting algorithm for composite optimization,” *arXiv preprint arXiv:2112.11928*, 2021.
- [123] A. Silveti-Falls, C. Molinari, and J. Fadili, “Generalized conditional gradient with augmented lagrangian for composite minimization,” *SIAM Journal on Optimization*, vol. 30, no. 4, pp. 2687–2725, 2020.
- [124] A. Silveti-Falls, C. Molinari, and J. Fadili, “Inexact and stochastic generalized conditional gradient with augmented lagrangian and proximal step,” *Journal of Non-smooth Analysis and Optimization*, vol. 2, no. Original research articles, 2021.
- [125] D. Soudry, E. Hoffer, M. S. Nacson, S. Gunasekar, and N. Srebro, “The implicit bias of gradient descent on separable data,” *The Journal of Machine Learning Research*, vol. 19, no. 1, pp. 2822–2878, 2018.
- [126] I. Steinwart and A. Christmann, *Support vector machines (Information Science and Statistics)*. Springer, New York, 2008, pp. xvi+601, ISBN: 978-0-387-77241-7.

- [127] T. Strohmer and R. Vershynin, “A randomized kaczmarz algorithm with exponential convergence,” *Journal of Fourier Analysis and Applications*, vol. 15, no. 2, pp. 262–278, 2009.
- [128] A. N. Tihonov, “Solution of incorrectly formulated problems and the regularization method,” *Soviet Math.*, vol. 4, pp. 1035–1038, 1963.
- [129] Y. Tsaig and D. L. Donoho, “Extensions of compressed sensing,” *Signal processing*, vol. 86, no. 3, pp. 549–571, 2006.
- [130] G. Vardi and O. Shamir, “Implicit regularization in relu networks with the square loss,” in *Conference on Learning Theory*, PMLR, 2021, pp. 4224–4258.
- [131] T. Vaskevicius, V. Kanade, and P. Rebeschini, “Implicit regularization for optimal sparse recovery,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [132] C. Vega, C. Molinari, L. Rosasco, and S. Villa, “Fast iterative regularization by reusing data,” *Journal of Inverse and Ill-posed Problems*, no. 0, 2023.
- [133] B. C. Vũ, “A splitting algorithm for dual monotone inclusions involving cocoercive operators,” *Adv. Comput. Math.*, vol. 38, no. 3, pp. 667–681, 2013, ISSN: 1019-7168. DOI: [10.1007/s10444-011-9254-8](https://doi.org/10.1007/s10444-011-9254-8).
- [134] Y. Wang, M. Chen, T. Zhao, and M. Tao, “Large learning rate tames homogeneity: Convergence and balancing effect,” in *The International Conference on Learning Representations*, 2022.
- [135] B. Woodworth et al., “Kernel and rich regimes in overparametrized models,” in *Conference on Learning Theory*, PMLR, 2020, pp. 3635–3673.
- [136] F. Wu and P. Rebeschini, “A continuous-time mirror descent approach to sparse phase retrieval,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 20 192–20 203, 2020.
- [137] F. Wu and P. Rebeschini, “Implicit regularization in matrix sensing via mirror descent,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 20 558–20 570, 2021.
- [138] L. Wu, C. Ma, et al., “How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective,” *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [139] X. Wu et al., “Implicit regularization and convergence for weight normalization,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 2835–2847, 2020.
- [140] L. Xiao, “Dual averaging methods for regularized stochastic learning and online optimization,” *J. Mach. Learn. Res.*, vol. 11, pp. 2543–2596, 2010, ISSN: 1532-4435.

- [141] Y. Yao, L. Rosasco, and A. Caponnetto, "On early stopping in gradient descent learning," *Constr. Approx.*, vol. 26, no. 2, pp. 289–315, 2007, ISSN: 0176-4276. DOI: [10.1007/s00365-006-0663-2](https://doi.org/10.1007/s00365-006-0663-2).
- [142] W. Yin, "Analysis and generalizations of the linearized bregman method," *SIAM Journal on Imaging Sciences*, vol. 3, no. 4, pp. 856–877, 2010.
- [143] W. Yin, S. Osher, D. Goldfarb, and J. Darbon, "Bregman iterative algorithms for ℓ_1 -minimization with applications to compressed sensing," *SIAM Journal on Imaging sciences*, vol. 1, no. 1, pp. 143–168, 2008.
- [144] C. You, Z. Zhu, Q. Qu, and Y. Ma, "Robust recovery via implicit bias of discrepant learning rates for double over-parameterization," *Advances in Neural Information Processing Systems*, vol. 33, pp. 17 733–17 744, 2020.
- [145] C. Yun, S. Krishnan, and H. Mobahi, "A unifying view on implicit bias in training linear neural networks," in *International Conference on Learning Representations*, 2020.
- [146] T. Zhang and B. Yu, "Boosting with early stopping: Convergence and consistency," *Ann. Statist.*, vol. 33, no. 4, pp. 1538–1579, 2005, ISSN: 0090-5364. DOI: [10.1214/009053605000000255](https://doi.org/10.1214/009053605000000255).
- [147] P. Zhao, Y. Yang, and Q.-C. He, "Implicit regularization via hadamard product over-parametrization in high-dimensional linear regression," *arXiv preprint arXiv:1903.09367*, 2019.
- [148] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *Journal of the royal statistical society: series B (statistical methodology)*, vol. 67, no. 2, pp. 301–320, 2005.
- [149] H. Zou and H. H. Zhang, "On the adaptive elastic-net with a diverging number of parameters," *Annals of statistics*, vol. 37, no. 4, p. 1733, 2009.
- [150] A. Zouzias and N. M. Freris, "Randomized extended kaczmarz for solving least squares," *SIAM Journal on Matrix Analysis and Applications*, vol. 34, no. 2, pp. 773–793, 2013.