

# **Consistency, Robustness and Sparsity for Learning Algorithms**

Konsistenz, Robustheit und Dünnbesetztheit von Lern-Algorithmen

Der Naturwissenschaftlichen Fakultät  
der  
Friedrich-Alexander-Universität Erlangen-Nürnberg  
  
zur Erlangung des Doktorgrades  
**Dr. rer. nat.**

vorgelegt von  
**Tim Roith**  
aus  
**Amberg**

Als Dissertation genehmigt von der Naturwissenschaftlichen Fakultät der  
Friedrich-Alexander-Universität Erlangen-Nürnberg

Tag der mündlichen Prüfung: —  
Vorsitzender des Promotionsorgans: —  
Gutachter\*in: Martin Burger  
Dejan Slepčev  
Franca Hoffmann

## **Acknowledgement**

Coming soon...

# Contents

Preface	vii
<b>I. Exposition</b>	<b>1</b>
1. Introduction	2
2. Learning Paradigms	5
2.1. Unsupervised Learning	6
2.2. Supervised Learning	6
2.3. Semi-Supervised Learning	7
3. Consistent Semi-Supervised Learning on Sparse Graphs	9
3.1. Graph-Based SSL and Consistency	10
3.1.1. Weighted Graphs	10
3.1.2. The $p$ -Laplacian: Continuum and Graph	13
3.1.3. Consistency for Graph-based SSL	18
3.2. Lipschitz Extensions and the Infinity Laplacian: Continuum and Graph	21
3.2.1. The Continuum Setting	21
3.2.2. Graph Lipschitz Extensions	31
3.3. Gamma Convergence: [LIP-I]	37
3.3.1. Setting and Preliminaries	38
3.3.2. $\Gamma$ -Convergence of the Discrete Functionals	41
3.3.3. Convergence of Minimizers	44
3.3.4. Application to Ground States	45
3.4. Uniform Convergence of AMLES: [LIP-II]	45
3.4.1. Setting	46
3.4.2. Convergence Results	47
3.4.3. Numerical Examples and Extensions	54
4. Robust and Sparse Supervised Learning	56
4.1. Setting	57
4.1.1. Network Architectures	58
4.1.2. Gradient Computation and Stochastic Gradient Descent	59
4.2. Adversarial Stability via Lipschitz Training: [CLIP]	60
4.2.1. Cheap Lipschitz Training	64
4.2.2. Analysis of Lipschitz Regularization	65
4.2.3. Numerical Results	67
4.3. Resolution Stability via FNOs: [FNO]	68
4.3.1. Fourier Neural Operators	70

## *Contents*

4.3.2.	Analytical Results for FNOs . . . . .	73
4.3.3.	Numerical Results . . . . .	76
4.4.	Sparsity via Bregman Iterations: [BREG-I] . . . . .	78
4.4.1.	Preliminaries on Convex Analysis and Bregman Iterations . .	79
4.4.2.	Linearized Bregman Iterations and Mirror Descent . . . . .	86
4.4.3.	Stochastic and Momentum Variants . . . . .	87
4.4.4.	Convergence of Stochastic Bregman Iterations . . . . .	88
4.4.5.	Numerical Results and Practical Considerations . . . . .	92

# List of Figures

3.1.	Dependence of the number of non-zero edges on the scaling parameter $\varepsilon$ . . . . .	13
3.2.	Solution to the Laplacian Learning problem for different number of data points. . . . .	19
3.3.	Solution to the Laplacian Learning problem for different number of data points. . . . .	21
3.4.	The domain in Example 3.16. . . . .	24
3.5.	Visualization for Example 3.20. . . . .	26
3.6.	Visualization of the relative boundary. . . . .	27
3.7.	Visualization of exterior and interior boundary on a graph. . . . .	35
3.8.	An example of a sharp internal corner, violating Eq. (3.20). . . . .	42
4.1.	Effects of applying convolutional filters with different resolutions. . . . .	73
4.2.	Visualization of the Bregman distance. . . . .	81
4.3.	Bregman iterations for image denoising in Example 4.21. . . . .	85

# Preface

This work is structured into two main parts, **Part I** the presentation and explanation of the topics and results presented in ??, the peer-reviewed articles.

Part I: Exposition	??: ??
Chapter 2: Learning Paradigms	—
Chapter 3: Consistent Semi-Supervised Learning on Sparse Graphs	????
Chapter 4: Robust and Sparse Supervised Learning	??????

**Part I** consists of five chapters, of which the first two give an introduction and explain the paradigms, *unsupervised*, *semi-supervised* and *supervised* learning. The next two chapters are split up thematically, concerning the topics semi-supervised and supervised learning respectively. Here, a short overview provides the necessary framework allowing us to explain the main contributions. The last chapter presents the conclusion. In ?? the following publications are reprinted:

- [LIP-I] T. Roith and L. Bungert. “Continuum limit of Lipschitz learning on graphs.” In: *Foundations of Computational Mathematics* (2022), pp. 1–39.
- [LIP-II] L. Bungert, J. Calder, and T. Roith. “Uniform convergence rates for Lipschitz learning on graphs.” In: *IMA Journal of Numerical Analysis* (Sept. 2022).
- [CLIP] L. Bungert, R. Raab, T. Roith, L. Schwinn, and D. Tenbrinck. “CLIP: Cheap Lipschitz training of neural networks.” In: *Scale Space and Variational Methods in Computer Vision: 8th International Conference, SSVM 2021, Proceedings*. Springer. 2021, pp. 307–319.
- [BREG-I] L. Bungert, T. Roith, D. Tenbrinck, and M. Burger. “A Bregman learning framework for sparse neural networks.” In: *Journal of Machine Learning Research* 23.192 (2022), pp. 1–43.
- [FNO] S. Kabri, T. Roith, D. Tenbrinck, and M. Burger. “Resolution-Invariant Image Classification based on Fourier Neural Operators.” In: *Scale Space and Variational Methods in Computer Vision: 9th International Conference, SSVM 2023, Proceedings*. Springer. 2023, pp. 307–319.

The following two works that are not part of this thesis but provide an additional insight.

- [LIP-III] L. Bungert, J. Calder, and T. Roith. *Ratio convergence rates for Euclidean first-passage percolation: Applications to the graph infinity Laplacian*. 2022. arXiv: [2210.09023 \[math.PR\]](#).
- [BREG-II] L. Bungert, T. Roith, D. Tenbrinck, and M. Burger. “Neural Architecture Search via Bregman Iterations.” In: (2021). arXiv: [2106.02479 \[cs.LG\]](#).

## TR’s Contribution

Here we list TR’s contribution to the publications included in the thesis.

**[LIP-I]:** This work builds upon the findings in TR’s masters thesis [Roi21]. It is however important to note that the results constitute a significant extension and are conceptually stronger than the ones in [Roi21], see [Section 3.3](#). TR adapted the continuum limit framework to the  $L^\infty$  case, worked out most of the proofs and wrote a significant part of the paper. In collaboration with LB, he identified the crucial domain assumptions that allow to work on non-convex domains and proved convergence for approximate boundary conditions.

**[LIP-II]:** In collaboration with LB, TR worked on the convergence proofs building upon the ideas of JC. Together with LB and JC he proved the main convergence result and the various lemmas leading up to it. Here, he was especially concerned with the adaptation of the theory for AMLEs to the graph case, which is a crucial element for the whole work. Furthermore, he contributed to the design and implementation of the numerical examples conducted in the paper.

**[CLIP]:** TR worked out the main algorithm proposed in the paper together with LB, based on LB’s idea. Together with LS, RR and DT he conducted the numerical examples and also wrote large parts of the source code. Furthermore, he wrote significant parts of the paper, where DT proofread and clarified the final document.

**[BREG-I]:** TR expanded LB’s ideas of employing Bregman iteration for sparse training, conceptualized by DT. Together with MB and LB he worked out the convergence analysis of stochastic Bregman iterations. Here, he also proposed a profound sparse initialization strategy. Furthermore, he conducted the numerical examples and wrote most of the source code.

**[FNO]:** This work is based on SK’s masters thesis, employing the initial ideas of MB for resolution invariance with FNOs. In the paper TR worked out the proofs for well-definedness and Fréchet-differentiability, together with SK. He wrote large parts of the paper and the source code, where DT helped with proofreading of the published version. Here, he conducted the numerical studies in collaboration with SK.

# **Part I.**

# **Exposition**

# Chapter 1

## Introduction

The field of *machine learning* emerged in the 1950s [Sam59; Ros58], motivated by the idea of letting a computer discover algorithms and patterns without having to explicitly arrange them by hand. After the initial phase and multiple “AI-winters” [SG96], numerous important developments—e.g. the rediscovery of the backpropagation algorithm, originally due to [Kel60; Ros+62] and then popularized in [RHW86], see e.g. [Sch22]—contributed to the relevance of learning methods. The advances in computer hardware together with the availability of large amounts of data, finally allowed the machine learning enthusiasm of the recent years to spark. While “deep” learning methods—i.e. techniques involving many stacked neural layers as originally proposed in [Ros58]—are the most prominent examples, there is a whole zoo of learning-based strategies that are actively applied in fields like computer vision [Cha+21], natural language processing [Khu+23] or healthcare [She+22]. In this work we mainly focus on data-driven approaches, applied to classification tasks, where the concrete modality of the given data determines our approach. Namely, we focus on supervised—the dataset consists only of input-output pairs, i.e. is fully labeled—and semi-supervised—the data is only partially labeled—learning tasks.

For both regimes especially the last 20 years have seen great success of these data-driven methods. However, the sometimes purely heuristic learning strategies also exhibit serious drawbacks. In the supervised setting one is usually interested in the generalization behavior of a learned classifier, i.e. how good is the performance on unseen inputs which are not part of the given training data. Unfortunately in [GSS14] it was discovered, that this performance can be completely corrupted, by small, seemingly invisible perturbations known as *adversarial attacks*. More generally this phenomenon leads us to the issue of *input robustness*. Given some input  $x$ , suppose that a human and some machine would classify this input to be of type  $c$ . In a rather vague but demonstrative formulation, the key implication we want to obtain for an input  $\bar{x}$  is

$$\left. \begin{array}{l} \bar{x} \text{ is close to } x, \\ \bar{x} \text{ is still classified as } c \text{ by a human} \end{array} \right\} \Rightarrow \text{the machine classifies } \bar{x} \text{ as } c.$$

Next to adversarial examples this also includes resolution changes of images, which do not change the classification by a human, if they are reasonably small. In any case, the existence of these perturbations exhibit critical flaws of learning methods and call for a better theoretical understanding of the employed models. This is where the mathematical foundation of the field becomes more relevant and properties apart from the classification performance come into play, which are discussed within this thesis.

For the semi-supervised setting we consider graph-based algorithms as originally proposed in [ZGL03] with the graph Laplacian. The main problem we highlight in this thesis was first observed in [NSZ09], namely that the classification performance deters significantly with increasing dimensionality of the data. In fact it turned out that solutions obtained by the standard graph Laplacian tend to be constant over the whole dataset, whenever the dimension is larger than two, which can be related to the Sobolev embedding theorem [AF03]. This issue is prevalent in the infinite data-limit, where a priori we consider the case, when the amount of unlabeled data points goes to infinity, which leads us to the question of *consistency* for semi-supervised algorithms.

An issue that is shared across the supervised and semi-supervised setting is the high demand for computational resources. Training a neural network usually involves the use of GPUs for long amounts of time. On the one hand this makes the process infeasible for less powerful machines or even mobile devices and on the other hand generates questionable amounts of CO<sub>2</sub> emissions [Hoe+21]. For graph-based semi-supervised learning one first needs to compute distances between many different data points, to obtain edge weights, which itself is a costly task. Furthermore, the computational complexity of various tasks on a given graph, scales with the number of edges. For example, the run time for Dijkstra's algorithm to compute shortest paths on a graph, already scales linearly with the amount of edges [Dij22]. In this thesis, the keyword to reduce the computational load in both cases, is *sparsity*. The concept of sparse matrices routes deeply into the field of numerical linear algebra [Lan52; GV13] and basically consists of exploiting zeros in a matrix to speed up the computation time. For neural networks this can be incorporated by enforcing the weight matrices of the layers to be sparse. For graphs, sparsity of the connectivity matrix simply means that we have only a small amount of active edges, which also reduces computational cost.

**Contributions in This Work** Taking up the previously mentioned subjects this thesis is concerned with *consistency*, *robustness* and *sparsity* of supervised and semi-supervised learning algorithms.

For the latter we mainly consider the so-called Lipschitz learning task [NSZ09] for which we prove convergence and convergence rates for discrete solutions to their continuum counterpart in the infinite data limit. Here, we always work in a framework that allows for very sparse and therefore computationally feasible graphs.

In the supervised regime we deal with input-robustness w.r.t. adversarial attacks and resolution changes. In the first case we propose an efficient algorithm, penalizing the Lipschitz constant [Lip77] of a neural network, which trains an adversarially robust network. In the multi-resolution setting we analyze the role of Fourier neural operators as proposed in [Li+20] and their connection to standard convolutional neural layers [Fuk80]. Concerning the computational complexity of neural network training, we propose an algorithm based on Bregman iterations [Osh+05] that allows for sparse weight matrices throughout the training. We also provide the convergence analysis for the stochastic adaption of the original Bregman iterations.

**Structure of The Exposition** In [Chapter 2](#) we introduce the learning paradigms and basic notions used throughout this thesis. We then present the topics on consistency for semi-supervised learning on graphs in [Chapter 3](#). After an explanatory introduction we then highlight the main contributions of [[LIP-I](#); [LIP-II](#)]. Here, we try to have as little redundancy to the prints in ?? as possible, while still allowing for an understandable context. In [Chapter 4](#) we then comment on the supervised part of this thesis. After an additional introduction the chapter contains three section presenting the works [[FNO](#); [CLIP](#); [BREG-I](#)] individually. Finally, in ?? we summarize the contents of the whole thesis and provide possible future directions.

# Chapter 2

## Learning Paradigms

Throughout this thesis, we assume to be given data  $\mathcal{X}_n \subset \mathcal{X} \subset \mathbb{R}^d$  consisting of  $n$  data points. We consider the task of *learning* a function  $f : \tilde{\mathcal{X}} \rightarrow \mathcal{Y}$  from the given data, where  $\mathcal{Y}$  denotes the output space. In our case the set  $\tilde{\mathcal{X}} \subset \mathcal{X}$  is usually chosen either as the set of data points  $\mathcal{X}_n$  or as the whole space  $\mathcal{X}$ . The two most important cases for us are listed below.

- **Classification:** The function  $f$  assigns a label to each  $x \in \tilde{\mathcal{X}}$  out of  $C \in \mathbb{N}$  possible classes, i.e.  $\mathcal{Y} = \{1, \dots, C\}$ . In some architectures the last layer of the neural network is given as a vector  $y \in \mathbb{R}^C$ . Typically, this vector is a probability vector, i.e.

$$y \in \Delta^C := \left\{ z \in [0, 1]^C : \sum_{i=1}^C z_i = 1 \right\}.$$

This can be enforced via the softmax function [Bri90]  $\text{softmax} : \mathbb{R}^C \rightarrow \mathbb{R}^C$

$$\text{softmax}(z)_i := \frac{\exp(z_i)}{\sum_{j=1}^C \exp(z_j)}$$

which was actually introduced by Boltzman in [Bol68]. This allows the interpretation that the  $i$ th entry of  $f_\theta(x) \in \Delta^C$  models the probability that  $x$  belongs to class  $i$ . In order to obtain a label one can simply choose the maximum entry, i.e.  $\text{argmax}_{i=1, \dots, C} f_\theta(x)_i$ .

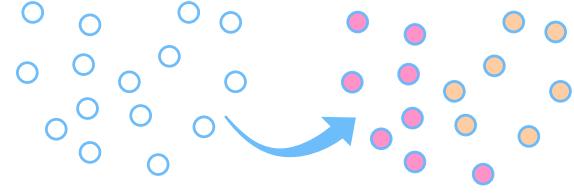
- **Image denoising:** The function  $f$  outputs a denoised version of an input image. Here we have  $\mathcal{X} = \mathcal{Y} = \mathbb{R}^{K \times N \times M}$ , where
  - $K \in \mathbb{N}$  is the number of color channels,
  - $N, M$  denote the width and height of the image.

The learning paradigms we consider in this thesis, differ by their usage of labeled data. We review the concepts in the following.

## 2.1. Unsupervised Learning

In the case of unsupervised learning we are not given any labeled data. In our context the most important application is data clustering. Other tasks involve dimensionality reduction or density estimation, see [ST14]. The clustering task consists of grouping data based on some similarity criterion. In this sense, clustering can also be interpreted as classification, i.e., the desired function is a mapping  $f : \tilde{\mathcal{X}} \rightarrow \{1, \dots, C\}$  where  $C \in \mathbb{N}$  denotes the number of clusters. Typically, one wants to obtain a clustering of the given data set, i.e.,  $\tilde{\mathcal{X}} = \mathcal{X}_n$ . We list some of the typical clustering methods below:

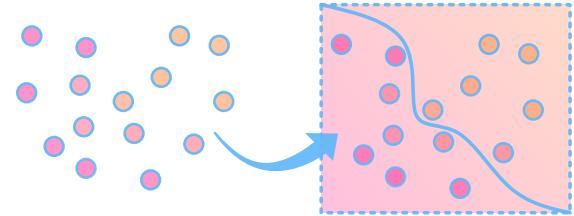
- K-means algorithm [Ste+56],
- expectation maximization [DLR77],
- Cheeger cuts [GS15; SB09; Gar+16; GMT22],
- spectral clustering [GS18; THH21; Hof+22].



Unsupervised learning is not the main focus of this present work. However, we note that especially the concepts developed in [GS15] for Cheeger cuts are crucial for the continuum limit framework in [Section 3.3](#).

## 2.2. Supervised Learning

In this setting, each data point  $x \in \mathcal{X}_n$  is labeled, via a given function  $g : \mathcal{X}_n \rightarrow \mathcal{Y}$  such that we have a finite training set  $\mathcal{T} = \{(x, g(x)) : x \in \mathcal{X}_n\}$ . The task is to infer a function defined on the underlying space,  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , i.e., we want to assign a label to unseen inputs  $x \in \mathcal{X}$  that are not necessarily part of the given data. Often, one models the problem via a joint probability function  $P_{\mathcal{X}, \mathcal{Y}}$  and assumes that the training data are independent and identically distributed (i.i.d) with respect to  $P_{\mathcal{X}, \mathcal{Y}}$ . In this interpretation, the aim is to approximate the conditional  $P_{\mathcal{X}, \mathcal{Y}}(y|x)$  for an input  $x \in \mathcal{X}$  and output  $y \in \mathcal{Y}$ .



In order to *learn* the function  $f$  from the given data, one needs to choose a parameterized class of functions, where typically each element can be described by a finite number of parameters. Among others, common methods or parametrizations include

- support vector machines [CV95; SS05],

- decision trees [MS63; Bre+84],
- neural networks [Tur04; Ros58; MP69].

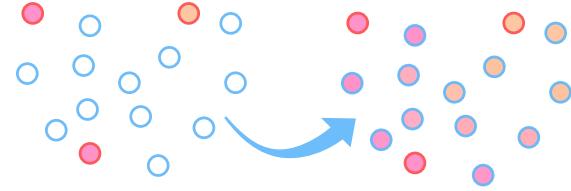
We refer to [Sch15] for an exhaustive historical overview.

In [Chapter 4](#) we exclusively focus on supervised learning algorithms employing neural networks, where the concrete setting is given at the start of the chapter.

### 2.3. Semi-Supervised Learning

In the semi-supervised setting we assume that only a fraction of the data  $\mathcal{X}_n$  is labeled, i.e., we are given a function  $g : \mathcal{O}_n \rightarrow \mathcal{Y}$  where  $\mathcal{O}_n \subset \mathcal{X}_n$  is the non-empty set of labeled data. Typically this constitutes only a small fraction of all available points, i.e.  $|\mathcal{O}_n| \ll |\mathcal{X}_n|$ . In this thesis we restrict ourselves to the *transductive setting*, i.e. we want to infer a function acting only on the data  $f : \mathcal{X}_n \rightarrow \mathcal{Y}$ . This is opposed to the inductive setting, where  $f$  also classifies unseen points  $x \in \mathcal{X}$ , [Zhu05]. Common algorithms and methods include

- expectation maximization and mixture models [DLR77; CCC+03],
- self-training and co-training [BM98],
- graph-based learning [Zhu05].



Mostly, we consider the extension task with  $\mathcal{Y}$  being chosen as  $\mathbb{R}$ . In application this can be seen as a binary classification task, where for  $o \in \mathcal{O}_n$  we have  $g(o) = 1$  if  $o$  belongs to a some class and  $g(o) = 0$  otherwise. The function  $f : \mathcal{X}_n \rightarrow \mathbb{R}$  then determines the probability that any vertex  $x \in \mathcal{X}_n$  belongs to this class, where we can binarize the output via some thresholding, e.g.,

$$x \text{ belongs to the class} \Leftrightarrow f(x) > 1/2.$$

This methodology can be extended to classification tasks beyond the binary case, via the so-called one-vs-all technique [ZGL03]. Given a classification problem with  $C \in \mathbb{N}$  possible classes, we assume that the labeling function  $g : \mathcal{O}_n \rightarrow \Delta^C$  outputs one-hot vectors, i.e.  $g(o)_c = 1$  if  $o$  belongs to class  $c$  and  $g(o)_c = 0$  otherwise, for every  $c = 1, \dots, C$ . We then perform the binary classification problem “ $x$  belongs to class  $c$ ” for every  $c = 1, \dots, C$ , by considering the extension task of

$$g_c : \mathcal{O}_n \rightarrow \mathbb{R} \quad g_c(o) = g(o)_c,$$

which yield functions  $f_c : \mathcal{X}_n \rightarrow [0, 1], c = 1, \dots, C$ . The final output can then either obtained by employing a Bayes classifier [DGL13], i.e.  $f : \mathcal{X}_n \rightarrow \{1, \dots, C\}$

$$f(x) := \operatorname{argmax}_{c=1, \dots, C} f_c(x)$$

or by applying a softmax to obtain a probability vector, i.e.  $f : \mathcal{X}_n \rightarrow \Delta^C$

$$f(x) := \text{softmax}(f_1(x), \dots, f_c(x)).$$

In [Chapter 3](#) we focus on graph-based learning algorithms, however we refer to [\[Zhu05\]](#) for an overview of semi-supervised learning algorithms.

# Chapter 2

## Consistent Semi-Supervised Learning on Sparse Graphs

This chapter considers the consistency of graph-based semi-supervised learning methods and contextualizes the topics provided in the prints [LIP-I; LIP-II]. We are given partially labeled data, where the task is to obtain a labeling on the whole graph. Considering learning algorithms we are especially interested in the following aspects:

- **Consistency:** what is the behavior in the infinite data limit?
- **Sparsity:** how does sparsity of the graph weights affect this behavior?

Section 3.1: Graph-Based SSL and Consistency

Section 3.2 Lipschitz Extensions and the Infinity Laplacian

Section 3.3: [LIP-I]

Section 3.4: [LIP-II]

In [LIP-I] we first consider  $\Gamma$ -convergence of discrete  $L^\infty$  functionals to their continuum counterpart. This yields convergence of minimizers and therefore consistency. The proofs allow for very sparse graphs, however they only directly apply to the Lipschitz extension task, which does not admit unique solutions. We refer to Section 3.3 for more details.

The issue of non-uniqueness is addressed in [LIP-II], which considers uniform convergence of graph infinity harmonic functions to their continuum counterpart. These functions are in fact special solutions of the extension task in [LIP-I], namely so-called absolutely minimizing Lipschitz extensions, which are—under certain assumptions—unique. Again, we are able to work with a very mild scaling assumption that allows for sparse graphs. Furthermore, we are able to prove the first quantitative convergence rates, where we provide the key insight, that a rate for graph distance functions transfers to a rate for infinity harmonic functions. The main contributions are presented in Section 3.4.

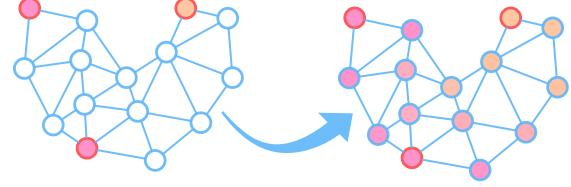
Before we present these works we first introduce the setting of consistency for graph-based SSL in Section 3.1, where we also discuss the  $p$ -Laplacian both in the continuum

and the graph setting. Similarly, we then introduce the problem we are actually concerned with, namely the  $\infty$ -Laplacian and Lipschitz extensions in [Section 3.2](#).

### 3.1. Graph-Based SSL and Consistency

In the semi-supervised learning setting of [Section 2.3](#), we are given a finite set  $\Omega_n \subset \tilde{\Omega}$  consisting of  $n$  points, where  $\tilde{\Omega}$  denotes the input space. We assume that a non-empty subset  $\mathcal{O}_n \subset \Omega_n$  is labeled, i.e., we are given a labeling function  $\mathbf{g} : \mathcal{O}_n \rightarrow \mathbb{R}$ . From now on we denote functions acting on a discrete set  $\Omega_n$  using bold symbols to distinguish them from functions acting on the continuum  $\tilde{\Omega}$ . In particular, the space of functions on  $\Omega_n$  can simply be identified with  $\mathbb{R}^n$ , since  $\{\mathbf{u} : \Omega_n \rightarrow \mathbb{R}\} \sim \mathbb{R}^n$ . The semi-supervised learning problem consists in finding an extension of  $\mathbf{g}$  from  $\mathcal{O}_n$  to the whole set  $\Omega_n$ , i.e.

$$\begin{aligned} & \text{find } \mathbf{u} : \Omega_n \rightarrow \mathbb{R}, \\ & \text{such that } \mathbf{u}(x) = \mathbf{g}(x) \text{ for all } x \in \mathcal{O}_n. \end{aligned} \tag{SSL}$$



In order to obtain meaningful solutions, one usually incorporates the *smoothness assumption* [[ST14](#)] which can be informally stated as follows:

“Points that are close together are more likely to share a similar label.”

**Remark 3.1 (Notational Remark).** In [Chapter 2](#) we employ the symbol  $\mathcal{X}$  to denote the input set, which is now replaced by  $\Omega$ . This is due to the fact, that [[LIP-I](#); [LIP-II](#); [LIP-III](#)] employ this notation for which we choose to use it here.  $\triangle$

Often, one assumes that the points  $x \in \Omega_n$  are i.i.d. w.r.t. a distribution  $\mu$  with density  $\rho : \Omega \rightarrow \mathbb{R}$ . However, as observed in [[Roi21](#)], the data distribution is not directly relevant for the works in [[LIP-I](#); [LIP-II](#)] and we can instead consider a more general class of point clouds. The appropriate generalization to study the limit of  $\Omega_n$  for  $n \rightarrow \infty$  is discussed in [Section 3.3](#).

#### 3.1.1. Weighted Graphs

In order to employ the smoothness assumption we need a notion of “closeness” on the set  $\Omega_n$ , for which we introduce weighted graphs. Namely we define a weighting function  $w$  that allows to compare points in  $\Omega_n$  and therefore induces the desired notion.

**Definition 3.2 (Weighted Graphs).** For a finite set  $\Omega_n$  and a weight function  $w_n : \Omega_n \times \Omega_n \rightarrow \mathbb{R}$ , the tuple  $(\Omega_n, w_n)$  is called a *weighted graph*.

**Other Notions of Graphs** Typically, a graph is defined as a pair  $(\Omega_n, E)$  where  $E$  denotes the set of edges. Here, one has two cases:

- *Undirected graph:*  $E = \{\{x, y\} : \text{there is an edge between } x \in \Omega_n \text{ and } y \in \Omega_n\}$ , i.e.,  $E \subset 2^{\Omega_n}$  and each edge is undirected, since  $\{x, y\} = \{y, x\}$ .
- *Directed graph:*  $E = \{(x, y) : \text{there is an edge from } x \text{ to } y\}$ , i.e.,  $E \subset \Omega_n \times \Omega_n$  and each edge is directed and  $(x, y) \neq (y, x)$ .

Additionally, one then considers a weight function  $W : E \rightarrow \mathbb{R}$  assigning a weight to each edge, and then defines the triple  $(\Omega_n, E, W)$  as a weighted graph. However, we note that all this information can be represented much more elegantly by a weight function  $w : \Omega_n \times \Omega_n \rightarrow \mathbb{R}$ . A directed edge set  $E \subset \Omega_n \times \Omega_n$  can be equivalently expressed by a weight function  $w : \Omega_n \times \Omega_n \rightarrow \mathbb{R}$  where  $w(x, y) \neq 0$  if and only if  $(x, y) \in E$ , a weighting  $W : E \rightarrow \mathbb{R}$  naturally transfers to  $w$ . In the case of an undirected graph, one simply requires the weight function to be symmetric, i.e.,  $w(x, y) = w(y, x)$  for all  $x, y \in \Omega_n$ . Furthermore, in the above definition the set  $\Omega_n$  does not enter the definition up to the ordering of its  $n \in \mathbb{N}$  elements. Therefore, a graph could be entirely represented by a weight function  $w : \{1, \dots, n\} \times \{1, \dots, n\} \rightarrow \mathbb{R}$ . However, since the definition of  $w$  will incorporate information about points  $\Omega_n \subset \mathbb{R}^d$  we use the notation  $(\Omega_n, w_n)$  for graphs in the following.

**Kernels and the Graph Scale** In most of our applications the data  $\Omega_n$  is given as a subset of  $\mathbb{R}^d$  and in fact we are interested in the limit  $n \rightarrow \infty$ , where  $\Omega_n$  fills out a domain  $\tilde{\Omega} \subset \mathbb{R}^d$ . In the continuum we consider *local* operators incorporating changes of functions  $u : \tilde{\Omega} \rightarrow \mathbb{R}$  at an infinitesimal small scale. Since interactions on a graph are inherently non-local in the Euclidean sense, we need to localize the interaction on the graph in the limit  $n \rightarrow \infty$ . A popular choice that guarantees this behavior is

$$w(x, y) = \eta_\varepsilon(|x - y|)$$

where  $\eta_\varepsilon : [0, \infty) \rightarrow [0, \infty]$  is a kernel function depending on a scaling parameter  $\varepsilon \in \mathbb{R}^+$ . The parameter  $\varepsilon$  is also referred to as the *graph scale* and informally speaking determines the scale of the graph interactions. I.e., the smaller  $\varepsilon$  the smaller the interaction radius of points in  $\Omega_n$  should be. In our specific setting of  $L^\infty$  problems it is typical to choose

$$\eta_\varepsilon(\cdot) = \frac{1}{c_\eta \varepsilon} \eta\left(\frac{\cdot}{\varepsilon}\right)$$

where  $\eta$  is a non-increasing kernel and  $c_\eta$  is a constant depending on the kernel. Typical examples of kernels include

- (constant weights)  $\eta(t) = 1_{[0,1]}(t)$ ,
- (exponential weights)  $\eta(t) = \exp(-t^2/(2\sigma^2))1_{[0,1]}(t)$ ,
- (non-integrable weights)  $\eta(t) = \frac{1}{t^p}1_{[0,1]}(t)$  with  $p \in (0, 1]$ .

**Remark 3.3.** In the continuum limit one needs to consider the value

$$\sigma_\eta = \sup_{t \in \mathbb{R}^+} t \eta(t),$$

which is assumed to be finite. Considering graph problems in  $L^p$  for  $p < \infty$  the corresponding value is given as

$$\sigma_\eta^{(p)} = \int_{\mathbb{R}^+} \eta(t) t^{d+p} dt \quad (3.1)$$

which was first employed in [GS15] for  $p = 1$  and then in [ST19] for general  $p < \infty$ . With the slight modification

$$\tilde{\sigma}_\eta^{(p)} = \int_{\mathbb{R}^+} \eta(t)^p t^{d+p} dt$$

we see that  $\tilde{\sigma}_\eta^{(p)}$  degenerates to  $\sigma_\eta$  in the limit  $p \rightarrow \infty$ .  $\triangle$

**Remark 3.4.** In [LIP-I] we choose  $c_\eta = 1$  and therefore  $\sigma_\eta$  then appears in the limit functional. In [LIP-II] we rescale the kernel, i.e.,  $c_\eta = \sigma_\eta$  which allows to work with unrescaled operators in the limit.  $\triangle$

**Remark 3.5.** In order to obtain continuum limits in the case  $p < \infty$  one usually employs a factor of  $1/\varepsilon^{p+d}$  in front of the graph weights [GS15; ST19]. In Section 3.2.2 we see that this factor degenerates to  $1/\varepsilon$  in the case  $p \rightarrow \infty$ . The intuition here is that problems in  $L^\infty$  do not treat mass in a quantitative but rather a qualitative manner. Therefore, factors like  $\varepsilon^d$  which appear because of integrals in  $\mathbb{R}^d$  do not contribute for  $p = \infty$ .  $\triangle$

**Sparse Graphs** An important practical aspect connected to the graph scale is the sparsity of the graph, which is given by the numbers of non-zero elements in the weight matrix  $W \in \mathbb{R}^{n \times n}$

$$W_{ij} := w(x_i, x_j) \quad x_i, x_j \in \Omega_n,$$

where we assume an ordering  $\Omega_n = \{x_1, \dots, x_n\}$ . In order to have a computationally feasible problem this matrix should have very few non-zero elements. Assuming that the kernel  $\eta$  has compact support, w.l.o.g.  $\text{supp}(\eta) \subset [0, 1]$  we observe that the sparsity is directly influenced by the graph scale  $\varepsilon$ . Namely,  $|x_i - x_j| > \varepsilon \Rightarrow W_{ij} = 0$ , see Fig. 3.1. Therefore, from a practical point of view  $\varepsilon$  should be chosen relatively small. However, showing consistency of graph-based SSI algorithms often requires scaling assumptions that permit the desired length scales [GS15; ST19; Cal19]. In Section 3.3 we give concrete examples of such scaling assumptions. One of the main goals of [LIP-I; LIP-II] was to show convergence and rates at the smallest possible length scale, which was indeed achieved.

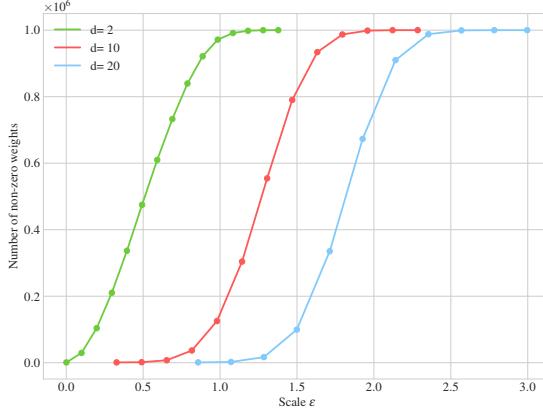


Figure 3.1.: Dependence of the number of non-zero edges on the scaling parameter  $\varepsilon$ . Here, we sample  $n = 10^3$  points uniformly in the cube  $[0, 1]^d$  and only connect vertices  $x_i, x_j \in \Omega_n$  if  $|x_i - x_j| < \varepsilon$ . The plots shows the behavior for  $d = 2, 10, 20$ .

**Remark 3.6.** In many practical applications graph weights connecting all points within an  $\varepsilon$ -ball of Euclidean distance are inferior to k-nearest neighbor (knn) graphs [Cal+20; FCL19; CT22]. While most consistency results employ the  $\varepsilon$ -ball setting, recently, the authors in [CT22] were able to show convergence rates for knn graphs. In this thesis we focus on the  $\varepsilon$ -ball setting, however it is an interesting open question to transfer the results of [LIP-I; LIP-II; LIP-III] to the knn setting.  $\triangle$

### 3.1.2. The $p$ -Laplacian: Continuum and Graph

The archetype of learning methods we consider in the following is the so-called *Laplacian learning*, which had one of its first appearances in [ZGL03]. The associated problem was given as

$$\begin{aligned} \min_{\mathbf{u}: \Omega_n \rightarrow \mathbb{R}} \sum_{x, y \in \Omega_n} w_n(x, y)^2 (\mathbf{u}(y) - \mathbf{u}(x))^2 \\ \text{subject to } \mathbf{u}(x) = \mathbf{g}(x) \text{ for all } x \in \mathcal{O}_n. \end{aligned} \tag{3.2}$$

Here, we always assume non-negative weights  $w_n : \Omega_n \times \Omega_n \rightarrow [0, \infty)$ . The intuitive idea behind this method is that it minimizes a discrete approximation of the Dirichlet energy and should therefore enforce a certain kind of “smoothness” of the solution  $\mathbf{u}$ . In [ZS05] it was proposed to generalize the idea of the graph Dirichlet energy in Eq. (3.2) to arbitrary  $1 \leq p < \infty$  motivated by the continuum counterpart. While our contributions in [LIP-I; LIP-II; LIP-III] mainly consider the case of  $p = \infty$ , the case  $p < \infty$  serves as a motivation for the whole discussion. Therefore, we briefly review the continuum setting and then discuss the situation on the graph.

**Continuum Setting** We follow the exposition in [Lin17]. Let  $\Omega \subset \mathbb{R}^d$  be a bounded domain, then we consider the  $p$ -Dirichlet energy for functions  $u \in W^{1,p}(\Omega)$ ,

$$\mathcal{E}_{p,\rho}(u) := \int_{\Omega} |\nabla u|^p \rho(x)^2 dx, \quad (3.3)$$

where  $\rho : \Omega \rightarrow \mathbb{R}$  is some given function. In the setting of [GS15; ST19]  $\rho$  denotes the density of data distribution as in Section 3.1. For  $\rho \equiv 1$  we simply write  $\mathcal{E}_{p,1} = \mathcal{E}_p$ . The associated variational problem is given in the following.

**Problem 3.7 (Variational Formulation).** For  $p \in [1, \infty)$  and  $V \subset W^{1,p}(\Omega)$  find  $u \in W^{1,p}(\Omega)$  such that

$$\mathcal{E}_{p,\rho}(u) \leq \mathcal{E}_{p,\rho}(v)$$

for all  $v$ , such that  $(u - v) \in W_0^{1,p}$ .

Assume for simplicity that  $\rho \equiv 1$  and that  $u \in V$  is a minimizer of the above problem, then its first variation must vanish, i.e., for all  $\phi \in C_0^\infty(\Omega)$  one has

$$\int_{\Omega} \langle |\nabla u|^p \nabla u, \nabla \phi \rangle dx = 0. \quad (3.4)$$

A function  $u \in V$  satisfying Eq. (3.4) is called a *weak solution* of the  $p$ -Laplace equation. In fact, if  $u$  is smooth enough one can infer that

$$\Delta_p u := \operatorname{div}(|\nabla u|^{p-2} \nabla u) = 0 \quad (3.5)$$

where  $\Delta_p$  is called the  $p$ -Laplacian. Boundary conditions on  $\partial\Omega$  given by a function  $g \in W^{1,p}(\Omega)$  can be incorporated by considering the set  $V_g := \{u \in W^{1,p}(\Omega) : u - g \in W_0^{1,p}(\Omega)\}$ . We state the following classical result, which can for example be found in [Lin17].

**Theorem 3.8 (Existence and Uniqueness).** For  $p \in (1, \infty)$  and  $g \in W^{1,p}(\Omega)$  there exists a unique minimizer  $u \in V_g$  of the  $p$ -Dirichlet energy, i.e.,

$$\operatorname{argmin}_{u \in V_g} \mathcal{E}_p(u) = u.$$

Moreover,  $u$  is a weak solution of the  $p$ -Laplace equation and there exists a function  $\tilde{u} \in C(\Omega)$  such that  $u = \tilde{u}$  a.e. in  $\Omega$ . If  $g \in C(\Omega)$  and  $\Omega$  is sufficiently smooth, then  $\tilde{u}|_{\partial\Omega} = g|_{\partial\Omega}$ .

*Proof.* The proof can be found in [Lin17, Thm. 2.16]. □

**Local Minimization Property** Let  $u$  solve [Problem 3.7](#) on the whole domain  $\Omega$  for given boundary values  $g \in W^{1,p}(\Omega)$  and let  $V \subset\subset \Omega$  be a sufficiently regular subset. Then we have that

$$\mathcal{E}_p(u) = \int_V |\nabla u|^p dx + \int_{\Omega \setminus V} |\nabla u|^p dx.$$

Now consider [Problem 3.7](#) restricted on  $V$  with boundary values given by  $u$  and denote by  $v$  the solution of this sub-problem on  $V$ . Since  $u = v$  on  $\partial V$  we can extend  $v$  onto  $\Omega$  by setting  $v = u$  in  $\Omega \setminus V$ , which gives  $v \in W^{1,p}(\Omega)$ . Therefore, we obtain

$$\mathcal{E}_p(v) = \int_V |\nabla v|^p dx + \int_{\Omega \setminus V} |\nabla u|^p dx \leq \int_V |\nabla u|^p dx + \int_{\Omega \setminus V} |\nabla u|^p dx = \mathcal{E}_p(u).$$

Since  $u$  was unique this yields  $u|_V = v|_V$ , and therefore we know that any solution on  $\Omega$  solves [Problem 3.7](#) also on any “nice” subset, with the boundary values given by itself. In [Section 3.2](#) we see that this *local minimization property* does not hold automatically and has to be enforced additionally. The underlying reason here is, that integrals are set-additive, but suprema are not, which is similarly explained in [\[ACJ04\]](#).

**Laplacian Learning** A natural extension of the problem given in [Eq. \(3.3\)](#) is obtained by considering the target functional

$$\mathbf{E}_p^{w_n}(\mathbf{u}) := \sum_{x,y \in \Omega_n} w_n(x,y)^p |\mathbf{u}(y) - \mathbf{u}(x)|^p,$$

which we refer to as the graph  $p$ -Dirichlet energy. Indeed, we notice structural similarities to the  $p$ -Dirichlet energy  $\mathcal{E}_p$  in [Eq. \(3.3\)](#), replacing the integral by a finite sum and derivatives by weighted finite differences

$$w_n(x,y)^p |\mathbf{u}(y) - \mathbf{u}(x)|^p.$$

This naturally leads to the following minimization problem.

**Problem 3.9 (Graph Energy Minimization).** Given a weighted graph  $(\Omega_n, w_n)$  and a labeling function  $\mathbf{g} : \mathcal{O}_n \rightarrow \mathbb{R}$ , for  $\mathcal{O}_n \subset \Omega_n$  we consider the problem

$$\min_{\mathbf{u}: \Omega_n \rightarrow \mathbb{R}} \mathbf{E}_p^{w_n}(\mathbf{u}) \text{ subject to } \mathbf{u}(x) = \mathbf{g}(x) \text{ for all } x \in \mathcal{O}_n.$$

Since every function  $\mathbf{u} : \Omega_n \rightarrow \mathbb{R}$  can be identified with a vector  $\mathbf{u} \in \mathbb{R}^n$ , the above problem is in fact an optimization problem in  $\mathbb{R}^n$ . The functional  $\mathbf{E}_p^{w_n} : \mathbb{R}^n \rightarrow \mathbb{R}$  is bounded from below by 0, continuous, convex and in the case  $p > 1$  even differentiable. However, one can prove unique existence with techniques very similar to the continuum case in [\[Lin17\]](#). We give the adapted proof below for completeness. The important property to establish uniqueness is that the graph is connected to the boundary, i.e., for every  $x \in \Omega_n$  there exists a path  $x = \gamma_1, \dots, \gamma_m = o \in \Omega_n$  with  $w(\gamma_{i+1}, \gamma_i) > 0$  connecting  $x$  to a point  $o \in \mathcal{O}_n$  in the boundary. This is very intuitive, since on any connected component that does not communicate with the boundary, an arbitrary constant minimizes the “graph gradient”.

**Theorem 3.10 (Existence and Uniqueness).** Let  $(\Omega_n, w_n)$  be a graph with non-negative weights, which is connected to its boundary  $\mathcal{O}_n \subset \Omega_n$ . Then Problem 3.9 admits a unique solution  $\mathbf{u} : \Omega_n \rightarrow \mathbb{R}$ .

*Proof.* The proof is an adaption from the continuum case in [Lin17, Thm. 2.16]. We start by proving uniqueness. Let  $E = \{(x, y) : w_n(x, y) > 0\}$  be the set of all active edges and denote by  $\nabla^{w_n} \mathbf{u} \in \mathbb{R}^{|E|}$

$$(\nabla^{w_n} \mathbf{u})_{(x,y)} := w_n(x, y) (\mathbf{u}(y) - \mathbf{u}(x))$$

the “graph gradient” on  $E$ , where we have that

$$\mathbf{E}_p^{w_n}(\mathbf{u}) = \sum_{(x,y) \in E} w_n(x, y)^p |\mathbf{u}(y) - \mathbf{u}(x)|^p = \sum_{(x,y) \in E} |(\nabla^{w_n} \mathbf{u})_{(x,y)}|^p = \|\nabla^{w_n} \mathbf{u}\|_p^p.$$

Let  $\mathbf{u}_1, \mathbf{u}_2$  be two solutions of Problem 3.9, i.e.  $\mathbf{E}_p^{w_n}(\mathbf{u}_1) = \mathbf{E}_p^{w_n}(\mathbf{u}_2)$ . Since  $\mathbf{u}_1, \mathbf{u}_2$  fulfill the boundary conditions, this also holds for  $1/2(\mathbf{u}_1 + \mathbf{u}_2)$  and therefore  $\mathbf{E}_p^{w_n}(\mathbf{u}_1) \leq \mathbf{E}_p^{w_n}(1/2(\mathbf{u}_1 + \mathbf{u}_2))$ .

Assume that there exists  $(\bar{x}, \bar{y}) \in E$  such that  $\nabla^{w_n} \mathbf{u}_1(\bar{x}, \bar{y}) \neq \nabla^{w_n} \mathbf{u}_2(\bar{x}, \bar{y})$ , then we have that

$$\begin{aligned} \mathbf{E}_p^{w_n}(\mathbf{u}_1) &\leq \mathbf{E}_p^{w_n} \left( \frac{\mathbf{u}_1 + \mathbf{u}_2}{2} \right) = \left\| \frac{\nabla^{w_n} \mathbf{u}_1(x, y) + \nabla^{w_n} \mathbf{u}_2(x, y)}{2} \right\|_p^p \\ &< \frac{1}{2} \sum_{(x,y) \in E \setminus \{(\bar{x}, \bar{y})\}} \left( |\nabla^{w_n} \mathbf{u}_1(x, y)|^p + |\nabla^{w_n} \mathbf{u}_2(x, y)|^p \right) \\ &\quad + \frac{1}{2} \left( |\nabla^{w_n} \mathbf{u}_1(\bar{x}, \bar{y})|^p + |\nabla^{w_n} \mathbf{u}_2(\bar{x}, \bar{y})|^p \right) \\ &= \frac{1}{2} \left( \mathbf{E}_p^{w_n}(\mathbf{u}_1) + \mathbf{E}_p^{w_n}(\mathbf{u}_2) \right) = \mathbf{E}_p^{w_n}(\mathbf{u}_1), \end{aligned}$$

which is a contradiction. Here, we employed the strict convexity of  $t \mapsto |t|^p$  which implies  $|t + \bar{t}|^p \leq 2^{p-1} (|t|^p + |\bar{t}|^p)$ , where the inequality is sharp if  $t \neq \bar{t}$ . Therefore, we obtain that  $\nabla^{w_n} \mathbf{u}_1 = \nabla^{w_n} \mathbf{u}_2$ , which yields

$$\mathbf{u}_1(y) - \mathbf{u}_1(x) = \mathbf{u}_2(y) - \mathbf{u}_2(x) \quad \text{for all } (x, y) \in E. \quad (3.6)$$

Since  $(\Omega_n, w_n)$  is connected to  $\mathcal{O}_n$  for any  $x \in \Omega_n$  we can find a path  $x = \gamma_1, \dots, \gamma_m = o$  connecting  $x$  to a point  $o \in \mathcal{O}_n$ . By definition  $w_n(\gamma_{i+1}, \gamma_i) > 0$  for all  $i = 1, \dots, m-1$  and therefore using Eq. (3.6) we obtain

$$\mathbf{u}_1(\gamma_i) - \mathbf{u}_1(\gamma_{i+1}) = \mathbf{u}_2(\gamma_i) - \mathbf{u}_2(\gamma_{i+1}) \quad \text{for all } i = 1, \dots, m-1.$$

which together with the boundary conditions  $\mathbf{u}_1(o) = \mathbf{u}_2(o) = \mathbf{g}(o)$  yields

$$\begin{aligned}\mathbf{u}_1(x) &= \mathbf{u}_1(\gamma_1) = \mathbf{u}_1(\gamma_1) - \mathbf{u}_1(\gamma_2) + \mathbf{u}_1(\gamma_2) \\ &= \mathbf{u}_1(\gamma_1) - \mathbf{u}_1(\gamma_2) + \mathbf{u}_1(\gamma_2) - \mathbf{u}_1(\gamma_3) + \mathbf{u}_1(\gamma_3) = \dots \\ &= \left[ \sum_{i=1}^{m-1} \mathbf{u}_1(\gamma_i) - \mathbf{u}_1(\gamma_{i+1}) \right] + \mathbf{u}_1(o) \\ &= \left[ \sum_{i=1}^{m-1} \mathbf{u}_2(\gamma_i) - \mathbf{u}_2(\gamma_{i+1}) \right] + \mathbf{u}_2(o) \\ &= \mathbf{u}_2(x).\end{aligned}$$

Since  $x \in \Omega_n$  was arbitrary, we get  $\mathbf{u}_1 = \mathbf{u}_2$ .

To show existence, we do not need to assume that the graph is connected to the boundary. On every connected component  $V \subset \Omega_n$  that is not connected to  $\mathcal{O}$  we can set  $\mathbf{u}(x) = c$  for all  $x \in V$ , where  $c \in \mathbb{R}$  is an arbitrary constant. We can select all edges in  $V$ , i.e.  $E_V := E \cap V^{\times 2}$  and since  $V$  was a connected component—i.e. it does not have active edges to vertices outside of  $V$ —we can also split up the functional

$$\mathbf{E}_p^{w_n}(\mathbf{u}) = \sum_{(x,y) \in E \setminus E_V} |\nabla^{w_n} \mathbf{u}(x, y)|^p + \underbrace{\sum_{(x,y) \in E_V} |\nabla^{w_n} \mathbf{u}(x, y)|^p}_{=0} = \sum_{(x,y) \in E \setminus E_V} |\nabla^{w_n} \mathbf{u}(x, y)|^p$$

where the contribution on the connected component is zero, since we set  $\mathbf{u}$  to be constant here. Therefore, w.l.o.g. we can assume that  $\Omega_n$  is connected to the boundary, since any component not connected to boundary does not contribute to the problem. Let now  $\mathbf{u} : \Omega_n \rightarrow \mathbb{R}$  be a function and for  $x \in \Omega_n$  let  $\gamma \in \Omega_n^{\times m}$  be a path to a point in the boundary  $o^x \in \mathcal{O}_n$ , then we have that

$$|\mathbf{u}(x)|^p \leq C \left( \sum_{i=1}^{m^x-1} |\mathbf{u}(\gamma_i) - \mathbf{u}(\gamma_{i+1})|^p + |\mathbf{g}(o^x)|^p \right) \leq C \left( \mathbf{E}_p^{w_n}(\mathbf{u}) + |\mathbf{g}(o^x)|^p \right)$$

where  $C > 0$  is a generic constant—not depending on  $\mathbf{u}$ —that also accounts for the factor  $(\min_{(x,y) \in E} w_n(x, y))^{-1}$ . Summing over all  $x \in \Omega_n$  yields

$$\begin{aligned}\|\mathbf{u}\|_p^p &= \sum_{x \in \Omega_n} |\mathbf{u}(x)|^p \leq C \left( n \mathbf{E}_p^{w_n}(\mathbf{u}) + \sum_{x \in \Omega_n} |\mathbf{g}(o^x)|^p \right) \\ &\leq C \left( \mathbf{E}_p^{w_n}(\mathbf{u}) + \|\mathbf{g}\|_p^p \right).\end{aligned}$$

Let now  $\mathbf{u}_j$  denote a sequence such that

$$\mathbf{E}_p^{w_n}(\mathbf{u}_j) \leq I + 1/j, \quad \text{where} \quad I := \inf_{\mathbf{u}: \mathbf{u} = \mathbf{g} \text{ on } \mathcal{O}_n} \mathbf{E}_p^{w_n}(\mathbf{u}).$$

Then we have that

$$\|\mathbf{u}_j\|_p^p \leq C (\mathbf{E}_p^{w_n}(\mathbf{u}_j) + \|\mathbf{g}\|_p^p) \leq C (I + 1/j + \|\mathbf{g}\|_p^p) \leq \tilde{C}$$

$j \in \mathbb{N}$ , i.e.  $\mathbf{u}_j$  is a uniformly bounded sequence of vectors in the finite dimensional space  $\mathbb{R}^n$ , and therefore there exists a subsequence—which we do not relabel—that converges to some  $\mathbf{u}^* \in \mathbb{R}^n$ . Since the functional  $\mathbf{E}_p^{w_n} : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuous we obtain that

$$\mathbf{E}_p^{w_n}(\mathbf{u}^*) = \lim_{j \rightarrow \infty} \mathbf{E}_p^{w_n}(\mathbf{u}_j) = I$$

which yields existence.  $\square$

**The Graph Laplacian** In the continuum case one considers the Euler–Lagrange equation for the functional  $\mathcal{E}_p$ , which yields the  $p$ -Laplacian, see [Section 3.1.2](#). Analogously, the optimality conditions for the graph  $p$ -Dirichlet energy  $\mathbf{E}_p^{w_n} : \mathbb{R}^n \rightarrow \mathbb{R}$  for  $p > 1$  yield

$$\nabla_{\mathbf{u}} \mathbf{E}_p^{w_n}(\mathbf{u}) = 0,$$

where for  $x \in \Omega_n$  we have

$$(\nabla_{\mathbf{u}} \mathbf{E}_p^{w_n}(\mathbf{u}))_x =: \Delta_p^{w_n} \mathbf{u}(x) = \sum_{y \in \Omega_n} w_n(x, y)^p |\mathbf{u}(y) - \mathbf{u}(x)|^{p-2} (\mathbf{u}(y) - \mathbf{u}(x)),$$

which is referred to as the graph  $p$ -Laplacian operator. This yields the following problem.

**Problem 3.11 (Graph  $p$ -Laplacian).** Given a weighted graph  $(\Omega_n, w_n)$  and a labeling function  $\mathbf{g} : \mathcal{O}_n \rightarrow \mathbb{R}$  with  $\mathcal{O}_n \subset \Omega_n$ , find a function  $\mathbf{u} : \Omega_n \rightarrow \mathbb{R}$  such that

$$\begin{aligned} \Delta_p^{w_n} \mathbf{u} &= 0, \text{ in } \Omega_n \setminus \mathcal{O}_n, \\ \mathbf{u} &= \mathbf{g} \text{ on } \mathcal{O}_n. \end{aligned}$$

Since the functional  $\mathbf{E}_p^{w_n}$  has a unique minimizer subject to the constraints given by  $\mathbf{g}$  and the graph  $p$ -Laplacian is derived via optimality conditions, one expects that the [Problem 3.9](#) and [Problem 3.11](#) are equivalent. This is formulated in the following theorem.

**Theorem 3.12 (Existence and Uniqueness).** Let  $(\Omega_n, w_n)$  be a graph with non-negative weights, that is connected to its boundary  $\mathcal{O}_n \subset \Omega_n$ . Then there exists a unique solution  $\mathbf{u} : \Omega_n \rightarrow \mathbb{R}$  to [Problem 3.11](#), which is also the unique minimizer of [Problem 3.9](#).

*Proof.* This statement can be proven similarly to [\[ETT15, Theorem 5.3\]](#).  $\square$

### 3.1.3. Consistency for Graph-based SSL

While the graph 2-Laplacian has been successfully employed for classification tasks, see e.g. [\[ZGL03; ZS05\]](#), it has been observed that solutions tend to degenerate for large amounts of data, whenever  $d \geq 2$  [\[NSZ09; AL11; El+16\]](#). We illustrate this phenomenon in [Example 3.13](#).

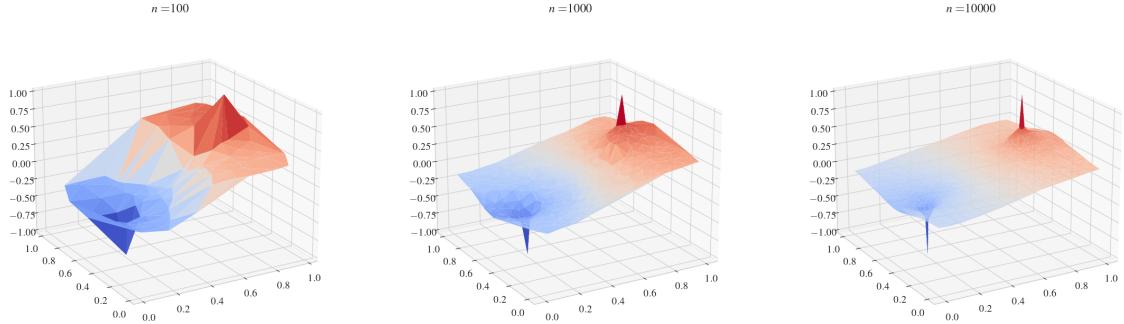


Figure 3.2.: Solution to the Laplacian Learning problem ( $p = 2$ ) for different number of data points  $n \in \{100, 1000, 10000\}$ .

**Example 3.13.** We sample different number of points  $n \in \{100, 1000, 10000\}$  in  $\Omega = (0, 1)^2$  and keep a fixed constraint on  $\mathcal{O}_n = \{o_1, o_2\} = \{(0.2, 0.5), (0.8, 0.5)\}$  with  $\mathbf{g}(o_1) = -1, \mathbf{g}(o_2) = 1$ . In Fig. 3.2 we observe that for increasing  $n$  the solution  $\mathbf{u}$  of Eq. (3.2) tends to be a constant equal to the average of the given constraints and spike at the given constraints.

This observation motivates the driving question of this chapter and the works presented:

*Are semi-supervised learning algorithms consistent in the infinite data limit?*

In our case “consistency” roughly ask the discrete solutions to converge to the continuum counterparts in the data limit  $n \rightarrow \infty$ . Here, we distinguish between point-wise consistency as in [VBB08; GK06; HAV05], consistency on a function space level [ST19; GS15; Cal19; LIP-I; LIP-II] and consistency of a discretized operator [Cal19; LIP-II]. Other related work also consider the limit of clustering methods [Hof+22; GS18] or the no-noise limit [Hof+20; Dun+20]. We detail the notion of consistency in the following.

**Consistency for the  $p$ -Laplacian** As observed in [NSZ09; AL11; El +16; CS20] the question of consistency is connected to the relation between  $p$  and  $d$ , where one considers the three regimes  $p < d$ ,  $p = d$  and probably the most relevant case  $p > d$ . The  $p$ -Dirichlet problem in the continuum can be formulated as a variational problem in  $W^{1,p}$ , however, in our setting we most likely consider pointwise constraints on a discrete set  $\mathcal{O}$  even in the continuum. Therefore, one requires that functions in  $W^{1,p}$  exhibit some continuity properties. By the Sobolev embedding theorem, this is the case when  $p > d$ , [AF03]. This leads to a first qualitative intuition that consistency can only be achieved in the case  $p > d$ . We note that  $p$ -harmonic functions are actually more regular also in the case  $p \leq d$ , beyond the implication of the Sobolev embedding theorem. However, this regularity does not help for the continuum limits [NSZ09; AL11; El +16].

The most relevant result for our contribution in [LIP-I] is the  $\Gamma$ -convergence statement of [GS15], which considers a clustering task, via the TV functional. Here, the authors

show

$$\frac{1}{\varepsilon_n n^2} \mathbf{E}_1^{w_n} \xrightarrow{\Gamma} \sigma_\eta \mathcal{E}_{1,\rho}$$

where the constant  $\sigma_\eta$  is defined as in [Eq. \(3.1\)](#). We review details on  $\Gamma$ -convergence in [Section 3.3](#). The important insight here, is that the graph scaling  $\varepsilon_n$  must not tend to zero too fast, which is guided by the following intuition:

- Problems on the graph are non-local, communication between points takes place on a finite scale.
- In the limit we want to obtain a local problem, communication takes place at a infinitesimal small scale.
- If the graph scale  $\varepsilon_n$  tends to zero too fast, there is not enough communication between vertices, or even worse the graph becomes disconnected.

The optimal scaling obtained in [\[GS15\]](#) for  $d \geq 3$  is given as

$$\lim_{n \rightarrow \infty} \left( \frac{\log n}{n} \right)^{1/d} \varepsilon_n^{-1} = 0, \quad (3.7)$$

i.e.  $\varepsilon_n$  must go to zero slower than  $\left( \frac{\log n}{n} \right)^{1/d}$ . This scaling is optimal, in the sense that the graph is disconnected with high probability if  $\varepsilon_n < \lambda \left( \frac{\log n}{n} \right)^{1/d}$ , see [\[GS15\]](#). The main difference here, is however that the considered clustering task does not directly involve a given labeling on the set  $\mathcal{O}_n \subset \Omega_n$ . This setting was considered in [\[ST19\]](#), where first the generalization of the  $\Gamma$ -convergence statement in [Eq. \(3.7\)](#) was provided. Furthermore, the authors then consider the constrained model by incorporating the boundary conditions into the functional

$$\mathbf{E}_p^{w_n, \text{cons}}(\mathbf{u}) := \begin{cases} \mathbf{E}_p^{w_n}(\mathbf{u}) & \text{if } \mathbf{u} = \mathbf{g} \text{ on } \mathcal{O}_n, \\ \infty & \text{else.} \end{cases}$$

In order to show  $\Gamma$ -convergence of the constrained functionals, an additional upper bound on the scaling has to be assumed. This means that the graph scaling must not tend to zero too slowly, namely

$$n \varepsilon_n^p \xrightarrow{n \rightarrow \infty} 0.$$

Together, with [Eq. \(3.7\)](#) this then implies

$$n^{-1/d} \ll \varepsilon_n \ll n^{-1/p}$$

which is only possible if  $p > d$ . Therefore, one again obtains well-posedness of the constrained problem, if  $p$  is large enough.

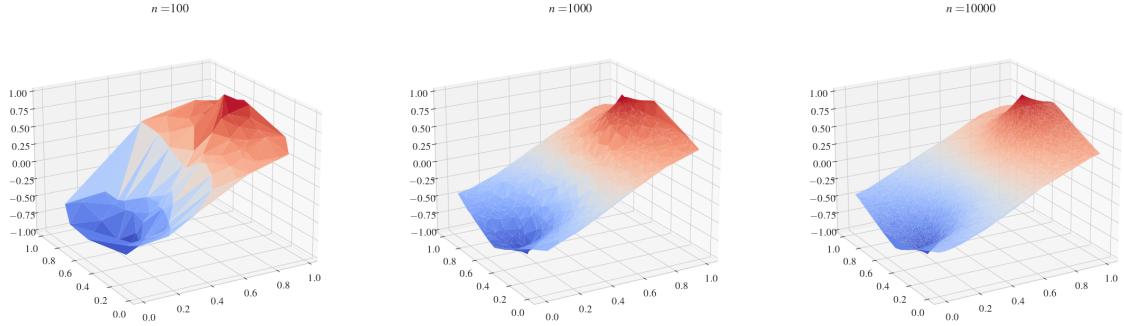


Figure 3.3.: Solution to the Laplacian Learning problem ( $p = \infty$ ) for different number of data points  $n \in \{100, 1000, 10000\}$ . The setup is otherwise copied from [Example 3.13](#).

**Sending  $p$  to Infinity** Since the assumption that  $p > d$  is crucial for asymptotic consistency, a natural idea is to send  $p$  to  $\infty$  and analyze the corresponding limit problem. This yields the so-called *Lipschitz learning* task [vLB04; Kyn+15], which is the main point of interest in this chapter. We formalize this problem in [Section 3.2](#). In [Fig. 3.3](#) we observe that for  $p = \infty$  we indeed obtain smoother solutions, compared to [Fig. 3.2](#). In this regard Lipschitz learning seems promising to overcome the consistency issue.

However, as already noticed in [El +16] being a pure  $L^\infty$  problem, Lipschitz learning does not respect the density of data in any way. Namely, the method is exclusively distance based. This behavior is a drawback in many machine learning applications. However, by carefully rescaling the graph weights one obtains a data sensitive problem, see [Cal19] and [Section 3.3](#). Here, we also want to mention an alternative method to overcome the spiking phenomenon for the  $p$ -Laplacian, by re-weighting the graph accordingly [CS20].

The contributions presented in this chapter are all connected to the Lipschitz learning task. Namely, we are able to show consistency and convergence rates under very mild scaling assumptions.

## 3.2. Lipschitz Extensions and the Infinity Laplacian: Continuum and Graph

Our main point of interest in this chapter is the Lipschitz learning task and its continuum limit. In this section we first motivate the continuum problem as considered in [[LIP-I](#); [LIP-II](#)] and then give some details on the discrete problem.

### 3.2.1. The Continuum Setting

We first recall, that for  $u \in W^{1,\infty}(\Omega)$  we have that

$$\lim_{p \rightarrow \infty} \mathcal{E}_{p,\rho}(u)^{1/p} = \operatorname{ess\,sup}_{x \in \operatorname{supp}(\rho)} |\nabla u(x)| =: \mathcal{E}_{\infty,\rho}(u),$$

see [Jen93]. Again, we see that that the weighting  $\rho$  only enters via its support, therefore we do not explicitly consider it in the following and instead assume  $\Omega = \text{supp}(\rho)$  and then only write  $\mathcal{E}_\infty$ . The functional  $\mathcal{E}_\infty$  is weak\*-lower semi-continuous over  $W^{1,\infty}(\Omega)$ , (see e.g. [BJW01, Thm. 2.6]). In the classical theory developed by Jensen in [Jen93] one considers the following problem, which tries to minimize the *sup-norm* of the gradient as described by Jensen.

**Problem 3.14 (Gradient Sup-Norm Problem).** For an open domain  $\Omega \subset \mathbb{R}^d$ , find a function  $u \in W^{1,\infty}(\Omega)$  such that

$$\|\nabla u\|_\infty \leq \|\nabla v\|_\infty \quad \text{for every } v \text{ with } (u - v) \in W_0^{1,\infty}(\Omega).$$

This problem can be seen as the limit problem for  $p \rightarrow \infty$  of [Problem 3.7](#). Additionally imposing boundary conditions with  $g : \partial\Omega \rightarrow \mathbb{R}$ , [Jen93] then draws the connection to so-called *Lipschitz extensions*, which are the driving concept in this section. We introduce a more general viewpoint—that does not require the notion of a gradient—later on, but first introduce the variational problem.

**The Lipschitz Constant** As noticed in [Jen93] working with the Lipschitz constant and the sup-norm of the gradient requires a careful treatment of the distance employed. Let  $\tilde{\Omega}$  be a set and let  $d(\cdot, \cdot)$  be a semi-metric on  $\tilde{\Omega}$ , that is  $d(\cdot, \cdot)$  fulfills the requirement of a metric up to triangle inequality. Then we define the Lipschitz constant of a function  $u : V \rightarrow \mathbb{R}$  on a subset  $V \subset \tilde{\Omega}$  as

$$\text{Lip}_d(u; V) := \sup_{x,y \in V, x \neq y} \frac{|u(x) - u(y)|}{d(x, y)}.$$

If  $d(\cdot, \cdot)$  denotes the Euclidean distance we omit the subscript. i.e.  $\text{Lip}_d = \text{Lip}$ . Additionally, we can introduce the space of Lipschitz functions  $\text{Lip}_d(V)$  on  $V$  via  $u \in \text{Lip}_d(V) \Leftrightarrow \text{Lip}_d(u; V) < \infty$ .

**Remark 3.15 (Lipschitz and Sobolev functions).** If  $\Omega \subset \mathbb{R}^d$  is sufficiently regular—e.g., it has Lipschitz boundary—then we have that

$$\text{Lip}(\Omega) = W^{1,\infty}(\Omega),$$

where this identity is of course to be understood in the sense of equivalence classes in  $L^p$  spaces. We refer to [Eva18] for a proof of this result.  $\triangle$

The above remark already relates Lipschitz with  $W^{1,\infty}$  functions. Often however, we need a quantitative comparison between the Lipschitz constant and the sup-norm of the gradient of a function. Here, it is essential which distance is chosen for the Lipschitz constant. For an open domain  $\Omega \subset \mathbb{R}^d$  and  $u \in W^{1,\infty}$  we have the inequality

$$\|\nabla u\|_\infty \leq \sup_{x \neq y} \frac{|u(x) - u(y)|}{|x - y|}$$

which can be proven via the definition of the gradient. For the reverse inequality, one has to take the geometry of the domain into account, namely for  $x, y \in \Omega$  we have that

$$|u(x) - u(y)| \leq \|\nabla u\|_\infty d_\Omega(x, y) \quad (3.8)$$

see, e.g., [BB11, Prop9.3, Rem. 7]. Here,

$$d_\Omega(x, y) = \inf \left\{ \int_0^1 |\dot{\gamma}(t)| dt : \gamma \in C^1([0, 1], \Omega) \text{ with } \gamma(0) = x, \gamma(1) = y \right\}$$

denotes the *geodesic distance* on  $\Omega$ . If  $\Omega$  is convex, we have that  $d_\Omega(x, y) = |x - y|$  for every  $x, y \in \Omega$  and therefore Eq. (3.8) yields  $\text{Lip}(u) = \|\nabla u\|_\infty$ . However, this situation changes for non-convex domains, see Example 3.16. Additionally it is often necessary to define a distance measure on the closure of  $\Omega \subset \mathbb{R}^d$ . In order to have a geodesic on  $\bar{\Omega}$  one can simply consider  $d_{\bar{\Omega}}$ , see e.g. [LIP-II], which then yields the length space  $(\bar{\Omega}, d_{\bar{\Omega}})$ . In the classical theory developed in [Jen93] one alternatively considers

$$\tilde{d}_{\bar{\Omega}}(x, y) := \liminf_{(\tilde{x}, \tilde{y}) \rightarrow (x, y)} d_\Omega(\tilde{x}, \tilde{y}).$$

The differences between these notions are demonstrated in the following example. Also note, that  $\tilde{d}_{\bar{\Omega}}$  is only a semi-metric on  $\bar{\Omega}$  since it is lacking a triangle inequality.

**Example 3.16.** For  $I = [-\pi, c] \cup [c, \pi]$  with  $c = \pi/6$  we consider the domain

$$\bigcup_{\theta \in I} B_1((\cos(\theta), \sin(\theta)))$$

which is visualized in Fig. 3.4 and the points  $x = (2c, 1), y = (2c, -1)$ . The line segment between  $x$  and  $y$  contains the point  $z = (2c, 0)$ , however  $z \notin \Omega$ . One can show that the geodesic has the length  $d_\Omega(x, y) = 4 \cos(\pi/6) + \pi \approx 6.606$  which is the length of the dotted path in Fig. 3.4. However, we observe that

$$\bar{\Omega} = \bigcup_{\theta \in I} \overline{B_1((\cos(\theta), \sin(\theta)))}$$

and in particular  $z \in \overline{B_1((c, c))}$ , therefore  $d_{\bar{\Omega}}(x, y) = 2$ .

**Solutions to the Gradient Sup-Norm Problem** Before generalizing the theory of Lipschitz extension to arbitrary metric spaces, we first note, that one can explicitly construct solutions of Problem 3.14. Namely, for given  $g \in \text{Lip}(\partial\Omega)$  the functions

$$\begin{aligned} \bar{g}(x) &:= \inf_{y \in \partial\Omega} g(y) + \text{Lip}_{\tilde{d}_{\bar{\Omega}}}(g; \partial\Omega) \cdot \tilde{d}_{\bar{\Omega}}(x, y) \\ \underline{g}(x) &:= \sup_{y \in \partial\Omega} g(y) - \text{Lip}_{\tilde{d}_{\bar{\Omega}}}(g; \partial\Omega) \cdot \tilde{d}_{\bar{\Omega}}(x, y) \end{aligned} \quad (3.9)$$

are solutions to the gradient sup-norm problem that coincide with  $g$  on  $\partial\Omega$ , see [Jen93, Th. 1.8]. The same concept of constructing solutions is applied in the following sections

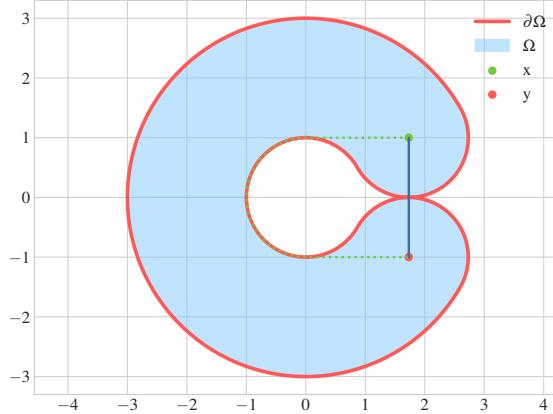


Figure 3.4.: The domain in Example 3.16.

in a more abstract setting. These solutions are then called Whitney and McShane or respectively maximal and minimal extensions, see Lemma 3.19.

One easily observes that there are cases where  $\bar{g} \neq g$  and therefore the problem does not admit a unique solution. A concrete example, to showcase this phenomena is given in [Jen93, p. 53].

**Lipschitz Extensions in Metric Spaces** The problem considered in the last section was motivated by a variational problem for the functional  $\mathcal{E}_\infty(u) = \|\nabla u\|_\infty$ . However, the theory of Lipschitz extensions provides a more general framework. Namely, it is not necessary that  $\Omega$  is a subset of  $\mathbb{R}^d$  but we can rather consider a metric space  $(\tilde{\Omega}, d)$  with  $\Omega \subset \tilde{\Omega}$ .

**Remark 3.17.** For applications within this thesis we have that  $\Omega \subset \mathbb{R}^d$  is an open bounded domain and then consider  $\tilde{\Omega} := \overline{\Omega}$ , i.e., the closure of  $\Omega$  within the topology induced by the Euclidean distance.  $\triangle$

A result originally due to Kierszbraun [Kir34] states that for two Hilbert spaces  $\mathcal{X}, \mathcal{Y}$ , a subset  $\mathcal{O} \subset \mathcal{X}$  and a function  $g : \mathcal{O} \rightarrow \mathcal{Y}$  there exists a function  $u : \mathcal{X} \rightarrow \mathcal{Y}$  such that

$$\begin{aligned} u &= g \text{ on } \mathcal{O}, \\ \text{Lip}(u; \mathcal{X}) &= \text{Lip}(g; \mathcal{O}). \end{aligned}$$

Here, the metrics for the respective Lipschitz constants are induced by the inner products of the Hilbert spaces. We refer to [Kir34] for the original proof and to [Sch69, Th. 1.31] for a proof of the version as stated above. In this work we only consider the case  $\mathcal{Y} = \mathbb{R}$  which allows for more general assumptions on the space  $\mathcal{X}$ . Below, we formulate the Lipschitz extension problem in our setting.

**Problem 3.18 (Lipschitz Extensions).** Let  $(\tilde{\Omega}, d)$  be a metric space and  $\mathcal{O} \subset \tilde{\Omega}$  be a bounded subset. For a given Lipschitz function  $g : \mathcal{O} \rightarrow \mathbb{R}$  find a Lipschitz

function  $u : \tilde{\Omega} \rightarrow \mathbb{R}$  such that

$$\text{Lip}_d(u; \tilde{\Omega}) = \text{Lip}_d(g; \mathcal{O}).$$

A function  $u : \tilde{\Omega} \rightarrow \mathbb{R}$  with this property is called *Lipschitz extension* of  $g$  to  $\tilde{\Omega}$ .

In this setting one can again explicitly construct solutions of the Lipschitz extension task. They are not unique, however one has an upper and a lower bound. In fact, conceptually these solutions are very similar to the functions in Eq. (3.9) and even coincide, whenever the sup-norm of the gradient is given as the Lipschitz constant.

**Lemma 3.19.** In the setting of Problem 3.18 we have that the

- **Whitney (or maximal) extension:**  $\bar{g}(x) := \inf_{y \in \mathcal{O}} g(y) + \text{Lip}_d(g; \mathcal{O}) \cdot d(x, y)$  and the
- **McShane (or minimal) extension:**  $\underline{g}(x) := \sup_{y \in \mathcal{O}} g(y) - \text{Lip}_d(g; \mathcal{O}) \cdot d(x, y)$

defined for  $x \in \tilde{\Omega}$  are Lipschitz extensions of  $g$  to  $\tilde{\Omega}$ . Moreover, let  $u : \tilde{\Omega} \rightarrow \mathbb{R}$  be any Lipschitz extension of  $g$ , the we have that

$$\underline{g} \leq u \leq \bar{g}.$$

*Proof.* We refer to [Whi92] and [McS34] for the proofs of the respective result.  $\square$

As demonstrated in Example 3.20, there are cases where  $\bar{g} \neq g$  and therefore, Lipschitz extensions are not unique in general. Furthermore, [ACJ04] points out that the Whitney and McShane extension do not allow for a comparison principle, which can also be observed in Example 3.20.

**Example 3.20.** Consider the set  $\tilde{\Omega} = [-1, 1]$  and  $\mathcal{O} = \{-1, 0, 1\}$  with

$$\begin{aligned} g_1(x) &:= 0, \\ g_2(x) &:= 1/2 (x - |x|), \\ g_3(x) &:= -g_2. \end{aligned}$$

Then we have that  $\underline{g}_2 \leq \underline{g}_1$  on  $\mathcal{O}$  but

$$\bar{g}_2 > \bar{g}_1 \text{ in } (0, 1),$$

see Fig. 3.5 for a visualization. Analogously, we have that  $\underline{g}_3 \geq \underline{g}_1$  on  $\mathcal{O}$  but

$$\bar{g}_3 < \bar{g}_1 \text{ in } (0, 1).$$

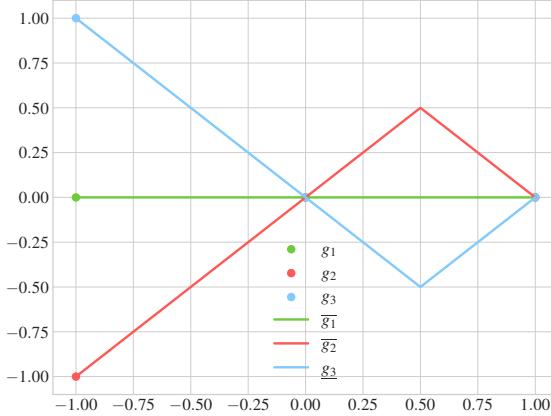


Figure 3.5.: The maximal extension does not admit a comparison principle, as demonstrated in Example 3.20.

**Absolutely Minimizing Extension** Sending  $p \rightarrow \infty$  in the variational formulation of the  $p$ -Laplace equation yields the Lipschitz extension task, which however does not admit unique solutions. So the question arises, which property is lost in the limit case. For  $p < \infty$  one has the local minimization property, as explained in Section 3.1.3. This lead Aronsson to introduce the concept of *absolutely minimizing Lipschitz extension* in [Aro67], by additionally enforcing the minimizing property on every subset. A function  $u \in W^{1,\infty}$  is called absolutely minimal, iff

$$\operatorname{ess\ sup}_{x \in V} |\nabla u| \leq \operatorname{ess\ sup}_{x \in V} |\nabla v| \quad \text{for every open } V \subset \Omega \quad (3.10)$$

and every function  $v$  such that  $(u - v) \in W_0^{1,\infty}$ . In fact in [Aro67] it is also shown, that for  $u_p$  denoting the solution of the corresponding  $p$ -Dirichlet problem, we have that  $u_p \xrightarrow{p \rightarrow \infty} u_\infty$ , which seems to validate the notion of absolute minimizers. In [ACJ04] it was shown, that one has an equivalent formulation involving the Lipschitz constant. For a given Lipschitz function  $g : \bar{\Omega} \rightarrow \mathbb{R}$  we have that  $u_\infty$  with  $(u_\infty - g) \in W_0^{1,\infty}(\Omega)$  fulfills Eq. (3.10) iff

$$\operatorname{Lip}(u_\infty; V) \leq \operatorname{Lip}(v; V) \text{ for every } V \subset \Omega$$

and every function  $v$  such that  $(u - v) \in W_0^{1,\infty}(V)$ , see [Aro67]. In this thesis we work with a notion of absolute minimizers, which is equivalent to the above formulation for convex domains in  $\mathbb{R}^d$ . However, for our applications it is more convenient to formulate the problem for abstract length spaces.

**Problem 3.21 (AMLEs).** Let  $(\tilde{\Omega}, d)$  be a length space,  $\mathcal{O} \subset \tilde{\Omega}$  a closed subset and  $g : \mathcal{O} \rightarrow \mathbb{R}$  a Lipschitz function. Find an extension  $u \in C(\tilde{\Omega})$  such that  $u = g$  on  $\mathcal{O}$  and

$$\operatorname{Lip}_d(u; \bar{V}) = \operatorname{Lip}_d(u, \partial V) \text{ for all open and connected sets } V \subset \tilde{\Omega} \setminus \mathcal{O}. \quad (3.11)$$

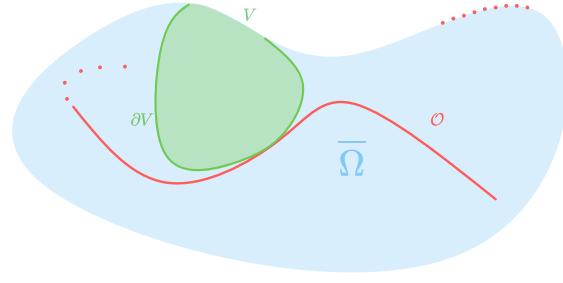


Figure 3.6.: A set  $V \subset \bar{\Omega}$  can be relatively open w.r.t. the metric space  $\bar{\Omega}$  although,  $V \cap \partial\bar{\Omega} \neq \emptyset$ , where  $\partial\bar{\Omega}$  is the boundary within the standard topology on  $\mathbb{R}^d$ . The relative boundary of  $\partial_{\bar{\Omega}}V$  does not include any parts of  $\partial\bar{\Omega}$ .

A function  $u$  fulfilling this property is called absolutely minimizing Lipschitz extension of  $g$ .

**Remark 3.22.** In our setting  $\Omega$  is an open subset of  $\mathbb{R}^d$  and we then choose  $\tilde{\Omega} = \bar{\Omega}$ . Here, it is important to note that the topological notions like boundary and interior are to be understood relative to  $\bar{\Omega}$ . A visualization of this concept can be found in Fig. 3.6. △

Fulfilling Eq. (3.11) also makes sense without enforcing boundary conditions on some set  $\mathcal{O}$ . Typically, one just says  $u$  is absolutely minimizing (AM) if it Eq. (3.11) holds, see [ACJ04]. In the convex and Euclidean case, existence of solutions to Problem 3.21 follow directly from Aronssons proof that the limit  $u_p \xrightarrow{p \rightarrow \infty} u_\infty$  fulfills Eq. (3.10), see [Aro67]. For general length spaces, it can be shown via an adapted Perron's method [Per23], which we state in the following result from [Juu02].

**Theorem 3.23 ([Juu02, Th. 4.3]).** Let  $(\tilde{\Omega}, d)$  be a separable length space, and suppose  $g : \mathcal{O} \rightarrow \mathbb{R}$  is Lipschitz. Then there exists an absolutely minimizing Lipschitz extension  $u$  of  $g$  to  $\tilde{\Omega}$ , i.e.  $u$  solves Problem 3.21.

We address the question of existence and uniqueness of solutions to Problem 3.21 in the following paragraphs.

**The Infinity Laplacian** In the Euclidean case one can derive an operator equation by considering the limit of the  $p$ -Laplacian operator. This yields the  $\infty$ -Laplacian which is defined as

$$\Delta_\infty u = \langle \nabla u, D^2 u \nabla u \rangle = \sum_{i,j} \partial_{x_i} u \partial_{x_j} u \partial_{x_i, x_j} u$$

for functions of class  $C^2$ . This operator had one of its first appearances in [Aro68], but we refer to [Lin17] for a detailed overview of the topic. The brief exposition here,

also follows [Lin17] to some extent. Intuitively the  $\infty$ -Laplace operator gives the second derivative of  $u$  into the direction of its gradient. The associated infinity Laplace equation takes the following form.

**Problem 3.24.** Let  $\Omega \subset \mathbb{R}^d$  be an open and bounded domain and  $g \in W^{1,\infty}$ , then the task is to find  $u$  such that

$$\begin{aligned}\Delta_\infty u &= 0 && \text{in } \Omega, \\ u &= g && \text{on } \partial\Omega.\end{aligned}$$

At first glance, the definition of the  $\infty$ -Laplacian only makes sense for functions of class  $C^2$ . However, as shown in [Yu06; Aro68] if  $u \in C^2(\Omega)$  fulfills  $\Delta_\infty u = 0$  in  $\Omega$ , then either  $\nabla u \neq 0$  in the whole domain or it is constant. However, one can easily construct boundary values  $g$  that are not constant but force the solution  $u$  to have critical points inside of the domain, see e.g. [Lin16]. Therefore, one requires a different concept, where one typically considers *viscosity solutions*. For the  $\infty$ -Laplacian this strategy was first employed in [BDM89], where the term *viscosity* is motivated by similar concepts for the Burgers' equation [Bur48]. First we consider subsolutions, where we take any  $\phi \in C^2(\Omega)$  that touches  $u$  from above at any  $x \in \Omega$ , i.e.

$$u < \phi \text{ in } \Omega \setminus \{x\} \quad \text{and } u(x) = \phi(x). \quad (3.12)$$

Since we  $\phi$  is  $C^2$  we can shift the application of  $\Delta_\infty$  onto this functions and then say  $u$  is a subsolution if for any  $\phi$  fulfilling Eq. (3.12) we have

$$-\Delta_\infty \phi(x) \leq 0.$$

Note, that the  $\infty$ -Laplacian is only evaluated at the touching point  $x \in \Omega$ . Analogously, we say  $u$  is a supersolution if for any  $\phi$  touching from below we have  $-\Delta_\infty \phi(x) \geq 0$ . If  $u$  is both a sub- and a supersolution, we say it is a solution in the viscosity sense.

**Remark 3.25.** The concept of viscosity solutions can similarly be applied to the  $p$ -Laplacian or even a more general class of differential operators, see [Lin17]. We also note, that one has the consistency result, that if  $u \in C^2$  is a viscosity solution, we also have that  $\Delta_\infty u = 0$  in the classical sense.  $\triangle$

Solving the  $\infty$ -Laplace equation in the viscosity sense is in fact equivalent to being absolutely minimizing, which we state in the following theorem taken from [ACJ04]

**Theorem 3.26 ([ACJ04, Theorem 4.13]).** A function  $u \in C(\overline{\Omega})$  fulfills  $\Delta_\infty u = 0$  in  $\Omega$  in the viscosity sense iff it is an absolutely minimizing function on  $\Omega$ .

In Problem 3.24 we only consider the case where boundary values are given on  $\partial\Omega$ . However, for our application we would rather have a Dirichlet condition on the set  $\mathcal{O} \subset \overline{\Omega}$ , since it often does not make sense to prescribe boundary values on  $\partial\Omega$ . Therefore, we

assume Neumann boundary conditions on this boundary and then consider the problem

$$\begin{aligned}\Delta_\infty u &= 0 \text{ in } \Omega \setminus \mathcal{O} \\ \frac{\partial u}{\partial \nu} &= 0 \text{ on } \partial\Omega \setminus \mathcal{O} \\ u &= g \text{ on } \mathcal{O}.\end{aligned}\tag{3.13}$$

In the case that  $\Omega$  is smooth and convex, from [ASS11, Lem. 3.1] we have that  $u \in C(\overline{\Omega})$  is an AMLE of  $g : \mathcal{O} \rightarrow \mathbb{R}$  to  $\overline{\Omega}$  iff it fulfills Eq. (3.13).

We can now also address the question of uniqueness, where we state the following famous result of Jensen [Jen93].

**Theorem 3.27** ([Jen93]). Let  $u, v \in C(\overline{\Omega})$  be two viscosity solutions of the  $\infty$ -Laplacian, then we have

$$\max_{\overline{\Omega}}(u - v) = \max_{\partial\Omega}(u - v).$$

Originally proven by Jensen in [Jen93], there exists a very simple proof by Amstrong and Smart [AS10] employing the *comparison with cones* which we detail in the following paragraph.

**Comparison with Cones** As shown in [ACJ04] the concept of absolutely minimizing extensions is equivalent to so-called *comparison with cones*. For some metric  $d(\cdot, \cdot)$  a cone function is defined as

$$x \mapsto a d(x, z) + c$$

where  $z \in \Omega$  denotes the origin or cone tip,  $a \geq 0$  its opening angle and  $c \in \mathbb{R}$  some offset. The property we consider in the following basically asks, if a function  $u$  is smaller than a cone on the boundary of some set  $V$  not including the tip  $z$ , then it should be smaller also in the interior of the domain. This means for any  $a \geq 0, c \in \mathbb{R}$  and  $z \notin V$  we have the implication

$$[u \leq a d(\cdot, z) + c \text{ on } \partial V] \quad \Rightarrow \quad [u \leq a d(\cdot, z) + c \text{ in } V].$$

One can omit the explicit use of the offset  $c \in \mathbb{R}$  and equivalently consider the property

$$\max_{\partial V}(u - a d(\cdot, z)) = \max_V(u - a d(\cdot, z))\tag{3.14}$$

which is reminiscent of the comparison principle in Theorem 3.26. In the Euclidean case cone functions are infinity harmonic away from their cone tip, since they are constant in

the direction of their gradient. We make the explicit computation below

$$\begin{aligned}\partial_{x_i} |x - z| &= \frac{x_i - z_i}{|x - z|} \\ \partial_{x_i x_j} |x - z| &= -\frac{(x_i - z_i)(x_j - z_j)}{|x - z|^3} \text{ for } i \neq j, \\ \partial_{x_i}^2 |x - z| &= \frac{1}{|x - z|^3} \sum_{j \neq i} (x_j - z_j)^2, \\ \Rightarrow \Delta_\infty |x - z| &= \sum_i \frac{(x_i - z_i)^2}{|x - z|^5} \sum_{j \neq i} (x_j - z_j)^2 \\ &\quad - \sum_{i \neq j} \frac{(x_i - z_i)^2 (x_j - z_j)^2}{|x - z|^5} = 0.\end{aligned}$$

In the Euclidean case as in [ACJ04] one says a function fulfills *comparison with cones* from above, if it fulfills Eq. (3.14) in  $\Omega$  for every subset  $V \subset \subset \Omega$ , every  $a \geq 0$  and  $z \in \mathbb{R}^n \setminus V$ . We say  $u$  fulfills *comparison with cones* from below if  $-u$  fulfills comparison from above. If it fulfills both comparisons from above and below, we say it fulfills comparison with cones. Here, [ACJ04, Prop. 2.1] shows the following equivalence to being absolutely minimizing, which however mainly focuses on the Euclidean case. Our setting slightly deviates from the Euclidean one, where in [LIP-II] we consider the following notion.

**Definition 3.28 ([LIP-II, Def. 4.1]).** We shall say that a upper semicontinuous function  $u \in USC(\bar{\Omega})$  satisfies CDF from above in  $\bar{\Omega} \setminus \mathcal{O}$ , if for each relatively open and connected subset  $V \subset \bar{\Omega} \setminus \mathcal{O}$ , any  $x_0 \in \bar{\Omega} \setminus V$  and  $a \geq 0$  we have

$$\max_{\bar{V}}(u - a d_{\bar{\Omega}}(x_0, \cdot)) = \max_{\partial^{\text{rel}} V}(u - a d_{\bar{\Omega}}(x_0, \cdot)).$$

Similarly, we say  $u \in LSC(\bar{\Omega})$  satisfies CDF from below in  $\bar{\Omega} \setminus \mathcal{O}$ , if for each relatively open and connected subset  $V \subset \bar{\Omega} \setminus \mathcal{O}$ , any  $x_0 \in \bar{\Omega} \setminus V$  and  $a \geq 0$  we have

$$\min_{\bar{V}}(u + a d_{\bar{\Omega}}(x_0, \cdot)) = \min_{\partial^{\text{rel}} V}(u + a d_{\bar{\Omega}}(x_0, \cdot)).$$

We say  $u \in C(\bar{\Omega})$  satisfies CDF if it satisfies CDF from above and below.

This definition is a special case of the notion in [JS06]. Therein, the authors also show the equivalence to the absolutely minimizing property, which we state in the following.

**Theorem 3.29 ([JS06, Prop. 4.1]).** A function  $u \in C(\bar{\Omega})$  is absolutely minimizing iff it fulfills comparison with cones in  $\Omega$ .

Therefore, we can rewrite Problem 3.21 employing the comparison with cones condition. Our metric space is given as  $(\bar{\Omega}, d_{\bar{\Omega}})$ , where  $\Omega \subset \mathbb{R}^n$  is some open set in the Euclidean sense, therefore we have a separable length space, where Theorem 3.23 yields existence

of solutions for [Problem 3.21](#). Concerning uniqueness, we refer to [\[Per+09\]](#) which were the first to prove uniqueness in the length space setting. For a more general result we refer to [\[NS12\]](#). In our work we obtain uniqueness by a careful adaption of the proof in [\[AS10\]](#), which we state here for completeness.

**Proposition 3.1 ([LIP-II, Prop. 4.11]).** Consider the metric space  $(\bar{\Omega}, d_{\bar{\Omega}})$ , then [Problem 3.21](#) has at most one solution.

### 3.2.2. Graph Lipschitz Extensions

We now consider the limit  $p \rightarrow \infty$  of [Problem 3.9](#) in the graph case. Analogously to [Section 3.2.1](#) we derive

$$\lim_{p \rightarrow \infty} \left( \mathbf{E}_p^{w_n}(\mathbf{u}) \right)^{1/p} = \max_{x, y \in \Omega_n} w_n(x, y) |\mathbf{u}(y) - \mathbf{u}(x)| =: \mathbf{E}_{\infty}^{w_n}(\mathbf{u})$$

which extends the graph  $p$ -Dirichlet energy to the case  $p = \infty$ . Again we notice structural similarities to the continuum version  $\mathcal{E}_{\infty}$ .

**Remark 3.30.** Informally speaking the functional  $\mathbf{E}_{\infty}^{w_n}$  combines elements of a gradient and a Lipschitz constant. Assuming that  $w_n(x, y)$  relates to  $1/|x - y|$  we see that the finite difference approximation resembles a Lipschitz constant. However, in the limit  $n \rightarrow \infty$  the weighting  $w_n(x, y)$  has a localizing property which fits the interpretation of a gradient better.  $\triangle$

This functional now yields the graph Lipschitz extension problem.

**Problem 3.31 (Graph Energy Minimization).** Given a weighted graph  $(\Omega_n, w_n)$  and a labeling function  $\mathbf{g} : \mathcal{O}_n \rightarrow \mathbb{R}$ , for  $\mathcal{O}_n \subset \Omega_n$  we consider the problem

$$\min_{\mathbf{u}: \Omega_n \rightarrow \mathbb{R}} \mathbf{E}_{\infty}^{w_n}(\mathbf{u}) \text{ subject to } \mathbf{u}(x) = \mathbf{g}(x) \text{ for all } x \in \mathcal{O}_n.$$

Since the weighting function  $w_n : \Omega_n \times \Omega_n \rightarrow \mathbb{R}_0^+$  does not induce a metric, [Problem 3.31](#) does not directly fit the framework of the abstract Lipschitz Extension in [Problem 3.18](#). However, we can consider paths in  $(\Omega_n, w_n)$ , connecting arbitrary  $x, y \in \Omega_n$  i.e.  $\gamma \in \Omega_n^{\times k}, x = \gamma_1, \dots, \gamma_k = y$  such that

$$w_n(\gamma_i, \gamma_{i+1}) > 0 \quad \text{for all } i = 1, \dots, k-1,$$

for which we define the length as

$$|\gamma| = \sum_{i=1}^{k-1} w_n(\gamma_i, \gamma_{i+1})^{-1}.$$

This yields the metric space  $(\Omega_n, d_{w_n})$ , where  $d_{w_n} : \Omega_n \times \Omega_n \rightarrow \mathbb{R}$  is defined as

$$d_{w_n}(x, y) := \min \{|\gamma| : \gamma \text{ is a path in } (\Omega_n, w_n) \text{ from } x \text{ to } y\}. \quad (3.15)$$

**Remark 3.32.** We note that it is important to only consider non-negative weights, otherwise any loop with a negative “length” would decrease the length of the whole path arbitrarily. However, restricting ourselves to non-negative weights we can easily see that the minimum in Eq. (3.15) is indeed attained.  $\triangle$

With this definition we can consider the Lipschitz extension task of  $g : \mathcal{O}_n \rightarrow \mathbb{R}$  to  $\Omega_n$  within the metric space  $(\Omega_n, d_{w_n})$ , i.e. within the setting of Problem 3.18. Therefore the question arises, whether the minimization problem in Problem 3.31 is equivalent to the metric Lipschitz extension problem for which we have the following lemma.

**Lemma 3.33.** For a graph  $(\Omega_n, w_n)$  with non-negative weights and a function  $\mathbf{u} : \Omega_n \rightarrow \mathbb{R}$  we have that

$$\mathbf{E}_{\infty}^{w_n}(\mathbf{u}) = \text{Lip}_{d_{w_n}}(\mathbf{u}).$$

Furthermore, for  $\mathcal{O}_n \subset \Omega_n$  and a function  $\mathbf{g} : \mathcal{O}_n \rightarrow \mathbb{R}$  and we have that

$$\mathbf{g} = \mathbf{u} \text{ on } \mathcal{O}_n \Rightarrow \text{Lip}_{d_{w_n}}(\mathbf{g}; \mathcal{O}_n) \leq \text{Lip}_{d_{w_n}}(\mathbf{u}).$$

*Proof.* **Step 1:** We show that  $\text{Lip}_{d_{w_n}}(\mathbf{u}) \leq \mathbf{E}_{\infty}^{w_n}(\mathbf{u})$ .

We can choose a path  $\gamma \in \Omega_n^{\times k}$  such that

$$\text{Lip}_{d_{w_n}}(\mathbf{u}) = \frac{\mathbf{u}(\gamma_1) - \mathbf{u}(\gamma_k)}{|\gamma|}.$$

The path  $\gamma$  allows to compare vertices  $\gamma_1, \gamma_k \in \Omega_n$  that are not necessarily neighbors in the graph. However, each consecutive vertices in the path are neighbors in the graph and therefore we have

$$w_n(\gamma_i, \gamma_{i+1}) |\mathbf{u}(\gamma_{i+1}) - \mathbf{u}(\gamma_i)| \leq \mathbf{E}_{\infty}^{w_n}(\mathbf{u}) \quad \text{for all } i = 1, \dots, k-1. \quad (3.16)$$

We now employ an elementary result for numbers  $a_i \in \mathbb{R}_0^+, b_i \in \mathbb{R}^+, i = 1, \dots, m \in \mathbb{N}$ , namely

$$[a_i \cdot b_i \leq c \in \mathbb{R} \quad \text{for } i = 1, \dots, m] \quad \Rightarrow \quad \frac{\sum_{i=1}^m a_i}{\sum_{i=1}^m b_i^{-1}} \leq c \quad (3.17)$$

which can be seen as follows

$$\begin{aligned} a_i \cdot b_i &\leq c \quad \text{for } i = 1, \dots, m, \\ \Rightarrow a_i &\leq b_i^{-1} \cdot c \quad \text{for } i = 1, \dots, m, \\ \Rightarrow \sum_{i=1}^m a_i &\leq \left( \sum_{i=1}^m b_i^{-1} \right) \cdot c, \\ \Rightarrow \frac{\sum_{i=1}^m a_i}{\sum_{i=1}^m b_i^{-1}} &\leq c. \end{aligned}$$

This then yields

$$\frac{|\mathbf{u}(\gamma_1) - \mathbf{u}(\gamma_k)|}{|\gamma|} \leq \frac{\sum_{i=1}^{k-1} |\mathbf{u}(\gamma_i) - \mathbf{u}(\gamma_{i+1})|}{|\gamma|} = \frac{\sum_{i=1}^{k-1} |\mathbf{u}(\gamma_i) - \mathbf{u}(\gamma_{i+1})|}{\sum_{i=1}^{k-1} w_n(\gamma_i, \gamma_{i+1})^{-1}} \leq \mathbf{E}_{\infty}^{w_n}(\mathbf{u})$$

where in the last inequality we employed Eq. (3.17) together with Eq. (3.16).

**Step 2:** We show that  $\text{Lip}_{d_{w_n}}(\mathbf{u}) \geq \mathbf{E}_{\infty}^{w_n}(\mathbf{u})$ .

Let  $x, y \in \Omega_n$ , then we know that  $d_w(x, y) \leq w_n(x, y)^{-1}$  and therefore

$$|\mathbf{u}(x) - \mathbf{u}(y)| w_n(x, y) \leq \frac{|\mathbf{u}(x) - \mathbf{u}(y)|}{d_w(x, y)} \leq \max_{\bar{x}, \bar{y} \in \Omega_n} \frac{|\mathbf{u}(\bar{x}) - \mathbf{u}(\bar{y})|}{d_w(\bar{x}, \bar{y})} = \text{Lip}_{d_{w_n}}(\mathbf{u}).$$

Since this holds for arbitrary  $x, y \in \Omega_n$  we have that

$$\mathbf{E}_{\infty}^{w_n}(\mathbf{u}) = \max_{x, y \in \Omega_n} |\mathbf{u}(x) - \mathbf{u}(y)| w_n(x, y) \leq \text{Lip}_{d_{w_n}}(\mathbf{u}).$$

**Step 3:** We show that  $\text{Lip}_{d_{w_n}}(\mathbf{g}; \mathcal{O}_n) \leq \text{Lip}_{d_{w_n}}(\mathbf{u})$ .

If  $\mathbf{g} = \mathbf{u}$  on  $\mathcal{O}$  this simply follows since the maximum for the Lipschitz constant of  $\mathbf{u}$  is taken over a larger set. Indeed, we have that

$$\begin{aligned} \text{Lip}_{d_{w_n}}(\mathbf{u}) &= \max_{x, y \in \Omega_n} \frac{|\mathbf{u}(x) - \mathbf{u}(y)|}{d_w(x, y)} \geq \max_{x, y \in \mathcal{O}_n} \frac{|\mathbf{u}(x) - \mathbf{u}(y)|}{d_w(x, y)} = \max_{x, y \in \mathcal{O}_n} \frac{|\mathbf{g}(x) - \mathbf{g}(y)|}{d_w(x, y)} \\ &= \text{Lip}_{d_{w_n}}(\mathbf{u}; \mathcal{O}_n). \end{aligned}$$

□

This lemma shows that the abstract Lipschitz extension task considered on the metric space  $(\Omega_n, d_{w_n})$  and the Graph  $\infty$ -Dirichlet minimization task are indeed equivalent. Therefore, we also have the Whitney and McShane extensions

$$\begin{aligned} \bar{\mathbf{g}}(x) &= \inf_{y \in \mathcal{O}_n} \mathbf{g}(y) + d_{w_n}(x, y) \\ \underline{\mathbf{g}}(x) &= \sup_{y \in \mathcal{O}_n} \mathbf{g}(y) - d_{w_n}(x, y) \end{aligned}$$

as solutions on the graph. Analogously, the problem does not admit unique solutions.

**Absolutely Minimizing Graph Extensions** Similarly to Section 3.2.1 we can now consider absolutely minimizing extensions. However, the problem in Problem 3.21 uses a notion of a boundary and it is not directly clear how to infer this concept to the discrete set  $\Omega_n$ . Therefore, we define the following what we mean by “boundary” on a graph.

**Definition 3.34.** Let  $(\Omega_n, w_n)$  be a weight graph and let  $V \subset \Omega_n$  be a subset, then we define

- the **exterior** boundary as  $\partial^{\text{ext}} := \{x \in \Omega_n \setminus V : w_n(x, y) > 0 \text{ for some } y \in V\}$ ,

- the **interior** boundary as  $\partial^{\text{int}} := \{x \in V : w_n(x, y) > 0 \text{ for some } y \in \Omega_n \setminus V\}$ .

The closure of  $V$  is then defined as  $\overline{V}^{\text{ext}} := V \cup \partial^{\text{ext}}$  and the interior as  $\overset{\circ}{V}^{\text{int}} := V \setminus \partial^{\text{int}}V$ .

We note that it is not possible to define a topology on  $\Omega_n$  that would yield the above notions. Namely, the only admissible topology in our case would be the discrete topology, i.e.,  $2^{\Omega_n}$ . However, in this topology the only closed sets are  $\emptyset$  and  $\Omega_n$  which is not useful for the applications in the following. Using the Kuratowski closure axioms [Kur22] we remark that the exterior closure on the graph is a so called pre- or Čech closure, [ČFK66]. For a set  $\mathcal{X}$ , a mapping  $\text{cl} : 2^{\mathcal{X}} \rightarrow 2^{\mathcal{X}}$  is called Čech closure, if the following conditions hold,

- $\text{cl}(\emptyset) = \emptyset$  ( $\text{cl}$  preserves the empty set),
- $V \subset \text{cl}(V)$  for all  $V \subset \mathcal{X}$  ( $\text{cl}$  is extensive),
- $\text{cl}(V_1 \cup V_2) = \text{cl}(V_1) \cup \text{cl}(V_2)$  for all  $V_1, V_2 \subset \mathcal{X}$  ( $\text{cl}$  preserves binary unions).

**Lemma 3.35.** The exterior closure on a weighted graph  $(\Omega_n, w_n)$  is a preclosure or Čech closure.

*Proof.* We first see that  $\overline{\emptyset}^{\text{ext}} = \emptyset$  and that  $V \subset \overline{V}^{\text{ext}}$  for every subset  $V \subset \Omega_n$ , i.e. the above defined closure preserves the empty set and is extensive. Furthermore, for two sets  $V_1, V_2 \subset \Omega_n$  we have that

$$\begin{aligned} & x \in \partial^{\text{ext}}(V_1 \cup V_2) \\ \Leftrightarrow & [x \notin V_1 \cup V_2] \wedge [\exists y \in V_1 \cup V_2 : w_n(x, y)] \neq 0 \\ \Leftrightarrow & [x \notin V_1 \cup V_2] \wedge \left( [\exists y \in V_1 : w_n(x, y) \neq 0] \vee [\exists y \in V_2 : w_n(x, y) \neq 0] \right) \\ \Leftrightarrow & \left[ x \in \partial^{\text{ext}}V_1 \setminus V_2 \right] \vee \left[ x \in \partial^{\text{ext}}V_2 \setminus V_1 \right] \\ \Leftrightarrow & x \in (\partial^{\text{ext}}V_1 \cup \partial^{\text{ext}}V_2) \setminus (V_1 \cup V_2). \end{aligned}$$

We have shown that  $\partial^{\text{ext}}(V_1 \cup V_2) = (\partial^{\text{ext}}V_1 \cup \partial^{\text{ext}}V_2) \setminus (V_1 \cup V_2)$ . Therefore, we have that

$$\begin{aligned} \overline{V_1 \cup V_2}^{\text{ext}} &= V_1 \cup V_2 \cup \partial^{\text{ext}}(V_1 \cup V_2) \\ &= V_1 \cup V_2 \cup ((\partial^{\text{ext}}V_1 \cup \partial^{\text{ext}}V_2) \setminus (V_1 \cup V_2)) \\ &= V_1 \cup \partial^{\text{ext}}V_1 \cup V_2 \cup \partial^{\text{ext}}V_2 \\ &= \overline{V_1}^{\text{ext}} \cup \overline{V_2}^{\text{ext}}. \end{aligned}$$

This shows that the closure preserves binary unions and therefore we have shown, that it is indeed a Čech closure.  $\square$

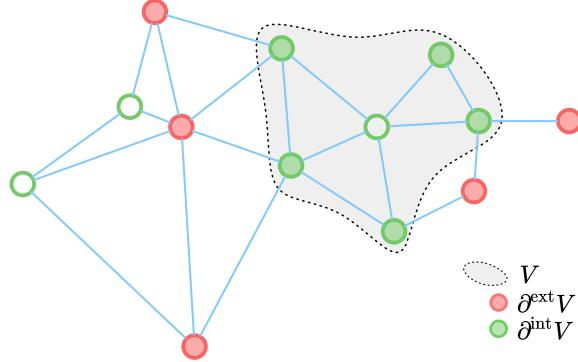


Figure 3.7.: Visualization of exterior and interior boundary on a graph.

The missing property, that inhibits the closure to induce a topology is the so-called idempotence. Namely, there are sets  $V \subset \Omega_n$  such that

$$\overline{V}^{\text{ext}} \neq \overline{\overline{V}^{\text{ext}}}^{\text{ext}}.$$

In the example visualized in Fig. 3.7 we see that  $\overline{\overline{V}^{\text{ext}}}^{\text{ext}} = \Omega_n \neq \overline{V}^{\text{ext}}$ . Since the closure we employ here does not induce a topology, we have a slightly modified notion of absolutely minimizers.

**Problem 3.36 (Graph AMLEs).** Given a connected weighted graph  $(\Omega_n, w_n)$ ,  $\mathcal{O}_n \subset \Omega_n$  and a function  $\mathbf{g} : \mathcal{O}_n \rightarrow \mathbb{R}$ , find a function  $\mathbf{u} : \Omega_n \rightarrow \mathbb{R}$  such that

$$\begin{aligned} \text{Lip}_{d_{w_n}}(\mathbf{u}; \overline{V}^{\text{ext}}) &= \text{Lip}_{d_{w_n}}(\mathbf{u}; \partial^{\text{ext}} V) \text{ for all connected } V \subset \Omega_n \setminus \mathcal{O}_n, \\ \mathbf{u} &= \mathbf{g} \text{ on } \mathcal{O}_n. \end{aligned}$$

This notion of absolute minimizers is employed in [LIP-II; LIP-III]. We see, that graph AMLEs are indeed special solutions of the basic Lipschitz extension problem on the graph, by choosing  $V = \Omega_n \setminus \mathcal{O}_n$  in the above problem.

**Comparison with Graph Distance functions** Analogously to the continuum case Section 3.2.1 we can also consider comparison with distance functions, but now on graphs. The main ingredients here are the graph distance function  $d_{w_n}$  and the notion of closure on a graph as developed in the last section.

**Definition 3.37.** For a weighted graph  $(\Omega_n, w_n)$  we say a function  $\mathbf{u} : \Omega_n \rightarrow \mathbb{R}$  fulfills comparison with distance function from above (CDFA) on a subset  $U \subset \Omega_n$  if for every  $V \subset U$  we have

$$\max_{\overline{V}^{\text{ext}}} (u - a d_{w_n}(\cdot, z)) = \max_{\partial^{\text{ext}} V} (u - a d_{w_n}(\cdot, z)) \quad (\text{CDFA})$$

for every  $z \in \Omega_n \setminus V$  and every  $a \geq 0$ . We say that  $\mathbf{u}$  fulfills comparison with distance function from below (CDFB) on a subset  $U \subset \Omega_n$  if for every  $V \subset U$  we have

$$\min_{V^{\text{ext}}} (u + a d_{w_n}(\cdot, z)) = \min_{\partial^{\text{ext}} V} (u + a d_{w_n}(\cdot, z)) \quad (\text{CDFB})$$

for every  $z \in \Omega_n \setminus V$  and every  $a \geq 0$ .

Analogously to the continuum case, we say that a function fulfills comparison with distance functions, if it fulfills both, (CDFA) and (CDFB).

**The Graph  $\infty$ -Laplacian** We can also obtain the limit of the Graph  $p$ -Laplace operator via the following formal calculation,

$$\begin{aligned} & \Delta_p^{w_n} \mathbf{u}(x) = 0 \\ \Leftrightarrow & \sum_{y \in \Omega_n} w_n(x, y)^p |\mathbf{u}(y) - \mathbf{u}(x)|^{p-2} (\mathbf{u}(y) - \mathbf{u}(x)) = 0 \\ \Leftrightarrow & \sum_{y: \mathbf{u}(x) \leq \mathbf{u}(y)} w_n(x, y)^p (\mathbf{u}(y) - \mathbf{u}(x))^{p-1} = \sum_{y: \mathbf{u}(x) > \mathbf{u}(y)} w_n(x, y)^p (\mathbf{u}(x) - \mathbf{u}(y))^{p-1}. \end{aligned}$$

Taking the terms on the left and right hand side to the power of  $1/p$  and then formally sending  $p \rightarrow \infty$  then yields

$$\begin{aligned} & \max_{y: \mathbf{u}(x) \leq \mathbf{u}(y)} w_n(x, y) (\mathbf{u}(y) - \mathbf{u}(x)) = \max_{y: \mathbf{u}(x) > \mathbf{u}(y)} w_n(x, y) (\mathbf{u}(x) - \mathbf{u}(y)) \\ \Leftrightarrow & \max_{y \in \Omega_n} w_n(x, y) (\mathbf{u}(y) - \mathbf{u}(x)) = - \min_{y \in \Omega_n} w_n(x, y) (\mathbf{u}(y) - \mathbf{u}(x)). \end{aligned}$$

This calculation motivates the definition of the graph infinity Laplacian

$$\Delta_\infty^{w_n} \mathbf{u}(x) := \max_{y \in \Omega_n} w_n(x, y) (\mathbf{u}(y) - \mathbf{u}(x)) + \min_{y \in \Omega_n} w_n(x, y) (\mathbf{u}(y) - \mathbf{u}(x)),$$

which then allows to formulate the associated problem as an extension of [Problem 3.11](#).

**Problem 3.38 (Graph  $\infty$ -Laplacian).** Given a weighted graph  $(\Omega_n, w_n)$  and a labeling function  $\mathbf{g} : \mathcal{O}_n \rightarrow \mathbb{R}$  with  $\mathcal{O}_n \subset \Omega_n$ , find a function  $\mathbf{u} : \Omega_n \rightarrow \mathbb{R}$  such that

$$\begin{aligned} & \Delta_\infty^{w_n} \mathbf{u} = 0, \text{ in } \Omega_n \setminus \mathcal{O}_n, \\ & \mathbf{u} = \mathbf{g} \text{ on } \mathcal{O}_n. \end{aligned}$$

To establish existence of solutions of the problem above one can employ the Perron method [Per23]. Alternatively, one can establish the equivalence to so-called lex-minimizers and the show existence as in [Kyn+15, Th. 3.3]. Uniqueness is a consequence of the following theorem from [Cal19], which is similar to the uniqueness result of Jensen in [Theorem 3.27](#).

**Theorem 3.39** ([Cal19, Th. 3.1]). Let  $(\Omega_n, w_n)$  be a graph connected to the boundary  $\mathcal{O}_n \subset \Omega_n$  and consider  $\mathbf{u}, \mathbf{v} : \Omega_n \rightarrow \mathbb{R}$  such that

$$\Delta_\infty^{w_n} \mathbf{u} \leq 0 \leq \Delta_\infty^{w_n} \mathbf{v} \quad \text{in } \Omega_n \setminus \mathcal{O}_n.$$

Then we know that

$$\max_{\Omega_n} (\mathbf{u} - \mathbf{v}) = \max_{\mathcal{O}_n} (\mathbf{u} - \mathbf{v})$$

**Relations Between the Graph Lipschitz Extensions** We now establish the connection between the different notions of Lipschitz extensions on the graph. Compared to the continuum case we do not establish the full equivalences but only the necessary implications required for the convergence proofs in [LIP-II]. Namely, in [LIP-II] we show that if  $\mathbf{u}$  solves the graph  $\infty$ -Laplace, then it is a graph AMLE and fulfills comparison with graph distance functions.

**Lemma 3.40** ([LIP-II, Th. 3.2, Prop. 3.8]). Let  $(\Omega_n, w_n)$  be a weighted connected graph and  $\mathbf{g} : \mathcal{O}_n \rightarrow \mathbb{R}$  be a given function for  $\mathcal{O}_n \subset \Omega_n$ . Furthermore, let  $\mathbf{u} : \Omega_n \rightarrow \mathbb{R}$  be graph infinity harmonic on  $\Omega_n \setminus \mathcal{O}_n$  with boundary conditions given by  $\mathbf{g}$ , i.e.,  $\mathbf{u}$  solves Problem 3.38 then we have that

- $\mathbf{u}$  is a graph AMLE, i.e.,  $\mathbf{u}$  solves Problem 3.36,
- $\mathbf{u}$  fulfills comparison with cones.

*Proof.* Both of the statements are proven in [LIP-II]. Let  $\mathbf{u}$  be such that

$$\begin{aligned} \Delta_\infty^{w_n} \mathbf{u} &= 0 \text{ in } \Omega_n \setminus \mathcal{O}_n \\ \mathbf{u} &= \mathbf{g} \text{ on } \mathcal{O}_n. \end{aligned}$$

From [LIP-II, Prop. 3.8] we have that  $\mathbf{u}$  is a graph AMLE. Furthermore, from [LIP-II, Th. 3.2] we have that  $\mathbf{u}$  fulfills comparison with cones. In fact, [LIP-II, Th. 3.2], shows a more refined statement, namely that

$$\begin{aligned} -\Delta_\infty^{w_n} \mathbf{u} \leq 0 &\Rightarrow \mathbf{u} \text{ fulfills (CDFA),} \\ -\Delta_\infty^{w_n} \mathbf{u} \geq 0 &\Rightarrow \mathbf{u} \text{ fulfills (CDFB).} \end{aligned}$$

□

### 3.3. Gamma Convergence: [LIP-I]

We now present the main results of [LIP-I]. This work considers the limit of the basic Lipschitz extension from the graph Problem 3.31 to the continuum case Problem 3.18. As

detailed in the previous sections both of these problems do not admit unique solutions, however it is still meaningful to study the convergence in the infinite data limit. We first start by specifying the setting and recalling the concept of  $\Gamma$ -convergence and then state the main results.

**Main Contributions in [LIP-I]** We first prove  $\Gamma$ -convergence of the functionals  $\mathbf{E}_\infty^{w_n}$  to their continuum counterpart  $\mathcal{E}_\infty$ . This result can be seen as the extension of the results in [GS15; ST19] to the case  $p = \infty$ . Here, we employ a different type of metric as we detail in the following. Furthermore, we identify the optimal graph scaling which in this case is derived in a deterministic manner. The key difference to [GS15; ST19] is that we can work with any point-clouds and not only i.i.d. points. The inherent reason is that  $L^\infty$  problems only consider mass in a qualitative not in a quantitative way. We then establish a compactness result, that allows us to infer convergence of minimizers. Finally, we apply this theory to the ground state problem.

Parts of this work was motivated by the findings in TR's master thesis [Roi21]. The convergence result in the latter, was however not fully established and therefore [LIP-I] constitutes a non-trivial expansion and also correction. We want to highlight two major differences:

- The domain  $\Omega$  in [Roi21] was assumed to be convex. This restriction was omitted in [LIP-I] by identifying a class of locally convex domains, that prevent internal sharp corners.
- We are able to work with asymptotic boundary conditions, i.e., the boundary conditions for each graph problem do need to be the same as in the continuum. They only need to approximate them in some sense.

### 3.3.1. Setting and Preliminaries

Here, we detail the concrete setting of [LIP-I] and provide some background on the employed notions.

**$\Gamma$ -Convergence** Originally, the concept of  $\Gamma$ -convergence dates back to De Giorgi [DF75] as a type of variational convergence. We refer to [Bra02; Dal12] for a detailed overview on this notion and related topics. While  $\Gamma$ -convergence was successfully employed in a pure continuum setting for a longer time (see e.g. [Mod77]), it was more recently used to prove convergence from a discrete to a continuum functional [CGL10; BY12; VB+12].

**Definition 3.41 ([DF75]:  $\Gamma$ -convergence).** Let  $X$  be a metric space and let  $F_n : X \rightarrow [-\infty, \infty]$  be a sequence of functionals. We say that  $F_n$   $\Gamma$ -converges to the functional  $F : X \rightarrow [-\infty, \infty]$  if

- (i) (**liminf inequality**) for every sequence  $(x_n)_{n \in \mathbb{N}} \subset X$  converging to  $x \in X$  we have that

$$\liminf_{n \rightarrow \infty} F_n(x_n) \geq F(x);$$

- (ii) (**limsup inequality**) for every  $x \in X$  there exists a sequence  $(x_n)_{n \in \mathbb{N}} \subset X$  converging to  $x$  and

$$\limsup_{n \rightarrow \infty} F_n(x_n) \leq F(x).$$

The most relevant reference for [LIP-I] was the disruptive work presented by García Trillos and Slepčev in [GS15], which is also commented on in Section 3.1. Among other important ideas and notions, we want to highlight two ingredients that directly influenced [LIP-I]:

- $\Gamma$ -convergence on a common metric space, that allows to compare graph functions with continuum functions.
- The proof strategy of first considering the discrete to non-local convergence and then non-local to continuum.

Employing  $\Gamma$ -convergence the authors in [ST19] were able to prove convergence of minimizers of the graph  $p$ -Dirichlet problem. Here, one can employ the convenient result that  $\Gamma$ -convergence implies convergence of minimizers.

**Lemma 3.42 ([Bra02, Thm. 1.21] Convergence of Minimizers).** Let  $X$  be a metric space and  $F_n : X \rightarrow [0, \infty]$  a sequence of functionals  $\Gamma$ -converging to  $F : X \rightarrow [0, \infty]$  which is not identically  $+\infty$ . If there exists a relatively compact sequence  $(x_n)_{n \in \mathbb{N}}$  such that

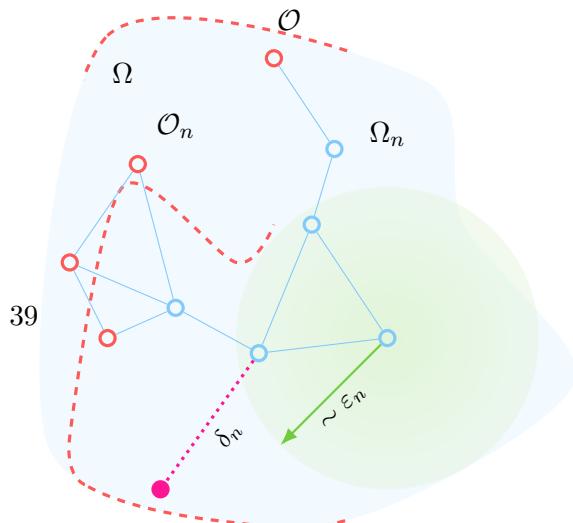
$$\lim_{n \rightarrow \infty} \left( F_n(x_n) - \inf_{x \in X} F_n(x) \right) = 0,$$

then we have that

$$\lim_{n \rightarrow \infty} \inf_{x \in X} F_n(x) = \min_{x \in X} F(x)$$

and any cluster point of  $(x_n)_{n \in \mathbb{N}}$  is a minimizer of  $F$ .

**The Kernel and the Graph Scale** In Section 3.1 we comment on the relevance of the kernel and graph scaling  $\varepsilon_n$  for continuum limits. We first state the concrete assumptions we make in [LIP-I] on the kernel  $\eta : \mathbb{R}_0^+ \rightarrow \mathbb{R}_0^+$ , namely we require



(K1)  $\eta$  is positive and continuous at 0,

(K2)  $\eta$  is non-increasing,

(K3)  $\text{supp}(\eta) \subset [0, t_\eta]$  for some  $t_\eta > 0$ .

As seen in [GS15; ST19] and Eq. (3.1)

we also need to consider the constant  $\sigma_\eta$   
which in our case is given by

$$\sigma_\eta = \underset{t \in \mathbb{R}^+}{\text{ess sup}} \eta(t) t = \lim_{p \rightarrow \infty} \left( \int_{\mathbb{R}^+} \eta(t)^p t^{d+p} dt \right)^{1/p}.$$

A main difference to the results for  $p < \infty$  is that we can work with point clouds  $\Omega_n$ , that do not necessarily need be constructed by an i.i.d. sampling process. However, this rises the question how we can ensure that  $\Omega_n \subset \bar{\Omega}$  fills out the continuum domain  $\bar{\Omega}$  sufficiently fast. In [LIP-I] we consider the maximum distance of any continuum point to its next vertex

$$\sup_{x \in \bar{\Omega}} \min_{y \in \Omega_n} |x - y|, \quad (3.18)$$

as a measure of the distance between  $\Omega_n$  and  $\bar{\Omega}$ . In fact, this is the Hausdorff distance ([Hau14]) between the two sets, since

$$d_H(\Omega_n, \bar{\Omega}) = \max \left\{ \sup_{x \in \bar{\Omega}} \inf_{y \in \Omega_n} |x - y|, \sup_{x \in \Omega_n} \inf_{y \in \bar{\Omega}} |x - y| \right\} = \sup_{x \in \bar{\Omega}} \min_{y \in \Omega_n} |x - y|,$$

where the second term vanishes since  $\Omega_n \subset \bar{\Omega}$ .

**Remark 3.43.** We remark on the historical correctness of the name ‘‘Hausdorff’’ distance, motivated by the review in [BT06]. This distance can actually be contributed to Pompeiu, who employed a similar notion in [Pom05], which also relates to concepts developed in Fréchet’s Ph.D. thesis [Fré06]. Hausdorff modified this distance into the presented form in [Hau14].  $\triangle$

We also allow the constraint set  $\mathcal{O}_n$  to vary with each  $n \in \mathbb{N}$ , therefore we also measure the distance between  $\mathcal{O}_n$  and the continuum constraint set  $\mathcal{O} \subset \bar{\Omega}$ , via

$$d_H(\mathcal{O}_n, \mathcal{O}) = \max \left\{ \sup_{x \in \mathcal{O}} \inf_{y \in \mathcal{O}_n} |x - y|, \max_{x \in \mathcal{O}_n} \inf_{y \in \mathcal{O}} |x - y| \right\}$$

where the second term does not vanish, because in general  $\mathcal{O}_n \not\subseteq \mathcal{O}$ . We assume that we have a Lipschitz continuous function  $g : \bar{\Omega} \rightarrow \mathbb{R}$  such that for each  $n$  the boundary conditions are given by the function  $\mathbf{g}_n = g|_{\mathcal{O}_n}$ . In total we then need to control the graph resolution

$$\delta_n := \max\{d_H(\Omega_n, \bar{\Omega}), d_H(\mathcal{O}_n, \mathcal{O})\}.$$

This then allows us to formulating our scaling assumption, namely we require that the graph resolution tends to zero faster than the graph scaling

$$\frac{\delta_n}{\varepsilon_n} \longrightarrow 0, \quad n \rightarrow \infty, \quad (3.19)$$

which has similarly employed in [Pen99a; Pen99b]. This is the weakest possible assumption that ensures connectivity of the graph in the limit  $n \rightarrow \infty$ . However, compared to Eq. (3.7) it is more of deterministic nature, since we do not assume anything on how  $\Omega_n$  is created. In case of a i.i.d. point cloud one can show that

$$\delta_n \sim \left( \frac{\log n}{n} \right)^{1/d}$$

see [Pen99b], which is the setting in [GS15].

**Assumptions on the Domain** As we mention in the introduction, we do not need to restrict ourselves to convex domains. However, for our convergence results, it is important that the geodesic distance

$$d_\Omega(x, y) = \inf \{ \text{len}(\gamma) : \gamma : [a, b] \rightarrow \Omega \text{ is a curve with } \gamma(a) = x, \gamma(b) = y \},$$

converges to the Euclidean distance if  $x$  goes to  $y$ . We formulate this in the following assumption

$$\limsup_{\delta \downarrow 0} \left\{ \frac{d_\Omega(x, y)}{|x - y|} : x, y \in \Omega, |x - y| < \delta \right\} \leq 1. \quad (3.20)$$

The inherent reason that we need to enforce this, is that edges in the graph are allowed to communicate via their direct Euclidean distance, even if this line does not lie in  $\Omega$ . While it would be easier to also consider the geodesic distance on the graph, this is unrealistic in many applications since we do not have access to the concrete shape of  $\Omega$ . Most importantly, Eq. (3.20) prohibits sharp internal corners in the domain, see Fig. 3.8. Furthermore, note that we consider the geodesic distance  $d_\Omega$  on the open domain  $\Omega$ , not on its closure. Therefore, we also exclude situations of a touching boundary as in Fig. 3.4. In [LIP-I] we also provide examples of domains fulfilling this condition.

### 3.3.2. $\Gamma$ -Convergence of the Discrete Functionals

In order to show convergence of discrete functionals defined on  $\Omega_n$  to continuum functions acting on  $\bar{\Omega}$ , one needs to define a common metric space. Employing results from [TS15] in [GS15] the authors introduced the space  $TL^p$

$$TL^p := \{(\mu, u) : \mu \in \mathcal{P}(\Omega), u \in L^p(\mu)\}$$

together with the transport distance

$$d_{TL^p}((\mu, u), (\nu, v)) = \inf_{\pi \in \Gamma(\mu, \nu)} \left( \int_{\Omega \times \Omega} |x - y|^p + |u(x) - v(y)|^p d\pi(x, y) \right)^{1/p}$$

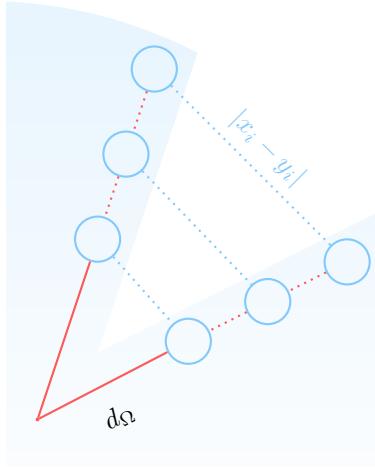


Figure 3.8.: An example of a sharp internal corner, violating Eq. (3.20). The geodesic distance between the pairs of points converging to the corner point is bigger than then Euclidean one by a constant factor.

where  $\Gamma(\mu, \nu)$  is the set of coupling between  $\mu$  and  $\nu$ . A sequence  $(\mu_n, u_n)$  converges to  $(\mu, u)$  in  $TL^p$  iff there exists a sequence of transportation maps  $T_n : \Omega \rightarrow \Omega$  with  $T_n \# \mu = \mu_n$  and

$$\int_{\Omega} |x - T_n(x)| d\mu(x) \rightarrow 0$$

such that  $u_n \circ T_n \xrightarrow{L^p} u$ , [GS15, Prop. 3.12]. Therefore, the maps  $T_n$  allow to employ standard convergence in  $L^p$ . In order to transfer this situation to  $L^\infty$  one could try to employ a  $\infty$ -Wasserstein distance. However, as argued in [Roi21] the strategy does not transfer directly, since convergence in  $W^\infty$  does not metrize weak convergence of measures [San15, Thm. 5.10]. However, as seen in [LIP-I] one can employ a more direct argument. For problems in  $L^p$  the conservation of mass was important for the maps  $T_n$  (i.e.  $T_n \# \mu = \mu_n$ ) such that the integrals could be transformed. This condition is irrelevant in  $L^\infty$ , namely we have the following analogous transformation rule in  $L^\infty$ .

**Lemma 3.44 ([LIP-I, Lem. 2]).** For two probability measures  $\mu, \nu \in \mathcal{P}(\overline{\Omega})$ , a measurable map  $T : \Omega \rightarrow \Omega$  which fulfills

- (i)  $\nu << T \# \mu$ ,
- (ii)  $T \# \mu << \nu$ ,

and for a measurable function  $u : \Omega \rightarrow \mathbb{R}$  we have that

$$\nu - \text{ess sup}_{x \in \Omega} u(x) = \mu - \text{ess sup}_{y \in \Omega} u(T(y)).$$

We want to compare the discrete measure  $\mu_n = \frac{1}{n} \sum_{x \in \Omega_n} \delta_x$  to the target measure  $\mu$ . In this setting a closest point projection  $p_n : \Omega \rightarrow \Omega_n$

$$p_n(x) \in \operatorname{argmin}_{y \in \Omega_n} |x - y|$$

fulfills the assumption of [Lemma 3.44](#). This allows us to extend the functional  $\mathbf{E}_\infty^{w_n}$  in [Problem 3.31](#) to  $L^\infty$  via

$$\mathbf{E}_\infty^{w_n}(u) = \begin{cases} \mathbf{E}_\infty^{w_n}(\mathbf{u}) & \text{if } u = \mathbf{u} \circ p_n, \text{ for some } \mathbf{u} : \Omega_n \rightarrow \mathbb{R}, \\ \infty & \text{else,} \end{cases} \quad (3.21)$$

which was similarly done in [\[GS15; ST19\]](#). Additionally, we incorporate the constraint on  $\mathcal{O}_n$  in [Problem 3.31](#) via

$$\mathbf{E}_\infty^{w_n, \text{cons}}(\mathbf{u}) := \begin{cases} \mathbf{E}(\mathbf{u}) & \text{if } \mathbf{u} = \mathbf{g} \text{ on } \mathcal{O}_n, \\ \infty & \text{else,} \end{cases}$$

with the analogous extension to  $L^\infty$  as in [Eq. \(3.21\)](#). We can now state the first main result of [\[LIP-I\]](#) that provides the  $\Gamma$ -convergence of the graph functionals.

**Theorem 3.45 (Discrete to continuum  $\Gamma$ -convergence).** Let  $\Omega \subset \mathbb{R}^d$  be a domain satisfying [\(3.20\)](#), let the kernel fulfill [\(K1\)-\(K3\)](#), then for any null sequence  $(\varepsilon_n)_{n \in \mathbb{N}} \subset (0, \infty)$  which satisfies the scaling condition [\(3.19\)](#) we have

$$\mathbf{E}_\infty^{n, \text{cons}} \xrightarrow{\Gamma} \sigma_\eta \mathcal{E}^{\text{cons}}. \quad (3.22)$$

The main proof strategy here is similar to the one in [\[GS15; ST19\]](#). Namely one defines the non-local functional

$$\mathcal{E}_\infty^\varepsilon(u) := \frac{1}{\varepsilon} \operatorname{ess\,sup}_{x, y \in \Omega} \{ \eta_\varepsilon(|x - y|) |u(x) - u(y)| \}, \quad \varepsilon > 0$$

for which we show that for any sequence  $\varepsilon_n \rightarrow 0$  we have

$$\mathcal{E}_\infty^{\varepsilon_n} \xrightarrow{\Gamma} \sigma_\eta \mathcal{E}_\infty,$$

see [\[LIP-I, Thm. 4\]](#). For the liminf inequality of the discrete functionals, one has to take special care of points  $x, y \in \Omega_n$  where  $\eta_{\varepsilon_n}(|p_n(x) - p_n(y)|) = 0$ . We want to bound  $\mathbf{E}_\infty^{w_n}$  from below by  $\mathcal{E}_\infty^{\varepsilon_n}$  for which we have to permit significant communication of  $x$  and  $y$  whenever  $p_n(x)$  and  $p_n(y)$  do not communicate. This can be done, (temporarily assuming  $\eta$  is constant on  $[0, t]$ ) by introducing a smaller length scale  $\tilde{\varepsilon}$  such that

- (i)  $|p_n(x) - p_n(y)| > t\varepsilon_n \Rightarrow |p_n(x) - p_n(y)| / \varepsilon_n < |x - y| / \tilde{\varepsilon}$ ,
- (ii)  $\lim_{n \rightarrow \infty} \tilde{\varepsilon} / \varepsilon = 1$ .

As shown in [LIP-I] the choice  $\tilde{\varepsilon} = \varepsilon - 2\delta/t$  fulfills (i). So in order to fulfill (ii) we obtain the scaling condition

$$\lim_{n \rightarrow \infty} \frac{\delta_n}{\varepsilon_n} = 0 \Rightarrow \lim_{n \rightarrow \infty} \frac{\varepsilon_n - 2\delta_n/t}{\varepsilon_n} = 1 - \frac{2}{t} \lim_{n \rightarrow \infty} \frac{\delta_n}{\varepsilon_n} = 1.$$

This then allows to prove the liminf inequality since  $\mathbf{E}_\infty^{w_n}(u_n) \geq \frac{\tilde{\varepsilon}_n}{\varepsilon_n} \mathcal{E}_\infty^{\tilde{\varepsilon}_n}(u_n)$  holds. The limsup inequality can then be shown by choosing the constant sequence, with some additional care for the changing constraint set  $\mathcal{O}_n$ .

### 3.3.3. Convergence of Minimizers

The convenient aspect of  $\Gamma$ -convergence is, that under additional compactness properties it directly shows convergence of minimizers, [Bra02, Thm. 8]. This yields the second main result in [LIP-I].

**Theorem 3.46 ([LIP-I, Thm. 2]).** Let  $\Omega \subset R^d$  be a domain satisfying (3.20), let the kernel fulfil (K1)-(K3), and  $(\varepsilon_n)_{n \in \mathbb{N}} \subset (0, \infty)$  be a null sequence which satisfies the scaling condition (3.19). Then any sequence  $(u)_{n \in \mathbb{N}} \subset L^\infty(\Omega)$  such that

$$\lim_{n \rightarrow \infty} \left( \mathbf{E}_\infty^{n, \text{cons}}(u_n) - \inf_{u \in L^\infty(\Omega)} \mathbf{E}_\infty^{n, \text{cons}}(u) \right) = 0$$

is relatively compact in  $L^\infty(\Omega)$  and

$$\lim_{n \rightarrow \infty} \mathbf{E}_\infty^{n, \text{cons}}(u_n) = \min_{u \in L^\infty(\Omega)} \sigma_\eta \mathcal{E}_\infty^{\text{cons}}(u).$$

Furthermore, every cluster point of  $(u_n)_{n \in \mathbb{N}}$  is a minimizer of  $\mathcal{E}_{\text{cons}}$ .

In order to show the compactness in the above theorem one employs the following lemma.

**Lemma 3.47 ([LIP-I, Lem. 4]).** Let  $(\Omega, \mu)$  be a finite measure space and  $K \subset L^\infty(\Omega; \mu)$  be a bounded set w.r.t.  $\|\cdot\|_{L^\infty(\Omega; \mu)}$  such that for every  $\varepsilon > 0$  there exists a finite partition  $\{V_i\}_{i=1}^n$  of  $\Omega$  into subsets  $V_i$  with positive and finite measure such that

$$\mu\text{-ess sup}_{x, y \in V_i} |u(x) - u(y)| < \varepsilon \quad \forall u \in K, i = 1, \dots, n, \tag{3.23}$$

then  $K$  is relatively compact.

**Remark 3.48.** The proof of this statement employs ideas from [DS88, Lem. IV.5.4] and appeared similarly in TR's master thesis. However, therein the statement was slightly wrong, which was corrected in [LIP-I].  $\triangle$

This lemma is used to show that if a sequence  $u_n$  fulfills

$$\sup_{n \in \mathbb{N}} \mathbf{E}_\infty^{w_n, \text{cons}}(u_n) < \infty$$

then it is relatively compact.

### 3.3.4. Application to Ground States

The convergence results of the previous paragraphs can be applied to the ground states problem. Namely, for a  $p$ -homogeneous functional  $\mathcal{E}$ —i.e.  $\mathcal{E}(cu) = |c|^p \mathcal{E}(u)$  for  $c \in \mathbb{R}$ —we consider the following problem

$$\min \left\{ \mathcal{E}(u) : \inf_{v \in \arg\min \mathcal{E}} \|u - v\| = 1 \right\},$$

as studied in [BKB20]. In our setting, if  $g \equiv 0$  then the functionals  $\mathbf{E}_\infty^{w_n, \text{cons}}$  and  $\mathcal{E}_\infty^{\text{cons}}$  are 1-homogeneous. In [LIP-I, Th. 5] we first show that the up to a global sign unique ground state of  $\mathcal{E}_\infty^{\text{cons}}$  is a multiple of the distance function

$$d_{\mathcal{O}}(x) = \inf_{y \in \mathcal{O}} d_\Omega(x, y).$$

We then employ the  $\Gamma$ -convergence results to show the convergence of discrete ground states to the above continuum one.

**Theorem 3.49 ([LIP-I, Thm. 6]: Convergence of Ground States).** Under the conditions of [LIP-I, Thm. 1] let the sequence  $(u_n)_{n \in \mathbb{N}} \subset L^\infty(\Omega)$  fulfill

$$\mathbf{u}_n \in \arg\min \left\{ \mathbf{E}_\infty^{w_n, \text{cons}}(\mathbf{u}) : \mathbf{u} \in L^\infty(\Omega), \|\mathbf{u}\|_{L^p} = 1 \right\}.$$

Then (up to a subsequence)  $\mathbf{u}_n \rightarrow u$  in  $L^\infty(\Omega)$  where

$$u \in \arg\min \left\{ \mathcal{E}_\infty^{\text{cons}}(u) : u \in L^\infty(\Omega), \|u\|_{L^p} = 1 \right\}$$

and it holds

$$\lim_{n \rightarrow \infty} \mathbf{E}_\infty^{w_n, \text{cons}}(\mathbf{u}_n) = \sigma_\eta \mathcal{E}_\infty^{\text{cons}}(u). \quad (3.24)$$

## 3.4. Uniform Convergence of AMLES: [LIP-II]

In the previous section we consider the limit of the Lipschitz extension task. The major limitations of this work are:

- The considered problem does not admit unique solutions. The  $\Gamma$ -convergence framework does not directly allow to select certain minimizers, e.g. AMLES, and consider the respective limit to continuum AMLES.

- We can only show qualitative convergence results. It is not directly possible to prove quantitative statements and convergence rates.

These drawbacks are the main motivation for the work presented in this section. We first state the main contributions and then provide details below.

**Main Contributions in [LIP-II]** We prove a quantitative convergence results for discrete AMLEs to their continuum counterpart. Convergence was already established in [Cal19], however employing a more restrictive scaling assumption and only allowing for smooth kernels  $\eta$ , which was due to the viscosity techniques employed. By using the AMLE characterization of infinity harmonic functions, we are able to show convergence rates under the weakest scaling assumption and allowing for very general kernels. This is done via a using a comparison principle on a homogenized length scale. The key insight we obtain in this work, is that convergence rates for graph distance functions imply rates for AMLEs. Finally, we examine some numerical convergence rates.

### 3.4.1. Setting

Conceptually, the basic setting is very similar to the previous one in Section 3.3. We highlight some minor differences in the paragraphs below.

**Assumptions on the Kernel** We employ (K2) and (K3) from [LIP-I] and set  $t_\eta = 1$  for simplicity. However we do not need to assume continuity in 0, i.e. (K1). Instead we additionally require the following.

(K4) We assume that there exists  $t_0 \in (0, 1]$ , chosen as the largest number with this property with  $\sigma_\eta = t_0\eta(t_0)$ .

Assuming (K2) to (K4) allows for the so-called singular weights

$$\eta(t) = \frac{1}{t^p} 1_{[0,1]}(t)$$

with  $p \in (0, 1]$ . As noted in [LIP-II] the above assumption is not redundant, since for example the function

$$\eta(t) = \frac{1}{t^{\frac{1}{t-1}+2}} 1_{[0,1]}(t), \quad t > 0,$$

violates (K4), while fulfilling (K2) and (K3).

**Assumptions on the Domain** Again we take the concept from [LIP-I], where we assume asymptotic convexity. We alter the Eq. (3.20) slightly to be more quantitative. We assume that there exists a function  $\phi : [0, \infty) \rightarrow [0, \infty)$  with  $\lim_{h \downarrow 0} \frac{\phi(h)}{h} = 0$  and  $r_\Omega > 0$  such that

$$d_{\bar{\Omega}}(x, y) \leq |x - y| + \phi(|x - y|), \quad \text{for all } x, y \in \Omega : |x - y| \leq r_\Omega. \quad (3.25)$$

Furthermore, for  $\varepsilon > 0$  we define  $\sigma_\phi(\varepsilon) = \sup_{0 < s \leq \varepsilon} \frac{\phi(s)}{s}$ . Note that we now employ the geodesic distance  $d_{\bar{\Omega}}$ , i.e. path can also lie on the boundary of  $\Omega$ .

**Assumptions on the Labeling Function** The assumption on the labeling for each sub-problem is weakened even further. Namely, we do not assume the labels for each discrete problem to be given by the continuum constraint  $g : \mathcal{O} \rightarrow \mathbb{R}$ . Instead we require

- (C1)  $\sup_{n \in \mathbb{N}} \text{Lip}_n(\mathbf{g}_n; \mathcal{O}_n) < \infty$  and  $\text{Lip}_\Omega(g; \mathcal{O}) < \infty$ ,
- (C2) there exists  $C > 0$  such that for all  $z \in \mathcal{O}$  it holds that  $|\mathbf{g}_n(\pi_{\mathcal{O}_n}(z)) - g(z)| \leq C\delta_n$ .

This assumption on the labeling function, is also connected to the no-noise and infinite data limit, see e.g. [Hof+20].

### 3.4.2. Convergence Results

In [LIP-I] the proof strategy relied on  $\Gamma$ -convergence. In [LIP-II] we use the comparison with cones characterization of AMLEs together with comparison principles for the infinity Laplace operator. An important strategy is to consider the homogenized operator

$$\Delta_\infty^\tau u(x) = \frac{1}{\tau^2} \left( \inf_{y \in B_\tau(x)} (u(y) - u(x)) + \sup_{y \in B_\tau(x)} (u(y) - u(x)) \right) \quad (3.26)$$

on the length scale  $\tau > 0$  which is typically larger than the graph scale  $\varepsilon$ . The results in the following are non-asymptotic, for which also our scaling assumption reads as follows

$$\begin{aligned} \varepsilon &\leq r_\Omega, \\ \frac{\varepsilon}{\tau} &< \frac{1}{2}, \\ \sigma_\phi(\varepsilon) + \frac{\delta_n}{\varepsilon} &\leq \frac{t_0}{4 + 2\sigma_\phi(\varepsilon)} \left(1 - 2\frac{\varepsilon}{\tau}\right). \end{aligned} \quad (3.27)$$

for some  $n \in \mathbb{N}$ . We detail the concrete appearance of the operator  $\Delta_\infty^\tau$  in the following paragraphs and first state the main general result. We employ the notation  $a \lesssim b$  which means that  $a \leq Cb$  for some constant  $C > 0$  depending only on  $g$  and  $\eta$ , but not on  $n, \varepsilon$  or  $\delta$ .

**Theorem 3.50 ([LIP-II, Th. 2.2]).** Let  $\tau > 0$ , let the kernel and data fulfill (K2) to (K4), (C1) and (C2) and let (3.25) and (3.27) hold. Let  $\mathbf{u}_n : \Omega_n \rightarrow \mathbb{R}$  solve Problem 3.36, and  $u : \bar{\Omega} \rightarrow \mathbb{R}$  solve Problem 3.21.

1. It holds

$$\sup_{\Omega_n} |u - \mathbf{u}_n| = \sup_{\bar{\Omega}} |u - u_n| \lesssim \tau + \sqrt[3]{\frac{\delta_n}{\varepsilon\tau} + \frac{\varepsilon}{\tau^2} + \frac{\phi(\varepsilon)}{\varepsilon\tau}}.$$

2. If  $\inf_{\bar{\Omega}} |\nabla u| > 0$ , then it even holds

$$\sup_{\Omega_n} |u - \mathbf{u}_n| = \sup_{\bar{\Omega}} |u - u_n| \lesssim \tau + \frac{\delta_n}{\varepsilon\tau} + \frac{\varepsilon}{\tau^2} + \frac{\phi(\varepsilon)}{\varepsilon\tau}.$$

Here,  $u_n : \bar{\Omega} \rightarrow \mathbb{R}$  denotes the piecewise constant extension of  $\mathbf{u}_n$ , as in (3.21).

We observe that in the case where  $|\nabla u|$  is bounded away from zero the rate improves, which is due to a perturbation we briefly mention in the following. The parameter  $\tau$  appears due to our proof strategy and does not have a concrete meaning for the final problem. An immediate consequence is that for an arbitrary  $\tau > 0$  we obtain

$$\sup_{\Omega_n} |u - \mathbf{u}_n| < \tau$$

by sending  $\delta_n, \varepsilon_n$  to zero as  $n \rightarrow \infty$  employing the weakest scaling assumption  $\delta_n \ll \varepsilon_n$ . Since  $\tau > 0$  was arbitrary this already yields convergence of discrete AMLEs under the weakest scaling assumption [LIP-II, Cor. 2.3]. In order to obtain the convergence rates involving only  $\delta$  and  $\varepsilon$  we can “optimize” over  $\tau$ . We first consider the small length scale regime, which in our case includes any graph scaling such that  $\varepsilon \lesssim \delta_n^{5/9}$ . In that case we can choose  $\tau = ((\delta_n + \phi(\varepsilon))/\varepsilon)^{1/4}$  to obtain

$$\sup_{\Omega_n} |u - \mathbf{u}_n| = \sup_{\bar{\Omega}} |u - u_n| \lesssim \left( \frac{\delta_n + \phi(\varepsilon)}{\varepsilon} \right)^{\frac{1}{4}}.$$

If  $|\nabla u|$  is bounded away from zero this can be improved to

$$\sup_{\Omega_n} |u - \mathbf{u}_n| = \sup_{\bar{\Omega}} |u - u_n| \lesssim \left( \frac{\delta_n + \phi(\varepsilon)}{\varepsilon} \right)^{\frac{1}{2}}.$$

where we consider the regime  $\varepsilon \lesssim \delta^{3/5}$ . In the convex case, i.e.  $\phi \equiv 0$  we therefore directly obtain the rates

$$\begin{aligned} \inf_{\bar{\Omega}} |\nabla u| &= 0 : \left( \frac{\delta_n}{\delta_n^{5/9}} \right)^{\frac{1}{4}} = \delta_n^{\frac{1}{9}}, \\ \inf_{\bar{\Omega}} |\nabla u| &> 0 : \left( \frac{\delta_n}{\delta_n^{3/5}} \right)^{\frac{1}{2}} = \delta_n^{\frac{1}{5}}, \end{aligned}$$

see [LIP-II, Cor. 2.4]. In fact the convexity assumption is too strict, instead we just have to require that  $\phi(\varepsilon) \lesssim \varepsilon^p$  for some power of  $p \in \mathbb{R}^+$  that does not deteriorate the rate. This is done in [LIP-II, Cor. 2.5, 2.6], where also a better rate is obtained by allowing for larger length scales. In the following we sketch the main ingredients of the proof. While in the  $\Gamma$ -convergence case we employ the hierarchy

“discrete  $\rightarrow$  non-local  $\rightarrow$  continuum”

our scheme now has the form

“discrete  $\rightarrow$  homogenized non-local  $\leftarrow$  continuum”.

We discuss the two directions below and then discuss how to “meet in the middle”.

**Non-Local to Continuum** In this paragraph we comment on how to relate continuum AMLEs to the homogenized operator in Eq. (3.26). Here, it is convenient to copy the following definitions from [LIP-II],

$$T^\tau u(x) := \sup_{\overline{B}_\Omega(x, \tau)} u, \quad T_\tau u(x) := \inf_{\overline{B}_\Omega(x, \tau)} u, \quad (3.28)$$

$$S_\tau^+ u := \frac{1}{\tau} (T^\tau u - u), \quad S_\tau^- u := \frac{1}{\tau} (u - T_\tau u), \quad (3.29)$$

so that we can write  $\Delta_\infty^\tau u = \frac{1}{\tau} (S_\tau^+ u - S_\tau^- u)$ . The astonishing property of this operator is the so-called “max-ball” lemma. If  $u \in USC(\overline{\Omega})$  fulfills comparison with distance functions from above, then we have that the function  $T^\tau u$  fulfills

$$-\Delta_\infty^\tau T^\tau u \leq 0, \quad \text{in } \Omega_\mathcal{O}^{2\tau}, \quad (3.30)$$

analogously if  $u \in LSC(\overline{\Omega})$  fulfills CDF form below we have

$$-\Delta_\infty^\tau T_\tau u \geq 0, \quad \text{in } \Omega_\mathcal{O}^{2\tau}. \quad (3.31)$$

Here, we denote by  $\Omega_\mathcal{O} = \overline{\Omega} \setminus \mathcal{O}$  and additional employ the inner parallel set,

$$\Omega_\mathcal{O}^\tau := \{x \in \overline{\Omega} : \text{dist}_\Omega(x, \mathcal{O}) > \tau\}. \quad (3.32)$$

Both, Eqs. (3.30) and (3.31) are shown in [LIP-II, Lem. 4.6], again borrowing concepts from [AS10]. In fact the proof is a simple consequence of the comparison with cones characterization. The strategy works as follows:

- Consider the ball  $B(x, 2\tau)$ .
- For any  $y \in B(x, 2\tau)$  we have that

$$\begin{aligned} u(y) &= u(x) + u(y) - u(x) = u(x) + \frac{u(y) - u(x)}{d_{\overline{\Omega}}(x, y)} d_{\overline{\Omega}}(x, y) \\ &\leq u(x) + \frac{T^{2\tau} u(x) - u(x)}{d_{\overline{\Omega}}(x, y)} d_{\overline{\Omega}}(x, y) \end{aligned}$$

- Consider the set  $V := B(x, 2\tau) \setminus \{x\}$  where on its boundary we can substitute the denominator by  $2\tau$  to obtain

$$u(y) \leq u(x) + \frac{T^{2\tau} u(x) - u(x)}{2\tau} d_{\overline{\Omega}}(x, y) \quad \text{for } y \in \partial V.$$

- Employing comparison with distance functions, this inequality holds on the whole domain  $V$ , therefor we can take the maximum over  $B(x, \tau)$  and see

$$T^\tau u(x) \leq \frac{T^{2\tau} u(x) - u(x)}{2\tau} \tau = \frac{1}{2} (T^{2\tau} u(x) + u(x))$$

- Rearranging the terms yields

$$\begin{aligned} 0 &\geq -\left(T^{2\tau}u + u - 2T^\tau u\right) = -(T^\tau T^\tau u + u - 2T^\tau u) \\ &\geq -(T^\tau T^\tau u + T_\tau T^\tau u - 2T^\tau u) = -\tau^2 \Delta_\infty^\tau T^\tau u \end{aligned}$$

It is interesting that the inequality is exact, i.e., we do not need to introduce an error term of order  $\tau$ . This observation allows us to pass from the continuum setting to the non-local one.

**Discrete to Non-Local** Compared to the previous paragraph passing from the discrete problem to the non-local one is more challenging. Analogously to the continuum case, where we considered the operators  $T^\tau, T_\tau$  we now define the functions

$$u_n^\tau(x) := \sup_{\overline{B}_\Omega(x, \tau) \cap \Omega_n} \mathbf{u}_n, \quad (u_n)_\tau(x) := \inf_{\overline{B}_\Omega(x, \tau) \cap \Omega_n} \mathbf{u}_n, \quad x \in \overline{\Omega}. \quad (3.33)$$

Concerning the connection between the discrete problem and the non-local operator, we state the following main result from [LIP-II].

**Theorem 3.51** ([LIP-II, Th. 5.13]). Assume that (K2) to (K4) and (3.25) and (3.27) hold. Let  $\mathbf{u}_n : \Omega_n \rightarrow \mathbb{R}$  solve the graph infinity Laplacian equation Problem 3.36. Then there exists a constant  $C > 0$  such that for all  $x_0 \in \Omega_\mathcal{O}^{2\tau+3\delta_n}$  it holds

$$-\Delta_\infty^\tau u_n^\tau(x_0) \leq \text{Lip}_n(\mathbf{g}_n) C \left( \frac{\delta_n}{\varepsilon\tau} + \frac{\varepsilon}{\tau^2} + \frac{\phi(\varepsilon)}{\varepsilon\tau} \right), \quad (3.34a)$$

$$-\Delta_\infty^\tau (u_n)_\tau(x_0) \geq -\text{Lip}_n(\mathbf{g}_n) C \left( \frac{\delta_n}{\varepsilon\tau} + \frac{\varepsilon}{\tau^2} + \frac{\phi(\varepsilon)}{\varepsilon\tau} \right). \quad (3.34b)$$

The above statement again only holds true on a inner parallel set, like [LIP-II, Lem. 4.6]. We later employ Lipschitz properties to obtain a result on the whole domain. Similarly, in the proof we first consider vertices  $\mathbf{x}_0 \in \Omega_\mathcal{O}^{2\tau+3\delta_n} \cap \Omega_n$  together with  $\mathbf{p}_n^\tau(\mathbf{x}_0), \mathbf{p}_n^{2\tau}$  such that

$$\begin{aligned} u_n^\tau(\mathbf{x}_0) &= \sup_{\overline{B}_\Omega(\mathbf{x}_0, \tau) \cap \Omega_n} \mathbf{u}_n = \mathbf{u}_n(\mathbf{p}_n^\tau(\mathbf{x}_0)), \\ u_n^{2\tau}(\mathbf{x}_0) &= \sup_{\overline{B}_\Omega(\mathbf{x}_0, 2\tau) \cap \Omega_n} \mathbf{u}_n = \mathbf{u}_n(\mathbf{p}_n^{2\tau}(\mathbf{x}_0)). \end{aligned}$$

Plugging in the definition of  $\Delta_\infty^\tau$  we then obtain the inequality

$$-\tau^2 \Delta_\infty^\tau u_n^\tau(\mathbf{x}_0) \leq 2\mathbf{u}_n(\mathbf{p}_n^\tau(\mathbf{x}_0)) - \mathbf{u}_n(\mathbf{p}_n^{2\tau}(\mathbf{x}_0)) - \mathbf{u}_n(\mathbf{x}_0).$$

Here, we would now like to employ comparison with graph distance functions for a similar strategy as in the continuum case. However, the choice of the set such that we have a desirable inequality on its boundary is not trivial. This is due the fact, that for

$u_n^\tau$  we take the supremum over a Euclidean ball, but we want to apply comparison with graph distance function. As in [LIP-I] we observe that our arguments need more care, whenever the Euclidean and the graph distance meet. We remark how the situation changes from the comparison with cones application in the continuum, where we assume that  $\phi \equiv 0$  for simplicity, the concrete details are given in [LIP-II, Sec. 5.2]:

- We know that  $u_n^{2\tau}$  maximizes on  $\overline{B}(2\tau, \mathbf{x}_0)$ .
- Since  $d_{w_n}(\mathbf{x}_0, \mathbf{w}) \leq |\mathbf{x}_0 - \mathbf{w}|$  we only need to shrink the graph ball by  $\varepsilon_n$  to ensure, that its closure lies in  $B(2\tau, \mathbf{x}_0)$ . This means we can consider

$$B = \{\mathbf{w} \in \Omega_n : d_{w_n}(\mathbf{x}_0, \mathbf{w}) \leq 2\tau - \varepsilon\}$$

and the we have that  $\overline{B} \subset B(2\tau, \mathbf{x}_0)$ .

- In the continuum case we knew that every point  $y \in \partial B(2\tau, \mathbf{x}_0)$  had exactly distance  $2\tau$  from  $x_0$ . This value is not as simple in the graph case, we need to consider the value

$$D = \inf_{\mathbf{y} \in \Omega_n \setminus \overline{B}(\mathbf{x}_0, 2\tau - \varepsilon)} d_{w_n}(\mathbf{x}_0, \mathbf{y}) \geq 2\tau - \varepsilon.$$

- Employing comparison with graph distance functions we then obtain

$$\mathbf{u}_n(\mathbf{w}) \leq \mathbf{u}_n(\mathbf{x}_0) + \frac{\mathbf{u}_n(\mathbf{p}_n^{2\tau}(\mathbf{x}_0)) - \mathbf{u}_n(\mathbf{x}_0)}{D} d_{w_n}(\mathbf{x}_0, \mathbf{w}) \quad (3.35)$$

for all vertices  $\mathbf{w} \in \Omega_n$  with  $d_{w_n}(\mathbf{x}_0, \mathbf{w}) \leq 2\tau - \varepsilon$ .

- We now choose  $\mathbf{w} = \mathbf{p}_n^\tau(\mathbf{x}_0)$ , however can not simply control the distance to  $\mathbf{x}_0$  by  $\tau$  instead we have the value

$$N = \sup_{\mathbf{y} \in \overline{B}(\mathbf{x}_0, \tau) \cap \Omega_n} d_{w_n}(\mathbf{x}_0, \mathbf{y})$$

and obtain

$$\begin{aligned} 2\mathbf{u}_n(\mathbf{p}_n^\tau(\mathbf{x}_0)) &\leq 2\mathbf{u}_n(\mathbf{x}_0) + (\mathbf{u}_n(\mathbf{p}_n^{2\tau}(\mathbf{x}_0)) - \mathbf{u}_n(\mathbf{x}_0)) \frac{2N}{D} \\ &= 2\mathbf{u}_n(\mathbf{x}_0) + (\mathbf{u}_n(\mathbf{p}_n^{2\tau}(\mathbf{x}_0)) - \mathbf{u}_n(\mathbf{x}_0)) \left(1 - 1 + \frac{2N}{D}\right) \\ &= \mathbf{u}_n(\mathbf{x}_0) + \mathbf{u}_n(\mathbf{p}_n^{2\tau}(\mathbf{x}_0)) + (\mathbf{u}_n(\mathbf{p}_n^{2\tau}(\mathbf{x}_0)) - \mathbf{u}_n(\mathbf{x}_0)) 2 \left(\frac{N}{D} - \frac{1}{2}\right). \end{aligned}$$

- Employing the Lipschitz estimates from [LIP-II, Sec. 5.1] we obtain

$$2\mathbf{u}_n(\mathbf{p}_n^\tau(\mathbf{x}_0)) \leq \mathbf{u}_n(\mathbf{x}_0) + \mathbf{u}_n(\mathbf{p}_n^{2\tau}(\mathbf{x}_0)) + C \operatorname{Lip}_n(\mathbf{g}_n) \left(\frac{N}{D} - \frac{1}{2}\right).$$

In our work we employ a concrete upper bound on the graph function [LIP-II, Lem. 5.5] which tells us that for any  $\mathbf{x}, \mathbf{y} \in \Omega_n$  we have

$$d_{w_n}(\mathbf{x}, \mathbf{y}) \leq \left(1 + \frac{4\delta_n}{t_0\varepsilon} + \frac{2\phi(\delta_n)}{t_0\varepsilon}\right) d_\Omega(\mathbf{x}, \mathbf{y}) + \tau_\eta \varepsilon$$

where

$$\tau_\eta := \sup_{0 < t \leq t_0} \left\{ \sigma_\eta \eta(t)^{-1} - t \right\}.$$

Plugging in this concrete error term, we then obtain the estimate

$$2\mathbf{u}_n(\mathbf{p}_n^\tau(\mathbf{x}_0)) \leq \mathbf{u}_n(\mathbf{x}_0) + \mathbf{u}_n(\mathbf{p}_n^{2\tau}(\mathbf{x}_0)) + \text{Lip}_n(\mathbf{g}_n)C \left( \frac{\delta_n \tau}{\varepsilon} + \varepsilon + \tau \frac{\phi(\varepsilon)}{\varepsilon} \right).$$

Rearranging this inequality then gives the first statement in [LIP-II, Th. 5.13]. The second one can be proven analogously. Quantitatively, we remark that the term

$$\frac{N}{D} - \frac{1}{2} = \frac{\sup_{\mathbf{y} \in \overline{B}(\mathbf{x}_0, \tau) \cap \Omega_n} d_{w_n}(\mathbf{x}_0, \mathbf{y})}{\inf_{\mathbf{y} \in \Omega_n \setminus \overline{B}(\mathbf{x}_0, 2\tau - \varepsilon)} d_{w_n}(\mathbf{x}_0, \mathbf{y})} - \frac{1}{2} \quad (3.36)$$

determines the leading term in the estimate and therefore the rate. This coined the phrase:

*“Rates for graph distance functions imply rates for graph AMLEs”.*

**Meeting in the Middle** The previous paragraphs allow us to relate the continuum and the discrete solution to the non-local operator. The question is now, how we can put these estimates together. Here, we employ the maximum principle for the operator  $\Delta_\infty^\tau$ , from [Sma10, Thm. 2.6.5]. Namely if for a constant  $C \geq 0$  two functions  $w, v : \Omega_\mathcal{O}^\tau \rightarrow \mathbb{R}$  satisfy

$$-\Delta_\infty^\tau w \leq C \leq -\Delta_\infty^\tau v \quad \text{in } V \subset \Omega_\mathcal{O}^\tau, \quad (3.37)$$

then it holds

$$\sup_{\Omega_\mathcal{O}^\tau} (w - v) = \sup_{\Omega_\mathcal{O}^\tau \setminus V} (w - v).$$

In our case the function  $w$  is chosen as  $u_n^\tau$ ,  $V = \Omega_\mathcal{O}^{2\tilde{\tau}_n}$  with  $\tilde{\tau}_n := \tau + \frac{3}{2}\delta_n$  and

$$\text{Lip}_n(\mathbf{g}_n)C \left| \frac{\delta_n}{\varepsilon\tau} + \frac{\varepsilon}{\tau^2} + \frac{\phi(\varepsilon)}{\varepsilon\tau} \right| =: C_{n,\tau},$$

in which case the left inequality in (3.37) holds true. We now would like to chose  $v = T_\tau u$  however, we now that  $-\Delta_\infty^\tau T_\tau u \geq 0$ . In order to obtain the desired inequality, we employ a perturbation trick. First we see that [Sma10, Lem. 2.6.4] tells us that for a bounded

function  $v : \Omega_{\mathcal{O}}^\tau \rightarrow \mathbb{R}$  with  $-\Delta_\infty^\tau v \geq 0$  on  $\Omega_{\mathcal{O}}^{2\tau}$ , there exists  $\delta_0 > 0$  such that for any  $0 \leq \delta \leq \delta_0$  the function  $w := v - \delta v^2$  satisfies

$$-\Delta_\infty^\tau w \geq -\Delta_\infty^\tau v + \delta(S_\tau^- v)^2 \quad \text{on } \Omega_{\mathcal{O}}^{2\tau}.$$

In the case that  $|\nabla u|$  is bounded away from 0 we can choose  $w := T_\tau u - \frac{C_{n,\tau}}{\alpha^2}(T_\tau u)^2$  with  $\alpha = \inf_{\bar{\Omega}} |\nabla u| > 0$ , which yields

$$-\Delta_\infty^\tau w \geq -\Delta_\infty^\tau T_\tau u + \frac{C_{n,\tau}}{\alpha^2}(S_\tau^- T_\tau u)^2 \quad \text{in } \Omega_{\mathcal{O}}^{2\tau}.$$

Since  $u$  is lower semi-continuous we can choose  $y_0$  with  $d_\Omega(x, y_0)$  such that

$$\begin{aligned} S_\tau^- T_\tau u(x) &= \frac{1}{\tau} (u(y_0) - T_\tau T_\tau u(x)) = \frac{1}{\tau} (u(y_0) - T_{2\tau} u(x)) \\ &\geq \frac{1}{\tau} \left( u(y_0) - (u(y_0) - \tau \inf_{\bar{\Omega}} |\nabla u|) \right) \\ &= \alpha, \end{aligned}$$

and therefore we have  $-\Delta_\infty^\tau w \geq C_{n,\tau}$ . Furthermore, we see that  $\|w - T_\tau u\|_\infty \leq \frac{c}{\alpha^2} C_{n,\tau}$  for some  $c > 0$ , which allows to make the following estimate

$$\begin{aligned} \sup_{\Omega_{\mathcal{O}}^{\tilde{\tau}_n}} (u_n^\tau - T_\tau u) &\leq \sup_{\Omega_{\mathcal{O}}^{\tilde{\tau}_n}} (u_n^\tau - w) + \frac{c}{\alpha^2} C_{n,\tau} \\ &\leq \sup_{\Omega_{\mathcal{O}}^{\tilde{\tau}_n} \setminus \Omega_{\mathcal{O}}^{2\tilde{\tau}_n}} (u_n^\tau - w) + \frac{c}{\alpha^2} C_{n,\tau} \\ &\leq \sup_{\Omega_{\mathcal{O}}^{\tilde{\tau}_n} \setminus \Omega_{\mathcal{O}}^{2\tilde{\tau}_n}} (u_n^\tau - T_\tau u) + 2 \frac{c}{\alpha^2} C_{n,\tau}. \end{aligned}$$

In the case that  $\alpha = 0$  we need to apply the additional perturbation trick from [Sma10, Lem. 2.6.3], which traditionally introduces a root of order 3. In total we then have the following result. Here we only include the results for the sub-solutions, the respective other inequality holds analogously and can be found in [LIP-II, Prop. 5.16].

**Proposition 3.2 ([LIP-II, Prop. 5.16]).** Under the previous assumptions there exists a constant  $C > 0$  such that

$$\sup_{\Omega_{\mathcal{O}}^{\tilde{\tau}_n}} (u_n^\tau - T_\tau u) \leq \sup_{\Omega_{\mathcal{O}}^{\tilde{\tau}_n} \setminus \Omega_{\mathcal{O}}^{2\tilde{\tau}_n}} (u_n^\tau - T_\tau u) + \sqrt[3]{\text{Lip}_n(\mathbf{g}_n) C \left( \frac{\delta_n}{\varepsilon \tau} + \frac{\varepsilon}{\tau^2} + \frac{\phi(\varepsilon)}{\varepsilon \tau} \right)}.$$

If additionally it holds that  $\alpha := \inf_{\bar{\Omega}} |\nabla u| > 0$ , then we even have

$$\sup_{\Omega_{\mathcal{O}}^{\tilde{\tau}_n}} (u_n^\tau - T_\tau u) \leq \sup_{\Omega_{\mathcal{O}}^{\tilde{\tau}_n} \setminus \Omega_{\mathcal{O}}^{2\tilde{\tau}_n}} (u_n^\tau - T_\tau u) + \text{Lip}_n(\mathbf{g}_n) C \left( \frac{\delta_n}{\varepsilon \tau} + \frac{\varepsilon}{\tau^2} + \frac{\phi(\varepsilon)}{\varepsilon \tau} \right).$$

Now we can employ the Lipschitz estimates from [LIP-II, Sec. 5.1] to extend the above result to the whole domain  $\Omega_{\mathcal{O}}$ , which then finally yields [LIP-II, Thm. 2.2].

### 3.4.3. Numerical Examples and Extensions

Here, we briefly comment on the numerical examples conducted in [LIP-II]. In order to evaluate the convergence rates numerically, we employ a known infinity harmonic function on  $\mathbb{R}^2$ , namely the Aronsson function

$$u(x_1, x_2) = |x_1|^{4/3} - |x_2|^{4/3}.$$

We choose the constraint set  $\mathcal{O} = \{(\pm 1, 0), (0, \pm 1)\}$  with the boundary values given by  $u$ . In order to ensure that  $u$  is the AMLE of its boundary values we choose the domain  $\Omega$  such that the Neumann boundary condition on  $\partial\Omega$  is fulfilled. In [LIP-II] we therefore introduce the Aronsson star

$$\Omega := \left\{ x \in [0, 1]^2 : |x_1|^{2/3} + |x_2|^{2/3} \leq 1 \right\},$$

which is non-convex, non-smooth but fulfills Eq. (3.20). In the discrete problems we consider all combinations of the following configurations,

- exact boundary vertices  $\mathcal{O}_n = \mathcal{O}$  and dilated boundary conditions with  $\mathcal{O}_n^{\text{dil}} = \{\mathbf{x} \in \Omega_n : \text{dist}(\mathbf{x}, \mathcal{O}) \leq \varepsilon_n\}$ ,
- unit weights and singular weights,
- scalings  $\varepsilon_n \sim \{\delta_n, \delta_n^{2/3}, \delta_n^{1/2}\}$ .

The results are displayed in [LIP-II, Fig. 1–Fig. 4], where we want to highlight the following observations:

- (i) The rates are generally better than the ones proven in [LIP-II, Thm. 2.2].
- (ii) Singular weights give better rates than unit weights.
- (iii) The best rate is achieved for the scaling  $\varepsilon_n \sim \delta_n$ , for which our results do not even proof convergence.

Concerning (i) we note that this does not directly imply that our analysis is not sharp. Typically, rates for the Aronsson function are better than the ones one can proof in general, see e.g. [Sma10]. Regarding point (iii) we want to mention the follow-up work [LIP-III], which deals with the case  $\varepsilon_n \sim \delta_n$ . Working in the setting of first passage percolation we are able to show a rate for the quantity

$$\left| \frac{\mathbb{E}[d_{\varepsilon_s, \mathcal{X}_s}(0, se_1)]}{\mathbb{E}[d_{\varepsilon_s, \mathcal{X}_s}(0, 2se_1)]} - \frac{1}{2} \right|$$



The code for all the experiments is available at [github.com/jwcalder/LipschitzLearningRates](https://github.com/jwcalder/LipschitzLearningRates).

is shown, [LIP-III, Thm. 2.1]. Here,  $\mathcal{X}_s$  denote modified Poisson point processes, over which the expectation is taken. The distance  $d_{\varepsilon_s, \mathcal{X}_s}(0, se_1)$  is basically a graph distance with unit weights, connecting 0 and the point  $se_1$ . We observe that this term is reminiscent of the dominating term Eq. (3.36) we have in [LIP-III]. In fact, employing the same strategy as displayed in the previous paragraphs, in [LIP-III] we are able to show a rate in the case  $\varepsilon_n \sim \delta_n$ .

# Chapter 4

## Robust and Sparse Supervised Learning

In this chapter we present the topics discussed in the works [CLIP; FNO; BREG-I] which are reprinted in ???. Compared to the previous chapter we are now in the setting of supervised learning as described in [Section 2.2](#). As already mentioned before, we focus on neural networks  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  parameterized by  $\theta \in \Theta$ , where  $\Theta$  is some parameter space. Here, our two main question arise, namely:

- **Input robustness:** how robust is  $x \mapsto f_\theta(x)$  w.r.t. input perturbations?
- **Parameter sparsity:** how can we obtain sparse parameters  $\theta \in \Theta$ ?

Section 4.1: Setting

Section 4.2: [CLIP]

Section 4.3: [FNO]

Section 4.4: [BREG-I]

Therefore, conceptually we again highlight the keyword **sparsity** and additionally **robustness**. In [CLIP] we consider robustness under adversarial perturbations. Here, the input is *attacked* on purpose to confuse a neural network classifier and therefore worsen its performance. In order to obtain a classifier that is less vulnerable against such attacks we propose an optimization strategy that selects parameters yielding a more robust network. In [Section 4.2](#) we comment on the topic and the contribution in more detail.

A different kind of input robustness is considered in [FNO]. In the setting of image classification, images are modeled as functions on a continuum domain that need to be discretized in order to represent them on a machine. However, this discretization is usually arbitrary and not inherent to the object of interest. Therefore, it is natural to assume that the output of the network should be independent of this discretization, also referred to as resolution, hence we consider robustness w.r.t. resolution changes. We remark on this and the publication in [Section 4.3](#).

Concerning computational performance and memory storage of the neural network we focus on sparsity of the parameters  $\theta \in \Theta$ . In [BREG-I] we propose a sparse optimization strategy based on Bregman iterations, which employs  $L^1$  type penalties to promote

sparsity. The conceptual difference between the existing algorithm and our proposal is the stochastic component in the gradient. Our theoretical results prove decay in loss and convergence of the iterates. Finally, we also conduct numerical experiments that demonstrate the efficiency of the method. We refer to [Section 4.4](#) for a detailed explanation.

Before, we start with the exposition on the mentioned works, we briefly review the common supervised learning framework. Building upon the basic observations in [Chapter 4](#) we give a slightly more detailed introduction in [Section 4.1](#).

## 4.1. Setting

We are given a finite training set  $\mathcal{T} \subset \mathcal{X} \times \mathcal{Y}$ . For a family of functions  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  parameterized by  $\theta \in \Theta$ , we consider the empirical minimization

$$\min_{\theta \in \Theta} \mathcal{L}(\theta),$$

where for a function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  we define

$$\mathcal{L}(\theta) := \frac{1}{|\mathcal{T}|} \sum_{(x,y) \in \mathcal{T}} \ell(f_\theta(x), y). \quad (4.1)$$

**Remark 4.1.** Assuming that  $\mathcal{T}$  is sampled from a joint distribution  $P_{\mathcal{X}, \mathcal{Y}}$  on  $\mathcal{X} \times \mathcal{Y}$  this approximates the computational infeasible population risk minimization problem

$$\inf_{\theta \in \Theta} \int_{\mathcal{X} \times \mathcal{Y}} \ell(f_\theta(x), y) d\pi(x, y).$$

△

In the following we provide two important choices for the function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ .

**Example 4.2 (MSE).** For image denoising problems, we often choose  $\mathcal{X} = \mathcal{Y} = [0, 1]^{N \times M}$ , assuming only one color channel for simplicity. In this context the mean squared  $L^2$  error (MSE), defined as

$$\ell(\bar{y}, y) := \frac{1}{N \cdot M} \|\bar{y} - y\|^2$$

is commonly employed. This loss function could however also be employed for classification problems.

**Example 4.3 (Cross-Entropy).** For classification problems one often chooses the *cross-entropy* or *negative log-likelihood* loss, [\[Goo52\]](#). For two discrete probability distributions,  $p, q : \{1, \dots, C\} \rightarrow \mathbb{R}$  one defines

$$H(p, q) := - \sum_{c=1}^C p_c \cdot \log(q_c)$$

see, e.g., [COR98]. Assuming that  $\mathcal{Y} = \Delta^C$  this allows to choose  $\ell(\bar{y}, y) := H(y, \bar{y})$ . If the network only maps to  $\mathbb{R}^C$  one often additionally inserts a soft-max function  $Q(y)_c := \exp(y_c) / \sum_c \exp(y_c)$  (see [Bol68]) and then sets

$$\ell(\bar{y}, y) := H(y, Q(\bar{y})).$$

In the case, where the output is given as labels  $y \in \{1, \dots, C\}$  one defines

$$\ell(\bar{y}, y) := H(y_{\text{oh}}, Q(\bar{y})) = -\log(Q(\bar{y}))$$

employing the one-hot notation from [Chapter 2](#).

#### 4.1.1. Network Architectures

In this thesis we focus on feed-forward neural networks, i.e., we consider layers of the form

$$\Phi(w, W, b)(z) := wz + \sigma(Wz + b) \quad (4.2)$$

where  $w \in \mathbb{R}$  models a residual connection,  $W \in \mathbb{R}^{n \times n}$  is a weight matrix,  $b \in \mathbb{R}^n$  a bias vector and  $z \in \mathbb{R}^n$  is the input. We consider a concatenation of  $L \in \mathbb{N}$  of such layers, which then forms a neural network

$$f_\theta = \Phi_L \circ \dots \circ \Phi_1$$

with parameters  $\theta = ((W_1, b_1, w_1), \dots, (W_L, b_L, w_L)) \in \Theta$  and layers  $\Phi_i := \Phi(w_i, W_i, b_i)$ . If there is no residual part, i.e.,  $w = 0$ , then we also allow dimensional changes in each layer, i.e.,  $z \in \mathbb{R}^n, W \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$ .

**MLP** In the easiest case we consider a perceptron [Ros58], which models a fully connected layer, i.e. every entry  $W_{ij}$  of the weight matrix is a parameter that is optimized in the training process.

**Convolutions** Especially important for visual tasks are convolutional layers. Here, we take a kernel  $k \in \mathbb{R}^{M \times M}$  and define the application of  $W = W(k)$  as

$$Wz = k * z$$

with  $*$  denoting the convolution. We refer to [Section 4.3](#) for the concrete definition of this operation. Typically, the input is of the form  $z \in \mathbb{R}^{K_{\text{in}} \times N \times M}$ , where  $K_{\text{in}}$  denotes the number of input channels. The layer is then a mapping  $\Phi : \mathbb{R}^{K_{\text{in}} \times N \times M} \rightarrow \mathbb{R}^{K_{\text{out}} \times N \times M}$ , where  $K_{\text{out}}$  denotes the number of output channels, which can be realized by different kernels  $k_{ij}$ . The application of  $W$  to an input  $z$  is defined as

$$(Wz)_{j,:,:} := \sum_{i=1}^{K_{\text{in}}} k_{ij} * z.$$

**ResNets** Often we also consider a residual component as displayed in Eq. (4.2) with the term  $w z$ . The idea of adding this residual component was first introduced in [SGS15] with a learnable parameter  $w \in \mathbb{R}$  and later popularized in [He+16a] by fixing  $w = 1$ , see also [He+16b], which then yields the celebrated ResNet architecture. In the following applications we both consider the case where  $w = 1$  is fixed, but also the possibility of learning the parameter  $w \in \mathbb{R}$  in [BREG-II].

#### 4.1.2. Gradient Computation and Stochastic Gradient Descent

Training a neural network requires solving an optimization problem w.r.t. to the parameters  $\theta \in \Theta$ . In this work we only focus on first order methods, however both zero [Rie22; Pin+17; Car+21] and second order methods [Mar10] have been successfully applied in this context. Employing first order methods, requires evaluating the gradient  $\nabla_{\theta}\mathcal{L}$ , however in this scenario it is not common to compute the full gradient but rather to have a gradient estimator. This estimator is usually obtained by randomly dividing the train set  $\mathcal{T}$  into disjoint mini-batches  $B_1 \cup \dots \cup B_b = \mathcal{T}$  and then successively computing the gradient of the mini-batch loss

$$\mathcal{L}(\theta; B) := \frac{1}{|B_i|} \sum_{(x,y) \in B_i} \ell(f_{\theta}(x), y).$$

Iterating over all batches  $i = 1, \dots, b$  is referred to as one epoch. From a mathematical point of view, this yields stochastic optimization methods, since in each step the true gradient is replaced by an estimator. In the abstract setting we let  $(\Omega, F, \mathbb{P})$  be a probability space and consider a function  $g : \Theta \times \Omega \rightarrow \Theta$  as an unbiased estimator of  $\nabla\mathcal{L}$ , i.e.

$$\mathbb{E}[g(\theta; \omega)] = \nabla\mathcal{L}(\theta) \text{ for all } \theta \in \Theta.$$

Most notably this method transforms the standard gradient descent update [Cau+47]

$$\theta^{(k+1)} = \theta^{(k)} - \tau^{(k)} \nabla\mathcal{L}(\theta^{(k)}),$$

to *stochastic* gradient descent [RM51]

draw  $\omega^{(k)}$  from  $\Omega$  using the law of  $\mathbb{P}$ ,

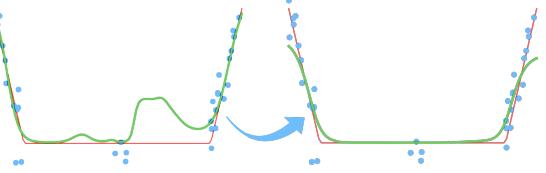
$$g^{(k)} := g(\theta^{(k)}; \omega^{(k)}),$$

$$\theta^{(k+1)} := \theta^{(k)} - \tau^{(k)} g^{(k)}.$$

## 4.2. Adversarial Stability via Lipschitz Training: [CLIP]

In the following we consider classification problems with  $C \in \mathbb{N}$  different classes and functions  $f : \mathcal{X} \rightarrow \Delta^C$  for which we denote by

$$f^{\text{MAP}}(x) := \operatorname{argmax}_c f(x)_c$$



the maximum a posteriori estimation for an input  $x \in \mathcal{X}$ . The capabilities of neural networks to perform classification tasks on unseen data—i.e. inputs  $x \in \mathcal{X}$  that are not part of the training data—are impressive and the reason for their popularity. However, it has been noticed in [GSS14] that it is relatively easy to “fool” classification networks in the following sense:

*Given a classification network and a human classifier  $f_\theta, h : \mathcal{X} \rightarrow \Delta^C$ , and an input  $x \in \mathcal{X}$  such that  $f_\theta^{\text{MAP}}(x) = h^{\text{MAP}}(x)$ . An adversarial example is an input  $\bar{x} \in \mathcal{X}$  that is close to  $x$ , but we have that*

$$f_\theta^{\text{MAP}}(\bar{x}) \neq h^{\text{MAP}}(\bar{x}) = h^{\text{MAP}}(x).$$

The vague concept of a *human classifier*, incorporates the intuition of the classification problem as a human would solve it. Assuming that we are given data  $\mathcal{T} \subset \mathcal{X} \times \mathcal{Y}$  that is sampled i.i.d. from a joint distribution  $P_{\mathcal{X}, \mathcal{Y}}$  this function could also be chosen as  $h(x)_c := P_{\mathcal{X}, \mathcal{Y}}(c|x)$ . The term “close” describes that the difference between the two images  $x, \bar{x}$  should be visually imperceptible.

**Remark 4.4.** Some authors argue that such instabilities are inherent to classification problems solved via a data-driven approach [Sha+18; FFF18]. From this point of view, trying to defend against these instabilities is not necessarily desirable. However, we do not consider this interpretation in the following.  $\triangle$

**What are Adversarial Examples Formally?** In order to formalize this idea, we omit any influence of the typically unavailable function  $h$  and instead consider  $\bar{x}$  adversarial if  $f_\theta^{\text{MAP}}(x) \neq f_\theta^{\text{MAP}}(\bar{x})$ . This interpretation only makes sense, if we assume that the output of  $f_\theta^{\text{MAP}}(x)$  is *correct*, which can only be easily checked on labeled data [Bun+23]. Furthermore, in this work we assume that  $\mathcal{X}$  models an image space  $\mathcal{X} = \mathbb{R}^{K \times N \times M}$ , with  $K, N, M \in \mathbb{N}$  and choose a distance  $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_0^+$  to measure the distance between points in  $\mathcal{X}$ .

**Remark 4.5.** In most of our examples, we choose  $d(\cdot, \cdot) = \|\cdot - \cdot\|_p$ . Employing a  $L^p$  norm in an image context—via pixel-wise comparisons—can be an unfavorable choice. Images that appear very different visually, can have a smaller  $L^p$  distance, than images that visually appear similar, see e.g. [Sta+21, Fig. 16]. However, it is easy to evaluate, which is crucial for most applications. In the absence of a better criterion, being close

to the original input in the  $L^p$  distance, is a practical way to decide if  $\bar{x}$  is an admissible adversarial example.  $\triangle$

Employing these simplifications we then say that  $\bar{x} \in \mathcal{X}$  is adversarial, if

$$d(x, \bar{x}) \leq \varepsilon \quad \text{and} \quad f_\theta^{\text{MAP}}(x) \neq f_\theta^{\text{MAP}}(\bar{x}),$$

where  $\varepsilon$  is called the *adversarial budget*. This parameter controls how far  $\bar{x}$  can be from the original point  $x$  to be still considered adversarial. Here, we also have to decide how much we trust the metric  $d(\cdot, \cdot)$  as a measure of closeness.

**Types of Adversarial Examples** Typically, adversarial examples are created from the clean image  $x$  via some distortion  $\delta \in \mathbb{R}^s$ . Together with an application map  $T : \mathcal{X} \times \mathbb{R}^s \rightarrow \mathcal{X}$  one then obtains  $\bar{x} = T(x, \delta)$  as the adversarial example. In this formulation one can alternatively employ the criterion  $\|\delta\| \leq \varepsilon$  to decide, whether  $T(x, \delta)$  is an adversarial example. We only list some of the approaches below:

- **Addition:** The most well-known examples are created by  $T(x, \delta) := x + \delta$ , i.e.  $\bar{x} = x + \delta$ . Here, it is important to note that typically images are assumed to have values between 0 and 1, i.e.  $\mathcal{X} = [0, 1]^{K \times N \times M}$ , therefore it is important to also ensure that  $\bar{x} \in \mathcal{X}$ .
- **Translation and Rotation:** Simple geometric transformations—that are unnoticeable for a human classifier—are quite effective to “fool” neural networks. Employing translations  $T_t : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{X}$  one has to choose the behavior on the boundary such that one obtains a valid image. The same holds true for rotations  $T_r : \mathcal{X} \times [-\pi, \pi] \rightarrow \mathcal{X}$ . In this case  $\delta \in [-\pi, \pi]$  models the angle of the rotation and yields an admissible adversarial example, if  $|\delta| \leq \varepsilon$ . Here, we see that this formulation loses some expressivity. Consider the MNIST dataset [LC10], where the task is to classify handwritten digits from 0 to 9:
  - We see that only  $\varepsilon < \pi/2$  makes sense, otherwise the number “6” could be always transformed to the number “9” and vice versa.
  - However, considering the number “0” rotations above the angle  $\pi/2$  can definitely yield proper adversarial examples  $\bar{x}$ .

We refer to [Eng+18] for a study on these types of adversarial examples.

- **Change of Basis:** As explored in [Guo+17] one can consider a different orthonormal basis of the image space  $\mathbb{R}^{K \times N \times M}$  and then perform the attack w.r.t. this basis. The map  $T : \mathcal{X} \times \mathbb{R}^{K \times N \times M} \rightarrow \mathcal{X}$  first obtains the different representation of  $x$ , then adds the coefficients of  $\delta$  and then maps back to the original basis. This is only meaningful, if one restricts certain coefficients of  $\delta$  to be zero in the alternative basis. For example, the discrete cosine transform ([ANR74]) has been applied successfully in this context [Guo+17].

In the following we restrict ourselves to the additive case and only consider examples  $\bar{x} = x + \delta$ .

**Finding Adversarial Examples** Employing a loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  the task of finding an adversarial examples  $\bar{x} = x + \delta \in \mathcal{X}$  can be relaxed via solving the problem

$$\max_{\bar{x} \in \mathcal{X}: d(x, \bar{x}) \leq \varepsilon} \ell(y, f_\theta(\bar{x})) = \max_{\delta: \|\delta\|_d \leq \varepsilon, x + \delta \in \mathcal{X}} \ell(y, f_\theta(x + \delta)), \quad (4.3)$$

where  $y = f_\theta(x)$ . Solving this problem is referred to as *attacking* the network  $f$ . More precisely, this is a so-called *untargeted* attack, since we do not prescribe  $\bar{x}$  to realize any special output as long as it confuses the network. Opposed to this, there are so called *targeted* attacks. Here, we pick  $c^{\text{adv}} \neq f_\theta^{\text{MAP}}(x)$  and then want to find  $\bar{x}$  such that  $f_\theta^{\text{MAP}}(\bar{x}) = c^{\text{adv}}$ , which yields the problem

$$\min_{\bar{x} \in \mathcal{X}: d(x, \bar{x}) \leq \varepsilon} \ell(c^{\text{adv}}, f_\theta(\bar{x})).$$

Conceptually, this problem is similar to Eq. (4.3) and differs only by a change of sign and a different reference label. A popular method to solve Eq. (4.3) is the projected sign gradient ascent iteration [KGB+16],

$$\bar{x}^{(k+1)} = \text{Proj}_d(\bar{x}^{(k)} + \tau \text{sign}(\nabla_{\bar{x}} \ell(f_\theta(x), f_\theta(\bar{x}^{(k)}))).$$

Here,  $\text{Proj}_d$  denotes the projection onto the set  $\mathcal{X} \cap B_{d, \varepsilon}(x)$ , where  $B_{d, \varepsilon}(x)$  is the ball with radius  $\varepsilon$  around  $x$ , w.r.t. to the distance  $d$ . Performing only one step of this iteration yields the fast gradient sign method [GSS14]. There is a wide variety of these so-called gradient-based open-box attacks [Yua+19], i.e., methods that assume that the gradient of the model is available. In more realistic scenarios, this might not be the case. Attacks that do not employ the actual gradient of the model to attack are called *closed box* attacks [Ily+18], which are not part of this thesis.

**Defending against Adversarial Attacks** We consider the question of finding parameters  $\theta \in \Theta$  such that corresponding model  $f_\theta$  is *adversarially robust*, i.e., is less vulnerable to attacks. Therefore, we want the attack problem in Eq. (4.3) to be hard to solve. This intuition leads to the optimization problem

$$\min_{\theta \in \Theta} \sum_{(x, y) \in \mathcal{T}} \max_{\bar{x}: d(x, \bar{x}) \leq \varepsilon} \ell(f_\theta(\bar{x}), y)$$

which is known as the *adversarial training* formulation [KGB16; Mad+17]. From a distributionally robust optimization point of view, this problem relates to

$$\min_{\theta \in \Theta} \max_{\tilde{P}: D(P_{\mathcal{X}, \mathcal{Y}}, \tilde{P}) \leq \varepsilon} \int_{\mathcal{X} \times \mathcal{Y}} \ell(f_\theta(x), y) d\tilde{P}(x, y)$$

where  $D$  denotes some distance on the space of probability distributions, see [BGM23]. Employing a batched gradient descent type iteration for the variable  $\theta$ , one then has to solve a problem of the form Equation (4.3) in every step, which can again be approximated via gradient ascent on the input variable  $x$ .

**Lipschitz Training of Neural Networks** Adversarial robustness is closely related to the Lipschitzness of a network  $f_\theta$ . Namely, for inputs  $x, \bar{x} \in \mathcal{X}$  that are close, we also want the outputs of  $f_\theta$  to be close. In other words we want to find a small constant  $L > 0$  such that

$$\|f_\theta(x) - f_\theta(\bar{x})\|_{\mathcal{Y}} \leq L \|x - \bar{x}\|_{\mathcal{X}},$$

where we typically employ an  $L^p$  norm for both  $\|\cdot\|_{\mathcal{X}}$  and  $\|\cdot\|_{\mathcal{Y}}$ . The smallest constant fulfilling this inequality is the Lipschitz constant  $\text{Lip}(f_\theta)$ . If this constant is small, one can expect that small input perturbations do not affect the classification output significantly.

**Remark 4.6.** Apart from adversarial robustness, having an upper bound on the Lipschitz constant of a neural network is also important for other applications, see, e.g., [Arn+23; Has+20; ACB17].  $\triangle$

In our work, we employ the Lipschitz constant as a regularizer and consider the problem

$$\min_{\theta} \mathcal{L}(\theta) + \lambda \text{Lip}(f_\theta) \quad (4.4)$$

for a parameter  $\lambda > 0$ . Computing the Lipschitz constant of neural networks is an NP-hard problem [SV18], therefore many works employ the estimate

$$\text{Lip}(f_\theta) \leq \prod_{l=1}^L \text{Lip}(\Phi_l) \leq C_\sigma \prod_{l=1}^L \|W_l\|, \quad (4.5)$$

where here  $\|\cdot\|$  denotes some matrix norm and  $C_\sigma > 0$  depends on the Lipschitz constants of the activation functions, see [ALG19; Gou+20; KMP20; RKH19]. This inequality is not sharp and usually overestimates the Lipschitz constant, as we see in the following example taken from [CLIP].

**Example 4.7.** We consider a feed-forward neural network  $f_\theta : \mathbb{R} \rightarrow \mathbb{R}$  with one hidden layer,

$$\begin{aligned} \Phi_1(z) &:= \text{ReLU}(W_1 z) := (\max\{z, 0\}, \max\{-z, 0\})^T, & W_1 &:= (1, -1), \\ \Phi_2(z) &:= W_2 z := z_1 + z_2, & W_2 &:= (1, 1)^T, \\ f_\theta &:= \Phi_2 \circ \Phi_1. \end{aligned}$$

For  $x \in \mathbb{R}$  we have that

$$\begin{aligned} x \geq 0 &\Rightarrow \Phi_1(x) = (x, 0)^T &\Rightarrow f_\theta(x) = x, \\ x \leq 0 &\Rightarrow \Phi_1(x) = (0, -x)^T &\Rightarrow f_\theta(x) = -x, \end{aligned}$$

and therefore  $f_\theta = |\cdot|$ , which yields that  $\text{Lip}(f_\theta) = 1$ . However employing the spectral norm for the weight matrices, we see that

$$\|W_1\| \cdot \|W_2\| = \sqrt{2} \sqrt{2} = 2.$$

Plugging the estimate of Eq. (4.5) into Eq. (4.4) therefore potentially over regularizes the problem. While this increases the stability of the network, it can worsen its classification performance.

**Contribution in [CLIP]** In [CLIP] we propose a strategy to solve Eq. (4.4) approximately, without the estimate in Eq. (4.5). The basic idea consists of approximating the Lipschitz constant on a finite set and using this approximation as a regularizer. Furthermore, we analyse the original model in Eq. (4.4) where we study existence of solutions, the influence of the parameter  $\lambda$  and the limits  $\lambda \rightarrow 0, \lambda \rightarrow \infty$ , see Section 4.2.2. Finally, we demonstrate the efficiency of the method by applying it to some simple toy problems, see Section 4.2.3.

#### 4.2.1. Cheap Lipschitz Training

We approximate the Lipschitz constant on a finite set  $\mathcal{X}_{\text{Lip}} \subset \mathcal{X} \times \mathcal{X}$  via

$$\text{Lip}(f_\theta; \mathcal{X}_{\text{Lip}}) := \max_{(x, \bar{x}) \in \mathcal{X}_{\text{Lip}}} \frac{\|f_\theta(x) - f_\theta(\bar{x})\|_{\mathcal{Y}}}{\|x - \bar{x}\|_{\mathcal{X}}} \approx \text{Lip}(f_\theta).$$

Disregarding the non-differentiable points of  $\theta \mapsto \text{Lip}(f_\theta; \mathcal{X}_{\text{Lip}})$  and employing the above approximation, allows us to solve the problem in Eq. (4.4) via stochastic gradient descent on the variable  $\theta$ .

**Remark 4.8.** Let  $f, g : \mathcal{X} \rightarrow \mathcal{Y}$  be two differentiable functions. If  $f(\bar{x}) > g(\bar{x})$  for some  $\bar{x} \in \mathcal{X}$  then we know that there exists some  $\epsilon > 0$  such that  $f(x) > g(x)$  for all  $x \in B_\epsilon(\bar{x})$ . We have that  $f \vee g = \max\{f, g\} = f$  in  $B_\epsilon(\bar{x})$  and therefore  $x \mapsto f(x) \vee g(x)$  is differentiable in  $\bar{x}$ . If  $f(\bar{x}) = g(\bar{x})$  then  $f \vee g$  could be non-differentiable at  $\bar{x}$ . However, in practice we employ automatic differentiation, where in this case one of the functions is chosen, say  $f$ , and the gradient at  $\bar{x}$  is computed as  $\nabla f$ .  $\triangle$

The strength of this approach is dependent on the quality of the set  $\mathcal{X}_{\text{Lip}}$ . In [CLIP] we propose to iteratively update the set via a gradient ascent type scheme. Namely, we initialize  $\mathcal{X}_{\text{Lip}}$  as a random perturbation of a subset of the given data. In each step we then update the points as follows:

- Consider  $L(x, \bar{x}) := \|f_\theta(x) - f_\theta(\bar{x})\| / \|x - \bar{x}\|$ ,  $x, \bar{x} \in \mathcal{X}$ .
- For each  $(x, \bar{x}) \in \mathcal{X}_{\text{Lip}}$  perform the update

$$x \leftarrow x + \tau L(x, \bar{x}) \nabla_x L(x, \bar{x}), \\ \bar{x} \leftarrow \bar{x} + \tau L(x, \bar{x}) \nabla_{\bar{x}} L(x, \bar{x}),$$

for a parameter  $\tau > 0$ .

This scheme performs gradient ascent on the Lipschitz constant, see [CLIP, Alg. 1]. For each mini-batch  $B \subset \mathcal{T}$  we first update the set  $\mathcal{X}_{\text{Lip}}$  and then update the parameters via

$$\theta \leftarrow \theta - \eta \nabla_\theta (\mathcal{L}(\theta; B) + \lambda \text{Lip}(f_\theta, \mathcal{X}_{\text{Lip}}))$$

which yields the algorithm presented in [CLIP, Alg. 2]. Similar to [Sha+19] one can reuse the gradients computed for  $\nabla_x$  for the computation of  $\nabla_\theta$ . This fact yields the attribute *cheap* in the Lipschitz training algorithm.

#### 4.2.2. Analysis of Lipschitz Regularization

In [CLIP] we analyse some basic properties of the original regularization problem in Eq. (4.4). We repeat the assumptions made therein.

**Assumption 1.** We assume that the loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R} \cup \{\infty\}$  satisfies:

- a)  $\ell(y, \bar{y}) \geq 0$  for all  $y, \bar{y} \in \mathcal{Y}$ ,
- b)  $y \mapsto \ell(y, \bar{y})$  is lower semi-continuous for all  $\bar{y} \in \mathcal{Y}$ .

**Assumption 2.** We assume that the map  $\theta \mapsto f_\theta(x)$  is continuous for all  $x \in \mathcal{X}$ .

**Assumption 3.** We assume that there exists  $\theta \in \Theta$  such that

$$\frac{1}{|\mathcal{T}|} \sum_{(x,y) \in \mathcal{T}} \ell(f_\theta(x), y) + \lambda \text{Lip}(f_\theta) < \infty.$$

Employing only Assumption 2 one can show that the map  $\theta \mapsto \text{Lip}(f_\theta)$  is lower semi-continuous, [CLIP, Lem. 1].

**Existence** If  $\Theta$  is compact or finite, one can show that there exist solutions of the problem in Eq. (4.4). In the general case, one needs to add a norm term, that ensures boundedness of a minimizing sequence. The following is the main existence result, which can be proven by the direct method in the calculus of variations [Dac07].

**Proposition 4.1** ([CLIP, Prop. 1]). Under Assumptions 1 to 3 the problem

$$\min_{\theta \in \Theta} \frac{1}{|\mathcal{T}|} \sum_{(x,y) \in \mathcal{T}} \ell(f_\theta(x), y) + \lambda \text{Lip}(f_\theta) + \mu \|\theta\|_\Theta$$

has a solution for all values  $\lambda, \mu > 0$ . Here,  $\|\cdot\|_\Theta$  denotes a norm on  $\Theta$ .

**Dependency on the Regularization Parameter** Assuming that for every  $\lambda > 0$  we have a solution  $\theta_\lambda \in \Theta$  of Eq. (4.4) one can show that

$$\begin{aligned} \lambda &\longmapsto \frac{1}{|\mathcal{T}|} \sum_{(x,y) \in \mathcal{T}} \ell(f_{\theta_\lambda}(x), y) \quad \text{is non-decreasing,} \\ \lambda &\longmapsto \text{Lip}(f_{\theta_\lambda}) \quad \text{is non-increasing.} \end{aligned}$$

This statement is the content of [CLIP, Prop. 2], however, the argument is exactly the same as in [BO13]. The intuition behind this result is that with increasing parameter  $\lambda$ , solutions of Eq. (4.4)—more precisely the corresponding networks—have smaller Lipschitz constants, i.e., tend to be more constant, which however also diminishes their expressivity. We formalize the limit cases in the following.

Assuming the realizability condition of [SB14] we know that there exists parameters  $\theta \in \Theta$  such that  $\mathcal{L}(\theta) = 0$ . This means there are parameters that fit the data perfectly. Considering the limit  $\lambda \rightarrow 0$  we obtain convergence to a solution of the unregularized problem with the smallest Lipschitz constant.

**Proposition 4.2 ([CLIP, Prop. 3]).** Let Assumptions 1 to 3 and the realizability assumption [SB14] be satisfied. If  $\theta_\lambda \rightarrow \theta^\dagger \in \Theta$  as  $\lambda \searrow 0$ , then

$$\theta^\dagger \in \operatorname{argmin} \left\{ \operatorname{Lip}(f_\theta) : \theta \in \Theta, \frac{1}{|\mathcal{T}|} \sum_{(x,y) \in \mathcal{T}} \ell(f_\theta(x), y) = 0 \right\}$$

if this problem admits a solution with  $\operatorname{Lip}(f_\theta) < \infty$ .

Furthermore, we study the effect of sending  $\lambda \rightarrow \infty$ , where we see that the network  $f_{\theta_\lambda}$  indeed tends to be constant as expected. We can explicitly characterize this constant as the closest point to barycenter of the output data, that is realizable by a neural network.

**Proposition 4.3 ([CLIP, Prop. 4]).** Let Assumptions 1 to 3 be satisfied and assume that

$$\mathcal{M} := \{y \in \mathcal{Y} : \exists \theta \in \Theta, f_\theta(x) = y, \forall x \in \mathcal{X}\} \neq \emptyset.$$

If  $\theta_\lambda \rightarrow \theta_\infty \in \Theta$  as  $\lambda \rightarrow \infty$ , then  $f_{\theta_\infty}(x) = \hat{y}$  for all  $x \in \mathcal{X}$  where

$$\hat{y} \in \operatorname{argmin}_{y' \in \mathcal{M}} \frac{1}{|\mathcal{T}|} \sum_{y \in \mathcal{T}_y} \ell(y', y).$$

### 4.2.3. Numerical Results

We briefly comment on the numerical results [CLIP]. All the experiments were implemented in Python [VD95] employing—among others—the PyTorch [Pas+19] package. We conduct experiments on the MNIST [LC10] and FashionMNIST [XRV17] datasets.

**Qualitative Example** We first apply the CLIP algorithm to regularize a one-dimensional regression problem. In [CLIP, Fig. 2] one observes the qualitative effect of the regularization. Namely, we are given noisy data from a ground truth function. The solution of the unregularized problem overfits the data and the resulting network has fluctuations and therefore a large Lipschitz constant in certain regions, where the ground truth function has a small Lipschitz constant. With increasing parameter  $\lambda$ , we observe that the networks  $f_{\theta_\lambda}$  are smoothed out and better approximate the ground truth function.

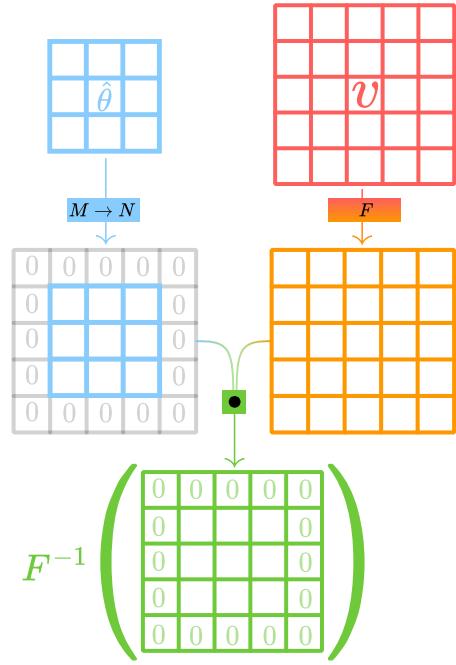
**Quantitative Evaluation of Adversarial Robustness** We additionally evaluate how robust a CLIP trained network is against adversarial attacks. Next to a standard learning rate scheduler, we also update the regularization parameter  $\lambda$  throughout the training process. Namely, we choose a target accuracy and evaluate the loss on a validation set after each epoch. If this accuracy is less than the target, we decrease  $\lambda$  and increase it if the accuracy is higher. In [CLIP, Tab. 1] we compare ourselves against a weight regularization scheme that is based on the estimate in Eq. (4.5). We see that CLIP outperforms this method in most cases.



The code for all the experiments is available at [github.com/TimRoith/CLIP](https://github.com/TimRoith/CLIP).

### 4.3. Resolution Stability via FNOs: [FNO]

In this section, we comment on the results in [FNO], concerning input robustness w.r.t. resolution changes. In practice we frequently employ the set  $\mathcal{X} = [0, 1]^{K \times N \times M}$  to represent images. However, from a modeling point of view it is more natural to assume that images are functions  $u : \Omega \rightarrow [0, 1]^K$  where  $\Omega \subset \mathbb{R}^2$  is some domain. In this sense, the space  $\mathcal{X}$  only constitutes a discretization of the space of all functions from  $\Omega$  to  $[0, 1]^K$ . The number of pixels, i.e.,  $N \cdot M$ , relates to the *resolution* of the image, where a higher number of pixels yields a higher resolution. In this sense, the resolution is merely an artifact of practical restrictions and should not be relevant for the classification. Therefore, one wants to obtain a resolution-independent classifier, which we study in the following.



**Images, Resolution and Scale** In the continuum setting we consider the  $d$ -dimensional torus  $\Omega = \mathbb{R}^d / \mathbb{Z}^d$ . An image is a function

$$u : \Omega \rightarrow [0, 1]^K$$

where  $K \in \mathbb{N}$  denotes the number of color channels, see e.g. [GW87]. In order to represent images on a computer, we discretize the domain  $\Omega$ , via a regular grid  $\Omega_N = \{x_0, \dots, x_N\}$  indexed by the set

$$\mathcal{J}_N := \{0, \dots, N - 1\}^d$$

with grid points  $x_j = j/(N - 1)$ , see e.g. [Kab22; KLM21]. For simplicity, we assume an equal number of discretization points in each dimension. Therefore,  $N^d$  is the resolution of an image discretized w.r.t.  $\mathcal{J}_N$ .

It is important to notice the differences to the concept of *scale*. A change in scale would be a change in the original image domain, for example, by zooming in on a certain area. In our scenario, we usually assume that the image  $u$  contains a certain entity to be classified. The scale roughly describes how “big” the entity is, or what percentage of the image it fills. We assume that the scale is fixed. Therefore, changing  $N$  directly influences the resolution of the discretization.

**How do Neural Networks react to Resolution Changes?** In many machine learning applications, one assumes a fixed input size and therefore a fixed resolution, assuming the same scale throughout the data. Formally, networks are defined as mappings  $f_\theta : \mathbb{R}^{K \times N \times N} \rightarrow \Delta^C$ , i.e. they only accept inputs of the fixed resolution  $N$ . In order to understand how this constraint can be weakened, we need to consider the concrete structure of our networks, which was similarly done in [KLM21; Kab22]. The architectures we consider are feed-forward networks of the form

$$f_\theta = \Phi^{\text{class}} \circ \mathcal{S} \circ \Phi^{\text{feature}}$$

where

- $\Phi^{\text{feature}}$  is the so-called feature extractor, which should be applicable independently from the input dimension,
- $\mathcal{S}$  is a function that maps inputs of any size to a fixed output dimension  $\mathbb{R}^s$ ,
- $\Phi^{\text{class}} : \mathbb{R}^s \rightarrow \Delta^C$  denotes the classification layer.

The feature extractor usually consists of convolutional layers. As in [Kab22, Ch. 2] we consider the discrete convolution of two discretized images  $u_N, v_N$  defined as

$$(u_N * v_N)(x_j) := \sum_{k \in \mathcal{J}_N} u_N(x_k) \cdot v_N(x_{j-k}) \quad (4.6)$$

where for negative indices  $j - k$  we set  $x_{j-k} := x_{j-k+N}$ . This assumes that both  $u_N$  and  $v_N$  live on the same discretization  $\mathcal{J}_N$ . For neural networks, we are interested in the convolution of an input image  $u$  and a kernel  $\theta \in \mathbb{R}^{\mathcal{J}_M}$ . Modeling spatial locality—and therefore a small support of the kernel—one usually chooses  $M$  much smaller than the resolution. This is for example motivated by the study in [HW62] which explores a similar methodology for the visual cortex of cats. However, if  $M < N$  one can not directly employ the convolution in Eq. (4.6), since there we assumed the same discretization for both inputs. In order to account for this dimension mismatch, we consider so-called *spatial zero-padding* for kernels  $\theta \in \mathbb{R}^{\mathcal{J}_M}$

$$\theta_k^{M \rightarrow N} = \begin{cases} \theta_k & \text{for } k \in \mathcal{J}_N \cap \mathcal{J}_M, \\ 0 & \text{for } k \in \mathcal{J}_N \setminus \mathcal{J}_M. \end{cases}$$

Using this method, one can define the convolution for inputs  $u_N$  of arbitrary input resolution  $N \geq M$  via

$$C(\theta)(u_N) = \theta^{M \rightarrow N} * u_N.$$

In [FNO] we refer to this as the *spatial implementation* of convolution. Up to the behavior on the boundary, this is in fact the standard implementation in most libraries, especially in PyTorch. Therefore, a feature extractor consisting of convolutional layers can take inputs of variable resolution. In fact, ignoring possible resolution changes within the

extractor—via pooling or strided convolutions—we have that  $\Phi^{\text{feature}}(u_N) \in \mathbb{R}_N^{\mathcal{J}}$  for any  $N \in \mathbb{N}$ .

The mapping  $\mathcal{S}$  can be realized as an adaptive pooling layer, see [Pas+19], which ensures a fixed output size. This methodology yields a *discretization invariant architecture*, see [Kab22; KLM21; Li+20]. This means that from a technical point of view, the network is able to produce outputs for inputs with arbitrary discretization. However, we are actually interested in a *discretization invariant functionality* (see [Kab22; KLM21; Li+20]), which also requires that the output similar for different resolutions.

**Input Interpolation** The technical possibility to handle different input sizes, as described above, usually does not perform well in practice. This is due to the fact that in the standard spatial implementation of convolution the support of the kernel changes with varying input dimension, hence the output differs, see Fig. 4.1. If a network is trained on a fixed input size, the filters are adapted to this size and therefore only create meaningful responses on this size.

A simple attempt to create a network, such that not only the architecture, but also the functionality is discretization independent—at least up to a certain degree—is input interpolation. In this case, our architecture is modified to

$$\tilde{f}_\theta = f_\theta \circ I$$

where  $I : \bigcup_{M \in \mathbb{N}} \mathbb{R}^{\mathcal{J}_M} \rightarrow \mathbb{R}^{\mathcal{J}_N}$  is an interpolation function, that maps inputs of arbitrary sizes to a fixed discretization  $\mathcal{I}_N$ . Typical choices here include nearest neighbor, bilinear or bicubic interpolation, see, e.g., [GW87]. Especially relevant in our case, is so-called *trigonometric interpolation*, where for  $v \in \mathbb{R}^{\mathcal{J}_M}$  we define

$$I^{\text{trigo}}(v) := v^M \xrightarrow{\Delta} N := F^{-1} \left( (Fv)^{M \rightarrow N} \right).$$

**Contribution in [FNO]** We study the connection between FNOs and CNNs for classification problems. We identify under which assumption the architectures are equivalent (see Lemma 4.9), but also where they are not, see Fig. 4.1. Here, we are especially interested in the multi-resolution case. We show that one layer of an FNO is equivariant with respect to trigonometric interpolation, Corollary 4.11. This is also underlined by numerical experiments, where we compare the FNO implementation to interpolation methods and a simple CNN implementation. Furthermore, we show that training equivalent CNN and FNO layers leads to different results, i.e., while they have the same forward pass, the gradients w.r.t. their parameters might differ, Lemma 4.10. Moreover, we show continuity and Fréchet-differentiability of abstract neural layers as operators between  $L^p$  spaces, see Section 4.3.2. Finally, we conduct numerical experiments supporting our theoretical findings Section 4.3.3.

### 4.3.1. Fourier Neural Operators

We want to obtain neural networks whose output does not depend on the discretization of the image  $u : \Omega \rightarrow \mathbb{R}^K$ . This raises the question whether it is possible to find a

formulation that allows us to work in the infinite dimensional setting. In [Kov+21] this issue was addressed in the setting of parametric PDEs, where the authors introduced the concept of *neural operators*. One layer of a neural operator is given as a mapping  $\mathcal{G} : L^p(\Omega) \rightarrow L^q(\Omega)$

$$\mathcal{G}(u)(x) = \sigma(\Psi(u)(x)) \quad \text{for a.e. } x \in \Omega, \quad (4.7)$$

with an affine linear part given by

$$\Psi(u) = Wu + \mathcal{K}u + b, \quad (4.8)$$

where

- $\mathcal{K} : u \mapsto \int_{\Omega} \kappa(\cdot, y) u(y) dy$  is a kernel integral operator with kernel  $\kappa : \Omega \times \Omega \rightarrow \mathbb{R}$ ,
- $W \in \mathbb{R}$  models a residual component,
- and  $b : \Omega \rightarrow \mathbb{R}$  models a bias.

By a slight abuse of notation, the activation function  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  acts as a Nemytskii operator, i.e.,

$$\sigma : v \mapsto \sigma(v(\cdot)), \quad (4.9)$$

see, e.g., [Trö10]. In Section 4.3.2 we analyze continuity and differentiability of a layer in this abstract form. However, the most relevant case for us, is when  $\mathcal{K}$  is a convolution operator, i.e.,  $\kappa(x, y) = \kappa(x - y)$  is a translation invariant kernel. In this special case  $\mathcal{G}$  is then known as layer of a *Fourier Neural Operator* (FNO) as introduced in [Li+20]. We can parameterize the kernel via its Fourier coefficients  $\hat{\theta}_k \in \mathbb{C}$

$$\kappa_{\hat{\theta}}(x) = \sum_{k \in \mathcal{I}} \hat{\theta}_k b_k(x), \quad (4.10)$$

where  $b_k(x) = \exp(2\pi i kx)$  denote the Fourier basis functions. In practice, we assume that  $\kappa_{\hat{\theta}}$  only has a finite amount of non-zero Fourier coefficients. We choose the set  $\mathcal{I}_N := \{-\lceil(N-1)/2\rceil, \dots, 0, \dots, \lfloor(N-1)/2\rfloor\}^d$  as the index set, as done in [Li+20]. Employing this set with odd  $N \in \mathbb{N}$ , we can easily enforce Hermitian symmetry  $\hat{\theta}_k = \hat{\theta}_{-k}$ , which ensures that  $\mathcal{K}_{\hat{\theta}}$  outputs real-valued functions. For now, we assume an odd number  $N$  and deal with the even case later. We note that this number of Fourier coefficients is completely independent of the spatial input discretization, which is the main advantage of FNOs.

**How to Perform Convolutions with FNOs?** In [FNO] we consider the discrete Fourier transform and its inverse

$$(Fv)_k = \frac{1}{\lambda} \sum_{j \in \mathcal{J}_N} v_j e^{-2\pi i \langle k, \frac{j}{N} \rangle}, \quad k \in \mathcal{I}_N,$$

$$\left(F^{-1}\hat{v}\right)_j = \frac{\lambda}{|\mathcal{J}_N|} \sum_{k \in \mathcal{I}_N} \hat{v}_k e^{2\pi i \langle k, \frac{j}{N} \rangle} \quad j \in \mathcal{J}_N,$$

with a normalization constant  $\lambda \in \{1, \sqrt{|\mathcal{J}_N|}, |\mathcal{J}_N|\}$ . We only consider parameters  $\hat{\theta} \in \mathbb{C}_{\text{sym}}^{\mathcal{I}_N} := F(\mathbb{R}^{\mathcal{J}_N})$ , where we can employ the convolution theorem (see e.g. [Gra14]) to define

$$K(\hat{\theta})(v) = F^{-1}(\hat{\theta} \cdot Fv) \quad \text{for } v \in \mathbb{R}^{\mathcal{J}_N}, \quad (4.11)$$

which is referred to as the FNO or spectral implementation of convolution.

**How do FNOs React to Resolution Changes?** The number of Fourier coefficients of  $\hat{\theta} \in \mathbb{C}_{\text{sym}}^{\mathcal{I}_M}$  is independent of the input resolution. Nevertheless, the point-wise multiplication in Eq. (4.11) is only defined for inputs  $v \in \mathbb{R}^{\mathcal{J}_M}$ , thus the question arises how FNOs can be adapted to dimension mismatch. Here, we use a conceptually similar idea, by employing zero-padding. However, the important and major difference is that this zero-padding is performed in the spectral domain. Assuming that  $N > M$  is odd we define

$$\hat{\theta}_k^{M \rightarrow N} = \begin{cases} \hat{\theta}_k & \text{for } k \in \mathcal{I}_N \cap \mathcal{I}_M, \\ 0 & \text{for } k \in \mathcal{I}_N \setminus \mathcal{I}_M, \end{cases}$$

and the spectral implementation of convolution for  $v_N \in \mathbb{R}^{\mathcal{J}_N}$  is given as

$$K(\hat{\theta})(v_N) := K(\hat{\theta}^{M \rightarrow N})(v_N).$$

**What Is The Difference Between Standard and FNO Implementation?** The difference between the spectral and spatial implementation is best explained in [FNO, Fig. 1], which we repeat in Fig. 4.1 for convenience. We are given a kernel  $\theta \in \mathbb{R}^{\mathcal{J}_M}$ , its unnormalized Fourier transform  $\hat{\theta} = \lambda F(\theta)$  and an input  $v_N \in \mathbb{R}^{\mathcal{J}_N}$ . If  $M = N$ , i.e., all dimensions match, we observe that spectral and spatial implementation are equivalent, see the middle row of Fig. 4.1. However, if we consider a higher resolution variant of the image with  $N > M$ , spatial zero-padding results in the kernel being localized in space and therefore the effect of convolving it with  $v_N$  changes. On the other hand for the spectral implementation we observe an equivariant behaviour. The resolution of the output changes but qualitatively the effect of the filter stays the same.

**Connection to Interpolation** As hinted in Fig. 4.1, when changing the resolution, the spectral implementation of convolution can be interpreted as a standard convolution with an interpolated kernel. In fact we observe that for  $\theta \in \mathbb{R}^{\mathcal{J}_M}, \hat{\theta} = \lambda F(\theta)$  and  $v_N \in \mathbb{R}^{\mathcal{J}_N}$  we have that

$$\begin{aligned} K(\hat{\theta})(v_N) &= F^{-1}(\hat{\theta}^{M \rightarrow N} \cdot Fv_N) \\ &= F^{-1}(F F^{-1}\hat{\theta}^{M \rightarrow N} \cdot Fv_N) \\ &= F^{-1}(F \theta^{M \xrightarrow{\Delta} N} \cdot Fv_N) = \theta^{M \xrightarrow{\Delta} N} * v_N \\ &= C(\theta^{M \xrightarrow{\Delta} N})(v_N). \end{aligned}$$

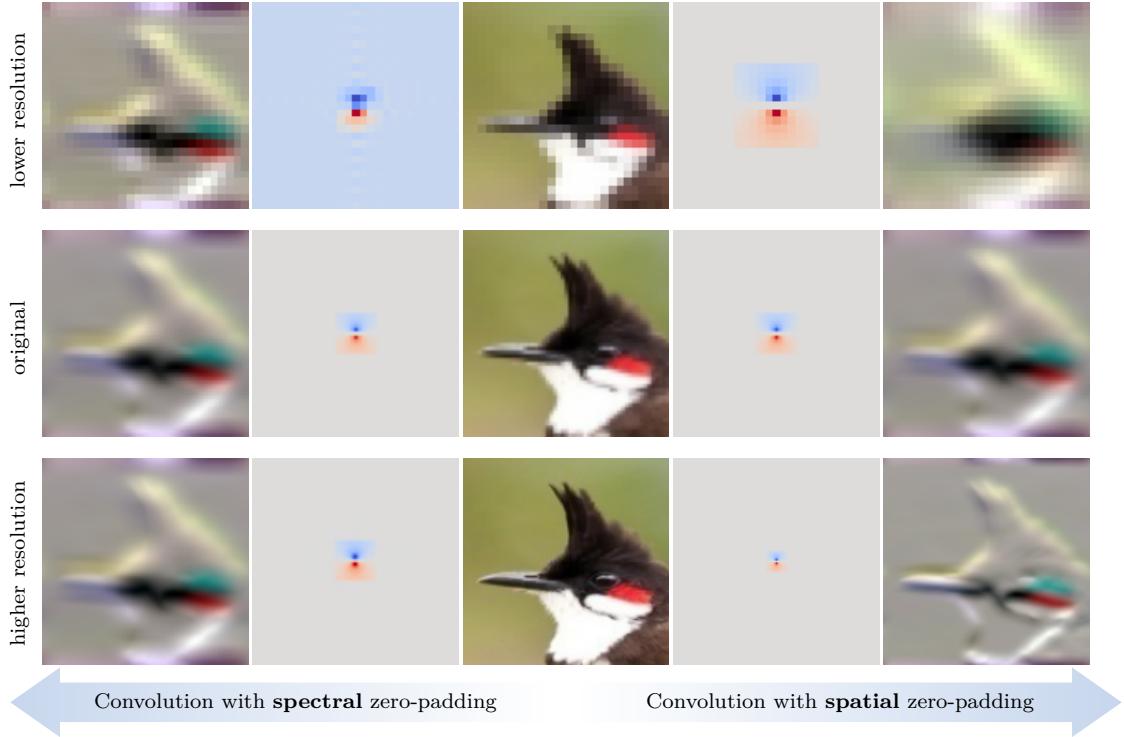


Figure 4.1.: The image is taken from [Kab22, Fig. 1]—depicting a red whiskered bulbul taken from the Birds500 dataset [Pio21]—and visualizes the different effects of spectral and spatial zero-padding.

Therefore, applying one FNO layer is equivalent to applying a standard CNN layer with a trigonometric interpolation of the kernel.

**Adaption to Even Dimensions** Zero-padding of the spectral coefficients only fulfills Hermitian symmetry in the case where  $M, N$  are odd. In order to adapt this to the even case, we employ so-called Nyquist splitting, see [Bri19]. In all our experiments, the implementation carefully employs this method, which ensures that the output of the spectral convolution is real valued. We also refer to [FNO, Sec. 3.3], where the details on this topic are given.

### 4.3.2. Analytical Results for FNOs

In this section, we comment on the theoretical findings in [FNO]. We first consider the abstract neural layer as in Eq. (4.7) for which we show continuity and Fréchet-differentiability.

**Continuity of Neural Layers** The results for neural layers in [FNO] mostly fall back to the theory of Nemytskii operators, see, e.g., [Trö10; AP93]. In order to show that the

layer  $\mathcal{G}$  is a well defined mapping from  $L^p$  to  $L^q$  one first needs to identify an exponent  $r \in [1, \infty]$  such that the affine part is a mapping  $\Psi : L^p \rightarrow L^r$ . For example if  $\kappa \in L^s$  with  $1/r + 1 = 1/p + 1/s$  it follows from Young's convolution inequality that  $\mathcal{K}$  maps to  $L^r$ , see [Gra14, Th. 1.2.12]. If then  $W = 0$  and  $b \in L^r$  we know that  $\Psi$  maps to  $L^r$ .

To ensure that the Nemytskii operator defines a mapping  $\sigma : L^r \rightarrow L^q$  one needs to assume a growth condition on  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ , see [FNO, Eq. 4] and originally [Trö10]. Under these assumptions, we obtain [FNO, Prop. 1]. Concrete examples fulfilling these assumptions are given in [FNO] borrowing concepts from [AP93; Trö10]. The most important activation function that is valid in this setting is the ReLU function

$$\text{ReLU}(x) = \max\{x, 0\}.$$

**Proposition 4.4 ([FNO, Prop. 1]).** For  $1 \leq p, q \leq \infty$  let  $\mathcal{L}$  be an operator layer given by (4.7) with an activation function  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ . If there exists  $r \geq 1$  such that

- (i) the affine part defines a mapping  $\Psi : L^p(\Omega) \rightarrow L^r(\Omega)$ ,
- (ii) the activation function  $\sigma$  generates a Nemytskii operator  $\sigma : L^r(\Omega) \rightarrow L^q(\Omega)$ ,

then it holds that  $\mathcal{L} : L^p(\Omega) \rightarrow L^q(\Omega)$ . If additionally  $\Psi$  is a continuous operator on the specified spaces and the function  $\sigma$  is continuous, or uniformly continuous in the case  $q = \infty$ , the operator  $\mathcal{L} : L^p(\Omega) \rightarrow L^q(\Omega)$  is also continuous.

**Differentiability of Neural Layers** We furthermore consider Fréchet differentiability of a neural layer w.r.t. the input variable. This can also be transferred to differentiability w.r.t. the parameters as we show in [FNO, Ex. 4]. Conceptually, the main result we repeat here is similar to the one on continuity of the last paragraph. The major difference is that we also need to assume differentiability of the activation function, also assuming a growth condition on its derivative. The ReLU function can therefore not be chosen in this setting, however the smooth approximation called GELU ([HG16])

$$\text{GELU}(x) := x \Phi(x)$$

where  $\Phi$  is the CDF of the standard normal distribution, can be employed.

**Proposition 4.5 ([FNO, Prop. 2]).** For  $1 \leq p, q \leq \infty$ , let  $\mathcal{L}$  be an operator layer given by (4.7) with affine part  $\Psi$  as in (4.8). If there exists  $r > q$ , or  $r = q = \infty$  such that

- (i) the affine part is a continuous operator  $\Psi : L^p(\Omega) \rightarrow L^r(\Omega)$ ,
- (ii) the activation function  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is continuously differentiable
- (iii) and the derivative of the activation function generates a Nemytskii operator  $\sigma' : L^r(\Omega) \rightarrow [L^r(\Omega) \rightarrow L^s(\Omega)]$  with  $s = rq/(r - q)$ ,

then it holds that  $\mathcal{L} : L^p(\Omega) \rightarrow L^q(\Omega)$  is Fréchet-differentiable in any  $v \in L^p(\Omega)$  with Fréchet-derivative  $D\mathcal{L}(v) : L^p(\Omega) \rightarrow L^q(\Omega)$

$$D\mathcal{L}(v)(h) = \sigma'(\Psi(v)) \cdot \tilde{\Psi}(h),$$

where  $\tilde{\Psi}$  denotes the linear part of  $\Psi$ , i.e.,  $\tilde{\Psi} = \Psi - b$ .

**Convertibility Between FNOs and CNNs** The following lemma formalizes the intuition that FNOs and CNNs are equivalent in certain settings. Given inputs of a fixed discretization  $J_N$  and parameters  $\theta \in \mathbb{R}^{J_M}$ , the standard convolution implementation is equivalent to the spectral one w.r.t. the parameters  $\hat{\theta} = \lambda F(\theta^{M \rightarrow N})$ . The subtle but important point here, is that the number of spectral parameters needs to be equal to the input size in order to achieve equivalence. In [FNO, Fig.3] we observe numerically that the number of used spectral coefficients actually needs to match the input size in order to achieve equivalence.

**Lemma 4.9 ([FNO, Lem. 3]).** Let  $M \leq N$  both be odd and let  $T : \mathbb{R}^{J_N} \rightarrow \mathbb{C}^{I_N}$  be defined for  $\theta \in \mathbb{R}^{J_N}$  as  $T(\theta) = \lambda F(\theta)$ . For any  $\theta \in \mathbb{R}^{J_M}$  and  $v \in \mathbb{R}^{J_N}$  it holds true that

$$C(\theta)(v) = K(T(\theta^{M \rightarrow N}))(v)$$

and for any  $\hat{\theta} \in \mathbb{C}_{\text{sym}}^{I_M}$  and  $v \in \mathbb{R}^{J_N}$  it holds true that

$$K(\hat{\theta})(v) = C(T^{-1}(\hat{\theta}^{M \rightarrow N}))(v).$$

Together with [FNO, Fig. 3] we see that in order to convert a CNN to a FNO we need a large number of spectral parameters. This is also connected to the fact that spatial locality can only be expressed using more spectral coefficients. Therefore, one might think that FNOs are infeasible due to memory requirements. However, it turns out that directly optimizing over the spectral parameters leads to a comparable performance, already for a low number of Fourier coefficients, which is reported in the blue curve in [FNO, Fig. 3]. This hints that training a FNO via gradient descent leads to different weights, that are not simply a conversion of a similarly trained CNN. The reason for these differences is that the gradients of spectral and standard convolutional layers, that have the same forward pass, are not directly convertible via the Fourier transformation. In fact, one obtains an additional factor, which we formalize in the following lemma.

**Lemma 4.10 ([FNO, Lem. 4]).** For odd  $N \in \mathbb{N}$  and  $v, \theta \in \mathbb{R}^{J_N}$  and  $\hat{\theta} = T(\theta)$  it holds true that

$$\nabla_{\hat{\theta}} K(\hat{\theta})(v) = \frac{1}{|J_N|} T \left( \nabla_{\theta} C(\theta)(v) \right).$$

**Interpolation Equivariance** The last result considers the main motivation of the whole section, namely, the resolution invariance of an FNO layer. However, we can not allow an arbitrary resizing operation of the input. Employing the spectral implementation of convolution layer, we can show equivariance w.r.t. trigonometric interpolation of the input.

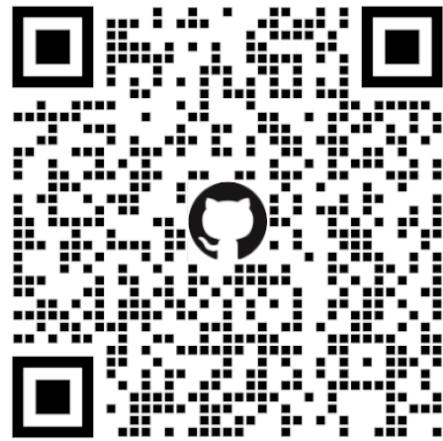
**Corollary 4.11.** For  $\hat{\theta} \in \mathbb{C}_{\text{sym}}^{I_M}$ ,  $v \in \mathbb{R}^{J_N}$ ,  $M \leq N$  it holds true for any  $L \geq M$  that

$$K(\hat{\theta})(v^N \xrightarrow{\Delta} L) = (K(\hat{\theta})(v))^N \xrightarrow{\Delta} L.$$

### 4.3.3. Numerical Results

In the numerical section of [FNO] we study the convertibility of CNNs to FNOs as discussed in Section 4.3.2 and the resolution invariance of the different proposed approaches. Again we employ—among others—the PyTorch package [Pas+19]. The experiments are conducted on the Fashion-MNIST [XRV17] and a former version of the BIRDS500 [Pio21] dataset. We remark that our implementation carefully treats the case of even kernel or input sizes, via Nyquist splitting. This allows us to apply all the derived results for the odd and also the even case.

**Convertibility and Training Differences** As already described in Section 4.3.2, the first experiment, displayed in [FNO, Fig. 3], studies how CNNs can be converted to FNOs. Here, we employ a network with two convolutional layers for feature extraction and a linear classification layer. We train this network—in the spatial implementation—on the FashionMNIST dataset [XRV17] employing varying spatial kernel sizes. Here, we see that for  $M = 5$  the spatial implementation already has the best performance and adding more parameters does not increase the performance. We then convert a set of spatial parameters to the Fourier formulation, employing again a varying number of spectral coefficients. It turns out that the requirement of [FNO, Lem. 3], that the number of coefficients must match the input dimension, is indeed relevant. We only obtain the full performance of the CNN if we use all spectral parameters. However, the example also visualizes [FNO, Lem. 4], namely that training a FNO conceptually leads to a different set of parameters. Optimizing over the spectral parameters yields a comparable performance already for a smaller number of coefficients.



The code for all the experiments is available at <https://github.com/samirak98/FourierImaging>.

**Resolution Invariance** In the second experiment we study the resolution invariance of the three discretization independent architectures we considered, namely the naive CNN adaptation, input interpolation and the FNO implementation. We first train the convolutional architecture as in the previous paragraph on the FashionMNIST dataset, that has an input size of  $28 \times 28$ . We resize the input data via trigonometric and bilinear interpolation—referred to as data sizing—to simulate multi-resolution data. One expects that the classification performance drops when the data is resized to a smaller size, since this step loses information. However, resizing the images to higher resolutions should yield the same performance.

In [FNO, Fig. 4] we see that the simple adaption does not perform well both in the lower and the higher resolution setting. Employing input interpolation improves the performance in the lower resolution setting. In the higher resolution regime, we obtain a constant performance, which is expected. Finally, we see that the FNO—which was converted from the CNN—performs as well as input interpolation, which justifies the resolution independence of this architecture.

In the second experiment, we trained a ResNet18 [He+16a] on a former version of the BIRDS500 dataset [Pio21], with an input size of  $112 \times 112$ . Concerning the naive adaption and the input interpolation, we observe the same behavior as in the previous example. However, the FNO variant performs slightly worse, which is surprising, especially in the higher resolution regime. In [FNO] we conclude that this is due the dimension changes within the ResNet architectures. These changes occur due to strided convolutions or pooling operations between the layers. In fact, in order to achieve the performance as displayed in [FNO, Fig. 4 (b)] we replaced every striding by a trigonometric interpolation, which fits better into the FNO framework and vastly improves the performance. Therefore, we summarize that architecture intern dimension changes—which are very common in practice—can potentially hinder resolution invariance.

## 4.4. Sparsity via Bregman Iterations: [BREG-I]

Starting from the first simple neural architectures (e.g. [Ros58]) the size of typical architectures has grown significantly in the last years [Hoe+21]. While this development allowed to solve increasingly difficult tasks with neural networks, it also lead to immense computational requirements, both for the training and evaluation of the net. Therefore, the question arises how the evaluation and training of large neural networks can be made more efficient. Some of the approaches here include (see [Gho+21] for a comprehensive overview):

- **Architecture Optimization:** Designing an architecture that can be trained to have the same performance as a comparable one, while having less parameters, e.g. [EMH19; How+17].
- **Quantization:** Lowering the precision of the machine numbers of the parameters, e.g. [Ban+18; CBD14].
- **Knowledge distillation:** Employing a trained larger network to learn a smaller network in a student-teacher approach, e.g. [Sch92; HVD15].
- **Pruning:** Parameters with small saliency, i.e., parameters that contribute little to the effective output of the network are removed or set to zero, in order to obtain a sparse weight matrix or a smaller architecture, see e.g. [LDS89; HSW93].

From the above methods, the pruning approach is most influential to our work in [BREG-I]. Namely, the concept of compressing the neural network by sparsifying the weight matrices is similarly employed. However, there are a few deviations from the classical pruning framework, which we remark in the following.

**Sparsity via Optimization** In [LDS89; HSW93] one assumes to be given a trained neural network, which is then compressed after training based on some criterion. The authors in [CFP97] employ an iterative pruning scheme, where a similar criterion is employed in order to throw away certain network weights after each step. In our work we want to employ a  $L^1$  type penalty on weight matrices  $W \in \mathbb{R}^{n \times m}$ ,

$$\|W\|_1 = \sum_{i=1}^n \sum_{j=1}^m |W_{ij}|$$

which has been widely employed to obtain sparse solution, e.g. [CM73]. In [Tib96] this yields to the so-called Lasso problem

$$\min_{\theta} \mathcal{L}(\theta) + \lambda \|\theta\|_1, \quad \lambda > 0,$$

where we extend the  $L^1$  norm to  $\Theta$  as discussed in [Section 4.4.5](#). This problem can for example be solved by a proximal gradient descent iteration

$$\theta^{(k+1)} = \text{soft shrinkage}(\theta - \tau \nabla \mathcal{L}(\theta^{(k)})) \quad (4.12)$$

where the shrinkage operator is defined in [Example 4.18](#). This iteration was employed in [Nit14; RVV20; Red+16] to train sparse neural networks.

**Sparse-to-Sparse Training** All of the sparsity-based methods mentioned before start with dense weight matrices and only decrease the number of parameters during or at the end of the training process. Our approach yields an iteration, where the network is sparse throughout the iteration. In fact we start with only very few non-zero parameters and only activate necessary weights during training. This paradigm is known as sparse-to-sparse or evolutionary training [Moc+18; DZ19; Evc+20; DYJ19; Fu+22; Hua+16; Liu+21].

**Contribution in [BREG-I]** Our work falls into the regime of sparse-to-sparse training. However, instead of relying on some heuristic growth strategy, we employ the concept of inverse scale flows (see [Section 4.4.1](#)) which allows us to obtain a optimization-driven framework, with a time-continuous interpretation. We propose a stochastic variant of linearized Bregman iterations ([Section 4.4.3](#)) and employ it to train a sparse neural network. We show monotonic decrease of the loss in the stochastic setting—which is not possible in the case of proximal gradient descent [Eq. \(4.12\)](#)—and convergence of the iterates under additional convexity assumptions, see [Section 4.4.4](#). Finally, we demonstrate the numerical efficiency of the method (see [Section 4.4.5](#)) and provide an interesting applications for neural architecture search, which was further developed in [BREG-II].

#### 4.4.1. Preliminaries on Convex Analysis and Bregman Iterations

We first review some necessary concepts from convex analysis that allow us to introduce the framework in [BREG-I]. We refer to [BB18; Roc97; BC11] for a more exhaustive introduction to the topics. The functional  $J$  is called lower semicontinuous if  $J(u) \leq \liminf_{n \rightarrow \infty} J(u_n)$  holds for all sequences  $(u_n)_{n \in \mathbb{N}} \subset \Theta$  converging to  $u$ .

**Definition 4.12.** Given a Hilbert space  $\Theta$  and a functional  $J : \Theta \rightarrow (-\infty, \infty]$ .

1. The functional  $J$  is called convex, if

$$J(\lambda \bar{\theta} + (1 - \lambda)\theta) \leq \lambda J(\bar{\theta}) + (1 - \lambda)J(\theta), \quad \forall \lambda \in [0, 1], \bar{\theta}, \theta \in \Theta. \quad (4.13)$$

2. The effective domain of  $J$  is defined as  $\text{dom}(J) := \{\theta \in \Theta : J(\theta) \neq \infty\}$  and  $J$  is called proper if  $\text{dom}(J) \neq \emptyset$ .

In the following we want to consider functionals  $J$  that are convex, but not necessarily differentiable. Therefor, we define the subdifferential.

**Definition 4.13.** of a convex and proper functional  $J : \Theta \rightarrow (-\infty, \infty]$  at a point  $\theta \in \Theta$  as

$$\partial J(\theta) := \left\{ p \in \Theta : J(\theta) + \langle p, \bar{\theta} - \theta \rangle \leq J(\bar{\theta}), \forall \bar{\theta} \in \Theta \right\}. \quad (4.14)$$

If  $J$  is differentiable, then the subdifferential coincides with the classical gradient (or Fréchet derivative). We denote  $\text{dom}(\partial J) := \{\theta \in \Theta : \partial J(\theta) \neq \emptyset\}$  and observe that  $\text{dom}(\partial J) \subset \text{dom}(J)$ .

The main algorithm in this section are so-called Bregman iterations, for which we first define the Bregman distance.

**Definition 4.14 (Bregman Distance).** Let  $J : \Theta \rightarrow (-\infty, \infty]$  be a proper, convex functional. Then we define for  $\theta \in \text{dom}(\partial J), \bar{\theta} \in \Theta$

$$D_J^p(\bar{\theta}, \theta) := J(\bar{\theta}) - J(\theta) - \langle p, \bar{\theta} - \theta \rangle, \quad p \in \partial J(\theta). \quad (4.15)$$

For  $p \in \partial J(\theta)$  and  $\bar{p} \in \partial J(\bar{\theta})$  we define the *symmetric* Bregman distance as

$$D_J^{\text{sym}}(\bar{\theta}, \theta) := D_J^p(\bar{\theta}, \theta) + D_J^{\bar{p}}(\theta, \bar{\theta}). \quad (4.16)$$

Intuitively, the Bregman distance  $D_J^p(\bar{\theta}, \theta)$ , measures the distance of  $J$  to its linearization around  $\theta$ , see Fig. 4.2. If  $J$  is differentiable, then the subdifferential is single valued—we can suppress the sup script  $p$ —and we have

$$D_J(\bar{\theta}, \theta) = J(\bar{\theta}) - J(\theta) - \langle \nabla J(\theta), \bar{\theta} - \theta \rangle.$$

**Example 4.15.** For  $\Theta = \mathbb{R}^n$  and  $J = \frac{1}{2} \|\cdot\|_2^2$  we see that  $\partial J(\theta) = \{\theta\}$  and therefore

$$\begin{aligned} D_J^p(\bar{\theta}, \theta) &= \frac{1}{2} \langle \bar{\theta}, \bar{\theta} \rangle - \frac{1}{2} \langle \theta, \theta \rangle - \langle \theta, \bar{\theta} - \theta \rangle \\ &= \frac{1}{2} \langle \bar{\theta}, \bar{\theta} \rangle + \frac{1}{2} \langle \theta, \theta \rangle - \langle \theta, \bar{\theta} \rangle \\ &= \frac{1}{2} \|\bar{\theta} - \theta\|_2^2 = J(\bar{\theta} - \theta). \end{aligned}$$

We can easily see that in general it is neither definite, symmetric nor fulfills the triangle inequality, hence it is not a metric. However, it fulfills the two distance axioms

$$D_J^p(\bar{\theta}, \theta) \geq 0, \quad D_J^p(\theta, \theta) = 0, \quad \forall \bar{\theta} \in \Theta, \theta \in \text{dom}(\partial J). \quad (4.17)$$

The same holds for the symmetric Bregman distance, where additionally—as the name suggests—the symmetry property is fulfilled. The last concept that is crucial in [BREG-I] is the so-called proximal operator.

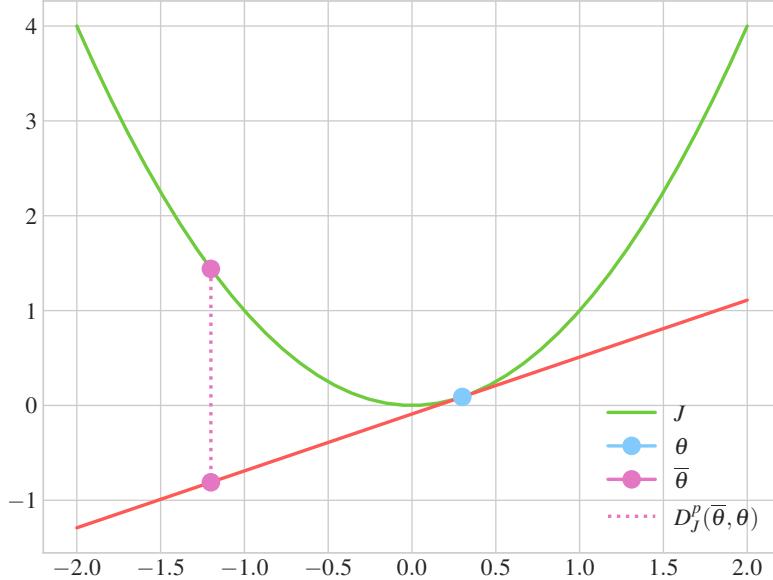


Figure 4.2.: Visualization of the Bregman distance.

**Definition 4.16.** Let  $J : \Theta \rightarrow (-\infty, \infty]$  be convex, proper and lower semicontinuous functional, then we define the *proximal operator* as

$$\text{prox}_J(\bar{\theta}) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{2} \|\theta - \bar{\theta}\|^2 + J(\theta).$$

If  $J$  is additionally a closed function, i.e., its sublevel sets

$$N_\alpha = \{\theta \in \text{dom } J : J(\theta) \leq \alpha\}$$

are closed for every  $\alpha \in \mathbb{R}$  then we have that the function  $\tilde{J} = \frac{1}{2} \|\theta - \cdot\|^2 + J(\theta)$  is closed, proper and *strongly convex* and therefore has a unique minimizer, see [Roc97, Thm. 27.1]. Additionally, one often considers a regularization parameter  $\lambda > 0$  and is then interested in  $\text{prox}_{\lambda J}$ .

**Remark 4.17.** The optimality conditions for  $\theta = \text{prox}_{\lambda J}(\bar{\theta})$  yield

$$\bar{\theta} - \theta \in \lambda \partial J(\theta).$$

For a proper, closed and convex function we obtain

$$\theta = (I + \lambda \partial J)^{-1}(\bar{\theta})$$

where  $(I + \lambda \partial J)^{-1}$  is called the *resolvent* and is a one-to-one mapping (see [PB+14, Ch. 3.2]) which justifies the equality in the above equation. If  $J$  is differentiable, then we have  $\partial J = \{\nabla J\}$  and therefore,

$$\text{prox}_{\lambda J} = (I + \lambda \nabla J)^{-1}.$$

△

In the following we list two relevant examples for the application in [BREG-I; BREG-II].

**Example 4.18.** If  $J = \|\cdot\|$  is a norm and  $\lambda > 0$  then we have that (see e.g. [PB+14])

$$\text{prox}_{\lambda J}(\bar{\theta}) = \bar{\theta} - \text{Proj}_{\|\cdot\|^*}(\bar{\theta}/\lambda)$$

where  $\text{Proj}_{\|\cdot\|^*}$  denotes the projection operator w.r.t. the dual norm  $\|\theta\|^* = \sup\{|\langle f, \theta \rangle| : f \in \Theta^*\}$ . In the case of  $\ell^p$  norms on  $\mathbb{R}^n$  we know that that

$$\|\theta\|_p^* = \|\theta\|_q$$

with  $1/p + 1/q = 1$  with the notational convention of  $1/\infty = 0$ . Especially relevant are the cases  $p \in \{1, 2\}$ . Here, we then have that

$$\text{prox}_{\lambda\|\cdot\|_2}(\bar{\theta}) = \bar{\theta} \left( 1 - \min \left\{ \frac{\lambda}{\|\bar{\theta}\|_2}, 1 \right\} \right) = \begin{cases} \bar{\theta}(1 - \lambda/\|\bar{\theta}\|_2) & \text{if } \|\bar{\theta}\|_2 \geq \lambda \\ 0 & \text{else} \end{cases}$$

and for  $i = 1, \dots, n$

$$\text{prox}_{\lambda\|\cdot\|_1}(\bar{\theta})_i = \text{sign}(\bar{\theta}_i) \max \left\{ |\bar{\theta}_i| - \lambda, 0 \right\} = \begin{cases} \bar{\theta}_i - \lambda & \text{if } \bar{\theta}_i > \lambda \\ 0 & \text{if } |\bar{\theta}_i| \leq \lambda \\ \bar{\theta}_i + \lambda & \text{if } \bar{\theta}_i < -\lambda \end{cases}$$

the so called *soft shrinkage operator*.

**Example 4.19 (Group Norms).** Another relevant functional  $J$  is the group norm  $\ell_{1,2}$  that—in the context of sparse neural networks—was first employed by [Sca+17]. Here, we assume that the parameters in  $\Theta$  can be grouped in a collection of parameters  $\mathcal{G}$ , for which we choose

$$J(\theta) = \sum_{g \in \mathcal{G}} \sqrt{\#\mathcal{G}} \|g\|_2$$

In this case the proximal operator is given as

$$\text{prox}_{\lambda J}(\bar{\theta})_g = g \max \left\{ 1 - \min \left\{ \frac{\lambda \sqrt{\#\mathcal{G}}}{\|g\|_2}, 1 \right\}, 0 \right\}$$

**Bregman Iterations** Our goal is to minimize a function  $\mathcal{L}$  while simultaneously obtaining a low value w.r.t. the functional  $J$ . One popular approach considers the regularized problem

$$\min_{\theta} \mathcal{L}(\theta) + \lambda J(\theta) \quad \lambda > 0,$$

see e.g. [Tik43; CP08; DDD04; FS06; FNW07; Cha04; CP11], which however influences the minimizers of the original problem  $\min_{\theta} \mathcal{L}(\theta)$ . In the derivation of the Bregman iterations one can take a different viewpoint. Assume that we want to employ an iterative scheme, where in each step we want minimize  $\mathcal{L}$  while penalizing the distance to the previous iterate. For a stepping parameter  $\tau$  and starting from some  $\theta^{(0)} \in \Theta$  this yields the update

$$\theta^{(k+1)} = \operatorname{argmin}_{\theta} \tau \mathcal{L}(\theta) + \frac{1}{2} \|\theta - \theta^{(k)}\|^2 = \operatorname{prox}_{\tau \mathcal{L}}(\theta^{(k)}). \quad (4.18)$$

At first sight this is either known as the proximal point algorithm [Bre67] or a minimizing movement scheme [De 93]. If  $\mathcal{L}$  is differentiable, this update can be rewritten as

$$\theta^{(k+1)} = (I + \tau \nabla \mathcal{L})^{-1} \theta^{(k)} \Leftrightarrow \frac{1}{\tau} (\theta^{(k+1)} - \theta^{(k)}) = -\nabla \mathcal{L}(\theta^{(k+1)})$$

which is a implicit Euler discretization ([Eul24]) of the time-continuos gradient flow

$$\partial_t \theta_t = -\nabla \mathcal{L}(\theta_t).$$

We see that the penalization term in Eq. (4.18) is in fact the Bregman distance of the functional  $\frac{1}{2} \|\cdot\|^2$ . In order to incorporate an arbitrary convex functional  $J$ —and therefore allow each iterate to only slightly deviate w.r.t. the Bregman distance of  $J$  to the previous iterate—we employ  $D_J^{p^{(k)}}(\cdot, \theta^{(k)})$  as a penalization term. In order to obtain a update scheme for the subgradients, we observe

$$\theta = \operatorname{argmin}_{\theta \in \Theta} D_J^{p^{(k)}}(\theta, \theta^{(k)}) + \tau \mathcal{L}(\theta) \quad (4.19)$$

$$\Leftrightarrow p^{(k)} + \tau \nabla \mathcal{L}(\theta) \in \partial J(\theta). \quad (4.20)$$

This finally yields *Bregman iteration* of [Osh+05]

$$\theta^{(k+1)} = \operatorname{argmin}_{\theta \in \Theta} D_J^{p^{(k)}}(\theta, \theta^{(k)}) + \tau \mathcal{L}(\theta), \quad (4.21a)$$

$$p^{(k+1)} = p^{(k)} - \tau \nabla \mathcal{L}(\theta^{(k+1)}) \in \partial J(\theta^{(k+1)}). \quad (4.21b)$$

The very nature of Bregman iterations means starting with a iterate  $\theta^{(0)}$  that has a low value in  $J$ —preferably  $J(\theta^{(0)}) = 0$ —and only increase  $J(\theta^{(k)})$  gradually as  $k$  increases.

**Remark 4.20.** Originally, the iterations were employed for solving inverse problems. Here, we are given a forward operator  $A : \Theta \rightarrow \tilde{\Theta}$  and a noisy measurement  $f = A\theta + \delta$  where  $\delta \in \tilde{\Theta}$  is additive noise. The loss function is then of the form

$$\mathcal{L} = \frac{1}{2} \|A \cdot - f\|_2^2$$

for which one can show that the Bregman iterations converge to a solution of

$$\min \{J(\theta) : A\theta = f\}, \quad (4.22)$$

see e.g. [Osh+05]. In comparison the concept of adding a regularizing term with parameter  $\lambda > 0$ , i.e. considering the problem

$$\min_{\theta} \mathcal{L}(\theta) + \lambda J(\theta)$$

actually modifies the minimizers. In this sense Bregman iterations do not introduce a bias.  $\triangle$

**Example 4.21.** In order to get an intuition about the behavior of Bregman iterations, we consider an image denoising task. I.e. we are given a noisy image  $\mathbb{R}^{n \times m} \ni f = u + \delta$  where  $\delta \in \mathbb{R}^{n \times m}$  is additive noise. In order to obtain  $u \in \mathbb{R}^{n \times n}$  from  $f$  we employ the TV functional [ROF92]

$$J(u) = TV(u) := \sum_{i,j} \sqrt{|u_{i+1,j} - u_{i,j}|^2 + |u_{i,j+1} - u_{i,j}|^2}$$

together with the loss function  $\mathcal{L}(u) := \frac{1}{2} \|u - f\|_2^2$ . We start with an image  $u^{(0)}$  such that  $TV(u^{(0)}) = 0$ , i.e. a constant image. In Fig. 4.3 we visualize the iteration. At lower iterations  $u^{(k)}$  only displays features on a larger scale, while at the end, the iteration converges back to smallest possible scale, the noisy data. In order to obtain an appropriate denoising, one needs to employ a early stopping here. This fits well to the insight from Eq. (4.22) since here the forward operator is the identity, i.e.

$$\left\{ u : \frac{1}{2} \|u - f\|^2 = 0 \right\} = \{f\}.$$

It should also be noted that this example only serves a explanatory purpose. In practice directly applying Eq. (4.21) for  $J = TV$  can become infeasible since the first minimization problem is expensive.

If  $J = \frac{1}{2} \|\cdot\|_2^2$  as in Example 4.15 then this amounts to the step

$$\theta^{(k+1)} = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{2} \|\theta - \theta^{(k)}\|_2^2 + \tau \mathcal{L}(\theta),$$

where the optimality conditions then yield

$$\theta^{(k+1)} - \theta^{(k)} + \tau \nabla \mathcal{L}(\theta^{(k+1)}) = 0 \Leftrightarrow \theta^{(k+1)} = \theta^{(k)} - \tau \nabla \mathcal{L}(\theta^{(k+1)})$$

which is a standard implicit Euler iteration. The time continuous flow for  $\tau \rightarrow 0$  is known as the *inverse scale space* flow [Bur+06; Bur+07],

$$\begin{cases} \dot{p}_t = -\nabla \mathcal{L}(\theta_t), \\ p_t \in \partial J(\theta_t), \end{cases}$$

where again for  $J = \frac{1}{2} \|\cdot\|_2^2$  we obtain that  $\partial J(\theta_t) = \theta_t$  and therefore obtain the standard gradient flow. Hence we see, that the inverse scale space flow is a generalization of the standard gradient flow.

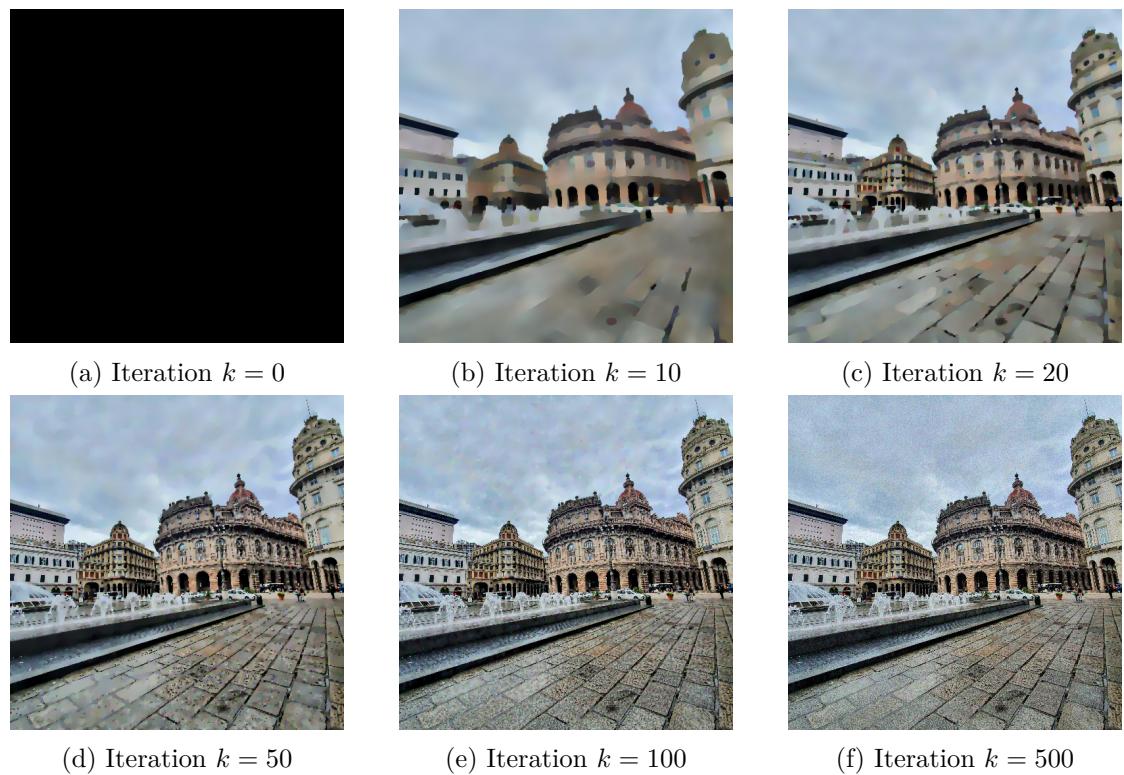


Figure 4.3.: Bregman iterations for image denoising in [Example 4.21](#).

#### 4.4.2. Linearized Bregman Iterations and Mirror Descent

The minimization step in Eq. (4.21) is infeasible for large scale applications, especially in our setting of neural networks. Therefore, we employ the idea introduced in [Yin+08; COS09]. Here, we first linearize the loss function around the previous iterate,

$$\mathcal{L}(\theta) \approx \mathcal{L}(\theta^{(k)}) + \langle \nabla \mathcal{L}(\theta^{(k)}), \theta - \theta^{(k)} \rangle.$$

The next step is to replace  $J$  with the strongly convex elastic net regularization

$$J_\delta := J + \frac{1}{2\delta} \|\cdot\|_2^2. \quad (4.23)$$

The minimization step then transforms to

$$\begin{aligned} & \operatorname{argmin}_{\theta \in \Theta} D_{J_\delta}^{p^{(k)}} (\theta, \theta^{(k)}) + \tau \langle \nabla \mathcal{L}(\theta^{(k)}), \theta \rangle \\ &= \operatorname{argmin}_{\theta \in \Theta} J(\theta) + \frac{1}{2\delta} \|\theta\|_2^2 - \langle p^{(k)}, \theta \rangle + \tau \langle \nabla \mathcal{L}(\theta^{(k)}), \theta \rangle \\ &= \operatorname{argmin}_{\theta \in \Theta} J(\theta) + \frac{1}{2\delta} \left\| \theta - \delta(p^{(k)} - \tau \nabla \mathcal{L}(\theta^{(k)})) \right\|_2^2 - \underbrace{\left\| p^{(k)} - \tau \nabla \mathcal{L}(\theta^{(k)}) \right\|_2^2}_{\text{constant in } \theta} \\ &= \operatorname{prox}_{\delta J} \left( \delta \left( p^{(k)} - \tau \nabla \mathcal{L}(\theta^{(k)}) \right) \right). \end{aligned} \quad (4.24)$$

Note that here  $p^{(k)}$  is a subgradient of  $J_\delta$  at  $\theta$  therefore we derive the subgradient update rule

$$p^{(k+1)} := p^{(k)} - \tau \nabla \mathcal{L}(\theta^{(k)}).$$

This finally yields the linearized Bregman iterations

$$p^{(k+1)} = p^{(k)} - \tau \nabla \mathcal{L}(\theta^{(k)}), \quad (4.25a)$$

$$\theta^{(k+1)} = \operatorname{prox}_{\delta J}(\delta p^{(k+1)}). \quad (4.25b)$$

The last line is equivalent to  $p^{(k+1)} \in \partial J_\delta(\theta^{(k+1)})$  for which we obtain the continuous linearized flow

$$\begin{aligned} \dot{p}_t &= -\nabla \mathcal{L}(\theta_t), \\ p_t &\in \partial J_\delta(\theta_t), \end{aligned}$$

see [Bur+06; Bur+07].

**Connections To Mirror Descent** As already noticed by [Vil+23] linearized Bregman iterations are equivalent to mirror descent in some situations. We show the equivalence in the following, where we employ similar arguments as in [BT03]. One assumes to be given a differentiable and strongly convex function  $h : \Theta \rightarrow \mathbb{R}$ , i.e.,

$$h(\bar{\theta}) - h(\theta) - \langle \nabla h(\theta), \bar{\theta} - \theta \rangle \geq \frac{1}{2} \|\bar{\theta} - \theta\|_2^2$$

for all  $\theta, \bar{\theta} \in \Theta$ . The mirror descent update then reads ([NY83; BT03])

$$\theta^{(k+1)} = \nabla h^* \left( \nabla h(\theta^{(k)}) - \tau \mathcal{L}(\theta^{(k)}) \right) \quad (4.26)$$

where  $h^*$  denotes the Fenchel conjugate

$$h^*(p) = \sup_{\theta} \langle p, \theta \rangle - h(\theta)$$

with the gradient (see [BV04])

$$\nabla h^*(p) = \operatorname{argmax}_{\theta} \{ \langle p, \theta \rangle - h(\theta) \}.$$

Therefore, we see that Eq. (4.26) can be written as

$$\begin{aligned} \theta^{(k+1)} &= \operatorname{argmax}_{\theta} \left\{ \langle \nabla h(\theta^{(k)}) - \tau \mathcal{L}(\theta^{(k)}), \theta \rangle - h(\theta) \right\} \\ &= \operatorname{argmax}_{\theta} \left\{ -D_h(\theta, \theta^{(k)}) - \tau \langle \mathcal{L}(\theta^{(k)}), \theta \rangle \right\} \\ &= \operatorname{argmin}_{\theta} \left\{ D_h(\theta, \theta^{(k)}) + \tau \langle \mathcal{L}(\theta^{(k)}), \theta \rangle \right\} \end{aligned}$$

which was our starting point to derive linearized Bregman iterations for  $h = J_\delta$  in Eq. (4.24). In fact, we can always find a convex functional  $J : \Theta \rightarrow \mathbb{R}$  such that  $h = J + \frac{1}{2} \|\cdot\|_2^2$  for which we see, that Eq. (4.25) is a more general formulation of Eq. (4.26).

#### 4.4.3. Stochastic and Momentum Variants

We want to employ linearized Bregman iterations to train a neural network. As mentioned in Section 4.1.2 we therefore do not compute the full gradient of  $\mathcal{L}$  but rather a mini-batched variant. This yields stochastic Bregman Iterations

$$\begin{aligned} &\text{draw } \omega^{(k)} \text{ from } \Omega \text{ using the law of } \mathbb{P}, \\ &g^{(k)} := g(\theta^{(k)}; \omega^{(k)}), \\ &v^{(k+1)} := v^{(k)} - \tau^{(k)} g^{(k)}, \\ &\theta^{(k+1)} := \operatorname{prox}_{\delta J}(\delta v^{(k+1)}), \end{aligned} \quad (4.27)$$

which we also abbreviate as the *LinBreg* algorithm in the following. The basic update scheme is given as the linearized Bregman iterations from [Osh+05], however the presence of a stochastic gradient estimator significantly complicates the convergence analysis, as observed in Section 4.4.4. However, this algorithm can now be efficiently employed to train a neural network. For the analogous stochastic mirror descent algorithm we refer to [Nem+09].

**Momentum Variant** Typically, the learning process of a neural network can be improved by introducing a momentum term (see e.g. [Nes83; Qia99]) in the optimizer. In our case this can be achieved, by replacing the gradient update on the subgradient variable. In [BREG-I] we first consider the inertia version of the gradient flow as

$$\begin{cases} \gamma \ddot{v}_t + \dot{v}_t = -\nabla \mathcal{L}(\theta_t), \\ v_t \in \partial J_\delta(\theta_t). \end{cases}$$

for which the discretization then reads

$$\begin{aligned} m^{(k+1)} &= \beta^{(k)} m^{(k)} + (1 - \beta^{(k)}) \tau^{(k)} g^{(k)}, \\ v^{(k+1)} &= v^{(k)} - m^{(k+1)}, \\ \theta^{(k+1)} &= \text{prox}_{\delta J}(\delta v^{(k+1)}). \end{aligned} \quad (4.28)$$

**Adamized Bregman Iteration** We shortly remark that one can replace the momentum update in Eq. (4.28) with an Adam update [KB14]. This then yields an Adamized version of linearized Bregman iterations as employed in [BREG-I].

#### 4.4.4. Convergence of Stochastic Bregman Iterations

While various previous works prove convergence of linearized Bregman iterations (see e.g. [Osh+05; COS09]), the stochastic setting requires special treatment. In [BREG-I] the first guarantees for the algorithm in Eq. (4.27) were proven. Other work on convergence of stochastic Bregman iterations, or mirror descent [DEH21; HR21; ZH18; DOr+21; AKL22] requires a differentiable functional  $J$ . However, since our main motivation is to a functional in the flavor of the  $\ell^1$  norm this is not applicable. Therefore, we present the novel convergence analysis of [BREG-I].

**Assumptions on the Gradient Estimator** In order to obtain convergence guarantees, we need to assume mainly two properties on the gradient estimator  $g(\cdot, \cdot)$ . First we assume unbiasedness, which means

$$\mathbb{E}[g(\theta; \omega)] = \nabla \mathcal{L}(\theta) \text{ for all } \theta \in \Theta.$$

The second assumption we need in the following is referred to as *bounded variance* of the estimator.

**Assumption 4.22 (Bounded variance).** There exists a constant  $\sigma > 0$  such that for any  $\theta \in \Theta$  it holds

$$\mathbb{E}[\|g(\theta; \omega) - \nabla \mathcal{L}(\theta)\|^2] \leq \sigma^2. \quad (4.29)$$

**Remark 4.23.** We want to remark, that this property is weaker than the bounded gradient assumption

$$\mathbb{E} [\|g(\theta; \omega)\|^2] \leq C$$

for some constant  $C > 0$ . In fact this condition can not be enforced together with a strong convexity assumption—which we employ in [Theorem 4.30](#)—as shown in [Ngu+18].  $\triangle$

**Assumptions on the Regularizer and on the Loss Function** The assumptions on the regularization functional  $J$  are mild and merely ensure the well-definedness of the proximal mapping,

**Assumption 4.24 (Regularizer).** We assume that  $J : \Theta \rightarrow (-\infty, \infty]$  is a convex, proper, and lower semicontinuous functional on the Hilbert space  $\Theta$ .

Our assumptions on the loss function  $\mathcal{L}$  are more restrictive. We require it to be bounded from below and differentiable, which are both standard assumptions. Additionally, we require Lipschitz continuity of the gradient, which also commonly employed in optimization literature.

**Assumption 4.25 (Loss function).** We assume the following conditions on the loss function:

- The loss function  $\mathcal{L}$  is bounded from below and without loss of generality we assume  $\mathcal{L} \geq 0$ .
- The function  $\mathcal{L}$  is continuously differentiable.
- The gradient of the loss function  $\theta \mapsto \nabla \mathcal{L}(\theta)$  is  $L$ -Lipschitz for  $L \in (0, \infty)$ :

$$\|\nabla \mathcal{L}(\tilde{\theta}) - \nabla \mathcal{L}(\theta)\| \leq L \|\tilde{\theta} - \theta\|, \quad \forall \theta, \tilde{\theta} \in \Theta. \quad (4.30)$$

If the loss function  $\mathcal{L}$  fulfills the previous assumptions we are able to prove loss decay of the iterates in [Theorem 4.29](#). However, in order to show convergence of the iterates we additionally need a convexity assumption. For a differentiable functional  $J$ , the authors in [DEH21] assumed

$$\nu D_J(\bar{\theta}, \theta) \leq D_{\mathcal{L}}(\bar{\theta}, \theta)$$

which for twice differentiable  $J, \mathcal{L}$  transfers to

$$\nu \nabla^2 J \lesssim \nabla^2 \mathcal{L}, \quad \forall \bar{\theta}, \theta \in \Theta.$$

Plugging in the definition of the Bregman dist  $D_{\mathcal{L}}$  we obtain

$$\nu D_J(\bar{\theta}, \theta) \leq \mathcal{L}(\bar{\theta}) - \mathcal{L}(\theta) - \langle \nabla \mathcal{L}(\theta), \bar{\theta} - \theta \rangle.$$

In this form one observes that this is in fact a convexity assumption on  $\mathcal{L}$  in a  $J$  dependent distance, as employed in [[BREG-I](#)].

**Assumption 4.26 (Strong convexity).** For a proper convex function  $H : \Theta \rightarrow \mathbb{R}$  and  $\nu \in (0, \infty)$ , we say that the loss function  $\theta \mapsto \mathcal{L}(\theta)$  is  $\nu$ -strongly convex w.r.t.  $H$ , if

$$\mathcal{L}(\bar{\theta}) \geq \mathcal{L}(\theta) + \langle \nabla \mathcal{L}(\theta), \bar{\theta} - \theta \rangle + \nu D_J^p(\bar{\theta}, \theta), \quad \forall \theta, \bar{\theta} \in \Theta, p \in \partial H(\theta). \quad (4.31)$$

**Remark 4.27.** We have two relevant cases for the choice of  $H$ . For  $H = \frac{1}{2} \|\cdot\|^2$  Assumption 4.26 reduces to standard strong  $\nu$ -convexity. The other relevant case, is  $H = J_\delta$ , i.e. we consider convexity w.r.t. to the functional  $J_\delta$ .  $\triangle$

**Remark 4.28.** In the setting of training a neural network, where we employ the empirical loss Eq. (4.1), this convexity assumption usually fails. While it is possible to enforce this conditions only locally around the minimum, this does not significantly improve the applicability. For future work, it would be desirable to enforce a Kurdyka–Łojasiewicz inequality, as in [Ben+21] for the deterministic case.  $\triangle$

**Loss Decay** The first convergence result considers the loss decay of the iterates. Here, we do not assume convexity of the loss function. Under this assumptions [Ben+21; BB18] were able to show the inequality

$$\begin{aligned} \mathbb{E} [\mathcal{L}(\theta^{(k+1)})] + \frac{1}{\tau^{(k)}} \mathbb{E} [D_J^{\text{sym}}(\theta^{(k+1)}, \theta^{(k)})] + \frac{C}{2\delta\tau^{(k)}} \mathbb{E} [\|\theta^{(k+1)} - \theta^{(k)}\|^2] \\ \leq \mathbb{E} [\mathcal{L}(\theta^{(k)})]. \end{aligned}$$

In our setting, employing a stochastic gradient estimator, we are able to prove a similar estimate. Here, we obtain an additional term scaled  $\sigma$  which controls the expected squared difference between the gradient estimator and the actual gradient. It should however be noted, that the proof is not only a trivial extension.

**Theorem 4.29 ([BREG-I, Th. 2]: Loss decay).** Assume that Assumptions 4.22, 4.24 and 4.25 hold true, let  $\delta > 0$ , and let the step sizes satisfy  $\tau^{(k)} \leq \frac{2}{\delta L}$ . Then there exist constants  $c, C > 0$  such that for every  $k \in \mathbb{N}$  the iterates of (4.27) satisfy

$$\begin{aligned} \mathbb{E} [\mathcal{L}(\theta^{(k+1)})] + \frac{1}{\tau^{(k)}} \mathbb{E} [D_J^{\text{sym}}(\theta^{(k+1)}, \theta^{(k)})] + \frac{C}{2\delta\tau^{(k)}} \mathbb{E} [\|\theta^{(k+1)} - \theta^{(k)}\|^2] \\ \leq \mathbb{E} [\mathcal{L}(\theta^{(k)})] + \tau^{(k)} \delta \frac{\sigma^2}{2c}, \end{aligned} \quad (4.32)$$

**Convergence of the Iterates** Here, we have two cases respectively proving convergence w.r.t. the  $L^2$  distance and the Bregman distance of  $J_\delta$ . The first assumes strong convexity with  $H = \frac{1}{2} \|\cdot\|^2$  in Assumption 4.26.

**Theorem 4.30 ([BREG-I, Th. 6]: Convergence in norm).** Assume that Assumptions 4.22, 4.24 and 4.25 and Assumption 4.26 for  $H = \frac{1}{2} \|\cdot\|^2$  hold true and let  $\delta > 0$ . Furthermore, assume that the step sizes  $\tau^{(k)}$  are such that for all  $k \in \mathbb{N}$ :

$$\tau^{(k)} \leq \frac{\mu}{2\delta L^2}, \quad \tau^{(k+1)} \leq \tau^{(k)}, \quad \sum_{k=0}^{\infty} (\tau^{(k)})^2 < \infty, \quad \sum_{k=0}^{\infty} \tau^{(k)} = \infty.$$

The function  $\mathcal{L}$  has a unique minimizer  $\theta^*$  and if  $J(\theta^*) < \infty$  the stochastic linearized Bregman iterations (4.27) satisfy the following:

- Letting  $d_k := \mathbb{E} [D_{J_\delta}^{v^{(k)}}(\theta^*, \theta^{(k)})]$  it holds

$$d_{k+1} - d_k + \frac{\mu}{4} \tau^{(k)} \mathbb{E} [\|\theta^* - \theta^{(k+1)}\|^2] \leq \frac{\sigma}{2} \left( (\tau^{(k)})^2 + \mathbb{E} [\|\theta^{(k)} - \theta^{(k+1)}\|^2] \right). \quad (4.33)$$

- The iterates possess a subsequence converging in the  $L^2$ -sense of random variables:

$$\lim_{j \rightarrow \infty} \mathbb{E} [\|\theta^* - \theta^{(k_j)}\|^2] = 0. \quad (4.34)$$

Here,  $J_\delta$  is defined as in (4.23).

For the second result we assume convexity w.r.t. the Bregman distance, i.e. we choose  $H = J_\delta$  in Assumption 4.26. This induces a relation between the Bregman distance of  $J$  and the loss function  $\mathcal{L}$ , which has been similarly employed in [DEH21].

**Theorem 4.31 ([BREG-I, Th. 11]: Convergence in the Bregman distance).**

Assume that Assumptions 4.22, 4.24 and 4.25 and Assumption 4.26 for  $H = J_\delta$  hold true and let  $\delta > 0$ . The function  $\mathcal{L}$  has a unique minimizer  $\theta^*$  and if  $J(\theta^*) < \infty$  the stochastic linearized Bregman iterations (4.27) satisfy the following:

- Letting  $d_k := \mathbb{E} [D_{J_\delta}^{v^{(k)}}(\theta^*, \theta^{(k)})]$  it holds

$$d_{k+1} \leq \left[ 1 - \tau^{(k)} \nu \left( 1 - \tau^{(k)} \frac{2\delta^2 L^2}{\nu} \right) \right] d_k + \delta (\tau^{(k)})^2 \sigma^2. \quad (4.35)$$

- For any  $\varepsilon > 0$  there exists  $\tau > 0$  such that if  $\tau^{(k)} = \tau$  for all  $k \in \mathbb{N}$  then

$$\limsup_{k \rightarrow \infty} d_k \leq \varepsilon. \quad (4.36)$$

- If  $\tau^{(k)}$  is such that

$$\lim_{k \rightarrow \infty} \tau^{(k)} = 0 \quad \text{and} \quad \sum_{k=0}^{\infty} \tau^{(k)} = \infty \quad (4.37)$$

then it holds

$$\lim_{k \rightarrow \infty} d_k = 0. \quad (4.38)$$

Here,  $J_\delta$  is defined as in (4.23).

#### 4.4.5. Numerical Results and Practical Considerations

Before briefly reviewing the numerical results in [BREG-I, Sec. 4], we remark on some practical considerations. In particular we comment on the parameter initialization strategy. All the experiments were implemented in Python [VD95] employing—among others—the PyTorch [Pas+19].

**Parameter Initialization** As already noticed in [GB10] parameter initialization has a significant impact on the training of the neural network. Here, in contrast to standard Bregman methods, we are not able to initialize the parameters of the neural network as  $\theta = 0$ . This is due to that fact, that a zero initialization induces symmetries in the network weights, for which one cannot utilize the full expressivity of the architecture [GBC16, Ch. 6]. Therefore, we rather employ the approach from [Liu+21; DZ19; Mar10] of sparsifying weight matrices  $\tilde{W}^l \in \mathbb{R}^{n_{l+1} \times n_l}$  up to a certain level, by a pointwise multiplication with a binary mask  $M^l \in \{0, 1\}^{n_{l+1} \times n_l}$

$$W^l := \tilde{W}^l \odot M^l.$$

Each entry in  $M^l$  is i.i.d. sampled from a Bernoulli distribution

$$M_{i,j}^l \sim \mathcal{B}(r).$$

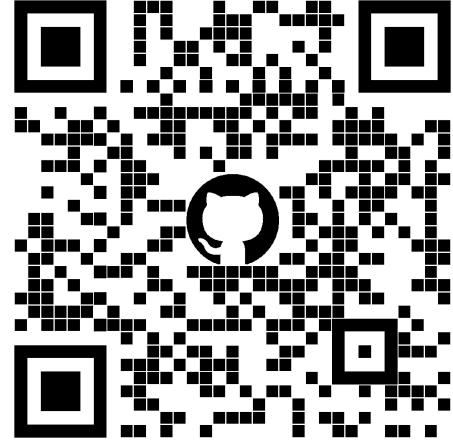
where the parameter  $r$  determines the sparsity,

$$N(W^l) := \frac{\|W^l\|_0}{n_l \cdot n_{l-1}} = 1 - S(W^l)$$

with  $N$  denoting the percentage of used parameters and  $S$  the sparsity. In [GB10] the authors advise to especially control the variance of the parameter initialization distribution, for which in [BREG-I] we derive

$$\text{Var} [\tilde{W}^l] = \frac{1}{r} \text{Var} [\tilde{W}^l \odot M^l] \quad (4.39)$$

and therefore scale the weights with the sparsity parameter  $r$  at initialization.



The code for all the experiments is available at [github.com/TimRoith/BregmanLearning](https://github.com/TimRoith/BregmanLearning).

**Choice of Regularizers** In all our experiments we choose a  $L^1$  type sparsity promoting regularization function functional  $J$ . We do not employ any coupling between weight matrices of different layers, and therefore for  $\theta = ((W^1, b^1), \dots, (W^L, b^L))$  we have

$$J(\theta) = \sum_{l=1}^L J^l(W^l)$$

where  $J^l$  is chosen according to the layer type. In the easiest case of a fully connected layer, we can choose

$$J^l(W^l) := \|W^l\|_1.$$

In the case of a convolutional layer we have that  $W^l$  is determined by convolutional kernels  $K_{i,j} \in \mathbb{R}^{k \times k}$ , see [Section 4.1.1](#). Here, we typically employ a group sparsity term in the form

$$J^l(W^l) = \|W^l\|_{2,1} = \sum_{i,j} \|K_{i,j}\|_2.$$

The outer sum acts as a  $L^1$  regularizer on the instances  $\|K_{i,j}\|$ . Sparsity in this sense, then amounts to having indices  $(i, j)$  for which  $\|K_{i,j}\|_2 = 0 \Leftrightarrow K_{i,j} = 0$ , i.e., we prune away whole convolutional filters. This effect is displayed in [[BREG-I](#), Fig. 1]. We can also employ group sparsity on fully connected layers, by considering row sparsity of  $W^l \in \mathbb{R}^{n_{l+1}, n_l}$

$$J^l(W^l) = \sum_{i=1}^{n_{l+1}} \|W_{i,:}\|_2 = \sum_{i=1}^{n_{l+1}} \sqrt{\sum_{j=1}^{n_l} W_{i,j}^2}.$$

In this setting we have a  $L^1$  penalty on the row norms  $\|W_{i,:}\|_2$  which therefore enforces whole rows to be zero. This is relevant, if we employ a layer architecture with  $\Psi^l(0) = 0$ , e.g., using no bias vectors and the ReLU activation function. In this setting, if the  $i$ th row of  $W^l$  is zero this effectively means, that the  $i$ th neuron in layer  $l+1$  is inactive. This observation allows the neural architecture search in one of the following paragraphs.

**Comments on the Numerical Results** We briefly remark on the numerical results as displayed in [[BREG-I](#), Sec. 4]. In the experiments we employed feed-forward networks with simple linear, convolutional and residual layers and tested on the three datasets [[Kri09](#); [XRV17](#); [LC10](#)].

The basic comparison between the algorithms SGD, ProxGD and LinBreg shows the qualitative behaviour of each iterations. Infantilizing sparse does not have any effect on SGD, since it does not preserve the sparsity in any way. ProxGD rather starts with many active parameters and reduces the this number during the iteration. Only the discretization of the inverse scale space flow—via Bregman iterations—shows the desired behaviour of gradually adding active parameters. Furthermore, in [[BREG-I](#),

Fig. 2] we can see, that the choice of  $\lambda$  in the regularizer  $J = \lambda \|\cdot\|_1$  changes the results significantly. In the light of Eq. (4.22) this is not expected for the standard Bregman iterations with a convex loss. It is therefore interesting to see, that in our non-convex and stochastic situation this effect changes.

The momentum variants, as discussed in Section 4.4.3 yield the desired effect of enhancing the validation accuracy, and respectively converging faster. However in each of the experiments, one can also observe that adding a momentum term has the effect that more parameters are added faster. On the one hand this could mean that the network actually requires more parameters to have a higher accuracy, for which a momentum variant is more likely to increase the number of needed parameters. However, the quantitative evaluation on the CIFAR10 dataset [Kri09] shows, that especially the Adamized version tends to increase the number of used parameters rather aggressively, while only slightly increasing the performance of the net. The performance here is very similar to the one of proximal gradient descent. However, the training of a residual network seems to be slightly better with a standard Lasso implementation. Here, one neglects the non-differentiability of the  $L^1$  norm and computes a derivate via automatic differentiation [Ral81; MDA15] (we employed the `autograd` library of the PyTorch package [Pas+19]). In order to obtain true zeros in the weight matrix one then has to employ a thresholding operation after the training. In some sense this method is not a proper sparse training approach, but rather a regularization method with an added pruning step at the end.

**Comments on Efficiency** One of the major advantages of the Bregman approach, is that the network is sparse already during the training time. As with all sparse-to-sparse training approaches this yields a very small number of active parameters over all training step. This sparsity can be easily exploited in each forward pass. However, it is not directly possible to achieve performance gains during the backward pass of the network, since in general

$$W_{ij}^l = 0 \not\Rightarrow \partial_{W_{ij}^l} \mathcal{L}(\theta) = 0.$$

In [BREG-I; BREG-II] there are no evaluations on the training time and memory consumption of the Bregman algorithm. Since the complexity does not increase in comparison to the standard SGD, one hopes to obtain a faster training time here. This is an interesting open question for future work. It should however be remarked that the computational complexity of the LinBreg algorithm does not increase significantly, compared to SGD, since the evaluation of the proximal operator is very efficient for  $L^1$  type functionals.

**Neural Architecture Search** An interesting aspect hinted in [BREG-I] is optimizing the architecture of a neural network via sparsity. In the example provided in [BREG-I, Fig. 4] one defines a super-architecture as multi-layer perceptron with an equal number of neurons in each layer. Training this task with the LinBreg algorithm reveals the well-known autoencoder structure [HZ93]. This idea was developed further in [BREG-II],

where also skip connection of a residual architecture were learned. Here, the super-architecture was given by a dense net [Hua+17], where each skip connection was scaled by a parameter, which was penalized with a sparsity term. The driving question here, if it is possible also learn an architecture similar to the U-Net as proposed in [RFB15].