

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	vi
I. Exposition	1
1. Introduction	2
2. Learning Paradigms	5
2.1. Unsupervised Learning	6
2.2. Supervised Learning	6
2.3. Semi-Supervised Learning	7
II. Prints	12

List of Figures

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.

Grafik
anpassen

Part I: Exposition	??: ??
Chapter 2: Learning Paradigms	—
??: ??	????
??: ??	????

Part I consists of three chapters, of which the first explains the paradigms, *unsupervised*, *semi-supervised* and *supervised* learning. The other chapters are the split up thematically, concerning the topics semi-supervised and supervised learning respectively. In each of these chapters a short introduction provides the necessary framework allowing us to explain the main contributions. The following publications are reprinted in **??**:

[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). DOI: [10.1093/imanum/drac048](https://doi.org/10.1093/imanum/drac048).

[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.

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 master 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 ???. TR adapted the continuum limit framework to the L^∞ 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. He contributed to both the numeric and the analysis 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 and RR he conducted the numerical examples and also wrote most of the source code. Furthermore, he wrote large parts of the paper.

[BREG-I]: TR expanded LB's ideas of employing Bregman iteration for sparse training. 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. 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 machine 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 of input-output pairs, i.e. is fully labeled—and the 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 that 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 and 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 for transformed input \bar{x} we want to obtain 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 heuristic 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 is 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 graph Laplacian tend to be constant over the whole dataset, whenever the dimension is bigger than 2. This issue is prevalent in the infinite data-limit, where a priori we consider the case, when the amount of unlabelled 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 are the computational resources required in both scenarios. 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 CO₂ emissions [Hoe+21]. For graph-based semi-supervised learning one first needs to compute distances between many different data points, to compute 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 scale 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 subject 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 [lipschitz1877lehrbuch] of a neural network, which trains an adversarially robust net. In the multi-resolution setting we analyze the role of Fourier neural operators as proposed by [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 [osher2005iterative] that allows

Chapter 1. Introduction

for sparse weight matrices throughout the training. We also provide the convergence analysis for the stochastic adaption of the original iterations.

Structure of The Exposition

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 task of *learning* a function $f : \tilde{\mathcal{X}} \rightarrow \mathcal{Y}$ from the given data, where the two most important cases for us are

- **classification:** f assigns a label to each $x \in \tilde{\mathcal{X}}$ out of a total 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.

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

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

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

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, $\arg\max_{i=1, \dots, C} f_\theta(x)_i$.

- **image denoising:** 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 set $\tilde{\mathcal{X}} \subset \mathbb{R}^d$ is usually either the set of data points \mathcal{X}_n or the whole space \mathcal{X} . The learning paradigms we consider in this thesis, differ by their usage of labeled data. We review the concepts in the following.

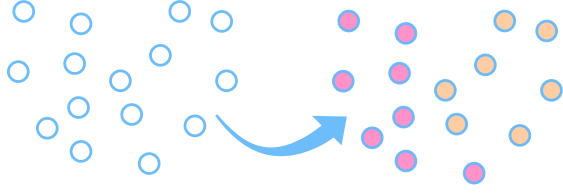
2.1. Unsupervised Learning

In this case 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 ??.



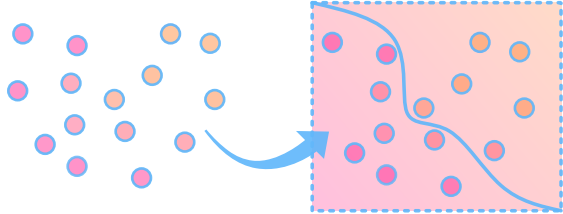
2.2. Supervised Learning

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 then to infer a function defined on the underlying space, i.e. $f : \mathcal{X} \rightarrow \mathcal{Y}$, i.e. we want to assign a label to unseen $x \in \mathcal{X}$ 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 i.i.d. w.r.t. $P_{\mathcal{X}, \mathcal{Y}}$. In this interpretation, a neural network can aim to model the conditional $P(y|x)$ 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 \mathcal{U} , where typically each element can be describe by a finite number of parameters. Among others, common methods or parametrizations include

- Support vector machines [CV95; SS05],
- decision Trees [MS63; Bre+84],

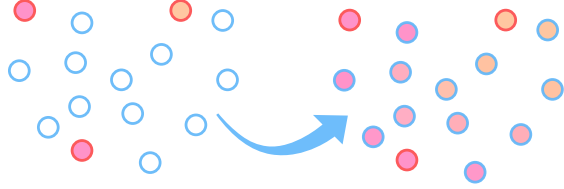


- neural networks [Tur04; Ros58; MP69].

In ?? we exclusively focus on supervised learning algorithms employing neural networks. We refer to [Sch15] for an exhaustive historical overview. The concrete setting and learning framework is given in ??.

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 set of labeled data. Typically the labeled data 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 afterwards via some thresholding, e.g.,

$$x \text{ belongs to the class} \Leftrightarrow f(x) > 0.5$$

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 yields a function $f_c(o)$. The final output can then either obtained by taking the argmax, 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 ?? we focus on graph-based learning algorithms, however we refer to [Zhu05] for a overview of semi-supervised learning algorithms.

Bibliography

Books

- [Ros+62] F. Rosenblatt et al. *Principles of neurodynamics: Perceptrons and the theory of brain mechanisms*. Vol. 55. Spartan books Washington, DC, 1962.
- [GV13] G. H. Golub and C. F. Van Loan. *Matrix computations*. JHU press, 2013.
- [Bre+84] L. Breiman, J. Friedman, R. Olshen, and C. Stone. *Classification and regression trees*. CRC Press, New York, 1984.
- [Zhu05] X. Zhu. *Semi-supervised learning with graphs*. Carnegie Mellon University, 2005.

Articles

- [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). DOI: [10.1093/imanum/drac048](https://doi.org/10.1093/imanum/drac048).
- [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.
- [Sam59] A. L. Samuel. “Some studies in machine learning using the game of checkers.” In: *IBM Journal of research and development* 3.3 (1959), pp. 210–229.
- [Ros58] F. Rosenblatt. “The perceptron: a probabilistic model for information storage and organization in the brain.” In: *Psychological review* 65.6 (1958), p. 386.
- [Kel60] H. J. Kelley. “Gradient theory of optimal flight paths.” In: *Ars Journal* 30.10 (1960), pp. 947–954.
- [RHW86] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. “Learning representations by back-propagating errors.” In: *nature* 323.6088 (1986), pp. 533–536.
- [Sch22] J. Schmidhuber. “Annotated history of modern AI and Deep learning.” In: *arXiv preprint arXiv:2212.11279* (2022).

- [Cha+21] J. Chai, H. Zeng, A. Li, and E. W. Ngai. “Deep learning in computer vision: A critical review of emerging techniques and application scenarios.” In: *Machine Learning with Applications* 6 (2021), p. 100134.
- [Khu+23] D. Khurana, A. Koli, K. Khatter, and S. Singh. “Natural language processing: State of the art, current trends and challenges.” In: *Multimedia tools and applications* 82.3 (2023), pp. 3713–3744.
- [She+22] M. Shehab, L. Abualigah, Q. Shambour, M. A. Abu-Hashem, M. K. Y. Shambour, A. I. Alsalibi, and A. H. Gandomi. “Machine learning in medical applications: A review of state-of-the-art methods.” In: *Computers in Biology and Medicine* 145 (2022), p. 105458.
- [GSS14] I. J. Goodfellow, J. Shlens, and C. Szegedy. “Explaining and harnessing adversarial examples.” In: *arXiv preprint arXiv:1412.6572* (2014).
- [NSZ09] B. Nadler, N. Srebro, and X. Zhou. “Statistical analysis of semi-supervised learning: The limit of infinite unlabelled data.” In: *Advances in neural information processing systems* 22 (2009).
- [Hoe+21] T. Hoefer, D. Alistarh, T. Ben-Nun, N. Dryden, and A. Peste. “Sparsity in Deep Learning: Pruning and growth for efficient inference and training in neural networks.” In: *J. Mach. Learn. Res.* 22.241 (2021), pp. 1–124.
- [Lan52] C. Lanczos. “Solution of systems of linear equations by minimized iterations.” In: *J. Res. Nat. Bur. Standards* 49.1 (1952), pp. 33–53.
- [Li+20] Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, and A. Anandkumar. “Fourier neural operator for parametric partial differential equations.” In: *arXiv preprint arXiv:2010.08895* (2020).
- [Fuk80] K. Fukushima. “Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position.” In: *Biological cybernetics* 36.4 (1980), pp. 193–202.
- [Bol68] L. Boltzmann. “Studien über das Gleichgewicht der lebenden Kraft.” In: *Wissenschaftliche Abhandlungen* 1 (1868), pp. 49–96.
- [ST14] A. Subramanya and P. P. Talukdar. “Graph-based semi-supervised learning.” In: *Synthesis Lectures on Artificial Intelligence and Machine Learning* 8.4 (2014), pp. 1–125.
- [Ste+56] H. Steinhaus et al. “Sur la division des corps matériels en parties.” In: *Bull. Acad. Polon. Sci* 1.804 (1956), p. 801.
- [DLR77] A. P. Dempster, N. M. Laird, and D. B. Rubin. “Maximum likelihood from incomplete data via the EM algorithm.” In: *Journal of the royal statistical society: series B (methodological)* 39.1 (1977), pp. 1–22.
- [GS15] N. García Trillos and D. Slepčev. “Continuum Limit of Total Variation on Point Clouds.” In: *Archive for Rational Mechanics and Analysis* 220.1 (2015), pp. 193–241. DOI: [10.1007/s00205-015-0929-z](https://doi.org/10.1007/s00205-015-0929-z).

- [SB09] A. Szlam and X. Bresson. “A total variation-based graph clustering algorithm for cheeger ratio cuts.” In: *UCLA Cam report* (2009), pp. 09–68.
- [Gar+16] N. García Trillos, D. Slepčev, J. Von Brecht, T. Laurent, and X. Bresson. “Consistency of Cheeger and ratio graph cuts.” In: *The Journal of Machine Learning Research* 17.1 (2016), pp. 6268–6313.
- [GMT22] N. García Trillos, R. Murray, and M. Thorpe. “From graph cuts to isoperimetric inequalities: Convergence rates of Cheeger cuts on data clouds.” In: *Archive for Rational Mechanics and Analysis* 244.3 (2022), pp. 541–598.
- [GS18] N. García Trillos and D. Slepčev. “A variational approach to the consistency of spectral clustering.” In: *Applied and Computational Harmonic Analysis* 45.2 (2018), pp. 239–281.
- [THH21] N. G. Trillos, F. Hoffmann, and B. Hosseini. “Geometric structure of graph Laplacian embeddings.” In: *The Journal of Machine Learning Research* 22.1 (2021), pp. 2934–2988.
- [Hof+22] F. Hoffmann, B. Hosseini, A. A. Oberai, and A. M. Stuart. “Spectral analysis of weighted Laplacians arising in data clustering.” In: *Applied and Computational Harmonic Analysis* 56 (2022), pp. 189–249.
- [CV95] C. Cortes and V. Vapnik. “Support-vector networks.” In: *Machine learning* 20 (1995), pp. 273–297.
- [MS63] J. N. Morgan and J. A. Sonquist. “Problems in the analysis of survey data, and a proposal.” In: *Journal of the American statistical association* 58.302 (1963), pp. 415–434.
- [MP69] M. Minsky and S. Papert. “An introduction to computational geometry.” In: *Cambridge tiass., HIT* 479 (1969), p. 480.
- [Sch15] J. Schmidhuber. “Deep learning in neural networks: An overview.” In: *Neural Networks* 61 (2015), pp. 85–117. DOI: <https://doi.org/10.1016/j.neunet.2014.09.003>.

Theses

- [Roi21] T. Roith. “Master thesis: Continuum limit of Lipschitz learning on graphs.” MA thesis. Friedrich-Alexander-Universität Erlangen-Nürnberg, 2021.

Part II.

Prints