

# Ergodic bilevel optimization

Christoph Brauer  
Dirk A. Lorenz

Institute of Analysis and Algebra, TU Braunschweig  
38092 Braunschweig, Germany  
Email: {ch.brauer, d.lorenz}@tu-braunschweig.de

**Abstract**—In image and signal processing and beyond, quantities of interest are often reconstructed from noisy measurements by means of suitable convex optimization problems. While model based approaches usually assume that the ingredients of the problem are a priori known, data driven approaches are motivated by situations where the objective function or constraints of the problem are partially unknown and shall be learned from data. This gives rise to bilevel optimization problems in which the original convex problem (hereafter referred to as the lower-level problem) appears as a constraint. Applying gradient based algorithms to bilevel optimization problems poses the difficulty to differentiate through the solution operator of the lower-level problem. In this contribution, we consider the approach to unroll a fixed number of update steps of the Chambolle-Pock algorithm applied to the lower-level problem in order to accomplish this kind of differentiation approximately. We investigate the asymptotic behavior of the resulting gradients and conclude that unrolling ergodic averages instead of ordinary iterates can have a positive effect upon the learning dynamics.

## I. INTRODUCTION

We consider the task to recover a ground truth  $\mathbf{x}^\dagger \in \mathbb{R}^n$  from degraded measurements  $\tilde{\mathbf{x}} \in \mathbb{R}^n$ . To that end, we assume that there exists an *unknown* analysis operator  $\mathbf{K} \in \mathbb{R}^{m \times n}$  such that with

$$S(\mathbf{K}, \tilde{\mathbf{x}}) := \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{K}\mathbf{x}) + G(\mathbf{x} - \tilde{\mathbf{x}}) \quad (1)$$

$\hat{\mathbf{x}} \in S(\mathbf{K}, \tilde{\mathbf{x}})$  is a good approximation of  $\mathbf{x}^\dagger$ , where  $F$  and  $G$  are a priori chosen proper, convex and l.s.c. functions. In a supervised learning environment, we have training data  $\mathcal{S}$  available which comprises a certain number of pairs  $(\tilde{\mathbf{x}}_t, \mathbf{x}_t^\dagger) \in \mathbb{R}^n \times \mathbb{R}^n$ . By means of a smooth loss function  $L : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ , we measure the discrepancy  $L(\hat{\mathbf{x}}, \mathbf{x}^\dagger)$  between a ground truth and the associated reconstruction. Therewith, our goal is to choose  $\mathbf{K}$  as a solution of the following bilevel optimization problem [1]:

$$\min_{\mathbf{K} \in \mathbb{R}^{m \times n}} \sum_{(\tilde{\mathbf{x}}_t, \mathbf{x}_t^\dagger) \in \mathcal{S}} L(\hat{\mathbf{x}}_t, \mathbf{x}_t^\dagger) \quad \text{s.t. } \forall t : \hat{\mathbf{x}}_t \in S(\mathbf{K}, \tilde{\mathbf{x}}_t) \quad (2)$$

## II. PLAIN UNROLLING SCHEME

After the change of variables  $\mathbf{r} := \mathbf{x} - \tilde{\mathbf{x}}$  the Chambolle-Pock iteration [2] for the lower-level problem (1) can be written as

$$\mathbf{z}_D^k = \mathbf{y}^{k-1} + \sigma \mathbf{K}(\tilde{\mathbf{x}} + \bar{\mathbf{r}}^{k-1}) \quad \mathbf{y}^k = \operatorname{prox}_{\sigma F^*}(\mathbf{z}_D^k) \quad (3)$$

$$\mathbf{z}_P^k = \mathbf{r}^{k-1} - \tau \mathbf{K}^\top \mathbf{y}^k \quad \mathbf{r}^k = \operatorname{prox}_{\tau G}(\mathbf{z}_P^k) \quad (4)$$

$$\bar{\mathbf{r}}^k = \mathbf{r}^k + \theta(\mathbf{r}^k - \mathbf{r}^{k-1}) \quad (5)$$

where we use  $\mathbf{r}^0 = \bar{\mathbf{r}}^0 = \mathbf{0}$  and  $\mathbf{y}^0 = \mathbf{0}$ , and  $\sigma, \tau$  and  $\theta$  are the usual step size and extrapolation parameters of the Chambolle-Pock algorithm. Additionally, we assume that both proximal operators are at least semismooth. Now, the original constraint  $\hat{\mathbf{x}} \in S(\mathbf{K}, \tilde{\mathbf{x}})$  is replaced by a relaxed version  $\hat{\mathbf{x}} = \tilde{\mathbf{x}} + \mathbf{r}^d(\mathbf{K}, \tilde{\mathbf{x}})$  which is computed according to (3)–(5) with some fixed  $d \in \mathbb{N}$ . The relaxed version of the bilevel problem then reads

$$\min_{\mathbf{K} \in \mathbb{R}^{m \times n}} \sum_{(\tilde{\mathbf{x}}_t, \mathbf{x}_t^\dagger) \in \mathcal{S}} L(\tilde{\mathbf{x}}_t + \mathbf{r}_t^d(\mathbf{K}, \tilde{\mathbf{x}}_t), \mathbf{x}_t^\dagger) \quad (6)$$

and because  $L$  is differentiable and the proximal operators are semismooth, the derivative of (6) can be computed explicitly. To that end, it is helpful to interpret the intermediate iterates  $\mathbf{y}^k$  and  $\mathbf{r}^k$  as activations of a recurrent neural network  $\mathbf{r}^d$  with nonlinearities  $\operatorname{prox}_{\sigma F^*}$  and  $\operatorname{prox}_{\tau G}$  and linear operator  $\mathbf{K}$ . Given that both proximal operators and the respective derivatives can be evaluated element-wise, one can use a backpropagation scheme [3] or an adjoint state method to derive the results presented in the next section.

## III. THEORETICAL RESULTS

**Lemma 1.** *Let  $L(\hat{\mathbf{x}}, \mathbf{x}^\dagger) = \frac{1}{2} \|\hat{\mathbf{x}} - \mathbf{x}^\dagger\|_2^2$ ,  $\delta_D^{d+1} = \bar{\delta}_D^{d+1} = \mathbf{0}$  and  $\delta_P^{d+1} = \tilde{\mathbf{x}} + \mathbf{r}^d(\mathbf{K}, \tilde{\mathbf{x}}) - \mathbf{r}^\dagger$ . Moreover, consider the reverse iteration*

$$\delta_P^k = \operatorname{prox}'_{\tau G}(\mathbf{z}_P^k) \odot (\delta_P^{k+1} + \sigma \mathbf{K}^\top \bar{\delta}_D^{k+1}) \quad (7)$$

$$\delta_D^k = \operatorname{prox}'_{\sigma F^*}(\mathbf{z}_D^k) \odot (\delta_D^{k+1} - \tau \mathbf{K} \delta_P^k) \quad (8)$$

$$\bar{\delta}_D^k = \delta_D^k + \theta(\delta_D^k - \delta_D^{k+1}) \quad (9)$$

for  $k = d, \dots, 1$ . Then, the gradient of  $L(\tilde{\mathbf{x}} + \mathbf{r}^d(\mathbf{K}, \tilde{\mathbf{x}}), \mathbf{x}^\dagger)$  with respect to  $\mathbf{K}$  is given by

$$\sum_{k=1}^d \sigma \delta_D^k (\tilde{\mathbf{x}} + \bar{\mathbf{r}}^{k-1})^\top - \tau \mathbf{y}^{k-1} (\delta_P^k)^\top. \quad (10)$$

Lemma 1 shows that the sequence generating  $\delta_P^k$  and  $\delta_D^k$  is almost another instance of Chambolle-Pock, except that the proximal operators are replaced by element-wise multiplication with the respective derivatives. Assuming that these remain constant after a certain number of iterations allows one to show convergence and draw conclusions about limits via the underlying optimization problems.

**Theorem 2.** *Let  $\mathbf{K}$ ,  $\sigma$  and  $\tau$  satisfy  $\sigma \tau \|\mathbf{K}\|^2 < 1$  and let  $\theta = 1$ . Moreover, suppose that there exists a  $k_0$  such that  $\operatorname{prox}'_{\tau G}(\mathbf{z}_P^k) = \operatorname{prox}'_{\tau G}(\mathbf{z}_P^{k_0}) \in \{0, 1\}^n$  and  $\operatorname{prox}'_{\sigma F^*}(\mathbf{z}_D^k) = \operatorname{prox}'_{\sigma F^*}(\mathbf{z}_D^{k_0}) \in \{0, 1\}^m$  is true for all  $k \geq k_0$ . Then, it holds that*

$$\lim_{d \rightarrow \infty} \delta_P^{k_0} \in \ker(\mathbf{K}) \quad \text{and} \quad \lim_{d \rightarrow \infty} \delta_D^{k_0} \in \ker(\mathbf{K}^\top).$$

## IV. MOTIVATION OF ERGODIC UNROLLING SCHEMES

Theorem 2 supports that  $\delta_P^k$  and  $\delta_D^k$  can vanish, especially in case  $d$  is large enough to ensure convergence of (3)–(5). However, this behavior may be undesirable, e.g. if we want to learn an operator  $\mathbf{K}$  that grants fast progress in early iterations of the forward scheme, or because the minimization of (6) can get stuck in bad local minima. This motivates us to adapt (6), this time using weighted ergodic averages  $e^d = \sum_{k=1}^d \alpha_k \mathbf{r}^k$  instead of only the final iterate. The effect on (7)–(9) can be considered a summation of several backpropagation sequences starting at each iterate  $\mathbf{r}^k$  with  $k \leq d$ . This modification theoretically prevents from the behavior described in Theorem 2, and our experiments show that it can also pay off in terms of substantially lower loss values on the training data.

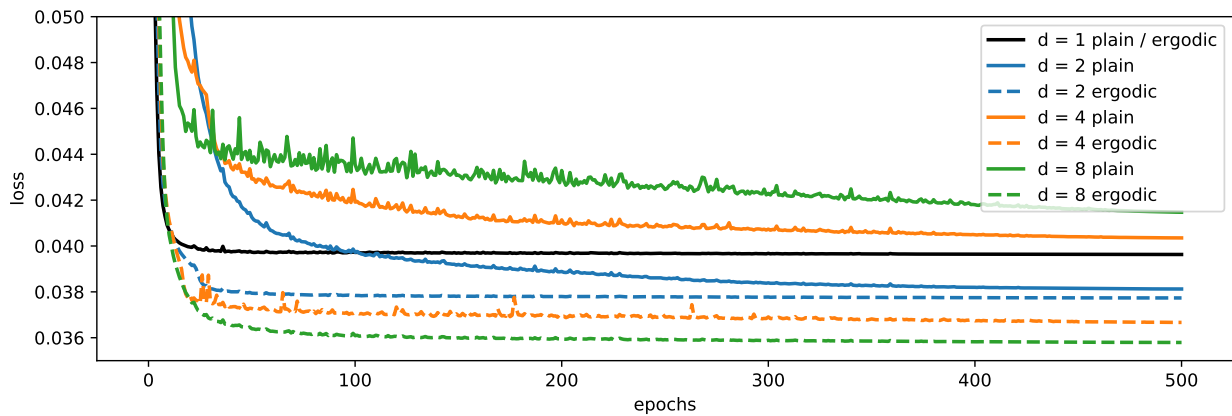


Fig. 1. We take up the experimental setting introduced in [4], [5], [6] in the context of speech dequantization, where  $F$  is the  $\ell_1$ -norm,  $G$  is an indicator function encoding the constraint  $\|\mathbf{x} - \bar{\mathbf{x}}\|_\infty \leq \eta$ , and the training data consists of 66.628 ground truth speech signal snippets from the IEEE corpus [7] and associated quantized versions. Both proximal operators are then projections onto sup norm balls satisfying the assumptions for Lemma 1 and Theorem 2. Throughout, we use stochastic gradient descent with decreasing step sizes to minimize the loss (6). For ergodic averaging we use weights  $\alpha_k = \mathcal{O}(1/k)$ . In case  $d = 1$ , we set  $\alpha_1 = 1$  such that the plain and the ergodic setting produce identical results. In the remaining cases with  $d \in \{2, 4, 8\}$  it can be observed that an increasing number of unrolled iterations leads to increasing loss values in the plain setting, whereas ergodic averaging leads to decreasing loss values induced by the learned analysis operators.

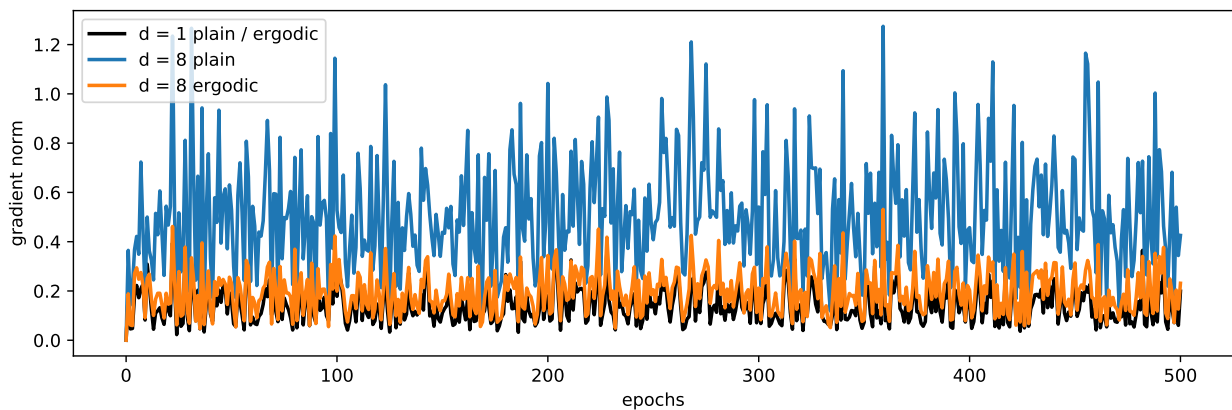


Fig. 2. In the same experimental setting as in Figure 1, it can be observed that the norms of the gradients (10) seem to be less dependent on the number of unrolled iterations in case one uses ergodic averaging. This behavior supports a more stable behavior of the training procedure and more robustness with respect to the utilized step sizes for stochastic gradient descent.

## REFERENCES

- [1] B. Colson, P. Marcotte and G. Savard, *An overview of bilevel optimization*, Annals of Operations Research 153(1), pp. 235–256, 2007.
- [2] A. Chambolle and T. Pock, *A First-Order Primal-Dual Algorithm for Convex Problems with Applications to Imaging*, Journal of Mathematical Imaging and Vision 40(1), pp. 120–145, 2011.
- [3] D. E. Rumelhart, G. E. Hinton and R. J. Williams, *Learning internal representations by error propagation*, California Univ San Diego La Jolla Inst for Cognitive Science, 1985.
- [4] C. Brauer, T. Gerkmann and D. Lorenz, *Sparse reconstruction of quantized speech signals*, In: 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5940–5944, 2016.
- [5] C. Brauer and D. Lorenz, *Primal-dual residual networks*, arXiv preprint arXiv:1806.05823, 2018.
- [6] C. Brauer, Z. Zhao, D. Lorenz and T. Fingscheidt, *Learning to dequantize speech signals by primal-dual networks: An approach for acoustic sensor networks*, To appear in: 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2019.
- [7] P. C. Loizou, *Speech Enhancement - Theory and Practice*, CRC Press, Taylor & Francis Group, 2007.
- [8] J. Adler and O. Öktem, *Learned Primal-dual Reconstruction*, IEEE Transactions on Medical Imaging 37(6), pp. 1322–1332, 2018.
- [9] A. Chambolle and T. Pock, *On the ergodic convergence rates of a first-order primal-dual algorithm*, Mathematical Programming 150(1-2), pp. 253–287, 2016.
- [10] Y. Chen, T. Pock and H. Bischof, *Learning  $\ell_1$ -based analysis and synthesis sparsity priors using bi-level optimization*, arXiv preprint arXiv:1401.4105, 2014.
- [11] K. Gregor and Y. LeCun, *Learning Fast Approximations of Sparse Coding*, In: Proceedings of the 27th International Conference on Machine Learning (ICML), pp. 399–406, 2010.
- [12] K. Kunisch and T. Pock, *A bilevel optimization approach for parameter learning in variational models*, SIAM Journal on Imaging Sciences 6(2), pp. 938–983, 2013.
- [13] P. Ochs, R. Ranftl, R. Brox and T. Pock, *Bilevel Optimization with Nonsmooth Lower Level Problems*, In: International Conference on Scale Space and Variational Methods in Computer Vision (SSVM), pp. 654–665, 2015.
- [14] P. Ochs, R. Ranftl, R. Brox and T. Pock, *Techniques for Gradient-Based Bilevel Optimization with Non-smooth Lower Level Problems*, Journal of Mathematical Imaging and Vision 56(2), pp. 175–194, 2016.
- [15] G. Riegler, D. Ferstl, M. Rüter and H. Bischof, *A Deep Primal-Dual Network for Guided Depth Super-Resolution*, In: Proceedings of the British Machine Vision Conference (BMVC), pp. 7.1–7.14, 2016.
- [16] G. Riegler, M. Rüter and H. Bischof, *ATGV-Net: Accurate Depth Super-Resolution*, In: European Conference on Computer Vision (ECCV), pp. 268–284, 2016.
- [17] S. Wang, S. Fidler and R. Urtasun, *Proximal Deep Structured Models*, In: Advances in Neural Information Processing Systems (NIPS), pp. 865–873, 2016.