



Technische
Universität
Braunschweig



Asymptotic analysis and truncated backpropagation for the unrolled primal-dual algorithm

Christoph Brauer and Dirk Lorenz, September 8, 2023

Introduction

Task

Recover ground truth $\mathbf{y} \in \mathbb{R}^n$
from noisy observation $\mathbf{x} \in \mathbb{R}^n$

via

Convex Problem

$$\hat{\mathbf{y}} \in S(\mathbf{K}, \mathbf{x}) := \underset{\mathbf{y}}{\operatorname{argmin}} F(\mathbf{K}\mathbf{y}) + G(\mathbf{y} - \mathbf{x})$$

use $\hat{\mathbf{K}}$

Bilevel Problem

$$\begin{aligned} \hat{\mathbf{K}} \in \underset{\mathbf{K}}{\operatorname{argmin}} & \sum_{i=1}^m \ell(\mathbf{y}_i, \hat{\mathbf{y}}_i) \\ \text{s.t.} & \forall i : \hat{\mathbf{y}}_i \in S(\mathbf{K}, \mathbf{x}_i) \end{aligned}$$

via

Training Data

$$\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_m, \mathbf{y}_m)\}$$

learn \mathbf{K}
from



Introduction

Bilevel Problem

$$\begin{aligned} \hat{\mathbf{K}} \in \operatorname{argmin}_{\mathbf{K}} & \quad \sum_{i=1}^m \ell(\mathbf{y}_i, \hat{\mathbf{y}}_i) \\ \text{s.t.} & \quad \forall i : \hat{\mathbf{y}}_i \in S(\mathbf{K}, \mathbf{x}_i) \end{aligned}$$

Convex Optimization Algorithm

$$A^1(\mathbf{K}, \mathbf{x}), \dots, A^L(\mathbf{K}, \mathbf{x})$$

Approximate Bilevel Problem

$$\hat{\mathbf{K}} \in \operatorname{argmin}_{\mathbf{K}} \sum_{i=1}^m \ell(\mathbf{y}_i, A^L(\mathbf{K}, \mathbf{x}_i))$$

Introduction

Approximate Bilevel Problem

$$\hat{\mathbf{K}} \in \operatorname{argmin}_{\mathbf{K}} \sum_{i=1}^m \ell(\mathbf{y}_i, A^L(\mathbf{K}, \mathbf{x}_i))$$

single
example

Approximate Gradient

$$\nabla_{\mathbf{K}} \ell(\mathbf{y}, A^L(\mathbf{K}, \mathbf{x}))$$

behavior
for $L \rightarrow \infty$

?

Outline

1. The algorithm A
2. Gradient backpropagation
3. Gradient of $\ell(\mathbf{y}, A^L(\mathbf{K}, \mathbf{x}))$ w.r.t. parameters \mathbf{K}
4. Limit of parameter gradient for $L \rightarrow \infty$
5. Speech dequantization
6. Truncated Backpropagation
7. Interpretability

The algorithm A

Algorithm 1 Chambolle-Pock

Choose $\sigma, \tau > 0$ and $\theta \in [0, 1]$
 Initialize $\mathbf{y}^0 = \bar{\mathbf{y}}^0 = \mathbf{x}$ and $\boldsymbol{\psi}^0 = \mathbf{0}$
for $l = 0, \dots, L - 1$ **do**

$$\mathbf{z}_D^{l+1} = \boldsymbol{\psi}^l + \sigma \mathbf{K} \bar{\mathbf{y}}^l$$

(Dual update)

$$\boldsymbol{\psi}^{l+1} = \text{prox}_{\sigma F^*}(\mathbf{z}_D^{l+1})$$

$$\mathbf{z}_P^{l+1} = \mathbf{y}^l - \tau \mathbf{K}^\top \boldsymbol{\psi}^{l+1}$$

(Prim. update)

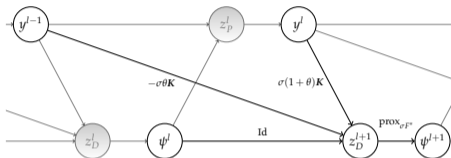
$$\mathbf{y}^{l+1} = \text{prox}_{\tau G}(\mathbf{z}_P^{l+1} - \mathbf{x}) + \mathbf{x}$$

$$\bar{\mathbf{y}}^{l+1} = \mathbf{y}^{l+1} + \theta(\mathbf{y}^{l+1} - \mathbf{y}^l)$$

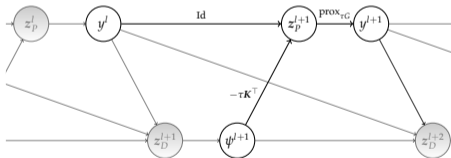
(Extrapolation)

end for

Dual update



Primal update



→ Apply standard backprop as for FCNNs

Gradient backpropagation

Algorithm 2 Backpropagated gradients

Choose σ, τ and θ as in Algorithm 1

Adopt $\mathbf{y}^L, \mathbf{z}_P^1, \dots, \mathbf{z}_P^L, \mathbf{z}_D^1, \dots, \mathbf{z}_D^L$ from Algorithm 1

Initialize $\delta_P^{L+1} = \nabla_{\mathbf{y}} \ell(\mathbf{y}, \mathbf{y}^L)$

Initialize $\delta_D^{L+1} = \bar{\delta}_D^{L+1} = \mathbf{0}$

for $l = L, \dots, 1$ **do**

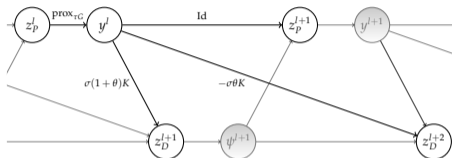
$$\delta_P^l = \mathcal{J}_{\text{prox}_{\tau G}}(\mathbf{z}_P^l - \mathbf{x})^\top (\delta_P^{l+1} + \sigma \mathbf{K}^\top \bar{\delta}_D^{l+1}) \quad (1)$$

$$\delta_D^l = \mathcal{J}_{\text{prox}_{\sigma F^*}}(\mathbf{z}_D^l)^\top (\delta_D^{l+1} - \tau \mathbf{K} \delta_P^l) \quad (2)$$

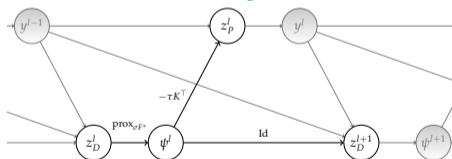
$$\bar{\delta}_D^l = \delta_D^l + \theta (\delta_D^l - \delta_D^{l+1}) \quad (3)$$

end for

Primal gradient $\delta_P^l := \nabla_{\mathbf{z}_P^l} \ell(\mathbf{y}, \mathbf{y}^L)$



Dual gradient $\delta_D^l := \nabla_{\mathbf{z}_D^l} \ell(\mathbf{y}, \mathbf{y}^L)$



→ Basically chain rule

Parameter gradient

Lemma

The gradients

$$\delta_P^l := \nabla_{\mathbf{z}_P^l} \ell(\mathbf{y}, \mathbf{y}^L) \quad \text{and}$$

$$\delta_D^l := \nabla_{\mathbf{z}_D^l} \ell(\mathbf{y}, \mathbf{y}^L)$$

can be computed recursively as outlined in Algorithm 2. Moreover, it holds that

$$\nabla_{\mathbf{K}} \ell(\mathbf{y}, \mathbf{y}^L) = \sum_{l=1}^L \sigma \delta_D^l (\bar{\mathbf{y}}^{l-1})^\top - \tau \boldsymbol{\psi}^l (\delta_P^l)^\top.$$

Proof.

\mathbf{K} affects the objective exactly through $\mathbf{z}_P^1, \dots, \mathbf{z}_P^L$ and $\mathbf{z}_D^1, \dots, \mathbf{z}_D^L$. \rightarrow Chain rule. □

Limit of parameter gradient

Theorem

Suppose that there exist constants $c \geq 0$ and $0 \leq \kappa < 1$ such that

$$\|\delta_D^{l,L}\| \leq c\kappa^{L-l} \quad \text{and} \quad \|\delta_P^{l,L}\| \leq c\kappa^{L-l}$$

hold for arbitrary L and $l \in \{1, \dots, L\}$. Then, the limits

$$\Delta_P := \lim_{L \rightarrow \infty} \sum_{l=1}^L \delta_P^{l,L} \quad \text{and} \quad \Delta_D := \lim_{L \rightarrow \infty} \sum_{l=1}^L \delta_D^{l,L}$$

exist and are finite. Moreover, it holds that

$$\lim_{L \rightarrow \infty} \nabla_{\mathbf{K}} \ell(\mathbf{y}, \mathbf{y}^L) = \sigma \Delta_D(\mathbf{y}^*)^\top - \tau \boldsymbol{\psi}^*(\Delta_P)^\top.$$

Limit of parameter gradient

Proof.

1. Show that Δ_p and Δ_D exist and are finite

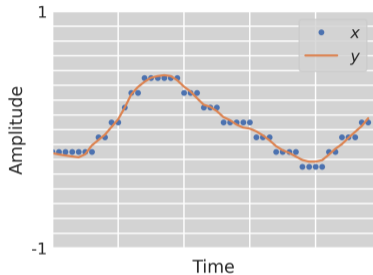
$$\lim_{L \rightarrow \infty} \sum_{l=1}^L \|\delta_p^{l,L}\| \leq \lim_{L \rightarrow \infty} \sum_{l=1}^L c\kappa^{L-l} = \lim_{L \rightarrow \infty} c \frac{1-\kappa^L}{1-\kappa} = \frac{c}{1-\kappa} < \infty$$

2. Rearrange gradient formula

$$\begin{aligned} \nabla_{\kappa} \ell(y, y^L) &= \sum_{l=1}^L \sigma \delta_D^{l,L} (\bar{y}^{l-1})^\top - \tau \psi^l (\delta_p^{l,L})^\top \\ &= \sum_{l=1}^L \sigma \delta_D^{l,L} (y^*)^\top - \tau \psi^* (\delta_p^{l,L})^\top - \sigma \delta_D^{l,L} (y^* - \bar{y}^{l-1})^\top + \tau (\psi^* - \psi^l) (\delta_p^{l,L})^\top \\ &= \underbrace{\sigma \left(\sum_{l=1}^L \delta_D^{l,L} (y^*)^\top \right)}_{\xrightarrow{L \rightarrow \infty} \sigma \Delta_D (y^*)^\top} - \underbrace{\tau \psi^* \left(\sum_{l=1}^L \delta_p^{l,L} \right)^\top}_{\xrightarrow{L \rightarrow \infty} \tau \psi^* (\Delta_p)^\top} - \underbrace{\sigma \sum_{l=1}^L \delta_D^{l,L} (y^* - \bar{y}^{l-1})^\top}_{\xrightarrow{L \rightarrow \infty} \mathbf{o}} + \underbrace{\tau \sum_{l=1}^L (\psi^* - \psi^l) (\delta_p^{l,L})^\top}_{\xrightarrow{L \rightarrow \infty} \mathbf{o}} \end{aligned}$$

□

Speech dequantization



$$\hat{\mathbf{y}} \in \underset{\mathbf{y}}{\operatorname{argmin}} \underbrace{\|\mathbf{K}\mathbf{y}\|_1}_{F(\mathbf{K}\mathbf{y})} \quad \text{s.t.} \quad \underbrace{\|\mathbf{y} - \mathbf{x}\|_\infty \leq \frac{\eta}{2}}_{G(\mathbf{y}-\mathbf{x}) = I_{\|\cdot\|_\infty \leq \frac{\eta}{2}}(\mathbf{y}-\mathbf{x})}$$

$$\mathbf{K} \in \mathbb{R}^{k \times n}$$

$$F^* = I_{\|\cdot\|_\infty \leq 1}$$

$$\operatorname{prox}_{\sigma F^*}(\mathbf{z}_D^l) = \min \left\{ 1, \max \left\{ -1, \mathbf{z}_D^l \right\} \right\}$$

$$\operatorname{prox}_{\tau G}(\mathbf{z}_P^l - \mathbf{x}) = \min \left\{ \frac{\eta}{2}, \max \left\{ -\frac{\eta}{2}, \mathbf{z}_P^l - \mathbf{x} \right\} \right\}$$

$$\delta_P^l = \overbrace{\operatorname{prox}'_{\tau G}(\mathbf{z}_P^l - \mathbf{x})}^{\in \{0,1\}^n} \odot (\delta_P^{l+1} + \sigma \mathbf{K}^\top \delta_D^{l+1})$$

$$\delta_D^l = \underbrace{\operatorname{prox}'_{\sigma F^*}(\mathbf{z}_D^l)}_{\in \{0,1\}^k} \odot (\delta_D^{l+1} - \tau \mathbf{K} \delta_P^l)$$

Speech dequantization

Theorem

Suppose that there exists an $l_0 \in \mathbb{N}$ such that $\text{prox}'_{\tau G}(\mathbf{z}_P^l - \mathbf{x}) = \text{prox}'_{\tau G}(\mathbf{z}_P^{l_0} - \mathbf{x})$ and $\text{prox}'_{\sigma F^*}(\mathbf{z}_D^l) = \text{prox}'_{\sigma F^*}(\mathbf{z}_D^{l_0})$ hold for all $l \geq l_0$.

Then, it holds for all $l \geq l_0$ that $\lim_{L \rightarrow \infty} \delta_P^l \in \ker(\mathbf{K})$ and $\lim_{L \rightarrow \infty} \delta_D^l \in \ker(\mathbf{K}^\top)$.

Proof.

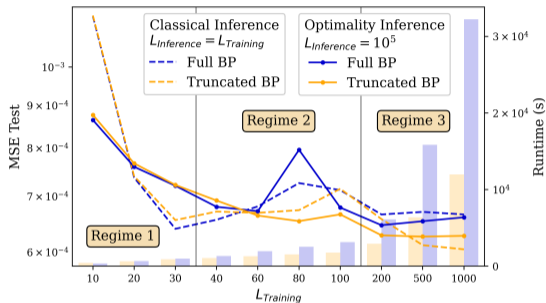
Set $\tilde{\mathbf{K}} := \text{prox}'_{\sigma F^*}(\mathbf{z}_D^{l_0}) \odot \mathbf{K} \odot \text{prox}'_{\tau G}(\mathbf{z}_P^{l_0} - \mathbf{x})$ and reverse engineer the optimization problem behind the Chambolle-Pock iteration to see that $\lim_{L \rightarrow \infty} \delta_D^l \in \text{argmin}_{\delta_D} \text{const. s.t. } \tilde{\mathbf{K}}^\top \delta_D = \mathbf{o}$ and analogously for the other case.

$$\delta_P^l = \delta_P^{l+1} + \sigma \tilde{\mathbf{K}}^\top \delta_D^{l+1}$$

$$\delta_D^l = \delta_D^{l+1} - \tau \tilde{\mathbf{K}} \delta_P^l$$

□

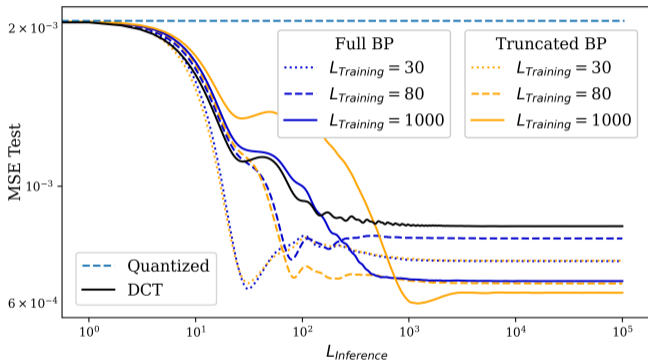
Truncated Backprop



- Regime 1. Full and truncated BP similar
- Regime 2. No significant performance increase
- Regime 3. Truncated BP outperforms full BP

- Full BP
 - Full backward pass according to Lemma
 - Requires $\sim 2L \times$ runtime of single iteration
 - Requires storage of $2L(k + n)$ variables
- Truncated BP
 - Use only $\mathbf{y}^L, \boldsymbol{\psi}^L, \mathbf{z}_P^L, \mathbf{z}_D^L$ and finite number of b backpropagated gradients to approximate $\lim_{L \rightarrow \infty} \nabla_{\mathbf{K}} \ell(\mathbf{y}, \mathbf{y}^L) = \sigma_{\Delta_D}(\mathbf{y}^*)^\top - \tau \boldsymbol{\psi}^*(\Delta_P)^\top$
 - Requires $\sim L + b \times$ runtime of single iteration
 - Requires storage of $2(k + n)$ variables

Interpretability



- **Bump at 30 iterations** illustrates gap between classical and optimality inference in Regime 1
- Classical unrolling is **prone to overfitting** to the number of unrolled iterations
- The best performing model in terms of the error here also features the slowest convergence speed

Summary

- **Asymptotic Analysis.** Limit of parameter gradient depends only on optimal solutions
- **Truncated backpropagation.** Outperforms usage of full gradients
- **Interpretability.** Limited when too few unrolled iterations during training
- **Future work.** Different algorithms, combine high interpretability and convergence speed

https://github.com/chrbraue/primal_dual_networks

Thank you!