

Neural Network Training as an Optimal Control Problem — An Augmented Lagrangian Approach —

Brecht Evens¹ Puya Latafat¹ Andreas Themelis² Johan Suykens¹ Panagiotis Patrinos¹

Abstract—Training of neural networks amounts to nonconvex optimization problems that are typically solved by using backpropagation and (variants of) stochastic gradient descent. In this work, we propose an alternative approach by viewing the training task as a nonlinear optimal control problem. Under this lens, backpropagation amounts to the sequential approach (single shooting) to optimal control, where the states variables have been eliminated. It is well known that single shooting may lead to ill-conditioning, and for this reason the simultaneous approach (multiple shooting) is typically preferred. Motivated by this hypothesis, an augmented Lagrangian algorithm is developed that only requires an approximate solution to the Lagrangian subproblems up to a user-defined accuracy. By applying this framework to the training of neural networks, it is shown that the inner Lagrangian subproblems are amenable to be solved using Gauss-Newton iterations. To fully exploit the structure of neural networks, the resulting linear least-squares problems are addressed by employing an approach based on forward dynamic programming. Finally, the effectiveness of our method is showcased on regression datasets.

I. INTRODUCTION

Feedforward deep neural networks (DNNs) are a prominent model for supervised learning, having a lot of success in various fields. The primary objective of this work is to devise a novel method for training DNNs with smooth activation functions; this task can be formally stated as follows.

Main problem. Given pairs $\{(a^{(\ell)}, b^{(\ell)}) \in \mathbb{R}^{d_0} \times \mathbb{R}^{d_{N+1}}\}_{\ell \in [m]}$, continuously differentiable functions $\{\Phi_j : \mathbb{R}^{d_j} \rightarrow \mathbb{R}^{d_{j+1}}\}_{j \in [N]}$ (operating in an element-wise fashion), and $\mu_w > 0$, find $\{W_j \in \mathbb{R}^{d_j \times d_{j-1}}\}_{j \in [N]}$ solutions to

$$\text{minimize}_{W_1, \dots, W_N} \frac{1}{2m} \sum_{\ell=1}^m \|\Phi_{N+1}(W_{N+1}x_N^{(\ell)}) - b^{(\ell)}\|^2 + \frac{\mu_w}{2} \sum_{j=1}^{N+1} \|W_j\|_F^2$$

$$\text{where } x_0^{(\ell)} := a^{(\ell)}, \quad \ell \in [m], \quad (1a)$$

$$x_j^{(\ell)} := \Phi_j(W_j x_{j-1}^{(\ell)}), \quad j \in [N], \ell \in [m]. \quad (1b)$$

Here, $(a^{(\ell)}, b^{(\ell)})$ are (given) training pairs, $N \in \mathbb{N}$ is the number of layers of the network, each one having d_i many

neurons/nodes and with Φ_i being the corresponding activation function, and μ_w is a regularization parameter for the weights W_i commonly used to avoid overfitting [15]. These optimization problems are typically solved using backpropagation [19] along with (variants of) stochastic gradient descent, due to their simplicity and effectiveness. However, these optimization methods suffer from various issues related to the challenging, highly nonconvex nature of the training task. First and foremost, due to the prominence of local minima and saddle points, trained DNN models tend to generalize poorly to test data. To alleviate this issue, various regularization methods have been introduced such as weight decay [15], batch normalization [13], and dropout [20], typically reducing the overfitting of the training data. More fundamentally, gradient-based methods are known to suffer from the vanishing gradient phenomenon [12], where the gradients in the output layers of DNNs decrease exponentially with the number of layers. Although recent studies have shown that piecewise affine activation functions such as ReLU, leaky ReLU [17], and Maxout unit [8] reduce the vanishing gradient problem by making the problem more sparse, the issue nevertheless persists especially in very deep networks.

To address these issues, in recent years a host of auxiliary variable methods have been introduced where the network structure is represented by equality constraints and the space of learning parameters is extended. By lifting the number of variables, these methods decompose the training task into a series of local subproblems which can be solved deterministically, typically using block coordinate descent (BCD) [4, 10, 24] or the alternating direction method of multipliers (ADMM) [21, 23, 25]. BCD and ADMM have been successful for this task due to their ability to convert the equality constrained optimization problems into unconstrained problems, which can then be solved more efficiently than their constrained counterparts. By increasing the dimension of the training problem, auxiliary variable methods can alleviate some of the issues from which classical gradient-based methods suffer. Most notably, it is observed that the vanishing gradient issue is alleviated as the auxiliary variables circumvent long-term dependencies between the network weights during training [25]. On the other hand, the increased dimensionality naturally makes the training task more challenging than when using classical gradient-based approaches.

The difference between traditional methods and auxiliary variable methods can be related to concepts from optimal control by viewing the training task as a nonlinear optimal control problem. Under this lens, auxiliary variable methods amount to the simultaneous approach (multiple shooting),

This work was supported by the Research Foundation Flanders (FWO) research projects G0A0920N, G086518N, G086318N, and PhD grant 1196820N; Research Council KU Leuven C1 project No. C14/18/068; Fonds de la Recherche Scientifique – FNRS and the Fonds Wetenschappelijk Onderzoek – Vlaanderen under EOS project no 30468160 (SeLMA); and the Japan Society for the Promotion of Science (JSPS) KAKENHI grant JP21K17710. Johan Suykens and Panagiotis Patrinos are affiliated to Leuven.AI - KU Leuven institute for AI, B-3000, Leuven, Belgium.

¹KU Leuven, Department of Electrical Engineering ESAT-STADIUS – Kasteelpark Arenberg 10, bus 2446, B-3001 Leuven, Belgium {brecht.evens,puya.latafat,johan.suykens,panos.patrinos}@kuleuven.be

²Kyushu University, Faculty of Information Science and Electrical Engineering (ISEE) – 744 Motoooka, Nishi-ku, Fukuoka 819-0395, Japan andreas.themelis@ees.kyushu-u.ac.jp

whereas backpropagation amounts to the sequential approach (single shooting), where the state variables are eliminated [16]. As it is well known that single shooting may lead to ill-conditioning of the optimization problem, it can be expected that multiple shooting methods can provide major advantages in the learning process of DNNs.

Motivated by this hypothesis, we develop a training methodology for neural networks based on an augmented Lagrangian framework that only requires finding approximate stationary points of the Lagrangian subproblems up to a user-defined accuracy. To fully exploit the structure of feed-forward neural networks, we additionally provide a computationally efficient approach to solve the inner subproblems based on forward dynamic programming. The overall approach leads to an efficient and provably convergent methodology for solving the highly nonconvex optimization problems emerging in the neural network training task.

A. Contributions

The contribution of this paper is twofold:

1) We introduce a novel augmented Lagrangian framework (ALM) for solving general nonconvex and nonsmooth equality constrained optimization problems. The framework is inspired by and extends [9, Alg. 1] by waiving smoothness assumptions and introducing a less conservative penalty update rule, yet preserving convergence to an approximate KKT point.

2) We apply this framework to the training of DNNs, which we address from an optimal control perspective. The resulting optimization problem's structure has a twofold benefit: first, the inner Lagrangian subproblems are amenable to be addressed with fast methods such as Gauss-Newton (GN); in turn, forward dynamic programming (FDP) can conveniently be employed to efficiently solve the resulting linear least-squares problems.

To reflect the modularity and the contribution of each component, the three procedures (outer ALM, inner GN, and FDP) are outlined in three standalone algorithms, each addressing a dedicated general problem.

B. Organization

The paper is organized as follows. The notation is introduced in the next subsection. An optimal control reformulation for the NN problem is presented in Section II. In Section III a novel augmented Lagrangian method (ALM) is proposed for general equality constrained nonlinear programs. The ALM method is specialized for the training of neural networks with smooth activation functions in Section IV, where a procedure based on the Gauss-Newton method and forward dynamic programming is proposed. The proofs of all the results are deferred to the appendix. Finally, numerical simulations showcasing the effectiveness of our proposed methodology on regression datasets are discussed in Section V.

C. Notation

We use $[N]$ to denote the set of indices $\{1, \dots, N\}$. We denote by \mathbb{R}^n the standard n -dimensional Euclidean space with inner product $\langle \cdot, \cdot \rangle$ and induced norm $\|\cdot\|$. The set of extended real numbers is defined as $\overline{\mathbb{R}} := \mathbb{R} \cup \{\infty\}$, and we say that an extended-real valued function $f : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ is proper if $\text{dom } f := \{x \in \mathbb{R}^n \mid f(x) < \infty\}$ is nonempty. The set of real n -by- m matrices is denoted by $\mathbb{R}^{n \times m}$. Given $A \in \mathbb{R}^{n \times m}$, $\|A\|_F$ is its Frobenius norm and $\text{vec}(A) \in \mathbb{R}^{nm}$ is the vector obtained by stacking the columns of A on top of one another. The sets of symmetric, symmetric positive semi-definite and symmetric positive definite n -by- n matrices are denoted by \mathbb{S}^n , \mathbb{S}_+^n and \mathbb{S}_{++}^n , respectively. For $V \in \mathbb{S}_{++}^n$ we define the scalar product $\langle x, y \rangle_V = \langle x, Vy \rangle$ and the induced norm $\|x\|_V = \sqrt{\langle x, x \rangle_V}$. The n -by- n identity matrix is denoted by I_n , or simply I when no ambiguity occurs. The vector of all zeros with dimension n and the n -by- m matrix of all zeros are denoted by 0_n and $0_{n \times m}$, respectively. The matrix Kronecker product is denoted by \otimes . The Jacobian of a differentiable function $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$, is denoted by $JF : \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$; $J_x F$ is a short-hand notation for the partial derivative $\frac{\partial F}{\partial x}$.

II. AN OPTIMAL CONTROL REFORMULATION

In the traditional approach, (1a) and (1b) are absorbed into the cost, thus forming an unconstrained minimization which is then solved by employing a stochastic (sub)gradient-type method. Here we take an alternative approach by viewing the minimization as an optimal control problem with N stages. To this end, (1) represents the dynamics of the problem and may compactly be written as

$$X_j = \Phi_j(W_j X_{j-1}), \quad j \in [N], \quad (2)$$

where $X_j \in \mathbb{R}^{d_j \times m}$ is a matrix whose i -th column is the vector $x_j^{(i)}$, for $i \in [m]$. By similarly letting $A \in \mathbb{R}^{d_0 \times m}$ and $Y \in \mathbb{R}^{d_{N+1} \times m}$ denote the input and output matrices (constructed using vectors $a^{(i)}$, $b^{(i)}$), the following compact reformulation of (1) is obtained

$$\begin{aligned} & \underset{(W_i)_{i \in [N+1]}, (X_i)_{i \in [N]}}{\text{minimize}} && \frac{1}{2m} \|\Phi_{N+1}(W_{N+1} X_N) - Y\|_F^2 + \frac{\mu_w}{2} \sum_{i=1}^{N+1} \|W_i\|_F^2 \\ & \text{subject to} && X_0 = A \\ & && X_{j+1} = \Phi_j(W_j X_j), \quad j \in [N]. \end{aligned} \quad (3)$$

A. Vectorized form

For simplicity of exposition and computational convenience, we condense the optimization variables W_i and X_i into a single long vector $\mathbf{z} = (\mathbf{w}, \mathbf{x})$ with

$$\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_{N+1}) \quad \text{and} \quad \mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$$

where, letting $w_{i,j} \in \mathbb{R}^{d_i-1}$ denote the j -th row of W_i and $x_i^{(j)}$ the j -th column of X_i ,

$$\mathbf{w}_i = (w_{i,1}, \dots, w_{i,d_i}) \in \mathbb{R}^{d_i d_i-1} \quad \text{and} \quad \mathbf{x}_i = (x_i^{(1)}, \dots, x_i^{(m)}) \in \mathbb{R}^{m d_i}.$$

In the vectorized notation, the cost function and the nonlinear constraints in (3) may be represented by f and $F(\mathbf{z}) = 0$ with

$$f(\mathbf{z}) = \frac{1}{2m} \|H_{N+1}(\mathbf{w}_{N+1}, \mathbf{x}_N) - \mathbf{y}\|^2 + \frac{\mu_w}{2} \|\mathbf{w}\|^2,$$

$$F(\mathbf{z}) = (\mathbf{x}_1 - H_1(\mathbf{w}_1, \mathbf{x}_0), \dots, \mathbf{x}_{N+1} - H_{N+1}(\mathbf{w}_{N+1}, \mathbf{x}_N)),$$

where $\mathbf{y} = \text{vec}(Y)$ and

$$H_j(\mathbf{w}_j, \mathbf{x}_{j-1}) = (\Phi_j(W_j x_{j-1}^{(1)}), \dots, \Phi_j(W_j x_{j-1}^{(m)})). \quad (4)$$

III. THE OUTER ALM ALGORITHM

With vectorized notation being adopted and as long as the activation functions Φ_j are locally Lipschitz, the minimization in (3) falls into the following general setting.

Problem I (General ALM framework). *For a proper, lower semicontinuous, lower bounded $f : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ and a locally Lipschitz $F : \mathbb{R}^n \rightarrow \mathbb{R}^p$ such that $\{z \in \text{dom } f \mid F(z) = 0\} \neq \emptyset$,*

$$\text{minimize}_{z \in \mathbb{R}^n} f(z) \quad \text{subject to } F(z) = 0. \quad (5)$$

This section proposes a conceptual algorithm for addressing **Problem I**, conceptual in the sense that, at this stage, no hint is given as to how the inner subproblems it involves can be solved. The algorithm will be concretized in the subsequent **Section IV**, where an implementable procedure for addressing these inner steps is detailed. The chosen method for the inner subproblems will ultimately require some additional structure and differentiability assumptions, which are nevertheless not needed for the (outer) ALM scheme presented in this section. For the sake of generality of the discussion and to well pinpoint where each requirement is invoked, the convergence proof of the outer scheme is given in this broader setting.

Equality constrained minimization problems as (5) are amenable to be addressed by means of augmented Lagrangian methods. For $\beta > 0$, we denote the corresponding β -augmented Lagrangian as

$$\begin{aligned} \mathcal{L}_\beta(z, \lambda) &:= f(z) + \langle \lambda, F(z) \rangle + \frac{\beta}{2} \|F(z)\|^2 \\ &= f(z) + \frac{\beta}{2} \|F(z)\|^2 + \lambda/\beta - \frac{1}{2\beta} \|\lambda\|^2, \end{aligned} \quad (6)$$

and we say that (z, λ) is an ε -KKT pair if

$$\|\nabla_z \mathcal{L}(z, \lambda)\|_\infty \leq \varepsilon, \quad \text{and} \quad (7a)$$

$$\|F(z)\|_\infty \leq \varepsilon, \quad (7b)$$

where $\mathcal{L} := \mathcal{L}_0$ is the (non-augmented) Lagrangian, and ∇_z denotes the gradient with respect to z or, in case of lack of differentiability, any vector in the subdifferential $\partial_z \mathcal{L}(z, \lambda)$.

Largely inspired by [9, Alg. 1], **Algorithm 1** hinges on the upper boundedness of the augmented Lagrangian along the iterates (see, e.g. [3, Ex. 4.12]) ensured by the initialization at a feasible point z^0 . Additionally to reducing the assumptions to the general setting of **Problem I**, we also introduced a less conservative update rule for the penalty parameter β , thus waiving the need for its unbounded growth which instead happens in [9] (unless the dual variable converges to 0).

Theorem 1. *Applied to **Problem I**, **Algorithm 1** terminates yielding an ε -KKT pair for (5).*

Note that the result can cope with rather general functions f and F , not necessarily derived from formulations as in (3). The optimal control structure will instead be exploited in the following **Section IV** where an iterative method for addressing the inner problems at **step 1.2** will be given.

Remark 2. When f is lower bounded, then so is $\mathcal{L}_\beta(\cdot, \lambda)$ for any λ , thus ensuring the existence of ε -stationary points as required in **step 1.2** for any $\varepsilon > 0$. In the setting of the optimal control problem (3), not only is this condition trivially

Algorithm 1 ALM for **Problem I**

REQUIRE Initial feasible point $z^0 \in \text{dom } f$ s.t. $F(z^0) = 0$, multiplier λ^0 , and penalty $\beta_0 > 0$
Parameters $0 < \gamma < 1 < \alpha, \xi$ and tolerance $\varepsilon > 0$

For $k = 0, 1, 2 \dots$

1.1: Set $\hat{z}^k = z^k$ if $\mathcal{L}_{\beta_k}(z^k, \lambda^k) \leq f(z^0)$, or $\hat{z}^k = z^0$ otherwise

1.2: Starting at \hat{z}^k , apply a descent method to compute an ε -stationary point z^{k+1} of

$$\text{minimize } \mathcal{L}_{\beta_k}(\cdot, \lambda^k), \quad (8)$$

i.e., a point z^{k+1} such that

$$\|\nabla_z \mathcal{L}_{\beta_k}(z^{k+1}, \lambda^k)\|_\infty \leq \varepsilon \quad (9)$$

1.3: Set $\lambda^{k+1} = \lambda^k + \beta_k F(z^{k+1})$

1.4: IF $\|F(z^{k+1})\|_\infty \leq \varepsilon$ THEN
RETURN ε -KKT pair (z^{k+1}, λ^{k+1})

1.5: Set $\beta_{k+1} = \beta_k$ if $\|F(z^{k+1})\|_\infty \leq \gamma \|F(z^k)\|_\infty$,
or $\beta_{k+1} = \max\{\xi\beta_k, \beta_0(k+1)^\alpha\}$ otherwise.

satisfied, but a feasible starting point z^0 can be obtained at virtually no cost by initializing the weights W_j^0 and unrolling the dynamics to generate the *state* variables X_j^0 . \square

IV. THE LAGRANGIAN SUBPROBLEM VIA GAUSS-NEWTON ITERATIONS

In this section we present a procedure for solving the inner minimization (8) in the setting of NNs with continuously differentiable activation functions. With the notational conventions of **Section II-A**, for a fixed multiplier $\lambda = (\lambda_1, \dots, \lambda_N)$ (λ_j being the one associated with the j -th dynamics) the Lagrangian subproblem associated with (3) is cast as follows.

Problem II (Lagrangian subproblem). *Given smooth functions $\{H_j\}_{j \in [N+1]}$, vectors λ, y of suitable sizes, and $\beta, \mu_w > 0$,*

$$\begin{aligned} \text{minimize}_{z=(w,x)} \mathcal{L}_\beta(z, \lambda) &:= \frac{1}{2m} \|H_{N+1}(w_{N+1}, x_N) - y\|^2 + \frac{\mu_w}{2} \|w\|^2 \\ &\quad - \frac{1}{2\beta} \|\lambda\|^2 + \frac{\beta}{2} \sum_{j=1}^N \|x_j - H_j(w_j, x_{j-1}) + \lambda_j/\beta\|^2. \end{aligned}$$

This smooth unconstrained least-squares problem is amenable to be solved by the Gauss-Newton (GN) method, which amounts to iteratively solving minimizations obtained after linearizing functions H_j around the last iterates, and then applying a standard line search to guarantee convergence. In the next subsection, we derive explicit expressions of the Jacobian matrices involved in the linearization.

A. Gauss-Newton linearization and update direction

Let $D^j : \mathbb{R}^{d_{j-1}} \rightarrow \mathbb{R}^{d_j \times d_j}$ be given by

$$D^j(v) := \text{diag}(\Phi'_j(\langle w_{j-1}, v \rangle), \dots, \Phi'_j(\langle w_{j,d_j}, v \rangle)),$$

(recall that Φ_j operates element-wise) and define

$$\mathcal{D}^j := \text{blkdiag}(D^j(x_{j-1}^{(1)}), \dots, D^j(x_{j-1}^{(m)})).$$

The Jacobians $J_{x_{j-1}} H_j \in \mathbb{R}^{m d_j \times m d_{j-1}}$ and $J_w H_j \in \mathbb{R}^{m d_j \times d_{j-1}}$ are then given by

$$J_{x_{j-1}} H_j(w_j, x_{j-1}) = \mathcal{D}^j(\mathbf{I}_m \otimes W_j), \quad \text{and}$$

Algorithm 2 Gauss-Newton procedure for **Problem II**

REQUIRE Initial point $\mathbf{z}^0 = (\mathbf{w}^0, \mathbf{x}^0)$ and $0 < \eta_1, \eta_2 < 1$ For $l = 0, 1, 2, \dots$ 2.1: [UPDATE DIRECTION] set $\mathbf{p}^l = \bar{\mathbf{z}}^l - \mathbf{z}^l$, where $\bar{\mathbf{z}}^l = (\bar{\mathbf{w}}^l, \bar{\mathbf{x}}^l)$ solves **Problem III** with $\mathcal{A}, \mathcal{B}, \mathbf{c}$ as in (10)2.2: [LINE SEARCH] set $\mathbf{z}^{l+1} = \mathbf{z}^l + \tau_l \mathbf{p}^l$, where τ_l is the largest number in $\{1, \eta_1, \eta_1^2, \dots\}$ such that

$$\mathcal{L}_\beta(\mathbf{z}^l + \tau_l \mathbf{p}^l, \lambda) \leq \mathcal{L}_\beta(\mathbf{z}^l, \lambda) - 2\eta_2 \tau_l \mathcal{G}_{\mathcal{A}, \mathcal{B}, 0}(\mathbf{p}^l)$$

with \mathcal{L}_β and $\mathcal{G}_{\mathcal{A}, \mathcal{B}, \mathbf{c}}$ as in **Problems II** and **III**2.3: IF $\|\nabla_{\mathbf{z}} \mathcal{L}_\beta(\mathbf{z}^{l+1}, \lambda)\|_\infty \leq \varepsilon$ THENRETURN $\mathbf{z}^{l+1} = (\mathbf{w}^{l+1}, \mathbf{x}^{l+1})$

$$\mathbf{J}_{\mathbf{w}_j} H_j(\mathbf{w}_j, \mathbf{x}_{j-1}) = \mathcal{D}^j(\mathbf{I}_{d_j} \otimes \mathbf{x}_{j-1}^{(1)}, \dots, \mathbf{I}_{d_j} \otimes \mathbf{x}_{j-1}^{(m)})^\top.$$

If $\mathbf{z}^l = (\mathbf{w}^l, \mathbf{x}^l)$ is the l -th iterate of a GN algorithm, denoting

$$\begin{aligned} A_{j+1} &= \begin{cases} 0_{m d_1 \times m d_0} & \text{if } j = 0 \\ \mathbf{J}_{\mathbf{x}_j} H_{j+1}(\mathbf{w}_{j+1}^l, \mathbf{x}_j^l) & \text{if } j \in [N] \end{cases} \\ B_j &= \mathbf{J}_{\mathbf{w}_j} H_j(\mathbf{w}_j^l, \mathbf{x}_{j-1}^l), \quad j \in [N+1] \\ \mathbf{c}_j &= \begin{cases} H_j(\mathbf{w}_j^l, \mathbf{x}_{j-1}^l) - A_j \mathbf{x}_{j-1}^l - B_j \mathbf{w}_j^l - \frac{1}{\beta} \lambda_j & \text{if } j \in [N] \\ H_j(\mathbf{w}_j^l, \mathbf{x}_{j-1}^l) - A_j \mathbf{x}_{j-1}^l - B_j \mathbf{w}_j^l - \mathbf{y} & \text{if } j = N+1, \end{cases} \end{aligned} \quad (10)$$

the linearized minimization yielding the l -th GN update direction reduces to the following problem.**Problem III** (GN direction). Given $\mathcal{A} = (A_1, \dots, A_{N+1})$, $\mathcal{B} = (B_1, \dots, B_{N+1})$ and $\mathbf{c} = (\mathbf{c}_1, \dots, \mathbf{c}_{N+1})$ with matrices A_j, B_j and vectors \mathbf{c}_j of suitable sizes, and given scalars $\beta, \mu_w > 0$,

$$\text{minimize } \mathcal{G}_{\mathcal{A}, \mathcal{B}, \mathbf{c}}(\mathbf{z}) \quad \text{where,}$$

denoting $(\delta_j, \rho_j) = (1, \beta)$ for $j \leq N$ and $(0, \frac{1}{m})$ otherwise,

$$\mathcal{G}_{\mathcal{A}, \mathcal{B}, \mathbf{c}}(\mathbf{z}) := \sum_{j=1}^{N+1} \left(\frac{\rho_j}{2} \|\delta_j \mathbf{x}_j - A_j \mathbf{x}_{j-1} - B_j \mathbf{w}_j - \mathbf{c}_j\|^2 + \frac{\mu_w}{2} \|\mathbf{w}_j\|^2 \right).$$

B. The Gauss-Newton algorithm

The structure of **Problem III** emphasizes how variables are weakly coupled, a phenomenon that owes to the stagewise structure of the optimal control problem (3). As a result, in spite of the large scale, **Problem III** admits a closed-form solution that is efficiently retrievable with a forward dynamic programming (FDP) approach detailed in the following **Section IV-C**. This routine may then be invoked by the GN method, synopsised in **Algorithm 2**, when computing the update directions at **step 2.1**.

In the next lemma we show that the GN method yields an ε -stationary solution for the original Lagrangian subproblem.

Lemma 3. *The limit points of the sequence $(\mathbf{z}^l)_{l \in \mathbb{N}}$ generated by **Algorithm 2** are stationary for **Problem II**. In particular, if $(\mathbf{z}^l)_{l \in \mathbb{N}}$ remains bounded (e.g., when $\mathcal{L}_\beta(\cdot, \lambda)$ is coercive), then **Algorithm 2** yields an ε -stationary solution.*

C. Forward dynamic programming

In this subsection we propose a recursive procedure for solving **Problem III** with given matrices $A_j \in \mathbb{R}^{r_j \times r_{j-1}}$, $B_j \in \mathbb{R}^{r_j \times s_j}$, and vectors $\mathbf{c}_j \in \mathbb{R}^{r_j}$, $j \in [N+1]$, thus providing an efficient routine for **step 2.1** of **Algorithm 2**. Inspired by

Algorithm 3 Recursive solution to **Problem III** with FDP

REQUIRE Initial state $\mathbf{x}_0 \in \mathbb{R}^{r_0}$

$$\text{set } M_1 = \frac{1}{\rho_1} \mathbf{I} + \frac{1}{\mu_w} B_1 B_1^\top, S_1 = \mathbf{I}_{r_0}, \mathbf{q}_0 = \mathbf{x}_0$$

3.1: [FORWARD RECURSION] For $j = 1, \dots, N$ (a) solve the linear system $M_j \tilde{\mathbf{c}}_j = \mathbf{c}_j$,(b) $\tilde{\mathbf{q}}_j = S_j \mathbf{q}_{j-1}$ ▷ as described in (18)(c) $\mathbf{q}_j = \rho_j G_j A_j \tilde{\mathbf{q}}_j + \tilde{\mathbf{c}}_j$ (d) $M_{j+1} = \frac{1}{\rho_{j+1}} \mathbf{I} + \frac{1}{\mu_w} B_{j+1} B_{j+1}^\top + A_{j+1} M_j A_{j+1}^\top$

3.2: [BACKWARD RECURSION]

(a) $\tilde{\mathbf{x}}_{N+1} = -S_{N+1} (A_{N+1}^\top G_{N+1} \mathbf{c}_{N+1})$ $\tilde{\mathbf{q}}_{N+1} = S_{N+1} \mathbf{q}_N$ ▷ as described in (18)(b) $\mathbf{x}_N = \tilde{\mathbf{q}}_{N+1} + \rho_{N+1} \tilde{\mathbf{x}}_{N+1}$ $\mathbf{w}_{N+1} = -E_{N+1} (A_{N+1} \mathbf{x}_N + \mathbf{c}_{N+1})$ (c) For $j = N, \dots, 2$: $\tilde{\mathbf{x}}_j = S_j (A_j^\top G_j (\mathbf{x}_j - \mathbf{c}_j))$ ▷ as described in (18) $\mathbf{x}_{j-1} = \tilde{\mathbf{q}}_j + \rho_j \tilde{\mathbf{x}}_j$ $\mathbf{w}_j = E_j (\mathbf{x}_j - A_j \mathbf{x}_{j-1} - \mathbf{c}_j)$ (d) $u_1 = E_1 (\mathbf{x}_1 - A_1 \mathbf{x}_0 - \mathbf{c}_1)$ RETURN $\mathbf{z} = (\mathbf{w}, \mathbf{x})$ with $\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_{N+1})$, $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$

the idea of forward dynamic programming, the minimization may be split into a series of simpler subproblems that are solved in a recursive manner:

$$V_1^*(\mathbf{x}_1) = \min_{\mathbf{w}_1} \left\{ \frac{\rho_1}{2} \|\mathbf{x}_1 - A_1 \mathbf{x}_0 - B_1 \mathbf{w}_1 - \mathbf{c}_1\|^2 + \frac{\mu_w}{2} \|\mathbf{w}_1\|^2 \right\} \quad (11)$$

$$\begin{aligned} V_j^*(\mathbf{x}_j) &= \min_{\mathbf{x}_{j-1}, \mathbf{w}_j} \left\{ V_{j-1}^*(\mathbf{x}_{j-1}) + \frac{\rho_j}{2} \|\mathbf{x}_j - A_j \mathbf{x}_{j-1} - B_j \mathbf{w}_j - \mathbf{c}_j\|^2 \right. \\ &\quad \left. + \frac{\mu_w}{2} \|\mathbf{w}_j\|^2 \right\}, \quad j = 2, \dots, N \end{aligned} \quad (12)$$

$$\begin{aligned} V_{N+1}^* &= \min_{\mathbf{x}_N, \mathbf{w}_{N+1}} \left\{ V_N^*(\mathbf{x}_N) + \frac{\rho_{N+1}}{2} \|A_{N+1} \mathbf{x}_N + B_{N+1} \mathbf{w}_{N+1} + \mathbf{c}_{N+1}\|^2 \right. \\ &\quad \left. + \frac{\mu_w}{2} \|\mathbf{w}_{N+1}\|^2 \right\}. \end{aligned} \quad (13)$$

Each stage consists of the minimization of the sum of the cost at the current stage and the optimal cost from the previous stage. The cost at the final stage V_{N+1}^* is equal to the optimal cost for **Problem III**. In order to obtain closed form solutions for each of the above minimizations, let $E_j \in \mathbb{R}^{s_j \times r_j}$ and $G_j, M_j, S_j \in \mathbb{R}^{r_j \times r_j}$, $j \in [N+1]$, be defined as

$$E_j = \left(\frac{\mu_w}{\rho_j} \mathbf{I} + B_j^\top B_j \right)^{-1} B_j^\top, \quad (14)$$

$$G_j = \mathbf{I} - B_j E_j, \quad (15)$$

$$M_j = \begin{cases} \frac{1}{\rho_1} \mathbf{I} + \frac{1}{\mu_w} B_1 B_1^\top & \text{if } j = 1 \\ \frac{1}{\rho_j} \mathbf{I} + \frac{1}{\mu_w} B_j B_j^\top + A_j M_{j-1} A_j^\top & \text{if } j > 1 \end{cases} \quad (16)$$

$$S_j = \begin{cases} \mathbf{I}_{r_0} & \text{if } j = 1 \\ M_{j-1} - M_{j-1} A_j^\top M_j^{-1} A_j M_{j-1} & \text{if } j > 1. \end{cases} \quad (17)$$

Note that matrices S_j need not be computed explicitly. Instead, given a vector $\mathbf{v} \in \mathbb{R}^{r_j}$, $S_j \mathbf{v}$ is computed as follows:

$$\begin{cases} (i) \text{ solve the linear system } M_j \bar{\mathbf{v}} = A_j (M_{j-1} \mathbf{v}) \\ (ii) \text{ set } S_j \mathbf{v} = M_{j-1} (\mathbf{v} - A_j^\top \bar{\mathbf{v}}). \end{cases} \quad (18)$$

The FDP procedure is presented in **Algorithm 3**. Other than matrix-vector products, the algorithm requires solving linear systems several times, which may be performed by

computing the *Cholesky factorization* of M_j and $\frac{\mu_w}{\rho_j} \mathbf{I} + B_j^\top B_j$ once, thus resulting in operations involving simple forward and backward substitution steps that substantially reduce the computational overhead.

Remark 4 (Positive definiteness). Since $\rho_j, \mu_w > 0$, it holds that $G_j, M_j \in \mathbb{S}_{++}^{r_j}$ for any j . Furthermore, using the Woodbury matrix identity and (16), the following alternative expression for S_{j+1} is obtained

$$S_{j+1} = (M_j^{-1} + \rho_{j+1} A_{j+1}^\top G_{j+1} A_{j+1})^{-1}, \quad (19)$$

establishing that also $S_{j+1} \in \mathbb{S}_{++}^{r_j}$. \square

The optimality of the solution obtained by the FDP procedure is established in the next lemma.

Lemma 5. *Suppose that $\mu_w > 0$. Then, $\mathbf{z} = (\mathbf{w}, \mathbf{x})$ generated by Algorithm 3 is the unique minimizer of Problem III.*

V. NUMERICAL EXPERIMENTS

A. Design of numerical experiments

We will generate training (and test) pairs $\{(a^{(\ell)}, b^{(\ell)})\}_{\ell \in [m]}$ for a three-layer neural network under the regression setting, analogous to the approach in [6], as follows:

$$b^{(\ell)} = W_3 \Phi(W_2 \Phi(W_1 a^{(\ell)})) + \delta,$$

where $a^{(\ell)} \sim \mathcal{N}(\mu, \Sigma)$ and $\delta \sim \delta_0 \mathcal{N}(0, 1)$. The mean $\mu \in \mathbb{R}^{d_0}$ and an additional random matrix $\Sigma_0 \in \mathbb{R}^{d_0 \times d_0}$ are generated by a normal distribution with standard deviation 0.2, and the covariance Σ is set to be $\Sigma_0^\top \Sigma_0$. The three-layer network consists of $N = 2$ hidden layers with respectively 20 and 5 neurons. As activation function the softplus function is used, i.e. $\Phi(x) := \ln(1 + \exp(x))$, a smooth approximation to the ReLU activation function which is often used in deep learning and known for its faster convergence. The weights W_i of the neural network are initialized according to Kaiming [11], which is a weight initialization procedure suitable for networks consisting of softplus activation functions, and we obtain a feasible starting point \mathbf{z}^0 by applying (2) recursively. All networks in this section are trained with regularization parameter $\mu_w = 0.1$. The following parameters for Algorithm 1 are used:

$$\begin{aligned} \lambda^0 &= 0, & \beta_0 &= 0.001 f(\mathbf{z}^0), & \gamma &= 0.5, \\ \alpha &= 2, & \varepsilon &= 10^{-3}, & \xi &= 2. \end{aligned}$$

Furthermore, to prevent solving the inner problems (8) up to an unnecessarily high tolerance ε in the first iterations, Eq. (9) is relaxed as follows:

$$\|\nabla_{\mathbf{z}} \mathcal{L}_{\beta_k}(\mathbf{z}^{k+1}, \lambda^k)\|_\infty \leq \varepsilon_k := \max(\bar{\varepsilon}, 0.5 \varepsilon_{k-1}) \quad (20)$$

with $\varepsilon_0 = 10^{-1}$ and $\bar{\varepsilon} = 10^{-2}$. Finally, the following parameters for the line search in Algorithm 2 are used:

$$\eta_1 = 0.8, \quad \eta_2 = 0.1.$$

The ALM framework and corresponding Gauss-Newton procedure are implemented using the SciPy sparse matrix library [22] in Python. The CHOLMOD library [5] is used to factorize $(\frac{\mu_w}{\rho_j} \mathbf{I} + B_j^\top B_j)$ and M_j . All experiments are conducted on an HP EliteBook 845 G7 with a 1.7GHz AMD Ryzen 7 PRO 4750U processor and 32 GB RAM.

B. Numerical results and discussion

The left-hand side of Table I shows the numerical results for training the previously introduced feedforward neural networks with varying input dimension d_0 and noise level δ_0 (averaged over 15 simulations) using our proposed ALM method, which yields an $\bar{\varepsilon}$ -KKT pair after only a couple of ALM iterations (as (7a) is satisfied for $\bar{\varepsilon}$ instead of ε).

All experiments are performed with a fixed sample size $m = 250$ for the training and test datasets. We should remark that the current implementation does not scale well with the sample size m both in terms of memory usage and computation time, as the matrices M_j in the FDP procedure become increasingly large. For this reason, our method would greatly benefit from a mini-batch implementation where the training set is split into smaller batches to compute the inner GN steps. This is considered as future work.

The typical performance of the ALM algorithm with low tolerance ($\varepsilon = \bar{\varepsilon} = 10^{-7}$) is visualized in Fig. 1 for a simulation with $d_0 = 15$ and $\delta_0 = 20\%$. In the earlier GN iterations mainly the loss is reduced, while in the final iterations the feasibility is recovered as the penalty parameter increases in the outer ALM iterations. For this reason, it makes sense to terminate our algorithm at higher tolerances, as in neural network training we are mainly interested in reducing the loss.

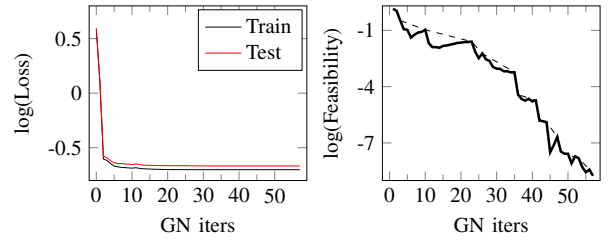


Fig. 1. Typical performance of the ALM algorithm. (Left) Training and test loss. (Right) Feasibility. The dashed line connects the points of the outer ALM iterations.

C. Comparison with first-order methods

We compare our method with two commonly used first-order methods for stochastic optimization, namely Adam [14] and stochastic gradient descent (SGD). We use the default implementations of these algorithms provided by the Keras library using the TensorFlow [1] backend with batch size 10, MSE loss function and ℓ_2 regularization with parameter μ_w .

The right portion of Table I shows the numerical results for training the three-layer network using Adam and SGD (averaged over 15 simulations) for 1000 epochs. No early stopping or other monitoring callbacks are used, minimizing the computation time per epoch. SGD is typically susceptible to stagnate at suboptimal points where it ceases to make significant progress, which explains its higher training MSE compared to Adam. When comparing with Adam and SGD it can be seen that our method tends to converge towards very good local optima, surpassing the performance

TABLE I

NUMERICAL RESULTS FOR TRAINING A THREE-LAYER NETWORK WITH VARYING INPUT DIMENSION AND NOISE LEVEL USING ALM, ADAM AND SGD.

| d_0 | δ_0 | ALM | | | | | | | Adam | | | SGD | | |
|-------|------------|--------------|----------|-------------------------------|--|-----------|----------|-------------|--------------|----------|-------------|--------------|----------|-------------|
| | | Training MSE | Test MSE | \mathcal{L}_{β_k} evals | $\nabla_x \mathcal{L}_{\beta_k}$ evals | ALM iters | GN iters | Time (m:ss) | Training MSE | Test MSE | Time (m:ss) | Training MSE | Test MSE | Time (m:ss) |
| 5 | 10% | 5.37e-2 | 5.04e-2 | 22 | 19 | 6 | 13 | 0:09 | 5.36e-2 | 5.05e-2 | 0:14 | 5.47e-2 | 5.12e-2 | 0:13 |
| 5 | 20% | 6.93e-2 | 6.62e-2 | 22 | 19 | 6 | 13 | 0:09 | 6.92e-2 | 6.63e-2 | 0:15 | 7.03e-2 | 6.68e-2 | 0:12 |
| 10 | 10% | 6.52e-2 | 6.47e-2 | 32 | 25 | 6 | 19 | 0:14 | 6.50e-2 | 6.44e-2 | 0:15 | 6.56e-2 | 6.48e-2 | 0:12 |
| 10 | 20% | 7.95e-2 | 8.08e-2 | 35 | 27 | 6 | 20 | 0:15 | 7.93e-2 | 8.07e-2 | 0:14 | 8.04e-2 | 8.15e-2 | 0:13 |
| 15 | 10% | 7.36e-2 | 7.98e-2 | 40 | 29 | 6 | 22 | 0:17 | 7.34e-2 | 7.96e-2 | 0:15 | 7.63e-2 | 8.35e-2 | 0:12 |
| 15 | 20% | 8.76e-2 | 9.49e-2 | 41 | 30 | 6 | 23 | 0:17 | 8.74e-2 | 9.47e-2 | 0:15 | 9.07e-2 | 9.93e-2 | 0:12 |

of SGD and occasionally even finding a better local minimum than Adam. Furthermore, the computation time of our methodology for training the introduced networks is reasonably similar to the ones of Adam and SGD. Overall, these results are encouraging as our method is expected to greatly benefit from a mini-batch implementation, further reducing the computation time and increasing scalability.

VI. CONCLUSIONS

In this paper a novel procedure for the training of neural networks was introduced that leverages an optimal control view and relies on three main components. First, a novel augmented Lagrangian method is presented for general non-smooth nonconvex equality constrained problems, which attains an ε -KKT solution. Second, when applied to the DNN problem we propose to solve the Lagrangian subproblems by employing Gauss-Newton iterations resulting in a series of linear least-squares problems. Third, owing to the stagewise structure in the optimal control formulation, we solve the linear least-squares GN problems through a simple recursive procedure based on forward dynamic programming. We observed encouraging results in comparison to fast first-order solvers such as Adam which are often used in a heuristic manner without theoretical guarantees. In the current implementation, our method is not competitive when using large numbers of training data. Future research directions include extending our scheme to mini-batch settings to tackle this issue. It is also interesting to extend the framework to allow for nonsmooth activations functions.

REFERENCES

- [1] M. Abadi, A. Agarwal, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv:1603.04467*, 2016.
- [2] D. Bertsekas. *Nonlinear Programming*. Athena Scientific, 2016.
- [3] E. Birgin and J. Martínez. *Practical Augmented Lagrangian Methods for Constrained Optimization*. SIAM, 2014.
- [4] M. Carreira-Perpinan and W. Wang. Distributed optimization of deeply nested systems. In *Artif. Intell. Stat.*, pages 10–19, 2014.
- [5] Y. Chen, T. Davis, et al. Algorithm 887: CHOLMOD, supernodal sparse Cholesky factorization and update/downdate. *ACM Trans Math Softw*, 35(3), oct 2008.
- [6] Y. Cui, Z. He, and J. Pang. Multicomposite nonconvex optimization for training deep neural networks. *SIAM J. Optim.*, 30(2):1693–1723, 2020.
- [7] F. Facchinei and J. Pang. *Finite-dimensional variational inequalities and complementarity problems*, volume II. Springer, 2003.
- [8] I. Goodfellow, D. Warde-Farley, et al. Maxout networks. In *Int. Conf. Mach. Learn.*, pages 1319–1327, 2013.
- [9] G. Grapiglia and Y. Yuan. On the complexity of an augmented Lagrangian method for nonconvex optimization. *IMA J. Numer. Anal.*, jul 2020.
- [10] F. Gu, A. Askari, and L. El Ghaoui. Fenchel lifted networks: A Lagrange relaxation of neural network training. In *Int. Conf. Artif. Intell. Stat.*, pages 3362–3371, 2020.
- [11] K. He, X. Zhang, et al. Deep residual learning for image recognition. In *2016 IEEE Conf. Comput. Vision Pattern Recogn. (CVPR)*, pages 770–778, Los Alamitos, CA, USA, jun 2016. IEEE Computer Society.
- [12] S. Hochreiter, Y. Bengio, et al. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies, 2001.
- [13] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proc. 32nd Int. Conf. Mach. Learn.*, volume 37 of *Proc. Mach. Learn. Res.*, pages 448–456, Lille, France, 07–09 Jul 2015. PMLR.
- [14] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014.
- [15] A. Krogh and J. Hertz. A simple weight decay can improve generalization. In *Proc. 4th Int. Conf. NIPS*, page 950–957. Morgan Kaufmann Publishers Inc., 1991.
- [16] Y. LeCun. *A theoretical framework for back-propagation*. IEEE Computer Society Press, 1992.
- [17] A. Maas, A. Hannun, and A. Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc ICML*, volume 30, page 3, 2013.
- [18] R. Rockafellar and R. Wets. *Variational Analysis*, volume 317. Springer, 2009.
- [19] D. Rumelhart, G. Hinton, and R. Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, 1986.
- [20] N. Srivastava, G. Hinton, et al. Dropout: a simple way to prevent neural networks from overfitting. *JMLR*, 15(1):1929–1958, 2014.
- [21] G. Taylor, R. Burmeister, et al. Training neural networks without gradients: A scalable ADMM approach. In *Int. Conf. Mach. Learn.*, pages 2722–2731, 2016.
- [22] P. Virtanen, R. Gommers, et al. SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nature Methods*, 17:261–272, 2020.
- [23] J. Wang, F. Yu, et al. ADMM for efficient deep learning with global convergence. In *Proc. 25th ACM SIGKDD Int. Conf. Knowledge Discov. Data Min.*, pages 111–119, 2019.
- [24] Z. Zhang and M. Brand. On the convergence of block coordinate descent in training dnns with Tikhonov regularization. In *Adv. Neural Inf. Process. Syst.*, pages 1719–1728, 2017.
- [25] Z. Zhang, Y. Chen, and V. Saligrama. Efficient training of very deep neural networks for supervised hashing. In *Proc. IEEE Conf. Comput. Vision Pattern Recogn.*, pages 1487–1495, 2016.

Lemma A.1. Suppose that $G : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is locally Lipschitz around a point \bar{z} at which $G(\bar{z}) = 0$. Then, $\phi(z) := \frac{1}{2}\|G(z)\|^2$ is strictly differentiable at \bar{z} (in the sense of [18, Def. 9.17]) with null gradient.

Proof. Let L be a Lipschitz constant for G in a neighborhood \mathcal{U} of \bar{z} . Then, for $z, z' \in \mathcal{U}$ we have

$$\begin{aligned} \frac{|\phi(z) - \phi(z') - \langle 0, z - z' \rangle|}{\|z - z'\|} &= \frac{\|G(z)\|^2 - \|G(z')\|^2}{2\|z - z'\|} \\ &= \frac{\|G(z) - G(z')\|^2 + 2\langle G(z'), G(z) - G(z') \rangle}{2\|z - z'\|} \\ &\leq \frac{L}{2}\|G(z) - G(z')\| + L\|G(z')\| \end{aligned}$$

which vanishes as $z, z' \rightarrow \bar{z}$, hence the claim. \square

Lemma A.2. Let $\mathcal{H}_i \in \mathbb{R}^{r_i \times r_i}$ be symmetric positive definite, $\mathcal{V}_i \in \mathbb{R}^{r_i \times p}$, and $v_i \in \mathbb{R}^{r_i}$, $i \in [N]$. If $\mathcal{U} := \sum_{i=1}^N \mathcal{V}_i^\top \mathcal{H}_i^{-1} \mathcal{V}_i$ is symmetric positive definite, then

$$\sum_{i=1}^N \|\mathcal{V}_i x - v_i\|_{\mathcal{H}_i^{-1}}^2 = \|\mathcal{U}x - d\|_{\mathcal{U}^{-1}}^2 - \|d\|_{\mathcal{U}^{-1}}^2 + \sum_{i=1}^N \|v_i\|_{\mathcal{H}_i^{-1}}^2,$$

where $d := \sum_{i=1}^N \mathcal{V}_i^\top \mathcal{H}_i^{-1} v_i$.

Proof. Let $q(x) = \sum_{i=1}^N \|\mathcal{V}_i x - v_i\|_{\mathcal{H}_i^{-1}}^2$. That $q(x) = \|x\|_{\mathcal{U}}^2 + \sum_{i=1}^N \|v_i\|_{\mathcal{H}_i^{-1}}^2 - 2\langle x, d \rangle$ is of immediate verification. Since q is quadratic, the Taylor expansion around its minimizer $x^* = \mathcal{U}^{-1}d$ is given by $q(x) = q(x^*) + \|x - x^*\|_{\mathcal{U}}^2$. Substituting x^* results in the claimed form. \square

Proof of Theorem 1. Owing to the update at step 1.3,

$$\mathcal{L}(z, \lambda^{k+1}) = \mathcal{L}_{\beta_k}(z, \lambda^k) + \frac{\beta_k}{2}\|F(z^{k+1})\|^2 - \frac{\beta_k}{2}\|F(z) - F(z^{k+1})\|^2.$$

The last term on the right-hand side is continuously differentiable (with null gradient) at z^{k+1} , owing to Lemma A.1. It then follows from [18, Ex. 8.8(c)] that $\partial_z \mathcal{L}(z^{k+1}, \lambda^{k+1}) = \partial_z \mathcal{L}_{\beta_k}(z^{k+1}, \lambda^k)$, hence that the pair (z^k, λ^k) satisfies condition (7a) for every $k \geq 1$, by virtue of Eq. (9) in step 1.2. It remains to show that (7b) too is eventually satisfied. Notice that, by definition of \hat{z}^k at step 1.1, $\mathcal{L}(z^{k+1}, \lambda^k) \leq f(z^0)$ holds for every k , which combined with (6) yields

$$\begin{aligned} \frac{1}{2\beta_k}\|\lambda^{k+1}\|^2 &= \frac{\beta_k}{2}\|F(z^{k+1}) + \lambda^k/\beta_k\|^2 \leq f(z^0) - f(z^{k+1}) + \frac{1}{2\beta_k}\|\lambda^k\|^2 \\ &\leq c + \frac{1}{2\beta_k}\|\lambda^k\|^2, \end{aligned}$$

where $c := f(z^0) - \inf f$ is a constant. Since $\beta_{k+1} \geq \beta_k$, it holds that $\frac{1}{2\beta_{k+1}}\|\lambda^{k+1}\|^2 \leq c + \frac{1}{2\beta_k}\|\lambda^k\|^2$, which leads to

$$\frac{1}{\beta_k}\|\lambda^k\|^2 \leq \frac{1}{\beta_0}\|\lambda^0\|^2 + 2kc \quad (\text{A.1a})$$

for every $k \in \mathbb{N}$. Moreover, since

$$\begin{aligned} \frac{1}{2}\|F(z^{k+1})\|^2 &\leq \|F(z^{k+1}) + \lambda^k/\beta_k\|^2 + \|\lambda^k/\beta_k\|^2 \\ &= \frac{1}{\beta_k}[\|\lambda^{k+1}\|^2 + \|\lambda^k\|^2], \end{aligned} \quad (\text{A.1b})$$

the β -update at step 1.5 implies that $\|F(z^{k+1})\|_\infty \rightarrow 0$ (\mathcal{Q} -linearly) if β_k is asymptotically constant, hence the claim. Otherwise, the set $\mathcal{K} := \{k \in \mathbb{N} \mid \beta_k = \max\{\xi\beta_{k-1}, \beta_0 k^\alpha\}\}$ is infinite. Then, for $k \in \mathcal{K}$, combining (A.1) yields

$$\begin{aligned} \frac{1}{2}\|F(z^{k+1})\|^2 &\leq \frac{\beta_{k+1}}{\beta_k^2} \left(\frac{1}{\beta_0}\|\lambda^0\|^2 + 2(k+1)c \right) + \frac{1}{\beta_k} \left(\frac{1}{\beta_0}\|\lambda^0\|^2 + 2kc \right) \\ &\leq \max \left\{ \frac{\xi}{\beta_k}, \frac{\beta_0(k+1)^\alpha}{\beta_k^2} \right\} \left(\frac{1}{\beta_0}\|\lambda^0\|^2 + 2(k+1)c \right) \end{aligned}$$

$$\begin{aligned} &+ \frac{1}{\beta_k} \left(\frac{1}{\beta_0}\|\lambda^0\|^2 + 2kc \right) \\ (\beta_k \geq k^\alpha) &\leq \max \left\{ \frac{\xi}{k^\alpha}, \frac{\beta_0(k+1)^\alpha}{k^{2\alpha}} \right\} \left(\frac{1}{\beta_0}\|\lambda^0\|^2 + 2(k+1)c \right) \\ &+ \frac{1}{k^\alpha} \left(\frac{1}{\beta_0}\|\lambda^0\|^2 + 2kc \right) \rightarrow 0 \end{aligned}$$

as $\mathcal{K} \ni k \rightarrow \infty$, owing to the fact that $\alpha > 1$. The second inequality uses the fact that, regardless of whether $k+1 \in \mathcal{K}$ or not, $\beta_{k+1} \leq \max\{\xi\beta_k, \beta_0(k+1)^\alpha\}$ holds (since $\xi > 1$). \square

Proof of Lemma 3. Note that matrices A_j , B_j and vectors c_j in Algorithm 2 depend on the current iterate z^l . Here, we use superscript l to emphasize this dependence. The linear least-squares Problem III solved at step 2.1 may equivalently be written as

$$\underset{z=(w,x)}{\text{minimize}} \quad \frac{1}{2}\|J(z^l)z - b(z^l)\|^2, \quad (\text{A.2})$$

where

$$J(z^l) = \begin{pmatrix} -\bar{B}^l & \bar{A}^l \\ \sqrt{\mu_w} \mathbf{I} & \end{pmatrix}$$

with

$$\bar{A} = \begin{pmatrix} \sqrt{\rho_1} \mathbf{I} & & & & \\ -\sqrt{\rho_2} A_2^l & \sqrt{\rho_2} \mathbf{I} & & & \\ & & \ddots & & \\ & & & -\sqrt{\rho_N} A_N^l & \sqrt{\rho_N} \mathbf{I} \\ & & & & -\sqrt{\rho_{N+1}} A_{N+1}^l \end{pmatrix},$$

$$\bar{B}^l = \text{blkdiag}(\sqrt{\rho_1} B_1^l, \dots, \sqrt{\rho_{N+1}} B_{N+1}^l),$$

and $b(z^l) = (\sqrt{\rho_1} c_1, \dots, \sqrt{\rho_{N+1}} c_{N+1}, 0_s)$.

In what follows we show that the eigenvalues of $J(z^l)^\top J(z^l)$ along a converging subsequence are bounded above and away from zero, and that step 2.2 is a restatement of the standard Armijo line search, at which point the claimed convergence results follow by standard arguments for gradient methods [2, 7]. For the latter, note that by the optimality conditions for (A.2), the solution \bar{z}^l satisfies $J(z^l)^\top J(z^l) \bar{z}^l = J(z^l)^\top b(z^l)$. Moreover, since $\nabla_z \mathcal{L}_\beta(z^l, \lambda) = J(z^l)^\top J(z^l) z^l - J(z^l)^\top b(z^l)$, combining the two equalities yields $\langle \nabla_z \mathcal{L}_\beta(z^l, \lambda), \bar{z}^l - z^l \rangle = -\|J(z^l)(\bar{z}^l - z^l)\|^2 = -2\mathcal{G}_{\text{st}, \mathcal{B}, 0}(p^l)$ establishing the claimed equivalence.

Suppose $(z^l)_{l \in \mathcal{K}}$ is a subsequence converging to a limit point z^* . Note that $J(z)^\top J(z)$ is nonsingular since $J(z)$ has full column rank for any z . Therefore, by continuity of $J(\cdot)$ and [7, Lem. 7.5.2] we have that $c_1 \mathbf{I} \leq J(z^l)^\top J(z^l) \leq c_2 \mathbf{I}$ for some $c_1, c_2 > 0$. The stationarity of z^* follows from [7, Prop. 8.3.7]. If $(z^l)_{l \in \mathcal{N}}$ remains bounded at least one converging subsequence exists establishing the claim. \square

Proof of Lemma 5. First, note that by (14)

$$\begin{aligned} G_j^\top G_j + \frac{\mu_w}{\rho_j} E_j^\top E_j &= \mathbf{I} - B_j E_j - E_j^\top B_j^\top + E_j^\top \left(B_j^\top B_j + \frac{\mu_w}{\rho_j} \mathbf{I} \right) E_j \\ &= \mathbf{I} - B_j E_j = G_j. \end{aligned} \quad (\text{A.3})$$

We proceed by induction to show that, for $j \in [N]$,

$$V_j^*(x_j) = \frac{1}{2}\|M_j^{-1} x_j - q_j\|_{M_j}^2 + C_j, \quad (\text{A.4})$$

where the term C_j does not depend on x_j, x_{j+1}, \dots, x_N . Here, we avoid deriving a recursion for C_j since it does not affect the computation of x_j and w_{j+1} in the next stages.

For the base case $j = 1$, by the first-order optimality condition for the minimization (11) the unique minimizer is computed as

$$\mathbf{w}_1^*(x_1) = E_1(\mathbf{x}_1 - A_1\mathbf{x}_0 - \mathbf{c}_1).$$

After substitution, by using (A.3) and simple algebra we obtain

$$V_1^*(\mathbf{x}_1) = \frac{1}{2}\|\mathbf{x}_1 - A_1\mathbf{x}_0 - \mathbf{c}_1\|_{\rho_1 G_1}^2 = \frac{1}{2}\|M_1^{-1}\mathbf{x}_1 - \mathbf{q}_1\|_{M_1}^2,$$

where $M_1 = \rho_1^{-1}G_1^{-1}$ and $\mathbf{q}_1 = M_1^{-1}(A_1\mathbf{x}_0 + \mathbf{c}_1)$.

Arguing by induction, suppose that (A.4) holds for some j such that $1 \leq j \leq N-1$. Let $\varphi(\mathbf{x}_j, \mathbf{w}_{j+1})$ denote the argument being minimized in (12). From direct computation

$$\nabla^2 \varphi(\mathbf{x}_j, \mathbf{w}_{j+1}) = \begin{pmatrix} M_j^{-1} + \rho_{j+1}A_{j+1}^\top A_{j+1} & \rho_{j+1}A_{j+1}^\top B_{j+1} \\ \rho_{j+1}B_{j+1}^\top A_{j+1} & \mu_w \mathbf{I} + \rho_{j+1}B_{j+1}^\top B_{j+1} \end{pmatrix}.$$

Since $M_j \in \mathbb{S}_{++}^r$, by forming its Schur complement and using (14) it follows that the Hessian is symmetric positive definite if and only if so is $M_j^{-1} + \rho_{j+1}A_{j+1}^\top G_{j+1}A_{j+1}$, which holds true. Hence, the subproblems have unique solutions. By the first-order optimality condition for (12), the solution pair $(\mathbf{x}_j^*, \mathbf{w}_{j+1}^*)$ satisfies

$$0 = M_j^{-1}\mathbf{x}_j^* - \mathbf{q}_j - \rho_{j+1}A_{j+1}^\top(\mathbf{x}_{j+1} - A_{j+1}\mathbf{x}_j^* - B_{j+1}\mathbf{w}_{j+1}^* - \mathbf{c}_{j+1})$$

and

$$0 = \mu_w \mathbf{w}_{j+1}^* - \rho_{j+1}B_{j+1}^\top(\mathbf{x}_{j+1} - A_{j+1}\mathbf{x}_j^* - B_{j+1}\mathbf{w}_{j+1}^* - \mathbf{c}_{j+1}).$$

The latter reads $\mathbf{w}_{j+1}^* = E_{j+1}(\mathbf{x}_{j+1} - A_{j+1}\mathbf{x}_j^* - \mathbf{c}_{j+1})$. After substituting \mathbf{w}_{j+1}^* into the former, using (14) and (A.3) we obtain

$$\mathbf{x}_j^* = S_{j+1}\mathbf{q}_j + P_j(\mathbf{x}_{j+1} - \mathbf{c}_{j+1}),$$

where $P_k = \rho_{k+1}S_{k+1}A_{k+1}^\top G_{k+1}$ and S_{k+1} is as in (17). Substituting the minimizer pair $(\mathbf{x}_j^*, \mathbf{w}_{j+1}^*)$ back in (12) and using (A.3) yields

$$\begin{aligned} V_{j+1}^*(\mathbf{x}_{j+1}) &= C_j + \frac{1}{2}\|M_j^{-1}\mathbf{x}_j^* - \mathbf{q}_j\|_{M_j}^2 \\ &\quad + \frac{1}{2}\|\mathbf{x}_{j+1} - A_{j+1}\mathbf{x}_j^* - \mathbf{c}_{j+1}\|_{\rho_{j+1}G_{j+1}}^2 \\ (\text{subs. } \mathbf{x}_j^*) &= C_j + \frac{1}{2}\|\mathcal{V}_1\mathbf{x}_{j+1} - \mathbf{v}_1\|_{\mathcal{H}_1}^2 \\ &\quad + \frac{1}{2}\|\mathcal{V}_2\mathbf{x}_{j+1} - \mathbf{v}_2\|_{\mathcal{H}_2}^2 \end{aligned}$$

where $\mathcal{V}_1 = M_j^{-1}P_j$, $\mathcal{H}_1 = M_j^{-1}$, $\mathcal{V}_2 = (\mathbf{I} - A_{j+1}P_j)$, $\mathcal{H}_2 = \rho_{j+1}^{-1}G_{j+1}^{-1}$,

$$\mathbf{v}_1 = (\mathbf{I} - M_j^{-1}S_{j+1})\mathbf{q}_j + M_j^{-1}P_j\mathbf{c}_{j+1}, \text{ and}$$

$$\mathbf{v}_2 = A_{j+1}S_{j+1}\mathbf{q}_j + (\mathbf{I} - A_{j+1}P_j)\mathbf{c}_{j+1}.$$

On the other hand, we have that

$$\begin{aligned} \mathcal{U} &= \sum_{i=1}^2 \mathcal{V}_i^\top \mathcal{H}_i^{-1} \mathcal{V}_i \\ &= P_j^\top M_j^{-1} P_j + \rho_{j+1} (\mathbf{I} - A_{j+1} P_j)^\top G_{j+1} (\mathbf{I} - A_{j+1} P_j) \quad (\text{A.5}) \\ &= P_j^\top S_{j+1}^{-1} P_j + \rho_{j+1} (\mathbf{I} - P_j^\top A_{j+1}^\top) G_{j+1} - \rho_{j+1} G_{j+1} A_{j+1} P_j \\ &= \rho_{j+1} G_{j+1} - \rho_{j+1}^2 G_{j+1} A_{j+1} S_{j+1} A_{j+1}^\top G_{j+1} \\ &= M_{j+1}^{-1}, \end{aligned}$$

where (19) was used in the second equality, and the Woodbury matrix identity in the last one. Therefore, we may apply Lemma A.2 to obtain

$$V_{j+1}^*(\mathbf{x}_{j+1}) = \frac{1}{2}\|M_{j+1}^{-1}\mathbf{x}_{j+1} - \mathbf{q}_{j+1}\|_{M_{j+1}}^2$$

$$- \frac{1}{2}\|\mathbf{q}_{j+1}\|_{M_{j+1}}^2 + \frac{1}{2}\sum_{i=1}^2 \|\mathbf{v}_i\|_{\mathcal{H}_i}^2 + C_j, \quad (\text{A.6})$$

with

$$\mathbf{q}_{j+1} = P_j^\top \mathbf{v}_1 + \rho_{j+1} (\mathbf{I} - A_{j+1} P_j)^\top G_{j+1} \mathbf{v}_2 \quad (\text{A.7})$$

$$= M_{j+1}^{-1} \mathbf{c}_{j+1} + \rho_{j+1} G_{j+1} A_{j+1} \mathbf{q}_j, \quad (\text{A.8})$$

where we used (19) and the alternative expression for M_{j+1} in (A.5). The last three terms in (A.6) are absorbed into C_{j+1} completing the induction argument.

It remains to solve (13). Arguing as before, from the first-order optimality condition the solution pair $(\mathbf{x}_N^*, \mathbf{w}_{N+1}^*)$ must satisfy

$$0 = M_N^{-1}\mathbf{x}_N^* - \mathbf{q}_N + \rho_{N+1}A_{N+1}^\top(A_{N+1}\mathbf{x}_N^* + B_{N+1}\mathbf{w}_{N+1}^* + \mathbf{c}_{N+1})$$

and

$$0 = \mu_w \mathbf{w}_{N+1}^* + \rho_{N+1}B_{N+1}^\top(A_{N+1}\mathbf{x}_N^* + B_{N+1}\mathbf{w}_{N+1}^* + \mathbf{c}_{N+1}).$$

The former equality is equivalent to the one given in step 3.2. After substituting \mathbf{w}_{N+1}^* back into the latter and using (A.3), the update for \mathbf{x}_N^* is obtained. \square