# Maximum Principle Based Algorithms for Deep Learning

Qianxiao Li LIQIX@IHPC.A-STAR.EDU.SG

Institute of High Performance Computing Agency for Science, Technology and Research 1 Fusionopolis Way, Connexis North, Singapore 138632

Long Chen XIDONGLC@PKU.EDU.CN

Peking University Beijing, China, 100080

Cheng Tai CHENGTAI@PKU.EDU.CN

Beijing Institute of Big Data Research and Peking University Beijing, China, 100080

Weinan E WEINAN@MATH.PRINCETON.EDU

Princeton University
Princeton, NJ 08544, USA,
Beijing Institute of Big Data Research
and Peking University
Beijing, China, 100080

Editor: Yoshua Bengio

#### Abstract

The continuous dynamical system approach to deep learning is explored in order to devise alternative frameworks for training algorithms. Training is recast as a control problem and this allows us to formulate necessary optimality conditions in continuous time using the Pontryagin's maximum principle (PMP). A modification of the method of successive approximations is then used to solve the PMP, giving rise to an alternative training algorithm for deep learning. This approach has the advantage that rigorous error estimates and convergence results can be established. We also show that it may avoid some pitfalls of gradient-based methods, such as slow convergence on flat landscapes near saddle points. Furthermore, we demonstrate that it obtains favorable initial convergence rate periteration, provided Hamiltonian maximization can be efficiently carried out - a step which is still in need of improvement. Overall, the approach opens up new avenues to attack problems associated with deep learning, such as trapping in slow manifolds and inapplicability of gradient-based methods for discrete trainable variables.

**Keywords:** deep learning, optimal control, Pontryagin's maximum principle, method of successive approximations

# 1. Introduction

Supervised learning using deep neural networks has become an increasingly successful tool in modern machine learning applications (Bengio, 2009; Schmidhuber, 2015; LeCun et al.,

©2018 Qianxiao Li, Long Chen, Cheng Tai, Weinan E.

License: CC-BY 4.0, see https://creativecommons.org/licenses/by/4.0/. Attribution requirements are provided at http://jmlr.org/papers/v18/17-653.html.

2015; Goodfellow et al., 2016). Efficient training methods of very deep neural networks, however, remain an active area of research. The most commonly applied training method is stochastic gradient descent (Robbins and Monro, 1951; Bottou, 2010) and its variants (Duchi et al., 2011; Zeiler, 2012; Kingma and Ba, 2014; Johnson and Zhang, 2013), where incremental updates to the trainable parameters are performed using gradient information computed via back-propagation (Kelley, 1960; Bryson, 1975). While efficient to implement, the incremental updates to the parameter tend to be slow, especially in the initial stages of the training. Moreover, other than the computation of gradients through back-propagation, the specific structure of deep neural networks is not exploited. These observations point to the question of whether there exists alternative training methods tailored to deep neural networks.

In a series of papers, we introduce an alternative approach by exploring the optimal control viewpoint of deep learning (E, 2017). Our focus will be on ideas and algorithms derived from the powerful Pontryagin's maximum principle (Boltyanskii et al., 1960; Pontryagin, 1987), which has two major components: the Hamiltonian dynamics and the condition that at each time the optimal parameters maximize the Hamiltonian. The second component suggests that optimization can be performed independently at different layers. One can also derive an explicit error control estimate based on the maximum principle (see Lemma 2 below).

In this first paper, we will consider the simplest context in which the deep neural networks are replaced by continuous (or discretized) dynamical systems, and devise numerical algorithms that are based on the optimality conditions in the Pontryagin's maximum principle. This leads to a new approach for training deep learning models that have certain advantages, such as fast initial descent and resilience to stalling in flat landscapes. An additional advantage is that one has a good control of the error through explicit estimates.

The rest of the paper is organized as follows. In Section 2, we present a dynamical systems viewpoint of function approximation and deep learning. We then discuss the necessary optimality conditions, which is the well-known Pontryagin's maximum principle. In Section 3 and 4, we discuss numerical methods to solve the necessary conditions and obtain error estimates and convergence guarantees. Using benchmarking examples, we then compare our method with traditional gradient-descent based methods for optimizing deep neural networks in Section 5. In Section 6, we discuss and compare our work with existing literature. Conclusion and outlook are given in Section 7.

# 2. Function Approximation by Dynamical Systems

We start with a description of the (continuous) dynamical systems approach to machine learning (see E 2017). The essential task of supervised learning is to approximate some function

$$F: \mathcal{X} \to \mathcal{Y}$$

which maps inputs in  $\mathcal{X} \subset \mathbb{R}^d$  (e.g. images, time-series) to labels in  $\mathcal{Y}$  (categories, numerical predictions). Given a collection of K sample input-label pairs  $\{x^i, y^i = F(x^i)\}_{i=1}^K$ , one aims to approximate F using these data points. In the dynamical systems framework, we consider the inputs  $x = (x^1, \dots, x^K) \in \mathbb{R}^{d \times K}$  as the initial condition of a system of ordinary

differential equations

$$\dot{X}_t^i = f(t, X_t^i, \theta_t), \quad X_0^i = x^i, \quad 0 \le t \le T,$$
 (1)

where  $\theta:[0,T]\to\Theta\subset\mathbb{R}^p$ , represents the control (training) parameters and  $X_t=(X_t^1,\ldots,X_t^K)\in\mathbb{R}^{d\times K}$  for all  $t\in[0,T]$ . The form of f is chosen as part of the machine learning model. For example, in deep learning, f is typically the composition (in either order) of a linear transformation and a component-wise nonlinear function (the activation function). For the solution to (1) to exist for any  $\theta$ , we shall assume hereafter that f and  $\nabla_x f$  are continuous in  $t, x, \theta$ . Note that weaker but more complicated conditions can be considered (Clarke, 2005). For the  $i^{\text{th}}$  input sample, the prediction of the "network" is a deterministic transformation of the terminal state  $g(X_T^i)$  for some  $g:\mathbb{R}^d\to\mathcal{Y}$ , which we can view collectively as a function of the initial state (input)  $x^i$  and the control parameters (weights)  $\theta$ . The dynamics (1) are decoupled across samples except for the dependence on the control  $\theta$ . We shall consider quite a general space of controls

$$\mathcal{U} := \{\theta : [0, T] \to \Theta : \theta \text{ is Lebesgue measurable}\}.$$

The aim is to select  $\theta$  from  $\mathcal{U}$  so that  $g(X_T^i)$  most closely resembles  $y^i$  for  $i=1,\ldots,K$ . To this end, we define a loss function  $\Phi: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$  which is minimized when its arguments are equal, and we consider minimizing  $\sum_i \Phi(g(x_T^i), y^i)$ . Since g is fixed, we shall absorb it into the definition of the loss function by defining  $\Phi_i(\cdot) := \Phi(\cdot, y_i)$ . Then, the supervised learning problem in our framework is

$$\min_{\theta \in \mathcal{U}} \sum_{i=1}^{K} \Phi_i(X_T^i) + \int_0^T L(\theta_t) dt, 
\dot{X}_t^i = f(t, X_t^i, \theta_t), \quad X_0^i = x^i, \quad 0 \le t \le T, \quad i = 1, \dots, K,$$
(2)

where  $L: \Theta \to \mathbb{R}$  is a running cost, or the regularizer<sup>1</sup>. We note here that alternatively, we can formulate the supervised learning problem more generally in terms of optimal control in function spaces, see Appendix A.

Problem (2) is a special case of a class of general optimal control problem for ordinary differential equations (Bertsekas, 1995; Athans and Falb, 2013). The advantage of this formulation is that we can write down and study the optimality conditions of (2) entirely in continuous time and derive numerical algorithms that can subsequently be discretized. In other words, we *optimize*, then discretize, as opposed to the traditional reverse approach in deep learning.

As was suggested in E (2017), deep residual networks (He et al., 2016) can be considered as the forward Euler discretization of the continuous approach described above. In this connection, the algorithms presented in this paper can also be formulated in the context of deep residual networks. For general deep neural networks, although one can also formulate similar algorithms, it is not clear at this moment that PMP holds and these algorithms are valid (e.g. converge to the right solution) in the general setting. This issue will be studied in future work.

<sup>1.</sup> We can also make L depend on  $X_t$ , but for simplicity of presentation and the fact that most current machine learning models do not regularize the states, we shall omit this general case.

The optimization problem (2) can be solved by first discretizing it into a discrete problem (a feed-forward neural network) and then applying back propagation and gradient descent approaches commonly used in deep learning. However, here we will present an alternative approach. Hereafter, for simplicity of notation we shall set K=1 drop the scripts i on all functions, noting that analogous results can be obtained in the general case since the dynamics and loss functions are decoupled across samples. Equivalently, we can think of this as effectively concatenating all K sample inputs into a single input vector of dimension  $d \times K$  and redefine our dynamics accordingly. Hence, all results remain valid if we perform full-batch training. The case of mini-batch training is discussed in Section 4.3.

#### 2.1 Pontryagin's Maximum Principle

In this section, we introduce a set of necessary conditions for optimal solutions of (2), known as the Pontryagin's Maximum Principle (PMP) (Boltyanskii et al., 1960; Pontryagin, 1987). This shall pave way for an alternative numerical algorithm to train (2) and its discrete-time counter-part.

To begin with, we define the Hamiltonian  $H:[0,T]\times\mathbb{R}^d\times\mathbb{R}^d\times\Theta\to\mathbb{R}$  given by

$$H(t, x, p, \theta) := p \cdot f(t, x, \theta) - L(\theta).$$

Theorem 1 (Pontryagin's Maximum Principle) Let  $\theta^* \in \mathcal{U}$  be an essentially bounded optimal control, i.e. a solution to (2) with  $ess\, sup_{t\in[0,T]}\|\theta_t^*\|_{\infty}<\infty$  (ess  $sup\ denotes\ the$ essential supremum). Denote by  $X^*$  the corresponding optimally controlled state process. Then, there exists an absolutely continuous co-state process  $P^*:[0,T]\to\mathbb{R}^d$  such that the Hamilton's equations

$$\dot{X}_{t}^{*} = \nabla_{p} H(t, X_{t}^{*}, P_{t}^{*}, \theta_{t}^{*}), \qquad X_{0}^{*} = x, 
\dot{P}_{t}^{*} = -\nabla_{x} H(t, X_{t}^{*}, P_{t}^{*}, \theta_{t}^{*}), \qquad P_{T}^{*} = -\nabla \Phi(X_{T}^{*}), \tag{4}$$

$$\dot{P}_t^* = -\nabla_x H(t, X_t^*, P_t^*, \theta_t^*), \qquad P_T^* = -\nabla \Phi(X_T^*), \tag{4}$$

are satisfied. Moreover, for each  $t \in [0,T]$ , we have the Hamiltonian maximization condition

$$H(t, X_t^*, P_t^*, \theta_t^*) \ge H(t, X_t^*, P_t^*, \theta) \text{ for all } \theta \in \Theta.$$
(5)

The proof of the PMP and its variants can be found in any optimal control theory reference, e.g. Athans and Falb (2013); Bertsekas (1995); Liberzon (2012). Some generalizations can be found in Clarke (2005) and references therein. For example, the requirement of the continuity of f with respect to t can be replaced by a much weaker measurability requirement if one assumes more conditions on  $\nabla_x f$ . In the statement of Theorem 1, we omitted a technicality involving an abnormal multiplier: the terminal condition for  $P^*$  should be  $P_T^* = -\lambda \nabla \Phi(X_T^*)$  and the Hamiltonian should be defined as  $H(t,x,p,\theta) = p \cdot f(t,x,\theta) - \lambda L(\theta)$  for some  $\lambda \geq 0$  (abnormal multiplier) that we can choose. When we are forced to always take  $\lambda = 0$ , the problem is singular and in a sense ill-posed (Athans and Falb, 2013). On the contrary, if we can take a positive  $\lambda$ , we can then rescale the equation for  $P^*$  so that we can take  $\lambda = 1$  without loss of generality. We shall hereafter assume that this is the case.

A few remarks are in order. First, Equation 3, 4 and 5 allow us to solve for the unknowns  $X^*, P^*, \theta^*$  simultaneously as a function of t. In this sense, the resulting optimal control  $\theta^*$  is open-loop and is not in a feed-back form  $\theta_t^* = \theta^*(X_t^*)$ . The latter is of closed-loop type and are typically obtained from dynamic programming and the Hamilton-Jacobi-Bellman formalism (Bellman, 2013). In this sense, the PMP gives a weaker control. However, open-loop solutions are sufficient for neural network applications, where the trained weights and biases are fixed and only depend on the layer number and not the inputs.

PMP can be regarded as a (highly non-trivial) generalization of the calculus of variations to non-smooth settings (since we only assume  $\theta^*$  to be measurable). Perhaps more familiar to the optimization community, the PMP is related to the Karush-Kuhn-Tucker (KKT) conditions for non-linear constrained optimization. Indeed, we can view (2) as a non-linear program over the function space  $\mathcal{U}$  where the constraint is the ODE (1). In this sense, the co-state process  $P^*$  plays the role of a continuous-time analogue of Lagrange multipliers. The key difference between the PMP and the KKT conditions (besides the lack of inequality constraints on the state) is the Hamiltonian maximization condition (5), which is stronger than a typical first-order condition that assumes smoothness with respect to  $\theta$  (e.g.  $\nabla_{\theta}H = 0$ ). In particular, the PMP says that H is not only stationary, but globally maximized at an optimal control - which is a much stronger statement if H is not concave. Moreover, the PMP makes minimal assumptions on the parameter space  $\Theta$ ; the PMP holds even when f is non-smooth with respect to  $\theta$ , or worse, when  $\Theta$  is a discrete subset of  $\mathbb{R}^p$ .

Last, we emphasize that the PMP is only a necessary condition, hence there can be cases where solutions to the PMP is not actually globally optimal for (2). Nevertheless, in practice the PMP is often strong enough to give good solution candidates, and when certain convexity assumptions are satisfied the PMP becomes sufficient (Bressan and Piccoli, 2007). In the next section, we will discuss numerical methods that can be used to solve the PMP.

## 3. Method of Successive Approximations

Now, our strategy is to devise numerical algorithms for training (2) via solving the PMP (Equation 3, 4 and 5). We derive and analyze algorithms entirely in continuous time, which allows us to characterize errors estimates and convergence in a more transparent fashion.

There are many methods for the numerical solution of the PMP, including two-point boundary value problem method (Bryson, 1975; Roberts and Shipman, 1972), and collocation methods (Betts, 1998) coupled with general non-linear programming techniques (Bertsekas, 1999; Bazaraa et al., 2013). See (Rao, 2009) for a more recent review. However, many of these methods concern small-scale problems typically encountered in control applications (e.g. trajectory optimization of spacecrafts) and do not scale well to modern machine learning problems with a large number of state and control variables. One exception is the method of successive approximations (MSA) (Chernousko and Lyubushin, 1982), which is an iterative method based on alternating propagation and optimization steps. We first introduce the simplest form of the MSA.

#### 3.1 Basic MSA

Observe that (3) is simply the equation

$$\dot{X}_t^* = f(t, X_t^*, \theta_t^*),$$

and is independent of the co-state  $P^*$ . Therefore, we may proceed in the following manner. First, we make an initial guess of the optimal control  $\theta^0 \in \mathcal{U}$ . For each  $k = 0, 1, 2, \ldots$ , we first solve (3)

$$\dot{X}_t^{\theta^k} = f(t, X_t^{\theta^k}, \theta_t^k), \quad X_0^{\theta^k} = x. \tag{6}$$

for  $X^{\theta^k}$ , which then allows us to solve (4)

$$\dot{P}_t^{\theta^k} = -\nabla_x H(t, X_t^{\theta^k}, P_t^{\theta^k}, \theta_t^k), \quad P_T^{\theta^k} = -\nabla \Phi(X_T^{\theta^k}), \tag{7}$$

to get  $P^{\theta^k}$ . Finally, we use the maximization condition (5) to set

$$\theta_t^{k+1} = \operatorname*{arg\,max}_{\theta \in \Theta} H(t, X_t^{\theta^k}, P_t^{\theta^k}, \theta),$$

for  $t \in [0, T]$ . The algorithm is summarized in Algorithm 1.

## Algorithm 1: Basic MSA

```
1 Initialize: \theta^{0} \in \mathcal{U};

2 for k = 0 to #Iterations do

3 | Solve \dot{X}_{t}^{\theta^{k}} = f(t, X_{t}^{\theta^{k}}, \theta_{t}^{k}), \quad X_{0}^{\theta^{k}} = x;

4 | Solve \dot{P}_{t}^{\theta^{k}} = -\nabla_{x}H(t, X_{t}^{\theta^{k}}, P_{t}^{\theta^{k}}, \theta_{t}^{k}), \quad P_{T}^{\theta^{k}} = -\nabla\Phi(X_{T}^{\theta^{k}});

5 | Set \theta_{t}^{k+1} = \arg\max_{\theta \in \Theta} H(t, X_{t}^{\theta^{k}}, P_{t}^{\theta^{k}}, \theta) for each t \in [0, T];

6 end
```

As is the case with the maximum principle, MSA consists of two major components: the forward-backward Hamiltonian dynamics and the maximization for the optimal parameters at each time. An important feature of MSA is that the Hamiltonian maximization step is decoupled for each  $t \in [0,T]$ . In the language of deep learning, the optimization step is decoupled for different network layers and only the Hamiltonian ODEs (Step 3,4 of Algorithm 1) involve propagation through the layers. This allows the parallelization of the maximization step, which is typically the most time-consuming step.

It has been shown that the basic MSA converges for a restricted class of linear quadratic regulators (Aleksandrov, 1968). However, in general it tends to diverge, especially if a bad initial  $\theta^0$  is chosen (Aleksandrov, 1968; Chernousko and Lyubushin, 1982). Our goal now is to modify the basic MSA to control its divergent behavior. Before we do so, it is important to understand why the MSA diverges, and in particular, the relationship between the maximization step in Algorithm 1 and the optimization problem (2).

#### 3.2 Error Estimate for the Basic MSA

For each  $\theta \in \mathcal{U}$ , let us denote

$$J(\theta) := \Phi(X_T^{\theta}) + \int_0^T L(\theta_t) dt,$$

where  $X^{\theta}$  satisfies (6). Our goal is to minimize  $J(\theta)$ . We show in the following Lemma the relationship between the values of J and the Hamiltonian maximization step. We start by making the following assumptions.

(A1)  $\Phi$  is twice continuously differentiable, with  $\Phi$  and  $\nabla \Phi$  satisfying a Lipschitz condition, i.e. there exists K > 0 such that

$$|\Phi(x) - \Phi(x')| + ||\nabla \Phi(x) - \nabla \Phi(x')|| \le K||x - x'||,$$

for all  $x, x' \in \mathbb{R}^d$ .

(A2)  $f(t, \cdot, \theta)$  is twice continuously differentiable in x, with  $f, \nabla_x f$  satisfying a Lipschitz condition in x uniformly in  $\theta$  and t, i.e. there exists K > 0 such that

$$||f(t,x,\theta) - f(t,x',\theta)|| + ||\nabla_x f(t,x,\theta) - \nabla_x f(t,x',\theta)||_2 \le K||x-x'||,$$

for all  $x, x' \in \mathbb{R}^d$  and  $t \in [0, T]$ . Note that  $\|\cdot\|_2$  denotes the induced 2-norm.

With these assumptions, we have the following estimate:

**Lemma 2** Suppose (A1)-(A2) holds. Then, there exists a constant C > 0 such that for any  $\theta, \phi \in \mathcal{U}$ ,

$$\begin{split} J(\phi) \leq &J(\theta) - \int_0^T \Delta_{\phi,\theta} H(t) dt \\ &+ C \int_0^T \|f(t,X_t^{\theta},\phi_t) - f(t,X_t^{\theta},\theta_t)\|^2 dt \\ &+ C \int_0^T \|\nabla_x H(t,X_t^{\theta},P_t^{\theta},\phi_t) - \nabla_x H(t,X_t^{\theta},P_t^{\theta},\theta_t)\|^2 dt, \end{split}$$

where  $X^{\theta}$ ,  $P^{\theta}$  satisfy Equations 6, 7 respectively and  $\Delta H_{\phi,\theta}$  denotes the change in Hamiltonian

$$\Delta H_{\phi,\theta}(t) := H(t, X_t^{\theta}, P_t^{\theta}, \phi_t) - H(t, X_t^{\theta}, P_t^{\theta}, \theta_t).$$

**Proof** See Appendix B for the proof and discussion on relaxing the assumptions.

In essence, Lemma 2 says that the Hamiltonian maximization step in MSA (step 5 in Algorithm 1) is in some sense the optimal descent direction for J. However, the last two terms on the right hand side indicates that this descent can be nullified if substituting  $\phi$  for  $\theta$  incurs too much error in the Hamiltonian dynamics (step 3,4 in Algorithm 1). In other words, the last two integrals measure the degree of satisfaction of the Hamiltonian dynamics (3), (4), which can be viewed as a feasibility condition, when one replaces  $\theta$  by  $\phi$ . Hence, we shall hereafter refer to these errors as feasibility errors. The divergence of the basic MSA happens when the feasibility errors blow up. Armed with this understanding, we can then modify the basic MSA to ensure convergence.

#### 3.3 Extended PMP and Extended MSA

As discussed previously in Lemma 2, the decrement of J is ensured if we can control the feasibility errors in the Hamiltonian dynamics in steps 3,4 of Algorithm 1. To this end, we employ a similar idea to augmented Lagrangians (Hestenes, 1969). Fix some  $\rho > 0$  and introduce the augmented Hamiltonian

$$\tilde{H}(t, x, p, \theta, v, q) := H(t, x, p, \theta) - \frac{1}{2}\rho \|v - f(t, x, \theta)\|^{2} - \frac{1}{2}\rho \|q + \nabla_{x}H(t, x, p, \theta)\|^{2}.$$
(8)

Then, we have the following set of alternative necessary conditions for optimality:

**Proposition 3 (Extended PMP)** Suppose that  $\theta^*$  is an essentially bounded solution to the optimal control problem (2). Then, there exists an absolutely continuous co-state process  $P^*$  such that the tuple  $(X_t^*, P_t^*, \theta_t^*)$  satisfies the necessary conditions

$$\dot{X}_{t}^{*} = \nabla_{p} \tilde{H}(t, X_{t}^{*}, P_{t}^{*}, \theta_{t}^{*}, \dot{X}_{t}^{*}, \dot{P}_{t}^{*}), \qquad X_{0}^{*} = x, \tag{9}$$

$$\dot{P}_t^* = -\nabla_x \tilde{H}(t, X_t^*, P_t^*, \theta_t^*, \dot{X}_t^*, \dot{P}_t^*), \qquad P_T^* = -\nabla_x \Phi(X_T^*), \qquad (10)$$

$$\tilde{H}(t, X_t^*, P_t^*, \theta_t^*, \dot{X}_t^*, \dot{P}_t^*) \ge \tilde{H}(t, X_t^*, P_t^*, \theta, \dot{X}_t^*, \dot{P}_t^*), \qquad \theta \in \Theta, t \in [0, T]. \tag{11}$$

**Proof** If  $\theta^*$  is optimal, then by the PMP there exists a co-state process  $P^*$  such that (3), (4) and (5) are satisfied. Then, for all  $t \in [0, T]$  and  $\theta \in \Theta$  we have

$$\nabla_{x} \tilde{H}(t, X_{t}^{*}, P_{t}^{*}, \theta, \dot{X}_{t}^{*}, \dot{P}_{t}^{*}) = \nabla_{x} H(t, X_{t}^{*}, P_{t}^{*}, \theta),$$

$$\nabla_{p} \tilde{H}(t, X_{t}^{*}, P_{t}^{*}, \theta, \dot{X}_{t}^{*}, \dot{P}_{t}^{*}) = \nabla_{p} H(t, X_{t}^{*}, P_{t}^{*}, \theta),$$

which implies that (9) and (10) are satisfied. Lastly, we can write

$$\begin{split} & \tilde{H}(t, X_t^*, P_t^*, \theta, \dot{X}_t^*, \dot{P}_t^*) \\ = & H(t, X_t^*, P_t^*, \theta) - \frac{1}{2} \rho \|\dot{X}_t^* - f(t, X_t^*, \theta)\|^2 - \frac{1}{2} \rho \|\dot{P}_t^* + \nabla_x H(t, X_t^*, P_t^*, \theta)\|^2. \end{split}$$

For each t,  $\theta^*$  maximizes all three terms on the RHS simultaneously, and hence (11) is also satisfied.

Compared with the usual PMP, the extended PMP is a weaker necessary condition. However, the advantage is that maximization of  $\tilde{H}$  naturally penalizes errors in the Hamiltonian dynamical equations, and hence we should expect MSA applied to the extended PMP to converge for large enough  $\rho$ . Note that the Hamiltonian equation steps do not change (since the added terms have no effect on optimal solutions) and the only change is the maximization step. The extended MSA (E-MSA) algorithm is summarized in Algorithm 2.

To establish convergence, define

$$\mu_k := \int_0^T \Delta H_{\theta^{k+1}, \theta^k}(t) dt \ge 0.$$

If  $\mu_k = 0$ , then from the Hamiltonian maximization step (11) we must have

$$0 = -\mu_k \le -\frac{1}{2}\rho \int_0^T \|f(t, X_t^{\theta^k}, \theta_t^{k+1}) - f(t, X_t^{\theta^k}, \theta_t^k)\|^2 dt$$
$$-\frac{1}{2}\rho \int_0^T \|\nabla_x H(t, X_t^{\theta^k}, P_t^{\theta^k}, \theta_t^{k+1}) - \nabla_x H(t, X_t^{\theta^k}, P_t^{\theta^k}, \theta_t^k)\|^2 dt. \le 0.$$

and so

$$\max_{\theta} \tilde{H}(X_t^{\theta^k}, P_t^{\theta^k}, \theta, \dot{X}_t^{\theta^k}, \dot{P}_t^{\theta^k}) = \tilde{H}(X_t^{\theta^k}, P_t^{\theta^k}, \theta_k, \dot{X}_t^{\theta^k}, \dot{P}_t^{\theta^k}).$$

i.e.  $(X^{\theta^k}, P^{\theta^k}, \theta^k)$  satisfy the extended PMP. In other words, the quantity  $\mu_k \geq 0$  measures the distance from a solution of the extended PMP, and if it equals 0, then we have a solution. We now prove the following result that guarantees the convergence of the extended MSA (Algorithm 2).

**Theorem 4** Let (A1)-(A2) be satisfied and  $\theta^0 \in \mathcal{U}$  be any initial measurable control with  $J(\theta^0) < +\infty$ . Suppose also that  $\inf_{\theta \in \mathcal{U}} J(\theta) > -\infty$ . Then, for  $\rho$  large enough, we have under Algorithm 2,

$$J(\theta^{k+1}) - J(\theta^k) \le -D\mu_k.$$

for some constant D > 0 and

$$\lim_{k \to 0} \mu_k = 0,$$

i.e. the extended MSA algorithm converges to the set of solutions of the extended PMP.

**Proof** Using Lemma 2 with  $\theta \equiv \theta^k$ ,  $\phi \equiv \theta^{k+1}$ , we have

$$J(\theta^{k+1}) - J(\theta^{k}) \le -\mu_{k} + C \int_{0}^{T} \|f(t, X_{t}^{\theta^{k}}, \theta_{t}^{k+1}) - f(t, X_{t}^{\theta^{k}}, \theta_{t}^{k})\|^{2} dt$$
$$+ C \int_{0}^{T} \|\nabla_{x} H(t, X_{t}^{\theta^{k}}, P_{t}^{\theta^{k}}, \theta_{t}^{k+1}) - \nabla_{x} H(t, X_{t}^{\theta^{k}}, P_{t}^{\theta^{k}}, \theta_{t}^{k})\|^{2} dt.$$

From the Hamiltonian maximization step in Algorithm 2, we know that

$$\begin{split} H(t, X_t^{\theta^k}, P_t^{\theta^k}, \theta_t^k) \leq & H(t, X_t^{\theta^k}, P_t^{\theta^k}, \theta_t^{k+1}) \\ & - \frac{1}{2} \rho \| f(t, X_t^{\theta^k}, \theta_t^{k+1}) - f(t, X_t^{\theta^k}, \theta_t^k) \|^2 \\ & - \frac{1}{2} \rho \| \nabla_x H(t, X_t^{\theta^k}, P_t^{\theta^k}, \theta_t^{k+1}) - \nabla_x H(t, X_t^{\theta^k}, P_t^{\theta^k}, \theta_t^k) \|^2. \end{split}$$

Hence, we have

$$J(\theta^{k+1}) - J(\theta^k) \le -(1 - \frac{2C}{\rho})\mu_k.$$

Pick  $\rho > 2C$ , then we indeed have  $J(\theta^{k+1}) - J(\theta^k) \leq -D\mu_k$  with  $D = (1 - \frac{2C}{\rho}) > 0$ . Moreover, we can rearrange and sum the above expression to get

$$\sum_{k=0}^{M} \mu_k \le D^{-1}(J(\theta^0) - J(\theta^{M+1})) \le D^{-1}(J(\theta^0) - \inf_{\theta \in \mathcal{U}} J(\theta)),$$

and hence  $\sum_{k=0}^{\infty} \mu_k < +\infty$ , which implies  $\mu_k \to 0$  and the extended MSA converges to a solution of the extended PMP.

# Algorithm 2: Extended MSA

```
1 Initialize: \theta^0 \in \mathcal{U}. Hyper-parameter: \rho;
2 for k=0 to #Iterations do
3 | Solve \dot{X}_t^{\theta^k} = f(t, X_t^{\theta^k}, \theta_t^k), \quad X_0^{\theta^k} = x;
4 | Solve \dot{P}_t^{\theta^k} = -\nabla_x H(t, X_t^{\theta^k}, P_t^{\theta^k}, \theta_t^k), \quad P_T^{\theta^k} = -\nabla \Phi(X_T^{\theta^k});
5 | Set \theta_t^{k+1} = \arg \max_{\theta \in \Theta} \tilde{H}(t, X_t^{\theta^k}, P_t^{\theta^k}, \theta, \dot{X}_t^{\theta^k}, \dot{P}_t^{\theta^k}) for each t \in [0, T];
6 end
```

# 4. Discrete-Time Formulation

In the previous section, we discussed the PMP and MSA in the continuous-time setting, where we showed that an appropriately extended version (E-MSA) converges to a solution of an extended PMP. Here, we shall discuss the discretized versions of PMP, MSA and E-MSA, as well as their connections to deep residual networks and back-propagation.

# 4.1 Discrete-Time PMP and Discrete-Time MSA

Applying Euler-discretization to Equation 1, we get

$$x_{n+1} = x_n + \delta f_n(x_n, \vartheta_n), \quad x_0 = x,$$

for n = 0, ..., N - 1, with  $\delta = T/N$  (step-size),  $x_n := X_{n\delta}$ ,  $\vartheta_n := \theta_{n\delta}$  and  $f_n(\cdot) := f(n\delta, \cdot)$ . Then, the discrete-time analogue of the control problem (2) is

$$\min_{\{\vartheta_0, \dots, \vartheta_{N-1}\} \in \Theta^N} \Phi(x_N) + \delta \sum_{n=0}^{N-1} L(\vartheta_n),$$

$$x_{n+1} = x_n + \delta f_n(x_n, \vartheta_n), \quad x_0 = x, \quad 0 \le n \le N - 1.$$
(12)

Observe that barring the constant  $\delta$ , this is exactly the supervised learning problem for deep residual networks<sup>2</sup>. Therefore, when suitably discretized, one expects that the E-MSA provides a means to train residual neural networks via the solution of the extended PMP.

We now write down formally the discretized form of the PMP. Let us use the shorthand  $g_n(x_n, \vartheta_n) := x_n + \delta f_n(x_n, \vartheta_n)$ . Define the scaled discrete Hamiltonian

$$H_n(x, p, \vartheta) = p \cdot g_n(x, \vartheta) - \delta L(\vartheta).$$

Then, a discrete-time PMP is the following set of conditions:

$$x_{n+1}^* = g_n(x_n^*, \vartheta_n^*), x_0^* = x,$$
  

$$p_n^* = \nabla_x H_n(x_n^*, p_{n+1}^*, \vartheta_n), p_N^* = -\nabla_x \Phi(x_N^*),$$
  

$$H_n(x_n^*, p_{n+1}^*, \vartheta_n^*) \ge H_n(x_n^*, p_{n+1}^*, \vartheta), \vartheta \in \Theta, n = 0, \dots, N-1.$$

<sup>2.</sup> If we pick ReLU activations (Hahnloser et al., 2000), then  $\delta$  can be absorbed into  $\vartheta$ 

The issue of whether the PMP holds for discrete time dynamical systems is a delicate one and there are known counterexamples (Butkovsky, 1963; Jackson and Horn, 1965; Nahorski et al., 1984). Nevertheless, they must hold approximately for small time step size and this is the situation we will consider in the current paper. We expect Lemma 2, which implies monotonicity of the E-MSA algorithm, to hold in the discrete-time case under appropriate conditions. We leave a rigorous analysis of these statements to future work. For numerical experiments presented in the next section, we shall almost always work with residual networks that can be regarded as discretizations of continuous networks so that the PMP holds approximately at least (Halkin, 1966).

For completeness, we summarize the discrete-time version of E-MSA in Algorithm 3. Note that for residual networks  $(g_n = x_n + \delta f_n)$ , this is equivalent to a forward Euler discretization on the state equation and a backward Euler discretization on the co-state equation in Algorithm 2. As before, the Hamiltonian maximization step is decoupled across layers and can be carried out in parallel.

# Algorithm 3: Discrete-time E-MSA

```
ı Initialize: Initialize: \vartheta_n^0 \in \Theta_n, n = 0, \dots, N-1. Hyper-parameter: \rho;
  2 for k = 0 to #Iterations do
               Set x_0^{\theta^k} = x;

for n = 0 to N - 1 do
 \begin{vmatrix} x_{n+1}^{\vartheta^k} = g_n(x_n^{\vartheta^k}, \vartheta_n^k) ; \\ \text{end} \end{vmatrix}
  3
  4
  \mathbf{5}
  6
                 Set p_N^{\vartheta^k} = -\nabla \Phi(x_N^{\vartheta^k}); for n = N-1 to 0 do
  7
  8
                  \begin{array}{c} p_n^{\vartheta^k} = \nabla_x H_n(x_n^{\vartheta^k}, p_{n+1}^{\vartheta^k}, \vartheta_n^k) \ ; \end{array}
  9
10
                  for n = 0 to N - 1 do
11
                          Set \vartheta_n^{k+1} = \arg\max_{\vartheta \in \Theta_n} H_n(x_n^{\vartheta^k}, p_{n+1}^{\vartheta^k}, \vartheta) - \frac{1}{2}\rho \|x_{n+1}^{\vartheta^k} - g_n(x_n^{\vartheta^k}, \vartheta)\|_2^2 - \frac{1}{2}\rho \|p_n^{\vartheta^k} - \nabla_x H_n(x_n^{\vartheta^k}, p_{n+1}^{\vartheta^k}, \vartheta)\|_2^2;
12
                  end
13
14 end
```

# 4.2 Relationship to Gradient Descent with Back-propagation

We note an interesting relationship of the MSA with classical gradient descent with backpropagation (Kelley, 1960; Bryson, 1975; LeCun et al., 1998). We have shown in Lemma 2 that the divergence of MSA can be attributed to the large errors in the Hamiltonian dynamics terms caused by the maximization step, which involve drastic changes in parameter values. Assuming each  $\Theta_n$  is a continuum and  $g_n$ , L are differentiable in  $\vartheta_n$ , a simple fix is to make the maximization step "soft": we replace step 12 in Algorithm 3 with a gradient ascent step:

$$\vartheta_n^{k+1} = \vartheta_n^k + \eta \nabla_{\vartheta} H_n(x_n^{\vartheta^k}, p_{n+1}^{\vartheta^k}, \vartheta_n^k), \tag{13}$$

for some small learning rate  $\eta$ . We now show that in the discrete-time setting, this is equivalent to the classical gradient descent with back-propagation.

**Proposition 5** The basic MSA in discrete-time (Algorithm 3 with  $\rho = 0$ ) with step 12 replaced by (13) is equivalent to gradient descent with back-propagation.

**Proof** Recall that the Hamiltonian is

$$H_n := p_{n+1} \cdot g_n(x_n, \vartheta_n) - \delta L(\vartheta_n),$$

and the total loss function is  $J(\vartheta) = \Phi(x_N) + \delta \sum_{n=0}^{N-1} L(\vartheta_n)$ . It is easy to see that  $p_n = -\nabla_{x_n} \Phi(x_N)$  by working backwards from n = N and the fact that  $\nabla_{x_n} x_{n+1} = \nabla_x g_n(x_n, \vartheta_n)$ . Then,

$$\begin{split} \nabla_{\vartheta_n} J(\vartheta) = & \nabla_{x_{n+1}} \Phi(x_N) \cdot \nabla_{\vartheta_n} x_{n+1} + \delta \nabla_{\vartheta_n} L(\vartheta_n) \\ = & - p_{n+1} \cdot \nabla_{\vartheta_n} g_n(x_n, \vartheta_n) + \delta \nabla_{\vartheta_n} L(\vartheta_n) \\ = & - \nabla_{\vartheta_n} H_n. \end{split}$$

Hence, (13) is simply the gradient descent step

$$\vartheta_n^{k+1} = \vartheta_n^k - \eta \nabla_{\vartheta_n} J(\vartheta^k).$$

As the proposition shows, gradient descent with back-propagation can be seen as a modification of the basic MSA by replacing the Hamiltonian maximization step with a gradient ascent step. However, we note that the PMP (and MSA convergence) holds, at least in continuous-time, even when differentiability with respect to  $\vartheta$  is not satisfied, and hence is more general than the classical back-propagation. In fact, the PMP formalism shows that the back-propagation of information through a deep network is handled by the co-state equation and there is no requirement or relationship to the gradients with respect to the trainable parameters. In other words, optimization is performed at each layer separately (with or without gradient information), and propagation is independent of optimization.

## 4.3 A Remark on Mini-batch Algorithms

So far, our discussion has focused on full-batch algorithms, where the input x represents the full set of training inputs. As modern supervised learning tasks typically involve a large number of training samples, usually the optimization problem has to be solved in mini-batches, where at each iteration we sub-sample m input-label pairs and optimize the parameters  $\theta$  (or  $\theta$  in discrete time) based on losses evaluated on these pairs. In the context of continuous-time PMP, we can write the batch version of the three necessary conditions as

$$\begin{split} \dot{X}_{t}^{i,*} &= \nabla_{p} H(t, X_{t}^{i,*}, P_{t}^{i,*}, \theta_{t}^{*}), & X_{0}^{i,*} &= x^{i}, \\ \dot{P}_{t}^{i,*} &= -\nabla_{x} H(t, X_{t}^{i,*}, P_{t}^{i,*}, \theta_{t}^{*}), & P_{T}^{i,*} &= -\nabla \Phi^{i}(X_{T}^{i,*}), \\ \theta_{t}^{*} &= \arg\max_{\theta \in \Theta} \sum_{i=1}^{M} H(t, X_{t}^{i,*}, P_{t}^{i,*}, \theta), & t \in [0, T], \end{split}$$

for samples  $i=1,\ldots,M$ . We omit for brevity the equivalent expressions for discrete-time. In particular, notice that the propagation steps are decoupled across samples, and hence can be carried out independently. The only difference is the maximization step, where in a mini-batch setting we would evaluate instead

$$\underset{\theta \in \Theta}{\operatorname{arg\,max}} \sum_{i=1}^{m} H(t, X_t^{i,*}, P_t^{i,*}, \theta).$$

If m is large enough and the samples are independently and identically drawn, then uniform law of large numbers (Jennrich, 1969) holds under fairly general conditions and ensures that the mini-batch mean of Hamiltonians converges uniformly in  $\theta$  to the full-batch sum. Hence, maximization performed on the mini-batch sum should be close to the actual maximization on the full Hamiltonian. Rigorous error estimates for the mini-batch version of our algorithm is out of the scope of the current work, and we use instead numerical results in Section 5 to demonstrate that the algorithm can also be carried out in a mini-batch fashion.

# 5. Numerical Experiments

In this section, we investigate the performance of E-MSA compared with the usual gradient-based approaches, namely stochastic gradient descent and its variants: Adagrad (Duchi et al., 2011) and Adam (Kingma and Ba, 2014). To illustrate key properties of E-MSA, we shall begin by investigating some synthetic examples. First, we consider a simple one-dimensional function approximation problem where we want to approximate  $F(x) = \sin(x)$  for  $x \in [-\pi, \pi]$  using a continuous time dynamical system. Let T = 5 and consider

$$\dot{X}_t = f(X_t, \theta_t) = \tanh(W_t X_t + b_t),$$

where  $\theta_t = (W_t, b_t) \in \mathbb{R}^{5 \times 5} \times \mathbb{R}^5$ , i.e. a continuous analogue of a fully connected feed forward neural networks with 5 nodes per layer. To match dimensions, we shall concatenate the input x to form a five dimensional vector of identical components, which is now the initial condition to the dynamical system on  $\mathbb{R}^d$ . The output of the network is  $\sum_{i=1}^5 X_T^i$ . and we define the loss function due to one sample to be  $\Phi(X_T) = (\sum_{i=1}^5 X_T^i - \sin(x))^2$ . For multiple samples, we average the loss function over all samples in the usual way. We apply E-MSA with discretization size  $\delta = 0.25$  (giving 20 layers) and compute the Hamiltonian maximization step using 10 iterations of limited memory BFGS method (L-BFGS) (Liu and Nocedal, 1989). In Figure 1(a), we compare the results with gradient descent based optimization approaches, where we observe that E-MSA has favorable convergence rate per-iteration. More interestingly, it is well-known that gradient descent may suffer slow convergence at flat regions or near saddle-points, where the gradients become very small and optimization may stall for a long time. This often occurs as a result of poor initialization of weights and biases (Sutskever et al., 2013). Here, we simulate this scenario by initializing all weights and biases  $(W_t, b_t)$  to be 0 and observe the optimization process. We see from Figure 1(b) that gradient descent based methods are more easily stalled at flat regions. We calculated numerically the eigenvalues of the Hessian at this region, which confirms that this is indeed very close to a saddle point. On the other hand, the Hamiltonian maximization in E-MSA can quickly escape the locally flat regions. One possible reason is that secondorder information employed by L-BFGS can off-set the small gradients and provide larger updates.

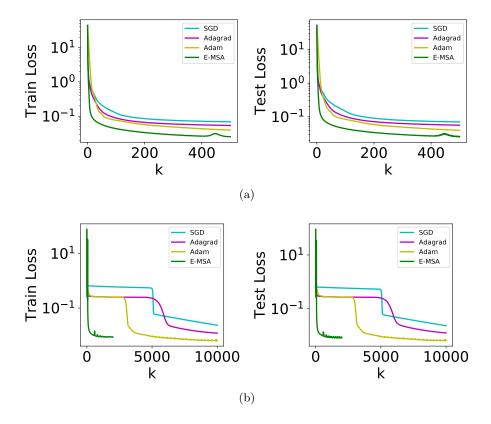


Figure 1: Comparison of E-MSA with gradient-based methods for approximating the sine function with a continuous, 5-dimensional dynamical system. A training and test set of 1000 samples each are used. (a) Loss function vs iterations for a good initialization, where weights are initialized with truncated random normal variables with standard deviation 0.1 and biases are initialized as constants equal to 0.1. We see that E-MSA has good convergence rate per iteration. (b) We use a poor initialization by setting all weights and biases to 0. We observe that gradient descent based methods tend to become stuck whereas E-MSA are better at escaping these slow manifolds, provided that  $\rho$  is well chosen (=1.0 in this case).

Next, we consider a familiar supervised learning test problem: the MNIST data set (Le-Cun, 1998) for handwritten digit recognition, with 55000 training samples and 10000 test samples. We employ a continuous dynamical system that resembles a (residual) convolution neural network (LeCun and Bengio, 1995) when discretized. More concretely, at each t we consider the map  $f(t, x, \theta) = \tanh(W \star x + b)$  where W is a  $3 \times 3$  convolution filter with 32 input and output channels. To match dimensions, we introduce two projection layers at the input (consisting of convolution, point-wise non-linearities followed by  $2 \times 2$  maxpooling). We also use a fully-connected classification layer as the final layer, with softmax cross-entropy loss. Note that the input projection layers and fully-connected output layers

are not of residual form, but we can nevertheless apply Algorithm 3 with the appropriate g. We use a total of 10 layers (2 projections, 1 fully-connected and 7 residual layers with  $\delta = 0.5$ , i.e T = 3.5). The model is trained with mini-batch sizes of 100 using E-MSA and gradient-descent based methods, namely SGD, Adagrad, and Adam. For E-MSA, we approximately solve the Hamiltonian maximization step using either 10 iterations of L-BFGS. Note that since we have decoupled the layers through the PMP, the L-BFGS step used to maximize H is tractable since it involves much fewer parameters than directly minimizing J. Figure 2 compares the performance of E-MSA with the other gradient-descent based methods, where we observe that E-MSA has good performance per-iteration, especially at early stages of training. However, we also show in Figure 3 that the wall-clock performance of our methods are not currently competitive, because the Hamiltonian maximization step is time consuming and the performance gains per iteration is outweighed by the running time. Note that wall-clock times are compared on the CPU for fairness since we did not use a GPU implementation of L-BFGS. As a further test, we train the same model on a different data set, the fashion MNIST data set (Xiao et al., 2017), where we again observe similar phenomena (see Figure 4). Experiments on more complex data sets such as ImageNet (Deng et al., 2009) with larger residual networks is a direction of future work. In particular, this may require further improvements to the Hamiltonian maximization step current handled by direct minimization with L-BFGS, which can be significantly slower (on a wall-clock basis) for larger networks and data sets.

# 6. Discussion and Related Work

We commence this section by highlighting the distinguishing features of E-MSA from traditional gradient-descent based training methods. First, the formulations of PMP and E-MSA do not involve gradient information with respect to the trainable parameters. In fact, Theorem 1 and Algorithm 2 remain valid even when the trainable parameters can only take values in a discrete set. Second, due to a more drastic argmax step taken at each iteration, E-MSA tends to have better convergence rates at the early steps of training, as observed in our numerical experiments (Section 5). Third, in the PMP formalism, the Hamiltonian equations for the state and co-state are the "forward and backward propagations", whereas given the state and co-state values, the optimization step is decoupled across layers. This allows one to potentially parallelize the often time-consuming optimization step. Moreover, from Lemma 2, we show that as long as the Hamiltonian is sufficiently increased in a layer without causing too much loss in the Hamiltonian dynamics feasibility conditions, we can ensure decrement of the loss function. This is the reason why we can use a small number of iterations of L-BFGS at each step. Moreover, this suggests that the argmax updates need not happen synchronously, i.e. the optimization in each layer can be a separate thread or process that computes the argmax and updates that layer's parameters independent of other layers. The propagation may also potentially be allowed to happen asynchronously as long as updates are sufficiently frequent. We leave a rigorous analysis of an asynchronous version of the current approach to future work. In summary, the main strength of the PMP (over e.g. solving the KKT conditions using gradient methods) is that PMP says that at the optimum, the Hamiltonian is not only stationary (KKT), but globally maximized. This hints that heuristic global optimization methods can be applied to H to obtain algorithms that

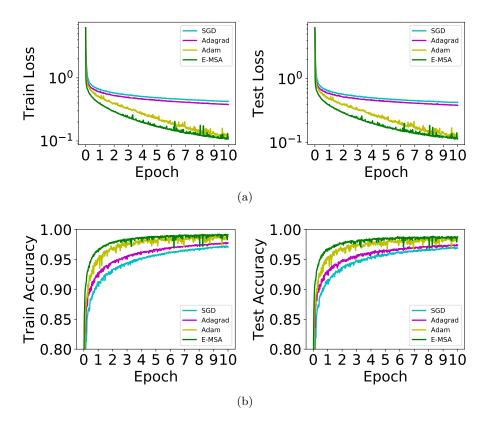


Figure 2: Comparison of E-MSA with gradient-based methods for the residual CNN on the MNIST data set. Mini-batch size of 100 is used so that each epoch of training consists of 550 iterations. (a) Train and test Loss vs epoch. (b) Train and test accuracy vs epoch. For each case, we tuned the associated hyper-parameters on a coarse grid for optimal performance. We observe that per-iteration, E-MSA performs favorably, at least at early times. This shows that if the augmented Hamiltonian can be efficiently maximized, we may obtain good performance.

are very different in behavior compared with gradient-descent based approaches. Again, Lemma 2 ensures that such heuristic global maximization need only be approximate.

As it currently stands, our experiments in Section 5 demonstrate that the Hamiltonian maximization step in E-MSA gives very different behavior compared with gradient-descent based methods. When the Hamiltonian is sufficiently maximized, we indeed obtain favorable performance compared with gradient descent based methods. Furthermore, we saw in Figure 1 that Hamiltonian maximization may avoid pitfalls such as a very flat landscape. Overall, the key to whether E-MSA (and other methods based on solving the PMP) will eventually constitute a replacement for gradient-descent based algorithm lies in the question of whether efficient Hamiltonian maximization can be performed at reasonable computational costs. Although this is still a non-convex optimization problem, it is much simpler than the original training problem because: (1) Optimization in the layers are decoupled and hence parameter space is greatly reduced; (2) The Hamiltonian is formally similar across

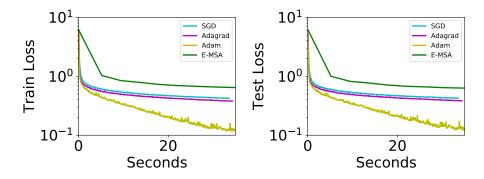


Figure 3: Comparison of E-MSA with gradient-based methods for the residual CNN on the MNIST data set on a wall-clock basis. We observe that currently, the gains per iteration is outweighed by the additional computational costs. Note that we did not use a GPU implementation for the L-BFGS algorithm used to maximize the augmented Hamiltonian, hence the wall-clock time for E-MSA is expected to be improved. Nevertheless, we expect that more efficient Hamiltonian maximization algorithms must be developed for E-MSA to out-perform gradient-based methods in terms of wall-clock efficiency.

different layers, loss functions and models, so specialized algorithms may be designed; (3) The Hamiltonian does not need to be maximized exactly, thus fast heuristic methods (Lee and El-Sharkawi, 2008) or learning (Andrychowicz et al., 2016; Jaderberg et al., 2016; Czarnecki et al., 2017) can potentially be used to perform this. All these are worthy of future exploration in order to make E-MSA truly competitive.

Next, we put our work in perspective by discussing related work in the optimal control, optimization and deep learning literature. First, the work on numerical algorithms for the solution of optimal control problem is abundant (see e.g. Rao 2009 for a survey). Many of the state-of-the-art techniques in the control theory literature assume a moderately small problem size, so that conventional non-linear programming techniques (Bertsekas, 1999; Bazaraa et al., 2013) as well as shooting (Roberts and Shipman, 1972) and collocation methods (Betts, 1998) produce efficient algorithms. This is usually not the case for largescale machine learning problems, where often, the only scalable approach is to rely on iterative updates to the parameters. This is the reason for our focus on the MSA algorithms (Chernousko and Lyubushin, 1982), as they are straight-forward to implement and typically have linear scaling in computational complexity with respect to the input and parameter sizes. The basic MSA is discussed in Krylov and Chernousko (1962), and a number of improved variants are discussed in Chernousko and Lyubushin (1982) and references therein. For example, a popular improvement is based on needle-perturbations, where controls are varied on small intervals at each iteration. While convergent, the main issue with the needle-perturbation approach is the requirement of a sufficiently fine mesh (i.e. many layers in the discretized network), which impacts computational speed. A possible solution is the use of adaptive meshes, which is a future direction we plan to investigate. Our variant of the MSA presented in this work differs from classical approaches (Chernousko and Lyubushin, 1982) mainly in the sense that we solve a weaker sufficient condition (extended

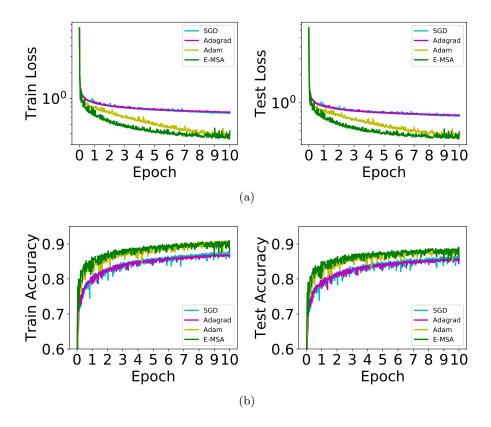


Figure 4: Comparison of E-MSA with gradient-based methods for the residual CNN on the fashion MNIST data set. We use the same network structure and mini-batch sizes as in Figure 2. The hyper-parameters have to be slightly re-tuned. (a) Train and test Loss vs epoch. (b) Train and test accuracy vs epoch. Again, we observe E-MSA performs favorably per-iteration.

PMP, Proposition 3), which then allows us to control errors in the Hamiltonian dynamical equations at every iteration without going into finer mesh-sizes. The regularization terms proportional to  $\rho$  is similar to the heuristic modifications suggested in Lyubushin (1982) by regularizing the distance between  $\theta^k$  and  $\theta^{k+1}$ , but we do not have to assume convexity of  $\Theta$  or that f is Lipschitz in  $\theta$ .

In the optimization literature, our work shares some similarity with the recently proposed ADMM methods (Taylor et al., 2016) for training deep neural networks, where the authors also considered necessary conditions with Lagrange multipliers that can decouple optimization across layers. The main difference in our work is that the PMP gives a stronger necessary condition (Hamiltonian maximization) that also applies to general parameter spaces (e.g., discrete, or bounded with non-linear constraints). Our modification of the basic MSA in terms of the augmented Hamiltonian is inspired by the method of augmented Lagrangians often applied in constrained optimization (Hestenes, 1969). The idea of viewing an initially discrete system as the discretization of a continuous-time system

has been explored in Li et al. (2017) in the form of stochastic optimization. Our current work is also in this flavor, but for neural network models.

In deep learning, there are a few works that share our perspective of deep neural networks as a discretization of a dynamical system. We note that the connection between the PMP and back-propagation has been pointed out qualitatively in LeCun (1988) and in the development of back-propagation (Bryson, 1975; Baydin et al., 2015), although to the best of our knowledge, this work is the first attempt to translate numerical algorithms for the PMP into training algorithms for deep learning that goes beyond gradient descent. The treatment of machine learning as function approximation via a dynamical system has been presented in E (2017). The recent work of Haber and Ruthotto (2017); Chang et al. (2017) also propose the dynamical systems viewpoint, and the authors used continuous-time tools to address stability issues. In contrast, our work focuses on the optimization aspects centered around the PMP. We also mention other recent approaches to decouple optimization in deep neural networks, such as synthetic gradients (Jaderberg et al., 2016; Czarnecki et al., 2017) and proximal back-propagation (Frerix et al., 2017).

#### 7. Conclusion and Outlook

In this paper, we discuss the viewpoint that deep residual neural networks can be viewed as discretization of a continuous-time dynamical system, and hence supervised deep learning can be regarded as solving an optimal control problem in continuous time. We explore a concrete consequence of this connection, by modifying the classical method of successive approximations for solving optimal control problems (in particular the PMP) into a method for solving a weaker sufficient condition (extended PMP). We prove the convergence of the resulting algorithm (E-MSA) and test it on various benchmark problems, where we observe that the E-MSA algorithm performs favorably on a per-iteration basis, especially at early stages of training, compared with gradient-based approaches such as SGD, Adagrad and Adam.

There are many avenues of future research. On the algorithmic side, it is necessary to further improve the computational efficiency of the E-MSA, in particular the Hamiltonian maximization step. Moreover, adaptive selection of  $\rho$  depending on iteration number and/or layer can be explored, e.g. by designing adaptive tuning schemes using control theoretic tools (Li et al., 2017). Also, it is desirable to formulate and analyze the PMP and E-MSA from a discrete-time perspective in order to broaden the method's application. From a modeling perspective, viewing deep neural networks as continuous-time dynamical systems is useful in the sense that it allows one to think of neural network architectures as dynamical objects. Indeed, at each training iteration of the E-MSA, we do not have to use the same discretization scheme to compute the Hamiltonian dynamical equations. Also, as the PMP and E-MSA assume little structure on the parameter space  $\Theta$ , it will also be interesting to apply the E-MSA to train neural networks that have discrete weights (e.g. those that can only take on binary values). Such networks have the advantage of fast inference speed and small memory requirement. However, training such networks is a challenge and most existing techniques rely on approximating or thresholding the derivatives (Courbariaux et al., 2015, 2016). With the PMP and MSA, we may be able to directly train discrete networks in a principled way.

# Acknowledgments

The work of W. E is supported in part by Major Program of NNSFC under grant 91130005, ONR grant N00014-13-1-0338, DOE grants DE-SC0008626 and DE-SC0009248 Q. Li is supported by the Agency for Science, Technology and Research, Singapore.

# Appendix A. Function Space Formulation

In this section, we give an alternative, non-rigorous formulation of the supervised learning problem as an optimal control problem on function spaces. This provides an alternative formulation of (continuous-time) deep learning that does not make reference to a specific set of input-outputs, but rather their conditional distributions. The idea is to consider the control of a continuity equation that describes the evolution of probability densities. Hereafter, we proceed formally by assuming all differentiability and integrability conditions are satisfied.

We would like to approximate, using a dynamical systems approach, some target joint probability density  $\rho(x,y)$ , where  $x \in \mathcal{X} \subset \mathbb{R}^d$  is a sample input and  $y \in \mathcal{Y}$  is the corresponding label. In the case where the labels are deterministically determined by the samples, i.e. there exists  $F: \mathcal{X} \to \mathcal{Y}$  such that y = F(x), we would have  $\rho(x,y) = \overline{\rho}(x)\delta(F(x) - y)$ . Here,  $\overline{\rho}(x)$  is the marginal density of  $\rho(x,y)$ . In general, we can write  $\rho(x,y) = \rho(y|x)\overline{\rho}(x)$ .

As before, the idea is to consider passing the inputs through a dynamical system

$$\dot{X}_t = f(t, X_t, \theta_t), \qquad X_0 = x. \tag{14}$$

We begin with a guess of a conditional density  $\rho_0(y|x)$  of y given x. In the deterministic case, we may set  $\rho_0(y|x) = \delta(y - F_0(x))$  for some  $F_0 : \mathcal{X} \to \mathcal{Y}$  (this is like the last layer of the neural network, be it a regressor or a classifier). Note that  $F_0$  is potentially very different from F, so that  $\rho_0(\cdot|x)$  is far from our target  $\rho(\cdot|x)$ .

To improve this approximation, we drive the initial condition by the controllable dynamical system (14). That is, we define the approximation at time t of  $\rho(y|x)$  to be  $\rho_t(y|x) := \langle \rho_0(y|\cdot), u_t \rangle$ , with  $u_t$  denoting the probability density of  $X_t$  at time t (pushforward distribution of  $X_t$  according to (14)). It is well-known that  $u_t$  follows the continuity equation, or Liouville equation (Gibbs, 2014); or forward Kolmogorov equation in stochastic processes, but with zero noise (Risken, 1996),

$$\frac{d}{dt}u_t = -\operatorname{div}(f(t, \cdot, \theta_t)u_t), \qquad u_0 = \delta_x, \tag{15}$$

where  $\operatorname{div} u = \sum_i \partial u / \partial x_i$  is the divergence operator and  $\delta_x(x') = \delta(x - x')$  is a point-mass at x. We shall assume that  $u_t \in \mathcal{H} \subset L^2(\mathbb{R}^d)$  for some function space  $\mathcal{H}$ , for all  $t \in (0, T]$ .

The goal now is to adjust  $\theta \in \mathcal{U}$  so that  $\rho_t(\cdot|x)$  is close to  $\rho(\cdot|x)$ . To this end, we define a differentiable loss function  $\Phi(\rho_1, \rho_2)$  that measures distances between two conditional densities  $\rho_1, \rho_2$  (e.g.,  $L^2$  loss, K-L divergence). Then, the learning problem can be formulated

as the following optimal control problem:

$$\min_{\theta \in \mathcal{U}} \mathbb{E}_{x \sim \overline{\rho}} \left[ \Phi(\rho_t(\cdot|x), \rho(\cdot|x)) + \int_0^T L(\theta_t) dt \right],$$

$$\frac{d}{dt} u_t = -\text{div}(f(t, \cdot, \theta_t) u_t), \qquad u_0 = \delta_x. \tag{16}$$

As before, L is a regularizer on the trainable parameters. Now, (16) is an optimal control problem on the function space  $\mathcal{H}$ .

We now write down formally a set of necessary conditions for optimality, in the form of the Pontryagin's maximum principle, for the present function-space control problem (16). Define the Hamiltonian functional  $H: [0,T] \times \mathcal{H} \times \mathcal{H} \times \Theta \to \mathbb{R}$ 

$$\begin{split} H(t,u,v,\theta) &:= -\langle v, \operatorname{div}(f(t,\cdot,\theta)u)\rangle - L(\theta) \\ &= -\int_{\mathbb{R}^d} v(x) \sum_{i=1}^d \frac{\partial}{\partial x_i} (f(t,x,\theta)u(x)) dx - L(\theta). \end{split}$$

Then, the Pontryagin's maximum principle for this system is expected to take the form: let  $\theta^* \in \mathcal{U}$  be an optimal control, then there exists a co-state process  $v_t \in \mathcal{H}$  such that

$$\begin{split} \frac{d}{dt}u_{t}^{*} &= D_{v}H(t, u_{t}^{*}, v_{t}^{*}, \theta_{t}^{*}), & u_{0}^{*} &= \delta_{x}, \\ \frac{d}{dt}v_{t}^{*} &= -D_{u}H(t, u_{t}^{*}, v_{t}^{*}, \theta_{t}^{*}), & v_{T}^{*} &= -D_{u}\Phi(\langle \rho_{0}, u_{T}^{*} \rangle, \rho(\cdot|x)) \\ \mathbb{E}_{x \sim \overline{\rho}}H(t, u_{t}^{*}, v_{t}^{*}, \theta_{t}^{*}) &\geq \mathbb{E}_{x \sim \overline{\rho}}H(t, u_{t}^{*}, v_{t}^{*}, \theta), & \theta \in \Theta, t \in [0, T], \end{split}$$

where D denotes the usual Fréchet derivative. Note that by definition, we have  $D_v H = -\text{div}(fu)$  and  $D_u H = f \cdot \nabla_x v$ . Observe that the co-state  $v^*$  satisfies the (time-reversed) adjoint Liouville's equation with a specified terminal condition. The PMP for similar functional optimal control problems has been studied in, among others, Pogodaev (2016); Roy and Borz (2017), albeit without the expectation over initial density.

In summary, the advantage of this formulation is that we make no explicit reference to the training data or target functions and formulate the entire problem as a control problem on probability densities. Of course, in practice, to implement an MSA-like algorithm, the terminal condition of the co-state will depend on the target joint density, which we can only access through the sampled data. A rigorous analysis of this function space control formulation and its consequences will be explored in future work.

# Appendix B. Proof of Lemma 2

First, observe that assumptions (A1)-(A2) in the main text implies that the second derivatives of f and  $\Phi$  are bounded by K. Provided that  $P_t^{\theta}$  is bounded, they also imply that the second derivatives of H with respect to x and p are bounded when evaluated on  $X_t^{\theta}$ ,  $P_t^{\theta}$ ,  $\theta_t$ . We first establish the boundedness of  $P_t^{\theta}$ .

**Lemma 6** Assume that (A1)-(A2) hold. Then, there exists a constant K' > 0 such that for any  $\theta$ ,

$$||P_t^{\theta}|| \le K',$$

for all  $t \in [0, T]$ .

**Proof** Using (7) and setting  $\tau := T - t$ ,  $\tilde{P}^{\theta}_{\tau} := P^{\theta}_{T-\tau}$  we get

$$\dot{\tilde{P}}_{\tau}^{\theta} = \tilde{P}_{\tau}^{\theta} \cdot \nabla_x f(t, X_{T-\tau}^{\theta}, T - \tau), \qquad \tilde{P}_{0}^{\theta} = -\nabla \Phi(X_T^{\theta}).$$

Using (A1)-(A2), we have  $||P_T^{\theta}|| = ||\nabla_x \Phi(X_T^{\theta})|| \le K$  and  $||\nabla_x f(t, X_t^{\theta}, \theta_t)||_2 \le K$ . Hence,

$$\|\dot{\tilde{P}}_{\tau}^{\theta}\| \le K \|\tilde{P}_{\tau}^{\theta}\|,$$

and by Gronwall's inequality,

$$\|\tilde{P}_{\tau}^{\theta}\| \le Ke^{KT} =: K'.$$

This proves the claim since it holds for any  $\tau$ .

We now prove Lemma 2. The approach here is similar to that employed in Rozonoer (1959).

**Proof** [Proof of Lemma 2] From (6) and the definition of the Hamiltonian, we have for any  $\theta \in \mathcal{U}$ ,

$$I(X^{\theta}, P^{\theta}, \theta) := \int_0^T P_t^{\theta} \cdot \dot{X}_t^{\theta} - H(t, X_t^{\theta}, P_t^{\theta}, \theta_t) - L(\theta_t) dt \equiv 0.$$

Denote  $\delta X_t = X_t^{\phi} - X_t^{\theta}$  and  $\delta P_t = P_t^{\phi} - P_t^{\theta}$ , then we have

$$0 \equiv I(X^{\phi}, P^{\phi}, \phi) - I(X^{\theta}, P^{\theta}, \theta)$$

$$= \int_{0}^{T} P_{t}^{\theta} \cdot \delta \dot{X}_{t} + \delta P_{t} \cdot \dot{X}_{t}^{\theta} + \delta P_{t} \cdot \delta \dot{X}_{t} dt$$

$$- \int_{0}^{T} H(t, X_{t}^{\phi}, P_{t}^{\phi}, \phi_{t}) - H(t, X_{t}^{\theta}, P_{t}^{\theta}, \theta_{t}) dt$$

$$- \int_{0}^{T} L(\phi_{t}) - L(\theta_{t}) dt.$$

$$(17)$$

Now, by integration by parts

$$\int_0^T P_t^{\theta} \cdot \delta \dot{X}_t dt = P_t^{\theta} \cdot \delta X_t \Big|_0^T - \int_0^T \dot{P}_t^{\theta} \cdot \delta X_t dt, \tag{18}$$

$$\int_{0}^{T} \delta P_{t} \cdot \delta \dot{X}_{t} dt = \delta P_{t} \cdot \delta X_{t} \Big|_{0}^{T} - \int_{0}^{T} \delta \dot{P}_{t} \cdot \delta X_{t} dt. \tag{19}$$

Using (6), (7) and (18), we have

$$\int_{0}^{T} P_{t}^{\theta} \cdot \delta \dot{X}_{t} + \delta P_{t} \cdot \dot{X}_{t}^{\theta} dt$$

$$= P_{t}^{\theta} \cdot \delta X_{t} \Big|_{0}^{T} + \int_{0}^{T} \left( f(t, X_{t}^{\theta}; \theta_{t}) \cdot \delta P + \nabla_{x} H(t, X_{t}^{\theta}, P_{t}^{\theta}, \theta_{t}) \cdot \delta X \right) dt$$

$$= P_{t}^{\theta} \cdot \delta X_{t} \Big|_{0}^{T} + \int_{0}^{T} \left( \nabla_{z} H(t, Z_{t}^{\theta}, \theta_{t}) \cdot \delta Z \right) dt. \tag{20}$$

where in the last line we defined  $Z^{\theta} := (X^{\theta}, P^{\theta})$ . Similarly, from (19) we get

$$\int_{0}^{T} \delta P_{t} \cdot \delta \dot{X}_{t} dt = \frac{1}{2} \int_{0}^{T} \delta P_{t} \cdot \delta \dot{X}_{t} dt + \frac{1}{2} \int_{0}^{T} \delta P_{t} \cdot \delta \dot{X}_{t} dt$$

$$= \frac{1}{2} \delta P_{t} \cdot \delta X_{t} \Big|_{0}^{T}$$

$$+ \frac{1}{2} \int_{0}^{T} \left( \left[ \nabla_{z} H(t, Z_{t}^{\phi}, \phi_{t}) - \nabla_{z} H(t, Z_{t}^{\theta}, \theta_{t}) \right] \cdot \delta Z_{t} \right) dt$$

$$= \frac{1}{2} \delta P_{t} \cdot \delta X_{t} \Big|_{0}^{T}$$

$$+ \frac{1}{2} \int_{0}^{T} \left[ \nabla_{z} H(t, Z_{t}^{\theta}, \phi_{t}) - \nabla_{z} H(t, Z_{t}^{\theta}, \theta_{t}) \right] \cdot \delta Z_{t} dt$$

$$+ \frac{1}{2} \int_{0}^{T} \delta Z_{t} \cdot \nabla_{z}^{2} H(t, Z_{t}^{\theta} + r_{1}(t) \delta Z_{t}, \phi_{t}) \cdot \delta Z_{t} dt. \tag{21}$$

where we have used Taylor's theorem in the last step with  $r_1(t) \in [0,1]$ . We now rewrite the boundary terms. Since  $\delta X_0 = 0$ , we have

$$(P_t^{\theta} + \frac{1}{2}\delta P_t) \cdot \delta X_t \Big|_0^T = (P_T^{\theta} + \frac{1}{2}\delta P_T) \cdot \delta X_T$$

$$= -\nabla \Phi(X_T^{\theta}) \cdot \delta X_T - \frac{1}{2}(\nabla \Phi(X_T^{\phi}) - \nabla \Phi(X_T^{\theta})) \cdot \delta X_T$$

$$= -\nabla \Phi(X_T^{\theta}) \cdot \delta X_T - \frac{1}{2}\delta X_T \cdot \nabla^2 \Phi(X_T^{\theta} + r_2 \delta X_T) \cdot \delta X_T$$

$$= -(\Phi(X_T^{\phi}) - \Phi(X_T^{\theta})) - \frac{1}{2}\delta X_T \cdot (\nabla^2 \Phi(X_T^{\theta} + r_2 \delta X_T) + \nabla^2 \Phi(X_T^{\theta} + r_3 \delta X_T)) \cdot \delta X_T, \quad (22)$$

for some  $r_2, r_3 \in [0, 1]$ . Lastly, for each  $t \in [0, T]$  we have

$$H(t, Z_t^{\phi}, \phi_t) - H(t, Z_t^{\theta}, \theta_t) = H(t, Z_t^{\theta}, \phi_t) - H(t, Z_t^{\theta}, \theta_t)$$

$$+ \nabla_z H(t, Z_t^{\theta}, \phi_t) \cdot \delta Z_t$$

$$+ \frac{1}{2} \delta Z_t \cdot \nabla_z^2 H(t, Z_t^{\theta} + r_4(t) \delta Z_t, \phi_t) \cdot \delta Z_t,$$

$$(23)$$

where  $r_4(t) \in [0, 1]$ .

Substituting (20), (21), (22), (23) into (17), we obtain

$$\left[\Phi(X_T^{\phi}) + \int_0^T L(\phi_t)\right] - \left[\Phi(X_T^{\theta}) + \int_0^T L(\theta_t)\right]$$

$$= \frac{1}{2}\delta X_T \cdot (\nabla^2 \Phi(X_T^{\theta} + r_2 \delta X_T) + \nabla^2 \Phi(X_T^{\theta} + r_3 \delta X_T)) \cdot \delta X_T$$

$$- \int_0^T \Delta H_{\phi,\theta}(t) dt$$

$$+ \frac{1}{2} \int_0^T (\nabla_z H(t, Z_t^{\theta}, \phi_t) - \nabla_z H(t, Z_t^{\theta}, \theta_t)) \cdot \delta Z_t dt$$

$$+ \frac{1}{2} \int_0^T \left(\delta Z_t \cdot [\nabla_z^2 H(t, Z_t^{\theta} + r_1(t) \delta Z_t, \phi_t) - \nabla_z^2 H(t, Z_t^{\theta} + r_4(t) \delta Z_t, \phi_t)] \cdot \delta Z_t\right) dt. \quad (24)$$

The left hand side is simply  $J(\phi) - J(\theta)$ , and so it remains to estimate the right hand side terms. First, let us estimate  $\delta X$  and  $\delta P$ . By definition,

$$\delta \dot{X}_t = f(t, X_t^{\phi}, \phi_t) - f(t, X_t^{\theta}, \theta_t).$$

Integrating, we get

$$\delta X_t = \int_0^t f(t, X_s^{\phi}, \phi_s) - f(t, X_s^{\theta}, \theta_s) ds,$$

and so

$$\|\delta X_{t}\| \leq \int_{0}^{t} \|f(t, X_{s}^{\phi}, \phi_{s}) - f(t, X_{s}^{\theta}, \theta_{s})\| ds$$

$$\leq \int_{0}^{t} \|f(t, X_{s}^{\phi}, \phi_{s}) - f(t, X_{s}^{\theta}, \phi_{s})\| ds$$

$$+ \int_{0}^{t} \|f(t, X_{s}^{\theta}, \phi_{s}) - f(t, X_{s}^{\theta}, \theta_{s})\| ds$$

$$\leq \int_{0}^{T} \|f(t, X_{s}^{\theta}, \phi_{s}) - f(t, X_{s}^{\theta}, \theta_{s})\| ds$$

$$+ K \int_{0}^{t} \|\delta X_{s}\| dt.$$
(25)

By Gronwall's inequality, we have

$$\|\delta X_t\| \le e^{KT} \int_0^T \|f(t, X_s^{\theta}, \phi_s) - f(t, X_s^{\theta}, \theta_s)\| ds.$$
 (26)

To estimate  $\delta P$ , we use the same substitution as in Lemma 6 with  $\tau = T - t$  and  $\tilde{\cdot}_{\tau} = \cdot_{T-t}$ . We get

$$\delta \tilde{P}_{\tau} = \delta \tilde{P}_{0} + \int_{0}^{\tau} \nabla_{x} H(t, \tilde{X}_{s}^{\phi}, \tilde{P}_{s}^{\phi}, \tilde{\phi}_{s}) - \nabla_{x} H(t, \tilde{X}_{s}^{\theta}, \tilde{P}_{s}^{\theta}, \tilde{\theta}_{s}) ds,$$

and hence using Lemma 6 and assumptions (A1)-(A2),

$$\begin{split} \|\delta\tilde{P}_{\tau}\| \leq &\|\delta\tilde{P}_{0}\| + \int_{0}^{\tau} \|\nabla_{x}H(t,\tilde{X}_{s}^{\phi},\tilde{P}_{s}^{\phi},\tilde{\phi}_{s}) - \nabla_{x}H(t,\tilde{X}_{s}^{\theta},\tilde{P}_{s}^{\theta},\tilde{\theta}_{s})\|ds \\ \leq &K\|\delta X_{T}\| + KK'\int_{0}^{T} \|\delta X_{s}\|ds + K\int_{0}^{\tau} \|\delta\tilde{P}_{s}\|ds \\ &+ \int_{0}^{T} \|\nabla_{x}H(t,X_{s}^{\theta},P_{s}^{\theta},\phi_{s}) - \nabla_{x}H(t,X_{s}^{\theta},P_{s}^{\theta},\theta_{s})\|ds \\ \leq &e^{KT}K(\|\delta X_{T}\| + K'\int_{0}^{T} \|\delta X_{s}\|ds) \\ &+ e^{KT}\int_{0}^{T} \|\nabla_{x}H(t,X_{s}^{\theta},P_{s}^{\theta},\phi_{s}) - \nabla_{x}H(t,X_{s}^{\theta},P_{s}^{\theta},\theta_{s})\|ds. \end{split}$$
 (27)

Using estimate (26), we obtain

$$\|\delta P_{t}\| \leq e^{2KT} K(1 + K'T) \int_{0}^{T} \|f(t, X_{s}^{\theta}, \phi_{s}) - f(t, X_{s}^{\theta}, \theta_{s})\| ds$$

$$+ e^{KT} \int_{0}^{T} \|\nabla_{x} H(t, X_{s}^{\theta}, P_{s}^{\theta}, \phi_{s}) - \nabla_{x} H(t, X_{s}^{\theta}, P_{s}^{\theta}, \theta_{s})\| ds.$$
(28)

Now, we substitute estimates (26) and (28) into (24) and rename constants for simplicity. Note that by assumptions (A1)-(A2) and Lemma 6, all the second derivative terms are bounded element-wise by some constant K''. Hence, we have  $|\delta Z_t \cdot A \cdot \delta Z_t| \leq K'' ||\delta Z||^2$  for each A being a second derivative matrix. Thus we obtain

$$\begin{split} J(\phi) - J(\theta) &\leq -\int_{0}^{T} \Delta H_{\phi,\theta}(t) dt \\ &+ \frac{1}{2} K'' \|\delta X_{T}\|^{2} \\ &+ K'' \int_{0}^{T} (\|\delta X_{t}\|^{2} + \|\delta P_{t}\|^{2}) dt \\ &+ \frac{1}{2} \int_{0}^{T} \|\delta X_{t}\| \|f(t, X_{t}^{\theta}, \phi_{t}) - f(t, X_{t}^{\theta}, \theta_{t})\| dt \\ &+ \frac{1}{2} \int_{0}^{T} \|\delta P_{t}\| \|\nabla_{x} H(t, X_{t}^{\theta}, P_{t}^{\theta}, \phi_{t}) - \nabla_{x} H(t, X_{t}^{\theta}, P_{t}^{\theta}, \theta_{t})\| dt \\ &\leq -\int_{0}^{T} \Delta H_{\phi,\theta}(t) dt \\ &+ C \left(\int_{0}^{T} \|f(t, X_{t}^{\theta}, \phi_{t}) - f(t, X_{t}^{\theta}, \theta_{t})\| dt\right)^{2} \\ &+ C \left(\int_{0}^{T} \|\nabla_{x} H(t, X_{t}^{\theta}, P_{t}^{\theta}, \phi_{t}) - \nabla_{x} H(t, X_{t}^{\theta}, P_{t}^{\theta}, \theta_{t})\|^{2} dt\right)^{2} \\ &\leq -\int_{0}^{T} \Delta H_{\phi,\theta}(t) dt \\ &+ C \int_{0}^{T} \|f(t, X_{t}^{\theta}, \phi_{t}) - f(t, X_{t}^{\theta}, \theta_{t})\|^{2} dt \\ &+ C \int_{0}^{T} \|\nabla_{x} H(t, X_{t}^{\theta}, P_{t}^{\theta}, \phi_{t}) - \nabla_{x} H(t, X_{t}^{\theta}, P_{t}^{\theta}, \theta_{t})\|^{2} dt. \end{split}$$

Remark 7 For applications, the global Lipschitz condition (A2) w.r.t. x on f may be restrictive. Note that this can be replaced by a local Lipschitz condition if we can show that  $X_t$ ,  $t \in [0,T]$  is bounded for all  $\theta \in \mathcal{U}$ . This is true if the parameter space  $\Theta$  is bounded, which we can safely assume in practice, as long as a suitable regularization is used that prevents the parameters from getting arbitrarily large. Alternatively, a projection step can be used to restrict the parameters to a bounded set. In either cases, this should not negatively affect the performance of the model.

#### References

- Vladimir V Aleksandrov. On the accumulation of perturbations in the linear systems with two coordinates. *Vestnik MGU*, 3, 1968.
- Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, and Nando de Freitas. Learning to learn by gradient descent by gradient descent. In *Advances in Neural Information Processing Systems*, pages 3981–3989, 2016.
- Michael Athans and Peter L Falb. Optimal control: an introduction to the theory and its applications. Courier Corporation, 2013.
- Atilim G Baydin, Barak A Pearlmutter, Alexey A Radul, and Jeffrey M Siskind. Automatic differentiation in machine learning: a survey. arXiv preprint arXiv:1502.05767, 2015.
- Mokhtar S Bazaraa, Hanif D Sherali, and Chitharanjan M Shetty. *Nonlinear programming:* theory and algorithms. John Wiley & Sons, 2013.
- Richard Bellman. Dynamic programming. Courier Corporation, 2013.
- Yoshua Bengio. Learning deep architectures for AI. Foundations and trends in Machine Learning, 2(1):1–127, 2009.
- Dimitri P Bertsekas. *Dynamic programming and optimal control*, volume 1. Athena scientific Belmont, MA, 1995.
- Dimitri P Bertsekas. Nonlinear programming. Athena scientific Belmont, 1999.
- John T Betts. Survey of numerical methods for trajectory optimization. *Journal of Guidance control and dynamics*, 21(2):193–207, 1998.
- Vladimir Grigor'evich Boltyanskii, Revaz Valer'yanovich Gamkrelidze, and Lev Semenovich Pontryagin. The theory of optimal processes. i. the maximum principle. Technical report, TRW SPACE TECHNOLOGY LABS LOS ANGELES CALIF, 1960.
- Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings* of COMPSTAT'2010, pages 177–186. Springer, 2010.
- Alberto Bressan and Benedetto Piccoli. Introduction to mathematical control theory. AIMS series on applied mathematics, Philadelphia, 2007.
- Arthur Earl Bryson. Applied optimal control: optimization, estimation and control. CRC Press, 1975.
- Anatolii B Butkovsky. Necessary and sufficient optimality conditions for sampled-data control systems. *Avtomat. i Telemekh*, 24(8):1056–1064, 1963.
- Bo Chang, Lili Meng, Eldad Haber, Lars Ruthotto, David Begert, and Elliot Holtham. Reversible architectures for arbitrarily deep residual neural networks. arXiv preprint arXiv:1709.03698, 2017.

- Felix L Chernousko and Alexey A Lyubushin. Method of successive approximations for solution of optimal control problems. *Optimal Control Applications and Methods*, 3(2): 101–114, 1982.
- Francis Clarke. The maximum principle in optimal control, then and now. Control and Cybernetics, 34(3):709, 2005.
- Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in Neural Information Processing Systems*, pages 3123–3131, 2015.
- Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to + 1 or 1. arXiv preprint arXiv:1602.02830, 2016.
- Wojciech M Czarnecki, Grzegorz Świrszcz, Max Jaderberg, Simon Osindero, Oriol Vinyals, and Koray Kavukcuoglu. Understanding synthetic gradients and decoupled neural interfaces. arXiv preprint arXiv:1703.00522, 2017.
- J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul): 2121–2159, 2011.
- Weinan E. A proposal on machine learning via dynamical systems. Communications in Mathematics and Statistics, 5(1):1–11, 2017.
- Thomas Frerix, Thomas Möllenhoff, Michael Moeller, and Daniel Cremers. Proximal back-propagation. arXiv preprint arXiv:1706.04638, 2017.
- J Willard Gibbs. Elementary principles in statistical mechanics. Courier Corporation, 2014.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. Deep learning. MIT press, 2016.
- Eldad Haber and Lars Ruthotto. Stable architectures for deep neural networks. arXiv preprint arXiv:1705.03341, 2017.
- Richard HR Hahnloser, Rahul Sarpeshkar, Misha A Mahowald, Rodney J Douglas, and H Sebastian Seung. Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. *Nature*, 405(6789):947, 2000.
- Hubert Halkin. A maximum principle of the pontryagin type for systems described by nonlinear difference equations. SIAM Journal on control, 4(1):90–111, 1966.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

- Magnus R Hestenes. Multiplier and gradient methods. Journal of optimization theory and applications, 4(5):303–320, 1969.
- R Jackson and F Horn. On discrete analogues of pontryagin's maximum principle. *International Journal of Control*, 1(4):389–395, 1965.
- Max Jaderberg, Wojciech M Czarnecki, Simon Osindero, Oriol Vinyals, Alex Graves, and Koray Kavukcuoglu. Decoupled neural interfaces using synthetic gradients. arXiv preprint arXiv:1608.05343, 2016.
- Robert I Jennrich. Asymptotic properties of non-linear least squares estimators. *The Annals of Mathematical Statistics*, 40(2):633–643, 1969.
- Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pages 315–323, 2013.
- Henry J Kelley. Gradient theory of optimal flight paths. Ars Journal, 30(10):947–954, 1960.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- Ivan A Krylov and Felix L Chernousko. On the method of successive approximations for solution of optimal control problems. *J. Comp. Mathem. and Mathematical Physics*, 2 (6), 1962.
- Yann LeCun. A theoretical framework for back-propagation. In *The Connectionist Models Summer School*, volume 1, pages 21–28, 1988.
- Yann LeCun. The MNIST database of handwritten digits. http://yann. lecun. com/exdb/mnist/, 1998.
- Yann LeCun and Yoshua Bengio. Convolutional networks for images, speech, and time series. The handbook of brain theory and neural networks, 3361(10):1995, 1995.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553): 436–444, 2015.
- Kwang Y Lee and Mohamed A El-Sharkawi. *Modern heuristic optimization techniques:* theory and applications to power systems, volume 39. John Wiley & Sons, 2008.
- Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic gradient algorithms. In *International Conference on Machine Learning*, pages 2101–2110, 2017.
- Daniel Liberzon. Calculus of variations and optimal control theory: a concise introduction. Princeton University Press, 2012.

- Dong C Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimization. *Mathematical programming*, 45(1):503–528, 1989.
- Alexey A Lyubushin. Modifications of the method of successive approximations for solving optimal control problems. *USSR Computational Mathematics and Mathematical Physics*, 22(1):29–34, 1982.
- Zbigniew Nahorski, Hans F Ravn, and René Victor Valqui Vidal. The discrete-time maximum principle: a survey and some new results. *International Journal of Control*, 40(3): 533–554, 1984.
- Nikolay Pogodaev. Optimal control of continuity equations. Nonlinear Differential Equations and Applications, 23(2):21, 2016.
- Lev S Pontryagin. Mathematical theory of optimal processes. CRC Press, 1987.
- Anil V Rao. A survey of numerical methods for optimal control. Advances in the Astronautical Sciences, 135(1):497–528, 2009.
- Hannes Risken. Fokker-planck equation. In *The Fokker-Planck Equation*, pages 63–95. Springer, 1996.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. The annals of mathematical statistics, pages 400–407, 1951.
- Sanford M Roberts and Jerome S Shipman. Two-point boundary value problems: shooting methods. SIAM Rev., 16(2):265266, 1972.
- Souvik Roy and Alfio Borz. Numerical investigation of a class of liouville control problems. *J Sci Comput*, 73:178, 2017.
- Lev I Rozonoer. The maximum principle of L.S. Pontryagin in optimal-system theory. *Automation and Remote Control*, 20(10):11, 1959.
- Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural networks*, 61: 85–117, 2015.
- Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147, 2013.
- Gavin Taylor, Ryan Burmeister, Zheng Xu, Bharat Singh, Ankit Patel, and Tom Goldstein. Training neural networks without gradients: A scalable ADMM approach. In *International Conference on Machine Learning*, pages 2722–2731, 2016.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. arXiv preprint arXiv:1708.07747, 2017.
- Matthew D Zeiler. Adadelta: an adaptive learning rate method. arXiv preprint arXiv:1212.5701, 2012.