

# Supplementary Material for the Manuscript “On the Implicit Bias of Adam”

November 22, 2023

## Contents

<a href="#">SA-1 Overview</a>	1
<a href="#">SA-2 RMSProp with <math>\varepsilon</math> outside the square root</a>	2
<a href="#">SA-3 RMSProp with <math>\varepsilon</math> inside the square root</a>	5
<a href="#">SA-4 Adam with <math>\varepsilon</math> outside the square root</a>	7
<a href="#">SA-5 Adam with <math>\varepsilon</math> inside the square root</a>	11
<a href="#">SA-6 Technical bounding lemmas</a>	12
<a href="#">SA-7 Proof of Theorem SA-2.3</a>	21
<a href="#">SA-8 Numerical experiments</a>	27
<a href="#">SA-9 Adam with <math>\varepsilon</math> inside the square root: informal derivation</a>	29

## SA-1 Overview

**SA-1.1.** This appendix provides some omitted details and proofs.

We consider two algorithms: RMSProp and Adam, and two versions of each algorithm (with the numerical stability  $\varepsilon$  parameter inside and outside of the square root in the denominator). This means there are four main theorems: [Theorem SA-2.4](#), [Theorem SA-3.4](#), [Theorem SA-4.4](#) and [Theorem SA-5.4](#), each residing in the section completely devoted to one algorithm. The simple induction argument taken from [1], essentially the same for each of these theorems, is based on an auxiliary result whose corresponding versions are [Theorem SA-2.3](#), [Theorem SA-3.3](#), [Theorem SA-4.3](#) and [Theorem SA-5.3](#). The proof of this result is also elementary but long, and it is done by a series of lemmas in [Section SA-6](#) and [Section SA-7](#), culminating in [Section SA-7.6](#). Out of these four, we only prove [Theorem SA-2.3](#) since the other three results are proven in the same way with obvious changes.

[Section SA-8](#) contains some details about the numerical experiments.

**SA-1.2 Notation.** We denote the loss of the  $k$ th minibatch as a function of the network parameters  $\theta \in \mathbb{R}^p$  by  $E_k(\theta)$ , and in the full-batch setting we omit the index and write  $E(\theta)$ . As usual,  $\nabla E$  means the gradient of  $E$ , and nabla with indices means partial derivatives, e.g.  $\nabla_{ij_s} E$  is a shortcut for  $\frac{\partial^3 E}{\partial \theta_i \partial \theta_j \partial \theta_s}$ .

The letter  $T > 0$  will always denote a finite time horizon of the ODEs,  $h$  will always denote the training step size, and we will replace  $nh$  with  $t_n$  when convenient, where  $n \in \{0, 1, \dots\}$  is the step number. We will use the same notation for the iteration of the discrete algorithm  $\{\theta^{(k)}\}_{k \in \mathbb{Z}_{\geq 0}}$ , the piecewise ODE solution  $\tilde{\theta}(t)$  and some auxiliary terms for each of the four algorithms: see [Definition SA-2.1](#), [Definition SA-](#)

3.1, Definition SA-4.1, Definition SA-5.1. This way, we avoid cluttering the notation significantly. We are careful to reference the relevant definition in all theorem statements.

## SA-2 RMSProp with $\varepsilon$ outside the square root

**Definition SA-2.1.** In this section, for some  $\boldsymbol{\theta}^{(0)} \in \mathbb{R}^p$ ,  $\nu^{(0)} = \mathbf{0} \in \mathbb{R}^p$ ,  $\rho \in (0, 1)$ , let the sequence of  $p$ -vectors  $\{\boldsymbol{\theta}^{(k)}\}_{k \in \mathbb{Z}_{\geq 0}}$  be defined for  $n \geq 0$  by

$$\begin{aligned} \nu_j^{(n+1)} &= \rho \nu_j^{(n)} + (1 - \rho) \left( \nabla_j E_n(\boldsymbol{\theta}^{(n)}) \right)^2, \\ \theta_j^{(n+1)} &= \theta_j^{(n)} - \frac{h}{\sqrt{\nu_j^{(n+1)} + \varepsilon}} \nabla_j E_n(\boldsymbol{\theta}^{(n)}). \end{aligned} \tag{SA-2.1}$$

Let  $\tilde{\boldsymbol{\theta}}(t)$  be defined as a continuous solution to the piecewise ODE

$$\begin{aligned} \dot{\tilde{\theta}}_j(t) &= - \frac{\nabla_j E_n(\tilde{\boldsymbol{\theta}}(t))}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon} \\ &+ h \left( \frac{\nabla_j E_n(\tilde{\boldsymbol{\theta}}(t)) \left( 2P_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \bar{P}_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \right)}{2 \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))} - \frac{\sum_{i=1}^p \nabla_{ij} E_n(\tilde{\boldsymbol{\theta}}(t)) \frac{\nabla_i E_n(\tilde{\boldsymbol{\theta}}(t))}{R_i^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon}}{2 \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon \right)} \right) \end{aligned} \tag{SA-2.2}$$

with the initial condition  $\tilde{\boldsymbol{\theta}}(0) = \boldsymbol{\theta}^{(0)}$ , where  $\mathbf{R}^{(n)}(\boldsymbol{\theta})$ ,  $\mathbf{P}^{(n)}(\boldsymbol{\theta})$  and  $\bar{\mathbf{P}}^{(n)}(\boldsymbol{\theta})$  are  $p$ -dimensional functions with components

$$\begin{aligned} R_j^{(n)}(\boldsymbol{\theta}) &:= \sqrt{\sum_{k=0}^n \rho^{n-k} (1 - \rho) \left( \nabla_j E_k(\boldsymbol{\theta}) \right)^2}, \\ P_j^{(n)}(\boldsymbol{\theta}) &:= \sum_{k=0}^n \rho^{n-k} (1 - \rho) \nabla_j E_k(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\boldsymbol{\theta})}{R_i^{(l)}(\boldsymbol{\theta}) + \varepsilon}, \\ \bar{P}_j^{(n)}(\boldsymbol{\theta}) &:= \sum_{k=0}^n \rho^{n-k} (1 - \rho) \nabla_j E_k(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \frac{\nabla_i E_n(\boldsymbol{\theta})}{R_i^{(n)}(\boldsymbol{\theta}) + \varepsilon}. \end{aligned}$$

### Assumption SA-2.2.

1. For some positive constants  $M_1, M_2, M_3, M_4$  we have

$$\begin{aligned} \sup_i \sup_k \sup_{\boldsymbol{\theta}} |\nabla_i E_k(\boldsymbol{\theta})| &\leq M_1, \\ \sup_{i,j} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ij} E_k(\boldsymbol{\theta})| &\leq M_2, \\ \sup_{i,j,s} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ijs} E_k(\boldsymbol{\theta})| &\leq M_3, \\ \sup_{i,j,s,r} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ijsr} E_k(\boldsymbol{\theta})| &\leq M_4. \end{aligned}$$

2. For some  $R > 0$  we have for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$

$$R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) \geq R, \quad \sum_{k=0}^n \rho^{n-k} (1 - \rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 \geq R^2,$$

where  $\tilde{\boldsymbol{\theta}}(t)$  is defined in Definition SA-2.1.

**Theorem SA-2.3** (RMSProp with  $\varepsilon$  outside: local error bound). *Suppose Assumption SA-2.2 holds. Then for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$*

$$\left| \tilde{\theta}_j(t_{n+1}) - \tilde{\theta}_j(t_n) + h \frac{\nabla_j E_n(\tilde{\theta}(t_n))}{\sqrt{\sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\theta}(t_k)) \right)^2} + \varepsilon} \right| \leq C_1 h^3$$

for a positive constant  $C_1$  depending on  $\rho$ .

The proof of [Theorem SA-2.3](#) is conceptually simple but very technical, and we delay it until [Section SA-7](#). For now assuming it as given and combining it with a simple induction argument gives a global error bound which follows.

**Theorem SA-2.4** (RMSProp with  $\varepsilon$  outside: global error bound). *Suppose Assumption SA-2.2 holds, and*

$$\sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\theta^{(k)}) \right)^2 \geq R^2$$

for  $\{\theta^{(k)}\}_{k \in \mathbb{Z}_{\geq 0}}$  defined in [Definition SA-2.1](#). Then there exist positive constants  $d_1, d_2, d_3$  such that for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$

$$\|\mathbf{e}_n\| \leq d_1 e^{d_2 n h} h^2 \quad \text{and} \quad \|\mathbf{e}_{n+1} - \mathbf{e}_n\| \leq d_3 e^{d_2 n h} h^3,$$

where  $\mathbf{e}_n := \tilde{\theta}(t_n) - \theta^{(n)}$ . The constants can be defined as

$$\begin{aligned} d_1 &:= C_1, \\ d_2 &:= \left[ 1 + \frac{M_2 \sqrt{\bar{\rho}}}{R + \varepsilon} \left( \frac{M_1^2}{R(R + \varepsilon)} + 1 \right) d_1 \right] \sqrt{\bar{\rho}}, \\ d_3 &:= C_1 d_2. \end{aligned}$$

*Proof.* We will show this by induction over  $n$ , the same way an analogous bound is shown in [\[1\]](#).

The base case is  $n = 0$ . Indeed,  $\mathbf{e}_0 = \tilde{\theta}(0) - \theta^{(0)} = \mathbf{0}$ . Then the  $j$ th component of  $\mathbf{e}_1 - \mathbf{e}_0$  is

$$\begin{aligned} [\mathbf{e}_1 - \mathbf{e}_0]_j &= [\mathbf{e}_1]_j = \tilde{\theta}_j(t_1) - \theta_j^{(0)} + \frac{h \nabla_j E_0(\theta^{(0)})}{\sqrt{(1-\rho) \left( \nabla_j E_0(\theta^{(0)}) \right)^2 + \varepsilon}} \\ &= \tilde{\theta}_j(t_1) - \tilde{\theta}_j(t_0) + \frac{h \nabla_j E_0(\tilde{\theta}(t_0))}{\sqrt{(1-\rho) \left( \nabla_j E_0(\tilde{\theta}(t_0)) \right)^2 + \varepsilon}}. \end{aligned}$$

By [Theorem SA-2.3](#), the absolute value of the right-hand side does not exceed  $C_1 h^3$ , which means  $\|\mathbf{e}_1 - \mathbf{e}_0\| \leq C_1 h^3 \sqrt{\bar{\rho}}$ . Since  $C_1 \sqrt{\bar{\rho}} \leq d_3$ , the base case is proven.

Now suppose that for all  $k = 0, 1, \dots, n-1$  the claim

$$\|\mathbf{e}_k\| \leq d_1 e^{d_2 k h} h^2 \quad \text{and} \quad \|\mathbf{e}_{k+1} - \mathbf{e}_k\| \leq d_3 e^{d_2 k h} h^3$$

is proven. Then

$$\begin{aligned} \|\mathbf{e}_n\| &\stackrel{(a)}{\leq} \|\mathbf{e}_{n-1}\| + \|\mathbf{e}_n - \mathbf{e}_{n-1}\| \leq d_1 e^{d_2(n-1)h} h^2 + d_3 e^{d_2(n-1)h} h^3 \\ &= d_1 e^{d_2(n-1)h} h^2 \left( 1 + \frac{d_3}{d_1} h \right) \stackrel{(b)}{\leq} d_1 e^{d_2(n-1)h} h^2 (1 + d_2 h) \end{aligned}$$

$$\stackrel{(c)}{\leq} d_1 e^{d_2(n-1)h} h^2 \cdot e^{d_2 h} = d_1 e^{d_2 n h} h^2,$$

where (a) is by the triangle inequality, (b) is by  $d_3/d_1 \leq d_2$ , in (c) we used  $1+x \leq e^x$  for all  $x \geq 0$ .

Next, combining [Theorem SA-2.3](#) with [\(SA-2.1\)](#), we have

$$\left| \mathbf{e}_{n+1} - \mathbf{e}_n \right|_j \leq C_1 h^3 + h \left| \frac{\nabla_j E_n(\tilde{\boldsymbol{\theta}}(t_n))}{\sqrt{A} + \varepsilon} - \frac{\nabla_j E_n(\boldsymbol{\theta}^{(n)})}{\sqrt{B} + \varepsilon} \right|, \quad (\text{SA-2.3})$$

where to simplify notation we put

$$\begin{aligned} A &:= \sum_{k=0}^n \rho^{n-k} (1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2, \\ B &:= \sum_{k=0}^n \rho^{n-k} (1-\rho) \left( \nabla_j E_k(\boldsymbol{\theta}^{(k)}) \right)^2. \end{aligned}$$

Using  $A \geq R^2$ ,  $B \geq R^2$ , we have

$$\left| \frac{1}{\sqrt{A} + \varepsilon} - \frac{1}{\sqrt{B} + \varepsilon} \right| = \frac{|A - B|}{(\sqrt{A} + \varepsilon)(\sqrt{B} + \varepsilon)(\sqrt{A} + \sqrt{B})} \leq \frac{|A - B|}{2R(R + \varepsilon)^2}. \quad (\text{SA-2.4})$$

But since

$$\begin{aligned} & \left| \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 - \left( \nabla_j E_k(\boldsymbol{\theta}^{(k)}) \right)^2 \right| \\ &= \left| \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) - \nabla_j E_k(\boldsymbol{\theta}^{(k)}) \right| \cdot \left| \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) + \nabla_j E_k(\boldsymbol{\theta}^{(k)}) \right| \\ &\leq 2M_1 \left| \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) - \nabla_j E_k(\boldsymbol{\theta}^{(k)}) \right| \leq 2M_1 M_2 \sqrt{p} \left\| \tilde{\boldsymbol{\theta}}(t_k) - \boldsymbol{\theta}^{(k)} \right\|, \end{aligned}$$

we have

$$|A - B| \leq 2M_1 M_2 \sqrt{p} \sum_{k=0}^n \rho^{n-k} (1-\rho) \left\| \tilde{\boldsymbol{\theta}}(t_k) - \boldsymbol{\theta}^{(k)} \right\|. \quad (\text{SA-2.5})$$

Combining [\(SA-2.4\)](#) and [\(SA-2.5\)](#), we obtain

$$\begin{aligned} & \left| \frac{\nabla_j E_n(\tilde{\boldsymbol{\theta}}(t_n))}{\sqrt{A} + \varepsilon} - \frac{\nabla_j E_n(\boldsymbol{\theta}^{(n)})}{\sqrt{B} + \varepsilon} \right| \\ &\leq \left| \nabla_j E_n(\tilde{\boldsymbol{\theta}}(t_n)) \right| \cdot \left| \frac{1}{\sqrt{A} + \varepsilon} - \frac{1}{\sqrt{B} + \varepsilon} \right| + \frac{\left| \nabla_j E_n(\tilde{\boldsymbol{\theta}}(t_n)) - \nabla_j E_n(\boldsymbol{\theta}^{(n)}) \right|}{\sqrt{B} + \varepsilon} \\ &\leq M_1 \cdot \frac{2M_1 M_2 \sqrt{p} \sum_{k=0}^n \rho^{n-k} (1-\rho) \left\| \tilde{\boldsymbol{\theta}}(t_k) - \boldsymbol{\theta}^{(k)} \right\|}{2R(R + \varepsilon)^2} + \frac{M_2 \sqrt{p} \left\| \tilde{\boldsymbol{\theta}}(t_n) - \boldsymbol{\theta}^{(n)} \right\|}{R + \varepsilon} \\ &= \frac{M_1^2 M_2 \sqrt{p}}{R(R + \varepsilon)^2} \sum_{k=0}^n \rho^{n-k} (1-\rho) \left\| \tilde{\boldsymbol{\theta}}(t_k) - \boldsymbol{\theta}^{(k)} \right\| + \frac{M_2 \sqrt{p}}{R + \varepsilon} \left\| \tilde{\boldsymbol{\theta}}(t_n) - \boldsymbol{\theta}^{(n)} \right\| \\ &\stackrel{(a)}{\leq} \frac{M_1^2 M_2 \sqrt{p}}{R(R + \varepsilon)^2} \sum_{k=0}^n \rho^{n-k} (1-\rho) d_1 e^{d_2 k h} h^2 + \frac{M_2 \sqrt{p}}{R + \varepsilon} d_1 e^{d_2 n h} h^2, \quad (\text{SA-2.6}) \end{aligned}$$

where in (a) we used the induction hypothesis and that the bound on  $\|\mathbf{e}_n\|$  is already proven.

Now note that since  $0 < \rho e^{-d_2 h} \leq \rho$ , we have  $\sum_{k=0}^n (\rho e^{-d_2 h})^k \leq \sum_{k=0}^{\infty} \rho^k = \frac{1}{1-\rho}$ , which is rewritten as

$$\sum_{k=0}^n \rho^{n-k} (1-\rho) e^{d_2 k h} \leq e^{d_2 n h}.$$

Then we can continue (SA-2.6):

$$\left| \frac{\nabla_j E_n(\tilde{\boldsymbol{\theta}}(t_n))}{\sqrt{A + \varepsilon}} - \frac{\nabla_j E_n(\boldsymbol{\theta}^{(n)})}{\sqrt{B + \varepsilon}} \right| \leq \frac{M_2 \sqrt{p}}{R + \varepsilon} \left( \frac{M_1^2}{R(R + \varepsilon)} + 1 \right) d_1 e^{d_2 n h} h^2 \quad (\text{SA-2.7})$$

Again using  $1 \leq e^{d_2 n h}$ , we conclude from (SA-2.3) and (SA-2.7) that

$$\|\mathbf{e}_{n+1} - \mathbf{e}_n\| \leq \underbrace{\left( C_1 + \frac{M_2 \sqrt{p}}{R + \varepsilon} \left( \frac{M_1^2}{R(R + \varepsilon)} + 1 \right) d_1 \right)}_{\leq d_3} \sqrt{p} e^{d_2 n h} h^3,$$

finishing the induction step.  $\square$

**SA-2.5 RMSProp with  $\varepsilon$  outside: full-batch.** In the full-batch setting  $E_k \equiv E$ , the terms in (SA-2.2) simplify to

$$\begin{aligned} R_j^{(n)}(\boldsymbol{\theta}) &= |\nabla_j E(\boldsymbol{\theta})| \sqrt{1 - \rho^{n+1}}, \\ P_j^{(n)}(\boldsymbol{\theta}) &= \sum_{k=0}^n \rho^{n-k} (1 - \rho) \nabla_j E(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E(\boldsymbol{\theta}) \sum_{l=k}^{n-1} \frac{\nabla_i E(\boldsymbol{\theta})}{|\nabla_i E(\boldsymbol{\theta})| \sqrt{1 - \rho^{l+1}} + \varepsilon}, \\ \bar{P}_j^{(n)}(\boldsymbol{\theta}) &= (1 - \rho^{n+1}) \nabla_j E(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E(\boldsymbol{\theta}) \frac{\nabla_i E(\boldsymbol{\theta})}{|\nabla_i E(\boldsymbol{\theta})| \sqrt{1 - \rho^{n+1}} + \varepsilon}. \end{aligned}$$

If  $\varepsilon$  is small and the iteration number  $n$  is large, (SA-2.2) simplifies to

$$\begin{aligned} \dot{\tilde{\boldsymbol{\theta}}}_j(t) &= -\text{sign} \nabla_j E(\tilde{\boldsymbol{\theta}}(t)) + h \frac{\rho}{1 - \rho} \cdot \frac{\sum_{i=1}^p \nabla_{ij} E(\tilde{\boldsymbol{\theta}}(t)) \text{sign} \nabla_i E(\tilde{\boldsymbol{\theta}}(t))}{|\nabla_j E(\tilde{\boldsymbol{\theta}}(t))|} \\ &= \left| \nabla_j E(\tilde{\boldsymbol{\theta}}(t)) \right|^{-1} \left[ -\nabla_j E(\tilde{\boldsymbol{\theta}}(t)) + h \frac{\rho}{1 - \rho} \nabla_j \left\| \nabla E(\tilde{\boldsymbol{\theta}}(t)) \right\|_1 \right]. \end{aligned}$$

### SA-3 RMSProp with $\varepsilon$ inside the square root

**Definition SA-3.1.** In this section, for some  $\boldsymbol{\theta}^{(0)} \in \mathbb{R}^p$ ,  $\boldsymbol{\nu}^{(0)} = \mathbf{0} \in \mathbb{R}^p$ ,  $\rho \in (0, 1)$ , let the sequence of  $p$ -vectors  $\{\boldsymbol{\theta}^{(k)}\}_{k \in \mathbb{Z}_{\geq 0}}$  be defined for  $n \geq 0$  by

$$\begin{aligned} \nu_j^{(n+1)} &= \rho \nu_j^{(n)} + (1 - \rho) \left( \nabla_j E_n(\boldsymbol{\theta}^{(n)}) \right)^2, \\ \theta_j^{(n+1)} &= \theta_j^{(n)} - \frac{h}{\sqrt{\nu_j^{(n+1)} + \varepsilon}} \nabla_j E_n(\boldsymbol{\theta}^{(n)}). \end{aligned} \quad (\text{SA-3.1})$$

Let  $\tilde{\boldsymbol{\theta}}(t)$  be defined as a continuous solution to the piecewise ODE

$$\begin{aligned} \dot{\tilde{\boldsymbol{\theta}}}_j(t) &= -\frac{\nabla_j E_n(\tilde{\boldsymbol{\theta}}(t))}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))} \\ &+ h \left( \frac{\nabla_j E_n(\tilde{\boldsymbol{\theta}}(t)) \left( 2P_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \bar{P}_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \right)}{2R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))^3} - \frac{\sum_{i=1}^p \nabla_{ij} E_n(\tilde{\boldsymbol{\theta}}(t)) \frac{\nabla_i E_n(\tilde{\boldsymbol{\theta}}(t))}{R_i^{(n)}(\tilde{\boldsymbol{\theta}}(t))}}{2R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))} \right). \end{aligned} \quad (\text{SA-3.2})$$

with the initial condition  $\tilde{\boldsymbol{\theta}}(0) = \boldsymbol{\theta}^{(0)}$ , where  $\mathbf{R}^{(n)}(\boldsymbol{\theta})$ ,  $\mathbf{P}^{(n)}(\boldsymbol{\theta})$  and  $\bar{\mathbf{P}}^{(n)}(\boldsymbol{\theta})$  are  $p$ -dimensional functions with components

$$\begin{aligned} R_j^{(n)}(\boldsymbol{\theta}) &:= \sqrt{\sum_{k=0}^n \rho^{n-k}(1-\rho)(\nabla_j E_k(\boldsymbol{\theta}))^2 + \varepsilon}, \\ P_j^{(n)}(\boldsymbol{\theta}) &:= \sum_{k=0}^n \rho^{n-k}(1-\rho) \nabla_j E_k(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\boldsymbol{\theta})}{R_i^{(l)}(\boldsymbol{\theta})}, \\ \bar{P}_j^{(n)}(\boldsymbol{\theta}) &:= \sum_{k=0}^n \rho^{n-k}(1-\rho) \nabla_j E_k(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \frac{\nabla_i E_n(\boldsymbol{\theta})}{R_i^{(n)}(\boldsymbol{\theta})}. \end{aligned} \tag{SA-3.3}$$

**Assumption SA-3.2.** For some positive constants  $M_1, M_2, M_3, M_4$  we have

$$\begin{aligned} \sup_i \sup_k \sup_{\boldsymbol{\theta}} |\nabla_i E_k(\boldsymbol{\theta})| &\leq M_1, \\ \sup_{i,j} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ij} E_k(\boldsymbol{\theta})| &\leq M_2, \\ \sup_{i,j,s} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ijs} E_k(\boldsymbol{\theta})| &\leq M_3, \\ \sup_{i,j,s,r} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ijsr} E_k(\boldsymbol{\theta})| &\leq M_4. \end{aligned}$$

**Theorem SA-3.3** (RMSProp with  $\varepsilon$  inside: local error bound). *Suppose Assumption SA-3.2 holds. Then for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$*

$$\left| \tilde{\theta}_j(t_{n+1}) - \tilde{\theta}_j(t_n) + h \frac{\nabla_j E_n(\tilde{\boldsymbol{\theta}}(t_n))}{\sqrt{\sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 + \varepsilon}} \right| \leq C_2 h^3$$

for a positive constant  $C_2$  depending on  $\rho$ , where  $\tilde{\boldsymbol{\theta}}(t)$  is defined in Definition SA-3.1.

The argument is the same as for Theorem SA-2.3.

**Theorem SA-3.4** (RMSProp with  $\varepsilon$  inside: global error bound). *Suppose Assumption SA-3.2 holds. Then there exist positive constants  $d_4, d_5, d_6$  such that for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$*

$$\|\mathbf{e}_n\| \leq d_4 e^{d_5 n h} h^2 \quad \text{and} \quad \|\mathbf{e}_{n+1} - \mathbf{e}_n\| \leq d_6 e^{d_5 n h} h^3,$$

where  $\mathbf{e}_n := \tilde{\boldsymbol{\theta}}(t_n) - \boldsymbol{\theta}^{(n)}$ ;  $\tilde{\boldsymbol{\theta}}(t)$  and  $\{\boldsymbol{\theta}^{(k)}\}_{k \in \mathbb{Z}_{\geq 0}}$  are defined in Definition SA-3.1. The constants can be defined as

$$\begin{aligned} d_4 &:= C_2, \\ d_5 &:= \left[ 1 + \frac{M_2 \sqrt{p}}{\sqrt{\varepsilon}} \left( \frac{M_1^2}{\varepsilon} + 1 \right) d_4 \right] \sqrt{p}, \\ d_6 &:= C_2 d_5. \end{aligned}$$

The argument is the same as for Theorem SA-2.4.

**SA-3.5 RMSProp with  $\varepsilon$  inside: full-batch.** In the full-batch setting  $E_k \equiv E$ , the terms in (SA-3.2) simplify to

$$R_j^{(n)}(\boldsymbol{\theta}) = \sqrt{|\nabla_j E(\boldsymbol{\theta})|^2 (1 - \rho^{n+1}) + \varepsilon},$$

$$P_j^{(n)}(\boldsymbol{\theta}) = \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E(\boldsymbol{\theta}) \sum_{l=k}^{n-1} \frac{\nabla_i E(\boldsymbol{\theta})}{\sqrt{|\nabla_i E(\boldsymbol{\theta})|^2 (1-\rho^{l+1}) + \varepsilon}},$$

$$\bar{P}_j^{(n)}(\boldsymbol{\theta}) = (1-\rho^{n+1}) \nabla_j E(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E(\boldsymbol{\theta}) \frac{\nabla_i E(\boldsymbol{\theta})}{\sqrt{|\nabla_i E(\boldsymbol{\theta})|^2 (1-\rho^{n+1}) + \varepsilon}}.$$

If the iteration number  $n$  is large, (SA-3.2) rapidly becomes

$$\dot{\tilde{\boldsymbol{\theta}}}_j(t) = -\frac{1}{\sqrt{|\nabla_j E(\tilde{\boldsymbol{\theta}}(t))|^2 + \varepsilon}} (\nabla_j E(\tilde{\boldsymbol{\theta}}(t)) + \text{bias}), \quad (\text{SA-3.4})$$

where

$$\text{bias} := \frac{h}{2} \left\{ -\frac{2\rho}{1-\rho} + \frac{1+\rho}{1-\rho} \cdot \frac{\varepsilon}{|\nabla_j E(\tilde{\boldsymbol{\theta}}(t))|^2 + \varepsilon} \right\} \nabla_j \|\nabla E(\tilde{\boldsymbol{\theta}}(t))\|_{1,\varepsilon}. \quad (\text{SA-3.5})$$

## SA-4 Adam with $\varepsilon$ outside the square root

**Definition SA-4.1.** In this section, for some  $\boldsymbol{\theta}^{(0)} \in \mathbb{R}^p$ ,  $\boldsymbol{\nu}^{(0)} = \mathbf{0} \in \mathbb{R}^p$ ,  $\beta, \rho \in (0, 1)$ , let the sequence of  $p$ -vectors  $\{\boldsymbol{\theta}^{(k)}\}_{k \in \mathbb{Z}_{\geq 0}}$  be defined for  $n \geq 0$  by

$$\begin{aligned} \nu_j^{(n+1)} &= \rho \nu_j^{(n)} + (1-\rho) \left( \nabla_j E_n(\boldsymbol{\theta}^{(n)}) \right)^2, \\ m_j^{(n+1)} &= \beta m_j^{(n)} + (1-\beta) \nabla_j E_n(\boldsymbol{\theta}^{(n)}), \\ \theta_j^{(n+1)} &= \theta_j^{(n)} - h \frac{m_j^{(n+1)} / (1-\beta^{n+1})}{\sqrt{\nu_j^{(n+1)} / (1-\rho^{n+1}) + \varepsilon}} \end{aligned}$$

or, rewriting,

$$\theta_j^{(n+1)} = \theta_j^{(n)} - h \frac{\frac{1}{1-\beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1-\beta) \nabla_j E_k(\boldsymbol{\theta}^{(k)})}{\sqrt{\frac{1}{1-\rho^{n+1}} \sum_{k=0}^n \rho^{n-k} (1-\rho) \left( \nabla_j E_k(\boldsymbol{\theta}^{(k)}) \right)^2 + \varepsilon}}. \quad (\text{SA-4.1})$$

Let  $\tilde{\boldsymbol{\theta}}(t)$  be defined as a continuous solution to the piecewise ODE

$$\begin{aligned} \dot{\tilde{\boldsymbol{\theta}}}_j(t) &= -\frac{M_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon} \\ &+ h \left( \frac{M_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \left( 2P_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \bar{P}_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \right)}{2 \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))} - \frac{2L_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \bar{L}_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))}{2 \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon \right)} \right). \end{aligned} \quad (\text{SA-4.2})$$

with the initial condition  $\tilde{\boldsymbol{\theta}}(0) = \boldsymbol{\theta}^{(0)}$ , where  $\mathbf{R}^{(n)}(\boldsymbol{\theta})$ ,  $\mathbf{P}^{(n)}(\boldsymbol{\theta})$ ,  $\bar{\mathbf{P}}^{(n)}(\boldsymbol{\theta})$ ,  $\mathbf{M}^{(n)}(\boldsymbol{\theta})$ ,  $\mathbf{L}^{(n)}(\boldsymbol{\theta})$ ,  $\bar{\mathbf{L}}^{(n)}(\boldsymbol{\theta})$  are

$p$ -dimensional functions with components

$$\begin{aligned}
R_j^{(n)}(\boldsymbol{\theta}) &:= \sqrt{\sum_{k=0}^n \rho^{n-k}(1-\rho)(\nabla_j E_k(\boldsymbol{\theta}))^2 / (1-\rho^{n+1})}, \\
M_j^{(n)}(\boldsymbol{\theta}) &:= \frac{1}{1-\beta^{n+1}} \sum_{k=0}^n \beta^{n-k}(1-\beta) \nabla_j E_k(\boldsymbol{\theta}), \\
L_j^{(n)}(\boldsymbol{\theta}) &:= \frac{1}{1-\beta^{n+1}} \sum_{k=0}^n \beta^{n-k}(1-\beta) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \sum_{l=k}^{n-1} \frac{M_i^{(l)}(\boldsymbol{\theta})}{R_i^{(l)}(\boldsymbol{\theta}) + \varepsilon}, \\
\bar{L}_j^{(n)}(\boldsymbol{\theta}) &:= \frac{1}{1-\beta^{n+1}} \sum_{k=0}^n \beta^{n-k}(1-\beta) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \frac{M_i^{(n)}(\boldsymbol{\theta})}{R_i^{(n)}(\boldsymbol{\theta}) + \varepsilon}, \\
P_j^{(n)}(\boldsymbol{\theta}) &:= \frac{1}{1-\rho^{n+1}} \sum_{k=0}^n \rho^{n-k}(1-\rho) \nabla_j E_k(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \sum_{l=k}^{n-1} \frac{M_i^{(l)}(\boldsymbol{\theta})}{R_i^{(l)}(\boldsymbol{\theta}) + \varepsilon}, \\
\bar{P}_j^{(n)}(\boldsymbol{\theta}) &:= \frac{1}{1-\rho^{n+1}} \sum_{k=0}^n \rho^{n-k}(1-\rho) \nabla_j E_k(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \frac{M_i^{(n)}(\boldsymbol{\theta})}{R_i^{(n)}(\boldsymbol{\theta}) + \varepsilon}.
\end{aligned} \tag{SA-4.3}$$

**Assumption SA-4.2.**

1. For some positive constants  $M_1, M_2, M_3, M_4$  we have

$$\begin{aligned}
\sup_i \sup_k \sup_{\boldsymbol{\theta}} |\nabla_i E_k(\boldsymbol{\theta})| &\leq M_1, \\
\sup_{i,j} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ij} E_k(\boldsymbol{\theta})| &\leq M_2, \\
\sup_{i,j,s} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ijs} E_k(\boldsymbol{\theta})| &\leq M_3, \\
\sup_{i,j,s,r} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ij sr} E_k(\boldsymbol{\theta})| &\leq M_4.
\end{aligned}$$

2. For some  $R > 0$  we have for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$

$$R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) \geq R, \quad \frac{1}{1-\rho^{n+1}} \sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 \geq R^2,$$

where  $\tilde{\boldsymbol{\theta}}(t)$  is defined in [Definition SA-4.1](#).

**Theorem SA-4.3** (Adam with  $\varepsilon$  outside: local error bound). *Suppose [Assumption SA-4.2](#) holds. Then for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$*

$$\left| \tilde{\boldsymbol{\theta}}_j(t_{n+1}) - \tilde{\boldsymbol{\theta}}_j(t_n) + h \frac{\frac{1}{1-\beta^{n+1}} \sum_{k=0}^n \beta^{n-k}(1-\beta) \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k))}{\sqrt{\frac{1}{1-\rho^{n+1}} \sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 + \varepsilon}} \right| \leq C_3 h^3$$

for a positive constant  $C_3$  depending on  $\beta$  and  $\rho$ .

The argument is the same as for [Theorem SA-2.3](#).

**Theorem SA-4.4** (Adam with  $\varepsilon$  outside: global error bound). *Suppose [Assumption SA-4.2](#) holds, and*

$$\frac{1}{1-\rho^{n+1}} \sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\boldsymbol{\theta}^{(k)}) \right)^2 \geq R^2$$

for  $\{\boldsymbol{\theta}^{(k)}\}_{k \in \mathbb{Z}_{\geq 0}}$  defined in [Definition SA-4.1](#). Then there exist positive constants  $d_7, d_8, d_9$  such that for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$

$$\|\mathbf{e}_n\| \leq d_7 e^{d_8 n h} h^2 \quad \text{and} \quad \|\mathbf{e}_{n+1} - \mathbf{e}_n\| \leq d_9 e^{d_8 n h} h^3,$$



where  $\mathbf{e}_n := \tilde{\boldsymbol{\theta}}(t_n) - \boldsymbol{\theta}^{(n)}$ . The constants can be defined as

$$\begin{aligned} d_7 &:= C_3, \\ d_8 &:= \left[ 1 + \frac{M_2\sqrt{p}}{R+\varepsilon} \left( \frac{M_1^2}{R(R+\varepsilon)} + 1 \right) d_7 \right] \sqrt{p}, \\ d_9 &:= C_3 d_8. \end{aligned}$$

*Proof.* Analogously to [Theorem SA-2.4](#), we will prove this by induction over  $n$ .

The base case is  $n = 0$ . Indeed,  $\mathbf{e}_0 = \tilde{\boldsymbol{\theta}}(0) - \boldsymbol{\theta}^{(0)} = \mathbf{0}$ . Then the  $j$ th component of  $\mathbf{e}_1 - \mathbf{e}_0$  is

$$\begin{aligned} [\mathbf{e}_1 - \mathbf{e}_0]_j &= [\mathbf{e}_1]_j = \tilde{\theta}_j(t_1) - \theta_j^{(0)} + \frac{h \nabla_j E_0(\boldsymbol{\theta}^{(0)})}{\left| \nabla_j E_0(\boldsymbol{\theta}^{(0)}) \right| + \varepsilon} \\ &= \tilde{\theta}_j(t_1) - \tilde{\theta}_j(t_0) + \frac{h \nabla_j E_0(\tilde{\boldsymbol{\theta}}(t_0))}{\sqrt{\left( \nabla_j E_0(\tilde{\boldsymbol{\theta}}(t_0)) \right)^2 + \varepsilon}}. \end{aligned}$$

By [Theorem SA-4.3](#), the absolute value of the right-hand side does not exceed  $C_3 h^3$ , which means  $\|\mathbf{e}_1 - \mathbf{e}_0\| \leq C_3 h^3 \sqrt{p}$ . Since  $C_3 \sqrt{p} \leq d_9$ , the base case is proven.

Now suppose that for all  $k = 0, 1, \dots, n-1$  the claim

$$\|\mathbf{e}_k\| \leq d_7 e^{d_8 k h} h^2 \quad \text{and} \quad \|\mathbf{e}_{k+1} - \mathbf{e}_k\| \leq d_9 e^{d_8 k h} h^3$$

is proven. Then

$$\begin{aligned} \|\mathbf{e}_n\| &\stackrel{(a)}{\leq} \|\mathbf{e}_{n-1}\| + \|\mathbf{e}_n - \mathbf{e}_{n-1}\| \leq d_7 e^{d_8(n-1)h} h^2 + d_9 e^{d_8(n-1)h} h^3 \\ &= d_7 e^{d_8(n-1)h} h^2 \left( 1 + \frac{d_9}{d_7} h \right) \stackrel{(b)}{\leq} d_7 e^{d_8(n-1)h} h^2 (1 + d_8 h) \\ &\stackrel{(c)}{\leq} d_7 e^{d_8(n-1)h} h^2 \cdot e^{d_8 h} = d_7 e^{d_8 n h} h^2, \end{aligned}$$

where (a) is by the triangle inequality, (b) is by  $d_9/d_7 \leq d_8$ , in (c) we used  $1 + x \leq e^x$  for all  $x \geq 0$ .

Next, combining [Theorem SA-4.3](#) with [\(SA-4.1\)](#), we have

$$\left| [\mathbf{e}_{n+1} - \mathbf{e}_n]_j \right| \leq C_3 h^3 + h \left| \frac{N'}{\sqrt{D'} + \varepsilon} - \frac{N''}{\sqrt{D''} + \varepsilon} \right|, \quad (\text{SA-4.4})$$

where to simplify notation we put

$$\begin{aligned} N' &:= \frac{1}{1 - \beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1 - \beta) \nabla_j E_k(\boldsymbol{\theta}^{(k)}), \\ N'' &:= \frac{1}{1 - \beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1 - \beta) \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)), \\ D' &:= \frac{1}{1 - \rho^{n+1}} \sum_{k=0}^n \rho^{n-k} (1 - \rho) \left( \nabla_j E_k(\boldsymbol{\theta}^{(k)}) \right)^2, \\ D'' &:= \frac{1}{1 - \rho^{n+1}} \sum_{k=0}^n \rho^{n-k} (1 - \rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2. \end{aligned}$$

Using  $D' \geq R^2$ ,  $D'' \geq R^2$ , we have

$$\left| \frac{1}{\sqrt{D'} + \varepsilon} - \frac{1}{\sqrt{D''} + \varepsilon} \right| = \frac{|D' - D''|}{(\sqrt{D'} + \varepsilon)(\sqrt{D''} + \varepsilon)(\sqrt{D'} + \sqrt{D''})} \leq \frac{|D' - D''|}{2R(R + \varepsilon)^2}. \quad (\text{SA-4.5})$$

But since

$$\begin{aligned}
& \left| \left( \nabla_j E_k(\boldsymbol{\theta}^{(k)}) \right)^2 - \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 \right| \\
&= \left| \nabla_j E_k(\boldsymbol{\theta}^{(k)}) - \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right| \cdot \left| \nabla_j E_k(\boldsymbol{\theta}^{(k)}) + \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right| \\
&\leq 2M_1 \left| \nabla_j E_k(\boldsymbol{\theta}^{(k)}) - \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right| \leq 2M_1 M_2 \sqrt{p} \left\| \boldsymbol{\theta}^{(k)} - \tilde{\boldsymbol{\theta}}(t_k) \right\|,
\end{aligned}$$

we have

$$|D' - D''| \leq \frac{2M_1 M_2 \sqrt{p}}{1 - \rho^{n+1}} \sum_{k=0}^n \rho^{n-k} (1 - \rho) \left\| \boldsymbol{\theta}^{(k)} - \tilde{\boldsymbol{\theta}}(t_k) \right\|. \quad (\text{SA-4.6})$$

Similarly,

$$\begin{aligned}
|N' - N''| &\leq \frac{1}{1 - \beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1 - \beta) \left| \nabla_j E_k(\boldsymbol{\theta}^{(k)}) - \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right| \\
&\leq \frac{1}{1 - \beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1 - \beta) M_2 \sqrt{p} \left\| \boldsymbol{\theta}^{(k)} - \tilde{\boldsymbol{\theta}}(t_k) \right\|.
\end{aligned} \quad (\text{SA-4.7})$$

Combining (SA-4.5), (SA-4.6) and (SA-4.7), we get

$$\begin{aligned}
& \left| \frac{N'}{\sqrt{D'} + \varepsilon} - \frac{N''}{\sqrt{D''} + \varepsilon} \right| \leq |N'| \cdot \left| \frac{1}{\sqrt{D'} + \varepsilon} - \frac{1}{\sqrt{D''} + \varepsilon} \right| + \frac{|N' - N''|}{\sqrt{D''} + \varepsilon} \\
&\leq \frac{1}{1 - \beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1 - \beta) M_1 \cdot \frac{2M_1 M_2 \sqrt{p}}{2R(R + \varepsilon)^2 (1 - \rho^{n+1})} \sum_{k=0}^n \rho^{n-k} (1 - \rho) \left\| \boldsymbol{\theta}^{(k)} - \tilde{\boldsymbol{\theta}}(t_k) \right\| \\
&\quad + \frac{M_2 \sqrt{p}}{(R + \varepsilon)(1 - \beta^{n+1})} \sum_{k=0}^n \beta^{n-k} (1 - \beta) \left\| \boldsymbol{\theta}^{(k)} - \tilde{\boldsymbol{\theta}}(t_k) \right\| \\
&= \frac{M_1^2 M_2 \sqrt{p}}{R(R + \varepsilon)^2 (1 - \rho^{n+1})} \sum_{k=0}^n \rho^{n-k} (1 - \rho) \left\| \boldsymbol{\theta}^{(k)} - \tilde{\boldsymbol{\theta}}(t_k) \right\| \\
&\quad + \frac{M_2 \sqrt{p}}{(R + \varepsilon)(1 - \beta^{n+1})} \sum_{k=0}^n \beta^{n-k} (1 - \beta) \left\| \boldsymbol{\theta}^{(k)} - \tilde{\boldsymbol{\theta}}(t_k) \right\| \\
&\stackrel{(a)}{\leq} \frac{M_1^2 M_2 \sqrt{p}}{R(R + \varepsilon)^2 (1 - \rho^{n+1})} \sum_{k=0}^n \rho^{n-k} (1 - \rho) d_7 e^{d_s k h} h^2 \\
&\quad + \frac{M_2 \sqrt{p}}{(R + \varepsilon)(1 - \beta^{n+1})} \sum_{k=0}^n \beta^{n-k} (1 - \beta) d_7 e^{d_s k h} h^2,
\end{aligned} \quad (\text{SA-4.8})$$

where in (a) we used the induction hypothesis and that the bound on  $\|\mathbf{e}_n\|$  is already proven.

Now note that since  $0 < \rho e^{-d_s h} < \rho$ , we have  $\sum_{k=0}^n (\rho e^{-d_s h})^k \leq \sum_{k=0}^n \rho^k = (1 - \rho^{n+1}) / (1 - \rho)$ , which is rewritten as

$$\frac{1}{1 - \rho^{n+1}} \sum_{k=0}^n \rho^{n-k} (1 - \rho) e^{d_s k h} \leq e^{d_s n h}.$$

By the same logic,

$$\frac{1}{1 - \beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1 - \beta) e^{d_s k h} \leq e^{d_s n h}.$$

Then we can continue (SA-4.8):

$$\left| \frac{N'}{\sqrt{D'} + \varepsilon} - \frac{N''}{\sqrt{D''} + \varepsilon} \right| \leq \frac{M_2 \sqrt{p}}{R + \varepsilon} \left( \frac{M_1^2}{R(R + \varepsilon)} + 1 \right) d_7 e^{d_s n h} h^2 \quad (\text{SA-4.9})$$

Again using  $1 \leq e^{dsnh}$ , we conclude from (SA-4.4) and (SA-4.9) that

$$\|\mathbf{e}_{n+1} - \mathbf{e}_n\| \leq \underbrace{\left( C_3 + \frac{M_2 \sqrt{p}}{R + \varepsilon} \left( \frac{M_1^2}{R(R + \varepsilon)} + 1 \right) d_7 \right)}_{\leq d_9} \sqrt{p} e^{dsnh} h^3,$$

finishing the induction step.  $\square$

## SA-5 Adam with $\varepsilon$ inside the square root

**Definition SA-5.1.** In this section, for some  $\boldsymbol{\theta}^{(0)} \in \mathbb{R}^p$ ,  $\boldsymbol{\nu}^{(0)} = \mathbf{0} \in \mathbb{R}^p$ ,  $\beta, \rho \in (0, 1)$ , let the sequence of  $p$ -vectors  $\left\{ \boldsymbol{\theta}^{(k)} \right\}_{k \in \mathbb{Z}_{\geq 0}}$  be defined for  $n \geq 0$  by

$$\begin{aligned} \nu_j^{(n+1)} &= \rho \nu_j^{(n)} + (1 - \rho) \left( \nabla_j E_n(\boldsymbol{\theta}^{(n)}) \right)^2, \\ m_j^{(n+1)} &= \beta m_j^{(n)} + (1 - \beta) \nabla_j E_n(\boldsymbol{\theta}^{(n)}), \\ \theta_j^{(n+1)} &= \theta_j^{(n)} - h \frac{m_j^{(n+1)} / (1 - \beta^{n+1})}{\sqrt{\nu_j^{(n+1)} / (1 - \rho^{n+1}) + \varepsilon}}. \end{aligned} \tag{SA-5.1}$$

Let  $\tilde{\boldsymbol{\theta}}(t)$  be defined as a continuous solution to the piecewise ODE

$$\begin{aligned} \dot{\tilde{\theta}}_j(t) &= - \frac{M_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))} \\ &+ h \left( \frac{M_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \left( 2P_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \bar{P}_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \right)}{2R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))^3} - \frac{2L_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \bar{L}_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))}{2R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))} \right). \end{aligned} \tag{SA-5.2}$$

with the initial condition  $\tilde{\boldsymbol{\theta}}(0) = \boldsymbol{\theta}^{(0)}$ , where  $\mathbf{R}^{(n)}(\boldsymbol{\theta})$ ,  $\mathbf{P}^{(n)}(\boldsymbol{\theta})$ ,  $\bar{\mathbf{P}}^{(n)}(\boldsymbol{\theta})$ ,  $\mathbf{M}^{(n)}(\boldsymbol{\theta})$ ,  $\mathbf{L}^{(n)}(\boldsymbol{\theta})$ ,  $\bar{\mathbf{L}}^{(n)}(\boldsymbol{\theta})$  are  $p$ -dimensional functions with components

$$\begin{aligned} R_j^{(n)}(\boldsymbol{\theta}) &:= \sqrt{\sum_{k=0}^n \rho^{n-k} (1 - \rho) (\nabla_j E_k(\boldsymbol{\theta}))^2 / (1 - \rho^{n+1}) + \varepsilon}, \\ M_j^{(n)}(\boldsymbol{\theta}) &:= \frac{1}{1 - \beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1 - \beta) \nabla_j E_k(\boldsymbol{\theta}), \\ L_j^{(n)}(\boldsymbol{\theta}) &:= \frac{1}{1 - \beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1 - \beta) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \sum_{l=k}^{n-1} \frac{M_i^{(l)}(\boldsymbol{\theta})}{R_i^{(l)}(\boldsymbol{\theta})}, \\ \bar{L}_j^{(n)}(\boldsymbol{\theta}) &:= \frac{1}{1 - \beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1 - \beta) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \frac{M_i^{(n)}(\boldsymbol{\theta})}{R_i^{(n)}(\boldsymbol{\theta})}, \\ P_j^{(n)}(\boldsymbol{\theta}) &:= \frac{1}{1 - \rho^{n+1}} \sum_{k=0}^n \rho^{n-k} (1 - \rho) \nabla_j E_k(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \sum_{l=k}^{n-1} \frac{M_i^{(l)}(\boldsymbol{\theta})}{R_i^{(l)}(\boldsymbol{\theta})}, \\ \bar{P}_j^{(n)}(\boldsymbol{\theta}) &:= \frac{1}{1 - \rho^{n+1}} \sum_{k=0}^n \rho^{n-k} (1 - \rho) \nabla_j E_k(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \frac{M_i^{(n)}(\boldsymbol{\theta})}{R_i^{(n)}(\boldsymbol{\theta})}. \end{aligned} \tag{SA-5.3}$$

**Assumption SA-5.2.** For some positive constants  $M_1, M_2, M_3, M_4$  we have

$$\sup_i \sup_k \sup_{\boldsymbol{\theta}} |\nabla_i E_k(\boldsymbol{\theta})| \leq M_1,$$

$$\begin{aligned} \sup_{i,j} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ij} E_k(\boldsymbol{\theta})| &\leq M_2, \\ \sup_{i,j,s} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ijs} E_k(\boldsymbol{\theta})| &\leq M_3, \\ \sup_{i,j,s,r} \sup_k \sup_{\boldsymbol{\theta}} |\nabla_{ijsr} E_k(\boldsymbol{\theta})| &\leq M_4. \end{aligned}$$

**Theorem SA-5.3** (Adam with  $\varepsilon$  inside: local error bound). *Suppose Assumption SA-5.2 holds. Then for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$*

$$\left| \tilde{\boldsymbol{\theta}}_j(t_{n+1}) - \tilde{\boldsymbol{\theta}}_j(t_n) + h \frac{\frac{1}{1-\beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1-\beta) \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k))}{\sqrt{\frac{1}{1-\rho^{n+1}} \sum_{k=0}^n \rho^{n-k} (1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 + \varepsilon}} \right| \leq C_4 h^3$$

for a positive constant  $C_4$  depending on  $\beta$  and  $\rho$ .

The argument is the same as for [Theorem SA-2.3](#).

**Theorem SA-5.4** (Adam with  $\varepsilon$  inside: global error bound). *Suppose Assumption SA-5.2 holds for  $\{\boldsymbol{\theta}^{(k)}\}_{k \in \mathbb{Z}_{\geq 0}}$  defined in Definition SA-5.1. Then there exist positive constants  $d_{10}$ ,  $d_{11}$ ,  $d_{12}$  such that for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$*

$$\|\mathbf{e}_n\| \leq d_{10} e^{d_{11} n h} h^2 \quad \text{and} \quad \|\mathbf{e}_{n+1} - \mathbf{e}_n\| \leq d_{12} e^{d_{11} n h} h^3,$$

where  $\mathbf{e}_n := \tilde{\boldsymbol{\theta}}(t_n) - \boldsymbol{\theta}^{(n)}$ . The constants can be defined as

$$\begin{aligned} d_{10} &:= C_4, \\ d_{11} &:= \left[ 1 + \frac{M_2 \sqrt{p}}{\sqrt{\varepsilon}} \left( \frac{M_1^2}{\varepsilon} + 1 \right) d_{10} \right] \sqrt{p}, \\ d_{12} &:= C_4 d_{11}. \end{aligned}$$

The argument is the same as for [Theorem SA-4.4](#).

## SA-6 Technical bounding lemmas

We will need the following lemmas to prove [Theorem SA-2.3](#).

**Lemma SA-6.1.** *Suppose Assumption SA-2.2 holds. Then*

$$\sup_{\boldsymbol{\theta}} \left| P_j^{(n)}(\boldsymbol{\theta}) \right| \leq C_5, \tag{SA-6.1}$$

$$\sup_{\boldsymbol{\theta}} \left| \bar{P}_j^{(n)}(\boldsymbol{\theta}) \right| \leq C_6, \tag{SA-6.2}$$

with constants  $C_5$ ,  $C_6$  defined as follows:

$$\begin{aligned} C_5 &:= p \frac{M_1^2 M_2}{R + \varepsilon} \cdot \frac{\rho}{1 - \rho}, \\ C_6 &:= p \frac{M_1^2 M_2}{R + \varepsilon}. \end{aligned}$$

*Proof of Lemma SA-6.1.* The proof is done in the following simple steps.

**SA-6.2 Proof of (SA-6.1).** This bound is straightforward:

$$\begin{aligned} \sup_{\boldsymbol{\theta}} \left| P_j^{(n)}(\boldsymbol{\theta}) \right| &= \sup_{\boldsymbol{\theta}} \left| \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E_k(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\boldsymbol{\theta})}{R_i^{(l)}(\boldsymbol{\theta}) + \varepsilon} \right| \\ &\leq p \frac{M_1^2 M_2}{R + \varepsilon} (1-\rho) \sum_{k=0}^n \rho^{n-k} (n-k) \leq p \frac{M_1^2 M_2}{R + \varepsilon} (1-\rho) \sum_{k=0}^{\infty} \rho^k k = C_5. \end{aligned}$$

**SA-6.3 Proof of (SA-6.2).** This bound is straightforward:

$$\begin{aligned} \sup_{\boldsymbol{\theta}} \left| \bar{P}_j^{(n)}(\boldsymbol{\theta}) \right| &= \sup_{\boldsymbol{\theta}} \left| \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E_k(\boldsymbol{\theta}) \sum_{i=1}^p \nabla_{ij} E_k(\boldsymbol{\theta}) \frac{\nabla_i E_n(\boldsymbol{\theta})}{R_i^{(n)}(\boldsymbol{\theta}) + \varepsilon} \right| \\ &\leq p \frac{M_1^2 M_2}{R + \varepsilon} (1-\rho) \sum_{k=0}^n \rho^{n-k} \leq p \frac{M_1^2 M_2}{R + \varepsilon} = C_6. \end{aligned}$$

This concludes the proof of [Lemma SA-6.1](#).  $\square$

**Lemma SA-6.4.** Suppose [Assumption SA-2.2](#) holds. Then the first derivative of  $t \mapsto \tilde{\theta}_j(t)$  is uniformly over  $j$  and  $t \in [0, T]$  bounded in absolute value by some positive constant, say  $D_1$ .

*Proof.* This follows immediately from  $h \leq T$ , [\(SA-6.1\)](#), [\(SA-6.2\)](#) and the definition of  $\tilde{\boldsymbol{\theta}}(t)$  given in [\(SA-2.2\)](#).  $\square$

**Lemma SA-6.5.** Suppose [Assumption SA-2.2](#) holds. Then

$$\sup_{t \in [0, T]} \sup_j \left| \left( \nabla_j E_n(\tilde{\boldsymbol{\theta}}(t)) \right) \right| \leq C_7, \quad (\text{SA-6.3})$$

$$\sup_{n, k} \sup_{t \in [t_n, t_{n+1}]} \left| \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t)) \left[ \dot{\tilde{\theta}}_i(t) + \frac{\nabla_i E_n(\tilde{\boldsymbol{\theta}}(t))}{R_i^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon} \right] \right| \leq C_8 h, \quad (\text{SA-6.4})$$

$$\sup_{k \leq n} \sup_{t \in [0, T]} \left| \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\boldsymbol{\theta}}(t))}{R_i^{(l)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon} \right| \leq (n-k) C_9, \quad (\text{SA-6.5})$$

$$\left| \left( P_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \right) \right| \leq C_{10} + C_{14}, \quad (\text{SA-6.6})$$

$$\left| \left( \bar{P}_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \right) \right| \leq C_{15}, \quad (\text{SA-6.7})$$

$$\left| \left( \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t)) \frac{\nabla_i E_n(\tilde{\boldsymbol{\theta}}(t))}{R_i^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon} \right) \right| \leq C_{13}, \quad (\text{SA-6.8})$$

$$\left| \left( \frac{\nabla_j E_n(\tilde{\boldsymbol{\theta}}(t)) \left( 2P_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \bar{P}_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \right)}{2 \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))} \right) \right| \leq C_{17}, \quad (\text{SA-6.9})$$

$$\left| \left( \frac{\sum_{i=1}^p \nabla_{ij} E_n(\tilde{\boldsymbol{\theta}}(t)) \frac{\nabla_i E_n(\tilde{\boldsymbol{\theta}}(t))}{R_i^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon}}{2 \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon \right)} \right) \right| \leq C_{18}, \quad (\text{SA-6.10})$$

with constants  $C_7, C_8, C_9, C_{10}, C_{11}, C_{12}, C_{13}, C_{14}, C_{15}, C_{16}, C_{17}, C_{18}$  defined as follows:

$$\begin{aligned}
C_7 &:= pM_2D_1, \\
C_8 &:= pM_2 \left[ \frac{M_1(2C_5 + C_6)}{2(R + \varepsilon)^2R} + \frac{pM_1M_2}{2(R + \varepsilon)^2} \right], \\
C_9 &:= p \frac{M_1M_2}{R + \varepsilon}, \\
C_{10} &:= D_1p^2 \frac{M_1M_2^2}{R + \varepsilon} \cdot \frac{\rho}{1 - \rho}, \\
C_{11} &:= \frac{D_1pM_1M_2}{R}, \\
C_{12} &:= D_1p^2 \frac{M_1M_3}{R + \varepsilon}, \\
C_{13} &:= C_{12} + pM_2 \left( \frac{D_1pM_2}{R + \varepsilon} + \frac{M_1}{(R + \varepsilon)^2} C_{11} \right) \\
&= \frac{D_1p^2}{R + \varepsilon} \left( M_1M_3 + M_2^2 + \frac{M_1^2M_2^2}{(R + \varepsilon)R} \right), \\
C_{14} &:= M_1C_{13} \frac{\rho}{1 - \rho}, \\
C_{15} &:= \frac{D_1p^2M_1M_2^2}{R + \varepsilon} + \frac{D_1p^2M_1^2M_3}{R + \varepsilon} + \frac{D_1p^2M_1M_2^2}{R + \varepsilon} + \frac{pM_1^2M_2C_{11}}{(R + \varepsilon)^2}, \\
C_{16} &:= \frac{2C_{11}}{R(R + \varepsilon)^3} + \frac{C_{11}}{(R + \varepsilon)^4}, \\
C_{17} &:= \frac{D_1pM_2 \cdot (2C_5 + C_6)}{2(R + \varepsilon)^2R} + \frac{M_1(2(C_{10} + C_{14}) + C_{15})}{2(R + \varepsilon)^2R} + \frac{M_1(2C_5 + C_6)C_{16}}{2}, \\
C_{18} &:= \frac{1}{2(R + \varepsilon)} \left( \frac{p^2D_1M_1M_3}{R + \varepsilon} + \frac{p^2D_1M_2^2}{R + \varepsilon} + \frac{pM_1M_2C_{11}}{(R + \varepsilon)^2} \right) + \frac{1}{2} \cdot \frac{pM_1M_2}{R + \varepsilon} \cdot \frac{C_{11}}{(R + \varepsilon)^2}.
\end{aligned}$$

*Proof of Lemma SA-6.5.* We divide this argument in several steps.

**SA-6.6 Proof of (SA-6.3).** This bound is straightforward:

$$\left| \left( \nabla_j E_n(\tilde{\boldsymbol{\theta}}(t)) \right) \right| = \left| \sum_{i=1}^p \nabla_{ij} E_n(\tilde{\boldsymbol{\theta}}(t)) \dot{\theta}_i(t) \right| \leq C_7.$$

**SA-6.7 Proof of (SA-6.4).** By (SA-2.2) we have for  $t = t_{n+1}^-$

$$\left| \dot{\theta}_j(t) + \frac{\nabla_j E_n(\tilde{\boldsymbol{\theta}}(t))}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon} \right| \leq h \left[ \frac{M_1(2C_5 + C_6)}{2(R + \varepsilon)^2R} + \frac{pM_1M_2}{2(R + \varepsilon)^2} \right],$$

giving (SA-6.4) immediately.

**SA-6.8 Proof of (SA-6.5).** This bound follows from the assumptions immediately.

**SA-6.9 Proof of (SA-6.6).** We will prove this by bounding the two terms in the expression

$$\begin{aligned}
& \frac{d}{dt} P_j^{(n)}(\tilde{\theta}(t)) \\
&= \sum_{k=0}^n \rho^{n-k} (1-\rho) \sum_{u=1}^p \nabla_{ju} E_k(\tilde{\theta}(t)) \dot{\theta}_u(t) \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \\
&+ \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E_k(\tilde{\theta}(t)) \sum_{i=1}^p \frac{d}{dt} \left\{ \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right\}.
\end{aligned} \tag{SA-6.11}$$

It is easily shown that the first term in (SA-6.11) is bounded in absolute value by  $C_{10}$ :

$$\begin{aligned}
& \left| \sum_{k=0}^n \rho^{n-k} (1-\rho) \sum_{u=1}^p \nabla_{ju} E_k(\tilde{\theta}(t)) \dot{\theta}_u(t) \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right| \\
&\leq D_1 p^2 \frac{M_1 M_2^2}{R + \varepsilon} (1-\rho) \sum_{k=0}^n \rho^k k \\
&\leq D_1 p^2 \frac{M_1 M_2^2}{R + \varepsilon} (1-\rho) \sum_{k=0}^{\infty} \rho^k k \\
&= C_{10}.
\end{aligned}$$

For the proof of (SA-6.6), it is left to show that the second term in (SA-6.11) is bounded in absolute value by  $C_{14}$ .

To bound  $\sum_{i=1}^p \frac{d}{dt} \left\{ \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right\}$ , we can use

$$\begin{aligned}
& \left| \sum_{i=1}^p \frac{d}{dt} \left\{ \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right\} \right| \\
&\leq \left| \sum_{i=1}^p \frac{d}{dt} \left\{ \nabla_{ij} E_k(\tilde{\theta}(t)) \right\} \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right| \\
&+ \left| \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{d}{dt} \left\{ \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right\} \right|
\end{aligned}$$

By the Cauchy-Schwarz inequality applied twice,

$$\begin{aligned}
& \left| \sum_{i=1}^p \frac{d}{dt} \left\{ \nabla_{ij} E_k(\tilde{\theta}(t)) \right\} \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right| \\
&\leq \sqrt{\sum_{i=1}^p \sum_{s=1}^p \left( \nabla_{ijs} E_k(\tilde{\theta}(t)) \right)^2} \sqrt{\sum_{u=1}^p \dot{\theta}_u(t)^2} \sqrt{\sum_{i=1}^p \left| \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right|^2} \\
&\leq M_{3p} \cdot D_1 \sqrt{p} \cdot \sqrt{\sum_{i=1}^p \left| \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right|^2} \leq (n-k) C_{12}.
\end{aligned}$$

Next, for any  $n$  and  $j$

$$\begin{aligned} \left| \frac{d}{dt} R_j^{(n)}(\tilde{\theta}(t)) \right| &= \frac{1}{R_j^{(n)}(\tilde{\theta}(t))} \left| \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E_k(\tilde{\theta}(t)) \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \dot{\theta}_i(t) \right| \\ &\leq \frac{1}{R_j^{(n)}(\tilde{\theta}(t))} D_1 p M_1 M_2 \sum_{k=0}^n \rho^{n-k} (1-\rho) \leq C_{11}. \end{aligned} \quad (\text{SA-6.12})$$

This gives

$$\begin{aligned} \left| \frac{d}{dt} \left\{ \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right\} \right| &\leq \frac{\left| \sum_{s=1}^p \nabla_{is} E_l(\tilde{\theta}(t)) \dot{\theta}_s(t) \right|}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} + \frac{\left| \nabla_i E_l(\tilde{\theta}(t)) \right| \cdot \left| \frac{d}{dt} R_i^{(l)}(\tilde{\theta}(t)) \right|}{\left( R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon \right)^2} \\ &\leq \frac{D_1 p M_2}{R + \varepsilon} + \frac{M_1}{(R + \varepsilon)^2} C_{11}. \end{aligned}$$

We have obtained

$$\left| \sum_{i=1}^p \frac{d}{dt} \left\{ \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right\} \right| \leq (n-k) C_{13}. \quad (\text{SA-6.13})$$

This gives a bound on the second term in (SA-6.11):

$$\begin{aligned} &\left| \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E_k(\tilde{\theta}(t)) \sum_{i=1}^p \frac{d}{dt} \left\{ \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right\} \right| \\ &\leq M_1 \sum_{k=0}^n \rho^{n-k} (1-\rho) (n-k) C_{13} \leq C_{14}, \end{aligned}$$

concluding the proof of (SA-6.6).

**SA-6.10 Proof of (SA-6.7).** We will prove this by bounding the four terms in the expression

$$\begin{aligned} &\frac{d}{dt} \left\{ \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E_k(\tilde{\theta}(t)) \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \frac{\nabla_i E_n(\tilde{\theta}(t))}{R_i^{(n)}(\tilde{\theta}(t)) + \varepsilon} \right\} \\ &= \text{Term1} + \text{Term2} + \text{Term3} + \text{Term4}, \end{aligned}$$

where

Term1

$$:= \sum_{k=0}^n \rho^{n-k} (1-\rho) \frac{d}{dt} \left\{ \nabla_j E_k(\tilde{\theta}(t)) \right\} \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \frac{\nabla_i E_n(\tilde{\theta}(t))}{R_i^{(n)}(\tilde{\theta}(t)) + \varepsilon},$$

Term2

$$:= \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E_k(\tilde{\theta}(t)) \sum_{i=1}^p \frac{d}{dt} \left\{ \nabla_{ij} E_k(\tilde{\theta}(t)) \right\} \frac{\nabla_i E_n(\tilde{\theta}(t))}{R_i^{(n)}(\tilde{\theta}(t)) + \varepsilon},$$

Term3



$$:= \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E_k(\tilde{\theta}(t)) \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \frac{\frac{d}{dt} \left\{ \nabla_i E_n(\tilde{\theta}(t)) \right\}}{R_i^{(n)}(\tilde{\theta}(t)) + \varepsilon},$$

Term4

$$:= - \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E_k(\tilde{\theta}(t)) \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \frac{\nabla_i E_n(\tilde{\theta}(t)) \frac{d}{dt} R_i^{(n)}(\tilde{\theta}(t))}{\left( R_i^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^2}.$$

To bound Term1, use  $\left| \frac{d}{dt} \left\{ \nabla_j E_k(\tilde{\theta}(t)) \right\} \right| \leq D_1 p M_2$ , giving

$$|\text{Term1}| \leq \frac{D_1 p^2 M_1 M_2^2}{R + \varepsilon} \sum_{k=0}^n \rho^{n-k} (1-\rho) \leq \frac{D_1 p^2 M_1 M_2^2}{R + \varepsilon}.$$

To bound Term2, use  $\left| \frac{d}{dt} \left\{ \nabla_{ij} E_k(\tilde{\theta}(t)) \right\} \right| \leq D_1 p M_3$ , giving

$$|\text{Term2}| \leq \frac{D_1 p^2 M_1^2 M_3}{R + \varepsilon} \sum_{k=0}^n \rho^{n-k} (1-\rho) \leq \frac{D_1 p^2 M_1^2 M_3}{R + \varepsilon}.$$

To bound Term3, use  $\left| \frac{d}{dt} \left\{ \nabla_i E_n(\tilde{\theta}(t)) \right\} \right| \leq D_1 p M_2$ , giving

$$|\text{Term3}| \leq \frac{D_1 p^2 M_1 M_2^2}{R + \varepsilon} \sum_{k=0}^n \rho^{n-k} (1-\rho) \leq \frac{D_1 p^2 M_1 M_2^2}{R + \varepsilon}.$$

To bound Term4, use (SA-6.12), giving

$$|\text{Term4}| \leq \frac{p M_1^2 M_2 C_{11}}{(R + \varepsilon)^2} \sum_{k=0}^n \rho^{n-k} (1-\rho) \leq \frac{p M_1^2 M_2 C_{11}}{(R + \varepsilon)^2}.$$

**SA-6.11 Proof of (SA-6.8).** This is proven in (SA-6.13).

**SA-6.12 Proof of (SA-6.9).** (SA-6.12) gives

$$\left| \frac{d}{dt} \left\{ \frac{1}{R_j^{(n)}(\tilde{\theta}(t))} \right\} \right| = \frac{\left| \frac{d}{dt} R_j^{(n)}(\tilde{\theta}(t)) \right|}{R_j^{(n)}(\tilde{\theta}(t))^2} \leq \frac{C_{11}}{R^2}, \quad (\text{SA-6.14})$$

$$\left| \frac{d}{dt} \left\{ \frac{1}{R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon} \right\} \right| = \frac{\left| \frac{d}{dt} R_j^{(n)}(\tilde{\theta}(t)) \right|}{\left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^2} \leq \frac{C_{11}}{(R + \varepsilon)^2}, \quad (\text{SA-6.15})$$

$$\left| \frac{d}{dt} \left\{ \frac{1}{\left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^2} \right\} \right| = \frac{2 \left| \frac{d}{dt} R_j^{(n)}(\tilde{\theta}(t)) \right|}{\left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^3} \leq \frac{2C_{11}}{(R + \varepsilon)^3}. \quad (\text{SA-6.16})$$

Combining two bounds above, we have

$$\left| \frac{d}{dt} \left\{ \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^{-2} R_j^{(n)}(\tilde{\theta}(t))^{-1} \right\} \right|$$

$$\leq \frac{\left| \frac{d}{dt} \left\{ \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^{-2} \right\} \right|}{R_j^{(n)}(\tilde{\theta}(t))} + \frac{\left| \frac{d}{dt} \left\{ R_j^{(n)}(\tilde{\theta}(t))^{-1} \right\} \right|}{\left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^2} \leq C_{16}.$$

We are ready to bound

$$\begin{aligned} & \left| \left( \frac{\nabla_j E_n(\tilde{\theta}(t)) \left( 2P_j^{(n)}(\tilde{\theta}(t)) + \bar{P}_j^{(n)}(\tilde{\theta}(t)) \right)}{2 \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\theta}(t))} \right) \right| \\ & \leq \left| \frac{\left( \nabla_j E_n(\tilde{\theta}(t)) \right) \left( 2P_j^{(n)}(\tilde{\theta}(t)) + \bar{P}_j^{(n)}(\tilde{\theta}(t)) \right)}{2 \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\theta}(t))} \right| + \\ & \quad + \left| \frac{\nabla_j E_n(\tilde{\theta}(t)) \left( 2P_j^{(n)}(\tilde{\theta}(t)) + \bar{P}_j^{(n)}(\tilde{\theta}(t)) \right)}{2 \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\theta}(t))} \right| \\ & \quad + \left| \frac{\nabla_j E_n(\tilde{\theta}(t)) \left( 2P_j^{(n)}(\tilde{\theta}(t)) + \bar{P}_j^{(n)}(\tilde{\theta}(t)) \right)}{2} \right| \\ & \quad \times \left| \left( \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^{-2} R_j^{(n)}(\tilde{\theta}(t))^{-1} \right) \right| \leq C_{17}. \end{aligned}$$

**SA-6.13 Proof of (SA-6.10).** Since

$$\left| \sum_{i=1}^p \nabla_{ij} E_n(\tilde{\theta}(t)) \frac{\nabla_i E_n(\tilde{\theta}(t))}{R_i^{(n)}(\tilde{\theta}(t)) + \varepsilon} \right| \leq \frac{pM_1M_2}{R + \varepsilon}$$

and, as we have already seen in the argument for (SA-6.7),

$$\left| \left( \sum_{i=1}^p \nabla_{ij} E_n(\tilde{\theta}(t)) \frac{\nabla_i E_n(\tilde{\theta}(t))}{R_i^{(n)}(\tilde{\theta}(t)) + \varepsilon} \right) \right| \leq \frac{p^2 D_1 M_1 M_3}{R + \varepsilon} + \frac{p^2 D_1 M_2^2}{R + \varepsilon} + \frac{pM_1M_2C_{11}}{(R + \varepsilon)^2},$$

we are ready to bound

$$\left| \left( \frac{\sum_{i=1}^p \nabla_{ij} E_n(\tilde{\theta}(t)) \frac{\nabla_i E_n(\tilde{\theta}(t))}{R_i^{(n)}(\tilde{\theta}(t)) + \varepsilon}}{2 \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)} \right) \right| \leq C_{18}.$$

The proof of Lemma SA-6.5 is concluded.  $\square$

**Lemma SA-6.14.** Suppose Assumption SA-2.2 holds. Then the second derivative of  $t \mapsto \tilde{\theta}_j(t)$  is uniformly over  $j$  and  $t \in [0, T]$  bounded in absolute value by some positive constant, say  $D_2$ .

*Proof.* This follows from the definition of  $\tilde{\boldsymbol{\theta}}(t)$  given in (SA-2.2),  $h \leq T$  and that the first derivatives of all three terms in (SA-2.2) are bounded by Lemma SA-6.5.  $\square$

**Lemma SA-6.15.** *Suppose Assumption SA-2.2 holds. Then*

$$\left| \left( \nabla_j E_n(\tilde{\boldsymbol{\theta}}(t)) \right)^{\cdot\cdot} \right| \leq C_{19}, \quad (\text{SA-6.17})$$

$$\left| \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \right)^{\cdot\cdot} \right| \leq C_{20}, \quad (\text{SA-6.18})$$

$$\left| \left( \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon \right)^{-2} \right)^{\cdot\cdot} \right| \leq C_{21}, \quad (\text{SA-6.19})$$

$$\left| \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))^{-1} \right)^{\cdot\cdot} \right| \leq C_{22}, \quad (\text{SA-6.20})$$

$$\left| \left( \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon \right)^{-2} R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))^{-1} \right)^{\cdot\cdot} \right| \leq C_{23}, \quad (\text{SA-6.21})$$

$$\left| \left( \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t)) \sum_{l=k}^{n-1} \frac{\nabla_l E_l(\tilde{\boldsymbol{\theta}}(t))}{R_i^{(l)}(\tilde{\boldsymbol{\theta}}(t)) + \varepsilon} \right)^{\cdot\cdot} \right| \leq (n-k)C_{24}, \quad (\text{SA-6.22})$$

with constants  $C_{19}, C_{20}, C_{21}, C_{22}, C_{23}, C_{24}$  defined as follows:

$$\begin{aligned} C_{19} &:= p^2 M_3 D_1^2 + p M_2 D_2, \\ C_{20} &:= \frac{C_{11}}{R^2} p M_1 M_2 D_1 + \frac{1}{R} p^2 M_2^2 D_1^2 + \frac{1}{R} p^2 M_1 M_3 D_1^2 + \frac{1}{R} p M_1 M_2 D_2, \\ C_{21} &:= \frac{6C_{11}^2}{(R+\varepsilon)^4} + \frac{2C_{20}}{(R+\varepsilon)^3}, \\ C_{22} &:= \frac{2C_{11}^2}{R^3} + \frac{C_{20}}{R^2}, \\ C_{23} &:= \frac{C_{21}}{R} + \frac{4C_{11}^2}{R^2(R+\varepsilon)^3} + \frac{C_{22}}{(R+\varepsilon)^2}, \\ C_{24} &:= p \left[ \frac{2C_{11}(D_1 M_2^2 p + D_1 M_1 M_3 p)}{(R+\varepsilon)^2} + M_1 M_2 \left( \frac{2C_{11}^2}{(R+\varepsilon)^3} + \frac{C_{20}}{(R+\varepsilon)^2} \right) \right. \\ &\quad \left. + \frac{2D_1^2 M_2 M_3 p^2 + M_2(D_1^2 M_3 p^2 + D_2 M_2 p) + M_1(D_1^2 M_4 p^2 + D_2 M_3 p)}{R+\varepsilon} \right]. \end{aligned}$$

*Proof of Lemma SA-6.15.* We divide this argument in several steps.

**SA-6.16 Proof of (SA-6.17).** This bound is straightforward:

$$\left| \left( \nabla_j E_n(\tilde{\boldsymbol{\theta}}(t)) \right)^{\cdot\cdot} \right| = \left| \sum_{i=1}^p \sum_{s=1}^p \nabla_{ijs} E_n(\tilde{\boldsymbol{\theta}}(t)) \dot{\theta}_s(t) \dot{\theta}_i(t) + \sum_{i=1}^p \nabla_{ij} E_n(\tilde{\boldsymbol{\theta}}(t)) \ddot{\theta}_i(t) \right| \leq C_{19}.$$

**SA-6.17 Proof of (SA-6.18).** Note that

$$\begin{aligned} \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t)) \right)^{\cdot\cdot} &= \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))^{-1} \right)^{\cdot} \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t)) \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t)) \dot{\theta}_i(t) \\ &\quad + R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t))^{-1} \sum_{k=0}^n \rho^{n-k} (1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t)) \right)^{\cdot} \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t)) \dot{\theta}_i(t) \end{aligned}$$

$$\begin{aligned}
& + R_j^{(n)}(\tilde{\theta}(t))^{-1} \sum_{k=0}^n \rho^{n-k}(1-\rho) \nabla_j E_k(\tilde{\theta}(t)) \sum_{i=1}^p \left( \nabla_{ij} E_k(\tilde{\theta}(t)) \right) \dot{\tilde{\theta}}_i(t) \\
& + R_j^{(n)}(\tilde{\theta}(t))^{-1} \sum_{k=0}^n \rho^{n-k}(1-\rho) \nabla_j E_k(\tilde{\theta}(t)) \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \ddot{\tilde{\theta}}_i(t),
\end{aligned}$$

giving by (SA-6.14)

$$\begin{aligned}
\left| \left( R_j^{(n)}(\tilde{\theta}(t)) \right)^{\cdot\cdot} \right| & \leq \frac{C_{11}}{R^2} p M_1 M_2 D_1 \sum_{k=0}^n \rho^{n-k}(1-\rho) + \frac{1}{R} p^2 M_2^2 D_1^2 \sum_{k=0}^n \rho^{n-k}(1-\rho) \\
& + \frac{1}{R} p^2 M_1 M_3 D_1^2 \sum_{k=0}^n \rho^{n-k}(1-\rho) + \frac{1}{R} p M_1 M_2 D_2 \sum_{k=0}^n \rho^{n-k}(1-\rho) \\
& \leq C_{20}.
\end{aligned}$$

**SA-6.18 Proof of (SA-6.19).** Note that

$$\left( \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^{-2} \right)^{\cdot\cdot} = \frac{6 \left( \left( R_j^{(n)}(\tilde{\theta}(t)) \right)^{\cdot} \right)^2}{\left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^4} - \frac{2 \left( R_j^{(n)}(\tilde{\theta}(t)) \right)^{\cdot\cdot}}{\left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^3},$$

giving by (SA-6.12) and (SA-6.18)

$$\left| \left( \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^{-2} \right)^{\cdot\cdot} \right| \leq C_{21}.$$

**SA-6.19 Proof of (SA-6.20).** The bound follows from (SA-6.12), (SA-6.18) and

$$\left( R_j^{(n)}(\tilde{\theta}(t))^{-1} \right)^{\cdot\cdot} = \frac{2 \left( \left( R_j^{(n)}(\tilde{\theta}(t)) \right)^{\cdot} \right)^2}{R_j^{(n)}(\tilde{\theta}(t))^3} - \frac{\left( R_j^{(n)}(\tilde{\theta}(t)) \right)^{\cdot\cdot}}{R_j^{(n)}(\tilde{\theta}(t))^2}.$$

**SA-6.20 Proof of (SA-6.21).** Putting  $a := \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^{-2}$ ,  $b := R_j^{(n)}(\tilde{\theta}(t))^{-1}$ , use

$$\begin{aligned}
|a| & \leq \frac{1}{(R+\varepsilon)^2}, \quad |b| \leq \frac{1}{R}, \\
|\dot{a}| & \leq \frac{2C_{11}}{(R+\varepsilon)^3}, \quad |\dot{b}| \leq \frac{C_{11}}{R^2}, \\
|\ddot{a}| & \leq C_{21}, \quad |\ddot{b}| \leq C_{22},
\end{aligned}$$

and

$$(ab)^{\cdot\cdot} = \ddot{a}b + 2\dot{a}\dot{b} + a\ddot{b}.$$

**SA-6.21 Proof of (SA-6.22).** Putting

$$\begin{aligned}
a & := \nabla_{ij} E_k(\tilde{\theta}(t)), \\
b & := \nabla_i E_l(\tilde{\theta}(t)), \\
c & := \left( R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon \right)^{-1},
\end{aligned}$$

we have

$$\begin{aligned} |a| &\leq M_2, & |\dot{a}| &\leq pM_3D_1, & |\ddot{a}| &\leq p^2M_4D_1^2 + pM_3D_2, \\ |b| &\leq M_1, & |\dot{b}| &\leq pM_2D_1, & |\ddot{b}| &\leq p^2M_3D_1^2 + pM_2D_2, \\ |c| &\leq \frac{1}{R+\varepsilon}, & |\dot{c}| &\leq \frac{C_{11}}{(R+\varepsilon)^2}, & |\ddot{c}| &\leq \frac{2C_{11}^2}{(R+\varepsilon)^3} + \frac{C_{20}}{(R+\varepsilon)^2}. \end{aligned}$$

(SA-6.22) follows.

The proof of Lemma SA-6.15 is concluded.  $\square$

**Lemma SA-6.22.** *Suppose Assumption SA-2.2 holds. Then the third derivative of  $t \mapsto \tilde{\theta}_j(t)$  is uniformly over  $j$  and  $t \in [0, T]$  bounded in absolute value by some positive constant, say  $D_3$ .*

*Proof.* By (SA-6.5), (SA-6.13) and (SA-6.22)

$$\begin{aligned} \left| \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right| &\leq (n-k)C_9, \\ \left| \left( \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right)' \right| &\leq (n-k)C_{13}, \\ \left| \left( \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t))}{R_i^{(l)}(\tilde{\theta}(t)) + \varepsilon} \right)'' \right| &\leq (n-k)C_{24}. \end{aligned}$$

From the definition of  $t \mapsto P_j^{(n)}(\tilde{\theta}(t))$ , it means that its derivatives up to order two are bounded. Similarly, the same is true for  $t \mapsto \bar{P}_j^{(n)}(\tilde{\theta}(t))$ .

It follows from (SA-6.19) and its proof that the derivatives up to order two of

$$t \mapsto \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^{-2} R_j^{(n)}(\tilde{\theta}(t))^{-1}$$

are also bounded.

These considerations give the boundedness of the second derivative of the term

$$t \mapsto \frac{\nabla_j E_n(\tilde{\theta}(t)) \left( 2P_j^{(n)}(\tilde{\theta}(t)) + \bar{P}_j^{(n)}(\tilde{\theta}(t)) \right)}{2 \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\theta}(t))}$$

in (SA-2.2). The boundedness of the second derivatives of the other two terms is shown analogously. By (SA-2.2) and since  $h \leq T$ , this means

$$\sup_j \sup_{t \in [0, T]} \left| \ddot{\tilde{\theta}}_j(t) \right| \leq D_3$$

for some positive constant  $D_3$ .  $\square$

## SA-7 Proof of Theorem SA-2.3

**Lemma SA-7.1.** *Suppose Assumption SA-2.2 holds. Then for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$ ,  $k \in \{0, 1, \dots, n-1\}$  we have*

$$\left| \nabla_j E_k(\tilde{\theta}(t_k)) - \nabla_j E_k(\tilde{\theta}(t_n)) \right| \leq C_7(n-k)h \quad (\text{SA-7.1})$$

*Proof.* (SA-7.1) follows from the mean value theorem applied  $n - k$  times.  $\square$

**Lemma SA-7.2.** *In the setting of Lemma SA-7.1, for any  $l \in \{k, k + 1, \dots, n - 1\}$  we have*

$$\left| \nabla_j E_k(\tilde{\theta}(t_l)) - \nabla_j E_k(\tilde{\theta}(t_{l+1})) - h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_n)) \frac{\nabla_i E_l(\tilde{\theta}(t_n))}{R_i^{(l)}(\tilde{\theta}(t_n)) + \varepsilon} \right| \leq (C_{19}/2 + C_8 + (n - l - 1)C_{13})h^2.$$

*Proof.* By the Taylor expansion of  $t \mapsto \nabla_j E_k(\tilde{\theta}(t))$  on the segment  $[t_l, t_{l+1}]$  at  $t_{l+1}$  on the left

$$\left| \nabla_j E_k(\tilde{\theta}(t_l)) - \nabla_j E_k(\tilde{\theta}(t_{l+1})) + h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_{l+1})) \dot{\tilde{\theta}}_i(t_{l+1}^-) \right| \leq \frac{C_{19}}{2} h^2.$$

Combining this with (SA-6.4) gives

$$\left| \nabla_j E_k(\tilde{\theta}(t_l)) - \nabla_j E_k(\tilde{\theta}(t_{l+1})) - h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_{l+1})) \frac{\nabla_i E_l(\tilde{\theta}(t_{l+1}))}{R_i^{(l)}(\tilde{\theta}(t_{l+1})) + \varepsilon} \right| \leq (C_{19}/2 + C_8)h^2. \quad (\text{SA-7.2})$$

Now applying the mean-value theorem  $n - l - 1$  times, we have

$$\begin{aligned} & \left| \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_{l+1})) \frac{\nabla_i E_l(\tilde{\theta}(t_{l+1}))}{R_i^{(l)}(\tilde{\theta}(t_{l+1})) + \varepsilon} - \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_{l+2})) \frac{\nabla_i E_l(\tilde{\theta}(t_{l+2}))}{R_i^{(l)}(\tilde{\theta}(t_{l+2})) + \varepsilon} \right| \leq C_{13}h, \\ & \dots \\ & \left| \sum_{i=1}^p \nabla_{ij} E_l(\tilde{\theta}(t_{n-1})) \frac{\nabla_i E_k(\tilde{\theta}(t_{n-1}))}{R_i^{(l)}(\tilde{\theta}(t_{n-1})) + \varepsilon} - \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_n)) \frac{\nabla_i E_l(\tilde{\theta}(t_n))}{R_i^{(l)}(\tilde{\theta}(t_n)) + \varepsilon} \right| \leq C_{13}h, \end{aligned}$$

and in particular

$$\left| \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_{l+1})) \frac{\nabla_i E_l(\tilde{\theta}(t_{l+1}))}{R_i^{(l)}(\tilde{\theta}(t_{l+1})) + \varepsilon} - \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_n)) \frac{\nabla_i E_l(\tilde{\theta}(t_n))}{R_i^{(l)}(\tilde{\theta}(t_n)) + \varepsilon} \right| \leq (n - l - 1)C_{13}h.$$

Combining this with (SA-7.2), we conclude the proof of Lemma SA-7.2.  $\square$

**Lemma SA-7.3.** *In the setting of Lemma SA-7.1,*

$$\left| \nabla_j E_k(\tilde{\theta}(t_k)) - \nabla_j E_k(\tilde{\theta}(t_n)) - h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_n)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t_n))}{R_i^{(l)}(\tilde{\theta}(t_n)) + \varepsilon} \right| \leq \left( (n - k)(C_{19}/2 + C_8) + \frac{(n - k)(n - k - 1)}{2} C_{13} \right) h^2.$$

*Proof.* Fix  $n \in \mathbb{Z}_{\geq 0}$ .

Note that

$$\left| \nabla_j E_k(\tilde{\theta}(t_k)) - \nabla_j E_k(\tilde{\theta}(t_n)) - h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_n)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\theta}(t_n))}{R_i^{(l)}(\tilde{\theta}(t_n)) + \varepsilon} \right|$$

$$\begin{aligned}
&= \left| \sum_{l=k}^{n-1} \left\{ \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_l)) - \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_{l+1})) - h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t_n)) \frac{\nabla_i E_l(\tilde{\boldsymbol{\theta}}(t_n))}{R_i^{(l)}(\tilde{\boldsymbol{\theta}}(t_n)) + \varepsilon} \right\} \right| \\
&\leq \sum_{l=k}^{n-1} \left| \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_l)) - \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_{l+1})) - h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t_n)) \frac{\nabla_i E_l(\tilde{\boldsymbol{\theta}}(t_n))}{R_i^{(l)}(\tilde{\boldsymbol{\theta}}(t_n)) + \varepsilon} \right| \\
&\stackrel{(a)}{\leq} \sum_{l=k}^{n-1} (C_{19}/2 + C_8 + (n-l-1)C_{13})h^2 = \left( (n-k)(C_{19}/2 + C_8) + \frac{(n-k)(n-k-1)}{2}C_{13} \right) h^2,
\end{aligned}$$

where (a) is by [Lemma SA-7.2](#).  $\square$

**Lemma SA-7.4.** *Suppose [Assumption SA-2.2](#) holds. Then for all  $n \in \{0, 1, \dots, \lfloor T/h \rfloor\}$*

$$\left| \sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 - R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))^2 \right| \leq C_{25}h \quad (\text{SA-7.3})$$

and

$$\left| \sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 - R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))^2 - 2hP_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) \right| \leq C_{26}h^2 \quad (\text{SA-7.4})$$

with  $C_{25}$  and  $C_{26}$  defined as follows:

$$\begin{aligned}
C_{25}(\rho) &:= 2M_1C_7 \frac{\rho}{1-\rho}, \\
C_{26}(\rho) &:= M_1|C_{19} + 2C_8 - C_{13}| \frac{\rho}{1-\rho} \\
&\quad + \left( M_1C_{13} + |C_{19} + 2C_8 - C_{13}|C_9 + \frac{(C_{19} + 2C_8 - C_{13})^2}{4} \right) \frac{\rho(1+\rho)}{(1-\rho)^2} \\
&\quad + \left( C_{13}C_9 + \frac{C_{13}}{2}|C_{19} + 2C_8 - C_{13}| \right) \frac{\rho(1+4\rho+\rho^2)}{(1-\rho)^3} + \frac{C_{13}^2}{4} \cdot \frac{\rho(1+11\rho+11\rho^2+\rho^3)}{(1-\rho)^4}.
\end{aligned}$$

*Proof.* Note that

$$\begin{aligned}
&\left| \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 - \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_n)) \right)^2 \right| \\
&\leq \left| \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) - \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_n)) \right| \cdot \left| \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) + \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_n)) \right| \\
&\stackrel{(a)}{\leq} C_7(n-k)h \cdot 2M_1,
\end{aligned}$$

where (a) is by [\(SA-7.1\)](#). Using the triangle inequality, we can conclude

$$\begin{aligned}
&\left| \sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 - R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))^2 \right| \\
&\leq 2M_1C_7h(1-\rho) \sum_{k=0}^n (n-k)\rho^{n-k} = 2M_1C_7h(1-\rho) \sum_{k=0}^n k\rho^k = 2M_1C_7 \frac{\rho}{1-\rho} h.
\end{aligned}$$

[\(SA-7.3\)](#) is proven.

We continue by showing

$$\begin{aligned}
& \left| \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 - \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_n)) \right)^2 \right. \\
& \quad \left. - 2 \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_n)) h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t_n)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\boldsymbol{\theta}}(t_n))}{R_i^{(l)}(\tilde{\boldsymbol{\theta}}(t_n)) + \varepsilon} \right| \\
& \leq 2M_1 \left( (n-k)(C_{19}/2 + C_8) + \frac{(n-k)(n-k-1)}{2} C_{13} \right) h^2 \\
& \quad + 2(n-k)C_9 \left( (n-k)(C_{19}/2 + C_8) + \frac{(n-k)(n-k-1)}{2} C_{13} \right) h^3 \\
& \quad + \left( (n-k)(C_{19}/2 + C_8) + \frac{(n-k)(n-k-1)}{2} C_{13} \right)^2 h^4.
\end{aligned} \tag{SA-7.5}$$

To prove this, use

$$|a^2 - b^2 - 2bKh| \leq 2|b| \cdot |a - b - Kh| + 2|K| \cdot h \cdot |a - b - Kh| + (a - b - Kh)^2$$

with

$$a := \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)), \quad b := \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_n)), \quad K := \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t_n)) \sum_{l=k}^{n-1} \frac{\nabla_i E_l(\tilde{\boldsymbol{\theta}}(t_n))}{R_i^{(l)}(\tilde{\boldsymbol{\theta}}(t_n)) + \varepsilon},$$

and bounding

$$\begin{aligned}
|a - b - Kh| & \stackrel{(a)}{\leq} \left( (n-k)(C_{19}/2 + C_8) + \frac{(n-k)(n-k-1)}{2} C_{13} \right) h^2, \\
|b| & \leq M_1, \quad |K| \leq (n-k)C_9,
\end{aligned}$$

where (a) is by Lemma SA-7.3. (SA-7.5) is proven.

We turn to the proof of (SA-7.4). By (SA-7.5) and the triangle inequality

$$\begin{aligned}
& \left| \sum_{k=0}^n \rho^{n-k} (1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 - R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))^2 - 2hP_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) \right| \\
& \leq (1-\rho) \sum_{k=0}^n \rho^{n-k} \left( \text{Poly}_1(n-k)h^2 + \text{Poly}_2(n-k)h^3 + \text{Poly}_3(n-k)h^4 \right) \\
& = (1-\rho) \sum_{k=0}^n \rho^k \left( \text{Poly}_1(k)h^2 + \text{Poly}_2(k)h^3 + \text{Poly}_3(k)h^4 \right),
\end{aligned}$$

where

$$\begin{aligned}
\text{Poly}_1(k) & := 2M_1 \left( k(C_{19}/2 + C_8) + \frac{k(k-1)}{2} C_{13} \right) = M_1 C_{13} k^2 + M_1 (C_{19} + 2C_8 - C_{13}) k, \\
\text{Poly}_2(k) & := 2kC_9 \left( k(C_{19}/2 + C_8) + \frac{k(k-1)}{2} C_{13} \right) = C_{13} C_9 k^3 + (C_{19} + 2C_8 - C_{13}) C_9 k^2, \\
\text{Poly}_3(k) & := \left( k(C_{19}/2 + C_8) + \frac{k(k-1)}{2} C_{13} \right)^2 \\
& = \frac{C_{13}^2}{4} k^4 + \frac{C_{13}}{2} (C_{19} + 2C_8 - C_{13}) k^3 + \frac{1}{4} (C_{19} + 2C_8 - C_{13})^2 k^2.
\end{aligned}$$

It is left to combine this with

$$\sum_{k=0}^n k \rho^k \leq \sum_{k=0}^{\infty} k \rho^k = \frac{\rho}{(1-\rho)^2},$$



$$\begin{aligned}
\sum_{k=0}^n k^2 \rho^k &\leq \sum_{k=0}^{\infty} k^2 \rho^k = \frac{\rho(1+\rho)}{(1-\rho)^3}, \\
\sum_{k=0}^n k^3 \rho^k &\leq \sum_{k=0}^{\infty} k^3 \rho^k = \frac{\rho(1+4\rho+\rho^2)}{(1-\rho)^4}, \\
\sum_{k=0}^n k^4 \rho^k &\leq \sum_{k=0}^{\infty} k^4 \rho^k = \frac{\rho(1+11\rho+11\rho^2+\rho^3)}{(1-\rho)^5}.
\end{aligned}$$

This gives

$$\begin{aligned}
&\left| \sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2 - R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))^2 - 2hP_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) \right| \\
&\leq \left( M_1 C_{13} \frac{\rho(1+\rho)}{(1-\rho)^2} + M_1 |C_{19} + 2C_8 - C_{13}| \frac{\rho}{1-\rho} \right) h^2 \\
&\quad + \left( C_{13} C_9 \frac{\rho(1+4\rho+\rho^2)}{(1-\rho)^3} + |C_{19} + 2C_8 - C_{13}| C_9 \frac{\rho(1+\rho)}{(1-\rho)^2} \right) h^3 \\
&\quad + \left( \frac{C_{13}^2}{4} \cdot \frac{\rho(1+11\rho+11\rho^2+\rho^3)}{(1-\rho)^4} + \frac{C_{13}}{2} |C_{19} + 2C_8 - C_{13}| \frac{\rho(1+4\rho+\rho^2)}{(1-\rho)^3} \right. \\
&\quad \left. + \frac{1}{4} (C_{19} + 2C_8 - C_{13})^2 \frac{\rho(1+\rho)}{(1-\rho)^2} \right) h^4 \\
&\stackrel{(a)}{\leq} \left[ M_1 |C_{19} + 2C_8 - C_{13}| \frac{\rho}{1-\rho} \right. \\
&\quad + \left( M_1 C_{13} + |C_{19} + 2C_8 - C_{13}| C_9 + \frac{(C_{19} + 2C_8 - C_{13})^2}{4} \right) \frac{\rho(1+\rho)}{(1-\rho)^2} \\
&\quad + \left( C_{13} C_9 + \frac{C_{13}}{2} |C_{19} + 2C_8 - C_{13}| \right) \frac{\rho(1+4\rho+\rho^2)}{(1-\rho)^3} \\
&\quad \left. + \frac{C_{13}^2}{4} \cdot \frac{\rho(1+11\rho+11\rho^2+\rho^3)}{(1-\rho)^4} \right] h^2,
\end{aligned}$$

where in (a) we used that  $h < 1$ . (SA-7.4) is proven.  $\square$

**Lemma SA-7.5.** *Suppose Assumption SA-2.2 holds. Then*

$$\begin{aligned}
&\left| \left( \sqrt{\sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) \right)^2} + \varepsilon \right)^{-1} - \left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) + \varepsilon \right)^{-1} \right. \\
&\quad \left. + h \frac{P_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))}{\left( R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))} \right| \leq \frac{C_{25}(\rho)^2 + R^2 C_{26}(\rho)}{2R^3(R+\varepsilon)^2} h^2.
\end{aligned}$$

*Proof.* Note that if  $a \geq R^2$ ,  $b \geq R^2$ , we have

$$\begin{aligned}
&\left| \frac{1}{\sqrt{a} + \varepsilon} - \frac{1}{\sqrt{b} + \varepsilon} + \frac{a-b}{2(\sqrt{b} + \varepsilon)^2 \sqrt{b}} \right| \\
&= \frac{(a-b)^2}{2\sqrt{b}(\sqrt{b} + \varepsilon)(\sqrt{a} + \varepsilon)(\sqrt{a} + \sqrt{b})} \underbrace{\left\{ \frac{1}{\sqrt{b} + \varepsilon} + \frac{1}{\sqrt{a} + \sqrt{b}} \right\}}_{\leq 2/R}
\end{aligned}$$

$$\leq \frac{(a-b)^2}{2R^3(R+\varepsilon)^2}.$$

By the triangle inequality,

$$\begin{aligned} \left| \frac{1}{\sqrt{a+\varepsilon}} - \frac{1}{\sqrt{b+\varepsilon}} + \frac{c}{2(\sqrt{b+\varepsilon})^2\sqrt{b}} \right| &\leq \frac{(a-b)^2}{2R^3(R+\varepsilon)^2} + \frac{|a-b-c|}{2(\sqrt{b+\varepsilon})^2\sqrt{b}} \\ &\leq \frac{(a-b)^2}{2R^3(R+\varepsilon)^2} + \frac{|a-b-c|}{2R(R+\varepsilon)^2} \end{aligned}$$

Apply this with

$$\begin{aligned} a &:= \sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\theta}(t_k)) \right)^2, \\ b &:= R_j^{(n)}(\tilde{\theta}(t_n))^2, \\ c &:= 2hP_j^{(n)}(\tilde{\theta}(t_n)) \end{aligned}$$

and use bounds

$$|a-b| \leq 2M_1C_7 \frac{\rho}{1-\rho} h, \quad |a-b-c| \leq C_{26}(\rho)h^2$$

by Lemma SA-7.4. □

**SA-7.6.** We are finally ready to prove Theorem SA-2.3.

*Proof of Theorem SA-2.3.* By (SA-6.9) and (SA-6.10), the first derivative of the function

$$t \mapsto \left( \frac{\nabla_j E_n(\tilde{\theta}(t)) \left( 2P_j^{(n)}(\tilde{\theta}(t)) + \bar{P}_j^{(n)}(\tilde{\theta}(t)) \right)}{2 \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\theta}(t))} - \frac{\sum_{i=1}^p \nabla_{ij} E_n(\tilde{\theta}(t)) \frac{\nabla_i E_n(\tilde{\theta}(t))}{R_i^{(n)}(\tilde{\theta}(t)) + \varepsilon}}{2 \left( R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon \right)} \right)$$

is bounded in absolute value by a positive constant  $C_{27} = C_{17} + C_{18}$ . By (SA-2.2), this means

$$\left| \ddot{\theta}_j(t) + \frac{d}{dt} \left( \frac{\nabla_j E_n(\tilde{\theta}(t))}{R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon} \right) \right| \leq C_{27}h.$$

Combining this with

$$\left| \tilde{\theta}_j(t_{n+1}) - \tilde{\theta}_j(t_n) - \dot{\theta}_j(t_n^+)h - \frac{\ddot{\theta}_j(t_n^+)}{2}h^2 \right| \leq \frac{D_3}{6}$$

by Taylor expansion, we get

$$\begin{aligned} &\left| \tilde{\theta}_j(t_{n+1}) - \tilde{\theta}_j(t_n) - \dot{\theta}_j(t_n^+)h + \frac{h^2}{2} \cdot \frac{d}{dt} \left( \frac{\nabla_j E_n(\tilde{\theta}(t))}{R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon} \right) \Big|_{t=t_n^+} \right| \\ &\leq \left( \frac{D_3}{6} + \frac{C_{27}}{2} \right) h^3. \end{aligned} \tag{SA-7.6}$$

Using

$$\left| \dot{\theta}_j(t_n) + \frac{\nabla_j E_n(\tilde{\theta}(t_n))}{R_j^{(n)}(\tilde{\theta}(t_n)) + \varepsilon} \right| \leq C_{28}h$$

with  $C_{28}$  defined as

$$C_{28} := \frac{M_1(2C_5 + C_6)}{2(R + \varepsilon)^2 R} + \frac{pM_1M_2}{2(R + \varepsilon)^2}$$

by (SA-2.2), and calculating the derivative, it is easy to show

$$\left| \frac{d}{dt} \left( \frac{\nabla_j E_n(\tilde{\theta}(t))}{R_j^{(n)}(\tilde{\theta}(t)) + \varepsilon} \right) \Big|_{t=t_n^+} - \text{FrDer} \right| \leq C_{29}h \quad (\text{SA-7.7})$$

for a positive constant  $C_{29}$ , where

$$\begin{aligned} \text{FrDer} &:= \frac{\text{FrDerNum}}{\left( R_j^{(n)}(\tilde{\theta}(t_n)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\theta}(t_n))} \\ \text{FrDerNum} &:= \nabla_j E_n(\tilde{\theta}(t_n)) \bar{P}_j^{(n)}(\tilde{\theta}(t_n)) \\ &\quad - \left( R_j^{(n)}(\tilde{\theta}(t_n)) + \varepsilon \right) R_j^{(n)}(\tilde{\theta}(t_n)) \sum_{i=1}^p \nabla_{ij} E_n(\tilde{\theta}(t_n)) \frac{\nabla_i E_n(\tilde{\theta}(t_n))}{R_i^{(n)}(\tilde{\theta}(t_n)) + \varepsilon}, \\ C_{29} &:= \left\{ \frac{pM_2}{R + \varepsilon} + \frac{M_1^2 M_2 p}{(R + \varepsilon)^2 R} \right\} C_{28}. \end{aligned}$$

From (SA-7.6) and (SA-7.7), by the triangle inequality

$$\left| \tilde{\theta}_j(t_{n+1}) - \tilde{\theta}_j(t_n) - \dot{\tilde{\theta}}_j(t_n^+)h + \frac{h^2}{2} \text{FrDer} \right| \leq \left( \frac{D_3}{6} + \frac{C_{27} + C_{29}}{2} \right) h^3,$$

which, using (SA-2.2), is rewritten as

$$\left| \tilde{\theta}_j(t_{n+1}) - \tilde{\theta}_j(t_n) + h \frac{\nabla_j E_n(\tilde{\theta}(t_n))}{R_j^{(n)}(\tilde{\theta}(t_n)) + \varepsilon} - h^2 \frac{\nabla_j E_n(\tilde{\theta}(t_n)) P_j^{(n)}(\tilde{\theta}(t_n))}{\left( R_j^{(n)}(\tilde{\theta}(t_n)) + \varepsilon \right)^2 R_j^{(n)}(\tilde{\theta}(t_n))} \right| \leq \left( \frac{D_3}{6} + \frac{C_{27} + C_{29}}{2} \right) h^3.$$

It is left to combine this with Lemma SA-7.5, giving the assertion of the theorem with

$$C_1 = \frac{D_3}{6} + \frac{C_{27} + C_{29}}{2} + M_1 \frac{C_{25}^2 + R^2 C_{26}}{2R^3(R + \varepsilon)^2}. \quad \square$$

## SA-8 Numerical experiments

**SA-8.1 Models.** We use small modifications of default Keras Resnet-50 and Resnet-101 architectures<sup>1</sup> for training on CIFAR-10 and CIFAR-100 (since image sizes are not the same as Imagenet), after verifying their correctness. The first convolution layer `conv1` has  $3 \times 3$  kernel, stride 1 and “same” padding. Then comes batch normalization, and relu. Max pooling is removed, and otherwise `conv2_x` to `conv5_x` are as described in [2], see Table 1 there (downsampling is performed by the first convolution of each bottleneck block, same as in this original paper, not the middle one as in version 1.5<sup>2</sup>; all convolution layers have learned biases). After `conv5` there is global average pooling, 10 or 100-way fully connected layer (for CIFAR-10 and CIFAR-100 respectively), and softmax.

<sup>1</sup><https://github.com/keras-team/keras/blob/v2.13.1/keras/applications/resnet.py>

<sup>2</sup>[https://catalog.ngc.nvidia.com/orgs/nvidia/resources/resnet\\_50\\_v1\\_5\\_for\\_pytorch](https://catalog.ngc.nvidia.com/orgs/nvidia/resources/resnet_50_v1_5_for_pytorch)

**SA-8.2 Data augmentation.** We subtract the per-pixel mean and divide by standard deviation, and we use the data augmentation scheme from [3], following [2], section 4.2. We take inspiration and some code snippets from [4] (though we do not use their models). During each pass over the training dataset, each  $32 \times 32$  initial image is padded evenly with zeros so that it becomes  $36 \times 36$ , then random crop is applied so that the picture becomes  $32 \times 32$  again, and finally random (probability 0.5) horizontal (left to right) flip is used.

**SA-8.3 Experiment details.** In experiments whose results are reported in Figures 4 and 5 of the main paper, we train for more than 3600 epochs and stop training when the train accuracy is near-perfect (Figure SA-1) and the testing accuracy does not significantly improve (Figure SA-2). Therefore, the maximal test accuracies are the final ones reached, and the maximal perturbed one-norms, after excluding the initial fall at the beginning of training, are at peaks of the “hills” on the norm curves (Figure SA-2).

Additional evidence (for ResNet-101 on CIFAR-100 and with hyperparameters different from the ones in Figures 4 and 5) is provided in Figures SA-3 and SA-4.

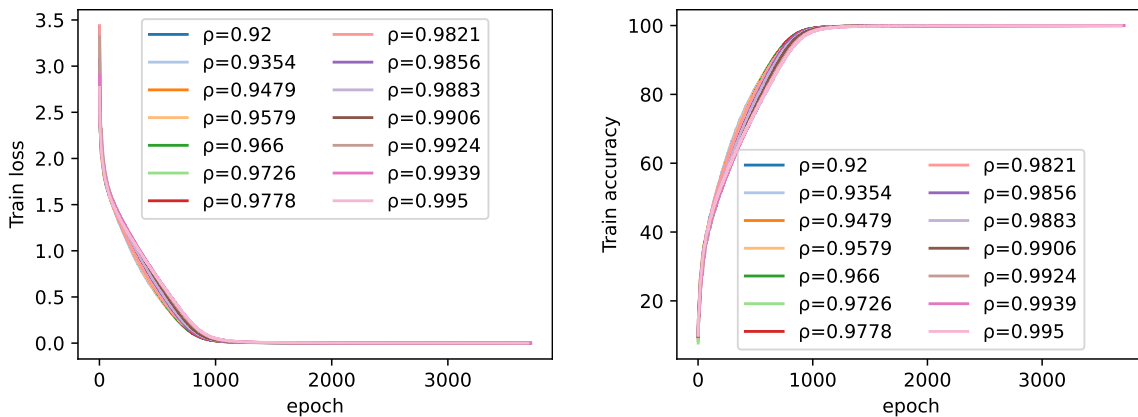


Figure SA-1: Train loss and train accuracy curves for full-batch Adam, ResNet-50 on CIFAR-10,  $\beta = 0.99$ ,  $\varepsilon = 10^{-8}$ ,  $h = 7.5 \cdot 10^{-5}$ .

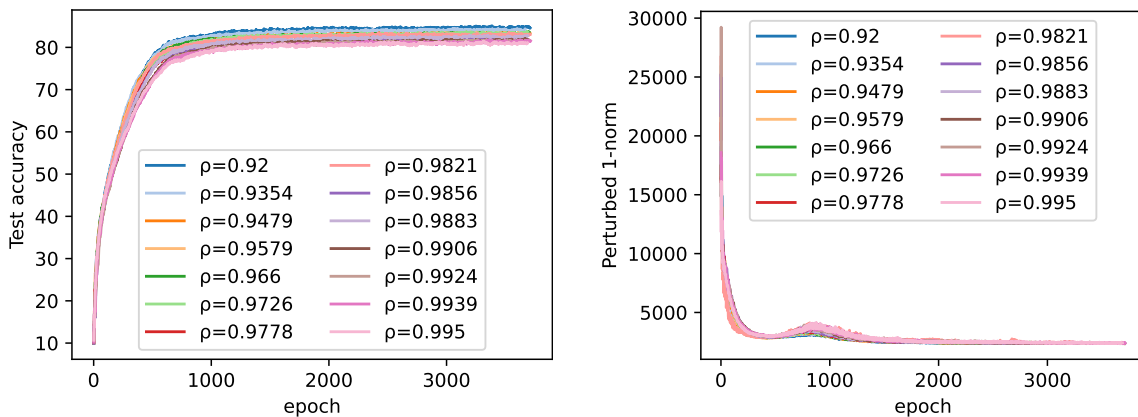


Figure SA-2: Test accuracy and  $\|\nabla E\|_{1,\varepsilon}$  after each epoch. The setting is the same as in Figure SA-1.

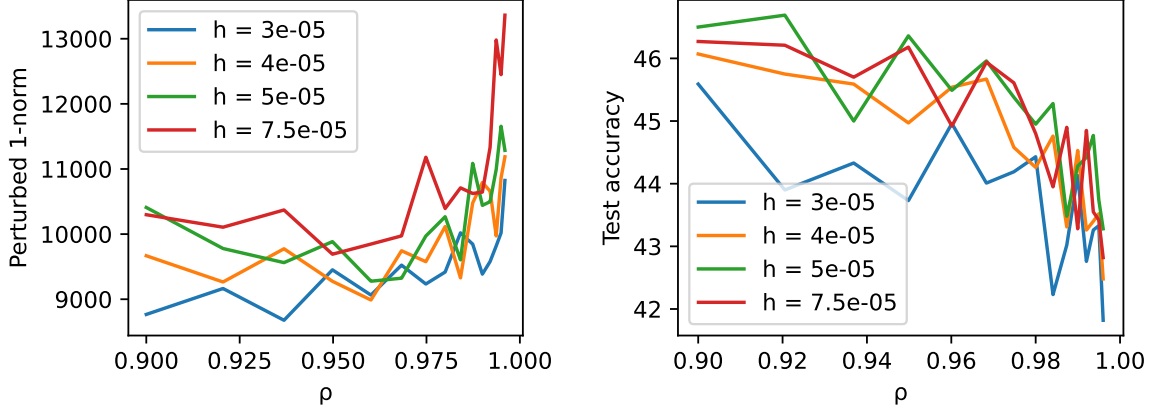


Figure SA-3: Resnet-101 on CIFAR-100 trained with full-batch Adam,  $\varepsilon = 10^{-8}$ ,  $\beta = 0.95$ . As  $\rho$  increases, the perturbed one-norm seems to rise and the test accuracy seems to fall (in the stable regime of training). Both metrics are calculated as in Figures 4 and 5 of the main paper.

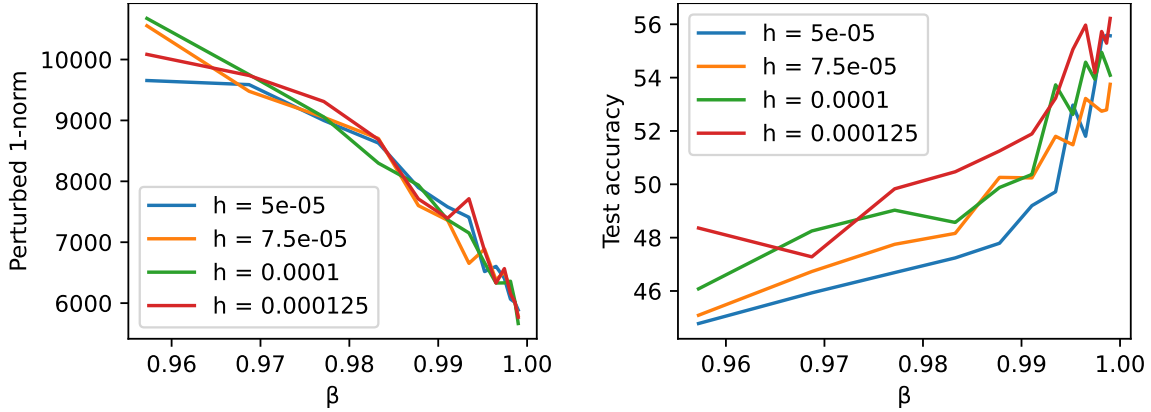


Figure SA-4: Resnet-101 on CIFAR-100 trained with full-batch Adam,  $\rho = 0.99$ ,  $\varepsilon = 10^{-8}$ . The perturbed one-norm seems to fall as  $\beta$  increases, and the test accuracy seems to rise. Both metrics are calculated as in Figures 4 and 5 of the main paper.

## SA-9 Adam with $\varepsilon$ inside the square root: informal derivation

**Result SA-9.1.** For  $n \in \{0, 1, 2, \dots\}$  we have

$$\begin{aligned} \tilde{\theta}_j(t_{n+1}) &= \tilde{\theta}_j(t_n) - h \frac{M_j^{(n)}(\tilde{\theta}(t_n))}{R_j^{(n)}(\tilde{\theta}(t_n))} \\ &\quad + h^2 \left( \frac{M_j^{(n)}(\tilde{\theta}(t_n)) P_j^{(n)}(\tilde{\theta}(t_n))}{R_j^{(n)}(\tilde{\theta}(t_n))^3} - \frac{L_j^{(n)}(\tilde{\theta}(t_n))}{R_j^{(n)}(\tilde{\theta}(t_n))} \right) + O(h^3). \end{aligned} \tag{SA-9.1}$$

*Derivation.* We take

$$\tilde{\theta}_j(t_{n+1}) = \tilde{\theta}_j(t_n) - h \frac{M_j^{(n)}(\tilde{\theta}(t_n))}{R_j^{(n)}(\tilde{\theta}(t_n))} + O(h^2)$$

for granted. Using this and the Taylor series, we can write

$$\nabla_j E_k(\tilde{\theta}(t_{n-1}))$$

$$\begin{aligned}
&= \nabla_j E_k(\tilde{\theta}(t_n)) + \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_n)) \left\{ \tilde{\theta}_i(t_{n-1}) - \tilde{\theta}_i(t_n) \right\} + O(h^2) \\
&= \nabla_j E_k(\tilde{\theta}(t_n)) + h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_n)) \frac{M_j^{(n-1)}(\tilde{\theta}(t_{n-1}))}{R_j^{(n-1)}(\tilde{\theta}(t_{n-1}))} + O(h^2) \\
&= \nabla_j E_k(\tilde{\theta}(t_n)) + h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_n)) \frac{M_j^{(n-1)}(\tilde{\theta}(t_n))}{R_j^{(n-1)}(\tilde{\theta}(t_n))} + O(h^2),
\end{aligned}$$

where in the last equality we just replaced  $t_{n-1}$  with  $t_n$  in the  $h$ -term since it only affects higher-order terms. Now doing this again for step  $n-1$  instead of step  $n$ , we will have

$$\begin{aligned}
&\nabla_j E_k(\tilde{\theta}(t_{n-2})) \\
&= \nabla_j E_k(\tilde{\theta}(t_{n-1})) + h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_{n-1})) \frac{M_j^{(n-2)}(\tilde{\theta}(t_{n-1}))}{R_j^{(n-2)}(\tilde{\theta}(t_{n-1}))} + O(h^2) \\
&= \nabla_j E_k(\tilde{\theta}(t_{n-1})) + h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_{n-1})) \frac{M_j^{(n-2)}(\tilde{\theta}(t_n))}{R_j^{(n-2)}(\tilde{\theta}(t_n))} + O(h^2),
\end{aligned}$$

where in the last equality we again replaced  $t_{n-1}$  with  $t_n$  since it only affects higher-order terms. Proceeding like this and adding the resulting equations, we have for  $n \in \{0, 1, \dots\}$ ,  $k \in \{0, \dots, n-1\}$  that

$$\begin{aligned}
&\nabla_j E_k(\tilde{\theta}(t_k)) \\
&= \nabla_j E_k(\tilde{\theta}(t_n)) + h \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_n)) \sum_{l=k}^{n-1} \frac{M_i^{(l)}(\tilde{\theta}(t_n))}{R_i^{(l)}(\tilde{\theta}(t_n))} + O(h^2),
\end{aligned}$$

where we ignored the fact that  $n-k$  is not bounded (we will get away with this because of exponential averaging). Hence, taking the square of this formal power series,

$$\begin{aligned}
&\rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\theta}(t_k)) \right)^2 = \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\theta}(t_n)) \right)^2 \\
&\quad + h \cdot 2\rho^{n-k}(1-\rho) \nabla_j E_k(\tilde{\theta}(t_n)) \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\theta}(t_n)) \sum_{l=k}^{n-1} \frac{M_i^{(l)}(\tilde{\theta}(t_n))}{R_i^{(l)}(\tilde{\theta}(t_n))} + O(h^2).
\end{aligned}$$

Summing up over  $k$ , we have

$$\frac{1}{1-\rho^{n+1}} \sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\theta}(t_k)) \right)^2 + \varepsilon = R_j^{(n)}(\tilde{\theta}(t_n))^2 + 2hP_j^{(n)}(\tilde{\theta}(t_n)) + O(h^2),$$

which, using the expression for the inverse square root  $(\sum_{r=0}^{\infty} a_r h^r)^{-1/2}$  of a formal power series  $\sum_{r=0}^{\infty} a_r h^r$ , gives us

$$\begin{aligned}
&\left( \sqrt{\frac{1}{1-\rho^{n+1}} \sum_{k=0}^n \rho^{n-k}(1-\rho) \left( \nabla_j E_k(\tilde{\theta}(t_k)) \right)^2 + \varepsilon} \right)^{-1} \\
&= \frac{1}{R_j^{(n)}(\tilde{\theta}(t_n))} - h \frac{P_j^{(n)}(\tilde{\theta}(t_n))}{R_j^{(n)}(\tilde{\theta}(t_n))^3} + O(h^2).
\end{aligned}$$

Similarly,

$$\begin{aligned}
& \frac{1}{1-\beta^{n+1}} \sum_{k=0}^n (1-\beta)\beta^{n-k} \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_k)) = \frac{1}{1-\beta^{n+1}} \sum_{k=0}^n (1-\beta)\beta^{n-k} \nabla_j E_k(\tilde{\boldsymbol{\theta}}(t_n)) \\
& + \frac{h}{1-\beta^{n+1}} \sum_{k=0}^n (1-\beta)\beta^{n-k} \sum_{i=1}^p \nabla_{ij} E_k(\tilde{\boldsymbol{\theta}}(t_n)) \sum_{l=k}^{n-1} \frac{M_i^{(l)}(\tilde{\boldsymbol{\theta}}(t_n))}{R_i^{(l)}(\tilde{\boldsymbol{\theta}}(t_n))} + O(h^2) \\
& = M_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) + hL_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) + O(h^2).
\end{aligned}$$

We conclude

$$\begin{aligned}
\tilde{\boldsymbol{\theta}}_j(t_{n+1}) &= \tilde{\boldsymbol{\theta}}_j(t_n) - h \left( M_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) + hL_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n)) + O(h^2) \right) \\
&\times \left( \frac{1}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))} - h \frac{P_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))^3} + O(h^2) \right) + O(h^3) \\
&= \tilde{\boldsymbol{\theta}}_j(t_n) - h \frac{M_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))} \\
&+ h^2 \left( \frac{M_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))P_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))^3} - \frac{L_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))}{R_j^{(n)}(\tilde{\boldsymbol{\theta}}(t_n))} \right) + O(h^3). \quad \square
\end{aligned}$$

**Result SA-9.2.** For  $t_n \leq t < t_{n+1}$ , the modified equation is (SA-5.2).

*Derivation.* Assume that the modified flow for  $t_n \leq t < t_{n+1}$  satisfies  $\dot{\tilde{\boldsymbol{\theta}}} = \tilde{\mathbf{f}}(\tilde{\boldsymbol{\theta}}(t))$  where

$$\tilde{\mathbf{f}}(\boldsymbol{\theta}) = \mathbf{f}(\boldsymbol{\theta}) + h\mathbf{f}_1(\boldsymbol{\theta}) + O(h^2).$$

By Taylor expansion, we have

$$\begin{aligned}
\tilde{\boldsymbol{\theta}}(t_{n+1}) &= \tilde{\boldsymbol{\theta}}(t_n) + h\dot{\tilde{\boldsymbol{\theta}}}(t_n^+) + \frac{h^2}{2}\ddot{\tilde{\boldsymbol{\theta}}}(t_n^+) + O(h^3) \\
&= \tilde{\boldsymbol{\theta}}(t_n) + h \left[ \mathbf{f}(\tilde{\boldsymbol{\theta}}(t_n)) + h\mathbf{f}_1(\tilde{\boldsymbol{\theta}}(t_n)) + O(h^2) \right] \\
&+ \frac{h^2}{2} \left[ \nabla \mathbf{f}(\tilde{\boldsymbol{\theta}}(t_n)) \mathbf{f}(\tilde{\boldsymbol{\theta}}(t_n)) + O(h) \right] + O(h^3) \tag{SA-9.2} \\
&= \tilde{\boldsymbol{\theta}}(t_n) + h\mathbf{f}(\tilde{\boldsymbol{\theta}}(t_n)) + h^2 \left[ \mathbf{f}_1(\tilde{\boldsymbol{\theta}}(t_n)) + \frac{\nabla \mathbf{f}(\tilde{\boldsymbol{\theta}}(t_n)) \mathbf{f}(\tilde{\boldsymbol{\theta}}(t_n))}{2} \right] + O(h^3).
\end{aligned}$$

Using Lemma SA-9.1 and equating the terms before the corresponding powers of  $h$  in (SA-9.1) and (SA-9.2), we obtain

$$\begin{aligned}
f_j(\boldsymbol{\theta}) &= -\frac{M_j^{(n)}(\boldsymbol{\theta})}{R_j^{(n)}(\boldsymbol{\theta})}, \\
f_{1,j}(\boldsymbol{\theta}) &= -\frac{1}{2} \sum_{i=1}^p \nabla_i f_j(\boldsymbol{\theta}) f_i(\boldsymbol{\theta}) + \frac{M_j^{(n)}(\boldsymbol{\theta})P_j^{(n)}(\boldsymbol{\theta})}{R_j^{(n)}(\boldsymbol{\theta})^3} - \frac{L_j^{(n)}(\boldsymbol{\theta})}{R_j^{(n)}(\boldsymbol{\theta})}. \tag{SA-9.3}
\end{aligned}$$

It is left to find  $\nabla_i f_j(\boldsymbol{\theta})$ . Using

$$\nabla_i R_j^{(n)}(\boldsymbol{\theta}) = \frac{\sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_{ij} E_k(\boldsymbol{\theta}) \nabla_j E_k(\boldsymbol{\theta})}{(1-\rho^{n+1})R_j^{(n)}(\boldsymbol{\theta})},$$

$$\nabla_i M_j^{(n)}(\boldsymbol{\theta}) = \frac{\sum_{k=0}^n \beta^{n-k} (1-\beta) \nabla_{ij} E_k(\boldsymbol{\theta})}{1-\beta^{n+1}}$$

we have

$$\begin{aligned} & \nabla_i \left( -\frac{M_j^{(n)}(\boldsymbol{\theta})}{R_j^{(n)}(\boldsymbol{\theta})} \right) \\ &= -\frac{\frac{R_j^{(n)}(\boldsymbol{\theta})^2}{1-\beta^{n+1}} \sum_{k=0}^n \beta^{n-k} (1-\beta) \nabla_{ij} E_k(\boldsymbol{\theta}) - \frac{M_j^{(n)}(\boldsymbol{\theta})}{1-\rho^{n+1}} \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_{ij} E_k(\boldsymbol{\theta}) \nabla_j E_k(\boldsymbol{\theta})}{R_j^{(n)}(\boldsymbol{\theta})^3} \\ &= -\frac{\sum_{k=0}^n \beta^{n-k} (1-\beta) \nabla_{ij} E_k(\boldsymbol{\theta})}{(1-\beta^{n+1}) R_j^{(n)}(\boldsymbol{\theta})} + \frac{M_j^{(n)}(\boldsymbol{\theta}) \sum_{k=0}^n \rho^{n-k} (1-\rho) \nabla_{ij} E_k(\boldsymbol{\theta}) \nabla_j E_k(\boldsymbol{\theta})}{(1-\rho^{n+1}) R_j^{(n)}(\boldsymbol{\theta})^3} \end{aligned}$$

Inserting this into (SA-9.3) concludes the proof.  $\square$

## References

- [1] Avrajit Ghosh, He Lyu, Xitong Zhang, and Rongrong Wang. “Implicit regularization in Heavy-ball momentum accelerated stochastic gradient descent”. In: *The Eleventh International Conference on Learning Representations*. 2023. URL: <https://openreview.net/forum?id=ZzdBhtEH9yB>.
- [2] 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*. 2016, pp. 770–778.
- [3] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. “Deeply-supervised nets”. In: *Artificial intelligence and statistics*. Pmlr. 2015, pp. 562–570.
- [4] Chia-Hung Yuan. *Training CIFAR-10 with TensorFlow2(TF2)*. <https://github.com/lionelmessi6410/tensorflow2-cifar>. 2021.