Stochastic Gradient Descent: Going As Fast As Possible But Not Faster

Alice Schoenauer-Sebag¹, Marc Schoenauer², Michèle Sebag² (1) Altschuler & Wu lab, UCSF (2)TAO, CNRS – INRIA – LRI – Université Paris-Sud



Jan. 9th, 2018





Position of the talk

Given a black-box optimization algorithm, and its hyper-parameters:

Offline setting

 Find good hp values, used during the whole run e.g. kernel; tree depth; regularization weights...

Online setting

- Adjust the hp values during the run
- ▶ NB: when considering a *fixed* schedule (e.g. $\eta_t = \frac{\eta}{\sqrt{t}}$), you are in the offline setting.

AUTO-ML

(ロ) (個) (E) (E) E のQ

Offline vs Online hyper-parameter setting

PROS/CONS

- More computationally expensive by orders of magnitude
- Potentially huge gains



Focusing on online adaptation

Large Scale Machine Learning

Increasing adoption of Stochastic Gradient Descent

$$\theta_{t+1} = \theta_t - \eta_t \nabla \mathcal{L}_\theta$$

• Adjusting η_t ?

Robbins & Munro 51; Jacobs 88; Nesterov 83 Sutskever et al. 13; Duchi et al. 13; Kingma & Ba 14

イロン イロン イヨン イヨン 三日

- Vanishing and Exploding gradients
- Most usually

$$\eta_t \sim rac{1}{\sqrt{t}} ext{ or } rac{1}{t}$$

Claims and Goal

Goal: Go as fast as possible

No the second second

Issue: There are explosions

Claim

- We cannot predict them
- ▶ We can observe them
- We can cure them

Remark

 If explosions are controlled, we can go faster speed-up and brake with reasonable cost

(ロ) (同) (E) (E) (E)

SALeRA: Safe Agnostic Learning Rate Adaptation

1st ingredient: Going faster

Compare gradient directions with an agnostic cumulative path

2nd ingredient: Catastrophic event manager

- Change detection test: Page-Hinkley
- Instant cure

 $heta_t = heta_{t-1}$ $\eta_t = rac{1}{2}\eta_{t-1}$

Claim: both questions can be addressed using statistical tests

Page 54, Hinkley 70

(ロ) (部) (注) (注) (注) ()

Background

Catastrophic events in Deep Learning SGD and learning rate adaptation

S-ALERA

Speeding up: Agnostic learning rate adaptation Catastrophic event management

イロト イロト イヨト イヨト 二日

7/46

Experiments

Catastrophic events in Deep Learning

There exist steep cliff structures in the derivative landscape

Goodfellow et al. 16



Catastrophic episodes

Avoiding cliff neighbourhood:

- Regularization: L₁, L₂, Max-norm
- Gradient clipping

Pascanu et al., 13

Batch normalization

loffe et Szegedy, 15

Weight initialization

Glorot et Bengio, 10

8/46

(ロ) (同) (E) (E) (E)

Gradient descent

Given $\mathcal{D} = \{(x_i, y_i), i \in [[1, N]]\},$ Minimize loss function \mathcal{L} on parameter space $\Theta \subset \mathcal{R}^d$

Gradient descent: Iterate

Compute full gradient

$$g_t = \frac{1}{n} \sum_{1}^{N} \nabla_{\theta} \mathcal{L}(x_i; \theta_t)$$

Update

$$\theta_{t+1} = \theta_t - \eta_t g_t$$

(ロ) (部) (目) (日) (日) (の)

9/46

with θ_t solution at step t, η_t learning rate

Stochastic gradient descent

Stochastic gradient descent (SGD)

- ▶ Sample a mini-batch MB (*) of size n < N (**) from D
- Compute gradient

$$g_t = \frac{1}{n} \sum_{MB} \nabla_{\theta} \mathcal{L}(x; \theta_t)$$

Update

$$\theta_{t+1} = \theta_t - \eta g_t$$

(*) About selecting the minibatch a.k.a curriculum learning

Start with easy examples Bengio et al. 09
Self-Paced Learning Kumar et al. 10
Prefer hardest examples Loshchilov et al. 16
It depends Lopedriza et al. 16
(**) n << N ? Keskar et al. 17

Learning rate deterministic schedule (1/2)

$$g_t = \frac{1}{n} \sum_{MB} \nabla_{\theta} \mathcal{L}(x; \theta_t)$$

Learning rate decay

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

11 / 46

$$\eta_t = \frac{\eta_0}{1 + \nu t}$$
$$_{t+1} = \theta_t - \eta_t g_t$$

 θ

With η_0 : initial learning rate, ν hyper-parameter

Relaxing the gradient (2/2)

$$g_t = \frac{1}{n} \sum_i \nabla_{\theta} \mathcal{L}(x_i; \theta_t - \gamma v_{t-1})$$

Nesterov accelerated gradient

Nesterov 83; Sutskever et al., 13

(ロ) (部) (目) (日) (日) (の)

12 / 46

$$v_t = \gamma v_{t-1} + \eta g_t$$
$$\theta_{t+1} = \theta_t - v_t$$

With η learning rate, γ hyper-parameter [.8, .999]

Learning rate adaptive schedule 1/2

Adagrad

Duchi et al., 13

◆□ > ◆□ > ◆三 > ◆三 > ・ 三 ・ のへで

13 / 46

For each coordinate $j \in [[1, d]]$

$$g_{t,j} = \frac{1}{n} \sum_{i} \frac{\partial \mathcal{L}}{\partial \theta_j} (x_i; \theta_t)$$

$$G_{t,j} = \sum_{k=1}^{t} g_{k,j}^2$$
$$\eta_{t,j} = \frac{\eta_0}{\sqrt{G_{t,j} + \epsilon}}$$
$$\theta_{t+1,j} = \theta_{t,j} - \eta_{t,j} g_{t,j}$$

With η_0 initial learning rate, $\epsilon << 1$

Learning rate adaptive schedule 2/2

Adam

Kingma & Ba, 14

For each coordinate $j \in [[1, d]]$

$$g_{t,j} = \frac{1}{n} \sum_{i} \frac{\partial \mathcal{L}}{\partial \theta_j}(x_i; \theta_t)$$

$$\begin{split} m_{t,j} &= \beta_1 m_{t-1} + (1 - \beta_1) g_{t,j} \\ v_{t,j} &= \beta_2 v_{t-1} + (1 - \beta_2) g_{t,j}^2 \\ \hat{m}_t &= m_t / (1 - \beta_1^t) \quad \hat{v}_t = v_t / (1 - \beta_2^t) \\ \theta_{t+1,j} &= \theta_{t,j} - \frac{\eta_0}{\sqrt{\hat{v}_{t,j}} + \epsilon} \hat{m}_{t,j} \end{split}$$

With η_0 : initial learning rate, $\epsilon \ll 1$ and $\beta_1 \ll \beta_2$ (default $\beta_1 = .9, \beta_2 = .999$)

The delta-bar-delta rule

Asymmetrical schedule

For each coordinate $j \in [[1, d]]$

$$g_{t,j} = \frac{1}{n} \sum_{MB} \frac{\partial \mathcal{L}}{\partial \theta_j}(x; \theta_t)$$

$$\begin{split} \mathbf{v}_{t,j} &= \gamma \mathbf{v}_{t-1,j} + (1-\gamma) g_{t,j} \\ \eta_{t,j} &= \eta_{t-1,j} + \beta \text{ if } \operatorname{sgn}(g_{t,j}) = \operatorname{sgn}(\mathbf{v}_{t,j}) \\ &= \eta_{t-1,j}/\nu \quad \text{otherwise} \\ \theta_{t+1} &= \theta_t - \eta_t g_t \end{split}$$

with ν , γ , β hyper-parameters

Additive increase Multiplicative decrease Jacobs, 88

<ロ > < 部 > < 言 > < 言 > 差 の Q (~ 15/46

Background

Catastrophic events in Deep Learning SGD and learning rate adaptation

S-ALERA

Speeding up: Agnostic learning rate adaptation Catastrophic event management

イロン イロン イヨン イヨン 三日

16 / 46

Experiments

The (μ, λ) -Evolution Strategy

Minimize $f : \mathbb{R}^n \to \mathbb{R}$

Initialize distribution parameters m, σ, C , set population size $\lambda \in \mathbb{N}$

While not terminate

1. Sample distribution $\mathcal{N}(\boldsymbol{m}, \sigma^2 \mathbf{C}) \rightarrow \boldsymbol{x}_1, \dots, \boldsymbol{x}_{\lambda} \in \mathbb{R}^n$

 $\mathbf{x}_i \sim \mathbf{m} + \sigma \, \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \qquad ext{for } i = 1, \dots, \lambda$

2. Evaluate x_1, \ldots, x_{λ} on f

Compute $f(x_1), \ldots, f(x_{\lambda})$

3. Update parameters

 $\boldsymbol{m}, \sigma, \boldsymbol{\mathsf{C}} \leftarrow \boldsymbol{\mathsf{F}}(\boldsymbol{m}, \sigma, \boldsymbol{\mathsf{C}}, \boldsymbol{x}_1, \dots, \boldsymbol{x}_{\lambda}, f(\boldsymbol{x}_1), \dots, f(\boldsymbol{x}_{\lambda}))$

◆□ → < 部 → < 差 → < 差 → 差 < つ へ ペ 17/46

The (μ, λ) -Evolution Strategy (2)

Gaussian Mutations

- mean vector $\boldsymbol{m} \in \mathbb{R}^n$ is the current solution
- the so-called step-size $\sigma \in \mathbb{R}_+$ controls the step length
- ▶ the covariance matrix $C \in \mathbb{R}^{n \times n}$ determines the shape of the distribution ellipsoid

How to update m, σ , and C?

$$\blacktriangleright m = \frac{1}{\mu} \sum_{i=1}^{i=\mu} x_{i:\lambda}$$

Selection of best μ samples

18 / 46

• Adaptive σ and C

Need for adaptation



Cumulative Step-Size Adaptation (CSA)

Measure the length of the evolution path

the pathway of the mean vector \boldsymbol{m} in the generation sequence



Rationale

If successive moves are 'less correlated' than a random walk, step size should be decreased

hovering around the optimum

Cumulative Step-Size Adaptation (CSA)

Measure the length of the evolution path

the pathway of the mean vector \boldsymbol{m} in the generation sequence



Rationale

If successive moves are 'more correlated' than a random walk, step size should be increased

still far from the optimum

Step-size Adaptation at Work





initial distribution, $\mathbf{C} = \mathbf{I}$

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^{\mathrm{T}}$
- ▶ ruling principle: the adaptation increases the probability of successful steps, y_w , to appear again



 \mathbf{y}_{w} , movement of the population mean \mathbf{m} (disregarding σ)

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \boldsymbol{y}_w \boldsymbol{y}_w^{\mathrm{T}}$
- > ruling principle: the adaptation increases the probability of successful steps, y_w , to appear again



mixture of distribution **C** and step y_w , $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times y_w y_w^{\mathrm{T}}$

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \boldsymbol{y}_w \boldsymbol{y}_w^{\mathrm{T}}$
- > ruling principle: the adaptation increases the probability of successful steps, y_w , to appear again



new distribution (disregarding σ)

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^{\mathrm{T}}$
- ▶ ruling principle: the adaptation increases the probability of successful steps, y_w , to appear again



movement of the population mean m

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^{\mathrm{T}}$
- ▶ ruling principle: the adaptation increases the probability of successful steps, y_w , to appear again



mixture of distribution C and step \boldsymbol{y}_w , C \leftarrow 0.8 \times C + 0.2 \times $\boldsymbol{y}_w \boldsymbol{y}_w^{\mathrm{T}}$

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^{\mathrm{T}}$
- > ruling principle: the adaptation increases the probability of successful steps, y_w , to appear again



- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^{\mathrm{T}}$
- > ruling principle: the adaptation increases the probability of successful steps, y_w , to appear again

Rank- μ **Update**



Remark: the old (sample) distribution shape has a great influence on the new distribution \longrightarrow iterations needed

CMA-ES in one page

Input:
$$\boldsymbol{m} \in \mathbb{R}^n$$
, $\sigma \in \mathbb{R}_+$, λ
Initialize: $\mathbf{C} = \mathbf{I}$, and $\boldsymbol{p}_{\mathbf{C}} = \mathbf{0}$, $\boldsymbol{p}_{\sigma} = \mathbf{0}$,
Set: $c_{\mathbf{C}} \approx 4/n$, $c_{\sigma} \approx 4/n$, $c_1 \approx 2/n^2$, $c_{\mu} \approx \mu_w/n^2$, $c_1 + c_{\mu} \leq 1$, $d_{\sigma} \approx 1 + \sqrt{\frac{\mu_w}{n}}$, and
 $w_{i=1...\lambda}$ such that $\mu_w = \frac{1}{\sum_{i=1}^{\mu} w_i^2} \approx 0.3 \lambda$

While not terminate

$$\begin{split} \mathbf{x}_{i} &= \mathbf{m} + \sigma \, \mathbf{y}_{i}, \quad \mathbf{y}_{i} \sim \mathcal{N}_{i}(\mathbf{0}, \mathbf{C}), \quad \text{for } i = 1, \dots, \lambda \\ \mathbf{m} \leftarrow \sum_{i=1}^{\mu} \mathbf{w}_{i} \, \mathbf{x}_{i:\lambda} &= \mathbf{m} + \sigma \, \mathbf{y}_{w} \quad \text{where } \, \mathbf{y}_{w} = \sum_{i=1}^{\mu} \mathbf{w}_{i} \, \mathbf{y}_{i:\lambda} \\ \mathbf{p}_{c} \leftarrow (1 - c_{c}) \, \mathbf{p}_{c} + \mathbbm{1}_{\{\|\mathbf{p}_{\sigma}\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_{c})^{2}} \sqrt{\mu_{w}} \, \mathbf{y}_{w} \\ \mathbf{p}_{\sigma} \leftarrow (1 - c_{\sigma}) \, \mathbf{p}_{\sigma} + \sqrt{1 - (1 - c_{\sigma})^{2}} \sqrt{\mu_{w}} \, \mathbf{C}^{-\frac{1}{2}} \, \mathbf{y}_{w} \\ \mathbf{C} \leftarrow (1 - c_{1} - c_{\mu}) \, \mathbf{C} + c_{1} \, \mathbf{p}_{c} \, \mathbf{p}_{c}^{\mathrm{T}} + c_{\mu} \sum_{i=1}^{\mu} \mathbf{w}_{i} \, \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^{\mathrm{T}} \\ \mathbf{\sigma} \leftarrow \sigma \times \exp\left(\frac{c_{\sigma}}{d_{\sigma}} \left(\frac{\|\mathbf{p}_{\sigma}\|}{\mathsf{E}\|\mathcal{N}(\mathbf{0},\mathbf{I})\|} - 1\right)\right) \\ \end{split}$$

Not covered on this slide:

termination, restarts, active or mirrored sampling, outputs, and boundaries

Invariances: Guarantee for Generalization

Invariance properties of CMA-ES

 Invariance to order preserving transformations in function space

like all comparison-based algorithms

 Translation and rotation invariance to rigid transformations of the search space

CMA-ES is almost parameterless

- Tuning of a small set of functions
- Default values generalize to whole classes
- Exception: population size for multi-modal functions

but see IPOP-CMA-ES Auger & Hansen 05 and BIPOP-CMA-ES Hansen 09





Hansen & Ostermeier 01

BBOB – Black-Box Optimization Benchmarking

- ACM-GECCO workshops
- Set of 25 benchmark functions, dimensions 2 to 40
- With known difficulties (ill-conditioning, non-separability, ...)
- Noisy and non-noisy versions

Competitors include

- BFGS (Matlab version),
- Fletcher-Powell,
- DFO (Derivative-Free Optimization, Powell 04)
- Differential Evolution
- Particle Swarm Optimization
- and many more



 \sim yearly since 2009

COCO platform

sketching the solution trajectory in 2d...



Intuition

If successive gradient directions are more correlated than random directions, learning rate can be increased.

Inspiration CMA-ES

Hansen and Ostermeier, 01; Ollivier et al., 17

(ロ) (同) (E) (E) (E)

sketching the solution trajectory in 2d...



Intuition

If successive gradient directions are less correlated than random directions, learning rate can be decreased.

Inspiration CMA-ES

Hansen and Ostermeier, 01; Ollivier et al., 17

(ロ) (同) (E) (E) (E)

Given memory parameter α **Define**

Gradient cumulative path

$$egin{aligned} p_0 &= 0 \ p_t &= lpha rac{g_t}{\|g_t\|} + (1-lpha) p_{t-1} \end{aligned}$$

Random cumulative path

$$\begin{aligned} \mathbf{r}_0 &= \mathbf{0} \\ u_t &\sim \mathbb{U}_{\mathcal{S}_d} \\ \mathbf{r}_t &= \alpha u_t + (1 - \alpha) \mathbf{r}_{t-1} \end{aligned}$$

Expectation and std deviation of random cumulative path norm

$$\mathbb{E}(\|\mathbf{r}_t\|^2) = \frac{\alpha}{2-\alpha} [1-(1-\alpha)^{2t}] \stackrel{\text{def}}{=} \mu(t)$$

$$std(\|\mathbf{r}_t\|^2) = \frac{1}{d} \frac{2\alpha^2(1-\alpha)^2}{(2-\alpha)^2((1-\alpha)^2+1)} [1-(1-\alpha)^t] [1-(1-\alpha)^{t-1}] \stackrel{def}{=} \sigma(d,t)$$

Agnostic momentum adaptation

$$\eta_{t+1} = \eta_t \exp\left(C\frac{\|p_t\|^2 - \mu(t)}{\sigma(d, t)}\right)$$

29 / 46

with C hyper-parameter

Background

Catastrophic events in Deep Learning SGD and learning rate adaptation

S-ALERA

Speeding up: Agnostic learning rate adaptation Catastrophic event management

(日) (圖) (E) (E) (E)

30 / 46

Experiments

Change detection test



Example: given time series $\ell_{t<100} \sim \mathcal{N}(0,1), \ \ell_{t>100} \sim \mathcal{N}(1,1)$

Page-Hinkley change detection test

Page 54; Hinkley 70

- Parameter Δ , alarm threshold
- Given a time series (ℓ_1, \ldots, ℓ_t) , compute

$$\begin{array}{l} \hline \ell_t = 1/t \sum_i^t \ell_i \\ \hline L_t = \sum_i^t \ell_i - \overline{\ell_i} \\ \hline L_{min} = min_i L_i, \\ \end{array}$$

$$\begin{array}{l} L_{max} = max_i L \\ L_{max} = max_i L \end{array}$$

► Test triggered if $PH_t = \{L_t - L_{min} > \Delta\}$ increase detected $PH_t = \{L_{max} - L_t > \Delta\}$ decrease detected

Catastrophic event detection: Meeting the cliff



Page-Hinkley change detection test

Page 54; Hinkley 70

 $\bar{\ell}_t = 1/t \sum_{1}^t \ell_i$ $L_t = \sum_{1}^t \ell_i - \bar{\ell}_i$

$$L_{min} = min_i L_i$$

$$\blacktriangleright \mathsf{PH}_t = \{L_t - L_{min} > \Delta\}$$

Monitoring the learning curve

Signal: weighted average of mini-batch loss

$$\ell_t = \rho \mathcal{L}(\theta, \mathsf{MB}_t) + (1 - \rho)\ell_{t-1}$$

- Only monitor loss increase
- Alarm threshold

Catastrophic event curation

Usual catastrophic episodes

Goodfellow et al., 2016



Curation: If PH_t:

Backtrack

$$\theta_{t+1} = \theta_t$$

Halve learning rate

$$\eta_{t+1} = \frac{1}{2}\eta_t$$

< □ > < 部 > < 差 > < 差 > 差 の Q (~ 33/46

S-ALERA

	Algorithm 1: S-ALERA: Agnostic LEarning Rate Adaptation and	Page-Hinkley change detection		
	Input: Model with loss function \mathcal{L} Parameters : Memory rate α , factor C Initial learning rate η_0 , mini-batch ratio ρ Initialize : $t \leftarrow 0; p \leftarrow 0; \eta \leftarrow \eta_0; \operatorname{init}(\theta)$ $L, L_{min}, \ell, \overline{\ell}, \leftarrow 0; \Delta \leftarrow \mathcal{L}(\theta, \operatorname{first mini-batch})/10$	//algorithm parameters // run parameters // initialization // initialize Page-Hinkley		
1 2 3	while stopping criterion not met do $MB \leftarrow$ new mini-batch ; $t \leftarrow t + 1$ $\ell = \rho \mathcal{L}(\theta, MB) + (1 - \rho)\ell$	// perform forward pass		
4	$ar{\ell} \leftarrow (\ell + tar{\ell})/(t+1)$	// empirical mean of batch losses		
5	$L \leftarrow L + (\ell - \bar{\ell})$ //	cumulated deviations from mean		
6	$L_{min} \leftarrow min(L_{min}, L)$	// lower bound of deviations		
7	if $L - L_{min} > \Delta$ then			
8	$ heta \leftarrow heta^{(b)}; \eta \leftarrow \eta/2$ // P	age-Hinkley triggered: backtrack		
9	$L, L_{min}, \bar{\ell}, \ell \leftarrow 0; t \leftarrow 0$	// and re-initialize Page-Hinkley ²		
10	else	B		
11	$ heta^{(b)} \leftarrow heta$	// save for possible backtracks		
12	$g \leftarrow \nabla_{\theta} \mathcal{L}(\theta, \mathbf{MB})$ // com	pute gradient with backward pass		
13	$p \leftarrow \alpha g / \ g\ + (1 - \alpha)p$ // exponential moving	average of normalized gradients		
14	$\eta \leftarrow \eta \exp\left(C(\ p\ _2^2 - \mu)/\sigma(d)\right)$	// agnostic learning rate update		
15	$ heta \leftarrow heta - \eta g$	// standard parameter update		
16	end			
17	end			

Background

Catastrophic events in Deep Learning SGD and learning rate adaptation

S-ALERA

Speeding up: Agnostic learning rate adaptation Catastrophic event management

イロト イロト イヨト イヨト 二日

35 / 46

Experiments

Experimental setting 1/2

- Datasets
 - MNIST
 - CIFAR-10

Le Cun et al., 1998 Krizhevsky, 2009

- Four architectures, activation ReLU
 - M0: softmax classification (no hidden layers)
 - M2: 2 fully connected layers
 - M2b: M2 with batch normalization (BN) in each layer
 - M4: 2 convolutional layers followed by 2 fully connected layers (BN in each layer, max-pooling).
- Environment 46 GPUs; Torch

Collobert et al., 2011

Avg classification error (mean and std. on 5 runs)

Experimental setting 2/2

Baselines:

- ▶ Nesterov Accelerated Gradient (NAG) $\gamma \in \{.8, .9, .99, .999\}$
- ADAGRAD • ADAM $\beta_1 \in \{.8, .9, .99\}, \beta_2 \in \{.4, .9, .99\}$
 - $\beta_1 \in \{.8, .9, .99\}, \beta_2 \in \{.99, .999, .9999\}$

- New algorithms:
 - S-ALERA
 - ALERA: S-ALERA w/o PH
 - ► AL-ADAM: ADAMwith ALERA

 $\begin{array}{c} \alpha \in \{.001,.01,.1,.25\},\\ \mathcal{C} \in \{3.10^{-8},3.10^{-7},3.10^{-6},3.10^{-5}\}\\ \text{idem}\\ \text{idem, with } \beta_1 = 0.9, \,\beta_2 = 0.999 \end{array}$

Results - Learning performances

			NAG	Adagrad	Adam	AG-ADAM	AGMA	S-AGMA
MNIST	M0	5ep.	7.75 (.14)	7.73 (.06)	7.72 (.12)	7.31 (.09)	7.51 (.09)	7.51 (.09)
		20ep.	7.59 (.07)	7.51 (.09)	7.43 (.10)	7.29 (.09)	7.43 (.03)	7.44 (.04)
	M2	5ep.	1.95 (.17)	2.00 (.06)	2.07 (.11)	1.93 (.17)	1.86 (.11)	1.87 (.05)
		20ep.	1.58 (.08)	1.71 (.08)	1.56 (.06)	1.57 (.04)	1.55 (.10)	1.59 (.09)
	M2b	5ep.	1.82 (.13)	1.72 (.07)	1.81 (.07)	1.66 (.08)	1.59 (.08)	1.59 (.08)
		20ep.	1.47 (.10)	1.48 (.06)	1.57 (.94)	1.53 (.05)	1.43 (.04)	1.48 (.09)
	M4b	5ep.	.85 (.08)	1.02 (.09)	.89 (.31)	.91 (.06)	.82 (.30)	.82 (.30)
		20ep.	.72 (.09)	.82 (.08)	.80 (.08)	.79 (.05)	.63 (.05)	.64 (.07)
CIFAR	M0	5ep.	60.37 (.55)	60.49 (.71)	60.60 (.45)	59.62 (.27)	59.89 (.19)	59.69 (.33)
		20ep.	59.73 (.19)	59.76 (.36)	59.81 (.24)	59.34 (.24)	59.31 (.11)	59.31 (.25)
	M2	5ep.	45.82 (.93)	44.81 (.62)	45.68 (.39)	44.91 (.42)	44.69 (17.58)	44.42 (.24)
		20ep.	45.08 (.32)	43.59 (.51)	44.43 (.50)	43.25 (.40)	43.19 (.21)	42.72 (.41)
	M2b	5ep.	45.01 (.84)	44.18 (.62)	44.30 (.96)	43.33 (.33)	43.08 (.17)	43.08 (.17)
		20ep.	42.50 (.48)	43.79 (.25)	43.60 (.64)	42.72 (.33)	42.12 (.15)	42.50 (.29)
	M4b	5ep.	27.74 (.48)	34.93 (.96)	28.50 (.68)	25.60 (.29)	28.61 (.50)	28.61 (.50)
		20ep.	27.45 (.39)	29.15 (.67)	27.84 (.59)	25.30 (.18)	26.35 (.64)	25.93 (.64)

Table 1: Best performances across all hyper-parameter settings at 5 and 20 epochs (average and standard deviation of test error on 5 runs). Bold: best result/line, red: less than 1 std.dev. away.

- ► AL-ADAM: outperforms ADAM
- ▶ S-ALERA: same as ALERA, less sensitive (less catastrophes)

Robust S-ALERA setting: $\alpha = .01$, $C = 3.10^{-6}$

ALERA vs S-ALERA on MNIST, M2, same seed



Comments First catastrophic episode around epoch 8: S-ALERA halves the learning rate for all three layers.

Further catastrophic events met are more rare and less severe for $\operatorname{S-ALERA}$ than $\operatorname{ALERA}.$

Analysis of the dividing factor

Considering the 1D parabola problem

$$F(heta) = rac{1}{2} a heta^2$$

Notations

- Optimal learning rate: $\eta^* = \frac{1}{a}$
- For $\eta > \eta_- = 2\eta^*$, loss increases

Phase 1: decreasing learning rate

1. Assume $\eta > \eta_-$

2. Thus divide η by ζ until falling in $\left[\frac{\eta_{-}}{\zeta}, \eta_{-}\right[$ Number of divisions $S = \log(\frac{\eta_{-}}{n})/\log(\zeta)$.

Phase 2: increasing learning rate

- Simplification:
 - If $\eta > \eta^*$, $\eta \to \eta (1 \varepsilon)$
 - Else $\eta \to \eta \left(1 + \varepsilon\right)$

Let ζ the dividing factor loss increase

Analysis of the dividing factor, 2

Case 1, $\zeta < 2$

• At the end of Phase 1, $\eta > \eta^*$; hence, multiplied by $(1 - \varepsilon)$ until reaching η^* .

Case 2, $\zeta > 2$

- At the end of Phase 1, η/η_{-} in $[\frac{1}{\zeta}, 1[$
- Case 2.1: η/η_- in $[\frac{1}{\zeta}, \frac{1}{2}]$
- Case 2.2: η/η₋ in [¹/₂, 1[

Case 2.1: Expectation of time to reach η^*

- Probability $\frac{1}{2}\frac{\zeta-2}{\zeta-1}$
- Expectation of number of multiplications by $(1 + \varepsilon)$

$$U = \frac{1}{\varepsilon} \int_{\frac{1}{\zeta}}^{\frac{1}{2}} -\log(2u) du = \frac{1}{\varepsilon} \left[-t \log(t) + t - t \log(2) \right]_{\frac{1}{\zeta}}^{\frac{1}{2}}$$
$$= \frac{1}{\varepsilon} \left(\frac{1}{2} - \frac{\log(\zeta)}{\zeta} - \frac{1 - \log(2)}{\zeta} \right)$$

Analysis of the dividing factor, 3

Case 2.2: Expectation of time to reach η^*

• Probability $\frac{1}{2} \frac{\zeta}{\zeta - 1}$

• Expectation of number of multiplications by $(1 - \varepsilon)$

$$T = \frac{1}{\varepsilon} \int_{\frac{1}{2}}^{1} \log(2u) du = \frac{1}{\varepsilon} [t \log(t) - t + t \log(2)]_{\frac{1}{2}}^{1}$$
$$= \frac{1}{\varepsilon} (\log(2) - \frac{1}{2})$$

Overall expected cost c_C : catastrophic regime (forward pass) c_S : standard regime (forward + backward pass)

$$c_{C}S + c_{S}\left(\frac{1}{2}\frac{\zeta-2}{\zeta-1}U + \frac{1}{2}\frac{\zeta}{\zeta-1}T\right)$$

Optimal ζ : minimizing

$$J(\zeta) = \frac{C}{\log(\zeta)} + \frac{1}{2}\frac{\zeta-2}{\zeta-1}\left(\frac{1}{2} - \frac{\log(\zeta)}{\zeta} - \frac{1 - \log(2)}{\zeta}\right) + \frac{1}{2}\frac{\zeta}{\zeta-1}\left(\log(2) - \frac{1}{2}\right)$$

with

$$C = \frac{c_C}{c_S} \varepsilon \log\left(\frac{\eta}{\eta_-}\right)$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Optimal dividing factor

Assuming C small ($c_C < c_S$, ε small)



・ロン ・回 と ・ ヨ と ・ ヨ と

3

43 / 46

Cost $J(\zeta)$ for different values of C: Optimal ζ in [3,5]

Conclusions

- A principled way to address catastrophic episodes
 - Detection: Page-Hinkley change detection test
 - Curation: backtrack and halve learning rates
- Enables a more aggressive learning rate adaptation scheme
 - Comparing the cumulative gradient path to that of random unit vectors

Remarks

 \blacktriangleright PH does reduce the catastrophic event rate (test error > 80%) by \approx 40%

イロン 不良 とくほど 不良 とうほう

44 / 46

Paper: https://arxiv.org/abs/1709.01427 Code: https://github.com/lalil⊘u/salera/

Perspectives

- More datasets
- Extension to Recurrent Neural Networks
- Adapt PH threshold Δ during learning
- Consider other learning signals

Acknowledgements

- Steve Altschuler and Lani Wu
- Yann Ollivier
- Sigurd Angenent