

UNIVERSITÄT LEIPZIG

Anti-Correlated Noise in Epoch-Based Stochastic Gradient Descent and its Implications

Prague, 10.04.2025 **Marcel Kühn** and Bernd Rosenow – Institute for Theoretical Physics





UNIVERSITAT LEIPZIG Institute for Theoretical Physics







UNIVERSITÄT LEIPZIG

• learning from examples

 $\{x_n, y_n \mid n=1, ..., N\}$... dataset x ... example y ... label

learning from examples •

> $\{x_n, y_n \mid n=1, ..., N\}$... dataset $x \ldots$ example $y \ldots$ label

- FIG 4. Has the UK been hit by a hurricane? 0: The Great Storm of 1987 was a violent extratropical P: **BoolQ Dataset** cyclone which caused casualties in England, France and the Channel Islands ... Yes. [An example event is given.] A:
- Does France have a Prime Minister and a President? **O**:
- ... The extent to which those decisions lie with the P: Prime Minister or President depends upon ...
- Yes. [Both are mentioned, so it can be inferred both A: exist.]
- Have the San Jose Sharks won a Stanley Cup? 0:
- ... The Sharks have advanced to the Stanley Cup fi-P: nals once, losing to the Pittsburgh Penguins in 2016
- A: No. [They were in the finals once, and lost.]

[Clark et al. (2019)]

learning from examples •

$$\{x_n, y_n \mid n=1, ..., N\}$$
 ... dataset x ... example

 $y \ldots$ label

FIG 4. Has the UK been hit by a hurricane? 0: The Great Storm of 1987 was a violent extratropical P: **BoolQ Dataset** cyclone which caused casualties in England, France and the Channel Islands ... Yes. [An example event is given.] A: **O**: Does France have a Prime Minister and a President? ... The extent to which those decisions lie with the P: Prime Minister or President depends upon ... Yes. [Both are mentioned, so it can be inferred both A: exist.] Have the San Jose Sharks won a Stanley Cup? 0: ... The Sharks have advanced to the Stanley Cup fi-P: nals once, losing to the Pittsburgh Penguins in 2016 A: No. [They were in the finals once, and lost.]

[Clark et al. (2019)]



[Krizhevsky and Hinton (2009)]

FIG 5. CIFAR10 Dataset

learning from examples •

$$\{x_n, y_n \mid n=1, ..., N\}$$
 ... dataset x ... example y ... label

FIG 4. Has the UK been hit by a hurricane? 0: The Great Storm of 1987 was a violent extratropical p. **BoolQ Dataset** cyclone which caused casualties in England, France Yes. [An example event is given.] A: 0: ... The extent to which those decisions lie with the **P**: Prime Minister or President depends upon ... A: Yes. [Both are mentioned, so it can be inferred both Have the San Jose Sharks won a Stanley Cup? 0: ... The Sharks have advanced to the Stanley Cup fi-P: nals once, losing to the Pittsburgh Penguins in 2016 A: No. [They were in the finals once, and lost.]



[Krizhevsky and Hinton (2009)]

FIG 5. CIFAR10 Dataset

Neural networks

 $f(oldsymbol{ heta},x)$... network $oldsymbol{ heta} \in \mathbb{R}^d$... network weights



UNIVERSITÄT LEIPZIG

Institute for Theoretical Physics

Neural networks

 $f(oldsymbol{ heta},x)$... network $oldsymbol{ heta} \in \mathbb{R}^d$... network weights





Institute for Theoretical Physics

$$\boldsymbol{\theta} \in \mathbb{R}^d \dots$$
 network weights
 $l_n(\boldsymbol{\theta}) \coloneqq l(f(\boldsymbol{\theta}, x_n), y_n) \dots$ example loss
 $L(\boldsymbol{\theta}) \coloneqq \frac{1}{N} \sum_{n=1}^N l_n(\boldsymbol{\theta}) \dots$ total loss

Quantify mismatch with loss •

$$\boldsymbol{\theta} \in \mathbb{R}^d \dots$$
 network weights
 $l_n(\boldsymbol{\theta}) \coloneqq l(f(\boldsymbol{\theta}, x_n), y_n) \dots$ example loss
 $L(\boldsymbol{\theta}) \coloneqq \frac{1}{N} \sum_{n=1}^N l_n(\boldsymbol{\theta}) \dots$ total loss

Quantify mismatch with loss •



FIG 8. Loss Landscape

$$\boldsymbol{\theta} \in \mathbb{R}^d \dots$$
 network weights
 $l_n(\boldsymbol{\theta}) \coloneqq l(f(\boldsymbol{\theta}, x_n), y_n) \dots$ example loss
 $L(\boldsymbol{\theta}) \coloneqq \frac{1}{N} \sum_{n=1}^N l_n(\boldsymbol{\theta}) \dots$ total loss

- Quantify mismatch with loss •
- Aim: minimize loss •



FIG 8. Loss Landscape

$$\boldsymbol{\theta} \in \mathbb{R}^d \dots$$
 network weights
 $l_n(\boldsymbol{\theta}) \coloneqq l(f(\boldsymbol{\theta}, x_n), y_n) \dots$ example loss
 $L(\boldsymbol{\theta}) \coloneqq \frac{1}{N} \sum_{n=1}^N l_n(\boldsymbol{\theta}) \dots$ total loss

- Quantify mismatch with loss
- Aim: minimize loss
- How? Gradient Descent!





Gradient Descent

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \boldsymbol{\nabla} L(\boldsymbol{\theta}_k)$$



Gradient Descent

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \boldsymbol{\nabla} L(\boldsymbol{\theta}_k)$$

Stochastic Gradient Descent (SGD) 1

$$\mathbf{g}_k(\boldsymbol{\theta}) = \frac{1}{S} \sum_{n \in \mathcal{B}_k} \boldsymbol{\nabla} l_n(\boldsymbol{\theta})$$

$$\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} - \eta \mathbf{g}_k(\boldsymbol{\theta}_{k-1})$$

$$k \ldots$$
 update step index

- η ... learning rate
- $S \ldots$ batch size



$$\mathcal{B}_k = \{n_1, ..., n_S\} \dots S \ll N$$
 randomly selected examples $(n_j \in \{1, ..., N\})$

UNIVERSITÄT Institute for Theoretical Physics

Gradient Descent

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \boldsymbol{\nabla} L(\boldsymbol{\theta}_k)$$

Stochastic Gradient Descent with Momentum (SGD) 1

$$\mathbf{g}_k(\boldsymbol{\theta}) = \frac{1}{S} \sum_{n \in \mathcal{B}_k} \boldsymbol{\nabla} l_n(\boldsymbol{\theta})$$
 (1a)

$$\mathbf{v}_k = -\eta \mathbf{g}_k(\boldsymbol{\theta}_{k-1}) + \beta \mathbf{v}_{k-1}$$
 (1b)

$$\boldsymbol{ heta}_k = \boldsymbol{ heta}_{k-1} + \mathbf{v}_k$$
 (1c)

- $k \ldots$ update step index
- η ... learning rate
- S ... batch size
- β ... momentum parameter

$$\mathcal{B}_k = \{n_1, ..., n_S\} \ldots S \ll N$$
 randomly selected examples

 $(n_i \in \{1, \dots, N\})$



How to choose examples?

- How to choose examples?
 - With or Without replacement (in epochs)



FIG 11. Epoch-based example selection

- How to choose examples?
 - > With or Without replacement (in epochs)



FIG 11. Epoch-based example selection

- How to choose examples?
 - With or Without replacement (in epochs)
- Reason for anti-correlations



FIG 11. Epoch-based example selection

• gradient noise:

 $\mathbf{g}_k(oldsymbol{ heta}) = rac{1}{S}\sum_{n\in\mathcal{B}_k}oldsymbol{
abla} l_n(oldsymbol{ heta})$ (1a)

$$egin{aligned} \delta \mathbf{g}_k \coloneqq \mathbf{g}_k(oldsymbol{ heta}) - oldsymbol{
abla} L(oldsymbol{ heta}) \ \mathbf{C} &= \mathop{\mathrm{cov}}_k(\delta \mathbf{g}_k, \delta \mathbf{g}_k) \end{aligned}$$

- gradient noise: • $\delta \mathbf{g}_k \coloneqq \mathbf{g}_k(\boldsymbol{\theta}) - \boldsymbol{\nabla} L(\boldsymbol{\theta})$ $\mathbf{g}_k(oldsymbol{ heta}) = rac{1}{S}\sum_{n\in\mathcal{B}_k}oldsymbol{
 abla} l_n(oldsymbol{ heta})$ (1a) $\mathbf{C} = \mathop{\mathrm{cov}}_k(\delta \mathbf{g}_k, \delta \mathbf{g}_k)$
- noise of one epoch has zero mean \rightarrow anti-correlations: ٠

$$\begin{split} &\frac{1}{M}\sum_{k\,\in\,\mathrm{epoch}\,e}\mathbf{g}_k(\boldsymbol{\theta}) = \boldsymbol{\nabla}L(\boldsymbol{\theta}) \quad M\coloneqq \frac{N}{S}\dots \text{batches per epoch} \\ &\Rightarrow \sum \quad \boldsymbol{\delta}\mathbf{g}_k = 0 \end{split}$$

 $k \in \operatorname{epoch} e$

- gradient noise: $\mathbf{g}_k(\theta) = \frac{1}{S} \sum_{n \in \mathcal{B}_k} \nabla l_n(\theta) \quad \text{(1a)} \quad \mathbf{C} = \underset{k}{\operatorname{cov}} (\mathbf{\delta} \mathbf{g}_k, \mathbf{\delta} \mathbf{g}_k)$
- noise of one epoch has zero mean \rightarrow anti-correlations:

$$\frac{1}{M} \sum_{k \in \text{epoch } e} \mathbf{g}_k(\boldsymbol{\theta}) = \boldsymbol{\nabla} L(\boldsymbol{\theta}) \quad M \coloneqq \frac{N}{S} \dots \text{ batches per epoch}$$
$$\Rightarrow \sum_{k \in \text{epoch } e} \boldsymbol{\delta} \mathbf{g}_k = 0 \qquad \Rightarrow -\sum_{j \in e/\{k\}} \boldsymbol{\delta} \mathbf{g}_j = \boldsymbol{\delta} \mathbf{g}_k$$

- gradient noise: •
 - $\mathbf{g}_k(oldsymbol{ heta}) = rac{1}{S}\sum_{n\in\mathcal{B}_k}oldsymbol{
 abla} l_n(oldsymbol{ heta})$ (1a)

$$egin{aligned} & oldsymbol{\delta} \mathbf{g}_k \coloneqq \mathbf{g}_k(oldsymbol{ heta}) - oldsymbol{
abla} L(oldsymbol{ heta}) \ & \mathbf{C} = \mathop{\mathrm{cov}}_k(oldsymbol{\delta} \mathbf{g}_k, oldsymbol{\delta} \mathbf{g}_k) \end{aligned}$$

noise of one epoch has zero mean \rightarrow anti-correlations: .

$$\frac{1}{M} \sum_{k \in \text{epoch } e} \mathbf{g}_k(\boldsymbol{\theta}) = \boldsymbol{\nabla} L(\boldsymbol{\theta}) \quad M \coloneqq \frac{N}{S} \dots \text{ batches per epoch}$$
$$\Rightarrow \sum_{k \in \text{epoch } e} \boldsymbol{\delta} \mathbf{g}_k = 0 \qquad \Rightarrow -\sum_{j \in e/\{k\}} \boldsymbol{\delta} \mathbf{g}_j = \boldsymbol{\delta} \mathbf{g}_k$$

average over chance for two batches from same epoch .

$$\operatorname{cov}_{k}(\boldsymbol{\delta}\mathbf{g}_{k}, \boldsymbol{\delta}\mathbf{g}_{k+h}) = \mathbf{C} \cdot \left(\delta_{h,0} - \mathbf{1}_{\{1,\dots,M\}}(|h|) \frac{M - |h|}{M(M - 1)}\right)$$
(2)

UNIVERSITÄT Institute for Theoretical Physics

LEIPZIG



FIG 12. Noise autocorrelation: Examples selected without replacement

- gradient noise: •
- $\mathbf{g}_k(\boldsymbol{\theta}) = rac{1}{S} \sum_{n \in \mathcal{B}_k} \boldsymbol{\nabla} l_n(\boldsymbol{\theta})$ (1a)

$$egin{aligned} & oldsymbol{\delta} \mathbf{g}_k \coloneqq \mathbf{g}_k(oldsymbol{ heta}) - oldsymbol{
abla} L(oldsymbol{ heta}) \ & \mathbf{C} = \mathop{\mathrm{cov}}_k(oldsymbol{\delta} \mathbf{g}_k, oldsymbol{\delta} \mathbf{g}_k) \end{aligned}$$

noise of one epoch has zero mean \rightarrow anti-correlations: .

$$\frac{1}{M} \sum_{k \in \text{epoch } e} \mathbf{g}_k(\boldsymbol{\theta}) = \boldsymbol{\nabla} L(\boldsymbol{\theta}) \quad M \coloneqq \frac{N}{S} \dots \text{ batches per epoch}$$
$$\Rightarrow \sum_{k \in \text{epoch } e} \boldsymbol{\delta} \mathbf{g}_k = 0 \qquad \Rightarrow -\sum_{j \in e/\{k\}} \boldsymbol{\delta} \mathbf{g}_j = \boldsymbol{\delta} \mathbf{g}_k$$

average over chance for two batches from same epoch

$$\operatorname{cov}_{k}(\boldsymbol{\delta}\mathbf{g}_{k}, \boldsymbol{\delta}\mathbf{g}_{k+h}) = \mathbf{C} \cdot \left(\delta_{h,0} - \mathbf{1}_{\{1,\dots,M\}}(|h|) \frac{M - |h|}{M(M - 1)}\right)$$
(2)





FIG 13. Noise autocorrelation: Examples selected with replacement

UNIVERSITÄT

- calculations done for static weights
 - o numerics also show anti-correlations at start of training



Questions about Anti-correlations?

Questions about Anti-correlations?

Implications for Weight Fluctuations

late phase of training: •

$$L(\boldsymbol{\theta}) \approx \frac{1}{2} \boldsymbol{\theta}^{\top} \mathbf{H} \boldsymbol{\theta}$$
 (3)

H ... Hessian matrix



late phase of training: .

$$L(\boldsymbol{\theta}) \approx \frac{1}{2} \boldsymbol{\theta}^{\top} \mathbf{H} \boldsymbol{\theta}$$
 (3)

H ... Hessian matrix

- analyze weight covariance •
 - resulting from stochasticity of SGD
 - measure for parameter exploration
 - insight about escape from minima

$$\boldsymbol{\Sigma} = \operatorname{cov}_k(\boldsymbol{ heta}_k, \boldsymbol{ heta}_k) \ \dots$$
 weight fluctuations (4)



FIG 16. Weight trajectory in toy loss

default case in physics: thermal noise

$$\mathbf{g}_k(\boldsymbol{\theta}) = rac{1}{S} \sum_{n \in \mathcal{B}_k} \boldsymbol{\nabla} l_n(\boldsymbol{\theta})$$
 (1a)

$$\mathbf{v}_{k} = -\eta \mathbf{g}_{k}(\boldsymbol{\theta}_{k-1}) + \beta \mathbf{v}_{k-1}$$
(1b)
$$\boldsymbol{\theta}_{k} = \boldsymbol{\theta}_{k-1} + \mathbf{v}_{k}$$
(1c)

 $\Sigma \propto \mathbf{H}^{-1} \dots$ (see Einstein Relation) for $\mathbf{C} \coloneqq \operatorname{cov}_k(\mathbf{g}_k, \mathbf{g}_k) \propto \mathbb{1} \dots$ isotropic noise

default case in physics: thermal noise •

$$\mathbf{g}_k(\boldsymbol{\theta}) = rac{1}{S} \sum_{n \in \mathcal{B}_k} \boldsymbol{\nabla} l_n(\boldsymbol{\theta})$$
 (1a)

$$\mathbf{v}_{k} = -\eta \mathbf{g}_{k}(\boldsymbol{\theta}_{k-1}) + \beta \mathbf{v}_{k-1}$$
(1b)
$$\boldsymbol{\theta}_{k} = \boldsymbol{\theta}_{k-1} + \mathbf{v}_{k}$$
(1c)

 $\mathbf{\Sigma} \propto \mathbf{H}^{-1}$... (see Einstein Relation) for $\mathbf{C} \coloneqq \operatorname{cov}_k(\mathbf{g}_k,\mathbf{g}_k) \propto \mathbb{1}$... isotropic noise

in neural networks: non-isotropic noise \rightarrow previously in literature: • (see Jastrzębski et al. (2018), Liu et al. (2021))

 $\mathbf{C} \propto \mathbf{H} \dots$ Hessian noise approximation

 $\Rightarrow \Sigma \propto 1$... isotropic weight fluctuations

default case in physics: thermal noise •

$$\mathbf{g}_k(\boldsymbol{\theta}) = rac{1}{S} \sum_{n \in \mathcal{B}_k} \boldsymbol{\nabla} l_n(\boldsymbol{\theta})$$
 (1a)

$$\mathbf{v}_k = -\eta \mathbf{g}_k(\boldsymbol{\theta}_{k-1}) + \beta \mathbf{v}_{k-1}$$
 (1b)
 $\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} + \mathbf{v}_k$ (1c)

 $\mathbf{\Sigma} \propto \mathbf{H}^{-1}$... (see Einstein Relation) for $\mathbf{C} \coloneqq \operatorname{cov}_k(\mathbf{g}_k,\mathbf{g}_k) \propto \mathbb{1} \dots$ isotropic noise

in neural networks: non-isotropic noise \rightarrow previously in literature: • (see Jastrzębski et al. (2018), Liu et al. (2021))

 $\mathbf{C} \propto \mathbf{H} \dots$ Hessian noise approximation

 $\Rightarrow \Sigma \propto 1$... isotropic weight fluctuations

recent empirical results contradict, more likely suggest • (see Feng and Tu (2022))

 $\Sigma \propto H \ldots$ anisotropic weight fluctuations

default case in physics: thermal noise

$$\mathbf{g}_k(\boldsymbol{\theta}) = \frac{1}{S} \sum_{n \in \mathcal{B}_k} \nabla l_n(\boldsymbol{\theta})$$
 (1a)

$$\mathbf{v}_k = -\eta \mathbf{g}_k(\boldsymbol{\theta}_{k-1}) + \beta \mathbf{v}_{k-1}$$
 (1b)
 $\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} + \mathbf{v}_k$ (1c)

 $\Sigma \propto \mathbf{H}^{-1} \dots$ (see Einstein Relation) for $\mathbf{C} \coloneqq \operatorname{cov}_k(\mathbf{g}_k, \mathbf{g}_k) \propto \mathbb{1} \dots$ isotropic noise

• in neural networks: non-isotropic noise → previously in literature: (see Jastrzębski et al. (2018), Liu et al. (2021))

 $\mathbf{C} \propto \mathbf{H}$... Hessian noise approximation

 $\Rightarrow \mathbf{\Sigma} \propto \mathbb{1} \ \dots$ isotropic weight fluctuations

• recent empirical results contradict, more likely suggest (see Feng and Tu (2022))

 $\mathbf{\Sigma} \propto \mathbf{H} \dots$ anisotropic weight fluctuations

solution: noise anti-correlations

weights covariance $\mathbf{\Sigma}\coloneqq\mathrm{cov}(\boldsymbol{ heta}_k, \boldsymbol{ heta}_k)$, similarly velocity covariance $\mathbf{\Sigma}_{\mathbf{v}}\coloneqq\mathrm{cov}(\mathbf{v}_k, \mathbf{v}_k)$ • $(\mathbf{v}_k = \boldsymbol{\theta}_k - \boldsymbol{\theta}_{k-1})$
- weights covariance $\Sigma \coloneqq \operatorname{cov}(\theta_k, \theta_k)$, similarly velocity covariance $\Sigma_{\mathbf{v}} \coloneqq \operatorname{cov}(\mathbf{v}_k, \mathbf{v}_k)$ $(\mathbf{v}_k = \theta_k - \theta_{k-1})$
- Assumption 1: Quadratic Approximation $L(\theta) = \frac{1}{2} \theta^{\top} \mathbf{H} \theta$ with Hessian \mathbf{H}

- weights covariance $\Sigma \coloneqq \operatorname{cov}(\boldsymbol{\theta}_k, \boldsymbol{\theta}_k)$, similarly velocity covariance $\Sigma_{\mathbf{v}} \coloneqq \operatorname{cov}(\mathbf{v}_k, \mathbf{v}_k)$ • $(\mathbf{v}_k = \boldsymbol{\theta}_k - \boldsymbol{\theta}_{k-1})$
- Assumption 1: Quadratic Approximation • $L(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^{\top} \mathbf{H} \boldsymbol{\theta}$ with Hessian \mathbf{H}
- Assumption 2: Anti-correlated Noise • Previously calculated SGD noise autocorrelation holds even for non static weights.

- weights covariance $\Sigma \coloneqq \operatorname{cov}(\theta_k, \theta_k)$, similarly velocity covariance $\Sigma_{\mathbf{v}} \coloneqq \operatorname{cov}(\mathbf{v}_k, \mathbf{v}_k)$ $(\mathbf{v}_k = \theta_k - \theta_{k-1})$
- Assumption 1: Quadratic Approximation $L(\theta) = \frac{1}{2} \theta^{\top} \mathbf{H} \theta$ with Hessian \mathbf{H}
- Assumption 2: Anti-correlated Noise
 Previously calculated SGD noise autocorrelation holds even for non static weights.
- Assumption 3: Hessian and noise covariance commute Covariance of noise commutes with Hessian, [C, H] = 0

- weights covariance $\Sigma \coloneqq \operatorname{cov}(\theta_k, \theta_k)$, similarly velocity covariance $\Sigma_{\mathbf{v}} \coloneqq \operatorname{cov}(\mathbf{v}_k, \mathbf{v}_k)$ $(\mathbf{v}_k = \theta_k - \theta_{k-1})$
- Assumption 1: Quadratic Approximation $L(\theta) = \frac{1}{2} \theta^{\top} \mathbf{H} \theta$ with Hessian \mathbf{H}
- Assumption 2: Anti-correlated Noise Previously calculated SGD noise autocorrelation holds even for non static weights.
- Assumption 3: Hessian and noise covariance commute Covariance of noise commutes with Hessian, [C, H] = 0
- with given assumptions: matrices commute, denote eigenvalues for common eigenvectors with i = 1, ..., d as: λ_i for $\mathbf{H}, \sigma_{\delta g, i}^2$ for $\mathbf{C}, \sigma_{\theta, i}^2$ for $\Sigma, \sigma_{v, i}^2$ for $\Sigma_{\mathbf{v}}$

Weight Fluctuations at a Minimum - Derivation

SGD with momentum in one dimension (with assumptions): •

$$g_k = \frac{\partial}{\partial \theta} \left[\frac{1}{2} \theta_{k-1} \lambda \theta_{k-1} \right] + \delta g_k , \quad v_k = -\eta g_k + \beta v_{k-1} , \quad \theta_k = \theta_{k-1} + v_k$$
(5)

$$\theta_k = (1 - \eta \lambda)\theta_{k-1} + \beta(\theta_{k-1} - \theta_{k-2}) - \eta \delta g_k$$
(6)

Weight Fluctuations at a Minimum - Derivation

• SGD with momentum in one dimension (with assumptions):

$$g_k = \frac{\partial}{\partial \theta} \left[\frac{1}{2} \theta_{k-1} \lambda \theta_{k-1} \right] + \delta g_k , \quad v_k = -\eta g_k + \beta v_{k-1} , \quad \theta_k = \theta_{k-1} + v_k$$
(5)

$$\theta_k = (1 - \eta \lambda)\theta_{k-1} + \beta(\theta_{k-1} - \theta_{k-2}) - \eta \delta g_k$$
(6)

• in matrix form:

$$\mathbf{x}_{k} = \mathbf{D}\mathbf{x}_{k-1} - \eta \delta g_{k} \mathbf{e}_{1}, \quad \mathbf{x}_{k} \coloneqq \begin{pmatrix} \theta_{k} \\ \theta_{k-1} \end{pmatrix}, \quad \mathbf{e}_{1} \coloneqq \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$
(7)
$$\mathbf{D} \coloneqq \begin{pmatrix} 1 + \beta - \eta \lambda & -\beta \\ 1 & 0 \end{pmatrix}$$
(8)

Weight Fluctuations at a Minimum - Derivation

SGD with momentum in one dimension (with assumptions): •

$$g_k = \frac{\partial}{\partial \theta} \left[\frac{1}{2} \theta_{k-1} \lambda \theta_{k-1} \right] + \delta g_k , \quad v_k = -\eta g_k + \beta v_{k-1} , \quad \theta_k = \theta_{k-1} + v_k$$
(5)

$$\theta_k = (1 - \eta\lambda)\theta_{k-1} + \beta(\theta_{k-1} - \theta_{k-2}) - \eta\delta g_k$$
(6)

in matrix form: •

LEIPZIG

$$\mathbf{x}_{k} = \mathbf{D}\mathbf{x}_{k-1} - \eta \delta g_{k} \mathbf{e}_{1}, \quad \mathbf{x}_{k} \coloneqq \begin{pmatrix} \theta_{k} \\ \theta_{k-1} \end{pmatrix}, \quad \mathbf{e}_{1} \coloneqq \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$
(7)
$$\mathbf{D} \coloneqq \begin{pmatrix} 1 + \beta - \eta \lambda & -\beta \\ 1 & 0 \end{pmatrix}$$
(8)

respective covariance matrix: •

$$\left\langle \mathbf{x}_{k} \mathbf{x}_{k}^{\top} \right\rangle = \begin{pmatrix} \sigma_{\theta}^{2} & \langle \theta_{k} \theta_{k-1} \rangle \\ \langle \theta_{k} \theta_{k-1} \rangle & \sigma_{\theta}^{2} \end{pmatrix}$$
 (9)

Weight Fluctuations at a Minimum - Derivation $\mathbf{x}_{k} = \mathbf{D}\mathbf{x}_{k-1} - \eta \delta g_{k} \mathbf{e}_{1}, \quad \mathbf{x}_{k} \coloneqq \begin{pmatrix} \theta_{k} \\ \theta_{k-1} \end{pmatrix}$ (7)

$$\left\langle \mathbf{x}_{k} \mathbf{x}_{k}^{\top} \right\rangle = \mathbf{D} \left\langle \mathbf{x}_{k-1} \mathbf{x}_{k-1}^{\top} \right\rangle \mathbf{D}^{\top} + \eta^{2} \left\langle \delta g_{k} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} - \eta \left(\mathbf{D} \left\langle \mathbf{x}_{k-1} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} + \left(\mathbf{D} \left\langle \mathbf{x}_{k-1} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} \right)^{\top} \right)$$
(10)

Weight Fluctuations at a Minimum - Derivation $\mathbf{x}_{k} = \mathbf{D}\mathbf{x}_{k-1} - \eta \delta g_{k} \mathbf{e}_{1}, \quad \mathbf{x}_{k} \coloneqq \begin{pmatrix} \theta_{k} \\ \theta_{k-1} \end{pmatrix}$ (7)

$$\left\langle \mathbf{x}_{k} \mathbf{x}_{k}^{\top} \right\rangle = \mathbf{D} \left\langle \mathbf{x}_{k-1} \mathbf{x}_{k-1}^{\top} \right\rangle \mathbf{D}^{\top} + \eta^{2} \left\langle \delta g_{k} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} - \eta \left(\mathbf{D} \left\langle \mathbf{x}_{k-1} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} + \left(\mathbf{D} \left\langle \mathbf{x}_{k-1} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} \right)^{\top} \right)$$
(10)

• notice $v_k = \theta_k - \theta_{k-1} \rightarrow \sigma_v^2 = 2\sigma_\theta^2 - 2\langle \theta_k \theta_{k-1} \rangle$, and apply on \mathbf{e}_1 :

$$\begin{pmatrix} \sigma_{\theta}^{2} \\ \sigma_{v}^{2} \end{pmatrix} = \mathbf{F} \left[\eta^{2} \sigma_{\delta g}^{2} \mathbf{e}_{1} \mathbf{e}_{1}^{\top} - \eta \left(\mathbf{D} \left\langle \mathbf{x}_{k-1} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} + \left(\mathbf{D} \left\langle \mathbf{x}_{k-1} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} \right)^{\top} \right) \right] \mathbf{e}_{1}$$
(11)

$$\mathbf{F} = \frac{1}{(1-\beta)\left(2(1+\beta)-\eta\lambda)\right)} \begin{pmatrix} \frac{1+\beta}{\eta\lambda} & \frac{2\beta(\eta\lambda-1-\beta)}{\eta\lambda}\\ 2 & 2(\eta\lambda-2) \end{pmatrix}$$
(12)

UNIVERSITAT LEIPZIG Institute for Theoretical Physics

mkuehn@itp.uni-leipzig.de

Weight Fluctuations at a Minimum - Derivation $\mathbf{x}_{k} = \mathbf{D}\mathbf{x}_{k-1} - \eta \delta g_{k} \mathbf{e}_{1}, \quad \mathbf{x}_{k} \coloneqq \begin{pmatrix} \theta_{k} \\ \theta_{k-1} \end{pmatrix} \quad (7)$

$$\left\langle \mathbf{x}_{k} \mathbf{x}_{k}^{\top} \right\rangle = \mathbf{D} \left\langle \mathbf{x}_{k-1} \mathbf{x}_{k-1}^{\top} \right\rangle \mathbf{D}^{\top} + \eta^{2} \left\langle \delta g_{k} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} - \eta \left(\mathbf{D} \left\langle \mathbf{x}_{k-1} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} + \left(\mathbf{D} \left\langle \mathbf{x}_{k-1} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} \right)^{\top} \right)$$
(10)

notice $v_k = \theta_k - \theta_{k-1} \rightarrow \sigma_v^2 = 2\sigma_\theta^2 - 2\langle \theta_k \theta_{k-1} \rangle$, and apply on \mathbf{e}_1 : •

$$\begin{pmatrix} \sigma_{\theta}^{2} \\ \sigma_{v}^{2} \end{pmatrix} = \mathbf{F} \left[\eta^{2} \sigma_{\delta g}^{2} \mathbf{e}_{1} \mathbf{e}_{1}^{\top} - \eta \left(\mathbf{D} \left\langle \mathbf{x}_{k-1} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} + \left(\mathbf{D} \left\langle \mathbf{x}_{k-1} \delta g_{k} \right\rangle \mathbf{e}_{1}^{\top} \right)^{\top} \right) \right] \mathbf{e}_{1}$$
(11)

$$\mathbf{F} = \frac{1}{(1-\beta)\left(2(1+\beta)-\eta\lambda)\right)} \begin{pmatrix} \frac{1+\beta}{\eta\lambda} & \frac{2\beta(\eta\lambda-1-\beta)}{\eta\lambda}\\ 2 & 2(\eta\lambda-2) \end{pmatrix}$$
(12)

iterate SGD equation, use anti-correlation assumption: •

$$\langle \mathbf{x}_{k-1} \delta g_k \rangle = \eta \sigma_{\delta g}^2 \frac{\mathbf{D}^M + (\mathbf{1} - \mathbf{D})M - \mathbf{1}}{(\mathbf{1} - \mathbf{D})^2 M (M - 1)} \mathbf{e}_1$$
(13)

Weight Fluctuations at a Minimum - Solution

- eigenvalues with i = 1, ..., d as: λ_i for **H**, $\sigma_{\delta g, i}^2$ for **C**, $\sigma_{\theta, i}^2$ for Σ , $\sigma_{v, i}^2$ for $\Sigma_{\mathbf{v}}$
- solution: $\left(\mathbf{e}_1 = (1 \ 0)^{\top}\right)$

$$\begin{pmatrix} \sigma_{\theta,i}^2 \\ \sigma_{v,i}^2 \end{pmatrix} = \eta^2 \sigma_{\delta g,i}^2 \mathbf{F}_{\mathbf{i}} \left[\mathbf{e}_1 - \left(\mathbf{E}_{\mathbf{i}} + \mathbf{E}_{\mathbf{i}}^\top \right) \mathbf{e}_1 \right]$$
(14)

$$\mathbf{F_i} = \frac{1}{(1-\beta)\left(2(1+\beta)-\eta\lambda_i\right)} \begin{pmatrix} \frac{1+\beta}{\eta\lambda_i} & \frac{2\beta(\eta\lambda_i-1-\beta)}{\eta\lambda_i}\\ 2 & 2(\eta\lambda_i-2) \end{pmatrix}$$
(15)

$$\mathbf{E}_{\mathbf{i}} \coloneqq \mathbf{D}_{\mathbf{i}} \frac{\mathbf{D}_{\mathbf{i}}^{M} + (\mathbf{1} - \mathbf{D}_{\mathbf{i}})M - \mathbf{1}}{(\mathbf{1} - \mathbf{D}_{\mathbf{i}})^{2}M(M - 1)} \mathbf{e}_{1}\mathbf{e}_{1}^{\top}$$
(16)



FIG 17. Empirical weight variance and velocity variance against Hessian eigenvalue

 $\mathbf{D}_{\mathbf{i}} \coloneqq \begin{pmatrix} 1+\beta-\eta\lambda_i & -\beta\\ 1 & 0 \end{pmatrix}$

1

UNIVERSITÄT

LEIPZIG

(17)

more comprehensive understanding with velocity-• correlation-time τ_i $\begin{pmatrix} \theta_{k,i} \coloneqq \boldsymbol{\theta}_k \cdot \mathbf{p}_i \\ v_{k,i} \coloneqq \mathbf{v}_k \cdot \mathbf{p}_i \end{pmatrix}$ for common eigenvector $\mathbf{p}_i \end{pmatrix}$

$$\tau_i \coloneqq \frac{2\sigma_{\theta,i}^2}{\sigma_{v,i}^2} = \frac{\sum_{h=1}^{\infty} h \cdot \operatorname{cov}(v_{k,i}, v_{k+h,i})}{\sum_{h=1}^{\infty} \operatorname{cov}(v_{k,i}, v_{k+h,i})}$$
(18)

• more comprehensive understanding with velocitycorrelation-time τ_i $\begin{pmatrix} \theta_{k,i} := \boldsymbol{\theta}_k \cdot \mathbf{p_i} \\ v_{k,i} := \mathbf{v}_k \cdot \mathbf{p_i} \end{pmatrix}$ for common eigenvector $\mathbf{p_i}$

$$\tau_i \coloneqq \frac{2\sigma_{\theta,i}^2}{\sigma_{v,i}^2} = \frac{\sum_{h=1}^{\infty} h \cdot \operatorname{cov}(v_{k,i}, v_{k+h,i})}{\sum_{h=1}^{\infty} \operatorname{cov}(v_{k,i}, v_{k+h,i})}$$
(18)

• two distinct regimes of Hessian eigenvalues (EVs), separated by $\lambda_{cross} \coloneqq \frac{3(1-\beta)S}{\eta N}$

 $\lambda_i > \lambda_{
m cross}$: can neglect anti-correlations

$$\to \tau_i \approx \frac{1+\beta}{\eta \lambda_i}$$

more comprehensive understanding with velocity- $\begin{pmatrix} \theta_{k,i} := \boldsymbol{\theta}_k \cdot \mathbf{p}_i \\ v_{k,i} := \mathbf{v}_k \cdot \mathbf{p}_i \\ \end{cases} \text{ for common eigenvector } \mathbf{p}_i \end{pmatrix}$ correlation-time au_i

$$\tau_i \coloneqq \frac{2\sigma_{\theta,i}^2}{\sigma_{v,i}^2} = \frac{\sum_{h=1}^{\infty} h \cdot \operatorname{cov}(v_{k,i}, v_{k+h,i})}{\sum_{h=1}^{\infty} \operatorname{cov}(v_{k,i}, v_{k+h,i})}$$
(18)

two distinct regimes of Hessian eigenvalues (EVs), • separated by $\lambda_{cross} \coloneqq \frac{3(1-\beta)S}{nN}$

 $\lambda_i > \lambda_{\rm cross}$: can neglect anti-correlations

$$\rightarrow \tau_i \approx \frac{1+\beta}{\eta \lambda_i}$$

 $\lambda_i < \lambda_{\rm cross}$: anti-correlations important

$$\rightarrow \tau_i \approx \frac{M}{3} \frac{1+\beta}{1-\beta} \eqqcolon \tau_{\text{SGD}}$$



more comprehensive understanding with velocity- $\begin{pmatrix} \theta_{k,i} := \boldsymbol{\theta}_k \cdot \mathbf{p_i} \\ v_{k,i} := \mathbf{v}_k \cdot \mathbf{p_i} \end{pmatrix} \text{ for common eigenvector } \mathbf{p_i} \end{pmatrix}$ correlation-time τ_i

$$\tau_i \coloneqq \frac{2\sigma_{\theta,i}^2}{\sigma_{v,i}^2} = \frac{\sum_{h=1}^{\infty} h \cdot \operatorname{cov}(v_{k,i}, v_{k+h,i})}{\sum_{h=1}^{\infty} \operatorname{cov}(v_{k,i}, v_{k+h,i})}$$
(18)

two distinct regimes of Hessian eigenvalues (EVs), separated by $\lambda_{cross} \coloneqq \frac{3(1-\beta)S}{nN}$

 $\lambda_i > \lambda_{\rm cross}$: can neglect anti-correlations

$$\rightarrow \tau_i \approx \frac{1+\beta}{\eta \lambda_i}$$

 $\lambda_i < \lambda_{\rm cross}$: anti-correlations important

$$\rightarrow \tau_i \approx \frac{M}{3} \frac{1+\beta}{1-\beta} \eqqcolon \tau_{\text{SGD}}$$



- two distinct regimes separated by $\lambda_{cross} \coloneqq \frac{3(1-\beta)S}{\eta N}$ •
- **Large Hessian EVs (** $\lambda_i > \lambda_{cross}$ **):** can neglect anti-correlations $\rightarrow \tau_i \approx \frac{1+\beta}{n\lambda_i}$ and •

$$\sigma_{\theta,i}^2 \approx \frac{\eta^2 \sigma_{\delta g,i}^2}{2(1-\beta)(1+\beta)} \cdot \frac{1+\beta}{\eta \lambda_i} , \ \sigma_{v,i}^2 \approx \frac{\eta^2 \sigma_{\delta g,i}^2}{(1-\beta)(1+\beta)}$$
(19)

- two distinct regimes separated by $\lambda_{cross} \coloneqq \frac{3(1-\beta)S}{\eta N}$
- Large Hessian EVs ($\lambda_i > \lambda_{cross}$): can neglect anti-correlations $\rightarrow \tau_i \approx \frac{1+\beta}{\eta \lambda_i}$ and

$$\sigma_{\theta,i}^2 \approx \frac{\eta^2 \sigma_{\delta g,i}^2}{2(1-\beta)(1+\beta)} \cdot \frac{1+\beta}{\eta \lambda_i} , \ \sigma_{v,i}^2 \approx \frac{\eta^2 \sigma_{\delta g,i}^2}{(1-\beta)(1+\beta)}$$
(19)

• Small Hessian EVs ($\lambda_i < \lambda_{cross}$): anti-correlations important, $\tau_i \approx \frac{M}{3} \frac{1+\beta}{1-\beta} \Rightarrow \tau_{SGD}$ and

$$\sigma_{\theta,i}^2 \approx \frac{\eta^2 \sigma_{\delta g,i}^2}{2(1-\beta)(1+\beta)} \cdot \frac{M}{3} \frac{1+\beta}{1-\beta} , \ \sigma_{v,i}^2 \approx \frac{\eta^2 \sigma_{\delta g,i}^2}{(1-\beta)(1+\beta)}$$
(20)

UNIVERSITAT LEIPZIG Institute for Theoretical Physics

approximate relation: •

$$\sigma_{v,i}^2 \propto \sigma_{\delta g,i}^2 \quad \text{or} \quad \boldsymbol{\Sigma}_{\mathbf{v}} \propto \mathbf{C} \ (\propto \mathbf{H})$$

$$\sigma_{\theta,i}^2 = \frac{1}{2} \tau_i \sigma_{v,i}^2 \quad \text{or} \quad \boldsymbol{\Sigma} = \frac{1}{2} \text{diag}(\tau_1, \dots, \tau_d) \boldsymbol{\Sigma}_{\mathbf{v}}$$

(in corresponding eigenbasis)

approximate relation: •

$$\sigma_{v,i}^2 \propto \sigma_{\delta g,i}^2 \quad \text{or} \quad \boldsymbol{\Sigma}_{\mathbf{v}} \propto \mathbf{C} \; (\propto \mathbf{H})$$

$$\sigma_{\theta,i}^2 = \frac{1}{2} \tau_i \sigma_{v,i}^2 \quad \text{or} \quad \boldsymbol{\Sigma} = \frac{1}{2} \text{diag}(\tau_1, \ldots, \tau_d) \boldsymbol{\Sigma}_{\mathbf{v}}$$

(in corresponding eigenbasis)

Hessian eigenvector subspace with $\lambda_i > \lambda_{cross}$ • $(\tau_i \propto \lambda_i^{-1})$:

$$\Sigma^{\scriptscriptstyle{(>)}} \propto 1$$

Hessian eigenvector subspace with $\lambda_i < \lambda_{cross}$ • $(\tau_i = \tau_{\text{SGD}})$:

$\Sigma^{\scriptscriptstyle{(<)}} \propto { m H}$



FIG 17. Empirical weight variance and velocity variance against Hessian eigenvalue

approximate relation:

$$\sigma_{v,i}^2 \propto \sigma_{\delta g,i}^2 \quad \text{or} \quad \boldsymbol{\Sigma}_{\mathbf{v}} \propto \mathbf{C} \; (\propto \mathbf{H})$$

$$\sigma_{\theta,i}^2 = \frac{1}{2} \tau_i \sigma_{v,i}^2 \quad \text{or} \quad \boldsymbol{\Sigma} = \frac{1}{2} \text{diag}(\tau_1, \ldots, \tau_d) \boldsymbol{\Sigma}_{\mathbf{v}}$$

(in corresponding eigenbasis)

Hessian eigenvector subspace with $\lambda_i > \lambda_{cross}$ • $(\tau_i \propto \lambda_i^{-1})$:

$\Sigma^{\scriptscriptstyle{(>)}} \propto 1$

Hessian eigenvector subspace with $\lambda_i < \lambda_{cross}$ • $(\tau_i = \tau_{\text{SGD}})$:

$\Sigma^{\scriptscriptstyle{(<)}} \propto {f H}$



FIG 17. Empirical weight variance and velocity variance against Hessian eigenvalue

- previous observation by Feng and Tu (2022) closer to $\mathbf{\Sigma} \propto \mathbf{H}^2$
- principal component analysis of Σ leads to finite size effects $ightarrow \mathbf{H}$ eigenbasis analysis superior

Varying Hyperparameters



FIG 20. Varying Batch Size: Empirical weight variance, velocity variance and correlation time against Hessian eigenvalue

Varying Hyperparameters



variance and correlation time against Hessian eigenvalue

1.0

0.5

0.0

0.01

0.02

Initial Learning Rate

Across

1.0

0.5

50

75

Batch Size

0.04

0.03

0.75

0.50

Momentum

× 1.0

0.5

0.00

0.25

125

×

100

Different Network





FIG 22. Improved Architecture - ResNet instead of LeNet: Empirical weight variance and velocity variance against Hessian eigenvalue FIG 23. Improved Architecture - ResNet instead of LeNet: Correlation time against Hessian eigenvalue

Questions about Weight Fluctuations?

Questions about Weight Fluctuations?

Implications for Training

Additional Error due to Fluctuations

$$\begin{split} L(\boldsymbol{\theta}) &= \frac{1}{2} \boldsymbol{\theta}^{\mathrm{T}} \mathbf{H} \boldsymbol{\theta} \\ \left\langle L(\boldsymbol{\theta}_{k}) \right\rangle_{k} &= \left\langle \sum_{i=1}^{d} \frac{1}{2} \theta_{i,k} \lambda_{i} \theta_{i,k} \right\rangle_{k} \qquad \qquad \sigma_{\theta,i}^{2} \begin{cases} = \sigma_{\theta,\max}^{2} & \text{if } \lambda_{i} > \lambda_{\mathrm{cross}} \\ \propto \lambda_{i} \text{, and } < \sigma_{\theta,\max}^{2} & \text{if } \lambda_{i} < \lambda_{\mathrm{cross}} \end{cases} \\ &= \sum_{i=1}^{d} \frac{1}{2} \lambda_{i} \sigma_{\theta,i}^{2} \end{split}$$

Additional Error due to Fluctuations

$$\begin{split} L(\boldsymbol{\theta}) &= \frac{1}{2} \boldsymbol{\theta}^{\mathrm{T}} \mathbf{H} \boldsymbol{\theta} \\ \left\langle L(\boldsymbol{\theta}_{k}) \right\rangle_{k} &= \left\langle \sum_{i=1}^{d} \frac{1}{2} \theta_{i,k} \lambda_{i} \theta_{i,k} \right\rangle_{k} \qquad \qquad \sigma_{\theta,i}^{2} \begin{cases} = \sigma_{\theta,\max}^{2} & \text{if } \lambda_{i} > \lambda_{\mathrm{cross}} \\ \propto \lambda_{i} \,, \text{ and } < \sigma_{\theta,\max}^{2} & \text{if } \lambda_{i} < \lambda_{\mathrm{cross}} \end{cases} \\ &= \sum_{i=1}^{d} \frac{1}{2} \lambda_{i} \sigma_{\theta,i}^{2} \end{split}$$

Reduction of weight fluctuations leads to 62% reduction of loss fluctuations •



FIG 24. Experimental weight variance and velocity variance against Hessian eigenvalue

 test accuracy of LeNet is 0.7% ± 0.2% higher for drawing examples in SGD without replacement (64.5%) than for SGD with replacement (63.8%)

- test accuracy of LeNet is 0.7% ± 0.2% higher for drawing examples in SGD without re-• placement (64.5%) than for SGD with replacement (63.8%)
- in agreement with Orvieto et al. (2022), where anti-correlated perturbed gradient descent • (Anti-PGD) was found beneficial for test accuracy and led to flatter minima

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \boldsymbol{\nabla} L(\boldsymbol{\theta}_k) + (\boldsymbol{\xi}_{k+1} - \boldsymbol{\xi}_k) \dots \text{Anti-PGD}$$

 $\boldsymbol{\xi}_{k,i} \sim \mathcal{N}(0, \sigma^2)$

- test accuracy of LeNet is 0.7% ± 0.2% higher for drawing examples in SGD without replacement (64.5%) than for SGD with replacement (63.8%)
- in agreement with Orvieto et al. (2022), where anti-correlated perturbed gradient descent (Anti-PGD) was found beneficial for test accuracy and led to flatter minima

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \boldsymbol{\nabla} L(\boldsymbol{\theta}_k) + (\boldsymbol{\xi}_{k+1} - \boldsymbol{\xi}_k) \dots \text{Anti-PGD}$$

 $\boldsymbol{\xi}_{k,i} \sim \mathcal{N}(0, \sigma^2)$

$$\mathbb{E}[\mathbf{z}_{k+1}|\mathbf{z}_k] \approx \mathbf{z}_k - \eta \nabla \tilde{L}(\mathbf{z}_k) \qquad (\mathbf{z}_k = \boldsymbol{\theta}_k - \boldsymbol{\xi}_k)$$
$$\tilde{L}(\mathbf{z}) \coloneqq L(\mathbf{z}) + \frac{\sigma^2}{2} \operatorname{Tr}(\mathbf{H}(\mathbf{z}))$$

- test accuracy of LeNet is 0.7% ± 0.2% higher for drawing examples in SGD without re-• placement (64.5%) than for SGD with replacement (63.8%)
- in agreement with Orvieto et al. (2022), where anti-correlated perturbed gradient descent • (Anti-PGD) was found beneficial for test accuracy and led to flatter minima

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \boldsymbol{\nabla} L(\boldsymbol{\theta}_k) + (\boldsymbol{\xi}_{k+1} - \boldsymbol{\xi}_k) \dots \text{Anti-PGD}$$

 $\xi_{k,i} \sim \mathcal{N}(0, \sigma^2)$

$$\mathbb{E}[\mathbf{z}_{k+1}|\mathbf{z}_k] \approx \mathbf{z}_k - \eta \nabla \hat{L}(\mathbf{z}_k) \qquad (\mathbf{z}_k = \boldsymbol{\theta}_k - \tilde{L}(\mathbf{z}_k)) = L(\mathbf{z}) + \frac{\sigma^2}{2} \operatorname{Tr}(\mathbf{H}(\mathbf{z}))$$



FIG 25. Gradient Descent with anti-correlated noise in a widening vallev

LEIPZIG

 $(\boldsymbol{\xi}_k)$



Conclusion

• SGD without replacement induces anti-correlated noise





Conclusion

- SGD without replacement induces anti-correlated noise
- Results in lower-than-expected weight variance in Hessian eigendirections with small EVs



Conclusion

- SGD without replacement induces anti-correlated noise
- Results in lower-than-expected weight variance in Hessian eigendirections with small EVs
- beneficial because gradients in flat directions then dominate fluctuations and lead network towards even flatter minima with improved generalization performance



Conclusion

- SGD without replacement induces anti-correlated noise
- Results in lower-than-expected weight variance in Hessian eigendirections with small EVs
- beneficial because gradients in flat directions then dominate fluctuations and lead network towards even flatter minima with improved generalization performance

Thank You!

Marcel Kühn and Bernd Rosenow

Institute for Theoretical Physics

mkuehn@itp.uni-leipzig.de

References

Berner et al. (2022)

Berner J, Grohs P, Kutyniok G and Petersen P 2022 The Modern Mathematics of Deep Learning. In: Grohs P, Kutyniok G, editors. Mathematical Aspects of Deep Learning. Cambridge: Cambridge University Press pp 1–111

Orvieto et al. (2022)

Orvieto A, Kersting H, Proske F, Bach F and Lucchi A 2022 Anticorrelated Noise Injection for Improved Generalization *Proceedings of the 39th International Conference on Machine Learning* vol 162 (PMLR) pp 17094–17116

Clark et al. (2019)

Clark C, Lee K, Chang M, Kwiatkowski T, Collins M, Toutanova K 2019 BoolQ: Exploring the Surprising Difficulty of Natural Yes/No Questions *Proceedings of the 2019* Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies vol 1 pp 2924--2936

Krizhevsky and Hinton (2009)

Krizhevsky A 2009 Learning multiple layers of features from tiny images Tech. rep. University of Toronto

Feng and Tu (2021)

Feng Y and Tu Y 2021 Proceedings of the National Academy of Sciences 118

He et al. (2016)

He K, Zhang X, Ren S and Sun J 2016 Deep Residual Learning for Image Recognition IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp 770–778

Jastrzębski et al. (2018)

Jastrzębski S, Kenton Z, Arpit D, Ballas N, Fischer A, Bengio Y and Storkey A 2018 Three Factors Influencing Minima in SGD (Arxiv Preprint 1711.04623)

Liu et al. (2021)

Liu K, Ziyin L and Ueda M 2021 Noise and Fluctuation of Finite Learning Rate Stochastic Gradient Descent Proceedings of the 38th International Conference on Machine Learning vol 139 (PMLR) pp 7045–7056

UNIVERSITÄT LEIPZIG
Training evolution



FIG 26. Evolution of loss (left) and accuracy (right) during training of Lenet