

# Optimization Algorithms in the Large

Exact Dynamics, Average-case Analysis, and Stepsize  
Criticality

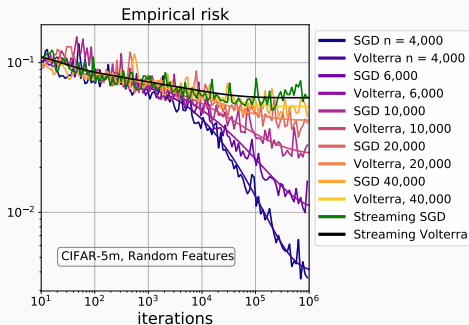
---

Courtney Paquette

Joint work: Fabian Pedregosa, Elliot Paquette, Bart van Merriënboer,  
Kiwon Lee, Jeffrey Pennington, Ben Adlam, and Andrew Cheng

Erice, May 2022

# Theory meets practice: CIFAR-5m



Using a random features model to predict CIFAR-5m (Nakkiran et al., '21) car/plane, the **Volterra equation** (using the Hessian spectra an input) gives **good predictions** for behavior of SGD.

# Typical Machine Learning Problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

High dimensional  $\Leftrightarrow$  large number of **features ( $d$ )** and **samples ( $n$ )**

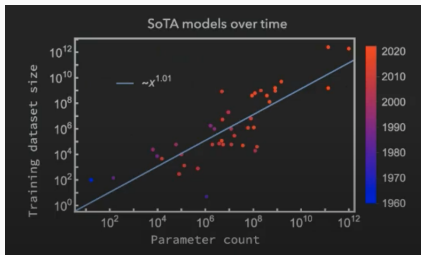
- ✓ State-of-the-art models have **millions/billions parameters**
  - Meena: 2.6 billion, Turing NLG: 17 billion, GPT-3: 175 billion

# Typical Machine Learning Problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

High dimensional  $\Leftrightarrow$  large number of **features** ( $d$ ) and **samples** ( $n$ )

- ✓ State-of-the-art models have **millions/billions parameters**
  - Meena: 2.6 billion, Turing NLG: 17 billion, GPT-3: 175 billion
- ✓ Ratio of features ( $d$ ) to samples ( $n$ ) are proportional



What's different about **high-dimensions**?

Input which generates worst complexity can be far from typical  
*"more room = more possibilities"*

What's different about **high-dimensions**?

Input which generates worst complexity can be far from typical  
*"more room = more possibilities"*

**How do we capture high-dimensional structure?**

**Probability distribution on the inputs**

**Remark:** results will hold for deterministic designs

Statistical learning (Mei & Montanari '19, Adlam & Pennington '21, Louart & Liao & Couillet '18)

Numerical Methods (Trogdon & Deift '19, Chandrasekher '21)

# Stochastic algorithms and random least squares

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 + \frac{\delta}{2} \|\mathbf{x}\|^2 = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) \stackrel{\text{def}}{=} \sum_{i=1}^n \underbrace{\frac{1}{2} (\mathbf{a}_i \mathbf{x} - b_i)^2 + \frac{\delta}{2n} \|\mathbf{x}\|^2}_{f_i(\mathbf{x})} \right\},$$

with *random*  $\mathbf{A} \in \mathbb{R}^{n \times d}$ ,  $\mathbf{b} \in \mathbb{R}^n$  *random* vector

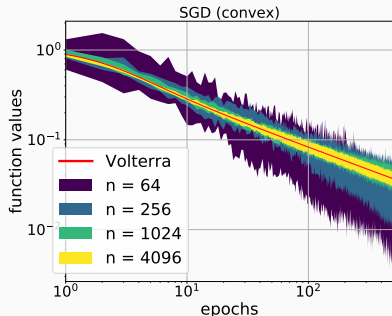
**Stochastic Gradient Descent (SGD)**     $\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma(t) \nabla f_{i_k}(\mathbf{x}_k)$

# Stochastic algorithms and random least squares

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \frac{\delta}{2} \|\mathbf{x}\|^2 = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) \stackrel{\text{def}}{=} \sum_{i=1}^n \underbrace{\frac{1}{2} (\mathbf{a}_i \mathbf{x} - b_i)^2 + \frac{\delta}{2n} \|\mathbf{x}\|^2}_{f_i(\mathbf{x})} \right\},$$

with *random*  $\mathbf{A} \in \mathbb{R}^{n \times d}$ ,  $\mathbf{b} \in \mathbb{R}^n$  *random* vector

**Stochastic Gradient Descent (SGD)**  $\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma(t) \nabla f_{i_k}(\mathbf{x}_k)$



For large models, as  $\frac{d}{n} \rightarrow r$ ,

- $f(\mathbf{x}_k) \xrightarrow{\text{Pr}}$  (smooth function)
- Analyze this smooth function
- Determined by the spectrum of the Hessian



$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 \quad \text{Hessian } \mathbf{H} = \mathbf{A}^T \mathbf{A}$$

## Distributions on data matrices $A$ and Random matrices

Many classes of random matrices share characteristics which depend on the low moments of their entries, called **Universality**

### Example

- **De-localization of eigenvectors of  $H$** : eigenvectors are not aligned with the unit vectors  
e.g., if  $A_{i,j} \sim N(0,1)$ , then eigenvectors of  $H \sim \text{Unif}(\mathbb{S}^{d-1})$

## Detour into random matrix theory...

$$\text{Hessian of least squares: } \mathbf{H} = \mathbf{A}^T \mathbf{A}$$

### Assumptions on data matrix (Bai & Silverstein '10, Benigni & Peche '19)

1. model size ( $d$ ) and # of samples ( $n$ ) polynomially related

$$d^\alpha \leq n \leq d^{1/\alpha} \quad \text{for some } \alpha \in (0, 1)$$

2. Mild assumptions on eigenvalues  $\lambda_{\max}$  and  $\lambda_{\min}$  of  $\mathbf{H}$ ,  $\|\mathbf{H}\|_2 < C$

# Detour into random matrix theory...

$$\text{Hessian of least squares: } \mathbf{H} = \mathbf{A}^T \mathbf{A}$$

## Assumptions on data matrix (Bai & Silverstein '10, Benigni & Peche '19)

1. model size ( $d$ ) and # of samples ( $n$ ) polynomially related

$$d^\alpha \leq n \leq d^{1/\alpha} \quad \text{for some } \alpha \in (0, 1)$$

2. Mild assumptions on eigenvalues  $\lambda_{\max}$  and  $\lambda_{\min}$  of  $\mathbf{H}$ ,  $\|\mathbf{H}\|_2 < C$

3. De-localization of eigenvectors of  $\mathbf{H}$ :

$\Omega$  contour enclosing eigenvalues of  $\mathbf{H}$ ,  $\theta < 1/2$

(i)  $\max_{z \in \Omega} \max_{1 \leq i \neq j \leq n} |\mathbf{e}_i^T R(z; \mathbf{A}\mathbf{A}^T) \mathbf{e}_j| \leq n^{\theta-1/2}$ .

(ii)  $\max_{z \in \Omega} \max_{1 \leq i \leq n} |\mathbf{e}_i^T R(z; \mathbf{A}\mathbf{A}^T) \mathbf{e}_i - \frac{1}{n} \text{tr} R(z; \mathbf{A}\mathbf{A}^T)| \leq n^{\theta-1/2}$ .

(iii) **Examples**

- **Isotropic features.** Entries of  $\mathbf{A} \sim N(0, 1)$
- **Sample covariance matrices.** independent samples w/ covariance between features
- **Random features.**  $\mathbf{A} = \sigma(\mathbf{X}\mathbf{W})$  where  $\sigma$  is an activation function

# Our framework

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \stackrel{\text{def}}{=} \underbrace{\frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2}_{\mathcal{L}(\mathbf{x})} + \frac{\delta}{2} \|\mathbf{x}\|^2, \quad \mathcal{L}(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2$$

## Model

- (random) **initialization**  $\mathbf{x}_0 \in \mathbb{R}^d$  is  $\|\mathbf{x}_0\|^2 \leq R$  (across dimensions)
- $\mathbf{b}$  is **target** vector,  $\|\mathbf{b}\|^2 \leq \tilde{R}$  (across dimensions)
- $\mathbf{b}$  and  $\mathbf{x}_0$  are "independent" of the eigenvectors of  $\mathbf{A}$

# Our framework

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \stackrel{\text{def}}{=} \underbrace{\frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2}_{\mathcal{L}(\mathbf{x})} + \frac{\delta}{2} \|\mathbf{x}\|^2, \quad \mathcal{L}(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2$$

## Model

- (random) **initialization**  $\mathbf{x}_0 \in \mathbb{R}^d$  is  $\|\mathbf{x}_0\|^2 \leq R$  (across dimensions)
- $\mathbf{b}$  is **target** vector,  $\|\mathbf{b}\|^2 \leq \tilde{R}$  (across dimensions)
- $\mathbf{b}$  and  $\mathbf{x}_0$  are "independent" of the eigenvectors of  $\mathbf{A}$

## Algorithmic considerations

- Small batch: batch size  $\rightarrow 0$  as  $n \rightarrow \infty$
- Multi-pass SGD
- Step size: bounded, continuous function s.t.  $\gamma(t) \rightarrow \tilde{\gamma}$   
(if  $\tilde{\gamma} = 0$ , then  $\int_0^\infty \gamma(s) ds = \infty$ , Robbins-Monro conditions)

# Exact Dynamics Idea: Diffusion Approximation

**Homogenized SGD** (C.P.-E. Paquette, NeurIPS '21 & Mori '21)

$$d\mathbf{X}_t = -\gamma(t)\nabla f(\mathbf{X}_t) dt + \sqrt{\frac{2}{n}\mathcal{L}(\mathbf{X}_t)\nabla^2\mathcal{L}(\mathbf{X}_t)} d\mathbf{B}_t$$

$\mathbf{X}_0 = \mathbf{x}_0$  and  $(\mathbf{B}_t : t \geq 0)$  is a  $d$ -dimen. standard Brownian motion

- New diffusion process (c.f. Li et al., Mandt et al.):
  - Dimension  $n \rightarrow \infty$  instead of stepsize  $\gamma \rightarrow 0$  to create it
  - Explicitly solvable and don't need to send stepsize to 0

# Exact Dynamics Idea: Diffusion Approximation

## Homogenized SGD (C.P.-E. Paquette, NeurIPS '21 & Mori '21)

$$d\mathbf{X}_t = -\gamma(t)\nabla f(\mathbf{X}_t) dt + \sqrt{\frac{2}{n}\mathcal{L}(\mathbf{X}_t)\nabla^2\mathcal{L}(\mathbf{X}_t)} d\mathbf{B}_t$$

$\mathbf{X}_0 = \mathbf{x}_0$  and  $(\mathbf{B}_t : t \geq 0)$  is a  $d$ -dimen. standard Brownian motion

- New diffusion process (c.f. Li et al., Mandt et al.):
  - Dimension  $n \rightarrow \infty$  instead of stepsize  $\gamma \rightarrow 0$  to create it
  - Explicitly solvable and don't need to send stepsize to 0

## Theorem: High dimensional equivalence of SGD

(C.P.-E. Paquette-B. Adlam-J.Pennington)

For any quadratic statistic  $q$ ,

$$\Pr\left(\sup_{0 \leq t \leq T} |q(\mathbf{x}_t) - \underbrace{q(\mathbf{X}_t)}_{\text{diffusion}}| > d^{-C}\right) \leq d^{-D}$$

concentration around mean:  $\underbrace{q(\mathbf{X}_t)}_{\text{diffusion}} \rightarrow \mathbb{E}[\underbrace{q(\mathbf{X}_t)}_{\text{diffusion}} \mid \mathbf{A}, \mathbf{b}, \mathbf{x}_0]$

# Dynamics of SGD: Loss function

**Theorem** (C.P.-Lee-E. Paquette-Pedregosa, COLT '21)

Let  $t = \#$  of passes through data. If  $\gamma(s) \rightarrow \tilde{\gamma} < 2 \left(\frac{1}{n} \text{tr}((\mathbf{H})^2(\mathbf{H} + \delta \mathbf{I})^{-1})\right)^{-1}$ ,

$$\Pr \left( \sup_{0 \leq t \leq T} |\mathcal{L}(\mathbf{x}_{\lfloor nt \rfloor}) - \psi_t| > d^{-C} \right) \leq d^{-D}$$

where  $\psi_t$  is solution to a **Volterra equation**

$$\psi_t = \underbrace{\mathcal{L}(\mathbf{X}_{\Gamma(t)}^{\text{gf}})}_{\text{gradient flow}} + \int_0^t \underbrace{\gamma^2(s) r h_2(t, s)}_{\text{noise term}} \psi_s \, ds$$

- $h_2(t, s) = \frac{1}{d} \sum_{i=1}^d \lambda_i^2 e^{-2(\Gamma(t) - \Gamma(s))(\lambda_i + \delta)}$ ,  $\lambda_i$  eigenvalues of  $\mathbf{H}$
- Integrated learning  $\Gamma(t) = \int_0^t \gamma(s) \, ds$



# Dynamics of SGD: Loss function

**Theorem** (C.P.-Lee-E. Paquette-Pedregosa, COLT '21)

Let  $t = \#$  of passes through data. If  $\gamma(s) \rightarrow \tilde{\gamma} < 2 \left(\frac{1}{n} \text{tr}((\mathbf{H})^2(\mathbf{H} + \delta \mathbf{I})^{-1})\right)^{-1}$ ,

$$\Pr \left( \sup_{0 \leq t \leq T} |\mathcal{L}(\mathbf{x}_{[nt]}) - \psi_t| > d^{-C} \right) \leq d^{-D}$$

where  $\psi_t$  is solution to a **Volterra equation**

$$\psi_t = \underbrace{\mathcal{L}(\mathbf{X}_{\Gamma(t)}^{\text{gf}})}_{\text{gradient flow}} + \int_0^t \underbrace{\gamma^2(s) r h_2(t, s)}_{\text{noise term}} \psi_s \, ds$$

- $h_2(t, s) = \frac{1}{d} \sum_{i=1}^d \lambda_i^2 e^{-2(\Gamma(t) - \Gamma(s))(\lambda_i + \delta)}$ ,  $\lambda_i$  eigenvalues of  $\mathbf{H}$
- Integrated learning  $\Gamma(t) = \int_0^t \gamma(s) \, ds$

## Rate of Convergence: Competition

- For **small**  $\gamma$ , same rate as **gradient descent**
- For **large**  $\gamma$ , **noise term** dominates rate

# Phase transition & Asymptotics: Fixed Stepsize

## Critical stepsize

$$\gamma_* = \frac{1}{\frac{r}{2} \int_0^\infty \frac{x^2}{x - \lambda_{\min}} d\mu(x)}$$

## Theorem

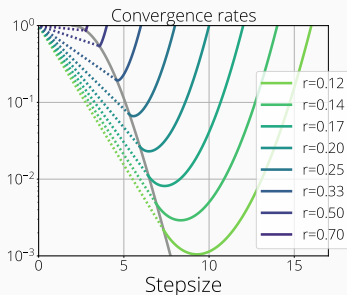
(C.P.-Lee-E. Paquette-Pedregosa, COLT '21)

For small  $\gamma < \gamma_*$ ,

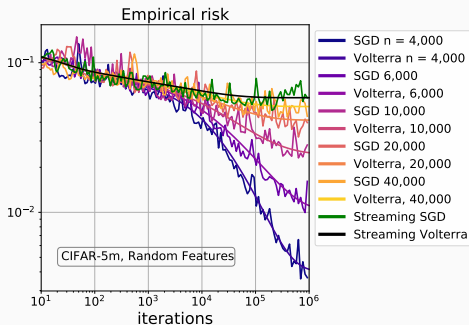
$$\psi_H(t) - \psi_H(\infty) \sim \frac{1}{t^\alpha} e^{-2\gamma t \lambda_{\min}}.$$

For large  $\gamma > \gamma_*$ ,  $\exists$  **non-linear**  $\lambda^*(\gamma)$

$$\text{and } \psi_H(t) - \psi_H(\infty) \sim \frac{1}{\gamma} e^{-2\gamma t \lambda^*(\gamma)}.$$



# Real world predictions: CIFAR-5m



Using a random features model to predict CIFAR-5m (Nakkiran et al., '21) car/plane, the **Volterra equation** (using the Hessian spectra as input) gives **good predictions** for behavior of SGD.

## Statistic: Expected risk

$$\mathcal{R}(\mathbf{x}_t) = \frac{1}{2} \mathbb{E}[(b - \mathbf{x}_t^T \mathbf{a})^2 | \mathbf{x}_t] \quad \text{where } (\mathbf{a}, b) \sim \mathcal{D}$$

**Theorem** (C.P.-Adlam-E. Paquette-Pennington)

Let  $t = \#$  of passes through data. If  $\gamma(t) \rightarrow \tilde{\gamma} < 2 \left( \frac{1}{n} \text{tr}((\mathbf{H})^2 (\mathbf{H} + \delta \mathbf{I})^{-1}) \right)^{-1}$ ,

$$\Pr \left( \sup_{0 \leq t \leq T} |\mathcal{R}(\mathbf{x}_{\lfloor tn \rfloor}) - \Omega_t| > d^{-C} \right) \leq d^{-D}$$

where  $\Omega_t$  equals

$$\Omega_t = \underbrace{\mathcal{R}(\mathbf{X}_{\Gamma(t)}^{\text{gf}})}_{\text{gradient flow}} + \int_0^t \underbrace{\gamma^2(s) \cdot K(t, s, \nabla^2 \mathcal{R}) \cdot \psi_s}_{\text{noise term}} ds,$$

and  $\psi_s$  limiting loss function  $\mathcal{L}$ ,  $K(t, s, \nabla^2 \mathcal{R})$  explicit kernel

# Statistic: Expected risk

$$\mathcal{R}(\mathbf{x}_t) = \frac{1}{2} \mathbb{E}[(b - \mathbf{x}_t^T \mathbf{a})^2 | \mathbf{x}_t] \quad \text{where } (\mathbf{a}, b) \sim \mathcal{D}$$

**Theorem** (C.P.-Adlam-E. Paquette-Pennington)

Let  $t = \#$  of passes through data. If  $\gamma(t) \rightarrow \tilde{\gamma} < 2 \left( \frac{1}{n} \text{tr}((\mathbf{H})^2 (\mathbf{H} + \delta \mathbf{I})^{-1}) \right)^{-1}$ ,

$$\Pr \left( \sup_{0 \leq t \leq T} |\mathcal{R}(\mathbf{x}_{\lfloor tn \rfloor}) - \Omega_t| > d^{-C} \right) \leq d^{-D}$$

where  $\Omega_t$  equals

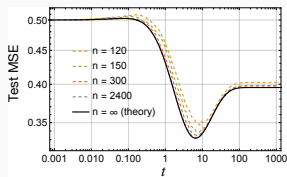
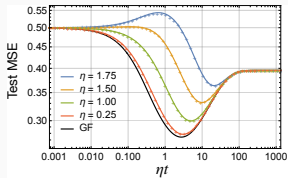
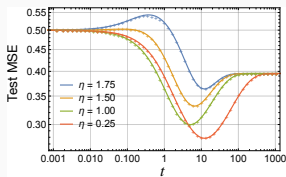
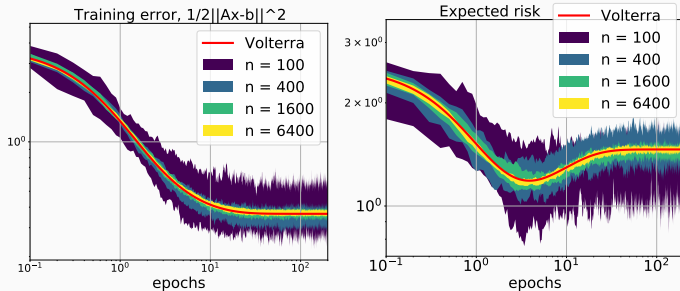
$$\Omega_t = \underbrace{\mathcal{R}(\mathbf{X}_{\Gamma(t)}^{\text{gf}})}_{\text{gradient flow}} + \int_0^t \underbrace{\gamma^2(s) \cdot K(t, s, \nabla^2 \mathcal{R}) \cdot \psi_s}_{\text{noise term}} ds,$$

and  $\psi_s$  limiting loss function  $\mathcal{L}$ ,  $K(t, s, \nabla^2 \mathcal{R})$  explicit kernel

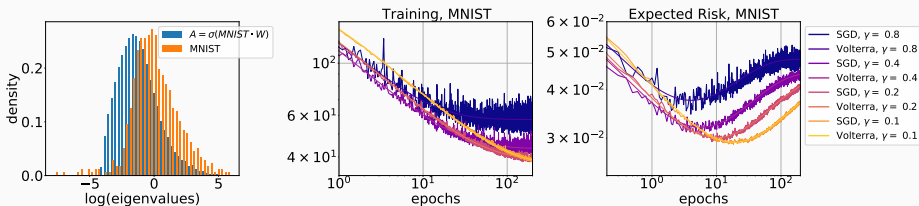
## Consequences

- Expected risk of SGD  $\geq$  Expected risk of gradient flow  
 $\Rightarrow$  no implicit regularization
- If  $\mathcal{L}(\mathbf{x}_t) \rightarrow 0$  or  $\gamma(t) \rightarrow 0$ , recover gradient flow at  $t \rightarrow \infty$

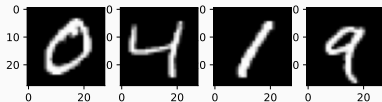
# Concentration effect



# Real world predictions: MNIST



Using a random features model to predict MNIST digits, the **Volterra equation** (using the Hessian spectra as input) gives **good predictions** for behavior of SGD.



**Question:** Can you go faster?



# SGD+M vanishing batch size

## Stochastic gradient descent with momentum (SGD+M)

$$\mathbf{y}_k = (1 - \theta)\mathbf{y}_{k-1} + \Gamma_1 \nabla f_{i_k}(\mathbf{x}_k)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{y}_k$$

✓  $\Gamma_1$  step size      ✓  $1 - \theta$  momentum parameter

## Homogenized SGD (Diffusion Process) (C.P.-E. Paquette, NeurIPS '21)

$$d\mathbf{X}_t = -\Gamma_1 \int_0^t e^{-n\theta(t-s)} d\mathbf{Z}_t$$

where  $d\mathbf{Z}_t = \nabla f(\mathbf{X}_t) dt + \sqrt{\frac{2}{n} \mathcal{L}(\mathbf{X}_t) \nabla^2 \mathcal{L}(\mathbf{X}_t)} d\mathbf{B}_t$

# SGD+M vanishing batch size

## Stochastic gradient descent with momentum (SGD+M)

$$\mathbf{y}_k = (1 - \theta)\mathbf{y}_{k-1} + \Gamma_1 \nabla f_{i_k}(\mathbf{x}_k)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{y}_k$$

✓  $\Gamma_1$  step size      ✓  $1 - \theta$  momentum parameter

## Homogenized SGD (Diffusion Process) (C.P.-E. Paquette, NeurIPS '21)

$$d\mathbf{X}_t = -\Gamma_1 \int_0^t e^{-n\theta(t-s)} d\mathbf{Z}_t$$

$$\text{where } d\mathbf{Z}_t = \nabla f(\mathbf{X}_t) dt + \sqrt{\frac{2}{n} \mathcal{L}(\mathbf{X}_t) \nabla^2 \mathcal{L}(\mathbf{X}_t)} d\mathbf{B}_t$$

Under **homogenized SGD**, **convolution-type Volterra equation**  $\psi_H$

(C.P.-E. Paquette, NeurIPS '21):

$$\psi_t = F(t) + \int_0^t \mathcal{I}(t-s)\psi_s ds.$$

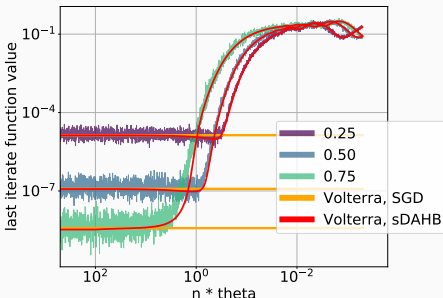
# Practical considerations: $\text{SGD}+\text{M} = \text{SGD}$ and $\text{SDAHB} \neq \text{SGD}$

## Experiment

- Fix matrix  $\mathbf{A}$ ; Run  $\text{SGD}+\text{M}$  for a fixed time and output function value
- Vary the  $\theta$  and  $\gamma$  holding  $\frac{\gamma}{\theta} = c$  fixed;  $c \in \{0.25, 0.5, 0.75\}$

## Conclusions

- For fixed  $\theta$ ,  $n$  large,  
$$\text{SGD}+\text{M}(\gamma, \theta) = \text{SGD}(\gamma/\theta)$$
- $\theta \sim 1/n$  barely faster than  $\text{SGD}$ , but not equivalent



# Batching + SGD+M

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \mathcal{L}(\mathbf{x}) \stackrel{\text{def}}{=} \sum_{i=1}^n \underbrace{\frac{1}{2} (\mathbf{a}_i \mathbf{x} - b_i)^2}_{f_i(\mathbf{x})} \right\},$$

with *batch*  $B \subset [n]$ , *batch fraction*  $\zeta \stackrel{\text{def}}{=} \frac{|B|}{n}$

**SGD+M with batches**     $\mathbf{y}_k = \Delta \cdot \mathbf{y}_{k-1} + \gamma \cdot \zeta \sum_{i_k \in B} \nabla f_{i_k}(\mathbf{x}_k)$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{y}_k$$

# Batching + SGD+M

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \mathcal{L}(\mathbf{x}) \stackrel{\text{def}}{=} \sum_{i=1}^n \underbrace{\frac{1}{2} (\mathbf{a}_i \mathbf{x} - b_i)^2}_{f_i(\mathbf{x})} \right\},$$

with *batch*  $B \subset [n]$ , *batch fraction*  $\zeta \stackrel{\text{def}}{=} \frac{|B|}{n}$

**SGD+M with batches**  $\mathbf{y}_k = \Delta \cdot \mathbf{y}_{k-1} + \gamma \cdot \zeta \sum_{i_k \in B} \nabla f_{i_k}(\mathbf{x}_k)$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{y}_k$$

## Results

- When  $\zeta \rightarrow 0$  as  $n \rightarrow \infty$ , SGD + M  $\equiv$  SGD
- When  $\zeta = 1$ , GD + M is faster than GD
  - $\mathcal{O}(\sqrt{\kappa})$  versus  $\mathcal{O}(\kappa)$ ,  $\kappa = \frac{\lambda_{\max}(\mathbf{H})}{\lambda_{\min}(\mathbf{H})}$

What happens in between?

# Concentration of SGD+M with batches

**Theorem** (C.P.-Lee-Cheng-E. Paquette)

$$\text{If } \gamma < \min \left\{ \frac{1+\Delta}{\zeta \lambda_{\max}(\mathbf{H})}, \frac{n(1-\Delta)}{(1-\zeta) \text{tr}(\mathbf{H})^{-1}} \right\},$$

$$\Pr \left( \sup_{0 \leq t \leq T} |\mathcal{L}(\mathbf{x}_t) - \psi(t)| > d^{-D} \right) \leq d^{-C}$$

where  $\psi(t)$  is solution to **(discrete) convolution-type Volterra**

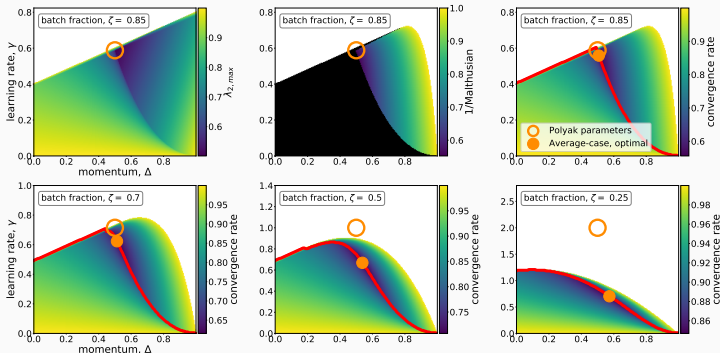
$$\psi(t+1) = \mathcal{L}(\mathbf{X}^{\text{gd}+\text{M}}) + \sum_{k=0}^t \underbrace{\gamma^2 \zeta (1-\zeta) H(t-k)}_{\text{noise}} \psi(k)$$

- $H(t-k)$  explicit function of eigenvalues of  $\mathbf{H}$
- $\Lambda$ , rate at which  $\mathcal{L}(\mathbf{X}^{\text{gd}+\text{M}})$  decrease
- For some  $\gamma$  and  $\Delta$ , **noise** term slows you down

# Convergence of SGD+M with batches

**Theorem** (C.P.-Lee-Cheng-E. Paquette)

$$\lim_{t \rightarrow \infty} (\mathcal{L}(x_t) - \mathcal{L}(\infty))^{1/t} = \max \left\{ \underbrace{\Lambda}_{\text{GD+M}}, \underbrace{\frac{-1}{3}}_{\text{noise}} \right\}$$



# Large vs Small batch: Convergence

$$\gamma = \frac{(1 - \sqrt{\Delta})^2}{\zeta \lambda_{\min}}, \quad \Delta = \max \left\{ \left( 1 - \frac{\zeta}{(1 - \zeta)\bar{\kappa}} \right)^2, \left( 1 - \frac{1}{\sqrt{\kappa}} \right)^2 \right\}$$

(average)  $\bar{\kappa} \stackrel{\text{def}}{=} \frac{\frac{1}{n} \text{tr}(\mathbf{H})}{\lambda_{\min}(\mathbf{H})}$ , *implicit conditioning ratio, ICR*  $\stackrel{\text{def}}{=} \frac{\bar{\kappa}}{\sqrt{\kappa}} = \frac{\text{average}}{\sqrt{\text{classic}}}$ .

## Phase transition

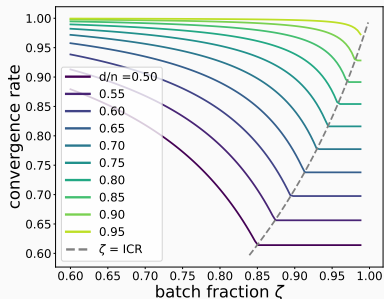
(C.P.-Lee-Cheng-E. Paquette)

- **Large batch:**  $\zeta \geq \text{ICR}$

SGD+M linearly at rate  $\mathcal{O}(1/\sqrt{\kappa})$   
and SGD+M accelerates

- **Small batch:**  $\zeta \leq \text{ICR}$

SGD+M linearly at rate  $\mathcal{O}(\zeta/\bar{\kappa})$   
SGD+M  $\Leftrightarrow$  SGD



**Saturating batch fraction** – after which increasing the batch fraction does not improve convergence.



# Thank you!

C. Paquette, E. Paquette, B. Adlam, J. Pennington. *Homogenization of SGD in high-dimensions: Exact dynamics and generalization properties*, [arxiv.org/pdf/2205.07069.pdf](https://arxiv.org/pdf/2205.07069.pdf)

K. Lee, A. Cheng, E. Paquette, C. Paquette. *Trajectory of Mini-Batch Momentum: Batch Size Saturation and Convergence in High Dimensions*, (submitted to NeurIPS '22)

C. Paquette, K. Lee, F. Pedregosa, E. Paquette. *SGD in the Large: Average-case Analysis, Asymptotics, and Step-size Criticality*, [arxiv.org/pdf/2102.04396.pdf](https://arxiv.org/pdf/2102.04396.pdf) (accepted at COLT 2021)

C. Paquette, E. Paquette. *Dynamics of Stochastic Momentum Methods on Large-scale, Quadratic Models*, [arxiv.org/pdf/2106.03696.pdf](https://arxiv.org/pdf/2106.03696.pdf) (accepted at NeurIPS 2021)