

Hitting the High-D(imensional) Notes:

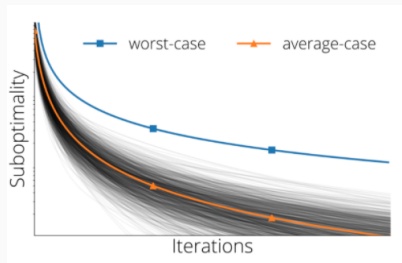
An ODE for SGD learning dynamics

Courtney Paquette (McGill & Google DeepMind)

Joint work: Elliot Paquette (McGill), Kiwon Lee (McGill), Elizabeth Collins-Woodfin (McGill), Inbar Seroussi (Tel-Aviv), Jeffrey Pennington (Google DeepMind), Ben Adlam (Google DeepMind), Andrew Cheng (PhD, Harvard)

Challenges/Problems for Optimization in Machine Learning

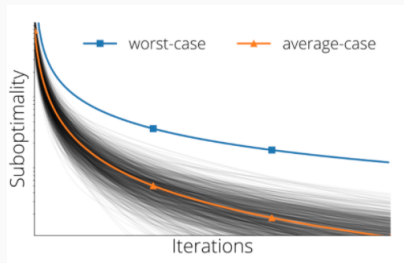
- ✓ Theory vs. Practice
- ✓ Mismatch with assumptions
→ too general in optimization theory
- ✓ Less convergence of algorithm



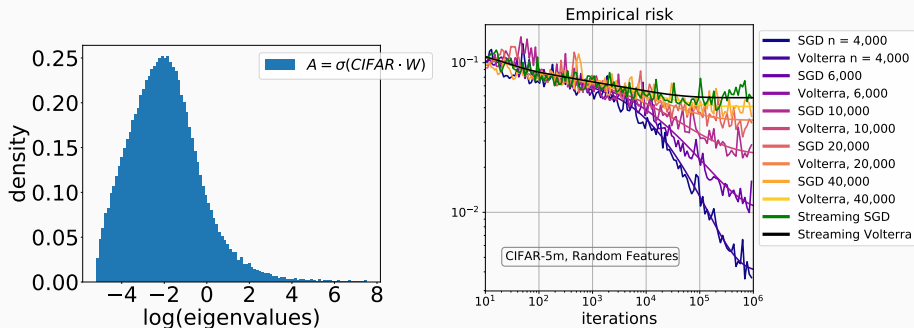
Challenges/Problems for Optimization in Machine Learning

✓ Theory vs. Practice

- ✓ Mismatch with assumptions
→ too general in optimization theory
- ✓ Less convergence of algorithm



Theory meets practice: CIFAR-5m



$$\min_x \|\sigma(\text{CIFAR} \cdot W)x - b\|^2$$

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.



Typical Machine Learning Problems

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

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

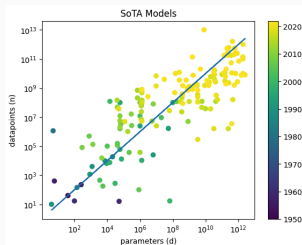
- ✓ State-of-the-art models have **millions/billions parameters**
 - Meena: 2.6 billion, Turing NLG: 17 billion, GPT-3: 175 billion

Typical Machine Learning Problems

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

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) is constant, $d^\alpha \leq n \leq d^{1/\alpha}$



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: Some 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)

High-dimensional linear composites

$$\min_{X \in \mathbb{R}^{d \times m}} \{ \mathcal{R}(X) = \mathbb{E}_{a, \epsilon} [f(a^T X; a^T X^*, \epsilon)] \}$$

- $a \sim N(0, K)$ data
- Covariance $K = \mathbb{E}[aa^T]$, $\|K\|_{\text{op}}$ bounded, independent of d
- ϵ is label noise, $X^* \in \mathbb{R}^{d \times m^*}$
- **Idea:** Think of $f : \mathbb{R}^m \rightarrow \mathbb{R}$ as pseudo-Lipschitz and **low dimensional**, i.e. $m, m^* \ll d$ **large**

High-dimensional linear composites

$$\min_{X \in \mathbb{R}^{d \times m}} \{ \mathcal{R}(X) = \mathbb{E}_{a, \epsilon} [f(a^T X; a^T X^*, \epsilon)] \}$$

- $a \sim N(0, K)$ data
- Covariance $K = \mathbb{E}[aa^T]$, $\|K\|_{\text{op}}$ bounded, independent of d
- ϵ is label noise, $X^* \in \mathbb{R}^{d \times m^*}$
- **Idea:** Think of $f : \mathbb{R}^m \rightarrow \mathbb{R}$ as pseudo-Lipschitz and **low dimensional**, i.e. $m, m^* \ll d$ **large**

What does this allow?

- GLMs, multi-index models
- e.g., multi-class logistic regression $m =$ number of classes
- X^* – ground truth signal, mean of data, fixed vector

One-pass SGD

$$\min_{X \in \mathbb{R}^{d \times m}} \{ \mathcal{R}(X) = \mathbb{E}_{a, \epsilon} [f(a^T X; a^T X^*, \epsilon)] \}$$

One-pass SGD: Generate new $(a_{k+1}, \epsilon_{k+1})$

$$X_{k+1} = X_k - \frac{\gamma_k}{d} a_{k+1} \otimes \nabla f(a_{k+1}^T X_k),$$

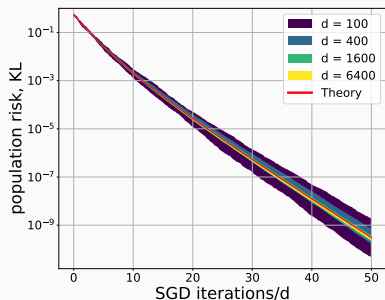
One-pass SGD

$$\min_{X \in \mathbb{R}^{d \times m}} \{ \mathcal{R}(X) = \mathbb{E}_{a, \epsilon} [f(a^T X; a^T X^*, \epsilon)] \}$$

One-pass SGD: Generate new $(a_{k+1}, \epsilon_{k+1})$

$$X_{k+1} = X_k - \frac{\gamma_k}{d} a_{k+1} \otimes \nabla f(a_{k+1}^T X_k),$$

For **large models**, as



$$\frac{\text{parameters}}{\text{samples}} = \frac{d}{n} \rightarrow r,$$

- $\mathcal{R}(X_k) \xrightarrow{\text{Pr}}$ (smooth function)
- Analyze this smooth function
- Determined by the spectrum of the covariance matrix $K = \mathbb{E}[aa^T]$

Why?

- High-dimensional compositional structure, $a^T X$ averages out d

Why?

- **High-dimensional compositional structure**, $a^T X$ averages out d
- Let $W \stackrel{\text{def}}{=} [X|X^*]$. Lots of statistics can be represented by

$$\varphi(X) = g(W^T q(K) W) = g \left(\begin{bmatrix} X^T q(K) X & X^T q(K) X^* \\ (X^*)^T q(K) X & (X^*)^T q(K) X^* \end{bmatrix} \right),$$

where q is a polynomial, $K = \mathbb{E}[aa^T]$.

Examples: Risk $\mathcal{R}(X)$, $\|\nabla \mathcal{R}\|^2$, distance to optimality $\|X - X^*\|^2$

Least squares:

$$\mathcal{R}(X) = \mathbb{E}_{a,\epsilon}[\text{tr}((a^T X - (a^T X^* + \epsilon))^2)] = \text{tr}((X - X^*)^T K (X - X^*)) + \mathbb{E}[\text{tr}(\epsilon \epsilon^T)]$$

Least squares:

$$\mathcal{R}(X) = \mathbb{E}_{a, \epsilon} [\text{tr}((a^T X - (a^T X^* + \epsilon))^2)] = \text{tr}((X - X^*)^T K (X - X^*)) + \mathbb{E} [\text{tr}(\epsilon \epsilon^T)]$$

Logistic Loss (binary, noiseless): Student-teacher generate targets

$y = \frac{\exp(a^T X^*)}{\exp(a^T X^*) + 1}$. Then the risk becomes

$$\mathcal{R}(X) = \mathbb{E}_a \left[-a^T X \cdot \frac{\exp(a^T X^*)}{\exp(a^T X^*) + 1} + \log(\exp(a^T X) + 1) \right].$$

Examples

Least squares:

$$\mathcal{R}(X) = \mathbb{E}_{a, \epsilon} [\text{tr}((a^T X - (a^T X^* + \epsilon))^2)] = \text{tr}((X - X^*)^T K (X - X^*)) + \mathbb{E} [\text{tr}(\epsilon \epsilon^T)]$$

Logistic Loss (binary, noiseless): Student-teacher generate targets

$y = \frac{\exp(a^T X^*)}{\exp(a^T X^*) + 1}$. Then the risk becomes

$$\mathcal{R}(X) = \mathbb{E}_a \left[-a^T X \cdot \frac{\exp(a^T X^*)}{\exp(a^T X^*) + 1} + \log(\exp(a^T X) + 1) \right].$$

Why? Let $a^T X \sim N(0, X^T K X)$

$$\mathbb{E}_a [\log(\exp(a^T X) + 1)] = \mathbb{E}_w [\log(\exp(\sqrt{X^T K X} w) + 1)], \quad w \sim N(0, 1)$$

Dynamics of Covariance Matrix

Goal: Understand the impact of SGD noise

Lots of statistics can be represented by

$$\varphi(X) = g(W^T q(K)W) = g \left(\begin{bmatrix} X^T q(K)X & X^T q(K)X^* \\ (X^*)^T q(K)X & (X^*)^T q(K)X^* \end{bmatrix} \right),$$

where q is a polynomial, $K = \mathbb{E}[aa^T]$.

Dynamics of Covariance Matrix

Goal: Understand the impact of SGD noise

Lots of statistics can be represented by

$$\varphi(X) = g(W^T q(K)W) = g \left(\begin{bmatrix} X^T q(K)X & X^T q(K)X^* \\ (X^*)^T q(K)X & (X^*)^T q(K)X^* \end{bmatrix} \right),$$

where q is a polynomial, $K = \mathbb{E}[aa^T]$.

Take away: To understand dynamics of SGD amounts to understanding **low-dimensional covariance matrix**

for any polynomial q , $W^T q(K)W \Rightarrow S(W; z) \stackrel{\text{def}}{=} W^T R(z; K)W$

where $R(z; K) = (z - K)^{-1}$ resolvent of K , $z \in \mathbb{C} \setminus (\text{spectrum } K)$.

Dynamics of Covariance Matrix

Goal: Understand the impact of SGD noise

Lots of statistics can be represented by

$$\varphi(X) = g(W^T q(K)W) = g \left(\begin{bmatrix} X^T q(K)X & X^T q(K)X^* \\ (X^*)^T q(K)X & (X^*)^T q(K)X^* \end{bmatrix} \right),$$

where q is a polynomial, $K = \mathbb{E}[aa^T]$.

Take away: To understand dynamics of SGD amounts to understanding **low-dimensional covariance matrix**

for any polynomial q , $W^T q(K)W \Rightarrow S(W; z) \stackrel{\text{def}}{=} W^T R(z; K)W$

where $R(z; K) = (z - K)^{-1}$ resolvent of K , $z \in \mathbb{C} \setminus (\text{spectrum } K)$.

Remark:

$$W^T q(K)W = \frac{1}{2\pi i} \oint_{\Gamma} q(z) W^T R(z; K)W \, dz$$

Main Result

Continuous time scale: iterates of SGD $k = td$, where $t \in \mathbb{R}$

$S(W; z) \stackrel{\text{def}}{=} W^T R(z; K) W$, where $R(z; K)$ is resolvent of K , $W = [X|X^*]$

Main Result

Continuous time scale: iterates of SGD $k = td$, where $t \in \mathbb{R}$

$S(W; z) \stackrel{\text{def}}{=} W^T R(z; K) W$, where $R(z; K)$ is resolvent of K , $W = [X|X^*]$

Theorem: High dimensional concentration of SGD

(C.P.-E.Collins-Woodfin-I. Seroussi-E. Paquette)

For a fixed $T > 0$,

$$\Pr \left(\sup_{0 \leq t \leq T} \left\| \underbrace{S(W_{\lfloor td \rfloor}; z)}_{\text{SGD}} - \underbrace{\mathcal{S}(t; z)}_{\text{deterministic function}} \right\| > d^{-C} \right) \leq d^{-D}$$

and, for any $\varphi(X) = g(W^T q(K) W)$,

$$\Pr \left(\sup_{0 \leq t \leq T} \left| \underbrace{\varphi(X_{\lfloor td \rfloor})}_{\text{SGD}} - \underbrace{\phi(t)}_{\text{deterministic function}} \right| > d^{-C} \right) \leq d^{-D}$$

Main Result

Continuous time scale: iterates of SGD $k = td$, where $t \in \mathbb{R}$

$S(W; z) \stackrel{\text{def}}{=} W^T R(z; K) W$, where $R(z; K)$ is resolvent of K , $W = [X|X^*]$

Theorem: High dimensional concentration of SGD

(C.P.-E. Collins-Woodfin-I. Seroussi-E. Paquette)

For a fixed $T > 0$,

$$\Pr \left(\sup_{0 \leq t \leq T} \left\| \underbrace{S(W_{\lfloor td \rfloor}; z)}_{\text{SGD}} - \underbrace{\mathcal{S}(t; z)}_{\text{deterministic function}} \right\| > d^{-C} \right) \leq d^{-D}$$

and, for any $\varphi(X) = g(W^T q(K) W)$,

$$\Pr \left(\sup_{0 \leq t \leq T} \left| \underbrace{\varphi(X_{\lfloor td \rfloor})}_{\text{SGD}} - \underbrace{\phi(t)}_{\text{deterministic function}} \right| > d^{-C} \right) \leq d^{-D}$$

- ✓ **Deterministic function** \mathcal{S} defined by ODE (see later)
- ✓ **Statistic**

$$\phi(t) = g \left(\oint_{\Gamma} q(z) \mathcal{S}(t; z) dz \right)^*$$

* ϕ will also satisfy an ODE.

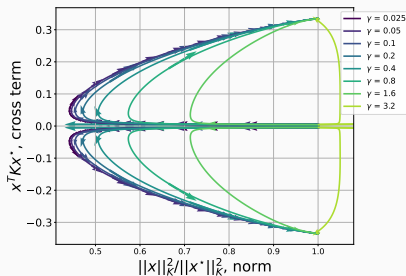
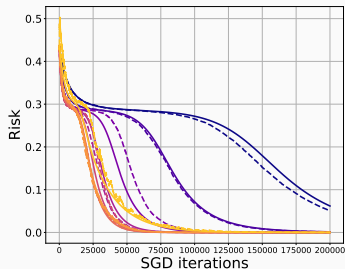
Example: Phase Retrieval

Problem: Find $x \in \mathbb{R}^d$ such that

$$(a^T X)^2 \approx (a^T X^*)^2, \quad a \in \mathbb{R}^d \quad \text{and} \quad X^* \in \mathbb{R}^d \quad \text{true signal}$$

Optimization formulation:

$$\min_{X \in \mathbb{R}^d} \{ \mathcal{R}(X) \stackrel{\text{def}}{=} \mathbb{E}_a [((a^T X)^2 - (a^T X^*)^2)^2] \}$$



Exact Dynamics Idea: Diffusion Approximation

Time scale: k iterates of SGD = td , where $t \in \mathbb{R}$

Homogenized SGD (C.P.-E. Collins-Woodfin-I. Seroussi-E. Paquette)

$$d\mathcal{X}_t = -\gamma_t \nabla_X \mathcal{R}(\mathcal{X}_t) dt + \gamma_t \langle \sqrt{K/d} \otimes \sqrt{\mathbb{E}_{a,\epsilon}[\nabla_x f(a^T X; a^T X^*, \epsilon)^{\otimes 2}]} \rangle, dB_t \rangle_{\mathcal{A} \otimes \mathcal{O}}$$

$\mathcal{X}_0 = X_0$ and $(B_t : t \geq 0)$ is a d -dimen. standard Brownian motion

- New diffusion process (Li et al., Mandt et al.)
- **Continuous time is made by $d \rightarrow \infty$** instead of stepsize $\gamma \rightarrow 0$, $t = 1$ means n SGD updates
- ~~\mathcal{X}_t mean/covariance same as $X_{[td]}^{\text{SGD}}$~~ , **Goal:** $\varphi(\mathcal{X}_t) \approx \varphi(X_{[td]}^{\text{SGD}})$

Exact Dynamics Idea: Diffusion Approximation

Time scale: k iterates of SGD = td , where $t \in \mathbb{R}$

Homogenized SGD (C.P.-E.Collins-Woodfin-I. Seroussi-E. Paquette)

$$d\mathcal{X}_t = -\gamma_t \nabla_X \mathcal{R}(\mathcal{X}_t) dt + \gamma_t \langle \sqrt{K/d} \otimes \sqrt{\mathbb{E}_{a,\epsilon}[\nabla_X f(a^T X; a^T X^*, \epsilon)^{\otimes 2}]} \rangle, dB_t \rangle_{\mathcal{A} \otimes \mathcal{O}}$$

$\mathcal{X}_0 = X_0$ and $(B_t : t \geq 0)$ is a d -dimen. standard Brownian motion

- New diffusion process (Li et al., Mandt et al.)
- **Continuous time is made by $d \rightarrow \infty$** instead of stepsize $\gamma \rightarrow 0$, $t = 1$ means n SGD updates
- ~~\mathcal{X}_t mean/covariance same as $X_{[td]}^{\text{SGD}}$, Goal: $\varphi(\mathcal{X}_t) \approx \varphi(X_{[td]}^{\text{SGD}})$~~

Theorem: High dimensional equivalence of SGD

(C.P.-E.Collins-Woodfin-I. Seroussi-E. Paquette)

For any $\varphi(X) = g(W^T q(K)W)$,

$$\Pr \left(\sup_{0 \leq t \leq T} |\varphi(X_{[td]}) - \underbrace{\varphi(\mathcal{X}_t)}_{\text{diffusion}}| > d^{-C} \right) \leq d^{-D}$$

Deterministic equivalent of φ

Assumption: Statistic $\varphi(X) = g(W^T q(K)W)$

Homogenized SGD

$$\begin{aligned}d\mathcal{X}_t &= -\gamma_t \nabla \mathcal{R}(\mathcal{X}_t) dt \\ &+ \gamma_t \langle \sqrt{K/d} \otimes \sqrt{\mathbb{E}_{a,\epsilon}[\nabla f(a^T \mathcal{X}_t)^{\otimes 2}]}, dB_t \rangle_{\mathcal{A} \otimes \mathcal{O}}\end{aligned}$$

Deterministic equivalent of φ

Assumption: Statistic $\varphi(X) = g(W^T q(K)W)$

Homogenized SGD

$$d\mathcal{X}_t = -\gamma_t \nabla \mathcal{R}(\mathcal{X}_t) dt \\ + \gamma_t \langle \sqrt{K/d} \otimes \sqrt{\mathbb{E}_{a,\epsilon}[\nabla f(a^T \mathcal{X}_t)^{\otimes 2}]}, dB_t \rangle_{\mathcal{A} \otimes \mathcal{O}}$$

Intuition: Apply Itô Calculus

$$d\varphi(\mathcal{X}_t) = -\gamma_t \langle \nabla \varphi(\mathcal{X}_t), \nabla \mathcal{R}(\mathcal{X}_t) \rangle dt \quad (\text{grad. flow}) \\ + \frac{\gamma_t^2}{2d} \langle \nabla^2 \varphi(\mathcal{X}_t), K \otimes \mathbb{E}_{a,\epsilon}[\nabla f(a^T \mathcal{X}_t)^{\otimes 2}] \rangle dt \quad (\text{SGD noise}) \\ + \text{incremental martingale} \\ \stackrel{\text{def}}{=} \mathcal{F}(S(\mathcal{W}_t; z)) + \text{incremental martingale}$$

All the quantities are functions

$$S(\mathcal{W}_t; z) \stackrel{\text{def}}{=} \mathcal{W}_t^T R(z; K) \mathcal{W}_t, \quad \mathcal{W}_t = [\mathcal{X}_t | X^*].$$

Deterministic equivalent

Idea: Set $\varphi(\mathcal{X}_t) = S(\mathcal{W}_t; z)$ in Ito (and drop martingale)

$$dS(\mathcal{W}_t; z) = \mathcal{G}(S(\mathcal{W}_t; z)) + \text{incremental martingale}$$

\Rightarrow ODE for deterministic equivalent $\mathcal{S}(t, z)$ for $S(\mathcal{W}_t, z)$

$$d\mathcal{S}(t, z) = \mathcal{G}(\mathcal{S}(t, z))$$

Deterministic equivalent

Idea: Set $\varphi(\mathcal{X}_t) = S(\mathcal{W}_t; z)$ in Ito (and drop martingale)
$$dS(\mathcal{W}_t; z) = \mathcal{G}(S(\mathcal{W}_t; z)) + \text{incremental martingale}$$
$$\Rightarrow \text{ODE for deterministic equivalent } \mathcal{S}(t, z) \text{ for } S(\mathcal{W}_t, z)$$
$$d\mathcal{S}(t, z) = \mathcal{G}(\mathcal{S}(t, z))$$

ODE for Deterministic Equivalent:

$$\begin{aligned} dS(\mathcal{W}_t; z) = \mathcal{G}(S(\mathcal{W}_t; z)) &\Rightarrow d\mathcal{S}(t; z) = \mathcal{G}(\mathcal{S}(t; z)) \quad \Leftarrow \text{solve numerically} \\ d\varphi(\mathcal{X}_t) = \mathcal{F}(S(\mathcal{W}_t; z)) &\Rightarrow d\phi(t) = \mathcal{F}(\mathcal{S}(t; z)) \end{aligned}$$

where ϕ is the deterministic equivalent of $\varphi(X)$

Theorem: High dimensional equivalence of SGD

(C.P.-E.Collins-Woodfin-I. Seroussi-E. Paquette)

For any $\varphi(X) = g(W^T q(K)W)$,

$$\Pr \left(\sup_{0 \leq t \leq T} |\varphi(X_{\lfloor td \rfloor}) - \underbrace{\varphi(\mathcal{X}_t)}_{\text{diffusion}}| > d^{-C} \right) \leq d^{-D}$$

and

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

Why Interesting?

Optimization Question

What choice of learning rate ensures distance to optimality decreases at each iteration of SGD?

→ Can't do this with SGD because of the stochasticity in the gradients

→ Can ask on the **deterministic equivalent** of the distance to optimality, $\|X - X^*\|^2$,

*What **stepsize** is needed for $\|X - X^*\|^2$ to be a decreasing function?*

Intuition-Critical Threshold

Deterministic equivalent of $\|X - X^*\|^2$:

$$\sup_{0 \leq t \leq T} \left| \|X_{\lfloor td \rfloor} - X^*\|^2 - \mathcal{D}^2(t) \right| \leq d^{-\varepsilon}$$

where

$$\frac{d}{dt} \mathcal{D}^2 = -2\gamma_t \mathbf{A}(\mathcal{J}) + \frac{\gamma_t^2}{d} \text{tr}(K) \mathbf{I}(\mathcal{J}), \quad \begin{cases} \mathbf{A}(\mathcal{J}) = \mathbb{E}_{a, \epsilon}[\langle x - x^*, \nabla f(x; x^*) \rangle], \\ \mathbf{I}(\mathcal{J}) = \mathbb{E}_{a, \epsilon}[\|\nabla f(x; x^*)\|^2], \\ \text{where } (x \oplus x^*) \sim N(0, W^T K W). \end{cases}$$

Intuition-Critical Threshold

Deterministic equivalent of $\|X - X^*\|^2$:

$$\sup_{0 \leq t \leq T} |\|X_{\lfloor td \rfloor} - X^*\|^2 - \mathcal{D}^2(t)| \leq d^{-\varepsilon}$$

where

$$\frac{d}{dt} \mathcal{D}^2 = -2\gamma_t A(\mathcal{S}) + \frac{\gamma_t^2}{d} \text{tr}(K) I(\mathcal{S}), \quad \begin{cases} A(\mathcal{S}) = \mathbb{E}_{a, \epsilon} [\langle x - x^*, \nabla f(x; x^*) \rangle], \\ I(\mathcal{S}) = \mathbb{E}_{a, \epsilon} [\|\nabla f(x; x^*)\|^2], \\ \text{where } (x \oplus x^*) \sim N(0, W^T K W). \end{cases}$$

Critical learning rate ($d\mathcal{D}^2 < 0$)

$$\gamma_t^{\text{critical}} = \frac{2A(\mathcal{S}(t; z))}{\frac{\text{tr}(K)}{d} I(\mathcal{S}(t; z))}$$

Intuition-Critical Threshold

Deterministic equivalent of $\|X - X^*\|^2$:

$$\sup_{0 \leq t \leq T} \left| \|X_{\lfloor td \rfloor} - X^*\|^2 - \mathcal{D}^2(t) \right| \leq d^{-\varepsilon}$$

where

$$\frac{d}{dt} \mathcal{D}^2 = -2\gamma_t A(\mathcal{S}) + \frac{\gamma_t^2}{d} \text{tr}(K) I(\mathcal{S}), \quad \begin{cases} A(\mathcal{S}) = \mathbb{E}_{a, \epsilon} [\langle x - x^*, \nabla f(x; x^*) \rangle], \\ I(\mathcal{S}) = \mathbb{E}_{a, \epsilon} [\|\nabla f(x; x^*)\|^2], \\ \text{where } (x \oplus x^*) \sim N(0, W^T K W). \end{cases}$$

Critical learning rate ($d\mathcal{D}^2 < 0$)

$$\gamma_t^{\text{critical}} = \frac{2A(\mathcal{S}(t; z))}{\frac{\text{tr}(K)}{d} I(\mathcal{S}(t; z))} \geq \frac{2q}{\frac{\text{tr}(K)}{d}}, \quad \text{where } \frac{A(\mathcal{S}(t; z))}{I(\mathcal{S}(t; z))} \geq q$$

- Functions $A(\mathcal{S})$ and $I(\mathcal{S})$ – don't carry K or d
- Lower bound A and I based on convexity/smoothness assumptions of f
- Critical stepsize depends on **average eigenvalue of K**

Theorem: Convergence of strongly convex

(C.P.-E. Collins-Woodfin-I. Seroussi-E. Paquette)

Suppose f is $\hat{\mu}$ -strongly convex and \hat{L} -Lipschitz gradients. for some $0 < \zeta < 1$, then for all $t \geq 0$

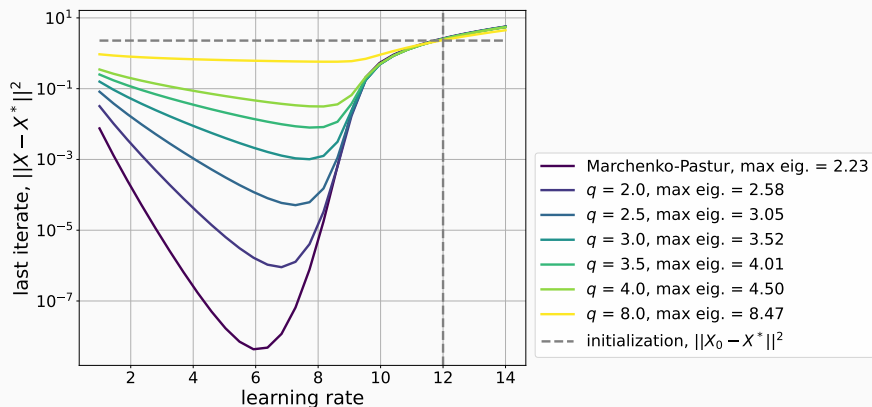
$$\mathcal{D}^2(t) \leq e^{-at} \mathcal{D}^2(0),$$

where $\gamma < \gamma^{\text{critical}}$

$$\text{(convergence rate)} \quad a = \frac{\hat{\mu}^2}{\hat{L}^2} \cdot \frac{\lambda_{\min}(K)}{\frac{1}{d} \text{tr}(K)}$$

Logistic regression

$$\gamma_t^{\text{critical}} = \frac{2A(\mathcal{L}(t; z))}{\frac{\text{tr}(K)}{d} l(\mathcal{L}(t; z))}$$



Caption: Covariance matrix $K = \text{diag}(\sigma_i^{2q} : i = 1, \dots, 1000)$, $\text{tr}(K)/d = 1$.

Focus on the **Least Squares** Problem with extensions
(e.g., multi-pass)

Our framework

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

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

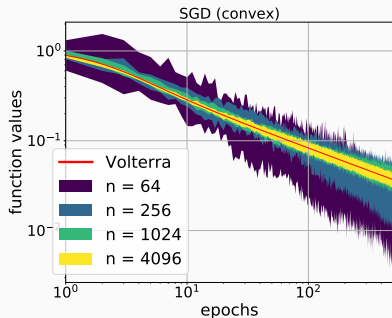
Multi-pass SGD Select index i_k $X_{k+1} = X_k - \gamma_k \nabla f_{i_k}(X_k)$

Our framework

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

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

Multi-pass SGD Select index i_k $X_{k+1} = X_k - \gamma_k \nabla f_{i_k}(X_k)$



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

- $\mathcal{L}(X_k) \xrightarrow{\text{Pr}} \mathcal{L}(t)$ (smooth function)
- Determined by the spectrum of the Hessian
- **Homogenized SGD (C.P.-E. Paquette, NeurIPS '21 & Mori '21)**

$$d\mathcal{X}_t = -\gamma(t) \nabla \mathcal{L}(\mathcal{X}_t) dt + \gamma_t \sqrt{\frac{2}{n} \mathcal{L}(\mathcal{X}_t) A^T A} dB_t$$

Remove Gaussian Assumption...

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

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 λ_{\max} and λ_{\min} of \mathbf{H}
3. De-localization of eigenvectors of $\mathbf{A}\mathbf{A}^T$: eigenvectors are not aligned with the unit vectors

e.g., if $A_{i,j} \sim N(0, 1)$, then eigenvectors of $\mathbf{H} \sim \text{Unif}(\mathbb{S}^{d-1})$

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

Real world predictions: MNIST

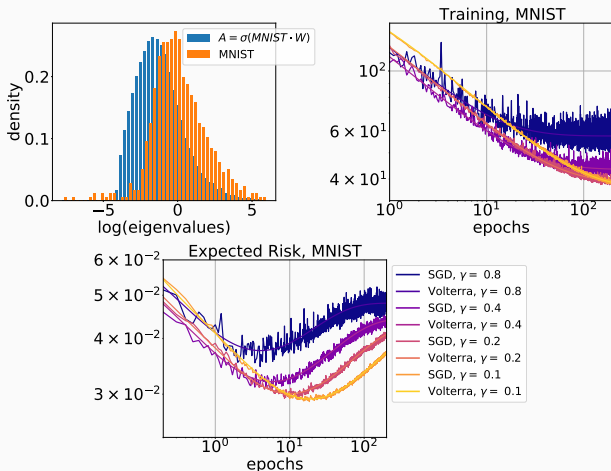
Expected Risk:

$$\mathcal{R}(X_k) = \frac{1}{2} \mathbb{E}[(b - X_k^T a)^2 | X_k] \quad \text{where } (a, b) \sim \mathcal{D}, \quad X_k = \text{SGD iterate on } \mathcal{L}$$

Real world predictions: MNIST

Expected Risk:

$$\mathcal{R}(X_k) = \frac{1}{2} \mathbb{E}[(b - X_k^T a)^2 | X_k] \quad \text{where } (a, b) \sim \mathcal{D}, \quad X_k = \text{SGD iterate on } \mathcal{L}$$



Phase transition & Asymptotics

Phase transition stepsize

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

Theorem

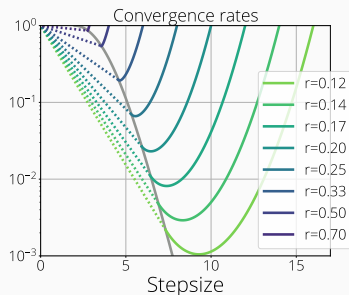
(C.P.-Lee-E. Paquette-Pedregosa, COLT '21)

For small $\gamma < \gamma_*$,

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

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

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



Large batch: SGD+M

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

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

Multi-pass SGD $Y_k = \Delta \cdot Y_{k-1} + \gamma \cdot \zeta \sum_{i_k \in B} \nabla f_{i_k}(X_k)$
+ momentum (SGD+M) $X_{k+1} = X_k - Y_k$

Large batch: SGD+M

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

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

$$\begin{array}{l} \text{Multi-pass SGD} \\ + \text{ momentum (SGD+M)} \end{array} \quad \begin{array}{l} Y_k = \Delta \cdot Y_{k-1} + \gamma \cdot \zeta \sum_{i_k \in B} \nabla f_{i_k}(X_k) \\ X_{k+1} = X_k - Y_k \end{array}$$

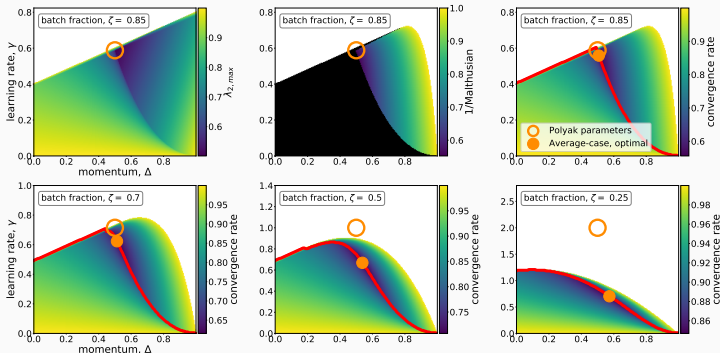
- Concentration of $\mathcal{L}(X_k) \rightarrow \mathcal{L}(k)$ deterministic, discrete **not** continuous

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

Convergence of SGD+M with batches

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

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



Condition numbers

$$\text{(average)} \quad \bar{\kappa} \stackrel{\text{def}}{=} \frac{\frac{1}{n} \text{tr}(A^T A)}{\lambda_{\min}(A^T A)} < \frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)} \stackrel{\text{def}}{=} \kappa \quad \text{(classic)}$$

Condition numbers

$$\text{(average)} \quad \bar{\kappa} \stackrel{\text{def}}{=} \frac{\frac{1}{n} \text{tr}(A^T A)}{\lambda_{\min}(A^T A)} < \frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)} \stackrel{\text{def}}{=} \kappa \quad \text{(classic)}$$

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

Suppose (near optimal parameters)

$$\gamma = \frac{(1 - \sqrt{\Delta})^2}{\zeta \lambda_{\min}(A^T A)}, \quad \Delta = \max \left\{ \underbrace{\left(1 - \frac{\zeta}{(1 - \zeta)\bar{\kappa}}\right)^2}_{\text{rate of SGD}}, \underbrace{\left(1 - \frac{1}{\sqrt{\bar{\kappa}}}\right)^2}_{\text{rate of GD+M}} \right\}$$

Then

$$\lim_{k \rightarrow \infty} (\mathcal{L}(k) - \mathcal{L}(\infty))^{1/k} = \Delta$$

Large vs Small batch: Convergence

$$\lim_{k \rightarrow \infty} (\mathcal{L}(k) - \mathcal{L}(\infty))^{1/k} = \max \left\{ \left(1 - \frac{\zeta}{(1 - \zeta)\bar{\kappa}} \right)^2, \left(1 - \frac{1}{\sqrt{\kappa}} \right)^2 \right\}$$

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

Large vs Small batch: Convergence

$$\lim_{k \rightarrow \infty} (\mathcal{L}(k) - \mathcal{L}(\infty))^{1/k} = \max \left\{ \left(1 - \frac{\zeta}{(1 - \zeta)\bar{\kappa}} \right)^2, \left(1 - \frac{1}{\sqrt{\kappa}} \right)^2 \right\}$$

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

Phase transition

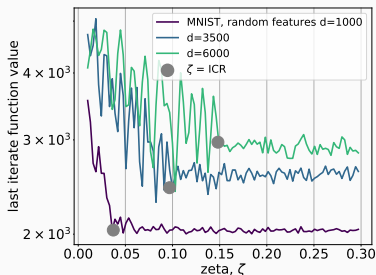
(C.P.-Lee-Cheng-E. Paquette)

- **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

K. Lee, A. Cheng, E. Paquette, C. Paquette. *Trajectory of Mini-Batch Momentum: Batch Size Saturation and Convergence in High Dimensions*,
arxiv.org/pdf/2206.01029.pdf (NeurIPS 2022)

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 (COLT 2021)

C. Paquette, E. Paquette. *Hitting the High-Dimensional Notes: An ODE for SGD Learning dynamics on GLMs and multi-index models*,
arxiv.org/pdf/2308.08977.pdf