

High-Dimensional Dynamics of SGD on Structured Data

Inbar Seroussi

Joint work with Elizabeth Colins-Woodfin (University of Oregon),

Courtney Paquette (McGill University) and Elliot Paquette (McGill University)

Statistical Physics and Machine Learning: Moving forward, Cargèse 2025, 6.8.2025



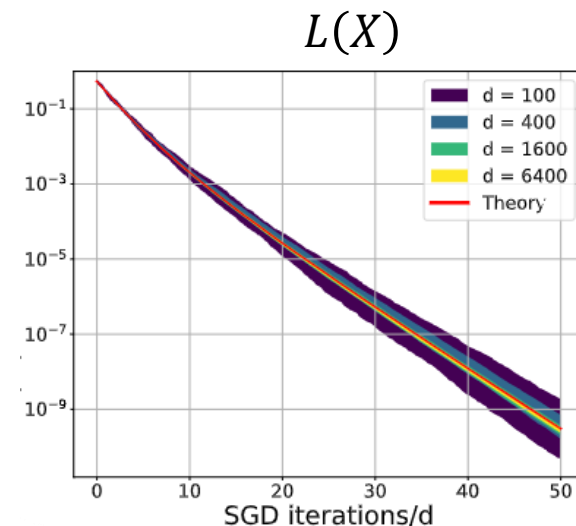
Goal and motivation for today

- Understand the **high-dimensional** dynamics of online **SGD** on **real** data
- What is the role of structure/high dimensionality in the dynamics?

Given data - $\{(\mathbf{a}_i; \mathbf{y}_i)\}_{i=1}^n$, $X \in \mathbb{R}^p$ is a set of learnable parameters with SGD

$$\min_{X \in \mathbb{R}^p} \{L(X) = \mathbb{E}_{(\mathbf{a}, \mathbf{y})}[f(X; \mathbf{a}, \mathbf{y})]\} \longleftrightarrow \min_{X \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n f(X; \mathbf{a}_i; \mathbf{y}_i)$$

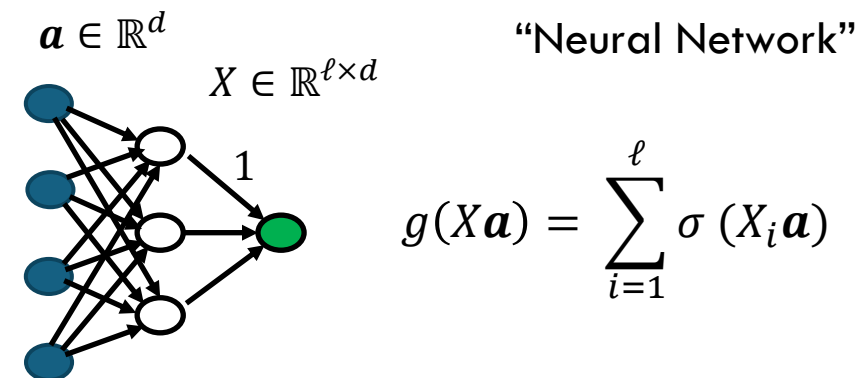
- Exact prediction of the dynamics
- Condition on feature learning, scaling, stability, classification capabilities



SGD dynamics - One class

Deterministic equivalence and more

Data and estimator setup



Target (Teacher) model:

$$\mathbf{y}_i = \phi(\mathbf{X}^* \mathbf{a}_i; \boldsymbol{\varepsilon}_i), \quad \boldsymbol{\varepsilon}_i \text{ - i.i.d. noise with bounded variance}$$

with a true matrix $\mathbf{X}^* \in \mathbb{R}^{\ell^* \times d}$ and $\phi: \mathbb{R}^{\ell^*} \rightarrow \mathbb{R}^m$, with $\ell^* = O_d(1)$

$$\mathbf{a}_i \sim \mathcal{N}(0, \mathbf{K}), \text{ with } \mathbf{K} \in \mathbb{R}^{d \times d}, \|\mathbf{K}\|_{\text{op}} \text{ bounded}$$



$$\xrightarrow{\quad} \{(\mathbf{a}_i, \mathbf{y}_i)\}_{i=1}^n$$

$$\phi(\mathbf{X}^* \mathbf{a}_i; \boldsymbol{\varepsilon}_i)$$

Estimator (Student) model Given $\{(\mathbf{a}_i, \mathbf{y}_i)\}_{i=1}^n$, choose a function $g: \mathbb{R}^{\ell} \rightarrow \mathbb{R}^m$, **estimate** the matrix

$$\mathbf{X} \in \mathbb{R}^{\ell \times d} \text{ with } \ell = O_d(1)$$



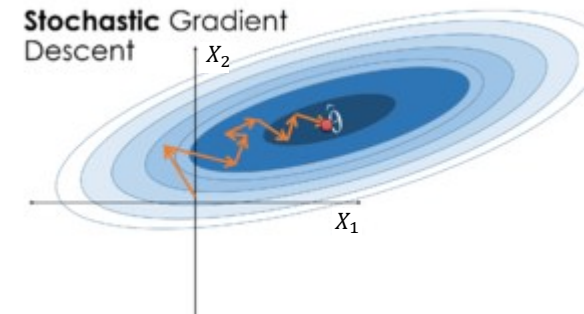
Optimization problem: $L(\mathbf{X}) = \mathbb{E}[\text{dist}(g(\mathbf{X}\mathbf{a}); \mathbf{y})] = \mathbb{E}_{(\mathbf{a}, \boldsymbol{\varepsilon})}[f(\mathbf{X}\mathbf{a}; \mathbf{X}^* \mathbf{a}, \boldsymbol{\varepsilon})]$

Stochastic Gradient Descent (SGD)

- One pass single batch (online learning) with fixed step size

$$X_{k+1} = X_k - \frac{\gamma_k}{d} \nabla_{X_k} f(X_k; \mathbf{a}_{k+1}, \mathbf{y}_{k+1})$$

- Initialization $X_0 \in \mathbb{R}^{\ell \times d}$ with a bounded norm
 - High-dimensional limit, n, d large, $\frac{n}{d} \rightarrow T \in (0, \infty)$
- Learning rate (step size)
- Samples (iteration) parameters



One class - No structure ($K = I_d$)

- Recall the features $\mathbf{a}_i \sim \mathcal{N}(0, K)$, $K = I_d$ - corresponds to isotropic data – **no structure**
- Iterates of SGD:

$$X_{k+1} = X_k - \frac{\gamma}{d} \nabla_r f \otimes \mathbf{a}_{k+1}, \text{ where } \nabla_r f = \nabla_{\mathbf{r}_k} f(\mathbf{r}_k; \mathbf{r}^*, \boldsymbol{\varepsilon}_{k+1})$$

$$\text{with } \mathbf{r}_k = X_k \mathbf{a}_{k+1}, \mathbf{r}^* = X^* \mathbf{a}_{k+1}$$

- Order Parameters: “**Norm**” $\langle X_k, X_k \rangle = X_k^\top X_k$, and “**Overlap**” $\langle X^*, X_k \rangle = (X^*)^\top X_k$

Theorem: (E. Collins-Woodfin, C&E. Paquette, **SI**): Fix $T = \frac{n}{d} \in [0, \infty)$ and for some $\varepsilon > 0$ with overwhelming probability,

$$\text{where } B_k = \begin{bmatrix} \langle X_k, X_k \rangle & \langle X_k, X^* \rangle \\ \langle X^*, X_k \rangle & \langle X^*, X^* \rangle \end{bmatrix}$$

$$\sup_{0 \leq t \leq T} \|\mathcal{B}(t) - B_{\lfloor td \rfloor}\| \leq d^{-\varepsilon}$$



Time scale: k iterates of SGD = td

continues time: $d \rightarrow \infty$ instead of $\gamma \rightarrow 0$

Limiting ODEs - No structure ($K = I_d$)

Given the deterministic B-matrix:

$$\mathcal{B}(t) = \begin{bmatrix} \mathcal{B}_{xx}(t) & \mathcal{B}_{x\star}(t) \\ \mathcal{B}_{x\star}(t) & \langle X^\star, X^\star \rangle \end{bmatrix}$$

Fisher matrix:

$$I(\mathcal{B}(t)) = \mathbb{E}[\nabla_{\mathbf{r}} f^{\otimes 2}]$$

Gradient of the loss:

$$H = \begin{bmatrix} \nabla_{\mathcal{B}_{xx}} \mathcal{L} & 0 \\ \nabla_{\mathcal{B}_{x\star}} \mathcal{L} & 0 \end{bmatrix}$$

The limiting ODEs:

$$\frac{d\mathcal{B}(t)}{dt} = \overset{\text{"Gradient" term}}{\boxed{-\gamma \left(\mathcal{B}(t) H(\mathcal{B}(t)) + H(\mathcal{B}(t))^\top \mathcal{B}(t) \right)}} + \overset{\text{"Noise" term}}{\boxed{\gamma^2 \begin{bmatrix} I(\mathcal{B}(t)) & 0 \\ 0 & 0 \end{bmatrix}}}$$

Order Parameters: “**Norm**” $\mathcal{B}_{xx,k} = \langle X_k, X_k \rangle$, and “**Overlap**” $\mathcal{B}_{x\star,k} = \langle X_k, X^\star \rangle$

Related literature - SGD in high dimension

Isotropic data ($K = I_d$):

- Two-layer neural net (Saad & Solla *Phys. Rev. E* '95, Riegler & Biehl *Physica A* '95...)
- Phase retrieval (Tan & Veryshynin *JMLR* '23, Mignacco et al. *NeurIPS* '20)
- Tensor PCA (Ben Arous et al. *NeurIPS* '22, Liang et al. *Inf. Inference* '23)
- Gaussian mixture models (Ben Arous et al. *ICLR* '24, 25')
- Generalized linear model (Gerbelot et al '22, Celentano et al. '21)
- Two-layer neural net (Goldt et al. *NeurIPS* '19)
-

Structured data (general K):

- Linear regression – Balasubramanian et al. '23, Wang et al. *J. Stat. Mech.* '19, Paquette et al. '22-25'
- Two-layer neural net – Yoshida et al. *NeurIPS* '19, Goldt et al. *PRX* '20

Structural data (general K)– Resolvent trick

Issue: One **cannot** write an autonomous set of equations,

$$\cancel{\frac{dB(t)}{dt} = F(B(t))}$$

Higher powers of K appears!

Solution: Random matrix theory trick!

Terms of the form
 $X^T K X, X^T K^2 X, \dots$

$$\text{Resolvent: } R(z; K) = (K - zI_d)^{-1} \text{ for } z \in \mathbb{C}$$

Some nice resolvent identities:

- $KR(z; K) = I_d + zR(z; K)$
- $R(z; K) = -\frac{1}{z} \left(I_d - \frac{K}{z} \right)^{-1} = \sum_{j=1}^{\infty} (K/z)^j$

This allows us to represent **any** polynomial of K !

$$p(K) = -\frac{1}{2\pi i} \oint_{\Gamma} p(z) R(z; K) dz$$

For any contour $\Gamma \subset \mathbb{C}$ enclosing the eigenvalues of K .

Structural data (general K)

Order Parameters:

“Resolvent Norm” $\langle X_k, X_k \rangle_R = X_k R(z; K) X_k^\top$

“Resolvent Overlap” $\langle X_k, X^* \rangle_R = X_k R(z; K) (X^*)^\top$

- Define the S matrix of “Order Parameters” : $S_k(z) = \begin{bmatrix} \langle X_k, X_k \rangle_R & \langle X_k, X^* \rangle_R \\ \langle X^*, X_k \rangle_R & \langle X^*, X^* \rangle_R \end{bmatrix} = \begin{bmatrix} X_k \\ X^* \end{bmatrix} R(z; K) \begin{bmatrix} X_k^\top & (X^*)^\top \end{bmatrix}$

Theorem (E. Collins-Woodfin, C&E. Paquette, **SI**): For any $T = \frac{n}{d} \in [0, \infty)$ and for some $\varepsilon > 0$ with

overwhelming probability

$$\sup_{0 \leq t \leq T} \|\mathcal{S}(t, z) - S_{[td]}(z)\| \leq d^{-\varepsilon}.$$

This then allows us to derive a limiting ODEs:

$$\frac{d\mathcal{S}(t, z)}{dt} = \mathcal{F}(z, \mathcal{S}(t, z))$$

Explicit risk curves

A large class of functions (statistics):

$$\varphi(X) = h(\langle [X, X^*]^{\otimes 2}, p(K) \rangle) \rightarrow h(-\frac{1}{2\pi i} \oint_{\Gamma} p(z) \mathcal{S}(t, z) dz)$$

h is α pseudo-Lipchitz function, and p is a polynomial

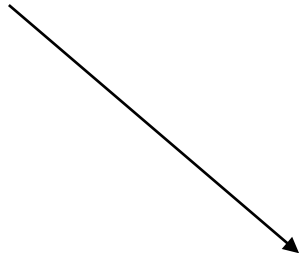
$$\varphi(X_k) = \mathcal{L}(X_k) = \mathbb{E}_{(\mathbf{a}, \varepsilon)} f(X\mathbf{a}, X^*\mathbf{a}; \varepsilon)$$

By our theorem



$$\mathcal{L}(X_k) = h(\langle [X, X^*]^{\otimes 2}, K \rangle) \rightarrow \mathcal{L}(t)$$

- Other functions: $\|X\|^2, \|X - X^*\|^2 \dots$


$$\begin{aligned} X_k \mathbf{a} &\sim \mathcal{N}(0, \langle X_k^{\otimes 2}, K \rangle), \\ X^* \mathbf{a} &\sim \mathcal{N}(0, \langle (X^*)^{\otimes 2}, K \rangle), \\ \mathbb{E}[\langle X^* \mathbf{a}, X_k \mathbf{a} \rangle] &= \langle X^* \otimes X_k, K \rangle \end{aligned}$$

Main result: Limiting process - Homogenized SGD

Homogenized SGD: The process \mathcal{X}_t satisfies the following SDE: “Noise” term

$$d\mathcal{X}_t = -\gamma \nabla \mathcal{L}(\mathcal{X}_t) dt + \boxed{\frac{\gamma}{\sqrt{d}} \sqrt{\mathbb{E}[(\nabla_r f)^{\otimes 2}] \otimes K} d\mathcal{W}_t}$$

Time scale:
 k iterates of SGD = td ,

where \mathcal{W}_t denotes d -dimensional Brownian motion
“Fisher matrix”

Theorem (E. Collins-Woodfin, C&E. Paquette, **SI**):

Fix $T = \frac{n}{d} \in [0, \infty)$, the process \mathcal{X}_t for $t \in [0, T]$ and some $\varepsilon > 0$ with overwhelming probability

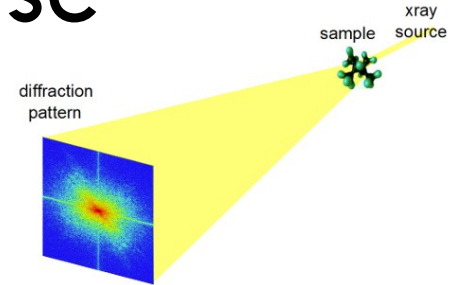
$$\sup_{0 \leq t \leq T} |\varphi(X_{\lfloor td \rfloor}) - \varphi(\mathcal{X}_t)| \leq d^{-\varepsilon}$$

Note $X_{\lfloor td \rfloor} \neq \mathcal{X}_t$

Recall SGD iterates: $X_{k+1} = X_k - \frac{\gamma}{d} \nabla_r f \otimes \mathbf{a}_{k+1}^\top$, with the population loss $\mathcal{L}(\mathcal{X}) = \mathbb{E}[f]$

Example 1: Phase retrieval – Hard phase

Candes et al., '11



- **Task:** Recover $X^* \in \mathbb{R}^{1 \times d}$, from modulo of projections on the vectors \mathbf{a} :

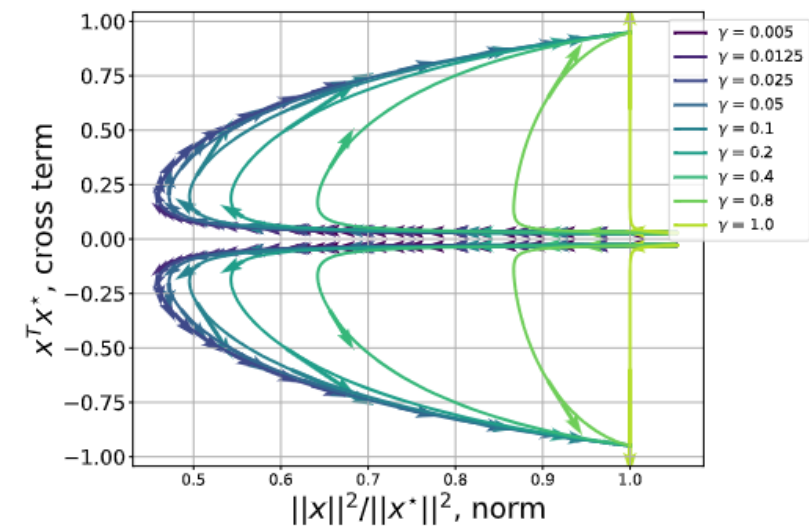
$$\mathcal{L}(X) = \mathbb{E}_{\mathbf{a}}[(|X\mathbf{a}| - |X^*\mathbf{a}|)^2]$$

Student: $g(X\mathbf{a}) = |X\mathbf{a}|$, and teacher: $\phi(X^*\mathbf{a}) = |X^*\mathbf{a}|$

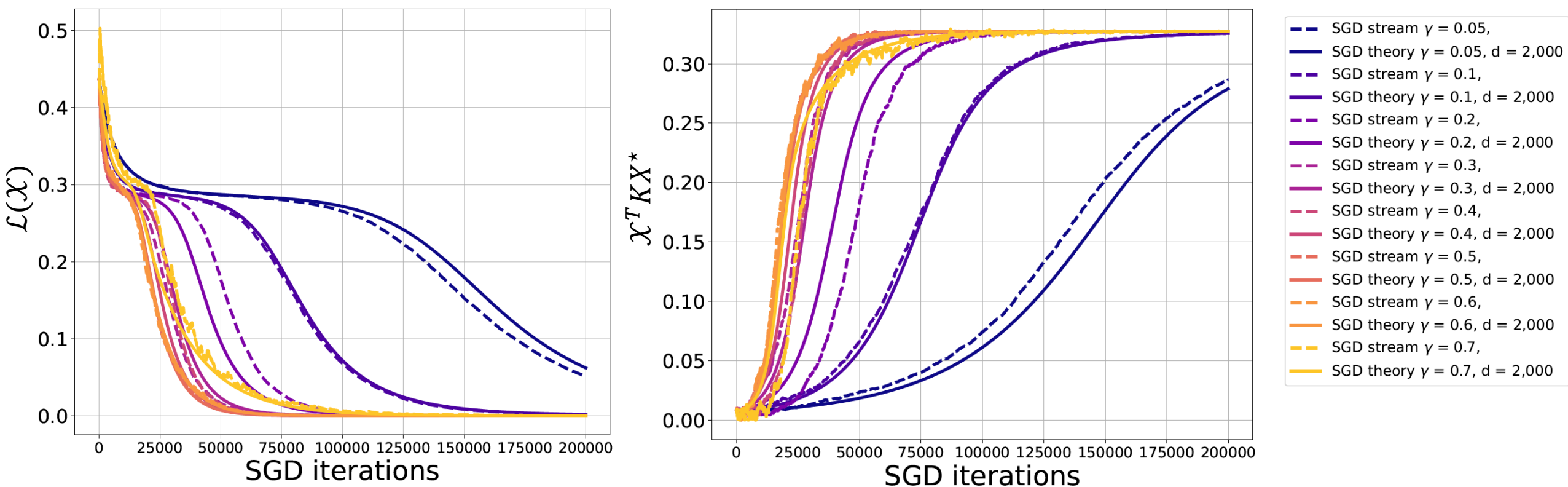
- Random initialization is problematic, suppose $X_0 \sim \mathcal{N}(0, I_d)$

$$\text{Initial overlap} = \langle X_0, X^* \rangle \sim \frac{1}{\sqrt{d}}$$

- If Initial overlap $\sim \frac{1}{\sqrt{d}}$ SGD converges in $n = O(d \log d)$
- If Initial overlap $\sim O(1)$ SGD converges in $n = O(d)$
- This can be seen directly from our equation of the norm and overlap



Example: Phase retrieval – risk and alignment



What learning rate ensures descent?

- **Distance to optimality**, by our theorem $\|X_{[td]} - X^*\|^2 \rightarrow \mathcal{D}(t)^2$:

$$\frac{d\mathcal{D}(t)^2}{dt} = -2\gamma A(t) + \frac{\gamma^2}{d} \text{Tr}(K) I(t)$$

where $A(t), I(t)$ are functions of the limiting norm and overlap

Thus, $\mathcal{D}(t)^2$ is decreasing when: $\gamma \leq \gamma_t^{\text{stable}} = \frac{2}{\frac{1}{d}\text{Tr}(K)} \frac{A(t)}{I(t)}$

- If for some $m > 0$, $mI(t) \leq A(t)$ (convexity and smoothness assumption):

$$\gamma \leq \frac{2m}{\frac{1}{d}\text{Tr}(K)}$$

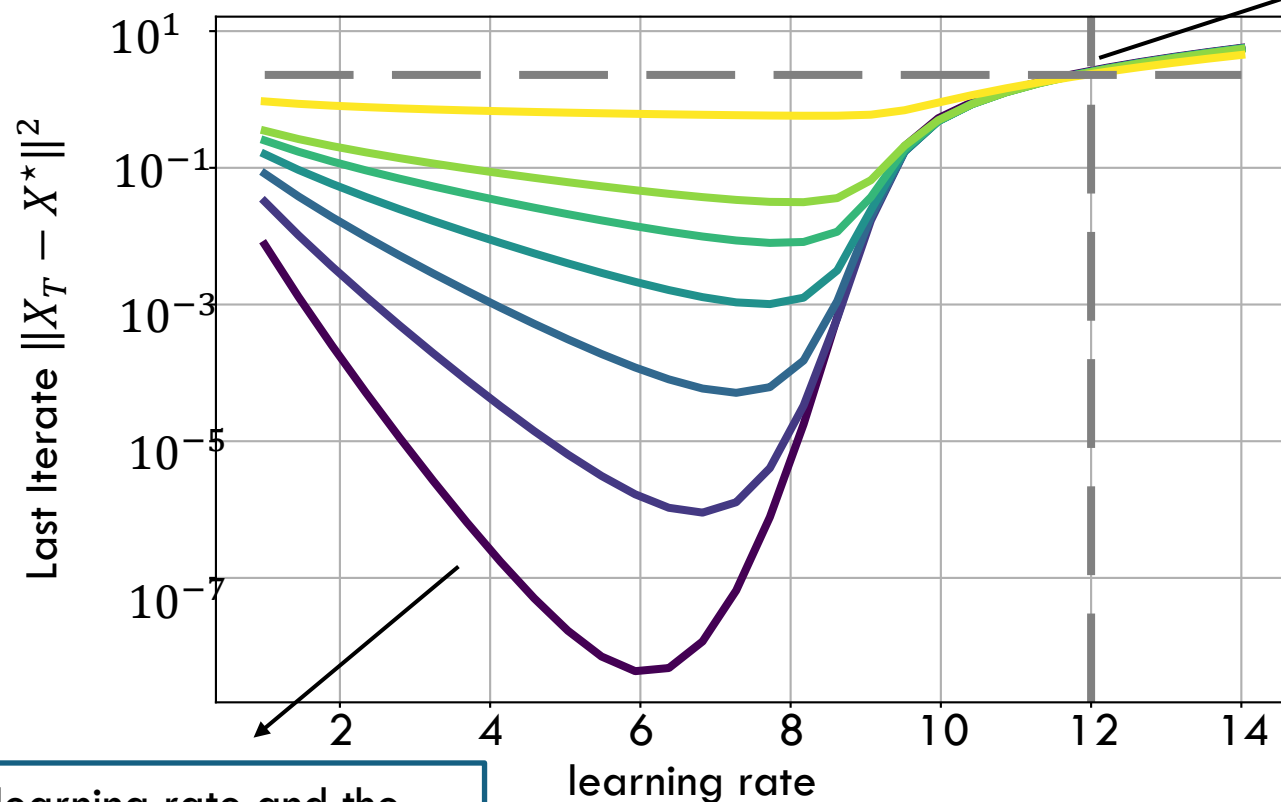
**Average
eigenvalue** rather
than largest!

- E..g. if ∇f is Lipschitz with constant L then $m = 1/2L$
- Convergence rate will now depend on $\frac{\lambda_{\min}(K)}{\frac{1}{d}\text{Tr}(K)}$

Dynamical threshold
Motivate ideas such as line
search, and Polyak step size

Descent and critical learning rate

$\mathcal{D}(t)^2$ in binary noiseless logistic regression for various γ and K



Threshold learning rate depends only on $\frac{1}{d} \text{Tr}(K)$

Average eigenvalue is fixed, $\frac{1}{d} \text{Tr}(K) = 1$

- Marchenko-Pastur, max eig. = 2.23
- $q = 2.0$, max eig. = 2.58
- $q = 2.5$, max eig. = 3.05
- $q = 3.0$, max eig. = 3.52
- $q = 3.5$, max eig. = 4.01
- $q = 4.0$, max eig. = 4.50
- $q = 8.0$, max eig. = 8.47
- Initialization, $\|X_0 - X^*\|^2$

$K = \text{diag}(\sigma_i^{2q}; 1, \dots, d)$, with $\sigma_i \sim \text{Unif}(1, 2)$

Optimal learning rate and the rate of convergence do vary as the max/min eigenvalue changes

Example 2: Stochastic adaptive methods - AdaGrad Norm

- Algorithm setup $X, X^* \in \mathbb{R}^d$, with $\gamma_0 = \frac{\gamma}{b_0} > 0$:

$$X_{k+1} = X_k - \frac{\gamma_k}{d} \nabla_{X_k} f(X_k; \mathbf{a}_{k+1}, \mathbf{y}_{k+1})$$
$$\gamma_k = \frac{\gamma}{\sqrt{b_0^2 + \sum_{j=1}^k \|\nabla_{X_k} f\|^2}}$$

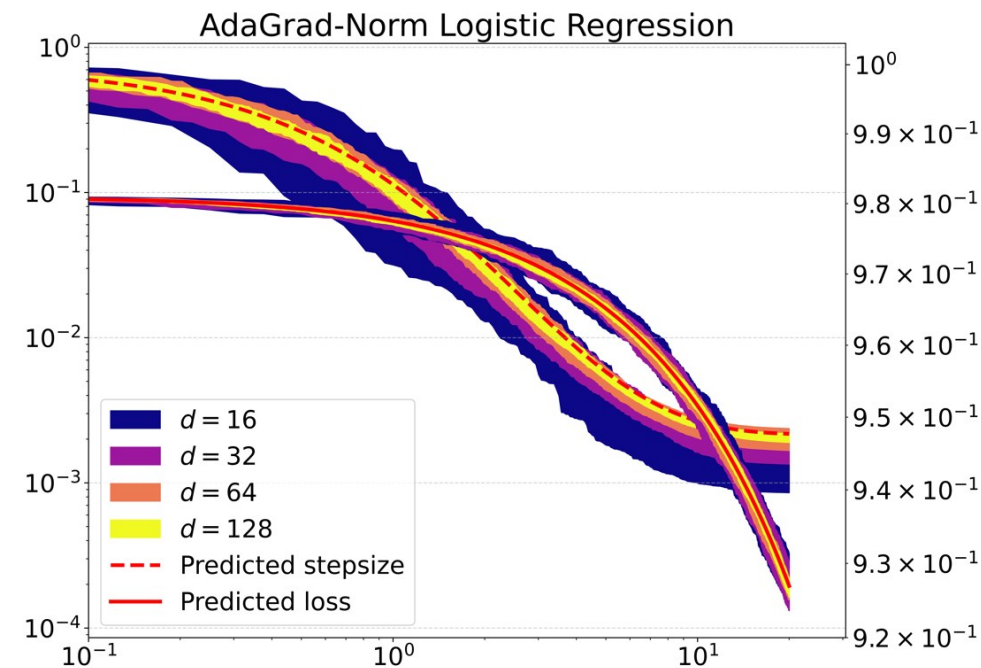
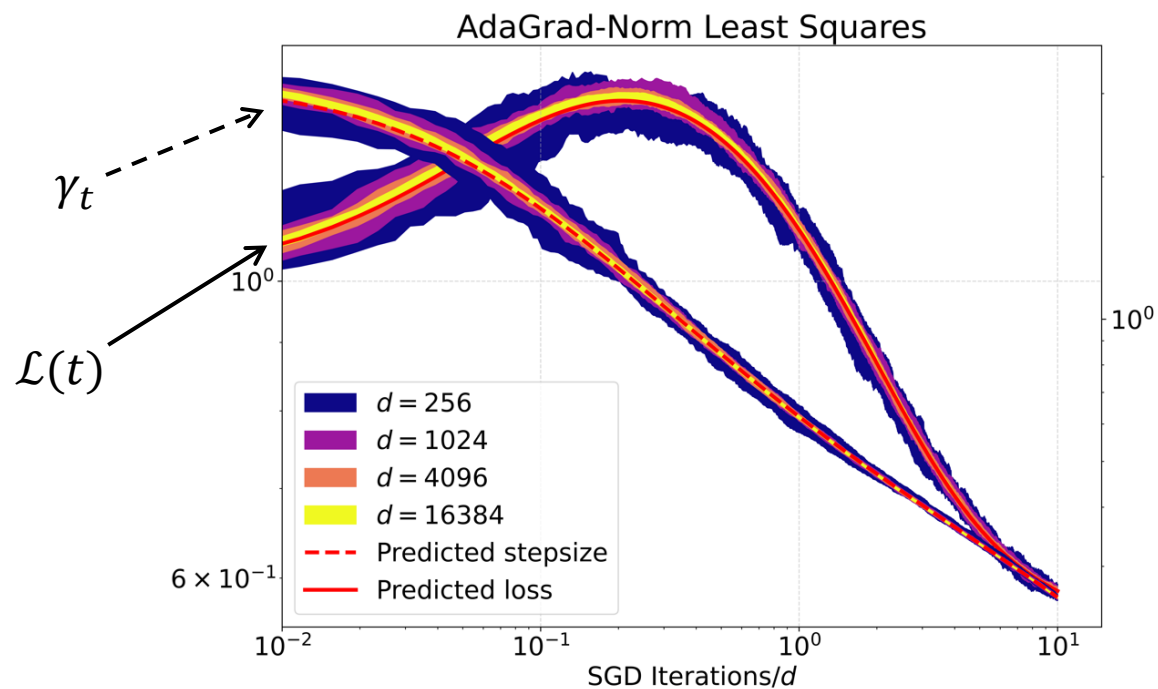
- Deterministic equivalence $\gamma_{\lfloor td \rfloor} \rightarrow \gamma_t$:

$$\gamma_t = \frac{\gamma}{\sqrt{b_0^2 + \frac{\text{Tr}(K)}{d} \int_0^t I(s) ds}}$$

with $I(s) = \mathbb{E}[f'(X^T a)^2]$

- For Least square: $I(s) = 2\mathcal{L}(s)$.

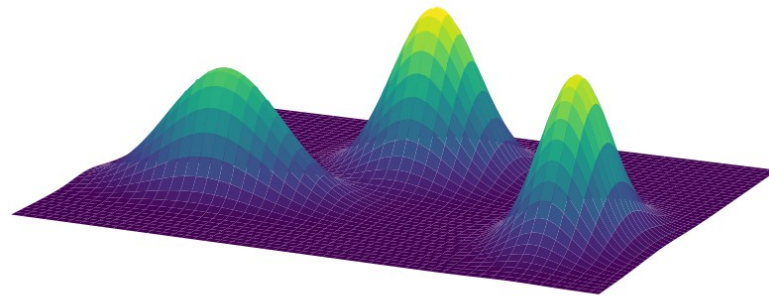
Stochastic adaptive methods – Exact Adagrad dynamics



How can we extend this to the multiclass setting?

Gaussian Mixture Model (GMM)

Joint work (in progress) with Elizabeth Colins-Woodfin



Model-setup - Data distribution



$$\rightarrow \{(\mathbf{a}_{i,I}, \mathbf{y}_{i,I})\}_{i=1}^n$$

Target (teacher) model: Data from ℓ^* classes:

$$\mathbf{a}_{i,I} \sim \mathcal{N}(\boldsymbol{\mu}_I, \mathbf{K}_I), \text{ with } \mathbf{K}_I \in \mathbb{R}^{d \times d}, \|\mathbf{K}_I\|_{\text{op}} \text{ bounded, all } \{\mathbf{K}_I\} \text{ commute}$$

$$c \in [\ell^*] = O(\log(d)), \text{ and } p_c = \mathbb{P}(c = I)$$

Our setup allow for the following two settings:

“Hard label” - $\mathbf{y}_{i,I} = \mathbf{I}$ or one - hot encoding of I

$$\text{“Soft label”} : \mathbf{y}_{i,I} = \phi_I(X^* \mathbf{a}_{i,I}; \boldsymbol{\varepsilon}_i) \text{ e.g. } y_i = \phi_I(X^* \mathbf{a}_{i,I}, \boldsymbol{\varepsilon}) = \text{softmax}(X^* \mathbf{a}_i) \quad \text{softmax}(\mathbf{r})_i = \frac{e^{r_i}}{\sum_j e^{r_j}}$$

Classifier and optimization problem GMM



Classifier (student) model:

Choose a function $g: \mathbb{R}^\ell \rightarrow \mathbb{R}^m$, **estimate using online SGD** the matrix $X \in \mathbb{R}^{\ell \times d}$

$$\ell = O(\log(d))$$

Optimization problem:

$$\min_{X \in \mathbb{R}^{d \times \ell}} L(X) = \mathbb{E}_{(\mathbf{a}, I)} [f_I(X\mathbf{a}_I; \mathbf{y}_I)]$$

Related work: Seddik et al ICLR 2020, Loureiro et al NIPS 2021, Mai & Liao 2019, Ben-Arous et al 2025

Main result - deterministic equivalence

Theorem (Collins-Woodfin, **SI** '25) Fix $T = \frac{n}{d} > 0$. For any $\epsilon \in (0, \frac{1}{2})$, with overwhelming probability,

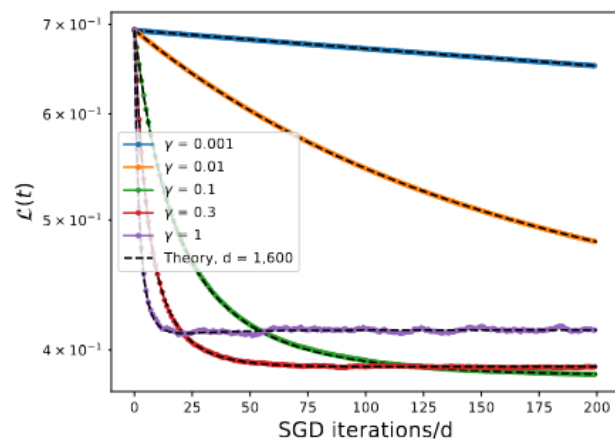
$$\sup_{0 \leq t \leq T} |L(X_{[td]}) - \mathcal{L}(t)| < Cd^{-\epsilon}$$

where $\mathcal{L}(t)$ is the “deterministic equivalent” of the risk, expressible in terms of a system of autonomous ODEs.

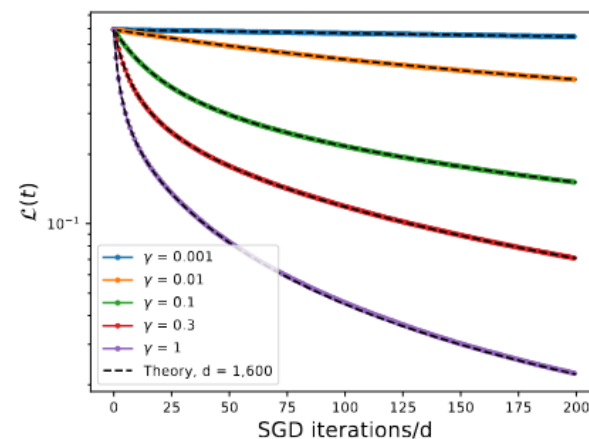
- This holds for other statistics of X , not just risk.

How does the structure of the classes affect
the SGD dynamics?

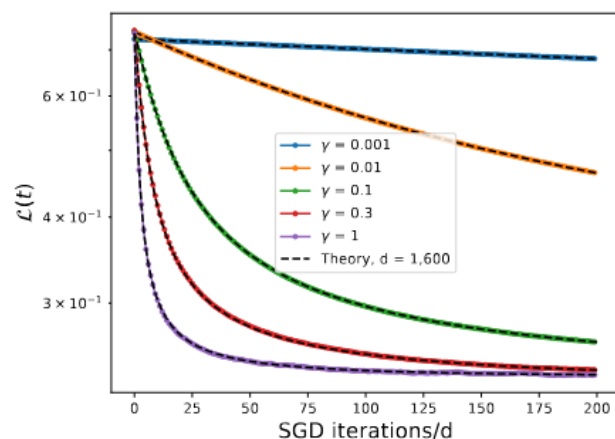
SGD vs theory for different data models



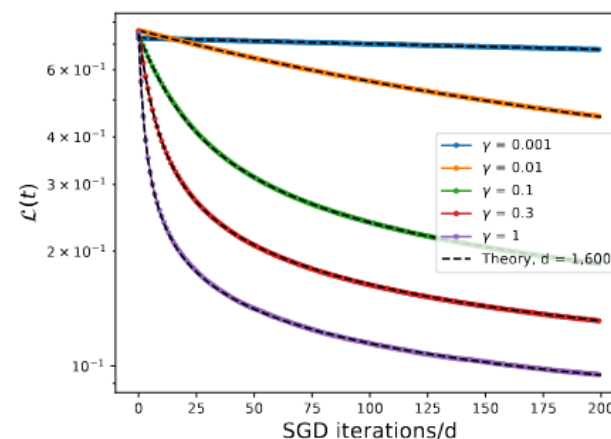
(a) Identity covariance



(b) Zero-one model

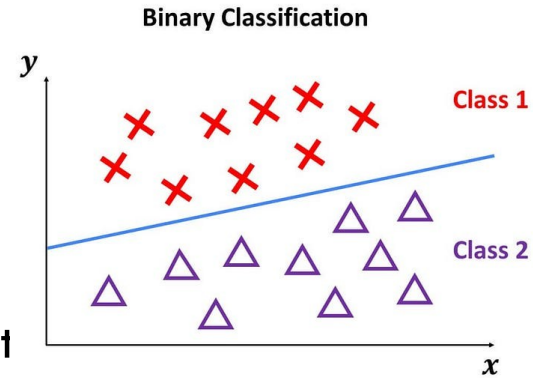


(c) Mild Power-law $\beta = 1$



(d) Extreme power-law $\beta = 0.2$

Binary logistic regression



Two classes: $y_i = 1_{i=1}$ and $a \mid i = 1 \sim N(\mu_1, K_1)$, and $a \mid i = 2 \sim N(\mu_2, K_2)$, with

$$L(X) = \mathbb{E}_{(a,y)} \left[-a^\top X y + \log(e^{a^\top X} + 1) \right]$$

For simplicity: $\mu_1 = -\mu_2 = \mu$ $K_1 = \text{diag}(\lambda_1^{(1)}, \dots, \lambda_d^{(1)})$ $K_2 = \text{diag}(\lambda_1^{(2)}, \dots, \lambda_d^{(2)})$

Identity model: $K_1 = K_2 = I_d$ ($|I_{11}| = d$)

Zero-One model - All eigenvalues in $\{0, 1\}$

• Partition indices $\{1, \dots, d\} = I_{00} \sqcup I_{01} \sqcup I_{10} \sqcup I_{11}$

• where $I_{jk} := \{i \leq d \mid \lambda_i^{(1)} = j, \lambda_i^{(2)} = k\}$

Example $d = 4$:

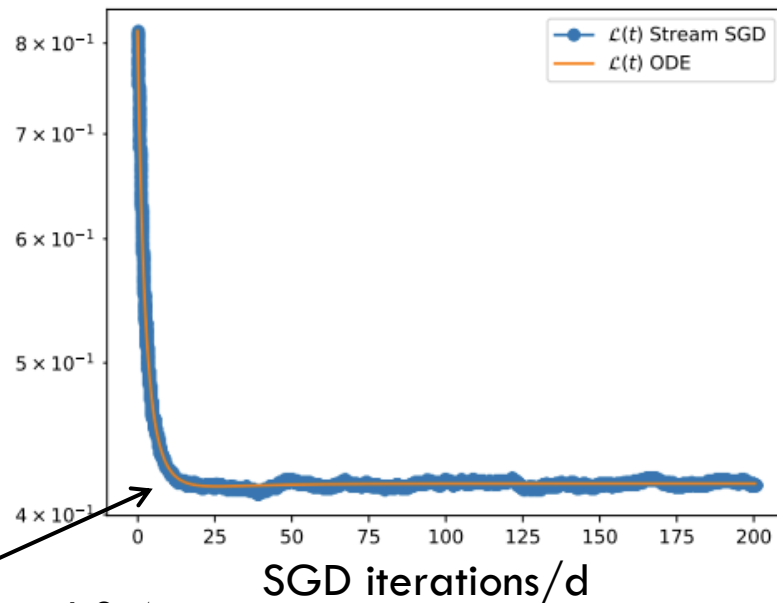
$$K_1 = \begin{bmatrix} 0 & & & \\ & 0 & & \\ & & 1 & \\ & & & 1 \end{bmatrix} \quad K_2 = \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 0 & \\ & & & 1 \end{bmatrix}$$

Comparing Identity and Zero-One models - Risk

Does SGD find the “perfect” subspace?

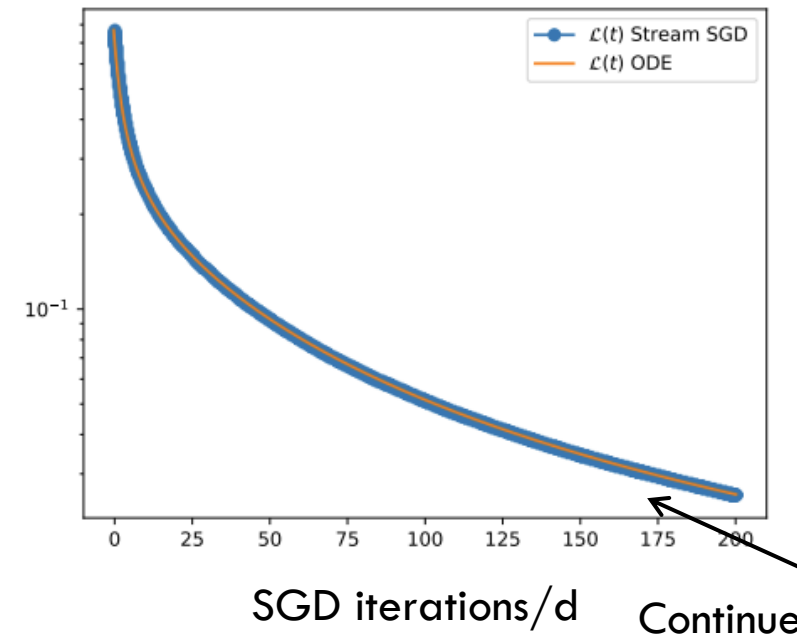
How “clean” directions (I_{00}) affect the classification?

Learning curve - identity



flattens out at around 0.4

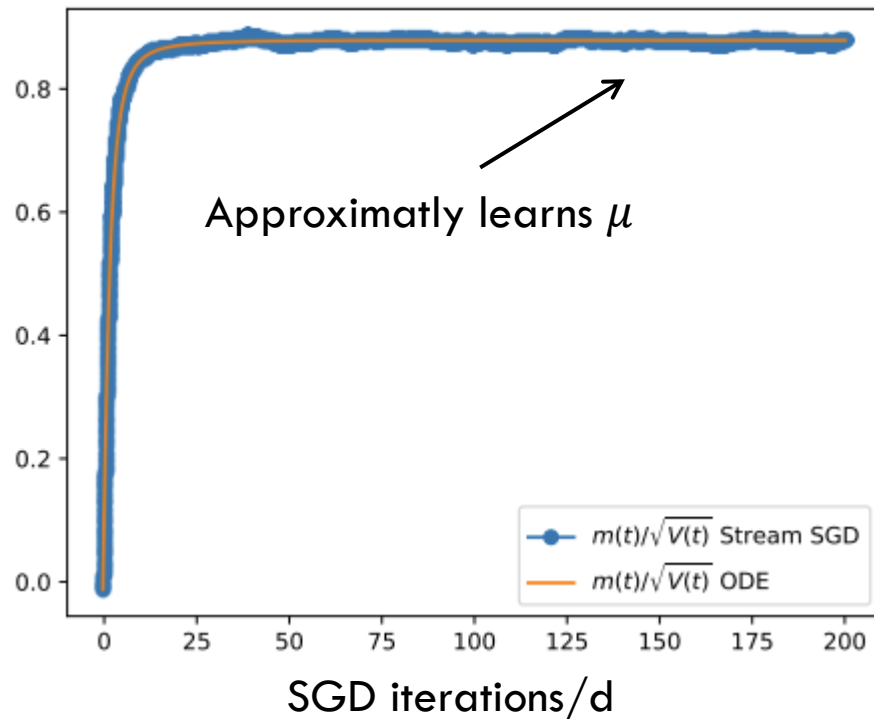
Learning curve - Zero-one



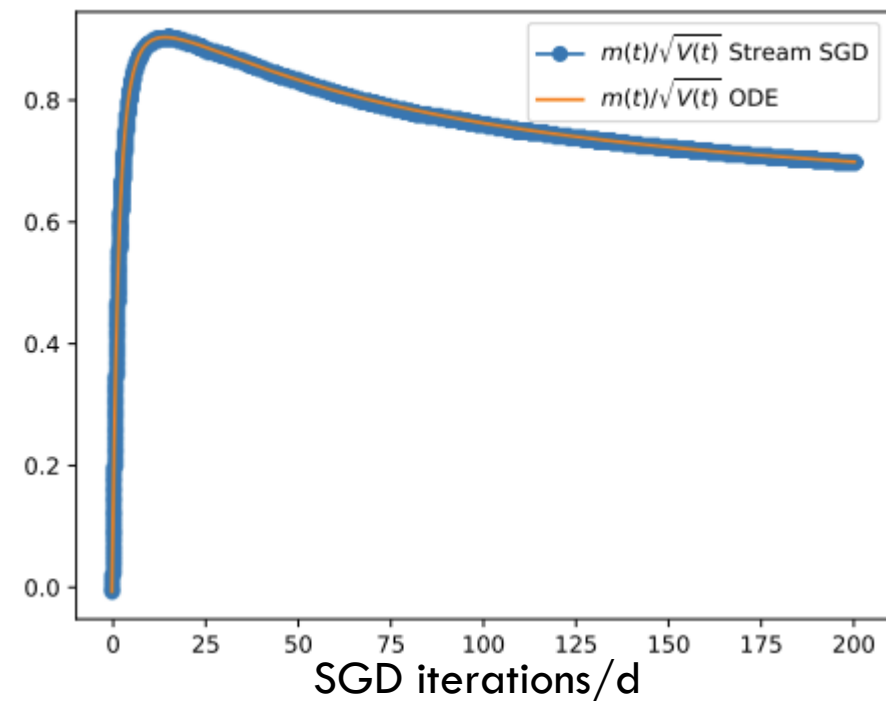
Comparing Identity and Zero-One models – Alignment

$$\text{Alignment} := \frac{\mu^\top X_{[td]}}{\|X_{[td]}\|} \rightarrow \frac{m(t)}{\sqrt{\mathcal{V}(t)}}$$

Alignment with μ - Identity



Alignment with μ - Zero-One



Zero-One asymptotic

Proposition (Collins-Woodfin, **SI** '25) : For $\gamma < 1, p_1 = \frac{1}{2}, |I_{00}| = |I_{01}| = |I_{10}| = |I_{11}| = d/4$. There exist $C_1(\gamma), C_2(\gamma)$ such that

$$t^{-C_1(\gamma)} \leq \mathcal{L}(t) \leq t^{-C_2(\gamma)}$$

where the alignment with μ :

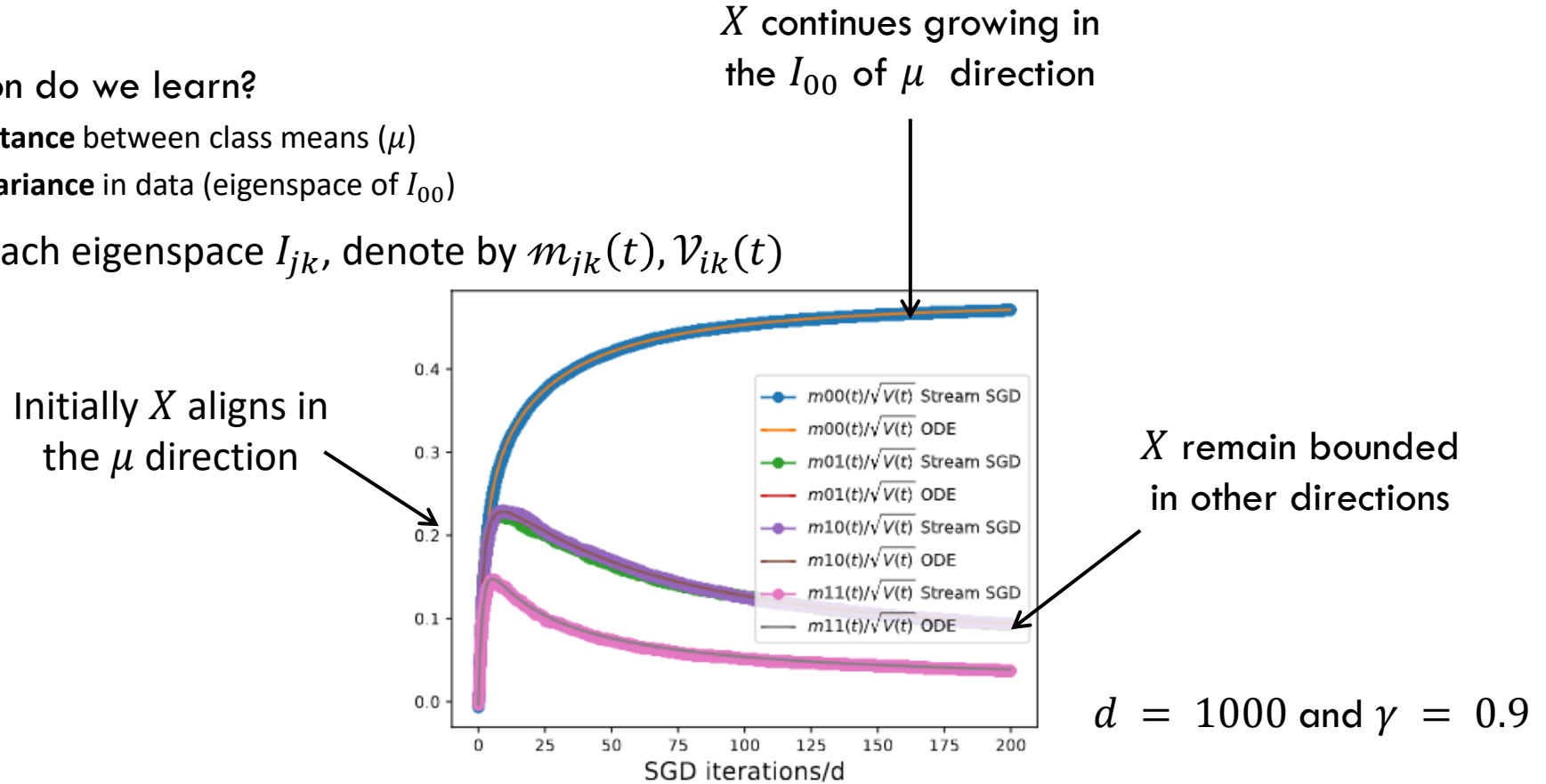
$$m(t) \asymp \log t, \quad \frac{m(t)}{\sqrt{\mathcal{V}(t)}} = \frac{1}{2} (1 + O((\log t)^{-1/2}))$$

Remarks:

- This implies analogous bounds on the **original loss** in high dimension by our Theorem
- $\mathcal{V}(t) \approx \|X_{[kd]}\|^2$ and the $m(t) \approx \mu^\top X_{[kd]}$ grows logarithmically with n ! (very different than the identity setting!)
- The covariance matrices has no power law structure.

Perfect classification vs clean directions

- What direction do we learn?
 - **Largest distance** between class means (μ)
 - **Smallest variance** in data (eigenspace of I_{00})
- Project into each eigenspace I_{jk} , denote by $m_{jk}(t), v_{jk}(t)$



In particular, we can prove that:

$$m_{00}(t) \asymp \log t,$$

$$m_{01}(t), m_{10}(t), m_{11}(t) = O(\sqrt{\log t}),$$

All in one - Power law covariance and mean

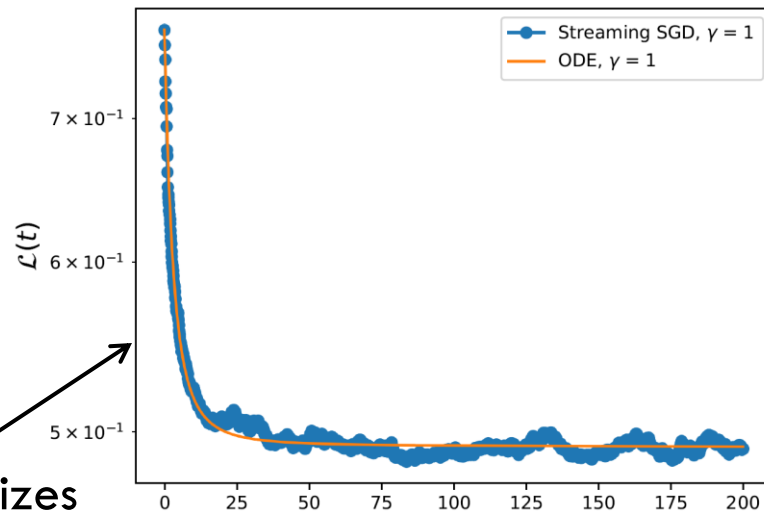
Power law model: $K_1 = K_2 = \text{diag}(\lambda_1, \dots, \lambda_d)$, $\mu_1 = -\mu_2 = \mu$

$$\lambda_i = \left(\frac{i}{d}\right)^\alpha \quad \text{and} \quad \mu_i^2 = \frac{1}{d} \left(\frac{i}{d}\right)^\beta, \quad \text{for } \beta \geq 0, \alpha > 1$$

Spectrum with
eigenvalues
accumulating near
zero

- There is a **phase transition** at $\alpha = 1 + \beta$

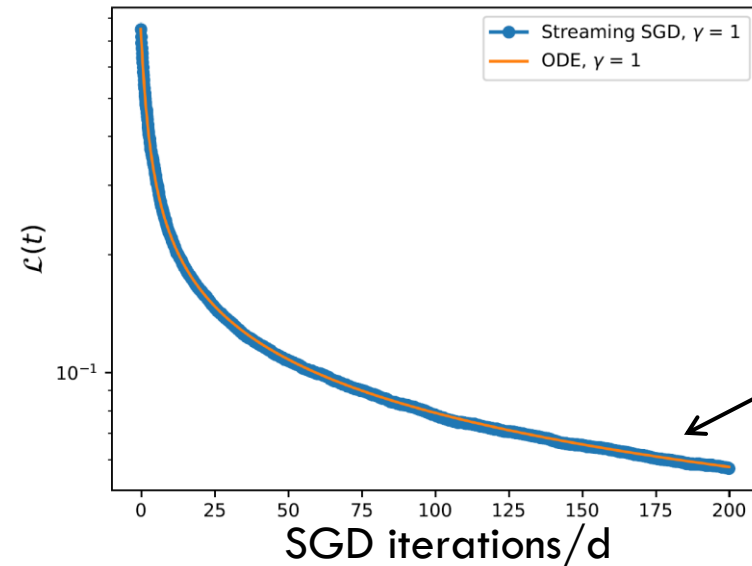
Mild power law



Risk stabilizes
as in the identity case

$$1 + \beta = 2.2 > \alpha = 1.2$$

Extreme power law



Risk continues
to go down

$$1 + \beta = 1 < \alpha = 1.2$$

Mild power law and identity regime

Proposition: Suppose $X_0 = 0$, $\gamma < 1$, $p_1 = \frac{1}{2}$. Then for $t \geq 1$,

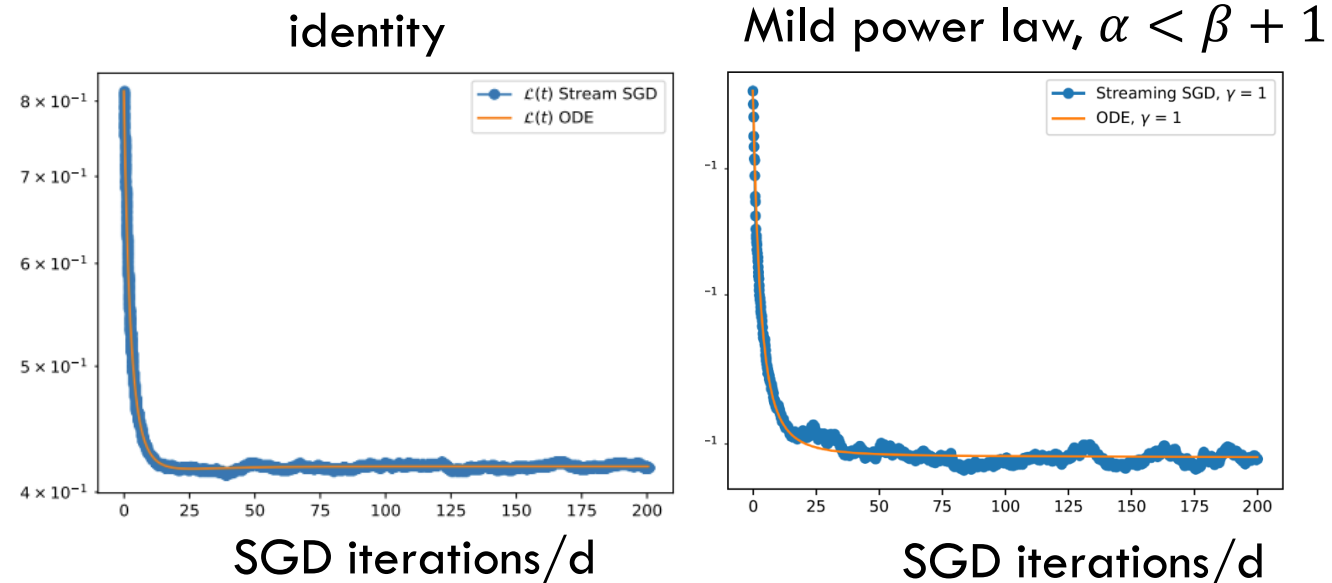
Mild power law ($\alpha < \beta + 1$)

- $m(t) \asymp \mu^\top [K]^{-1} \mu$
- $\mathcal{L}(t) \asymp \mathcal{L}_{\min} > 0$

Identity covariance: $K = I_d$ with $\|\mu\| = O(1)$

- $m(t) \asymp \mu^\top [K]^{-1} \mu$
- $\mathcal{L}(t) \asymp \mathcal{L}_{\min} > 0$

Rate of convergence are different!



Extreme power law ($\alpha \geq \beta + 1$)

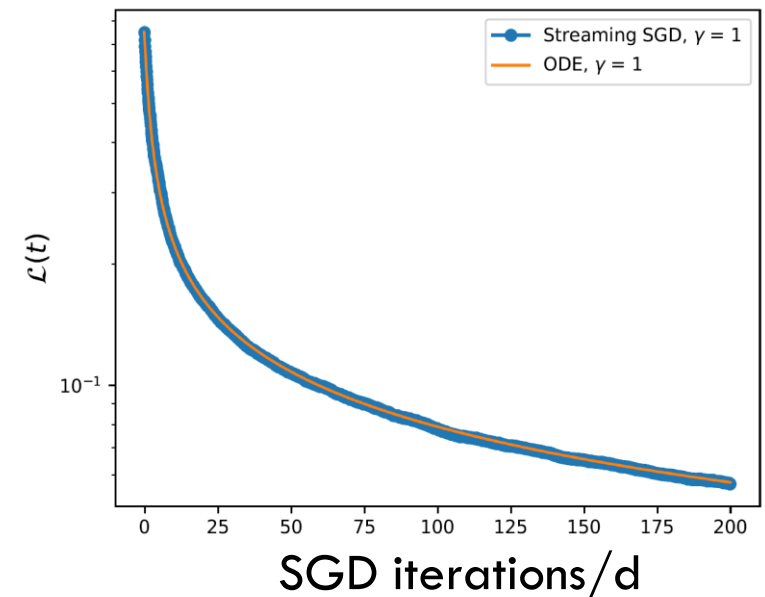
Proposition: Suppose $X_0 = 0$, $\gamma < 1$, $p_1 = \frac{1}{2}$. Then for $t \geq 1$,

- $m(t)$ **grows** with t at a polylog rate
- $\mathcal{L}(t) \rightarrow 0$ faster than polynomial decay,
but still slower than exponential decay.

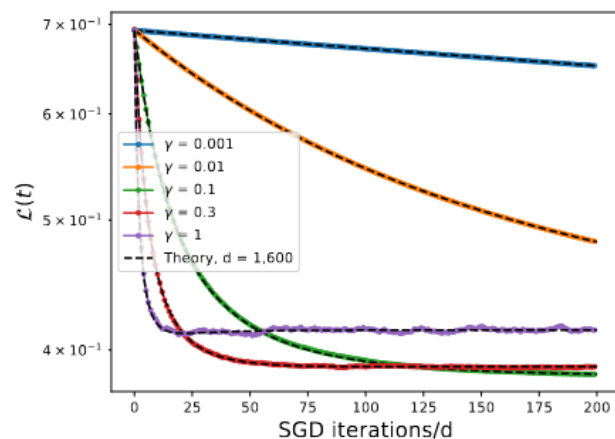
Remarks:

- Closely related to the Zero-One thought rates are different!
- **Small variance** directions contribute the **most to the learning**.

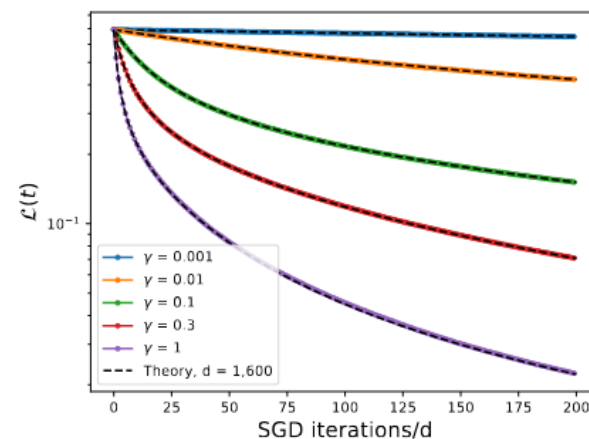
Extreme power law, $\alpha > \beta + 1$



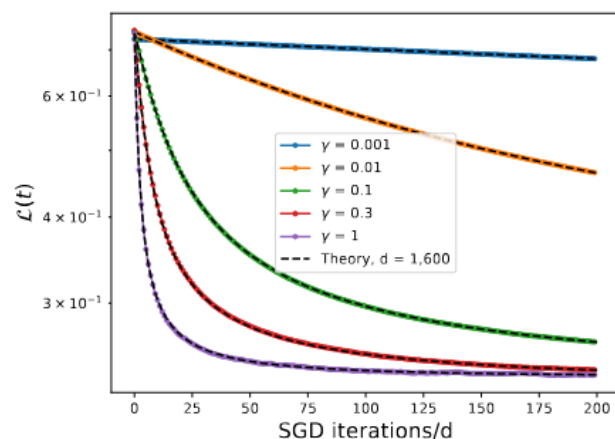
SGD vs theory for different data models



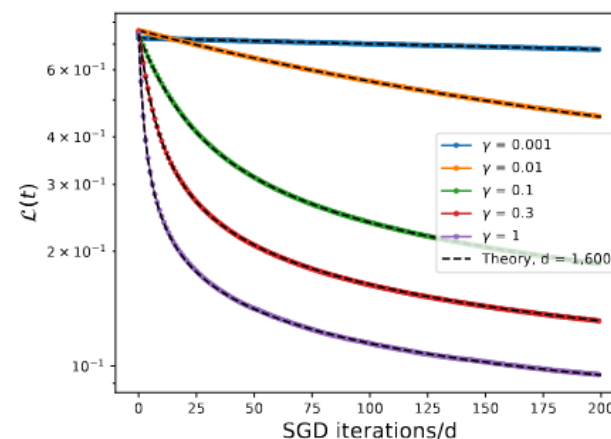
(a) Identity covariance



(b) Zero-one model



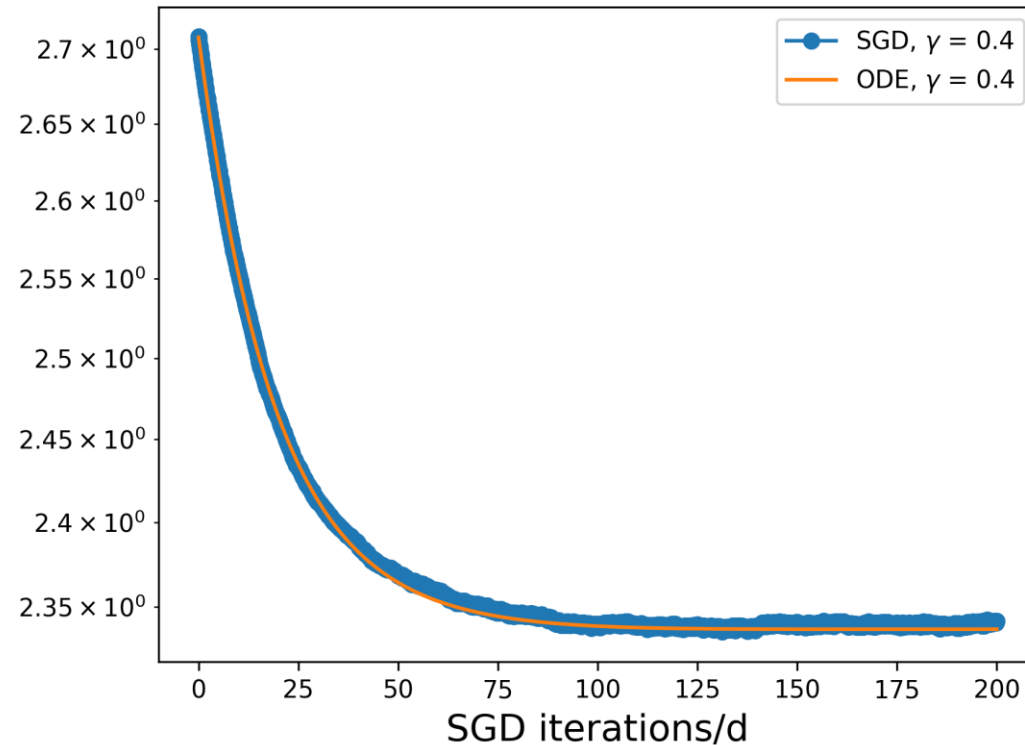
(c) Mild Power-law $\beta = 1$



(d) Extreme power-law $\beta = 0.2$

Large number of classes

- We allow the number of classes to grow as $\ell = \ell^* = O(\log(d))$

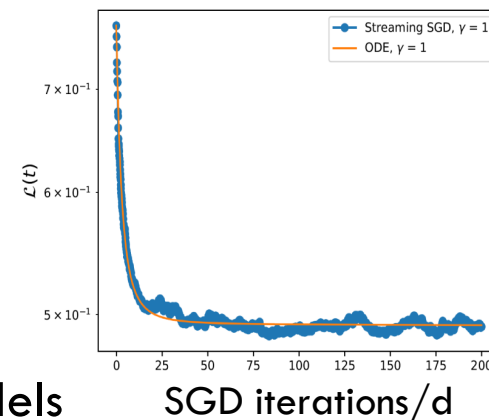


$$\ell = \ell^* = 15, d = 500$$

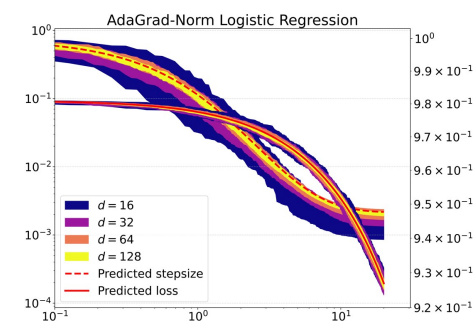
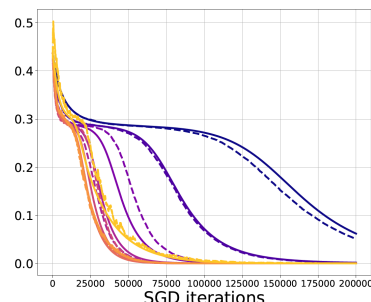
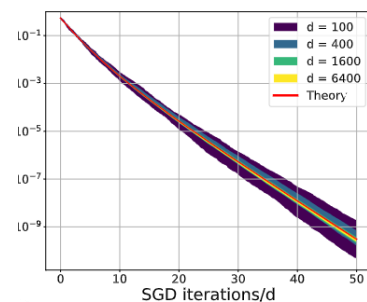
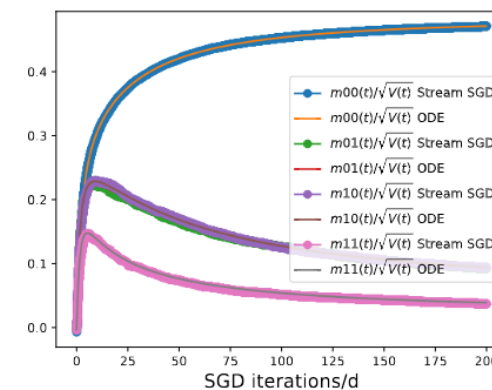
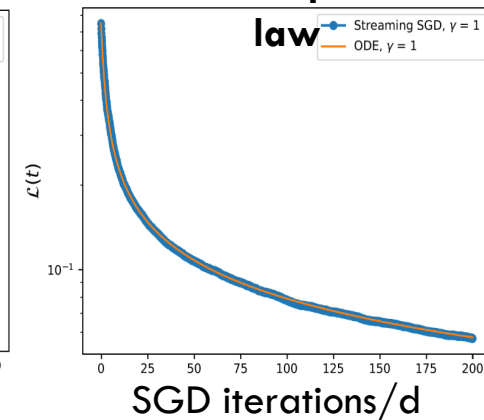
Conclusions

- An **exact** asymptotic theory of online SGD for multi-index models
- Applies to **stochastic adaptive methods**, such as Adagrad
- Extension of the theory for **nonisotropic** Gaussian mixture models
- Algorithm-dependent scaling laws and **phase transition** as a function of the structure
- Asymptotic analysis show the exact scaling behavior of the loss and other statistics
- Allow for growing number of classes $\ell^* = O(\log d)$

Mild power law



Extreme power law



Thank You!

Questions?

- Collins-Woodfin, E & **Seroussi, I** "SGD dynamics for Gaussian Mixtures models with non-isotropic Covariance and mean" (To appear soon!)
- Collins-Woodfin, E., Paquette, C., Paquette, E., & **Seroussi, I.** (2024). Hitting the high-dimensional notes: An ode for SGD learning dynamics on GLMs and multi-index models. *Information and Inference: A Journal of the IMA*, 13(4), iaae028.
- Collins-Woodfin, E., **Seroussi, I.**, Malaxechebarría, B.G., Mackenzie, A.W., Paquette, E. and Paquette, C., 2024. The High Line: Exact Risk and Learning Rate Curves of Stochastic Adaptive Learning Rate Algorithms. arXiv preprint arXiv:2405.19585. *NeurIPS 2024*