
Implicit Regularization or Implicit Conditioning? Exact Risk Trajectories of SGD in High Dimensions

Courtney Paquette

Google Research, Brain and McGill University
courtney.paquette@mcgill.ca

Elliot Paquette

McGill University
elliott.paquette@mcgill.ca

Ben Adlam

Google Research, Brain

Jeffrey Pennington

Google Research, Brain

Abstract

Stochastic gradient descent (SGD) is a pillar of modern machine learning, serving as the go-to optimization algorithm for a diverse array of problems. While the empirical success of SGD is often attributed to its computational efficiency and favorable generalization behavior, neither effect is well understood and disentangling them remains an open problem. Even in the simple setting of convex quadratic problems, worst-case analyses give an asymptotic convergence rate for SGD that is no better than full-batch gradient descent (GD), and the purported implicit regularization effects of SGD lack a precise explanation. In this work, we study the dynamics of multi-pass SGD on high-dimensional convex quadratics and establish an asymptotic equivalence to a stochastic differential equation, which we call homogenized stochastic gradient descent (HSGD), whose solutions we characterize explicitly in terms of a Volterra integral equation. These results yield precise formulas for the learning and risk trajectories, which reveal a mechanism of implicit conditioning that explains the efficiency of SGD relative to GD. We also prove that the noise from SGD negatively impacts generalization performance, ruling out the possibility of any type of implicit regularization in this context. Finally, we show how to adapt the HSGD formalism to include streaming SGD, which allows us to produce an exact prediction for the excess risk of multi-pass SGD relative to that of streaming SGD (bootstrap risk).

1 Introduction

Stochastic gradient descent (SGD) is the algorithm of choice for optimization in modern machine learning and has been hailed as a major reason for deep learning's success [11, 21]. Explanations for the effectiveness of SGD typically refer to its computational efficiency and to its favorable generalization properties, but theoretical understanding of these purported benefits is far from complete.

The efficiency of SGD has been the subject of extensive research, dating back to the original work of Robbins and Monro [68] and extending to modern large-scale machine learning applications (see e.g. [10, 12]). However, despite its widespread adoption and algorithmic simplicity, surprisingly little is known about how SGD performs in the types of high-dimensional optimization problems that occur in practice. Part of the challenge in deriving robust high-level conclusions about the efficiency of SGD is simply that those conclusions can depend on precisely which quantities are measured and what assumptions are leveraged. For example, in the extreme setting where the samples are one-hot vectors, running SGD on a quadratic function is actually identical to running full-batch gradient descent; as such, any statements about the two algorithms' relative efficiency must be data-dependent.

36th Conference on Neural Information Processing Systems (NeurIPS 2022).

Furthermore, the majority of prior analyses focus on the streaming or single-pass setting, where each sample is seen a single time. While this setting is appropriate when the number of samples n is much larger than the dimensionality d , it does not adequately describe the practically-relevant overparameterized or high-dimensional settings where $d \gtrsim n$.

Moreover, the practical success of SGD has been so remarkable in recent years that a growing body of literature has suggested that its benefit to generalization extends beyond what any improved efficiency might reasonably afford [76, 31, 14, 72, 74]. Some of the myriad explanations for SGD's favorable generalization properties include the local geometry of minimizers [32, 26, 82, 24], connections to approximate Bayesian inference [49], and the regularization properties of noise [75], among many others. While some of these perspectives are intuitive and compelling, they are often difficult to rigorously establish from either an empirical or a theoretical perspective. Empirically, simulations at large scale command significant computational resources, and it can be challenging to push to sufficiently late times or sufficiently large batches to establish the appropriate baselines [72, 75]. Theoretically, the strongest existing results are again in the single-pass setting, for which a number of works have established excess risk bounds for quadratic problems [6, 18, 20, 81]. Much less is known in the multi-pass setting, though stability results were established by [25], and some recent works have begun examining generalization [39].

In this work, we study the dynamics of multi-pass SGD on high-dimensional convex quadratic functions and derive exact asymptotic predictions for the learning and risk trajectories. Our analysis establishes an asymptotic equivalence to a stochastic differential equation, which we call homogenized stochastic gradient descent (HSGD), whose solutions we characterize explicitly in terms of a Volterra integral equation. These results allow us to define a precise data-dependent implicit-conditioning ratio (ICR) that determines whether SGD is more efficient than its full-batch cousins. The ICR favors SGD for many practical datasets, providing some explanation for the observed superior efficiency of SGD; interestingly, we also highlight settings for which SGD is less efficient than full-batch momentum gradient descent, underscoring the data-dependence of the conclusions. Moreover, our results also show that SGD does not improve generalization performance, whether measured in-distribution or out-of-distribution, and therefore that SGD does not offer any form of implicit regularization in this setting. We emphasize that our results do not rule out possible benefits for non-convex problems, but they do provide some of the first explicit negative results in the convex quadratic case.

1.1 Contributions

Our primary contributions are to:

1. Establish the equivalence of quadratic statistics computed on the iterates of SGD and on a particular stochastic Langevin diffusion process called homogenized SGD (Theorem 1);
2. Exactly characterize the asymptotic training and risk trajectories as the solutions of a deterministic Volterra integral equation (Theorem 2);
3. Prove that the noise from SGD negatively impacts generalization performance, both in- and out-of-distribution (Section 3.1), but explain why the impact is often minimal in practice;
4. Introduce the implicit-conditioning ratio that describes when and by how much SGD accelerates convergence relative to the best full-batch methods (Section 3.2);
5. Analyze the limit of streaming SGD to show its inability to capture many salient features of the dynamics of multi-pass SGD (Appendix C).

2 Preliminaries and background

Problem setting. We consider high-dimensional ℓ^2 -regularized least squares problems defined by,

$$\min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \frac{\delta}{2} \|\mathbf{x}\|_2^2 = \sum_{i=1}^n \underbrace{\frac{1}{2} \left((\mathbf{a}_i \mathbf{x} - b_i)^2 + \frac{\delta}{n} \|\mathbf{x}\|_2^2 \right)}_{\stackrel{\text{def}}{=} f_i(\mathbf{x})} \right\}, \quad (1)$$

where $\delta \geq 0$ is the ridge-regularization parameter. We denote the ridgeless empirical risk as

$$\mathcal{L}(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2. \quad (2)$$

On the problem (1), the steps taken by gradient decent (GD) can be written recursively as

$$\mathbf{x}_{k+1}^{\text{m-gd}} = \mathbf{x}_k^{\text{m-gd}} - \gamma_k \nabla f(\mathbf{x}_k^{\text{m-gd}}) = \mathbf{x}_k^{\text{m-gd}} - \gamma_k \mathbf{A}^T (\mathbf{A} \mathbf{x}_k^{\text{m-gd}} - \mathbf{b}) - \gamma_k \delta \mathbf{x}_k^{\text{m-gd}} + \Delta (\mathbf{x}_k^{\text{m-gd}} - \mathbf{x}_{k-1}^{\text{m-gd}}), \quad (3)$$

where $\Delta > 0$ is the momentum parameter, γ_k is the learning rate schedule, and $\mathbf{x}_0 \in \mathbb{R}^d$ is an initial vector assumed to be independent of all other randomness and having norm at most 1. When A is large, computing these updates can be expensive, so an unbiased estimator for the true gradient is often used, where a subset of the data points are selected uniformly at random. We focus on the setting with batch size equal to one and without momentum, which we refer to as stochastic gradient descent (SGD), and for which the iterates can be written recursively as

$$\mathbf{x}_{k+1}^{\text{sgd}} = \mathbf{x}_k^{\text{sgd}} - \gamma_k \nabla f_{i_k}(\mathbf{x}_k^{\text{sgd}}) = \mathbf{x}_k^{\text{sgd}} - \gamma_k \mathbf{A}^T \mathbf{e}_{i_k} \mathbf{e}_{i_k}^T (\mathbf{A} \mathbf{x}_k^{\text{sgd}} - \mathbf{b}) - \frac{\gamma_k \delta}{n} \mathbf{x}_k^{\text{sgd}}, \quad (4)$$

where the $i_k \sim \text{Unif}([n])$ iid. While it would also be possible to consider mini-batch SGD, previous work has shown that batch sizes that are vanishingly small as a fraction of the number of samples are equivalent to the single-batch analysis, after appropriately adjusting the time by a factor of the batch size [60, Theorem 1]; similarly, we do not consider high-dimensional SGD with momentum as it degenerates to SGD [58]. See also [30].

Diffusion approximations and homogenized SGD. A common paradigm for understanding SGD is through stochastic Langevin diffusions (SLD), i.e. solutions of equations of the form

$$d\mathbf{X}_t = -\gamma(\nabla f(\mathbf{X}_t) dt + \sqrt{\Sigma_t} dB_t), \quad (5)$$

where γ is the step size of SGD, f is the loss function, B_t is a d -dimensional standard Brownian motion, and the matrix $0 \preceq \Sigma_t \in \mathbb{R}^{d \times d}$ models the noise covariance. In many analyses, no concrete connection between SGD and SLD is developed, and the diffusion is merely used to build intuition. A common example is the isotropic case ($\Sigma_t \propto \mathbf{I}_d$), for which the Fokker-Planck equation implies that the dynamics are reversible with respect to a density proportional to $e^{-C_\gamma f(x)}$ with $C_\gamma > 0$ some constant. Consequently, the process can escape local minima, exhibiting a trade-off between the entropy and depth of minima and thereby highlighting a possible mechanism of implicit regularization. In the general anisotropic case, describing the stationary distribution is more difficult; nonetheless, the local geometry near minima of f can be analyzed, see [14, 35].

While this type of implicit entropic regularization might ultimately underlie the generalization benefits of SGD for nonconvex problems, currently we lack a precise connection between a concrete SLD and a practical nonconvex learning problem. As such, the implicit regularization effects of SGD on nonconvex losses remains a largely unsolved problem.

For convex quadratics, however, the implications of Eq. (5) are quite clear: there is no notion of implicit regularization as the noise in SLD negatively impacts generalization performance. Note that because the noise is mean zero, any SLD is centered around *gradient flow* (GF) $\mathcal{X}_t^{\text{gf}}$, which solves

$$d\mathcal{X}_t^{\text{gf}} = -\nabla f(\mathcal{X}_t^{\text{gf}}), \quad (6)$$

leading to the following conclusion for generalization (see also [88]):

Lemma 1. *Suppose the objective function is $f(\mathbf{x}) = \frac{1}{2}(\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \delta\|\mathbf{x}\|^2)$. Suppose $(\mathbf{X}_t : t \in [0, \infty))$ is an SLD (i.e. \mathbf{X}_t solves (5)) with $\|\Sigma_t\|_{\text{op}}$ almost surely bounded by some $C < \infty$. Suppose the population risk $\mathcal{R} : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function and denote $\mathbf{x}_* \stackrel{\text{def}}{=} \lim_{t \rightarrow \infty} \mathcal{X}_t^{\text{gf}}$, then*

$$\underbrace{\mathbb{E}[\mathcal{R}(\mathbf{X}_t)]}_{\text{pop. risk of SLD}} \geq \underbrace{\mathcal{R}(\mathcal{X}_{\gamma t}^{\text{gf}})}_{\text{pop. risk of gradient flow}} \quad \text{for all } t \geq 0 \text{ and hence,} \quad \underbrace{\liminf_{t \rightarrow \infty} \mathbb{E}[\mathcal{R}(\mathbf{X}_t)]}_{\text{limiting pop. risk of SLD}} \geq \underbrace{\mathcal{R}(\mathbf{x}_*)}_{\text{limiting pop. risk gradient flow}}.$$

If in addition \mathcal{R} is strictly convex, and $\Sigma_t \rightarrow \Sigma_\infty$ with $\Sigma_\infty \succ 0$, then the inequality is strict.

Proof. The mean $\mathbb{E}[\mathbf{X}_t]$, by the linearity of the gradient ∇f , is GF. Under the conditions given, the law of \mathbf{X}_t converges to a Gaussian variable centered at \mathbf{x}_* . Hence by Fatou's lemma and Jensen's inequality, the inequality follows. \square

We emphasize that this conclusion applies even under general distribution shifts, so long as the risk remains a convex function. Still, the utility of Lemma 1 may not be immediately clear, as it pertains

to SLD and we have not yet established any concrete connection between SLD and the process of interest, SGD. Nor is it evident what form such a connection should take—the agreement between SGD and an SLD cannot occur at the level of individual states since the randomness from each process is not assumed to be coupled. Instead, the most we can hope for is that statistics of the processes agree. Specifically, we might hope that matching the noise structure of SGD with a careful choice of SLD will cause relevant statistics, like the population risk, to be equal.

It turns out that such a choice of SLD exists for convex quadratic problems in high dimensions, and is given by *homogenized SGD* (HSGD), introduced simultaneously in [52, 58]. Both the empirical and population risks (\mathcal{L} , \mathcal{R} resp.) of HSGD agree with the same of SGD in the high-dimensional limit (see Thm. 1). Mathematically, HSGD is the strong solution of the stochastic differential equation:

$$d\mathbf{X}_t \stackrel{\text{def}}{=} -\gamma(t)\nabla\mathcal{L}(\mathbf{X}_t) dt + \gamma(t)\sqrt{\frac{2}{n}\mathcal{L}(\mathbf{X}_t)\nabla^2\mathcal{L}(\mathbf{X}_t)} d\mathbf{B}_t, \quad \text{for quadratic } \mathcal{L}, \quad (7)$$

where again \mathbf{B}_t is a d -dimensional standard Brownian motion, $\gamma(t)$ is the learning rate schedule, and the initial condition is $\mathbf{X}_0 = \mathbf{x}_0$. Roughly, HSGD is a diffusion approximation to SGD that gains explanatory power when the *dimensionality* is large. In particular, it does not require the step size γ to be small, in contrast to the usual paradigm of SLD approximations. Note that as with other universality results, the details of the noise distribution are not relevant and only the second-order correlations contribute, which are carefully matched by HSGD to SGD.

The precise sense of the comparison requires us to evaluate low-dimensional statistics of the high-dimensional dynamics; “low-dimensional” must be effective, in that the univariate statistics of the SGD iterates concentrate around the same statistic evaluated on HSGD. For understanding generalization or implicit regularization properties, a important statistic is the population risk, \mathcal{R} .

Assumptions. For all parts of our analysis to hold, the pair (\mathbf{A}, \mathbf{b}) of the data matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ and target vector $\mathbf{b} \in \mathbb{R}^n$ must satisfy some *quasi-random* assumptions—a set of deterministic conditions on the pair (\mathbf{A}, \mathbf{b}) that are satisfied with high probability by natural classes of random matrix-vector pairs (see Appendix B for specifics). We use the convention that the target and initialization vectors are bounded independent of n , $\|\mathbf{b}\|_2 \leq C$ and $\|\mathbf{x}_0\|_2 \leq C$, respectively.

We illustrate some examples below that we have shown to satisfy the quasi-random assumptions.

- Gaussian linear regression. Here the rows of \mathbf{A} are iid and drawn from a Gaussian with norm-bounded covariance Σ and the target \mathbf{b} is drawn from a generative model, $\mathbf{b} = \mathbf{A}\tilde{\mathbf{x}} + \boldsymbol{\eta}$ for some unknown signal $\tilde{\mathbf{x}} \in \mathbb{R}^d$ and independent noise $\boldsymbol{\eta} \in \mathbb{R}^n$.
- Subgaussian linear designs. In the example above, we can relax the Gaussian assumption to be of the form $\mathbf{x}\Sigma^{1/2}$ for \mathbf{x} a vector of iid centered subgaussian random variables [88, 30].
- Gaussian random features with a linear ground truth [51, 3, 66, 4]. Suppose \mathbf{A} is given by $\sigma(\mathbf{X}\mathbf{W})$ for an iid standard Gaussian weight matrix \mathbf{W} and Gaussian data matrix \mathbf{X} . Suitable assumptions on the activation function σ and the covariance Σ of \mathbf{X} added.

Assumption 1. *The population risk $\mathcal{R} : \mathbb{R}^d \rightarrow \mathbb{R}$ is a quadratic, that is, it is a degree-2 polynomial or, equivalently, can be represented by*

$$\mathcal{R}(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T\mathbf{T}\mathbf{x} + \mathbf{u}^T\mathbf{x} + c$$

for some $d \times d$ symmetric matrix \mathbf{T} , vector $\mathbf{u} \in \mathbb{R}^d$, and scalar $c \in \mathbb{R}$. We further assume that $\|\nabla^2\mathcal{R}\|_{op} \leq C$, $\|\nabla\mathcal{R}(0)\|_2 \leq C$, and $|\mathcal{R}(0)| \leq C$.

A natural population risk is given by $\mathcal{R}(\mathbf{x}) = \frac{1}{2}\mathbb{E}[(\mathbf{a} \cdot \mathbf{x} - b)^2]$ where $(\mathbf{a}, b) \sim \mathcal{D}$. This distribution \mathcal{D} may or may not be the same as the distribution that generated the data $[\mathbf{A} \mid \mathbf{b}]$ used in training.

As we work in the high-dimensional limit, we suppose that $\gamma_k = \gamma(k/n)$ for a smooth, bounded function $\gamma(\cdot)$ such that $\gamma(t) \rightarrow \gamma \in [0, \infty)$ and $\hat{\gamma} \stackrel{\text{def}}{=} \sup_{t \geq 0} \gamma(t) < \infty$.

3 Main results

Our main results are analyzable (non-asymptotic) expressions for the empirical risk \mathcal{L} and the population risk \mathcal{R} of SGD at any time t for the high-dimensional least squares problem (1). To begin, we first establish the following equivalence between SGD and HSGD.

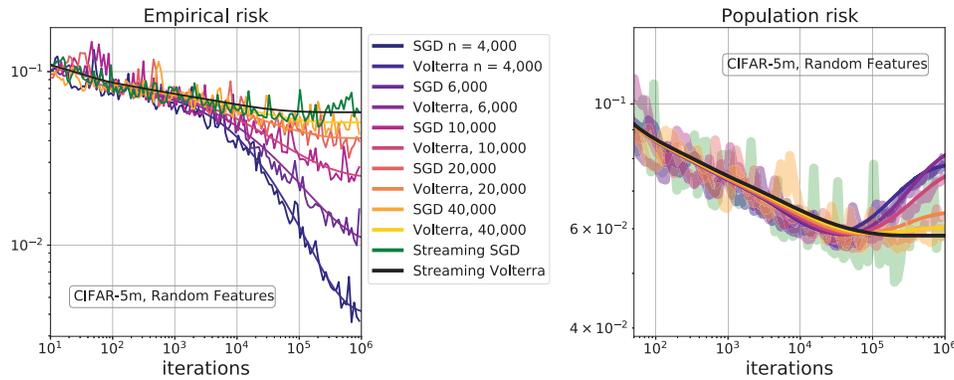


Figure 1: **Single runs of SGD vs. HSGD (Volterra) in streaming** on standardized CIFAR-5M [54] with car/plane class vector (1,000,000 samples); a standardized ReLu (74) random features model (see Appendix B.1) was applied with increasing number of samples n and fixed $d = 6000$. The predicted behavior from HSGD (denoted by Volterra) matches the performance of single runs of SGD for finite n and streaming ($n = \infty$). Shaded region (right) is the moving average of a single run of SGD. Empirical risk (left) increases monotonically with n to its limit while population risk generally decreases with n . Streaming corresponds to $n = \infty$ (see Appendix C). For consistency across sample sizes, time is measured in iterations. Additional details in App G.

Theorem 1 (Equivalence of SGD and HSGD). *Suppose the pair $(\mathbf{A}, \mathbf{b}) \in \mathbb{R}^{n \times d} \times \mathbb{R}^d$ satisfy the quasi-random assumptions with $d^\epsilon \leq n \leq d^{1/\epsilon}$ for some $\epsilon \in (0, 1]$. Let the iterates $\mathbf{x}_t = \mathbf{x}_{[t]}^{sgd}$ be generated from multi-pass SGD Eq. (4) and \mathbf{X}_t be the solution of Eq. (7). Then for any deterministic $T > 0$ and any $D > 0$, there is a $C > 0$ such that*

$$\Pr \left[\sup_{0 \leq t \leq T} \left\| \begin{pmatrix} \mathcal{L}(\mathbf{x}_{[tn]}) \\ \mathcal{R}(\mathbf{x}_{[tn]}) \end{pmatrix} - \begin{pmatrix} \mathcal{L}(\mathbf{X}_t) \\ \mathcal{R}(\mathbf{X}_t) \end{pmatrix} \right\|_2 > d^{-\epsilon/2} \right] \leq Cd^{-D}.$$

The rigorous proof is given in [61]. For the rest of this paper, we will use homogenized SGD to analyze the behavior of multi-pass SGD. See also Appendix C where we heuristically extend this to the case of streaming SGD.

While the comparison of SGD to HSGD requires relatively strong assumptions on \mathbf{A} and \mathbf{b} , the analysis of HSGD can be performed under weaker assumptions (no quasirandomness assumptions are needed). It suffices to suppose the problem is high dimensional in the following sense:

Assumption 2. *The empirical risk \mathcal{L} satisfies $\text{tr} \nabla^2 \mathcal{L} = n$ and $0 \preceq \nabla^2 \mathcal{L} \preceq nd^{-\epsilon}$ for some $\epsilon > 0$.*

This corresponds to the normalization where $\nabla^2 \mathcal{L} = \mathbf{A}^T \mathbf{A}$ and each row of \mathbf{A} is length 1 and hence $\text{tr} \nabla^2 \mathcal{L} = n$.

Under Assumptions 1 and 2, the dynamics of the empirical and population risk under HSGD concentrate around a deterministic dynamical system driven by a Volterra integral equation:

Volterra Dynamics, Multi-pass. The following deterministic dynamical system is the high-dimensional limit for $\mathcal{L}(\mathbf{X}_t)$ and $\mathcal{R}(\mathbf{X}_t)$, respectively

$$\Psi_t = \mathcal{L}(\mathcal{X}_{\Gamma(t)}^{\text{gf}}) + \int_0^t K(t, s; \nabla^2 \mathcal{L}) \Psi_s \, ds \quad (\text{Empirical risk}) \quad (8)$$

$$\Omega_t = \mathcal{R}(\mathcal{X}_{\Gamma(t)}^{\text{gf}}) + \int_0^t K(t, s; \nabla^2 \mathcal{R}) \Psi_s \, ds \quad (\text{Population risk}) \quad (9)$$

where the *integrated learning rate* Γ and *kernel* K , for any $d \times d$ matrix \mathbf{P} , respectively are

$$\Gamma(t) = \int_0^t \gamma(s) \, ds, \quad K(t, s; \mathbf{P}) = \frac{\gamma^2(s)}{n} \text{tr}((\nabla^2 \mathcal{L}) \mathbf{P} \exp(-2(\nabla^2 \mathcal{L} + \delta \mathbf{I}_d)(\Gamma(t) - \Gamma(s))))). \quad (10)$$

Theorem 2 (Concentration of HSGD around Volterra dynamics). *Under Assumptions 1 and 2, for any $T > 0$ and for any $D > 0$ there exists sufficiently large $C > 0$ such that for all $d > 0$*

$$\Pr \left[\sup_{0 \leq t \leq T} \left\| \begin{pmatrix} \mathcal{L}(\mathbf{X}_t) \\ \mathcal{R}(\mathbf{X}_t) \end{pmatrix} - \begin{pmatrix} \Psi_t \\ \Omega_t \end{pmatrix} \right\| > d^{-\epsilon/2} \right] \leq C d^{-D},$$

where Ψ_t and Ω_t solve (8) and (9).

We give a formal proof of the concentration result in Appendix D.1 in Theorem 11.

3.1 No implicit regularization from SGD

From (9), for convex \mathcal{R} we observe immediately that the population risk Ω_t is only larger than the population risk of GF. Moreover, we have an explicit formula for the excess risk due to SGD noise,

$$\underbrace{\Omega_t - \mathcal{R}(\mathcal{X}_{\Gamma(t)}^{\text{gf}})}_{\text{excess risk due to SGD}} \stackrel{\text{def}}{=} \int_0^t K(t, s; \nabla^2 \mathcal{R}) \times \underbrace{\Psi_s}_{\text{limiting loss } \mathcal{L}} \, ds.$$

Note that the population risk of SGD tracks that of GF. If GF overfits, SGD overfits as well; there is no statistical regularization due to the noise of SGD applied to empirical risk minimization (ERM).

We can further analyze the long-time behavior of SGD with exact limiting values for this excess risk.

Theorem 3 (Time infinity risk values). *If $\gamma(t) \rightarrow 0$ as $t \rightarrow \infty$ but $\Gamma(t) \rightarrow \infty$ (i.e. the usual Robbins-Monro setting), then the excess population risk of SGD over GF tends to 0. If on the other hand $\gamma(t) \rightarrow \gamma \in (0, 2(\frac{1}{n} \text{tr} \{ \frac{(\mathbf{A}^T \mathbf{A})^2}{\mathbf{A}^T \mathbf{A} + \delta \mathbf{I}_d} \})^{-1})$, then with Ψ_∞ given by the limiting empirical risk,*

$$\Psi_\infty = \mathcal{L}(\mathcal{X}_\infty^{\text{gf}}) \times \left(1 - \frac{\gamma}{2n} \text{tr} \left\{ \frac{(\nabla^2 \mathcal{L})^2}{\nabla^2 \mathcal{L} + \delta \mathbf{I}_d} \right\} \right)^{-1}$$

the excess risk due to SGD converges to

$$\Omega_t - \mathcal{R}(\mathcal{X}_{\Gamma(t)}^{\text{gf}}) \rightarrow \frac{\gamma \Psi_\infty}{2n} \times \text{tr} \left\{ \frac{(\nabla^2 \mathcal{R})(\nabla^2 \mathcal{L})}{\nabla^2 \mathcal{L} + \delta \mathbf{I}_d} \right\}.$$

There are a few conclusions to draw directly from this. In the interpolation regime, that is where $\Psi_\infty = 0$, there is no excess risk due to SGD and there is no need to send γ to 0. Moreover, if the empirical risk Ψ_∞ is small, the excess risk due to SGD is proportional to $\gamma \Psi_\infty$, and hence it is frequently orders of magnitude smaller than other potential sources of error. Furthermore, the excess risk is affected by how similar the population and empirical risks are, in the large directions. Ridge regularization can substantially reduce excess risk due to SGD in cases where population risk has many small eigenvalues. In summary, either by sending $\gamma \rightarrow 0$, working in the interpolation regime, or otherwise in a regime Ψ_∞ is small, the excess risk incurred by running SGD is minimal.

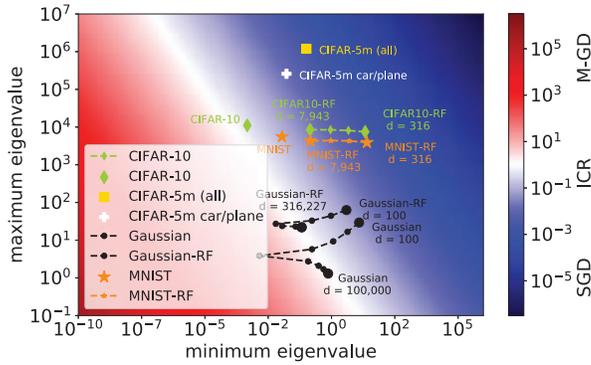


Figure 2: ICR as function of largest and smallest eigenvalues of \mathbf{A} with trace normalized to be 1; blue (smaller, SGD favored) and red (larger, full batched MGD favored). Points indicate ICR for image datasets (MNIST, CIFAR-10, CIFAR-5m) as well as their images under ReLU random feature maps of various dimensions (lines). Gaussian random features ($n = 2000$, $n_0 = 100$, various d) and Gaussian data ($n = 2000$, various d) shows ICR for over- and under-parameterized models. Natural datasets tend to favor SGD.

3.2 Implicit conditioning of SGD

In contrast, the algorithmic advantages of SGD are substantial. To simplify the discussion, we consider only the case of constant learning rate γ . In this case, the kernel in (8) and (9) simplifies to a convolution kernel, which has a much simpler theory. To characterize the rates, we define λ_{\min} as the smallest non-zero eigenvalue of $\nabla^2 \mathcal{L}$. Then for generic initial conditions, (in particular almost surely if \mathbf{X}_0 is nonzero isotropic), GF has the following convergence rate

$$\lim_{t \rightarrow \infty} \left(\mathcal{L}(\mathbf{x}_{\gamma t}^{\text{gf}}) - \mathcal{L}(\mathbf{x}_{\infty}^{\text{gf}}) \right)^{1/t} = \begin{cases} e^{-\gamma(\lambda_{\min}(\nabla^2 \mathcal{L}) + \delta)}, & \text{if } \delta > 0, \\ e^{-2\gamma\lambda_{\min}(\nabla^2 \mathcal{L})}, & \text{otherwise.} \end{cases}$$

Here we use the notation that $\lambda_{\min}(\mathbf{H})$ and $\lambda_{\max}(\mathbf{H})$ are the smallest and largest eigenvalues of the matrix \mathbf{H} . The rate of convergence of Ψ_t to Ψ_{∞} can be no faster than the underlying GF, given by the rate above. On the other hand, for larger γ the Volterra term in (8) can frustrate the convergence. The *Malthusian exponent* of the convolution Volterra equation is given by

$$\lambda_* = \inf \left\{ x : 1 = \int_0^{\infty} e^{xt} K(t; \nabla^2 \mathcal{L}) dt \stackrel{\text{def}}{=} \gamma^2 \int_0^{\infty} e^{xt} \text{tr}((\nabla^2 \mathcal{L})^2 \exp(-2\gamma(\nabla^2 \mathcal{L} + \delta \mathbf{I}_d)t)) dt \right\}. \quad (11)$$

As $\nabla^2 \mathcal{L}$ is finite dimensional, we have that $\lambda_* \leq 2\gamma(\lambda_{\min}(\nabla^2 \mathcal{L}) + \delta)$, owing to the divergence of the integral as x approaches this value from below. Note that in principal the Malthusian exponent can be negative, in which case SGD is *divergent*. The Malthusian exponent gives the effective rate of convergence of constant learning rate SGD. Define

$$\Xi(\gamma) \stackrel{\text{def}}{=} \begin{cases} \min\{\gamma(\lambda_{\min}(\nabla^2 \mathcal{L}) + \delta), \lambda_*(\gamma)\} & \text{if } \delta > 0, \\ \lambda_*(\gamma) & \text{if } \delta = 0. \end{cases} \quad (12)$$

Theorem 4 (SGD convergence rates, average-case). *Then the rates of convergence of both the empirical and population risk are controlled by this parameter*

$$\lim_{t \rightarrow \infty} (\Psi_t - \Psi_{\infty})^{1/t} = e^{-\Xi(\gamma)} = \lim_{t \rightarrow \infty} (\Omega_t - \Omega_{\infty})^{1/t}.$$

Furthermore, when $\gamma = n(\text{tr}(\mathbf{A}^T \mathbf{A}))^{-1}$, we have the rate guarantee $\Xi(\gamma) \geq \frac{\lambda_{\min}(\nabla^2 \mathcal{L}) + \delta}{2}$.

The major difference between SGD and full batch methods such as momentum gradient descent (MGD; see Appendix F.2 for definitions) is that they have different sensitivities to the Hessian spectrum of the empirical risk $\mathcal{L}(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$. Define the condition numbers

$$\kappa \stackrel{\text{def}}{=} \frac{\lambda_{\max}(\nabla^2 \mathcal{L}) + \delta}{\lambda_{\min}(\nabla^2 \mathcal{L}) + \delta} \quad \text{and} \quad \bar{\kappa} \stackrel{\text{def}}{=} \frac{\frac{1}{n} \text{tr}(\nabla^2 \mathcal{L})}{\lambda_{\min}(\nabla^2 \mathcal{L}) + \delta}.$$

The first of these is the classical condition number of the ridge problem, while the second is the averaged condition number that regulates the behavior of SGD in the high-dimensional limit. MGD has been long established to have a rate of convergence, with proper tuning, controlled by the square root of the condition number [65], which is known to be optimal amongst first order algorithms.

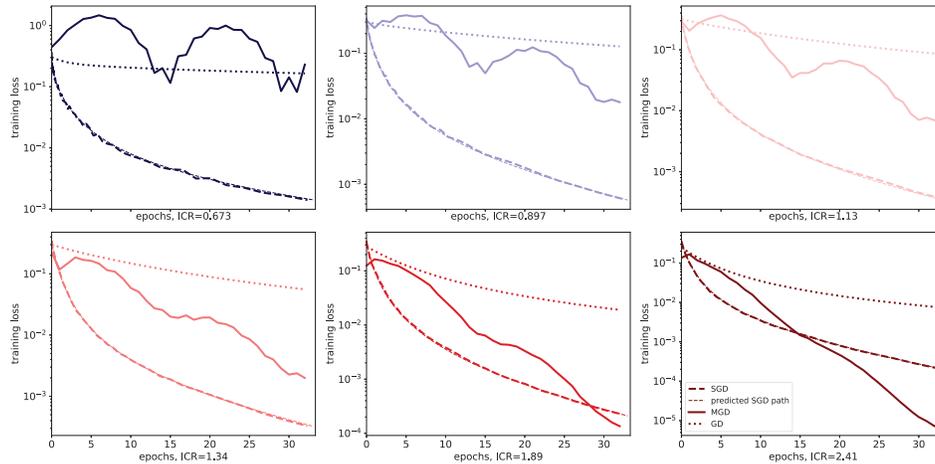


Figure 3: **ICR effect on full batch MGD versus SGD** in a synthetic least-squares setting. We consider minimizing, for an $n \times d$ matrix \mathbf{A} (with $n = 2400, d = 3600$) $\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$, where \mathbf{A} is a Gaussian matrix, with correlated rows and \mathbf{b} is given by $\mathbf{A}\boldsymbol{\beta} + \boldsymbol{\eta}$ for a ground truth, which is isotropic normal of expected norm-square 1, and $\boldsymbol{\eta}$ is isotropic normal of expected norm-square 0.02.

Theorem 5 (Convergence rates for MGD). *For isotropic random initialization \mathbf{x}_0 or noisy \mathbf{b} , $\delta > 0$, and strictly convex population risk \mathcal{R}*

$$(\mathcal{L}(\mathbf{x}_k^{m-gd}) - \mathcal{L}(x_*))^{1/k} \xrightarrow[k \rightarrow \infty]{\text{a.s.}} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right) \quad \text{and} \quad (\mathcal{R}(\mathbf{x}_k^{m-gd}) - \mathcal{R}(x_*))^{1/k} \xrightarrow[k \rightarrow \infty]{\text{a.s.}} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right).$$

See Appendix F.2 for elaboration.

In light of Theorems 4 and 5, we can define the *implicit-conditioning ratio* as

$$\text{ICR} \stackrel{\text{def}}{=} \frac{\bar{\kappa}}{\sqrt{\kappa}} \approx \log \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right) \bar{\kappa},$$

which measures the efficiency of SGD over MGD in that SGD with constant learning rate $n / \text{tr}(\nabla^2 \mathcal{L})$ trains in an ICR-multiple of the number of epochs that MGD requires (lower is better for SGD).

Problems favor SGD when there are large outlier eigenvalues, a common feature of Hessian spectra in practice [69, 70, 1]. Indeed, if the largest eigenvalues are on the same order as the *unnormalized trace*, individual SGD iterates are as effective as full-batch gradient. In contrast, when the Hessian spectrum is tightly packed, which is less common in practice but can occur after some preprocessing techniques or e.g. for uncorrelated Gaussian samples, then MGD is favored. See Fig. 2.

3.3 Numerical results on ICR

In Fig. 3, we illustrate how ICR affects the relative performance of full batch MGD versus single batch SGD in a synthetic least squares setting where we have tight control over the all the particulars of the problem and in a neural network setting. We control the ICR by controlling the covariance singular value spectrum of the rows of \mathbf{A} , which we take as Pareto distributed with exponent s for varying $s > 2$. This choice allows us to affect the ICR of the problem without changing the minimum curvature of the Hessian. The results, comparing SGD, MGD, and GD, are given in Figure 3.

We note a few key qualitative observations. First, even in MGD favored configurations, SGD will outperform MGD on short time scales. When optimizing the hyperparameters in MGD for long-time performance, minimal curvature (which in this case is just the minimal eigenvalue of $\mathbf{A}\mathbf{A}^T$) plays a major role in the choices; being tuned for long-time performance, MGD typically performs suboptimally at initialization. In contrast, the learning rate in SGD only depends on average curvature, and so it generally performs better at initialization on problems with a larger interval of Hessian spectra.

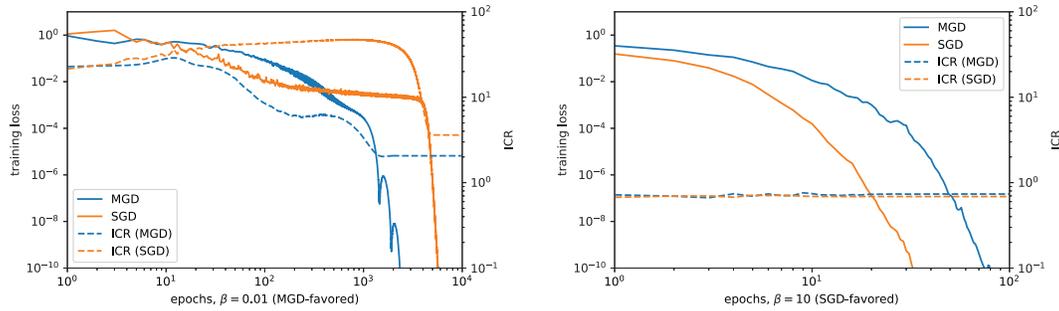


Figure 4: **ICR and full-batch MGD versus SGD** with a fully connected 2-layer neural network with activation function $f(x) = \frac{1}{\sqrt{\beta}} \operatorname{erf}(\sqrt{\beta}x)$ on a 500-sample subset of CIFAR-10 with targets car/plane. Momentum hyper-parameters were tuned empirically to give best performance. We note that the loss curve of MGD also displays the characteristic cycloidal oscillations of tuned momentum on inhomogeneous problems. *Left*: $\beta = 0.01$. ICR varies over the course of training but stays above 1, favoring momentum as illustrated in the figure. *Right*: $\beta = 10$. ICR varies slightly but always stays below 1 throughout the training and SGD is favored over full-batch momentum.

Second, we note that the problem setup was chosen to hold the minimum curvature roughly constant while varying s . When s tends to 2 from above, the largest eigenvalues of $\mathbf{A}\mathbf{A}^T$ grows with feature dimension d , but the average and minimum eigenvalue stays bounded with feature dimension. Hence we can send the ICR to 0 by choosing an s above 2 and increasing d (or n).¹ On the other hand, by sending $s \rightarrow \infty$, we send the covariance matrix to the identity, which tends to be momentum favored.²

In Fig. 4, we run SGD on a fully-connected 2-layer neural network on a subset of CIFAR-10 in order to examine the dynamics of the ICR for a non-trivial problem and to see how our insights might play out in practice. Owing to the non-convexity of this problem, we define the ICR in terms of the Gauss-Newton approximation to the Hessian, or equivalently in terms of the Neural Tangent Kernel [29]. By changing the activation function of the network, we can vary the initial ICR from an SGD-favored to a momentum-favored value. While the ICR does change over the course of training, we find that, at least in this setting, the initial ICR can nevertheless predict the relative performance of SGD versus MGD. Indeed, for activation functions for which the ICR remains above 1.0, the training remains MGD-favored over sufficiently long times, and we observe that MGD with optimal parameters does converge faster than SGD. In contrast, when the ICR remains below 1.0, we find that SGD outperforms MGD.

4 Conclusion.

Using a specific type of SLD (called HSGD) that matches the second-order correlations in the noise of SGD, we demonstrated that their empirical and population risks match in the high-dimensional limit. Moreover, the risks of HSGD behavior deterministically, as described by a Volterra equation. With this connection, we investigated the benefits of SGD on a convex objective. While there is no statistical benefit to generalization from the noise of SGD, in overparameterized, interpolating settings little is lost compared to GD. Moreover, when computational restrictions are imposed, SGD can be radically faster than GD because of its dependence on a different condition number of the Hessian. We characterized this speed up using the ICR, which when calculated for datasets common in deep learning clearly favors SGD. This should highlight the difficulty in studying implicit regularization for SGD empirically: any experiment necessarily has a finite computational budget and may find

¹The relevance of $s > 2$ is that the Pareto has moments up to and including the second moment. For values less than 2, the covariance spectra is sufficiently heavy that the maximum eigenvalue of $\mathbf{A}\mathbf{A}^T$ dominates the trace. In that regime, the problem becomes effectively sparse, with an intrinsic dimension depending only on s , and it should be expected that GD/SGD are approximately equivalent and outperform MGD.

²Furthermore, the 'average curvature' speedup of SGD on short time scales becomes muted, owing to all curvatures being the same.

lower population risks with SGD simply via its improved conditioning. Finally, we demonstrated limitations in using streaming SGD alone as a tool for studying generalization.

As future work, a major outstanding problem (both theoretically and empirically) is extending the analysis above to non-quadratic losses, both train and test, and especially to other high-dimensional problems not in the kernel regime. Finally, data augmentation can naturally be considered by randomly augmenting each sample from $\hat{\mathcal{D}}_n$ in Eq. (30).

Acknowledgments and Disclosure of Funding

C. Paquette's research was supported by CIFAR AI Chair, MILA, a Discovery Grant from the Natural Science and Engineering Council (NSERC), and the FRQNT New University Researcher's Start-up Program. Research by E. Paquette was supported by a Discovery Grant from the Natural Science and Engineering Council (NSERC). Additional revenues related to this work: C. Paquette has part-time employment at Google Research, Brain Team, Montreal, QC.

References

- [1] A. Abdel-Gawad and S. Ratner. Adaptive optimization of hyperparameters in L2-regularised logistic regression. *Technical report*, 2007.
- [2] B. Adlam and J. Pennington. The Neural Tangent Kernel in High Dimensions: Triple Descent and a Multi-Scale Theory of Generalization. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119, pages 74–84, 13–18 Jul 2020.
- [3] B. Adlam and J. Pennington. Understanding Double Descent Requires A Fine-Grained Bias-Variance Decomposition. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pages 11022–11032, 2020.
- [4] B. Adlam, J.A. Levinson, and J. Pennington. A Random Matrix Perspective on Mixtures of Nonlinearities in High Dimensions. In *International Conference on Artificial Intelligence and Statistics*, pages 3434–3457. PMLR, 2022.
- [5] S. Asmussen. *Applied probability and queues*, volume 51 of *Applications of Mathematics (New York)*. Springer-Verlag, New York, second edition, 2003. Stochastic Modelling and Applied Probability.
- [6] F. Bach and E. Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate $O(1/n)$. *Advances in neural information processing systems (NeurIPS)*, 26, 2013.
- [7] D. Barrett and B. Dherin. Implicit Gradient Regularization. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=3q5IqUrkcF>.
- [8] P.L. Bartlett, P.M. Long, G. Lugosi, and A. Tsigler. Benign overfitting in linear regression. *Proc. Natl. Acad. Sci. USA*, 117(48):30063–30070, 2020. doi: 10.1073/pnas.1907378117.
- [9] B. Bordelon and C. Pehlevan. Learning Curves for SGD on Structured Features. In *International Conference on Learning Representations (ICLR)*, 2022.
- [10] L. Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010.
- [11] L. Bottou. Stochastic gradient descent tricks. In *Neural networks: Tricks of the trade*, pages 421–436. Springer, 2012.
- [12] L. Bottou and O. Bousquet. The tradeoffs of large scale learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 20, 2007.
- [13] J.-P. Bouchaud and A. Georges. Anomalous diffusion in disordered media: statistical mechanisms, models and physical applications. *Physics reports*, 195(4-5):127–293, 1990.
- [14] P. Chaudhari and S. Soatto. Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. In *International Conference on Learning Representations (ICLR)*, pages 1–20, 2018.
- [15] X. Cheng, N. Chatterji, Y. Abbasi-Yadkori, P. Bartlett, and M. Jordan. Sharp convergence rates for Langevin dynamics in the nonconvex setting. *arXiv preprint arXiv:1805.01648*, 2018.
- [16] A. Defossez and F. Bach. Averaged Least-Mean-Squares: Bias-Variance Trade-offs and Optimal Sampling Distributions. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 38, pages 205–213, 2015.
- [17] M. Dereziński, F. T. Liang, and M. W. Mahoney. Exact expressions for double descent and implicit regularization via surrogate random design. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pages 5152–5164, 2020.
- [18] A. Dieuleveut, N. Flammarion, and F. Bach. Harder, better, faster, stronger convergence rates for least-squares regression. *The Journal of Machine Learning Research*, 18(1):3520–3570, 2017.

- [19] M. Engeli, Th. Ginsburg, H. Rutishauser, and E. Stiefel. Refined iterative methods for computation of the solution and the eigenvalues of self-adjoint boundary value problems. *Mitt. Inst. Angew. Math. Zürich*, 8:107, 1959.
- [20] R. Ge, S.M. Kakade, R. Kidambi, and P. Netrapalli. The step decay schedule: A near optimal, geometrically decaying learning rate procedure for least squares. *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 2019.
- [21] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [22] G. Gripenberg. On the resolvents of nonconvolution Volterra kernels. *Funkcial. Ekvac.*, 23(1): 83–95, 1980.
- [23] S. Gunasekar, J. Lee, D. Soudry, and N. Srebro. Characterizing Implicit Bias in Terms of Optimization Geometry. *Proceedings of the 35th International Conference on Machine Learning (ICML)*, 2018.
- [24] M. Gurbuzbalaban, U. Simsekli, and L. Zhu. The Heavy-Tail Phenomenon in SGD. In *International Conference on Learning Representations (ICLR)*, 2020.
- [25] M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, volume 48, pages 1225–1234, 2016.
- [26] S. Hochreiter and J. Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997.
- [27] E. Hoffer, I. Hubara, and D. Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017.
- [28] W. Hu, C. Li, L. Li, and J. Liu. On the diffusion approximation of nonconvex stochastic gradient descent. *arXiv preprint arXiv:1705.07562*, 2017.
- [29] A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems (NeurIPS)*, 2018.
- [30] P. Jain, S. Kakade, R. Kidambi, P. Netrapalli, and A. Sidford. Accelerating Stochastic Gradient Descent for Least Squares Regression. In *Proceedings of the 31st Conference On Learning Theory (COLT)*, volume 75, pages 545–604, 2018.
- [31] S. Jastrzebski, Z. Kenton, D. Arpit, N. Ballas, A. Fischer, Y. Bengio, and A. Storkey. Three factors influencing minima in sgd. *arXiv preprint arXiv:1711.04623*, 2017.
- [32] N. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. Tang. On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima. *arXiv preprint arXiv:1609.04836*, 2016.
- [33] D. Kobak, J. Lomond, and B. Sanchez. The Optimal Ridge Penalty for Real-world High-dimensional Data Can Be Zero or Negative due to the Implicit Ridge Regularization. *Journal of Machine Learning Research*, 21(169):1–16, 2020.
- [34] A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, U. of Toronto, 2009.
- [35] D. Kunin, J. Sagastuy-Brena, L. Gillespie, E. Margalit, H. Tanaka, S. Ganguli, and D. Yamins. Rethinking the limiting dynamics of SGD: modified loss, phase space oscillations, and anomalous diffusion. *arXiv preprint arXiv:2107.09133*, 2021.
- [36] H. Kushner and G.G. Yin. *Stochastic approximation and recursive algorithms and applications*, volume 35. Springer Science & Business Media, 2003.
- [37] Y. LeCun, C. Cortes, and C. Burges. "mnist" handwritten digit database, 2010. URL <http://yann.lecun.com/exdb/mnist>.

- [38] J. Lee, Y. Bahri, R. Novak, S. Schoenholz, J. Pennington, and J. Sohl-Dickstein. Deep Neural Networks as Gaussian Processes. In *International Conference on Learning Representations (ICLR)*, 2018.
- [39] Y. Lei, T. Hu, and K. Tang. Generalization Performance of Multi-pass Stochastic Gradient Descent with Convex Loss Functions. *J. Mach. Learn. Res.*, 22:25–1, 2021.
- [40] A. Lewkowycz, Y. Bahri, E. Dyer, J. Sohl-Dickstein, and G. Gur-Ari. The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*, 2020.
- [41] C. Li, H. Farkhoor, R. Liu, and J. Yosinski. Measuring the intrinsic dimension of objective landscapes. *arXiv preprint arXiv:1804.08838*, 2018.
- [42] Q. Li, C. Tai, and W. E. Stochastic Modified Equations and Adaptive Stochastic Gradient Algorithms. In *Proceedings of the 34th International Conference on Machine Learning (ICLR)*, volume 70, pages 2101–2110, 2017.
- [43] Q. Li, C. Tai, and W. E. Stochastic Modified Equations and Dynamics of Stochastic Gradient Algorithms I: Mathematical Foundations. *Journal of Machine Learning Research*, 20(40):1–47, 2019.
- [44] Z. Liao, R. Couillet, and M. Mahoney. A Random Matrix Analysis of Random Fourier Features: Beyond the Gaussian Kernel, a Precise Phase Transition, and the Corresponding Double Descent. *arXiv preprint arXiv:2006.05013*, 2020.
- [45] J. Lin and L. Rosasco. Optimal Rates for Multi-pass Stochastic Gradient Methods. *Journal of Machine Learning Research*, 18(97):1–47, 2017.
- [46] L. Ljung. Analysis of recursive stochastic algorithms. *IEEE Trans. Automatic Control*, AC-22(4):551–575, 1977.
- [47] L. Ljung, G. Pflug, and H. Walk. *Stochastic approximation and optimization of random systems*, volume 17 of *DMV Seminar*. Birkhäuser Verlag, Basel, 1992. doi: 10.1007/978-3-0348-8609-3.
- [48] S. Mandt, M. Hoffman, and D. Blei. A variational analysis of stochastic gradient algorithms. In *International conference on machine learning (ICML)*, 2016.
- [49] S. Mandt, M.D. Hoffman, and D. Blei. Stochastic Gradient Descent as Approximate Bayesian Inference. *Journal of Machine Learning Research*, 18(134):1–35, 2017.
- [50] V.A. Marčenko and L.A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Mathematics of the USSR-Sbornik*, 1967.
- [51] S. Mei and A. Montanari. The generalization error of random features regression: precise asymptotics and the double descent curve. *Comm. Pure Appl. Math.*, 75(4):667–766, 2022.
- [52] T. Mori, L. Ziyin, K. Liu, and M. Ueda. Logarithmic landscape and power-law escape rate of sgd. *arXiv preprint arXiv:2105.09557*, 2021.
- [53] E. Moulines and F. Bach. Non-Asymptotic Analysis of Stochastic Approximation Algorithms for Machine Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 24, 2011.
- [54] P. Nakkiran, B. Neyshabur, and H. Sedghi. The Deep Bootstrap Framework: Good Online Learners are Good Offline Generalizers. In *International Conference on Learning Representations (ICLR)*, 2021.
- [55] R.M. Neal. *Priors for Infinite Networks*, pages 29–53. Springer New York, New York, NY, 1996. doi: 10.1007/978-1-4612-0745-0_2. URL https://doi.org/10.1007/978-1-4612-0745-0_2.
- [56] D. Needell, N. Srebro, and R. Ward. Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm. *Math. Program.*, 155(1-2, Ser. A):549–573, 2016.

- [57] B. Neyshabur, R. Tomioka, and N. Srebro. In Search of the Real Inductive Bias: On the Role of Implicit Regularization in Deep Learning. *arXiv preprint arXiv:1412.6614*, 2014.
- [58] C. Paquette and E. Paquette. Dynamics of Stochastic Momentum Methods on Large-scale, Quadratic Models. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, 2021.
- [59] C. Paquette, K. Lee, F. Pedregosa, and E. Paquette. SGD in the Large: Average-case Analysis, Asymptotics, and Stepsize Criticality. In *Proceedings of Thirty Fourth Conference on Learning Theory (COLT)*, volume 134, pages 3548–3626, 2021.
- [60] C. Paquette, B. van Merriënboer, and F. Pedregosa. Halting Time is Predictable for Large Models: A Universality Property and Average-case Analysis. *Foundations of Computational Mathematics*, 2022. URL <https://doi.org/10.1007/s10208-022-09554-y>.
- [61] E. Paquette, C. Paquette, B. Adlam, and J. Pennington. Homogenization of SGD in high-dimensions: exact dynamics and generalization properties. *arXiv preprint arXiv:2205.07069*, 2022.
- [62] F. Pedregosa. A Hitchhiker’s Guide to Momentum, 2021. URL <http://fa.bianp.net/blog/2021/hitchhiker/>.
- [63] S. Pesme, L. Pillaud-Vivien, and N. Flammarion. Implicit bias of sgd for diagonal linear networks: a provable benefit of stochasticity. *Advances in Neural Information Processing Systems (NeurIPS)*, 34, 2021.
- [64] L. Pillaud-Vivien, A. Rudi, and F. Bach. Statistical Optimality of Stochastic Gradient Descent on Hard Learning Problems through Multiple Passes. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 31, 2018.
- [65] B.T. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 04, 1964.
- [66] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1177–1184, 2008.
- [67] S. Resnick. *Adventures in stochastic processes*. Birkhäuser Boston, Inc., Boston, MA, 1992.
- [68] H. Robbins and S. Monro. A Stochastic Approximation Method. *Ann. Math. Statist.*, 1951.
- [69] L. Sagun, L. Bottou, and Y. LeCun. Eigenvalues of the hessian in deep learning: Singularity and beyond. *arXiv preprint arXiv:1611.07476*, 2016.
- [70] L. Sagun, U. Evci, V. Guney, Y. Dauphin, and L. Bottou. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.
- [71] S. Shalev-Shwartz and T. Zhang. Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization. *Math. Program.*, 155(1-2, Ser. A):105–145, 2016. doi: 10.1007/s10107-014-0839-0.
- [72] C.J. Shallue, J. Lee, J. Antognini, J. Sohl-Dickstein, R. Frostig, and G.E. Dahl. Measuring the Effects of Data Parallelism on Neural Network Training. *Journal of Machine Learning Research*, 20:1–49, 2019.
- [73] V. Shankar, A. Fang, W. Guo, S. Fridovich-Keil, J. Ragan-Kelley, L. Schmidt, and B. Recht. Neural Kernels Without Tangents. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119, pages 8614–8623, 2020.
- [74] S. Smith and Q. Le. A bayesian perspective on generalization and stochastic gradient descent. In *International Conference on Learning Representations (ICLR)*, pages 1–13, 2018.
- [75] S. Smith, E. Elsen, and S. De. On the Generalization Benefit of Noise in Stochastic Gradient Descent. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119, pages 9058–9067, 2020.

- [76] S.L. Smith, B. Dherin, D. Barrett, and S. De. On the origin of implicit regularization in stochastic gradient descent. In *International Conference on Learning Representations (ICLR)*, pages 1–14, 2021.
- [77] N. Tripuraneni, B. Adlam, and J. Pennington. Covariate Shift in High-Dimensional Random Feature Regression. *arXiv preprint arXiv:2111.08234*, 2021.
- [78] A. Tsigler and P. L. Bartlett. Benign overfitting in ridge regression. *arXiv preprint arXiv:2009.14286*, 2020.
- [79] A. Vardhan Varre, L. Pillaud-Vivien, and N. Flammarion. Last iterate convergence of SGD for Least-Squares in the Interpolation regime. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pages 21581–21591, 2021.
- [80] S. Wojtowysch. Stochastic gradient descent with noise of machine learning type. Part II: Continuous time analysis. *arXiv preprint arXiv:2106.02588*, 2021.
- [81] J. Wu, D. Zou, V. Braverman, Q. Gu, and S. Kakade. Last Iterate Risk Bounds of SGD with Decaying Stepsize for Overparameterized Linear Regression. *arXiv preprint arXiv:2110.06198*, 2021.
- [82] L. Wu, C. Ma, and W. E. How SGD Selects the Global Minima in Over-parameterized Learning: A Dynamical Stability Perspective. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 31, 2018.
- [83] S. Yaida. Fluctuation-dissipation relations for stochastic gradient descent. In *International Conference on Learning Representations (ICLR)*, 2019.
- [84] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021. doi: 10.1145/3446776.
- [85] Z. Zhu, J. Wu, B. Yu, L. Wu, and J. Ma. The anisotropic noise in stochastic gradient descent: Its behavior of escaping from sharp minima and regularization effects. *arXiv preprint arXiv:1803.00195*, 2018.
- [86] L. Ziyin, K. Liu, T. Mori, and M. Ueda. Strength of Minibatch Noise in SGD. In *International Conference on Learning Representations (ICLR)*, 2022.
- [87] D. Zou, J. Wu, V. Braverman, Q. Gu, D.P. Foster, and S. Kakade. The Benefits of Implicit Regularization from SGD in Least Squares Problems (NeurIPS). In *Advances in Neural Information Processing Systems*, volume 34, pages 5456–5468, 2021.
- [88] D. Zou, J. Wu, V. Braverman, Q. Gu, and S. M. Kakade. Risk Bounds of Multi-Pass SGD for Least Squares in the Interpolation Regime. *arXiv preprint arXiv:2203.03159*, 2022.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]** We make sure to clearly state our assumptions particular that the results only hold for the ℓ^2 regularized least squares problem.
 - (b) Did you describe the limitations of your work? **[Yes]** We ensure the reader is aware that many of our results follow from another paper. We also make clear which results are proven and which results speculative, but have numerical support.
 - (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** We address the potential negative societal impacts in the appendix. We do not anticipate any ethical or societal issues. The results presented in this paper concern the analysis of existing methods.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]** The work presented is purely theoretical.

2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes] The full set of assumptions on the data matrix \mathbf{A} , signal β , initialization x_0 , and target vector \mathbf{b} are listed in the theorems statements with a complete description in Appendix B. Moreover we state our setting, namely that we are working on the ℓ^2 regularized least squares problem, clearly in the introduction of the main paper.
 - (b) Did you include complete proofs of all theoretical results? [Yes] All proofs of theoretical results are presented and, if not, referenced to where the proof can be found. We also clearly state conjectures which we believe should be true (see e.g., streaming SGD in Appendix C).
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] We provide code in our supplemental material and describe our experimental simulations in detail in Appendix G. We state all parameters in the captions of the figures.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] All parameter choices are clearly stated under each figure and in Appendix G. For example, we give the exact dimensions of the matrices used as well as the learning rates on SGD.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] Variance across the runs of SGD were reported.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] We indicated the type of resources in Appendix G. All our experiments were done using Colabs and only required a single standard GPU.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes] We use three data sets to verify our claims. All three data sets (MNIST [37], CIFAR-10 [34], and CIFAR-5m) are open source and available in TensorFlow. We cited the creators in the main text.
 - (b) Did you mention the license of the assets? [Yes] All three data sets used are open source and available on TensorFlow.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] We include the code used our simulations as part of the supplemental material.
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A] The data sets are open source.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A] The data sets used do not contain any personally identifiable information or offensive content. The data sets are standard data sets used in benchmarking.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] The work is purely theoretical and does not require any crowdsourcing or human subjects.
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] The work is purely theoretical and does not require any crowdsourcing or human subjects.
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] The work is purely theoretical and does not require any crowdsourcing or human subjects.