

LEARNING CURVES FOR STOCHASTIC GRADIENT DESCENT ON STRUCTURED FEATURES

Anonymous authors
Paper under double-blind review

ABSTRACT

The generalization performance of a machine learning algorithm such as a neural network depends in a non-trivial way on the structure of the data distribution. To analyze the influence of data structure on test loss dynamics, we study an exactly solvable model of stochastic gradient descent (SGD) which predicts test loss when training on features with arbitrary covariance structure. We solve the theory exactly for both Gaussian features and arbitrary features and we show that the simpler Gaussian model accurately predicts test loss of nonlinear random-feature models and deep neural networks trained with SGD on real datasets such as MNIST and CIFAR-10. We show that the optimal batch size at a fixed compute budget is typically small and depends on the feature correlation structure, demonstrating the computational benefits of SGD with small batch sizes. Lastly, we extend our theory to the more usual setting of stochastic gradient descent on a fixed subsampled training set, showing that both training and test error can be accurately predicted in our framework on real data.

1 INTRODUCTION

Due to the challenge of modeling the structure of realistic data, theoretical studies of generalization often attempt to derive data-agnostic generalization bounds or study the typical performance of the algorithm on simple data distributions. The first set of theories derive bounds based on the complexity or capacity of the function class and often struggle to explain the success of modern learning systems which generalize well on real data but are sufficiently powerful to fit random noise (Mohri et al., 2012; Zhang et al., 2017). Rather than exploring data-agnostic bounds, it is often useful to analyze how algorithms generalize typically or on average over a stipulated data distribution (Engel & Van den Broeck, 2001). In this style of analysis, the data distribution is usually assumed to be highly symmetric, stipulating that input data follows a factorized probability distribution across input variables (Advani et al., 2013). For example, *spherical cow* models treat data vectors as drawn from the isotropic Gaussian distribution or uniformly from the sphere while Boolean hypercube models treat data as random binary vectors. Rather than being distributed isotropically throughout the entire set of ambient dimensions, realistic datasets often lie on low dimensional structures embedded in high dimensional ambient spaces (Pope et al., 2021). For example, MNIST and CIFAR-10 lie on surfaces with intrinsic dimension of ~ 14 and ~ 35 respectively (Spigler et al., 2020). To understand the average-case performance of SGD in more realistic learning problems, incorporating structural information about realistic data distributions is necessary.

In this paper, we first explore the minimal improvement on the spherical cow approximation by studying an *elliptical cow* model, where the image of the data under a possibly nonlinear feature map is treated as a Gaussian with certain covariance. We express the generalization error in terms of the induced distribution of nonlinear features. Using the structure of the data, we study the evolution of the expected test loss during SGD. We derive test error dynamics throughout SGD in terms of the correlation structure in a feature space. We analyze SGD on random feature models and neural networks using MNIST and CIFAR-10 and accurately predict test loss scalings on these datasets.

We then analyze the general case where the feature distribution is arbitrary and provide an exact solution for the expected test loss dynamics. This result requires not only the second moment structure but also all of the fourth moments of the features. For MNIST and CIFAR-10, we empirically observe

that the Gaussian model provides an excellent approximation to the true dynamics due to negligible non-Gaussian effects.

We explore in detail the effect of minibatch size, m , on learning dynamics. By varying m , we can interpolate our theory between single sample SGD ($m = 1$) and gradient descent on the population loss ($m \rightarrow \infty$). To explore the computational advantages SGD compared to standard full batch gradient descent we analyze the loss achieved at a fixed compute budget $C = tm$ for different minibatch size m and number of steps t , trading off the number of parameter update steps for denoising through averaging. We show that generally, the optimal batch size is small, with the precise optimum dependent on the learning rate and structure of the features.

Overall, our theory shows how learning rate, minibatch size and data structure interact with the structure of the learning problem to determine generalization dynamics. It provides a predictive account of training dynamics in wide neural networks.

2 PROBLEM DEFINITION AND SETUP

We study stochastic gradient descent on a linear model with parameters \mathbf{w} and feature map $\psi(\mathbf{x}) \in \mathbb{R}^N$ (with N possibly infinite). Some interesting examples of linear models are random feature models, where $\psi(\mathbf{x}) = \phi(\mathbf{G}\mathbf{x})$ for random matrix \mathbf{G} and point-wise nonlinearity ϕ (Rahimi & Recht 2008; Mei & Montanari 2020). Another interesting linearized setting is wide neural networks with neural tangent kernel (NTK) parameterization (Jacot et al. 2020; Lee et al. 2020). Here the features are parameter gradients of the neural network function $\psi(\mathbf{x}) = \nabla_{\theta} f(\mathbf{x}, \theta)|_{\theta_0}$ at initialization. We will study both of these special cases in experiments.

We optimize the set of parameters \mathbf{w} by SGD to minimize a population loss of the form

$$L(\mathbf{w}) = \left\langle (\mathbf{w} \cdot \psi(\mathbf{x}) - y(\mathbf{x}))^2 \right\rangle_{\mathbf{x} \sim p(\mathbf{x})}, \quad (1)$$

where \mathbf{x} are input data vectors associated with a probability distribution $p(\mathbf{x})$, $\psi(\mathbf{x})$ is a nonlinear feature map and $y(\mathbf{x})$ is a target function which we can evaluate on training samples. We assume that the target function is square integrable $\langle y(\mathbf{x})^2 \rangle_{\mathbf{x}} < \infty$ over $p(\mathbf{x})$. Our aim is to elucidate how this population loss evolves during stochastic gradient descent on \mathbf{w} . We derive a formula in terms of the eigendecomposition of the feature correlation matrix and the target function

$$\Sigma = \langle \psi(\mathbf{x})\psi(\mathbf{x})^\top \rangle_{\mathbf{x}} = \sum_{k=1}^N \lambda_k \mathbf{u}_k \mathbf{u}_k^\top, \quad y(\mathbf{x}) = \sum_k v_k \mathbf{u}_k^\top \psi(\mathbf{x}) + y_{\perp}(\mathbf{x}), \quad (2)$$

where $\langle y_{\perp}(\mathbf{x})\psi(\mathbf{x}) \rangle = 0$. We justify this decomposition of $y(\mathbf{x})$ in the Appendix A using an eigendecomposition and show that it is general for target functions and features with finite variance.

During learning, parameters \mathbf{w} are updated to estimate a target function y which, as discussed above, can generally be expressed as a linear combination of features $y = \mathbf{w}^* \cdot \psi + y_{\perp}$. At each time step t , the weights are updated by taking a stochastic gradient step on a fresh mini-batch of m examples

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{m} \sum_{\mu=1}^m \psi_{t,\mu} (\mathbf{w}_t \cdot \psi_{t,\mu} - y_{t,\mu}), \quad (3)$$

where each of the vectors $\psi_{t,\mu}$ are sampled independently and $y_{t,\mu} = \mathbf{w}^* \cdot \psi_{t,\mu}$. The learning rate η controls the gradient descent step size while the batch size m gives a empirical estimate of the gradient at timestep t . At each timestep, the test-loss, or generalization error, has the form

$$L_t = \left\langle (\mathbf{w} \cdot \psi - \mathbf{w}^* \cdot \psi)^2 \right\rangle_{\psi} = (\mathbf{w}_t - \mathbf{w}^*)^\top \Sigma (\mathbf{w}_t - \mathbf{w}^*) + \langle y_{\perp}(\mathbf{x})^2 \rangle, \quad (4)$$

which quantifies exactly the test error of the vector \mathbf{w}_t . Note, however, that L_t is a random variable since \mathbf{w}_t depends on the precise history of sampled feature vectors $\mathcal{D}_t = \{\psi_{t,\mu}\}$. Our theory, which generalizes the recursive method of (Werfel et al. 2004) allows us to compute the *expected* test loss by averaging over all possible sequences, to obtain $\langle L_t \rangle_{\mathcal{D}_t}$. We show that our calculated learning curves are not limited to the one-pass setting, but rather can accommodate sampling minibatches from a finite training set with replacement and testing on a separate test set which we address in Section 3.4

In summary, we will develop a theory that predicts the expected test loss $\langle L_t \rangle_{\mathcal{D}_t}$ averaged over training sample sequences \mathcal{D}_t in terms of the quantities $\{\lambda_k, v_k, \langle y_{\perp}(\mathbf{x})^2 \rangle_{\mathbf{x}}\}$. This will reveal how the structure in the data and the learning problem influence test error dynamics during SGD. This theory is a quite general analysis of linear models on square loss, analyzing the performance of linearized models on arbitrary data distributions, feature maps ψ , and target functions $y(\mathbf{x})$.

3 ANALYTIC FORMULAE FOR LEARNING CURVES

3.1 LEARNABLE AND NOISE FREE PROBLEMS

Before studying the general case, we first analyze the setting where the target function is *learnable*, meaning that there exist weights \mathbf{w}^* such that $y(\mathbf{x}) = \mathbf{w}^* \cdot \psi(\mathbf{x})$. For many cases of interest, this is a reasonable assumption, especially when applying our theory to real datasets by fitting an atomic measure on P points $\frac{1}{P} \sum_{\mu} \delta(\mathbf{x} - \mathbf{x}^{\mu})$. We will further assume that the induced feature distribution is Gaussian so that all moments of ψ can be written in terms of the covariance Σ . We will remove these assumptions in later sections.

Theorem 3.1. *Suppose the features ψ follow a Gaussian distribution $\psi \sim \mathcal{N}(0, \Sigma)$ and the target function is learnable in these features $y = \mathbf{w}^* \cdot \psi$. After t steps of SGD with minibatch size m and learning rate η , the expected (over possible sample sequences \mathcal{D}_t) test loss $\langle L_t \rangle_{\mathcal{D}_t}$ has the form*

$$\langle L_t \rangle_{\mathcal{D}_t} = \boldsymbol{\lambda}^{\top} \mathbf{A}^t \mathbf{v}^2, \quad \mathbf{A} = (\mathbf{I} - \eta \text{diag}(\boldsymbol{\lambda}))^2 + \frac{\eta^2}{m} \text{diag}(\boldsymbol{\lambda}^2) + \frac{\eta^2}{m} \boldsymbol{\lambda} \boldsymbol{\lambda}^{\top} \quad (5)$$

where $\boldsymbol{\lambda}$ is a vector containing the eigenvalues of Σ and \mathbf{v}^2 is a vector containing elements $(\mathbf{v}^2)_k = v_k^2 = (\mathbf{u}_k \cdot \mathbf{w}^*)^2$ for eigenvectors \mathbf{u}_k of Σ . The function $\text{diag}(\cdot)$ constructs a diagonal matrix with the argument vector placed along the diagonal.

Proof. See Appendix B for the full derivation. We will provide a brief sketch of the proof here. The strategy of the proof relies on the fact that $\langle L_t \rangle = \text{Tr} \Sigma \mathbf{C}_t$ where $\mathbf{C}_t = \langle (\mathbf{w}_t - \mathbf{w}^*)(\mathbf{w}_t - \mathbf{w}^*)^{\top} \rangle_{\mathcal{D}_t}$. We derive the following recursion relation for this error matrix

$$\mathbf{C}_{t+1} = (\mathbf{I} - \eta \Sigma) \mathbf{C}_t (\mathbf{I} - \eta \Sigma) + \frac{\eta^2}{m} [\Sigma \mathbf{C}_t \Sigma + \Sigma \text{Tr}(\Sigma \mathbf{C}_t)] \quad (6)$$

The loss only depends on the quantities $c_{k,t} = \mathbf{u}_k^{\top} \mathbf{C}_t \mathbf{u}_k$. As a vector the recurrence for these quantities can be solved $\mathbf{c}_t = \mathbf{A}^t \mathbf{v}^2$. Using the fact that $L_t = \sum_k \lambda_k \mathbf{u}_k^{\top} \mathbf{C}_t \mathbf{u}_k = \sum_k c_{k,t} \lambda_k = \boldsymbol{\lambda}^{\top} \mathbf{A}^t \mathbf{v}^2$, we obtain the desired result. \square

Below we provide some immediate interpretations of this result.

- The matrix \mathbf{A} contains two components; a matrix $(\mathbf{I} - \eta \text{diag}(\boldsymbol{\lambda}))^2$ which represents the time-evolution of the loss under *average gradient updates*. The remaining matrix $\frac{\eta^2}{m} (\text{diag}(\boldsymbol{\lambda}^2) + \boldsymbol{\lambda} \boldsymbol{\lambda}^{\top})$ arises due to fluctuations in the gradients, a consequence of the stochastic sampling process.
- The test loss obtained when training directly on the population loss can be obtained by taking the minibatch size $m \rightarrow \infty$. In this case, $\mathbf{A} \rightarrow (\mathbf{I} - \eta \text{diag}(\boldsymbol{\lambda}))^2$ and one obtains the population loss $L_t^{\text{pop}} = \sum_k v_k^2 \lambda_k (1 - \eta \lambda_k)^{2t}$. This population loss can also be obtained by considering small learning rates, i.e. the $\eta \rightarrow 0$ limit, where $\mathbf{A} = (\mathbf{I} - \eta \text{diag}(\boldsymbol{\lambda}))^2 + O(\eta^2)$.
- For general $\boldsymbol{\lambda}$ and $\eta^2/m > 0$, \mathbf{A} is non-diagonal, indicating that the components $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$ are not learned independently as t increases like for L_t^{pop} , but rather interact during learning due to non-trivial coupling across eigenmodes at large η^2/m . This is unlike offline theory for learning in feature spaces (kernel regression) where errors across eigenmodes were shown to decouple and are learned at different rates (Bordelon et al., 2020; Canatar et al., 2020). This observation of mixing across covariance eigenspaces agrees with a recent analysis of SGD, which introduced recursively defined “mixing terms” that couple each mode’s evolution (Varre et al., 2021).
- Though increasing m always improves generalization at fixed time t (proof given in Appendix D), learning with a fixed compute budget (number of gradient evaluations) $C = tm$, can favor smaller batch sizes. We provide an example of this in the next sections and Figure 1(d)-(f).

- The lower bound $\langle L_t \rangle \geq \lambda^\top v^2 \left[(1 - \eta)^2 + \frac{\eta^2}{m} |\lambda|^2 \right]^t$ can be used to find necessary stability conditions on m, η . This bound implies that $\langle L_t \rangle$ will diverge if $m < \frac{\eta |\lambda|^2}{2 - \eta}$. The learning rate must be sufficiently small and the batch size sufficiently large to guarantee convergence. This stability condition depends on the features through $|\lambda|^2 = \sum_k \lambda_k^2$. One can derive heuristic optimal batch sizes and optimal learning rates through this lower bound. See Figure 2 and Appendix C.

3.1.1 SPECIAL CASE 1: UNSTRUCTURED ISOTROPIC FEATURES

This special case was previously analyzed by Werfel et al. (2004) which takes $\Sigma = \mathbf{I} \in \mathbb{R}^{N \times N}$ and $m = 1$. We extend their result for arbitrary m , giving the following learning curve

$$\langle L_t \rangle_{\mathcal{D}_t} = \left((1 - \eta)^2 + \frac{1 + N}{m} \eta^2 \right)^t \|\mathbf{w}^*\|^2, \quad \langle L_t^* \rangle_{\mathcal{D}_t} = \left(1 - \frac{m}{m + N + 1} \right)^t \|\mathbf{w}^*\|^2, \quad (7)$$

where the second expression has optimal η . First, we note the strong dependence on the ambient dimension N : as $N \gg m$, learning happens at a rate $\langle L_t \rangle \sim e^{-tm/N}$. Increasing the minibatch size m improves the exponential rate by reducing the gradient noise variance. Second, we note that this feature model has the same rate of convergence for every learnable target function y . At small m , the convergence at any learning rate η is much slower than the convergence of the $m \rightarrow \infty$ limit, $L_{pop} = (1 - \eta)^{2t} \|\mathbf{w}^*\|^2$ which does not suffer from a dimensionality dependence due to gradient noise. Lastly, for a fixed compute budget $C = tm$, the optimal batch size is $m^* = 1$; see Figure 1(d). This can be shown by differentiating $\langle L_{C/m} \rangle$ with respect to m (see Appendix E). In Figure 1(a) we show theoretical and simulated learning curves for this model for varying values of N at the optimal learning rate and in Figure 1(d), we show the loss as a function of minibatch size for a fixed compute budget $C = tm = 100$.

3.1.2 SPECIAL CASE 2: POWER LAWS AND EFFECTIVE DIMENSIONALITY

Realistic datasets such as natural images or audio tend to exhibit nontrivial correlation structure, which often results in power-law spectra when the data is projected into a feature space, such as a randomly initialized neural network (Spigler et al., 2020; Canatar et al., 2020; Bahri et al., 2021). In the $\frac{\eta^2}{m} \ll 1$ limit, if the feature spectra and task spectra follow power laws, $\lambda_k \sim k^{-b}$ and $\lambda_k v_k^2 \sim k^{-a}$ with $a, b > 1$, then Theorem 3.1 implies that generalization error also falls with a power law: $\langle L_t \rangle \sim Ct^{-\beta}$, $\beta = \frac{a-1}{b}$ where C is a constant. See Appendix G for a derivation. Notably, these predicted exponents we recovered as a special case of our theory agree with prior work on SGD with power law spectra, which give exponents in terms of the feature correlation structure (Berthier et al., 2020; Dieuleveut et al., 2016; Velikanov & Yarotsky, 2021; Varre et al., 2021). Further, our power law scaling appears to accurately match the qualitative behavior of wide neural networks trained on realistic data (Hestness et al., 2017; Bahri et al., 2021), which we study in Section 4.

We show an example of such a power law scaling with synthetic features in Figure 1(b). Since the total variance approaches a finite value as $N \rightarrow \infty$, the learning curves are relatively insensitive to N , and are rather sensitive to the eigenspectrum through terms like $|\lambda|^2$ and $\mathbf{1}^\top \lambda$, etc. In Figure 1(c), we see that the scaling of the loss is more similar to the power law setting than the isotropic features setting in a random features model of MNIST, agreeing excellently with our theory.

For this model, we find that there can exist optimal batch sizes when the compute budget $C = tm$ is fixed (Figure 1(e) and (f)). In Appendix C.1 we heuristically argue that the optimal batch size for power law features should scale as, $m^* \approx \frac{\eta^2}{(2b-1)(1-\eta)^2}$. Figure 2 shows a test of this result and related experiments.

We provide further evidence of the existence of power law structure on realistic data in Figure 3 (a)-(c), where we provide spectra and test loss learning curves for MNIST and CIFAR-10 on ReLU random features. The eigenvalues $\lambda_k \sim k^{-b}$ and the task power tail sums $\sum_{n=k}^\infty \lambda_n v_n^2 \sim k^{-a+1}$ both follow power laws, generating power law test loss curves. These learning curves are contrasted with isotropically distributed data in \mathbb{R}^{784} passed through the same ReLU random feature model and we see that structured data distributions allow much faster learning than the unstructured data. Again, our theory predicts experimental curves accurately across variations in learning rate, batch size and noise (Figure 3).

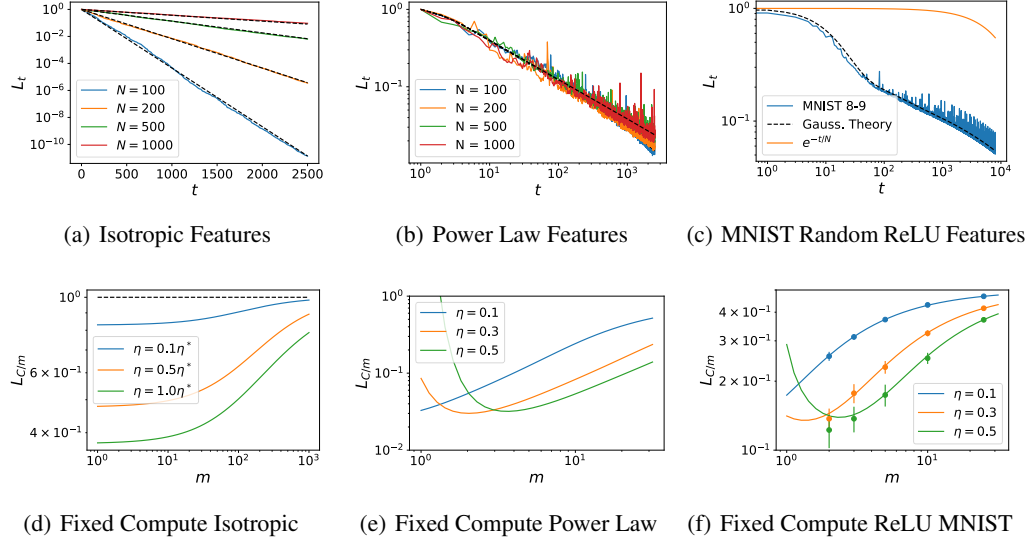


Figure 1: Isotropic features generated as $\psi \sim \mathcal{N}(0, \mathbf{I})$ have qualitatively different learning curves than power-law features observed in real data. Black dashed lines are theory. (a) Online learning with N -dimensional isotropic features gives a test loss which scales like $L_t \sim e^{-t/N}$ for any target function, indicating that learning requires $t \sim N$ steps of SGD, using the optimal learning rates $\eta^* = \frac{m}{N+m+1}$. (b) Power-law features $\psi \sim \mathcal{N}(0, \Lambda)$ with $\Lambda_{kl} = \delta_{k,l} k^{-2}$ have non-extensive give a power-law scaling $L_t \sim t^{-\beta}$ with exponent $\beta = O_N(1)$. (c) Learning to discriminate MNIST 8’s and 9’s with $N = 4000$ dimensional random ReLU features (Rahimi & Recht, 2008), generates a power law scaling at large t , which is both quantitatively and qualitatively different than the scaling predicted by isotropic features $e^{-t/N}$. (d)-(f) The loss at a fixed compute budget $C = tm = 100$ for (d) isotropic features, (e) power law features and (f) MNIST ReLU random features with simulations (dots average and standard deviation for 30 runs). Intermediate batch sizes are preferable on real data.

3.2 ARBITRARY INDUCED FEATURE DISTRIBUTIONS: THE GENERAL SOLUTION

The result in the previous section was proven exactly in the case of Gaussian vectors (see Appendix B). For arbitrary (possibly non-Gaussian) distributions, we obtain a slightly more involved result (see Appendix F).

Theorem 3.2. Let $\psi(\mathbf{x}) \in \mathbb{R}^N$ be an arbitrary feature map with covariance matrix $\Sigma = \sum_k \lambda_k \mathbf{u}_k \mathbf{u}_k^\top$. After diagonalizing the features $\phi_k(\mathbf{x}) = \mathbf{u}_k^\top \psi(\mathbf{x})$, introduce the fourth moment tensor $\kappa_{ijkl}^4 = \langle \phi_i \phi_j \phi_k \phi_l \rangle$. The expected loss is exactly $\langle L_t \rangle = \sum_k \lambda_k c_k(\lambda, \kappa, \mathbf{v})$.

We provide an exact formula for c_k in the Appendix F. We see that the test loss dynamics depends only on the second and fourth moments of the features through quantities λ_k and κ_{ijkl} respectively. We recover the Gaussian result as a special case when κ_{ijkl} is a simple weighted sum of these three products of Kronecker tensors $\kappa_{ijkl}^{Gauss} = \lambda_i \lambda_j \delta_{ik} \delta_{jl} + \lambda_i \lambda_k \delta_{ij} \delta_{kl} + \lambda_i \lambda_j \delta_{il} \delta_{jk}$. As an alternative to the above closed form expression for $\langle L_t \rangle$, a recursive formula which tracks N mixing coefficients has also recently been utilized to analyze the test loss dynamics for arbitrary distributions (Varre et al., 2021).

Next we show that a mild regularity condition, similar to those assumed in other recent works (Berthier et al., 2020; Varre et al., 2021), on the fourth moment structure of the features allows derivation of an upper bound which is qualitatively similar to the Gaussian theory.

Theorem 3.3. If the fourth moments satisfy $\langle \psi \psi^\top \mathbf{G} \psi \psi^\top \rangle \preceq (\alpha + 1) \Sigma \mathbf{G} \Sigma + \alpha \text{Tr} \Sigma \mathbf{G}$ for any positive-semidefinite \mathbf{G} , then

$$L_t \leq \lambda^\top \mathbf{A}^t \mathbf{v}^2, \quad \mathbf{A} = (\mathbf{I} - \eta \text{diag}(\lambda))^2 + \frac{\alpha \eta^2}{m} [\text{diag}(\lambda^2) + \lambda \lambda^\top]. \quad (8)$$

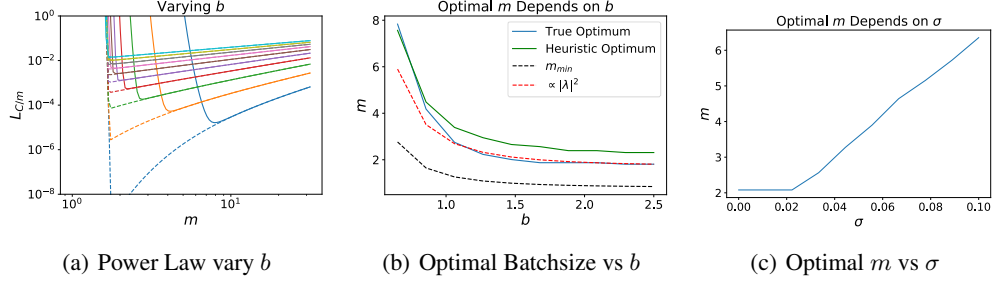


Figure 2: Optimal batch size depends on feature structure and noise level. (a) For power law features $\lambda_k \sim k^{-b}$, $\lambda_k v_k^2 \sim k^{-a}$, the m dependence of the loss $L_{C/m}$ depends strongly on the feature exponent b . Each color is a different b value, evenly spaced in $[0.6, 2.5]$ with $a = 2.5$, $C = 500$. The solid lines show the exact theory while dashed lines show the error predicted by first order perturbation theory where the mode coupling term $\frac{\eta^2}{m} \lambda \lambda^\top$ is replaced with a decoupled term $\frac{\eta^2}{m} \text{diag}(\lambda^2)$. This shows that mode coupling is necessary to accurately predict optimal m . (b) The optimal m scales proportionally with $|\lambda|^2 \approx \frac{1}{2b-1}$. We plot the lower bound m_{\min} (black), the heuristic optimum (m which optimizes the lower bound for L , shown in green) and $\frac{\eta^2}{(1-\eta)^2} |\lambda|^2$ (red). (c) In noisy problems, the optimal batch size also depends on noise level σ , scaling roughly as $m^* \propto \sigma$.

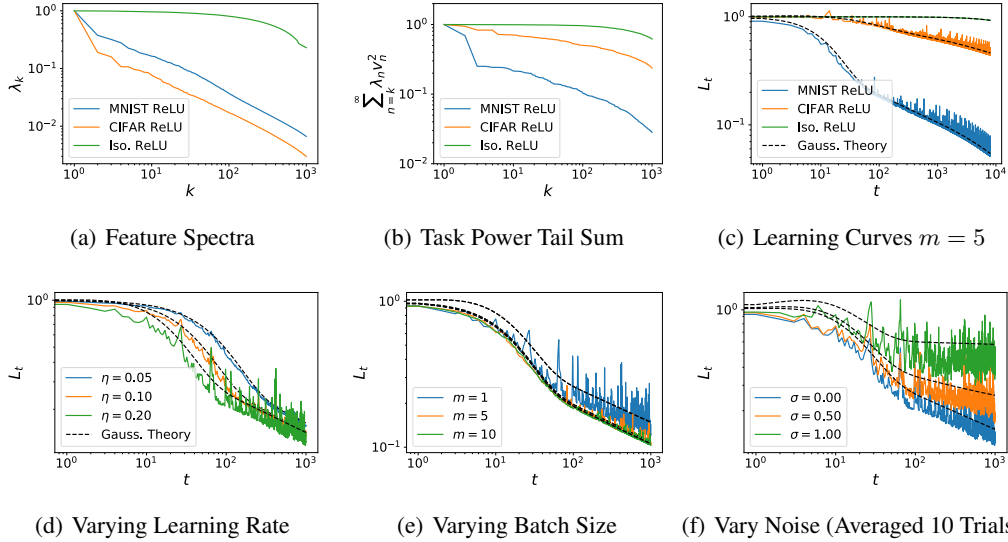


Figure 3: Structure in the data distribution, nonlinearity, batchsize and learning rate all influence learning curves. (a) ReLU random feature embedding in $N = 4000$ dimensions of MNIST and CIFAR images have very different eigenvalue scalings than spherically isotropic vectors in 784 dimensions. (b) The task power spectrum decays much faster for MNIST than for random isotropic vectors. (c) Learning curves reveal the data-structure dependence of test error dynamics. Dashed lines are theory curves derived from equation. (d) Increasing the learning rate increases the initial speed of learning but induces large fluctuations in the loss and can be worse at large t . (e) Increasing the batch size alters both the average test loss L_t and the variance. (f) Noise in the target values during training produces an asymptotic error L_∞ which persists even as $t \rightarrow \infty$.

We provide this proof in Appendix [F.1](#). We note that the assumed bound on the fourth moments is tight for Gaussian features with $\alpha = 1$, recovering our previous theory. Thus, if this condition on the fourth moments is satisfied, then the loss for the non-Gaussian features is upper bounded by the Gaussian test loss theory with the batch size effectively altered $\tilde{m} = m/\alpha$.

The question remains whether the Gaussian approximation will provide an accurate model on *realistic data*. We do not provide a proof of this conjecture, but verify its accuracy in empirical experiments on MNIST and CIFAR-10 as shown in Figure 3. In Appendix Figure F.1, we show that the fourth moment matrix for a ReLU random feature model and its projection along the eigenbasis of the feature covariance is accurately approximated by the equivalent Gaussian model.

3.3 UNLEARNABLE OR NOISE CORRUPTED PROBLEMS

In general, the target function $y(\mathbf{x})$ may depend on features which cannot be expressed as linear combinations of features $\psi(\mathbf{x})$, $y(\mathbf{x}) = \mathbf{w}^* \cdot \psi(\mathbf{x}) + y_\perp(\mathbf{x})$. Let $\langle y_\perp(\mathbf{x})^2 \rangle_{\mathbf{x}} = \sigma^2$. Note that y_\perp does not have to be a deterministic function of \mathbf{x} , but can also be a stochastic process which is uncorrelated with $\psi(\mathbf{x})$.

Theorem 3.4. *For a target function with unlearnable variance $\langle y_\perp^2 \rangle = \sigma^2$, the expected test loss has the form*

$$\langle L_t \rangle = \boldsymbol{\lambda}^\top \mathbf{A}^t \mathbf{v}^2 + \frac{1}{m} \eta^2 \sigma^2 \boldsymbol{\lambda}^\top (\mathbf{I} - \mathbf{A})^{-1} (\mathbf{I} - \mathbf{A}^t) \boldsymbol{\lambda} \quad (9)$$

which has an asymptotic, irreducible error $\langle L_\infty \rangle = \frac{1}{m} \eta^2 \sigma^2 \boldsymbol{\lambda}^\top (\mathbf{I} - \mathbf{A})^{-1} \boldsymbol{\lambda}$ as $t \rightarrow \infty$.

See Appendix H for the proof. The convergence to the asymptotic error takes the form $\langle L_t - L_\infty \rangle = \boldsymbol{\lambda}^\top \mathbf{A}^t (\mathbf{v}^2 - \frac{1}{m} \eta^2 \sigma^2 (\mathbf{I} - \mathbf{A})^{-1} \boldsymbol{\lambda})$. We note that this quantity is not necessarily monotonic in t and can exhibit local maxima for sufficiently large σ^2 , as in Figure 3(f).

3.4 TEST/TRAIN SPLITS

Rather than interpreting our theory as a description of the average test loss during SGD in a one-pass setting, where data points are sampled from the a distribution at each step of SGD, our theory can be suitably modified to accommodate multiple random passes over a finite training set. To accomplish this, one must first recognize that the training and test distributions are different.

Theorem 3.5. *Let $\hat{p}(\mathbf{x}) = \frac{1}{M} \sum_{\mu} \delta(\mathbf{x} - \mathbf{x}^\mu)$ be the empirical distribution on the M training data points and let $\hat{\Sigma} = \langle \psi(\mathbf{x}) \psi(\mathbf{x})^\top \rangle_{\mathbf{x} \sim \hat{p}(\mathbf{x})} = \sum_k \hat{\lambda}_k \mathbf{u}_k \mathbf{u}_k^\top$ be the induced feature correlation matrix on this training set. Further, let $p(\mathbf{x})$ be the test distribution Σ its corresponding feature correlation. Then we have*

$$\begin{aligned} \langle L_{\text{train}} \rangle &= \text{Tr} [\hat{\Sigma} \mathbf{C}_t] , \quad \langle L_{\text{test}} \rangle = \text{Tr} [\Sigma \mathbf{C}_t] \\ \mathbf{u}_k^\top \mathbf{C}_t \mathbf{u}_k &= [\mathbf{A}^t \mathbf{v}^2]_k , \quad \mathbf{u}_k^\top \mathbf{C}_t \mathbf{u}_\ell = \left(1 - \eta \hat{\lambda}_k - \eta \hat{\lambda}_\ell + \eta^2 \left(1 + \frac{1}{m} \right) \hat{\lambda}_k \hat{\lambda}_\ell \right)^t v_k v_\ell \end{aligned} \quad (10)$$

where $\mathbf{A} = \left((\mathbf{I} - \text{diag}(\hat{\lambda}))^2 + \frac{\eta^2}{m} [\text{diag}(\hat{\lambda}^2) + \hat{\lambda} \hat{\lambda}^\top] \right)$.

We provide the proof of this theorem in Appendix I. The interpretation of this result is that it provides the expected training and test loss if, at each step of SGD, m points from the training set $\{\mathbf{x}^1, \dots, \mathbf{x}^M\}$ are sampled uniformly with replacement and used to calculate a stochastic gradient. Note that while Σ can be full rank, the rank of $\hat{\Sigma}$ has rank upper bounded by M , the number of training samples. As a consequence learning will only occur along the M dimensional subspace spanned by the data. Thus, the test error will have an irreducible component at large time, as evidenced in Figure 4. While the training errors continue to go to zero, the test errors saturate at a M -dependent final loss. This result can also allow one to predict errors on other test distributions.

4 COMPARING NEURAL NETWORK FEATURE MAPS

We can utilize our theory to compare how wide neural networks of different depths generalize when trained with SGD on a real dataset. With a certain initialization of the model parameters, infinite width networks behave as linear functions of their parameters $f(\mathbf{x}, \boldsymbol{\theta}) \approx f(\mathbf{x}, \boldsymbol{\theta}_0) + \nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}_0) \cdot (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$ (Lee et al., 2020). To predict test loss dynamics with our theory, it therefore suffices to characterize the geometry of the gradient features $\psi(\mathbf{x}) = \nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta})$. In Figure 5, we show the Neural Tangent

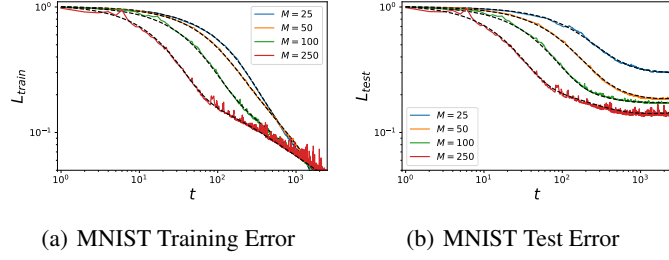


Figure 4: Training and test errors of a model trained on a training set of size M can be computed with the C_t matrix. Dashed black lines are theory. (a) The training error for MNIST random feature model approaches zero asymptotically. (b) The test error saturates to a quantity dependent on M .

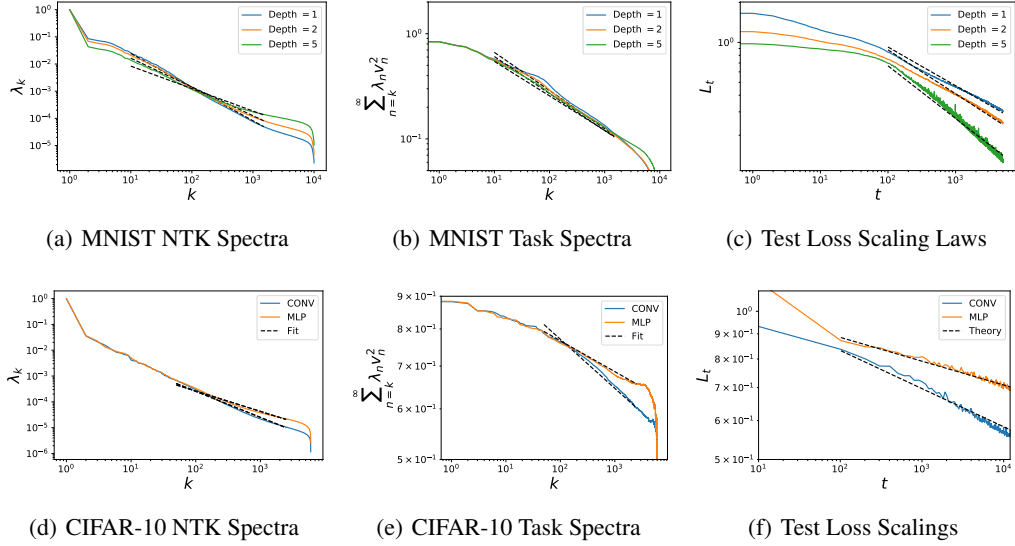


Figure 5: ReLU neural networks of depth D and width 500 are trained with SGD on full MNIST. (a)-(b) Feature and spectra are estimated by diagonalizing the infinite width NTK matrix on the training data. We fit a simple power law to each of the curves $\lambda_k \sim k^{-b}$ and $v_k^2 \sim k^{-a}$. (c) Experimental test loss during SGD (color) compared to theoretical power-law scalings $t^{-\frac{a-1}{b}}$ (dashed color). Deeper networks train faster due to their slower decay in their feature eigenspectra λ_k , though they have similar task spectra. (d)-(f) The spectra and test loss for convolutional and fully connected networks on CIFAR-10. The CNN obtains a better convergence exponent due to its faster decaying task spectra. The predicted test loss scalings (dashed black) match those observed in experiments (color).

Kernel (NTK) eigenspectra and task-power spectra for fully connected neural networks of varying depth, calculated with the Neural Tangents API (Novak et al., 2020). We compute the kernel on a subset of 10,000 randomly sampled MNIST images and estimate the power law exponents for the kernel and task spectra λ_k and v_k^2 . We find that, accross architectures, the task spectra v_k^2 are highly similar, but that the kernel eigenvalues λ_k decay more slowly for deeper models, corresponding to a smaller exponent b . As a consequence, deeper neural network models train more quickly during stochastic gradient descent as we show in Figure 5 (c). After fitting power laws to the spectra $\lambda_k \sim k^{-b}$ and the task power $v_k^2 \sim k^{-a}$, we compared the true test loss dynamics (color) for a width-500 neural network model with the predicted power-law scalings $\beta = \frac{a-1}{b}$ from the fit exponents a, b . The predicted scalings from NTK regression accurately describe trained networks at finite width. On CIFAR-10, we compare the scalings of the CNN model and a standard MLP and find that the CNN obtains a better exponent due to its faster decaying tail sum $\sum_{n=k}^{\infty} \lambda_n v_n^2$.

5 RELATED WORK AND DISCUSSION

The analysis of stochastic gradient descent has a long history dating back to seminal works of Polyak & Juditsky (1992); Ruppert (1988), who showed that the time-averaged iterates provably converge to the optimum in noisy convex problems with the optimal minimax rate of $O(\sigma^2 t^{-1})$ in the large t limit. Many more works have examined a similar setting, identifying how averaged and accelerated versions of SGD perform asymptotically (Flammarion & Bach, 2015; 2017; Jain et al., 2018; Duchi & Ruan, 2021; Shapiro, 1989; Robbins & Monro, 1951; Chung, 1954; Duchi & Ruan, 2021; Yu et al., 2020; Anastasiou et al., 2019; Gurbuzbalaban et al., 2020).

Recent studies have also analyzed the asymptotics of noise-free MSE problems with arbitrary feature structure. A series of works have established that, under mild regularity conditions of the features, power law test loss curves emerge with exponents which are better than the $1/t$ rates in the noisy problem (Berthier et al., 2020; Pillaud-Vivien et al., 2018; Dieuleveut et al., 2016; Varre et al., 2021; Dieuleveut & Bach, 2016; Ying & Pontil, 2008; Fischer & Steinwart, 2020). Their predicted exponent in the small learning rate limit agrees with the exponent we derive with the saddle point approximation in Section 3.1.2. In the most recent of these studies, (Varre et al., 2021) show that the error signals across covariance eigenspaces mix at finite learning rate, an observation which also occurs in our theory since the derived \mathbf{A} matrix is non-diagonal. This fact can result in qualitatively different behavior (such as different optimal hyperparameters) than what arises if one assumes the gradient covariance and feature covariance are simultaneously diagonalizable (Zhang et al., 2019). Our work improves upon these prior results by examining the effect of minibatch size, test train splits, and providing exact expressions for both Gaussian and arbitrary features.

Several famous works have analyzed average case online learning in shallow and two-layer neural networks for unstructured data (Heskes & Kappen, 1991; Biehl & Riegler, 1994; Mace & Coolen, 1998; Saad & Solla, 1999; Cun et al., 1991; Goldt et al., 2019). However, a more recent analysis has demonstrated that structured data has significant influence in the two-layer student-teacher setting (Goldt et al., 2020) and for regression generalization (Mel & Ganguli, 2021). While we focus on discrete time in this work, some other recent works have analyzed shallow classification in a continuous time dynamical field theory approach (Mignacco et al., 2020).

Studying the simple setting of least squares regression on isotropic Gaussian features, Werfel et al. (2004) computed average case learning curves for SGD with minibatch size of one. Their results are non-asymptotic and exact, though the assumptions on the features are highly restrictive. We generalize their result and method so that it can describe structured features and arbitrary batch sizes.

Understanding the computational benefit that SGD provides over full-batch gradient descent requires understanding how test loss dynamics depend on batch-size m . Ma et al. (2018) study the tradeoff between taking many steps of gradient descent at small m and taking a small number of steps at large m . They show that for small m , doubling the batch size and cutting in half the number of steps give roughly the same loss. After a critical m , however, they observe a saturation effect where making m larger does not reduce the loss significantly since denoising estimated gradients provides diminishing returns. Our results improve upon this initial study since we provide an exact analysis of SGD with varying batch size at fixed compute budget and show how the optimal batch sizes depend on the feature covariance.

6 CONCLUSION

By studying a simple model of stochastic gradient descent, we were able to uncover how the geometry of the data in an induced feature space governs the dynamics of the test loss. We derived average learning curves $\langle L_t \rangle$ for both Gaussian and general non-Gaussian features and showed the conditions under which the Gaussian approximation is accurate. The proposed model allowed us to explore the role of the data distribution and neural network architecture on the learning curves, demonstrating how the power-law spectra observed in wide neural networks on real data allow an escape of the curse of dimensionality during SGD. We verified our theory with experiments on MNIST and CIFAR-10. In addition, we explored the role of batch size, learning rate, and label noise level on generalization. We found that for a fixed compute budget small minibatch sizes give the lowest expected loss, providing a quantitative demonstration of a benefit of SGD over large batch gradient descent.

REFERENCES

- Madhu Advani, Subhaneil Lahiri, and Surya Ganguli. Statistical mechanics of complex neural systems and high dimensional data. *Journal of Statistical Mechanics: Theory and Experiment*, 2013(03):P03014, Mar 2013. ISSN 1742-5468. doi: 10.1088/1742-5468/2013/03/p03014. URL <http://dx.doi.org/10.1088/1742-5468/2013/03/P03014>.
- Andreas Anastasiou, Krishnakumar Balasubramanian, and Murat A. Erdogdu. Normal approximation for stochastic gradient descent via non-asymptotic rates of martingale clt. In Alina Beygelzimer and Daniel Hsu (eds.), *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pp. 115–137, Phoenix, USA, 25–28 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v99/anastasiou19a.html>.
- Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural scaling laws. *arXiv preprint arXiv:2102.06701*, 2021.
- Carl Bender and Steven Orszag. *Advanced Mathematical Methods for Scientists and Engineers: Asymptotic Methods and Perturbation Theory*, volume 1. 01 1999. ISBN 978-1-4419-3187-0. doi: 10.1007/978-1-4757-3069-2.
- Raphaël Berthier, Francis Bach, and Pierre Gaillard. Tight nonparametric convergence rates for stochastic gradient descent under the noiseless linear model, 2020.
- M Biehl and P Riegler. On-line learning with a perceptron. *Europhysics Letters (EPL)*, 28(7): 525–530, dec 1994. doi: 10.1209/0295-5075/28/7/012. URL <https://doi.org/10.1209/0295-5075/28/7/012>.
- Blake Bordelon, Abdulkadir Canatar, and Cengiz Pehlevan. Spectrum dependent learning curves in kernel regression and wide neural networks. In *International Conference on Machine Learning*, pp. 1024–1034. PMLR, 2020.
- Abdulkadir Canatar, B. Bordelon, and C. Pehlevan. Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks. *Nature Communications*, 12:1–12, 2020.
- K. L. Chung. On a Stochastic Approximation Method. *The Annals of Mathematical Statistics*, 25 (3):463 – 483, 1954. doi: 10.1214/aoms/1177728716. URL <https://doi.org/10.1214/aoms/1177728716>.
- Yann Le Cun, Ido Kanter, and Sara A. Solla. Eigenvalues of covariance matrices: Application to neural-network learning. *Phys. Rev. Lett.*, 66:2396–2399, May 1991. doi: 10.1103/PhysRevLett.66.2396. URL <https://link.aps.org/doi/10.1103/PhysRevLett.66.2396>.
- Aymeric Dieuleveut and Francis Bach. Nonparametric stochastic approximation with large step-sizes. *The Annals of Statistics*, 44(4):1363 – 1399, 2016. doi: 10.1214/15-AOS1391. URL <https://doi.org/10.1214/15-AOS1391>.
- Aymeric Dieuleveut, Nicolas Flammarion, and Francis Bach. Harder, better, faster, stronger convergence rates for least-squares regression. *Journal of Machine Learning Research*, 18, 02 2016.
- John C. Duchi and Feng Ruan. Asymptotic optimality in stochastic optimization. *The Annals of Statistics*, 49(1):21 – 48, 2021. doi: 10.1214/19-AOS1831. URL <https://doi.org/10.1214/19-AOS1831>.
- A. Engel and C. Van den Broeck. *Statistical Mechanics of Learning*. Cambridge University Press, 2001. doi: 10.1017/CBO9781139164542.
- Simon Fischer and Ingo Steinwart. Sobolev norm learning rates for regularized least-squares algorithms. *Journal of Machine Learning Research*, 21(205):1–38, 2020. URL <http://jmlr.org/papers/v21/19-734.html>.
- Nicolas Flammarion and Francis Bach. From averaging to acceleration, there is only a step-size, 2015.

- Nicolas Flammarion and Francis Bach. Stochastic composite least-squares regression with convergence rate $o(1/n)$. In Satyen Kale and Ohad Shamir (eds.), *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pp. 831–875. PMLR, 07–10 Jul 2017. URL <https://proceedings.mlr.press/v65/flammarion17a.html>.
- Sebastian Goldt, Madhu Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborová. Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/cab070d53bd0d200746fb852a922064a-Paper.pdf>.
- Sebastian Goldt, Marc Mézard, Florent Krzakala, and Lenka Zdeborová. Modeling the influence of data structure on learning in neural networks: The hidden manifold model. *Phys. Rev. X*, 10: 041044, Dec 2020. doi: 10.1103/PhysRevX.10.041044. URL <https://link.aps.org/doi/10.1103/PhysRevX.10.041044>.
- Mert Gurbuzbalaban, Umut Simsekli, and Lingjiong Zhu. The heavy-tail phenomenon in sgd. *arXiv preprint arXiv:2006.04740*, 2020.
- Tom M. Heskes and Bert Kappen. Learning processes in neural networks. *Phys. Rev. A*, 44:2718–2726, Aug 1991. doi: 10.1103/PhysRevA.44.2718. URL <https://link.aps.org/doi/10.1103/PhysRevA.44.2718>.
- Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md. Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable, empirically, 2017.
- Arthur Jacot, Berfin Şimşek, Francesco Spadaro, Clément Hongler, and Franck Gabriel. Kernel alignment risk estimator: Risk prediction from training data, 2020.
- Prateek Jain, S. Kakade, Rahul Kidambi, Praneeth Netrapalli, and Aaron Sidford. Accelerating stochastic gradient descent for least squares regression. In *COLT*, 2018.
- Jaehoon Lee, Lechao Xiao, Samuel S Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(12): 124002, Dec 2020. ISSN 1742-5468. doi: 10.1088/1742-5468/abc62b. URL <http://dx.doi.org/10.1088/1742-5468/abc62b>.
- Siyuan Ma, Raef Bassily, and Mikhail Belkin. The power of interpolation: Understanding the effectiveness of sgd in modern over-parametrized learning. In *ICML*, pp. 3331–3340, 2018. URL <http://proceedings.mlr.press/v80/ma18a.html>.
- C. Mace and A. Coolen. Statistical mechanical analysis of the dynamics of learning in perceptrons. *Statistics and Computing*, 8:55–88, 1998.
- Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve, 2020.
- Gabriel Mel and Surya Ganguli. A theory of high dimensional regression with arbitrary correlations between input features and target functions: sample complexity, multiple descent curves and a hierarchy of phase transitions. In *International Conference on Machine Learning*, pp. 7578–7587. PMLR, 2021.
- Francesca Mignacco, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborová. Dynamical mean-field theory for stochastic gradient descent in gaussian mixture classification, 2020.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. The MIT Press, 2012. ISBN 026201825X.

- Roman Novak, Lechao Xiao, Jiri Hron, Jaehoon Lee, Alexander A. Alemi, Jascha Sohl-Dickstein, and Samuel S. Schoenholz. Neural tangents: Fast and easy infinite neural networks in python. In *International Conference on Learning Representations*, 2020. URL <https://github.com/google/neural-tangents>.
- Loucas Pillaud-Vivien, Alessandro Rudi, and Francis Bach. Statistical optimality of stochastic gradient descent on hard learning problems through multiple passes, 2018.
- B. Polyak and A. Juditsky. Acceleration of stochastic approximation by averaging. *Siam Journal on Control and Optimization*, 30:838–855, 1992.
- Phil Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic dimension of images and its impact on learning. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=XJk19XzGq2J>.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In J. Platt, D. Koller, Y. Singer, and S. Roweis (eds.), *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2008. URL <https://proceedings.neurips.cc/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf>.
- Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. *The Annals of Mathematical Statistics*, 22(3):400 – 407, 1951. doi: 10.1214/aoms/1177729586. URL <https://doi.org/10.1214/aoms/1177729586>.
- David Ruppert. Efficient estimations from a slowly convergent robbins-monro process. 02 1988.
- David Saad and Sara Solla. Dynamics of on-line gradient descent learning for multilayer neural networks. 04 1999.
- Alexander Shapiro. Asymptotic Properties of Statistical Estimators in Stochastic Programming. *The Annals of Statistics*, 17(2):841 – 858, 1989. doi: 10.1214/aos/1176347146. URL <https://doi.org/10.1214/aos/1176347146>.
- Stefano Spigler, Mario Geiger, and Matthieu Wyart. Asymptotic learning curves of kernel methods: empirical data versus teacher–student paradigm. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(12):124001, Dec 2020. ISSN 1742-5468. doi: 10.1088/1742-5468/abc61d. URL <http://dx.doi.org/10.1088/1742-5468/abc61d>.
- Aditya Varre, Loucas Pillaud-Vivien, and Nicolas Flammarion. Last iterate convergence of sgd for least-squares in the interpolation regime, 2021.
- Maksim Velikanov and Dmitry Yarotsky. Universal scaling laws in the gradient descent training of neural networks, 2021.
- Justin Werfel, Xiaohui Xie, and H. Seung. Learning curves for stochastic gradient descent in linear feedforward networks. In S. Thrun, L. Saul, and B. Schölkopf (eds.), *Advances in Neural Information Processing Systems*, volume 16. MIT Press, 2004. URL <https://proceedings.neurips.cc/paper/2003/file/f8b932c70d0b2e6bf071729a4fa68dfc-Paper.pdf>.
- Yiming Ying and Massimiliano Pontil. Online gradient descent learning algorithms. *Found. Comput. Math.*, 8(5):561–596, October 2008. ISSN 1615-3375.
- Lu Yu, Krishnakumar Balasubramanian, Stanislav Volgushev, and Murat A. Erdogdu. An analysis of constant step size sgd in the non-convex regime: Asymptotic normality and bias, 2020.
- C. Zhang, S. Bengio, Moritz Hardt, B. Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *ArXiv*, abs/1611.03530, 2017.
- Guodong Zhang, Lala Li, Zachary Nado, James Martens, Sushant Sachdeva, George E. Dahl, Christopher J. Shallue, and Roger Grosse. Which algorithmic choices matter at which batch sizes? insights from a noisy quadratic model, 2019.