
The alignment property of SGD noise and how it helps select flat minima: A stability analysis

Lei Wu

School of Mathematical Sciences
Peking University

Mingze Wang

School of Mathematical Sciences
Peking University

WeiJie J. Su

Department of Statistics and Data Science
University of Pennsylvania

Abstract

The phenomenon that stochastic gradient descent (SGD) favors flat minima has played a critical role in understanding the implicit regularization of SGD. In this paper, we provide an explanation of this striking phenomenon by relating the particular noise structure of SGD to its *linear stability* (Wu et al., 2018). Specifically, we consider training over-parameterized models with square loss. We prove that if a global minimum θ^* is linearly stable for SGD, then it must satisfy $\|H(\theta^*)\|_F \leq O(\sqrt{B}/\eta)$, where $\|H(\theta^*)\|_F, B, \eta$ denote the Frobenius norm of Hessian at θ^* , batch size, and learning rate, respectively. Otherwise, SGD will escape from that minimum *exponentially* fast. Hence, for minima accessible to SGD, the sharpness—as measured by the Frobenius norm of the Hessian—is bounded *independently* of the model size and sample size. The key to obtaining these results is exploiting the particular structure of SGD noise: The noise concentrates in sharp directions of local landscape and the magnitude is proportional to loss value. This alignment property of SGD noise provably holds for linear networks and random feature models (RFMs), and is empirically verified for nonlinear networks. Moreover, the validity and practical relevance of our theoretical findings are also justified by extensive experiments on CIFAR-10 dataset.

1 Introduction

Modern machine learning (ML) models are often operated with far more unknown parameters than training examples, a regime referred to as over-parameterization. In this regime, there are many global minima, all of which have zero training loss but their test performance can be significantly different [47]. Fortunately, it is often observed that SGD converges to those generalizable ones, even without needing any explicit regularizations [49]. This suggests there must exist certain “implicit regularization” mechanism at work [32, 14, 47, 6].

More mysteriously, SGD solutions often generalize better than gradient descent (GD) solutions [21, 38]. Therefore, the SGD noise must play a critical role in implicit regularization. The most popular explanation is that SGD favors flatter minima [21] and flatter minima generalize better [16]. This flat-minima principle has been extensively and successfully adopted in practice to tune the hyperparameters of SGD [44, 21] and to design new optimizers [17, 11, 43] for improving generalization. Therefore, understanding how SGD noise biases SGD towards flatter minima is of paramount importance, which is the main focus of this paper.

The works [44, 10, 48] show that SGD noise is highly anisotropic; [30, 41] find the magnitude of SGD noise to be loss dependent. Both structures are shown to be critical for SGD picking flat minima.

However, these works all make unrealistic (even wrong) over-simplifications of SGD noise (see the related work section for more details) in their analysis. In addition, instead of studying SGD, they all consider the continuous-time stochastic differential equation (SDE), which is a good modeling of SGD only in finite time and small learning rate (LR) regime [25]. It is generally unclear how this SDE modeling is relevant for understanding SGD with a large LR—a regime preferred in practice. Consequently, these works only provide intuitive and empirical analyses, lacking a quantitative characterization of when and how SGD favors flat minima.

Another line of works [46, 29] relate the selection bias of SGD to the *dynamical stability*. In over-parameterized case, all global minima are fixed points of SGD but their dynamical stabilities can be very different. At unstable minima, a small perturbation will drive SGD to leave away, whereas, for stable minima, SGD can stay around and even converge back after initial perturbations. Thus SGD prefers stable minima over unstable ones. Specifically, [46, 29] analyze the linear stability [1] of SGD, showing that a linearly stable minimum must be flat and uniform. Different from SDE-based analysis, this stability-based analysis is relevant for large-LR SGD and is even empirically accurate in predicting the properties of minima selected by SGD [46, 20, 8].

In this work, we follow the linear stability analysis in [46, 29] but take the particular structure of SGD noise into consideration. We establish a direct connection between linear stability and flatness, which allows us to obtain a quantitative characterization of how the learning rate and batch size affect the flatness of minima accessible to SGD. In contrast, [46, 29] have to introduce the another quantity: non-uniformity together with flatness to characterize linear stability because of neglecting the noise structure.

Setup Let $\{(x_i, y_i)\}_{i=1}^n$ with $x_i \in \mathbb{R}^d, y_i \in \mathbb{R}^K$ be the training set and $f(\cdot; \theta)$ with $\theta \in \mathbb{R}^p$ be our model. The *model size* is defined to be the number of parameters p and in this paper. For simplicity, we will always assume $K = 1$ and the extension to the case of $K > 1$ is straightforward. Let $L_i(\theta) = \frac{1}{2}|f(x_i; \theta) - y_i|^2$ be the fitting error at the i -th sample and $L(\theta) = \frac{1}{n} \sum_{i=1}^n L_i(\theta)$ be the empirical risk. We shall focus on the over-parameterized case in the sense that $\inf_{\theta} L(\theta) = 0$. To minimize $L(\cdot)$, we consider the mini-batch SGD:

$$\theta_{t+1} = \theta_t - \frac{\eta}{B} \sum_{i \in I_t} \nabla L_i(\theta_t), \quad (1)$$

where η and B are the learning rate and batch size, respectively. This SGD can be rewritten as $\theta_{t+1} = \theta_t - \eta(\nabla L(\theta_t) + \xi_t)$, where ξ_t is the noise, satisfying $\mathbb{E}[\xi_t] = 0$ and $\mathbb{E}[\xi_t \xi_t^T] = \Sigma(\theta_t)/B$. Here the noise covariance $\Sigma(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla L_i(\theta) \nabla L_i(\theta)^T - \nabla L(\theta) \nabla L(\theta)^T$. To characterize the local geometry of loss landscape, we consider the Fisher matrix: $G(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla f(x_i; \theta) \nabla f(x_i; \theta)^T$ and the Hessian matrix: $H(\theta) = G(\theta) + \frac{1}{n} \sum_{i=1}^n (f(x_i; \theta) - y_i) \nabla^2 f(x_i; \theta)$. When the loss value is small, $H(\theta) \approx G(\theta)$ and in particular, if θ^* is an global minimum, $H(\theta^*) = G(\theta^*)$.

Notations. For a vector a , let $\|a\| = \sqrt{a^T a}$ and $\|a\|_W = \sqrt{a^T W a}$. For a matrix A , denote by $\{\lambda_j(A)\}$ the eigenvalue of A in a decreasing order. For other notations, we refer to Appendix C.

Our main contributions are summarized as follows.

- We first show that for many ML models, the SGD noise is geometry aware: 1) the noise magnitude is proportional to the loss value; 2) the noise covariance aligns well with the Fisher matrix. Specifically, to quantify the alignment strength, we define a loss-scaled alignment factor $\mu(\theta)$, which is proved to be bounded from below, i.e., there exists a *size-independent* positive constant μ_0 such that $\mu(\theta) \geq \mu_0$, for linear networks (Proposition 2.3) and RFMs (Proposition 2.5), and is also empirically justified for nonlinear networks. Moreover, we identify that it is the uniformity of model gradient norms $\{\|\nabla f(x_i; \theta)\|_{G(\theta)}\}_i$ that accounts for this *alignment property* of SGD noise.
- We then provide a thorough analysis of the linear stability of SGD by exploiting the alignment property of noise. We prove in Theorem 3.3 that if a global minimum θ^* is linearly stable, then $\|H(\theta^*)\|_F \leq \eta^{-1} \sqrt{B/\mu_0}$. Here the constant μ_0 quantifies the alignment strength of SGD noise. Hence, for minima accessible to SGD, the Hessian’s Frobenius norm—the flatness perceived by SGD—is bounded independently of the model size and sample size. Moreover, if a minimum is too sharp, violating the preceding stability condition, SGD will escape from it *exponentially fast* (Theorem 3.5). Together, we obtain a quantitative characterization of when and how much SGD dislikes sharp minima.

- Our theoretical findings are also corroborated with well-designed experiments on a variety of models including linear networks, RFMs, convolutional networks, and fully-connected networks. In particular, the practical relevance is demonstrated in Section 4 by extensive experiments on classifying full CIFAR-10 dataset with VGG nets and ResNets.

1.1 Related work

Noise structures. [52, 19, 27] consider the Hessian-based approximation: $\Sigma(\theta) \approx \sigma^2 H(\theta)$, where σ is a small constant. [53] proposes an improved version: $\Sigma(\theta) \approx 2L(\theta)H(\theta)$. But these approximations in general cannot be accurate since Hessian is not semi-positive definite (SPD) in non-convex region. More recently, [30] and [41] study SGD by assuming $\Sigma(\theta) = 2L(\theta)H(\theta^*)$, where θ^* is a minimum of interest, and $\Sigma(\theta) = \sigma^2 L(\theta)I_p$, respectively. These assumptions completely ignore the state-dependence of noise direction. In contrast, we assume $\Sigma(\theta) = 2L(\theta)C(\theta)$ with $C(\theta)$ having a nontrivial alignment with the Fisher matrix $G(\theta)$, which does not impose any explicit structural assumption on $C(\theta)$. As a result, our assumption is much weaker and can be rigorously justified for popular ML models both theoretically and empirically. More importantly, we show that this weak alignment property is sufficient for analyzing the linear stability of SGD. We anticipate that our alignment assumption can be also adopted to analyze other properties of SGD.

Escape from sharp minima. The escape behavior of SGD was first studied in [52, 46], as an indicator of how much SGD dislikes sharp minima. One of the most mysterious observation is that the escape happens in an unreasonably efficient way. However, the theoretical analysis there assumes the noise to be state-independent, and consequently, the derived escape time depends polynomially on the loss barrier. Later [48, 30] attempt to study this issue using the classical diffusion-based framework [12] (Itô-SDE), which cannot explain the unreasonable escape efficiency at all since the resulting escape rates depend on the loss barrier exponentially. See also [30] for an improved analysis. [37, 51] argues that the SGD noise is heavy-tailed and thus SGD should be modeled as Lévy-SDE instead of Itô-SDE. Moreover, it is shown that the heavy-tailedness can ensure the escape rate depends on the basin volume instead of the loss barrier. Unfortunately, the volume in high dimensions always scales with the dimension exponentially and consequently, this does not explain the escape efficiency in high dimensions. Moreover, whether SGD noise is really heavy-tailed and whether the heavy-tailedness is really important for generalization are still debatable for neural networks [40, 26]. In contrast, we show that the unreasonable escape efficiency comes from the particular geometry-aware structure of SGD noise, regardless of whether the noise is heavy- or light-tailed.

Flatness metrics In the literature, a variety of flatness metrics have been adopted, such as the largest eigenvalue of Hessian [21], the trace of Hessian [9, 6], the basin volume [51], and the ones scaled by parameter norms [28, 39] in order to achieve the scaling-invariance for ReLU nets. These metrics are proposed for either computation easiness or bounding generalization gaps. It is unclear if they are perceivable to SGD, let alone how the boundedness of them depends on the batch size, learning rate, as well as the model size and sample size. We show that for SGD solutions, the Frobenius norm of Hessian—a flatness perceived by SGD through the linear stability—is bounded by a size-independent quantity. Note that a similar stability argument also applies to GD but only yielding the boundedness of the largest eigenvalue of Hessian [46, 31].

Lastly, we particularly mention the work [33], which provides a fine-grained analysis of the implicit bias of training two-layer diagonal linear networks. This work is related to ours since we both consider the magnitude and direction structure of SGD noise simultaneously. However, the analysis in [33] is limited to the specific toy model but ours is relevant for general models.

2 The alignment property of SGD noise

Since $\nabla L_i(\theta) = (f(x_i; \theta) - y_i)\nabla f(x_i; \theta)$, we have the following intuitive approximation [30]:

$$\begin{aligned} \Sigma(\theta) &= \frac{2}{n} \sum_{i=1}^n L_i(\theta) \nabla f(x_i; \theta) \nabla f(x_i; \theta)^T - \nabla L(\theta) \nabla L(\theta)^T \stackrel{(i)}{\approx} \frac{2}{n} \sum_{i=1}^n L_i(\theta) \nabla f(x_i; \theta) \nabla f(x_i; \theta)^T \\ &\stackrel{(ii)}{\approx} 2 \left(\frac{1}{n} \sum_{i=1}^n L_i(\theta) \right) \frac{1}{n} \sum_{i=1}^n \nabla f(x_i; \theta) \nabla f(x_i; \theta)^T = 2L(\theta)G(\theta), \end{aligned} \quad (2)$$

where (i) assumes that the full-batch gradient $\nabla L(\theta)$ to be negligible compared with the sample gradients $\{\nabla L_i(\theta)\}$; (ii) assumes that $\{\nabla f(x_i; \theta)\}_i$ and $\{L_i(\theta)\}_i$ are nearly decoupled. This approximation cannot be true in general but tells us that 1) The noise magnitude is proportional to the loss value; 2) the noise covariance aligns with the Fisher matrix.

Motivated by (2), we define $\alpha(\theta) = \frac{\text{Tr}(G(\theta)\Sigma(\theta))}{\|G(\theta)\|_F \|\Sigma(\theta)\|_F}$, $\beta(\theta) = \frac{\|\Sigma(\theta)\|_F}{2L(\theta)\|G(\theta)\|_F}$. Here $\alpha(\theta)$ quantifies the similarity between $\Sigma(\theta)$ and $G(\theta)$, characterizing how much the noise concentrates in sharp directions of local landscape. $\beta(\theta)$ characterizes the relative non-degeneracy of noise (with respect to the loss value). Then we define the loss-scaled alignment factor:

$$\mu(\theta) = \alpha(\theta)\beta(\theta) = \frac{\text{Tr}(\Sigma(\theta)G(\theta))}{2L(\theta)\|G(\theta)\|_F^2}, \quad (3)$$

which characterizes the direction and magnitude alignment simultaneously. Intuitively speaking, if $\mu(\theta)$ is bounded below, SGD noise is non-degenerate in sharp directions of local landscape. In particular, $\alpha(\theta) = \beta(\theta) = \mu(\theta) = 1$ if the approximation (2) holds.

We say the noise satisfies the μ -alignment if $\mu(\theta) > 0$. Compared with the decoupling approximation (2), the μ -alignment is a much weaker condition. Note that this specific weak quantification of alignment is inspired for analyzing the linear stability of SGD, which is the focus of this paper. Specifically, Theorem 3.3 shows that $\mu(\theta)$ along with the Frobenius norm of Hessian determines the linear stability of SGD. One may define other metrics to quantify alignment strength for studying other properties of SGD, but this is beyond the scope of this paper.

A relaxed alignment. Let $\Sigma_1(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla L_i(\theta) \nabla L_i(\theta)^T$, $\Sigma_2(\theta) = \nabla L(\theta) \nabla L(\theta)^T$. Then $\Sigma(\theta) = \Sigma_1(\theta) - \Sigma_2(\theta)$. It is often believed that the full-batch gradient ∇L is relatively small compared to the sample gradients $\{\nabla L_i\}_i$. As a result, the influence of $\Sigma_2(\theta)$ should be negligible compared to $\Sigma_1(\theta)$. To disentangle the influences of them, we define

$$\mu_1(\theta) = \frac{\text{Tr}(\Sigma_1(\theta)G(\theta))}{2L(\theta)\|G(\theta)\|_F^2}, \quad \mu_2(\theta) = \frac{\text{Tr}(\Sigma_2(\theta)G(\theta))}{2L(\theta)\|G(\theta)\|_F^2}.$$

Then $\mu(\theta) = \mu_1(\theta) - \mu_2(\theta)$. Our linear stability analysis in Section 3 show that $\mu_1(\theta) \geq \mu_1 > 0$, a condition we refer to as μ_1 -alignment, is sufficient to ensure that SGD only selects flat minima.

2.1 Why does the alignment property hold?

Definition 2.1 (Norm uniformity of model gradients). Let $g_i(\theta) = \nabla f(x_i; \theta)$, $\chi_i(\theta) := \|g_i(\theta)\|_{G(\theta)}^2 = g_i^T(\theta)G(\theta)g_i(\theta)$, $\bar{\chi}(\theta) = \frac{1}{n} \sum_{i=1}^n \chi_i(\theta)$. Define $\gamma(\theta) := \min_{i \in [n]} \frac{\chi_i(\theta)}{\bar{\chi}(\theta)}$.

The quantity $\gamma(\theta)$ measures the uniformity of model gradient norms and this property can guarantee the μ_1 -alignment as stated below.

Lemma 2.2. $\mu_1(\theta) \geq \gamma(\theta)$

Proof. Noticing $\bar{\chi}(\theta) = \frac{1}{n} \sum_{i=1}^n g_i(\theta)^T G(\theta) g_i(\theta) = \text{Tr}(G(\theta) \frac{1}{n} \sum_{i=1}^n g_i(\theta) g_i(\theta)^T) = \|G(\theta)\|_F^2$, we have

$$\begin{aligned} \text{Tr}(\Sigma_1(\theta)G(\theta)) &= \frac{2}{n} \sum_{i=1}^n L_i(\theta) \text{Tr}(g_i(\theta)g_i(\theta)^T G(\theta)) \\ &= \frac{2}{n} \sum_{i=1}^n L_i(\theta) \chi_i \geq \frac{2}{n} \sum_{i=1}^n L_i(\theta) \gamma \bar{\chi} = 2\gamma L(\theta) \|G(\theta)\|_F^2 \end{aligned}$$

Thus $\mu_1(\theta) = \text{Tr}(\Sigma_1(\theta)G(\theta))/(2L(\theta)\|G(\theta)\|_F^2) \geq \gamma(\theta)$. \square

The above proof suggests that the ‘‘decoupling’’ approximation in (2) holds in a weak sense if $\{\|\nabla f(x_i; \theta)\|_{G(\theta)}\}_i$ are uniform. One can apply a similar argument by assuming the uniformity of the fitting errors $\{L_i(\theta)\}$, which, unfortunately, we find never hold in practice. In contrast, we will show that the norm uniformity of model gradients provably holds for linear networks and RFMs, and can be empirically justified for nonlinear networks.

Over-parameterized linear models. Consider an over-parameterized linear model (OLM): $f(x; \theta) = F(\theta)^T x$, where $F : \Omega \mapsto \mathbb{R}^d$ denotes a general re-parameterization function. Note that $f(\cdot; \theta)$ only represents linear functions but the corresponding landscape can be highly non-convex. Typical examples include the linear network: $F(\theta) = W_1 W_2 \cdots W_L$ and the diagonal linear network: $F(\theta) = (\alpha_1^2 - \beta_1^2, \dots, \alpha_d^2 - \beta_d^2)^T$. Both have attracted a lot of attention in analyzing the implicit bias of GD and SGD [3, 42, 33, 13, 4]. The following proposition provides a precise characterization of the noise covariance for OLM models, whose proof is deferred to Appendix D.

Proposition 2.3. *Denote by $\mathcal{N}(0, S)$ the Gaussian distribution with mean zero and covariance matrix S . Suppose $f(\cdot; \theta)$ is a general OLM and $x \sim \mathcal{N}(0, S)$. Consider the online SGD setting, i.e., $n = \infty$. Then, $\Sigma(\theta) = \nabla L(\theta) \nabla L(\theta)^T + 2L(\theta)G(\theta)$ and $\mu_1(\theta) \geq \mu(\theta) \geq 1$*

This proposition shows that the alignment property holds in the entire parameter space and moreover, the alignment strength is independent of model size. Here the alignment is only proved for the infinite-sample case. Similar results should hold for finite-sample cases by concentration inequalities as long as n is relatively large, but this straightforward extension does not bring any new insights. It is more interesting to consider the low-sample regime (i.e., $n < d$), where the alignments indeed hold (at least in typical regions explored by SGD) as demonstrated empirically in Figure 1. In addition, the above proposition provides a closed-form expression of the noise covariance, which might be useful for analyzing other properties of SGD beyond the linear stability. A comprehensive analysis of these issues is left to future work.

Feature-based models. Consider a feature-based model $f(x; \theta) = \sum_{j=1}^m \theta_j \varphi_j(x) = \langle \theta, \Phi(x) \rangle$. In this case, the model gradients: $g_i = g_i(\theta) = \Phi(x_i)$ and the Hessian and Fisher matrix: $G = H = \frac{1}{n} \sum_{i=1}^n g_i g_i^T$ all are constant. But the noise covariance $\Sigma(\theta)$ is still state-dependent. The norm uniformity of model gradients also becomes constant: $\gamma = \min_i \chi_i / (\frac{1}{n} \sum_{i=1}^n \chi_i)$ with $\chi_i = \|g_i\|_G^2$.

Lemma 2.4. $\mu_1(\theta) \geq \gamma, \mu_2(\theta) \leq \tau(G) := \lambda_1^2(G) / \sum_j \lambda_j^2(G)$, and $\mu(\theta) \geq \gamma - \tau(G)$.

The above lemma suggests that the $\mu_2(\theta)$ term is negligible as long as the Fisher matrix is not nearly rank-1. By bounding the γ and $\tau(G)$, we can prove the μ_1 - and μ -alignment for random ReLU feature models as stated in the following proposition. The proof is presented in Appendix F, where similar results for general RFMs are provided (see Proposition F.7).

Proposition 2.5. *If $\varphi_j(x) = \text{ReLU}(w_j^T x)$ with $w_j \stackrel{iid}{\sim} \text{Unif}(\sqrt{d}\mathbb{S}^{d-1})$ and $x \sim \text{Unif}(\mathbb{S}^{d-1})$. Then, for any $\delta \in (0, 1)$, if $m \geq n \gtrsim d^5 \log(1/\delta)$, then w.p. at least $1 - \delta$, $\mu_1(\theta) \geq 1$ and $\mu(\theta) \gtrsim d^{-1}$.*

Although feature-based models are linear, their analysis is still applicable to understand nonlinear models, as long as the nonlinear model *locally* behaves like the linearized one: $f_{\text{lin}}(x; \theta) := f(x; \theta^*) + \langle \theta - \theta^*, \nabla f(x; \theta^*) \rangle$ with $\nabla f(x; \theta^*)$ learned from data. Hence, Proposition 2.5 can explain why the alignment property holds in a local region around global minima. Note that this is sufficient for characterizing the linear stability of SGD, which is a local property in nature.

2.2 Empirical validations

Figure 1a reports the values of $\alpha(\theta_t), \beta(\theta_t), \mu(\theta_t)$ during the SGD training of four types of models, including the linear networks and RFMs analyzed above, and fully-connected networks (FCN) and convolutional neural networks (CNN). First, one can see that $\alpha(\theta_t)$'s are quite close to 1 during the entire training, suggesting the strong concentration of SGD noise in sharp directions of local landscape. Second, $\beta(\theta_t)$'s keep bounded below, implying the noise magnitudes are sufficiently large with respect to the training loss. As a result, $\mu(\theta_t)$'s are significantly positive for all models examined. In particular, for the linear networks, the alignment holds in the low-sample regime, where $n < d$.

The size independence. Figure 1b further examines how the extent of over-parameterization affects the alignment strength. One can see clearly that for linear networks and RFMs, $\mu(\theta)$'s are independent of the model size, which confirms our theoretical findings proved above. In addition, we also observe that for nonlinear networks, the alignment strength is also (nearly) independent of the model size. For instance, for CNN, the value of $\mu(\theta)$ only decreases from around 1.05 to 1.0 as the model size is increased by more than two orders of magnitude.

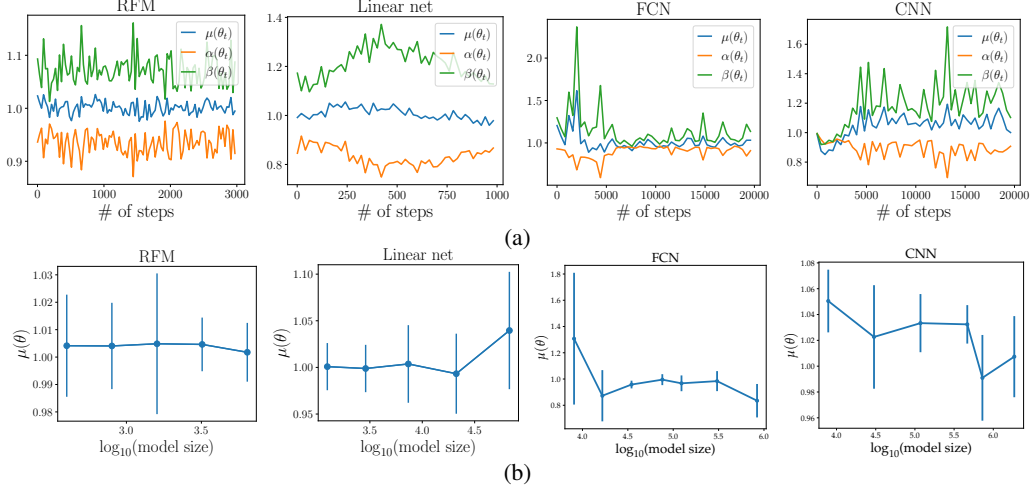


Figure 1: **The alignment property of SGD noise.** Four types of models, including the RFM, linear network, fully-connected network (FCN), and convolutional neural network (CNN), are examined. We refer to Appendix A for the experimental setup. Note that the linear network is trained in a low-sample regime ($n = 100, d = 50$). (a) The alignment factors during training. (b) How the alignment strength of convergent solution changes with the over-parameterization. The error bar corresponds to the standard deviation over 5 runs.

The norm uniformity of model gradients. Figure 2 shows the norm uniformity of model gradients, where we report the values of $\gamma(\theta)$ in the entire SGD trajectory. One can see that $\gamma(\theta)$ is indeed bounded below, implying that the SGD noise satisfies the μ_1 -alignment according to Lemma 2.2.

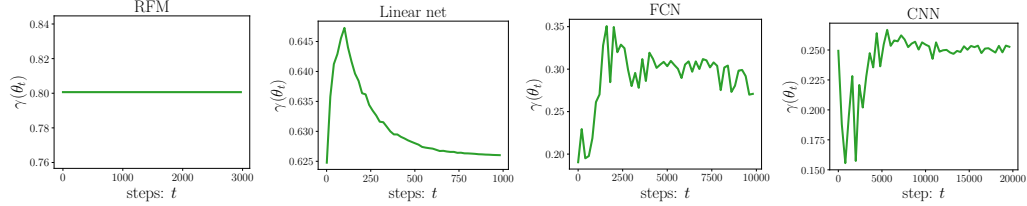


Figure 2: **The norm uniformity of model gradients.** The values of $\gamma(\theta_t)$ in the entire SGD trajectory are reported. It is shown that the norms of model gradient norms are uniform during the entire training process.

Note that in experiments, we only show that the alignment property is satisfied in typical regions explored by SGD, including the random initialization and the convergent region. In contrast, for OLMs and RFMs, we in fact prove the alignment property for the entire parameter space.

3 The linear stability analysis

Let θ^* be a global minimum of $L(\cdot)$. When θ_t is close to θ^* , the local dynamical behavior of SGD can be characterized by linearizing the dynamics around θ^* :

$$\tilde{\theta}_{t+1} = \tilde{\theta}_t - \frac{\eta}{B} \sum_{i \in I_t} \nabla^2 L_i(\theta^*) (\tilde{\theta}_t - \theta^*), \quad (4)$$

where $\nabla^2 L_i(\theta^*) = \nabla f(x_i; \theta^*) \nabla f(x_i; \theta^*)^T$. This corresponds to the local quadratic approximation of the loss $L(\cdot)$ or the local linearization of the model around θ^* :

$$f_{\text{lin}}(x; \theta) = f(x; \theta^*) + \langle \theta - \theta^*, \nabla f(x; \theta^*) \rangle. \quad (5)$$

Specifically, (4) is exactly the SGD of training the linearized model (5).

Definition 3.1 (Linear stability). *A global minimum θ^* is said to be linearly stable if there exists a $C > 0$ such that it holds for the linearized dynamics (4) that $\mathbb{E}[L(\tilde{\theta}_t)] \leq C \mathbb{E}[L(\tilde{\theta}_0)], \forall t \geq 0$.*

It is well-known in dynamical system that the local behavior of the original nonlinear dynamics can be characterized by the linearized one if the local quadratic approximation is non-degenerate. However, in over-parameterized case, the local quadratic approximation is degenerate in flat directions. Consequently, one may be concerned about the relevance of local quadratic approximation and the

resulting linear stability analysis. Fortunately, the stability in Definition 3.1 is particularly measured with the change of loss value. Thus the instability mostly comes from noise perturbations in sharp directions and the alignment property guarantees that noise mostly concentrates in sharp directions. Consequently, the flat directions contribute little to the instability. In sharp directions, the local quadratic approximation is always valid, thereby explaining the relevance of linear stability analysis. The rigorous formulation of this intuition is left to future work and we instead resort to numerical experiments to demonstrate the validity in this paper.

For simplicity, we will use θ_t to denote $\tilde{\theta}_t$; let $\theta^* = 0$ and $g_i = \nabla f(x_i; \theta^*)$. For the linearized model $f_{\text{lin}}(\cdot; \theta)$, we have $L(\theta) = \frac{1}{2n} \sum_{i=1}^n |\theta^T g_i|^2 = \frac{1}{2} \theta^T H \theta$, $G = H = \frac{1}{n} \sum_{i=1}^n g_i g_i^T$, where we omit the dependence on θ^* for simplicity. Note that the Fisher and Hessian matrix are constant but the noise covariance $\Sigma(\theta) = \frac{1}{n} \sum_{i=1}^n |g_i^T \theta|^2 g_i g_i^T - H \theta \theta^T H$ is still state-dependent.

Before considering the specific linearized SGD (4), we first have a general result.

Lemma 3.2. *Consider a general SGD: $\theta_{t+1} = \theta_t - \eta(\nabla L(\theta_t) + \xi_t)$ for the linearized model (5), where $(\xi_t)_{t \geq 1}$ are any noises satisfying $\mathbb{E}[\xi_t] = 0$, $\mathbb{E}[\xi_t \xi_t^T] = S(\theta_t)$. Then we have*

$$\mathbb{E}[L(\theta_{t+1})] = \mathbb{E}[r(\theta_t)L(\theta_t) + \eta^2 v(\theta_t)], \quad (6)$$

where $v(\theta) = \text{Tr}(HS(\theta))/2$ and $r(\theta) \geq 0$. Moreover, if $\eta \leq 2/\lambda_1(H)$, then $r(\theta) \leq 1$.

Proof. Using the fact $\mathbb{E}[\xi_t] = 0$ and $\mathbb{E}[\xi_t \xi_t^T] = S(\theta_t)$, we have

$$\begin{aligned} \mathbb{E}[L(\theta_{t+1})] &= \mathbb{E}\left[\frac{1}{2}(\theta_t - \eta \nabla L(\theta_t) + \eta \xi_t)^T H (\theta_t - \eta \nabla L(\theta_t) + \eta \xi_t)\right] \\ &= \mathbb{E}[L(\theta_t) - \eta \nabla L(\theta_t)^T H \theta_t + \frac{\eta^2}{2} \nabla L(\theta_t)^T H \nabla L(\theta_t)] + \frac{\eta^2}{2} \mathbb{E}[\text{Tr}(HS(\theta_t))] \\ &= \mathbb{E}[r(\theta_t)L(\theta_t) + \eta^2 v(\theta_t)], \end{aligned} \quad (7)$$

where $r(\theta) = 1 - 2\eta \frac{\theta^T H^2 \theta}{\theta^T H \theta} + \eta^2 \frac{\theta^T H^3 \theta}{\theta^T H \theta}$ since $\nabla L(\theta) = H\theta$. By Lemma G.2, $r(\theta) \geq 0$ and if $\eta \leq 2/\lambda_1(H)$, then $r(\theta) \leq 1$. \square

The two terms $r(\theta_t)L(\theta_t)$ and $\eta^2 v(\theta_t)$ denote the contributions from the full-batch gradient $\nabla L(\theta_t)$ and the noise ξ_t , respectively. The stability is affected by both terms simultaneously. It is well-known that if θ^* is linearly stable for GD, then $\lambda_1(H(\theta^*)) \leq 2/\eta$ (see, e.g., [46, 31]). This also holds for SGD but SGD imposes a stricter condition because of the extra $\eta^2 v(\theta_t)$ term. Specifically,

$$\mathbb{E}[L(\theta_{t+1})] = \mathbb{E}[r(\theta_t)L(\theta_t) + \eta^2 v(\theta_t)] \geq \eta^2 \mathbb{E}[v(\theta_t)] = 0.5\eta^2 \text{Tr}(HS(\theta)). \quad (8)$$

Therefore, the more $S(\theta)$ aligns with H , the more unstable that minimum is. Specifically, let $S(\theta) = 2\sigma^2 L(\theta)C(\theta)$. Then, $\mathbb{E}[L(\theta_{t+1})] \geq \eta^2 \sigma^2 \mathbb{E}[L(\theta_t) \text{Tr}(HC(\theta_t))]$. Thus to ensure a stable convergence, we should roughly have $\text{Tr}(HC(\theta)) \leq \frac{1}{\sigma^2 \eta^2}$. We next show that this can lead to a flatness control by utilizing the alignment between $C(\theta)$ and H .

3.1 The linear stability imposes size-independent flatness constraints

For the mini-batch SGD, the following theorem characterizes how the batch size and learning rate affect the flatness—as measured by the Frobenius norm of Hessian—of minima accessible to SGD.

Theorem 3.3. *Let θ^* be a global minimum that is linearly stable. Denote by $\mu(\theta)$ the alignment factors for the linearized SGD (4) and model (5). If $\mu(\theta) \geq \mu_0$, then $\|H(\theta^*)\|_F \leq \frac{1}{\eta} \sqrt{\frac{B}{\mu_0}}$.*

Proof. By (8) and the definition of $\mu(\theta)$, we have

$$\mathbb{E}[L(\theta_{t+1})] \geq \frac{\eta^2}{2B} \mathbb{E}[\text{Tr}(H\Sigma(\theta_t))] \geq \frac{\eta^2 \|H\|_F^2}{B} \mathbb{E}[\mu(\theta_t)L(\theta_t)] \geq \frac{\mu_0 \eta^2 \|H\|_F^2}{B} \mathbb{E}[L(\theta_t)]. \quad (9)$$

To ensure the stability, we must have $\mu_0 \eta^2 \|H\|_F^2 / B \leq 1$, leading to $\|H\|_F \leq \sqrt{B/\mu_0}/\eta$. \square

We have shown in Section 2 that μ_0 is (nearly) size-independent, and thus the obtained upper bound of flatness is also (nearly) size-independent. As a comparison, for GD, the linear stability can only ensure $\lambda_{\max}(H(\theta^*)) \leq 2/\eta$. This gives a bound of the Hessian's Frobenius norm: $\|H(\theta^*)\|_F \leq 2\sqrt{p}/\eta$,

depending on the model size explicitly. The comparison of two bounds partially explains why SGD tends to select flatter minima than GD.

We show below that the μ -alignment can be further relaxed to the μ_1 -alignment. The proof is similar to the one of Theorem 3.3 and deferred to Appendix G.

Proposition 3.4. *Under the setting of Theorem 3.3, if the noise of linearized SGD satisfies $\mu_1(\theta) \geq \mu_1$, then $\|H(\theta^*)\|_F \leq \min\left(\frac{B}{\sqrt{(B-1)\mu_1}}, \frac{2B}{\mu_1}\right) \frac{1}{\eta}$.*

When $B \gg 1$, the bound becomes $B/(\eta\sqrt{(B-1)\mu_1}) \approx \sqrt{B/\mu_1}/\eta$, which is the same as the case of μ -alignment. Thus the influence of ∇L is indeed negligible compared with $\{\nabla L_i\}_i$.

Note that the linear stability is local in nature and hence our analysis essentially only needs the μ_1 -alignment to hold locally around minima of interest. Experiments in Figure 2 shows that $\gamma(\theta^*)$ is always bounded below, i.e., the norm uniformity of model gradients is satisfied. Combining with Proposition F.7, we can conclude that the alignment property assumed in Proposition 3.4 holds.

Tightness of our analysis. In the analysis above, we only inspect the instability caused by the noise, with the full-batch gradient completely ignored. Therefore, we anticipate that our bound is tighter in small-batch regime, where the noise term dominates the full-batch term. We will see that our numerical experiments indeed confirm this tightness in small-batch regime. However, for obtaining the tightest bound, one may need to consider both components simultaneously; this is much more complicated and left to future work.

Numerical validations. Figure 3a numerically shows how the Frobenius norm of Hessian (not only the upper bound) changes with extent of over-parameterization, where the trace of Hessian is also plotted for comparison. One can see that the Frobenius norm indeed keeps almost unchanged as increasing the model size but the Hessian trace increases significantly. Figure 3b further shows the ratio between the Frobenius norm and our upper bound in the training process and two batch sizes are examined. We have two observations. First, the correctness of our bound holds for the entire SGD trajectory, suggesting that the linear stability analysis is relevant for the entire training process. Second, as expected, the theoretical bound is indeed tighter for the case with a smaller batch size.

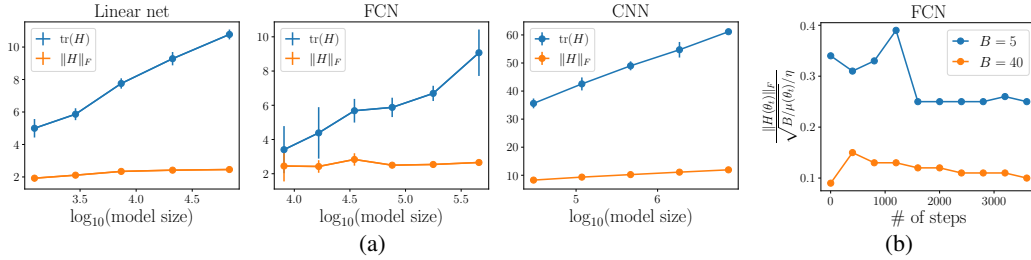


Figure 3: (a) The Frobenius norm and trace of Hessian vs. model size. The error bar corresponds to the standard deviation estimated over 5 runs. (b) The ratio between the Frobenius norm and our upper bound in the training process. Here $B = 5$ and $B = 40$ are examined. For more experiment details, we refer to Appendix A.

3.2 SGD escapes from sharp minima exponentially fast

The following theorem shows that the pure noise-driven escape from a sharp minimum is *exponentially fast*, whose proof follows trivially from (9).

Theorem 3.5. *Under the setting of Theorem 3.3, if $\|H(\theta^*)\|_F > \frac{1}{\eta}\sqrt{\frac{B}{\mu_0}}$, then the linearized SGD satisfies $\mathbb{E}[L(\theta_t)] \geq \gamma_0^t \mathbb{E}[L(\theta_0)]$ with $\gamma_0 = \frac{\eta^2 \mu_0}{B} \|H(\theta^*)\|_F^2 > 1$.*

Hence, linearized SGD takes roughly $\log_{\gamma_0}(1/\varepsilon)$ steps to escape from a $O(\varepsilon)$ -loss region to a $O(1)$ -loss region. The escape time depends on the loss barrier only logarithmically and is independent of the parameter space dimension. Due to the local closeness between linearized SGD and the original SGD, this partially explains the unreasonable escape efficiency of SGD for training big models. In contrast, the escape rates of existing works [48, 52, 51, 30] are either exponentially slow with respect to the loss barrier or suffer from the curse of dimensionality.

Figure 4 shows the trajectories of SGD escaping from sharp minima. It is demonstrated that the escape is indeed exponentially fast and specifically, 10 steps are enough for SGD escaping to a high-loss region for all the models examined. In addition, we observe that the escape is still exponentially fast in the high-loss region, although our analysis only applies to a local region. How can we explain this nonlocal escape behavior? We leave the study of this interesting phenomenon to future work.

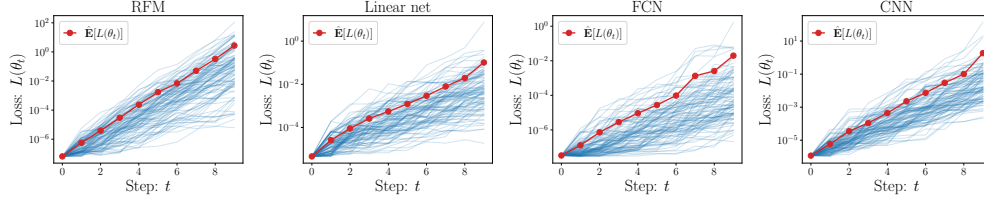


Figure 4: **The exponentially fast escape from sharp minima.** The blue curves are 200 trajectories of SGD; the red curve corresponds to the average. The sharp minimum is found by GD. When GD nearly converge, we switch to SGD with the same learning rate. This choice ensures that the minimum is stable for GD, and thus the escape is purely driven by SGD noise. For more experimental details, we refer to Appendix A.

3.3 The importance of the noise structure

The magnitude structure. Theorem 3.5 together with its proof suggests that the loss dependence of noise magnitude is critical for obtaining the exponentially fast escape. The intuition is as follows. When θ_t is perturbed by noise to θ_{t+1} where $L(\theta_{t+1}) > L(\theta_t)$, the noise magnitude becomes larger there and thus θ_{t+1} is easier to be perturbed to a larger-loss region. This positive feedback drives SGD to leave exponentially fast. On the contrary, the following lemma shows that if the noise is uniformly bounded, the noise-driven escape is at most linear in time.

Lemma 3.6. *Under the setting of Lemma 3.2, assume $\eta \leq 2/\lambda_1(H)$ and $\mathbb{E}[HS(\theta)] \leq 2\sigma^2$. Then $\mathbb{E}[L(\theta_t) - L(\theta_0)] \leq \eta^2\sigma^2 t$.*

We set $\eta \leq 2/\lambda_1(H)$ to avoid the exponential escape caused by the full-batch component.

Proof. By Lemma G.2, when $\eta \leq 2/\lambda_1(H)$, $r(\theta) \leq 1$. Thus Lemma 3.2 implies $\mathbb{E}[L(\theta_{t+1})] \leq \mathbb{E}[L(\theta_t)] + \eta^2\sigma^2$, which implies $\mathbb{E}[L(\theta_t)] \leq \mathbb{E}[L(\theta_0)] + \eta^2\sigma^2 t$. \square

The direction structure. We now turn to consider the impact of direction structure. Consider general SGDs: $\theta_{t+1} = \theta_t - \eta(\nabla L(\theta_t) + \xi_t)$ with $\mathbb{E}[\xi_t \xi_t] = S(\theta_t)/B$ for the linearized model (5). We compare two type of (unrealistic) noises:

- *Geometry-aware noise:* $S_1(\theta) = 2L(\theta)H$.
- *Isotropic noise:* $S_2(\theta) = 2\sigma^2 L(\theta)I_p$ with $\sigma^2 = \text{Tr}(H)/p$. Here, the value of σ^2 is chosen to ensure two types of noises have the same total variance for a fair comparison [52].

Note that p denotes the model size. Analogous to Theorem 3.3, for the second isotropic SGD,

$$\mathbb{E}[L(\theta_{t+1})] \geq \frac{\eta^2}{2B} \text{Tr}(HS(\theta)) \geq \mathbb{E}[L(\theta_t)] \frac{\sigma^2 \eta^2}{B} \text{Tr}(H) = \mathbb{E}[L(\theta_t)] \frac{\eta^2}{pB} \text{Tr}(H)^2.$$

Hence, the instability decreases with the parameter-space dimension and the resulting flatness constraint is $\text{Tr}(H(\theta^*)) \leq \sqrt{pB}/\eta$, depending on the model size explicitly. In contrast, for the first noise, Theorem 3.3 implies $\|H(\theta^*)\|_F \leq \sqrt{B}/\eta$, independent of the model size. This difference can be intuitively explained as follows. The isotropic noise wastes most energy in perturbing SGD along flat directions, which barely affects the instability. In contrast, the geometry-aware noise focuses most energy on perturbing SGD along sharp directions, causing much more instability.

4 Larger-scale experiments

We have provided small-scale experiments to justify the validity of our theoretical findings for a variety of ML models. Here we turn to demonstrate the practical relevance by consider the classification of the CIFAR-10 dataset [22] with VGG nets [36] and ResNets [15]. In training, all explicit regularizations are removed to keep consistent with our theoretical analysis. More details of the experimental setup can be found in Appendix A.

The alignment property and escaping behavior. Figure 5a reports the alignment strength of SGD noise during training for VGG-19 and ResNet-110. One can see that the alignment factors are significantly positive and similar results are also observed for a variety of VGG nets and ResNets of different depths and can be found in Figure 8 of Appendix B. One can see that the alignment strength is nearly independent of the model size. Figure 5b shows the behavior of escaping from sharp minima for VGG-19 and ResNet-110. One can still observe that the escape is exponentially fast and similar observation for other ResNets and VGG nets can be found in Figure 9 in the appendix. These observations suggest that our theoretical findings also hold for this practical setting.

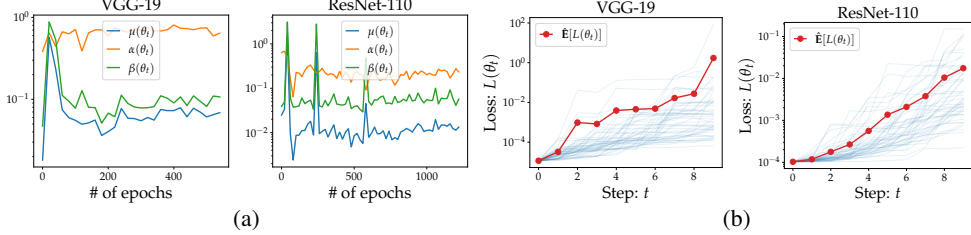


Figure 5: (a) The alignment factors and the escaping behavior for VGG-19. Similar results are also observed for all examined ResNets and VGG nets, which can be found in Appendix B.2. (b) The actual flatness of SGD solution and the corresponding theoretical upper bound (Theorem 3.3). (c) The upper bound becomes tighter as decreasing the batch size.

The flatness and upper bound. Figure 6a reports the flatness of convergent solution and the corresponding upper bound predicted by our theory for ResNets and VGG nets. It is again observed that the flatness is nearly independent of the model size. A surprising observation in Figure 6a is that our upper bounds are rather tight, see, e.g., VGG-16 and VGG-19. This tightness suggests that SGD runs (nearly) at the edge of stability [46, 8]. Moreover, as expected, Figure 6b shows that the our bound becomes tighter as decreasing batch size. This is consistent with what we observe in small-scale experiment in Figure 3b.

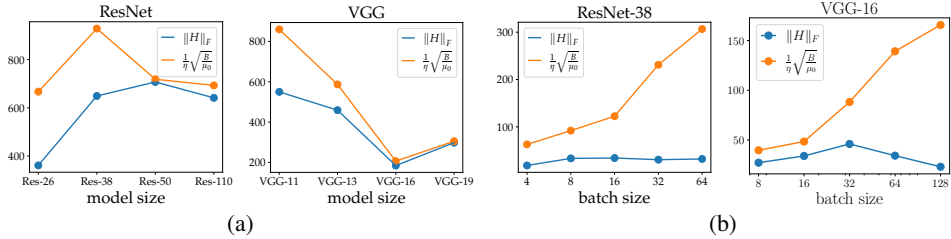


Figure 6: (a) The flatness and upper bound for various ResNets and VGG nets. (b) How the tightness of upper bound changes with decreasing batch size.

5 Conclusion

We provide a stability-based explanation of why SGD selects flat minima. Our current understanding is as follows. 1) For popular ML models, the SGD noise aligns very well with local landscape. 2) This alignment property ensures that the flatness of stable minima must be size-independent. This understanding is made rigorous and quantitative by introducing a loss-scaled alignment factor to characterize the alignment strength and analyze the linear stability. Obviously, many questions remains open. For example, can we understand what roles the stability plays in the whole dynamic process of SGD instead of only around global minima? Can we establish the connection between the Hessian’s Frobenius norm and generalization? Can we provide a fine-grained characterization of the noise structure and how the structure is related to implicit regularization of SGD? We leave the discussion of these important questions to future work.

Acknowledgements

We thank Zhiqin Xu and the anonymous reviewers for helpful suggestions. The work of Lei Wu is supported by a startup fund from Peking University. The work of Weijie J. Su is supported in part by NSF Grants CAREER DMS-1847415 and an Alfred Sloan Research Fellowship.

References

- [1] Vladimir Igorevich Arnold. *Geometrical methods in the theory of ordinary differential equations*, volume 250. Springer Science & Business Media, 2012. 2
- [2] Nachman Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, 68(3):337–404, 1950. 22
- [3] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019. 5
- [4] Shahar Azulay, Edward Moroshko, Mor Shpigel Nacson, Blake E Woodworth, Nathan Srebro, Amir Globerson, and Daniel Soudry. On the implicit bias of initialization shape: Beyond infinitesimal mirror descent. In *International Conference on Machine Learning*, pages 468–477. PMLR, 2021. 5
- [5] Francis Bach. Breaking the curse of dimensionality with convex neural networks. *The Journal of Machine Learning Research*, 18(1):629–681, 2017. 25
- [6] Guy Blanc, Neha Gupta, Gregory Valiant, and Paul Valiant. Implicit regularization for deep neural networks driven by an Ornstein-Uhlenbeck like process. In *Conference on learning theory*, pages 483–513. PMLR, 2020. 1, 3
- [7] Youngmin Cho and Lawrence K Saul. Kernel methods for deep learning. In *Proceedings of the 22nd International Conference on Neural Information Processing Systems*, pages 342–350, 2009. 25
- [8] Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2020. 2, 10
- [9] Alex Damian, Tengyu Ma, and Jason D Lee. Label noise SGD provably prefers flat global minimizers. *Advances in Neural Information Processing Systems*, 34:27449–27461, 2021. 3
- [10] Yu Feng and Yuhai Tu. The inverse variance–flatness relation in stochastic gradient descent is critical for finding flat minima. *Proceedings of the National Academy of Sciences*, 118(9), 2021. 1
- [11] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2020. 1
- [12] Crispin Gardiner. *Stochastic methods*, volume 4. springer Berlin, 2009. 3
- [13] Jeff Z HaoChen, Colin Wei, Jason Lee, and Tengyu Ma. Shape matters: Understanding the implicit bias of the noise covariance. In *Conference on Learning Theory*, pages 2315–2357. PMLR, 2021. 5
- [14] Hangfeng He and Weijie Su. The local elasticity of neural networks. In *International Conference on Learning Representations*, 2020. 1
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 9, 15
- [16] S. Hochreiter and J. Schmidhuber. Flat minima. *Neural Computation*, 9(1):1–42, 1997. 1
- [17] P Izmailov, AG Wilson, D Podoprikin, D Vetrov, and T Garipov. Averaging weights leads to wider optima and better generalization. In *34th Conference on Uncertainty in Artificial Intelligence 2018, UAI 2018*, pages 876–885, 2018. 1
- [18] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018. 21
- [19] Stanisław Jastrzebski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. Three factors influencing minima in SGD. *arXiv preprint arXiv:1711.04623*, 2017. 3
- [20] Stanisław Jastrzebski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun Cho, and Krzysztof Geras. The break-even point on optimization trajectories of deep neural networks. In *International Conference on Learning Representations*, 2019. 2

- [21] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations (ICLR)*, 2017. 1, 3
- [22] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images, 2009. 9
- [23] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 15
- [24] Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces: Isoperimetry and processes*. Springer Science & Business Media, 2013. 20
- [25] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic gradient algorithms. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 2101–2110. PMLR, Aug 2017. 2
- [26] Zhiyuan Li, Sathika Malladi, and Sanjeev Arora. On the validity of modeling SGD with stochastic differential equations (SDEs). In *Advances in Neural Information Processing Systems*, volume 34, 2021. 3
- [27] Zhiyuan Li, Tianhao Wang, and Sanjeev Arora. What happens after SGD reaches zero loss? –a mathematical framework. In *International Conference on Learning Representations*, 2022. 3
- [28] Tengyuan Liang, Tomaso Poggio, Alexander Rakhlin, and James Stokes. Fisher-Rao metric, geometry, and complexity of neural networks. In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89, pages 888–896. PMLR, 2019. 3
- [29] Chao Ma and Lexing Ying. On linear stability of SGD and input-smoothness of neural networks. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [30] Takashi Mori, Liu Ziyin, Kangqiao Liu, and Masahito Ueda. Logarithmic landscape and power-law escape rate of SGD. *arXiv preprint arXiv:2105.09557*, 2021. 1, 3, 8
- [31] Rotem Mulayoff, Tomer Michaeli, and Daniel Soudry. The implicit bias of minima stability: A view from function space. *Advances in Neural Information Processing Systems*, 34, 2021. 3, 7
- [32] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. *arXiv preprint arXiv:1412.6614*, 2014. 1
- [33] Scott Pesme, Loucas Pillaud-Vivien, and Nicolas Flammarion. Implicit bias of SGD for diagonal linear networks: A provable benefit of stochasticity. *Advances in Neural Information Processing Systems*, 34, 2021. 3, 5
- [34] Dino Sejdinovic and Arthur Gretton. What is an RKHS? *Lecture Notes*, 2012. 25
- [35] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014. 20
- [36] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 9, 15
- [37] Umut Simsekli, Levent Sagun, and Mert Gurbuzbalaban. A tail-index analysis of stochastic gradient noise in deep neural networks. In *International Conference on Machine Learning*, pages 5827–5837. PMLR, 2019. 3
- [38] Weijie Su. Neurashed: A phenomenological model for imitating deep learning training. *arXiv preprint arXiv:2112.09741*, 2021. 1
- [39] Yusuke Tsuzuku, Issei Sato, and Masashi Sugiyama. Normalized flat minima: Exploring scale invariant definition of flat minima for neural networks using PAC-Bayesian analysis. In *International Conference on Machine Learning*, pages 9636–9647. PMLR, 2020. 3
- [40] Xingyu Wang, Sewoong Oh, and Chang-Han Rhee. Eliminating sharp minima from SGD with truncated heavy-tailed noise. In *International Conference on Learning Representations*, 2022. 3
- [41] Stephan Wojtowytsch. Stochastic gradient descent with noise of machine learning type. part II: Continuous time analysis. *arXiv preprint arXiv:2106.02588*, 2021. 1, 3
- [42] Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, pages 3635–3673. PMLR, 2020. 5

- [43] Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization. *Advances in Neural Information Processing Systems*, 33, 2020. [1](#)
- [44] Jingfeng Wu, Wenqing Hu, Haoyi Xiong, Jun Huan, Vladimir Braverman, and Zhanxing Zhu. On the noisy gradient descent that generalizes as SGD. In *International Conference on Machine Learning*, pages 10367–10376. PMLR, 2020. [1](#)
- [45] Lei Wu and Jihao Long. A spectral-based analysis of the separation between two-layer neural networks and linear methods. *Journal of Machine Learning Research*, 23(119):1–34, 2022. [25](#)
- [46] Lei Wu, Chao Ma, and Weinan E. How SGD selects the global minima in over-parameterized learning: A dynamical stability perspective. *Advances in Neural Information Processing Systems*, 31:8279–8288, 2018. [2](#), [3](#), [7](#), [10](#)
- [47] Lei Wu, Zhanxing Zhu, and Weinan E. Towards understanding generalization of deep learning: Perspective of loss landscapes. *arXiv preprint arXiv:1706.10239*, 2017. [1](#)
- [48] Zeke Xie, Issei Sato, and Masashi Sugiyama. A diffusion theory for deep learning dynamics: Stochastic gradient descent exponentially favors flat minima. In *International Conference on Learning Representations*, 2020. [1](#), [3](#), [8](#)
- [49] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017. [1](#)
- [50] Hongyi Zhang, Yann N. Dauphin, and Tengyu Ma. Residual learning without normalization via better initialization. In *International Conference on Learning Representations*, 2019. [15](#)
- [51] Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Chu Hong Hoi, et al. Towards theoretically understanding why SGD generalizes better than Adam in deep learning. *Advances in Neural Information Processing Systems*, 33, 2020. [3](#), [8](#)
- [52] Zhanxing Zhu, Jingfeng Wu, Bing Yu, Lei Wu, and Jinwen Ma. The anisotropic noise in stochastic gradient descent: Its behavior of escaping from sharp minima and regularization effects. In *International Conference on Machine Learning*, pages 7654–7663. PMLR, 2019. [3](#), [8](#), [9](#)
- [53] Liu Ziyin, Kangqiao Liu, Takashi Mori, and Masahito Ueda. Strength of minibatch noise in SGD. In *International Conference on Learning Representations*, 2022. [3](#)

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? [\[Yes\]](#) See Section ??.
- Did you include the license to the code and datasets? [\[No\]](#) The code and the data are proprietary.
- Did you include the license to the code and datasets? [\[N/A\]](#)

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
- (b) Did you describe the limitations of your work? [\[Yes\]](#) Local analysis
- (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#)
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)

2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#)
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#)
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)? [\[Yes\]](#) See Appendix A.1
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) See Appendix A.1
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#) Estimated over 5 runs
 - (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)? [\[No\]](#)
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? [\[Yes\]](#)
 - (b) Did you mention the license of the assets? [\[N/A\]](#)
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[N/A\]](#)
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [\[N/A\]](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)
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? [\[N/A\]](#)
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)

A Experiment setup

We consider training the following models in the over-parameterized regime. In training, all explicit regularizations (including weight decay, dropout, data augmentation, batch normalization, learning rate decay) are removed, and a simple constant-LR SGD is used to train our models.

Small-scale models: In this case, we set the sample size to be particularly small on purpose, which allows interested readers to quickly reproduce these experiments on their own computer. Note that this choice does not impact our conclusions since we also discuss in details the influence of changing the extent of over-parameterization and provide larger-scale experiments.

- **Random feature model (RFM).** The inputs $\{x_i\}_{i=1}^n$ are drawn from $\mathcal{N}(0, I_d)$ with $d = 10, n = 200$. The labels are generated by

$$f^*(x) = 0.2x_1 + \left(\sum_{i=2}^d x_i - 1 \right)^2 / 3 + \sin \left(\sum_{i=1}^{d-2} x_i x_{i+2} / 4 \right).$$

The model is $f(x; \theta) = \sum_{j=1}^m \theta_j \text{ReLU}(w_j^T x)$ with $m = 2000$ and $\{w_j\}_{j=1}^m$ independently drawn from $N(0, I_d)$ at initialization and fixed during the training. This model is trained by SGD with learning rate $\eta = 0.003$ and batch size $B = 5$.

- **Linear networks.** The inputs $\{x_i\}_{i=1}^n$ drawn are drawn from $\mathcal{N}(0, I_d)$ with $d = 100, n = 50$. Here we set $n < d$ to examine the low-sample regime. The labels are generated by $f^*(x) = \sum_{i=1}^d x_i / d$. The model is a four-layer linear network: $d \rightarrow m \rightarrow m \rightarrow m \rightarrow 1$ with $m = 50$. This model is trained by SGD with learning rate $\eta = 0.1$ and batch size $B = 5$. The default LeCun initialization is used.
- **Fully-connected networks (FCN).** We randomly sample n data from the MNIST training set and label $\{1, 2, 3, 4, 5\}$ to 0 and $\{5, 6, 7, 8, 9, 10\}$ to 1 to form our new training set. The model is a ReLU-activated fully-connected network with the architecture: $784 \rightarrow m \rightarrow m \rightarrow 1$. Except for studying the influence of over-parameterization, we always set $m = 30, n = 1000$, and the model size is $p = 25441$. This model is trained by SGD with $\eta = 0.05, B = 5$ and the LeCun initialization is used.
- **Convolutional neural networks (CNN).** This is a small LeNet-type CNN [23] whose architecture is given in Table 1. The training set is the same as the one constructed above for FCN. The LeCun initialization is used and the model is trained by SGD with $\eta = 0.1, B = 5$.

Table 1: The architecture of CNN with m controlling the network width. We always set $m = 20$ except for studying the impact of over-parameterization, where we vary the value of m .

Layer	Output size
input	$28 \times 28 \times 1$
$3 \times 3 \times m$, conv	$28 \times 28 \times m$
$3 \times 3 \times 2m$, conv	$28 \times 28 \times 2m$
2×2 , avgpool	$14 \times 14 \times 2m$
$3 \times 3 \times 2m$, conv	$14 \times 14 \times 2m$
$3 \times 3 \times m$, conv	$14 \times 14 \times m$
2×2 , avgpool	$7 \times 7 \times m$
flatten	$49m$
$49m \rightarrow 1$, linear	1

Larger-scale models: For these experiments, the full CIFAR-10 dataset are used to train our model.

- **VGG nets.** The models are the standard VGG nets (including VGG-11, VGG-13, VGG-16, and VGG-19) for classifying CIFAR-10 proposed in [36].
- **ResNets.** The residual networks for CIFAR-10 proposed in [15] are considered and the models include ResNets-26, ResNet-38, ResNet-50, and ResNet-100. For ResNets, we follow [50] to use the fixup initialization in order to ensure that the model can be trained without batch normalization.

Both VGG nets and ResNets are trained by SGD with learning rate $\eta = 0.1$ and batch size $B = 64$ until the training loss becomes smaller than 10^{-4} .

Hyperparameter choices. In all experiments, the default model size, sample size, learning rate, and batch size described above are used unless explicitly specified, e.g., studying the influence of over-parameterization.

Efficient computations of the alignment factors and flatness. Let $g_i = \nabla f(x_i; \theta)$, $e_i = f(x_i; \theta) - y_i$. Let $Q = (g_1, \dots, g_n) \in \mathbb{R}^{p \times n}$ and $S = (e_1 g_1, \dots, e_n g_n) \in \mathbb{R}^{p \times n}$. Then,

$$G = \frac{1}{n} \sum_{i=1}^n g_i g_i^T = \frac{1}{n} Q Q^T \in \mathbb{R}^{p \times p}$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n e_i^2 g_i g_i^T - \left(\frac{1}{n} \sum_{i=1}^n e_i g_i \right) \left(\frac{1}{n} \sum_{i=1}^n e_i g_i \right)^T = \frac{1}{n} S P S^T \in \mathbb{R}^{p \times p},$$

where $P = I_{n \times n} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \in \mathbb{R}^{n \times n}$. Here $\mathbf{1} \in \mathbb{R}^n$ denotes the all-one vector.

Notice that the computation of $\alpha(\theta)$, $\beta(\theta)$, $\mu(\theta)$, $\|H(\theta)\|_F$ can be reduced to computing $\|G\|_F$, $\|\Sigma\|_F$, and $\text{Tr}(G\Sigma)$. The time complexity of naively computing them is on the order of $O(p^2 n)$, which is prohibitive for large-scale models where $p \gg n$. A more efficient way is to use the following equations

$$\begin{aligned} \|G\|_F^2 &= \frac{1}{n^2} \text{Tr}(Q Q^T Q Q^T) = \frac{1}{n^2} \|Q^T Q\|_F^2 \\ \|\Sigma\|_F^2 &= \frac{1}{n^2} \text{Tr}(S P S^T S P S^T) = \frac{1}{n^2} \|P S^T S\|_F^2 \\ \text{Tr}(G\Sigma) &= \frac{1}{n^2} \text{Tr}(Q Q^T S P S^T) = \frac{1}{n^2} \|P S^T Q\|_F^2, \end{aligned} \tag{10}$$

where all the matrices are $n \times n$. Hence, using these equations, the computation complexity becomes $O(n^2 p)$. This is much smaller than $O(p^2 n)$ when $p \gg n$.

- For small-scale experiments, the equations in (10) are directly used.
- For the large-scale models, we need further approximations since the computation complexity $O(n^2 p)$ is still prohibitive in this case. Notice that the formulations in Eq. (10) are all in the form of sample average, which allows us to perform Monte-Carlo approximation. Specifically, we randomly choose B samples from x_1, \dots, x_n and still use (10) to estimate these quantities. But now the computation complexity becomes $O(B^2 p)$. For the experiments on CIFAR-10, we test B 's with different values and find that $B = 50$ is sufficient to obtain a reliable approximation of the original full-data quantity. Hence, for all large-scale experiments in this paper, we use $B = 50$ to speed up the computation of alignment factors and flatness. We clarify that the models are still trained on the full dataset.

B Extra experiment results

B.1 Small-scale experiments

Figure 7 shows the relative norm of model gradient, i.e., $\chi_i(\theta)/\bar{\chi}(\theta)$ of each samples when the model converges. It is shown that $\gamma(\theta) = \min_i \chi_i(\theta)/\bar{\chi}(\theta)$ is indeed bounded below. This explains why the SGD noise satisfies the alignment property locally according to Lemma 2.4.

B.2 Larger-scale experiments

Figure 8 reports the values of alignment factors during the training for ResNet-26, ResNet-38, ResNet-50, VGG-13, and VGG-19. We see that the alignment factors are always significantly positive. In particular, for VGG-16, $\alpha(\theta)$ is close to 1 and $\mu(\theta)$ is roughly on the order of 0.1. In addition, Figure 9 shows that the noise-driven escape is indeed exponentially fast as suggested by our theoretical analysis.

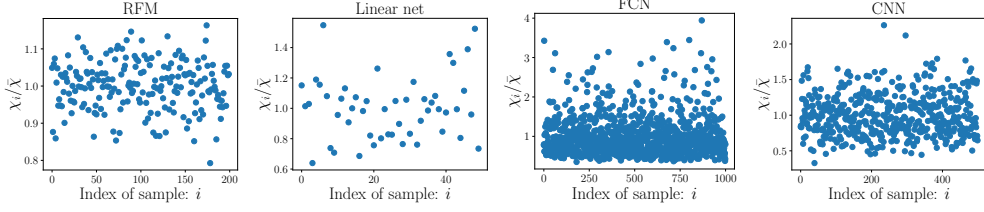


Figure 7: The model gradient norms of SGD solutions for the RFM, linear network, FCN, and CNN.

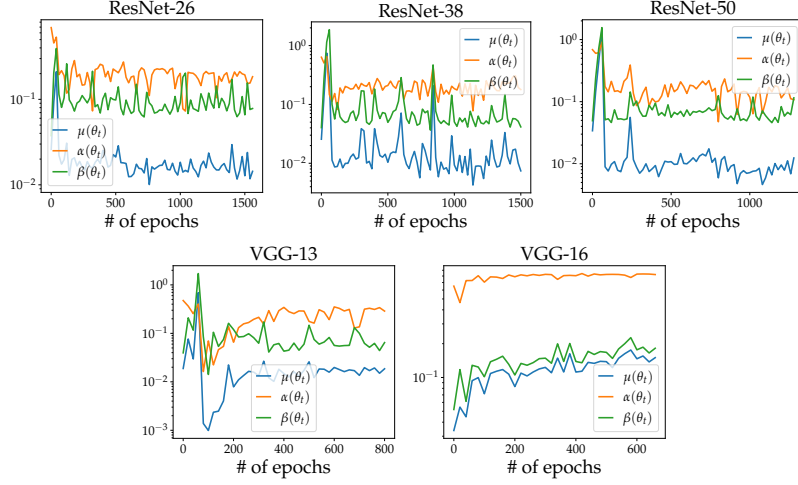


Figure 8: The alignment factors during the SGD training for classifying full CIFAR-10 dataset with VGG nets and ResNets .

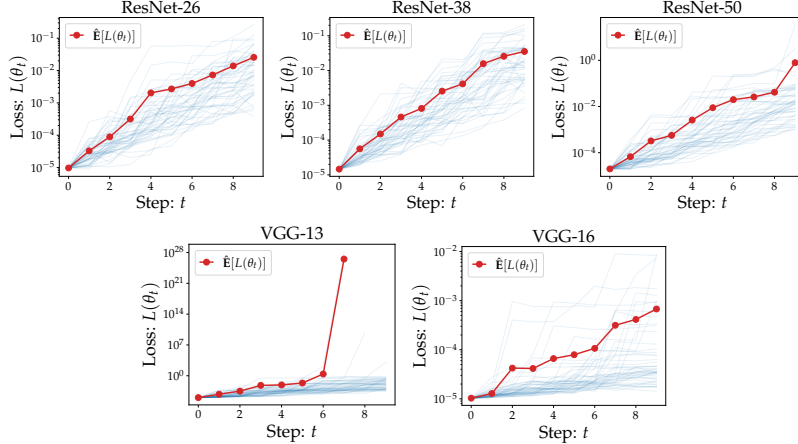


Figure 9: **The exponentially fast escape from sharp minima for training CIFAR-10 dataset.** The blue curves are 50 trajectories of SGD; the red curve corresponds to the average. The sharp minimum is found by SGD with $B = 64$ and $\eta = 0.1$. When it nearly converges, we switch to SGD with the same learning rate but a smaller batch size $B = 4$. This choice ensures that the escape is purely driven by SGD noise.

However, it should be stressed that the loss-scaled alignment factors for these larger-scale models are clearly smaller than the ones reported in small-scale experiments. In particular, one can see that the values of $\mu(\theta)$ for various ResNets are only the order of 0.01. To understand the underlying reason, we additionally examine the same models but for classifying a two-class subset of CIFAR-10 with the results reported in Figure 10. It is shown that in this case, the loss-scaled alignment factors are roughly on the order of 0.1 for all the models, significantly larger than the case of classifying the full CIFAR-10 dataset. It is also very surprising to observe that for VGG nets, the standard alignment

factor $\alpha(\theta)$'s are nearly 1, suggesting that the noise covariance completely aligns with the geometry of local landscape.

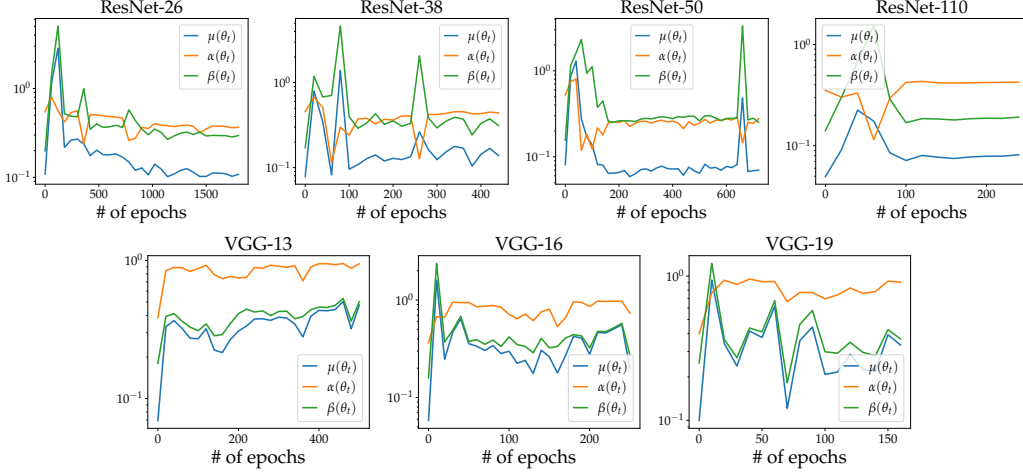


Figure 10: **The alignment factors during the SGD training for classifying a two-class subset of CIFAR-10 with VGG nets and ResNets .** In this experiment, we only pick the data of the class “0” and “1” from CIFAR-10 to train the model. It is shown that for this simpler problem, the alignment factors are significantly larger than the full-data case (see Figure 5a and 8) for all the models examined, regardless how over-parameterized the model is. In particular, we surprisingly observe that $\alpha(\theta)$'s are pretty close to 1 for all VGG nets, suggesting a complete alignment between the noise covariance and Fisher matrix. We leave to the detailed analysis of this unreasonable alignment to future work.

C Notations.

For a vector v , let $\|v\|_p = (\sum_i v_i^p)^{1/p}$ and $\hat{v} = v/\|v\|_2$. When $p = 2$, we omit the subscript for simplicity. For a matrix $A = (a_{i,j})$, denote by $\|A\|_F = (\sum_{i,j} a_{i,j}^2)^{1/2}$ the Frobenius norm. For a matrix or linear operator A , denote by $\{\lambda_i(A)\}_{i \geq 1}$ the eigenvalues of A in a non-increasing order. Let $\mathbb{S}^{d-1} = \{x \in \mathbb{R}^d \mid \|x\| = 1\}$, $r\mathbb{S}^{d-1} = \{x \in \mathbb{R}^d \mid \|x\| = r\}$, and $\tau_{d-1} = \text{Unif}(\mathbb{S}^{d-1})$. For any distribution ρ , let $\|f\|_{L_2(\rho)}^2 = \mathbb{E}_{x \sim \rho}[f^2(x)]$. We use $X \lesssim Y$ to indicate $X \leq CY$ for an absolute constant $C > 0$ and $X \gtrsim Y$ is defined similarly. Moreover, we use $X \sim Y$ to mean $cY \leq X \leq CY$ for some absolute constants $c, C > 0$. We also use C to denote an absolute constant, whose value may change from line to line. For a vector a and a SPD matrix W , let $\|a\|_W = \sqrt{a^T W a}$.

D Proof of Proposition 2.3

We first need the following lemma.

Lemma D.1. *Suppose $z \sim \mathcal{N}(0, S)$. Then, $\mathbb{E}_z[|v^T z|^2 z z^T] = 2S v v^T S + \|S^{1/2} v\|^2 S$.*

Proof. First assume $S = I_d$. Notice that

$$\mathbb{E}[|v^T z|^2 z_i z_j] = \mathbb{E}_z\left[\sum_{k,l=1}^d v_k v_l z_k z_l z_i z_j\right] = \begin{cases} 2v_i v_j & \text{if } i \neq j \\ 2v_i^2 + \sum_{k=1}^d v_k^2 & \text{if } i = j. \end{cases}$$

Therefore,

$$\mathbb{E}[|v^T z|^2 z z^T] = 2v v^T + \|v\|^2 I_d \quad (11)$$

For general S , let $x = S^{-1/2}z$. Then $x \sim \mathcal{N}(0, I_d)$. Then we have

$$\begin{aligned}\mathbb{E}[|v^T z|^2 z z^T] &= \mathbb{E}[|v^T S^{1/2} x| S^{1/2} z z^T S^{1/2}] \\ &= S^{1/2} \mathbb{E}[|v^T S^{1/2} x| z z^T] S^{1/2} \\ &= S^{1/2} \left(2S^{1/2} v v^T S^{1/2} + \|S^{1/2} v\|^2 I_d \right) S^{1/2} \\ &= 2S v v^T S + \|S^{1/2} v\|^2 S,\end{aligned}$$

where the third step follows from (11). \square

Now we start to prove Proposition 2.3. Let $\nabla F : \mathbb{R}^p \mapsto \mathbb{R}^d$ be the Jacobian matrix of F . Then $\nabla f(x; \theta) = \nabla F(\theta)^T x$. Recall that we assume $y = f(x; \theta^*)$, i.e., there is no label noise. Thus, we have

$$\begin{aligned}L(\theta) &= \frac{1}{2} \mathbb{E}_x[|(F(\theta) - F(\theta^*))^T x|^2] = \frac{1}{2} \|u(\theta)\|_S^2 \\ G(\theta) &= \mathbb{E}_x[\nabla F(\theta)^T x x^T \nabla F(\theta)] = \nabla F(\theta)^T S \nabla F(\theta) \\ H(\theta) &= G(\theta) + (F(\theta) - F(\theta^*)) \nabla^2 F(\theta),\end{aligned}$$

where $u(\theta) = F(\theta) - F(\theta^*)$. For the noise covariance, we have

$$\Sigma_2(\theta) = \nabla L(\theta) \nabla L(\theta)^T = \nabla F(\theta)^T S u(\theta) u(\theta)^T S \nabla F(\theta),$$

and

$$\begin{aligned}\Sigma_1(\theta) &= \mathbb{E}_x[|F(\theta)^T x - F(\theta^*)^T x|^2 \nabla F(\theta)^T x x^T \nabla F(\theta)] \\ &= \nabla F(\theta)^T \mathbb{E}_x[|(F(\theta) - F(\theta^*))^T x|^2 x x^T] \nabla F(\theta) \\ &= \nabla F(\theta)^T (2S u(\theta) u(\theta)^T S + \|u(\theta)\|_S^2 S) \nabla F(\theta) \\ &= 2\nabla F(\theta)^T S u(\theta) u(\theta)^T S \nabla F(\theta) + \|u(\theta)\|_S^2 \nabla F(\theta)^T S \nabla F(\theta) \\ &= 2\nabla L(\theta) \nabla L(\theta)^T + 2L(\theta) G(\theta),\end{aligned}$$

where the third step follows from Lemma D.1. Consequently, $\mu_1(\theta) \geq \mu(\theta) \geq 1$. \square

E Proof of Lemma 2.4

Recall the definition of feature-based model: $f(x; \theta) = \sum_{j=1}^m \theta_j \varphi_j(x) = \langle \theta, \Phi(x) \rangle$. Here $\Phi(x) \in \mathbb{R}^m$ denotes the feature of x . In this case, the Hessian and Fisher matrix are both constant but the SGD noise is still state-dependent. Let $g_i = \Phi(x_i)$ be the sample feature of x_i . Then,

$$\Sigma(\theta) = \frac{2}{n} \sum_{i=1}^n L_i(\theta) g_i g_i^T - \nabla L(\theta) \nabla L(\theta)^T, \quad G(\theta) = H(\theta) = \frac{1}{n} \sum_{i=1}^n g_i g_i^T.$$

Proof. Noticing $\bar{\chi} = \frac{1}{n} \sum_{i=1}^n g_i^T G g_i = \|G\|_F^2$, then

$$\begin{aligned}\text{Tr}(\Sigma_1(\theta) G) &= \frac{1}{n} \sum_{i=1}^n |g_i^T \theta|^2 \text{Tr}(g_i g_i^T G) \\ &= \frac{1}{n} \sum_{i=1}^n |g_i^T \theta|^2 \chi_i \geq \frac{\gamma \bar{\chi}}{n} \sum_{i=1}^n |g_i^T \theta|^2 = 2\gamma L(\theta) \|G\|_F^2.\end{aligned}$$

Thus $\mu_1(\theta) \geq \gamma$. In addition,

$$\begin{aligned}\mu_2(\theta) &= \frac{\nabla L(\theta)^T G \nabla L(\theta)}{2L(\theta) \|G\|_F^2} \\ &= \frac{\theta^T G^3 \theta}{\theta^T G \theta \|G\|_F^2} \leq \frac{\lambda_1^2(G)}{\|G\|_F^2},\end{aligned}$$

where the last step follows from Lemma G.1. \square

F Proof of Proposition 2.5

F.1 Rademacher complexity

We shall use the Rademacher complexity to bound the difference between empirical quantities and the corresponding population ones.

Definition F.1 (Rademacher complexity). *Let \mathcal{F} be a set of functions. The (empirical) Rademacher complexity of \mathcal{F} is defined as $\widehat{\text{Rad}}_n(\mathcal{F}) = \mathbb{E}_\xi[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \xi_i f(z_i)]$, where $\{\xi_i\}_{i=1}^n$ are i.i.d. random variables satisfying $\mathbb{P}(\xi_i = +1) = \mathbb{P}(\xi_i = -1) = \frac{1}{2}$.*

In particular, the following classic result will be frequently used, which is a restatement of [35, Theorem 26.5].

Theorem F.2. *Consider a function class \mathcal{F} and assume $|f| \leq B$. Then for any $\delta \in (0, 1)$, w.p. at least $1 - \delta$ over the choice of (z_1, z_2, \dots, z_n) , we have,*

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(z_i) - \mathbb{E}_z[f(z)] \right| \leq 4\widehat{\text{Rad}}_n(\mathcal{F}) + B\sqrt{\frac{2 \ln(2/\delta)}{n}}.$$

Lemma F.3 (Contraction property [24]). *Let $\phi : \mathbb{R} \mapsto \mathbb{R}$ be β -Lispchitz continuous and $\phi \circ \mathcal{F} = \{\phi \circ f : f \in \mathcal{F}\}$. Then, $\widehat{\text{Rad}}_n(\phi \circ \mathcal{F}) \leq \beta \widehat{\text{Rad}}_n(\mathcal{F})$.*

Lemma F.4. *Let $\mathcal{F} = \{u^T x : u \in \mathbb{S}^{d-1}\}$ be the linear class. Then $\widehat{\text{Rad}}_n(\mathcal{F}) \leq \sqrt{\frac{\sum_{i=1}^n \|x_i\|^2}{n^2}}$.*

Lemma F.5. *Let $\phi : \mathcal{X} \mapsto H$ be a feature map and H be a Hilbert space. Define $k(x, y) = \langle \phi_x, \phi_y \rangle_H$ and $\mathcal{H} = \{f(x) = \langle \phi_x, h \rangle_H \mid \|h\|_H \leq 1\}$. Then, $\widehat{\text{Rad}}_n(\mathcal{H}) \leq \sqrt{\sum_{i=1}^n k(x_i, x_i)}/n$.*

Proof. By the definition, we have

$$\begin{aligned} n\widehat{\text{Rad}}_n(\mathcal{H}) &= \mathbb{E}_\xi \sup_{\|h\|_H \leq 1} \sum_{i=1}^n \xi_i \langle \phi_{x_i}, h \rangle_H \\ &= \mathbb{E}_\xi \sup_{\|h\|_H \leq 1} \left\langle \sum_{i=1}^n \xi_i \phi_{x_i}, h \right\rangle_H = \mathbb{E}_\xi \left\| \sum_{i=1}^n \xi_i \phi_{x_i} \right\| \\ &\stackrel{(i)}{\leq} \sqrt{\mathbb{E}_\xi \left\| \sum_{i=1}^n \xi_i \phi_{x_i} \right\|^2} = \sqrt{\sum_{i=1}^n \langle \phi_{x_i}, \phi_{x_i} \rangle_H} = \sqrt{\sum_{i=1}^n k(x_i, x_i)}, \end{aligned}$$

where (i) follows from the Jensen's inequality. \square

For the proof of the above classic lemmas, we refer to [35]. The following lemma concerns the Rademacher complexity of the product of two function classes.

Lemma F.6. *Let \mathcal{F} and \mathcal{G} be two function classes. Suppose that $\sup_{f \in \mathcal{F}} \|f\|_\infty \leq A$ and $\sup_{g \in \mathcal{G}} \|g\|_\infty \leq B$. Define $\mathcal{F} * \mathcal{G} = \{f(x)g(x) : \mathcal{X} \mapsto \mathbb{R} : f \in \mathcal{F}, g \in \mathcal{G}\}$. Then,*

$$\widehat{\text{Rad}}_n(\mathcal{F} * \mathcal{G}) \leq (A + B)(\widehat{\text{Rad}}_n(\mathcal{F}) + \widehat{\text{Rad}}_n(\mathcal{G})).$$

Proof. By the definition of Rademacher complexity,

$$\begin{aligned}
n\widehat{\text{Rad}}_n(\mathcal{F} * \mathcal{G}) &= \mathbb{E}_\xi \left[\sup_{f \in \mathcal{F}, g \in \mathcal{G}} \sum_{i=1}^n f(x_i)g(x_i)\xi_i \right] \\
&= \mathbb{E}_\xi \left[\sup_{f \in \mathcal{F}, g \in \mathcal{G}} \sum_{i=1}^n \frac{(f(x_i) + g(x_i))^2}{4} \xi_i - \sum_{i=1}^n \frac{(f(x_i) - g(x_i))^2}{4} \xi_i \right] \\
&\leq \mathbb{E}_\xi \left[\sup_{f \in \mathcal{F}, g \in \mathcal{G}} \sum_{i=1}^n \frac{(f(x_i) + g(x_i))^2}{4} \xi_i \right] + \mathbb{E}_\xi \left[\sup_{f \in \mathcal{F}, g \in \mathcal{G}} \sum_{i=1}^n \frac{(f(x_i) - g(x_i))^2}{4} \xi_i \right] \\
&\stackrel{(i)}{\leq} \frac{A+B}{2} \left(\mathbb{E}_\xi \left[\sup_{f \in \mathcal{F}, g \in \mathcal{G}} \sum_{i=1}^n (f(x_i) + g(x_i)) \xi_i \right] + \mathbb{E}_\xi \left[\sup_{f \in \mathcal{F}, g \in \mathcal{G}} \sum_{i=1}^n (f(x_i) - g(x_i)) \xi_i \right] \right) \\
&\leq (A+B)n(\widehat{\text{Rad}}_n(\mathcal{F}) + \widehat{\text{Rad}}_n(\mathcal{G})),
\end{aligned}$$

where (i) follows from the Lemma F.3 and the fact that $t^2/4$ is $(A+B)/2$ Lipschitz continuous since $|f| \leq A, |g| \leq B$. \square

F.2 Results for general random feature models

Consider a random feature model (RFM):

$$f(x; \theta) = \frac{1}{\sqrt{m}} \sum_{i=1}^m \theta_i \varphi(x; w_i),$$

with $\{w_j\}_{j=1}^m \stackrel{iid}{\sim} \pi$ and $x \sim \rho$. Here $\varphi : \mathcal{X} \times \Omega \mapsto \mathbb{R}$ is an arbitrary parametric feature function. The scaling factor $m^{-1/2}$ is added only to ease taking the limit, which does not affect the noise structure. The associated kernel and kernel operator are given by $k(x, x') = \mathbb{E}_{w \sim \pi} [\sigma(w^T x) \sigma(w^T x')]$ and $\mathcal{K} : L_2(\rho) \mapsto L_2(\rho), \mathcal{K}u = \mathbb{E}_{x' \sim \rho} [k(\cdot, x') u(x')]$.

In this case, $g_i = m^{-1/2}(\varphi(x_i; w_1), \varphi(x_i; w_2), \dots, \varphi(x_i; w_m))^T \in \mathbb{R}^m$. Define the empirical kernel matrix: $\hat{k}(x, x') = \frac{1}{m} \sum_{s=1}^m \varphi(x; w_s) \varphi(x'; w_s)$. Then $g_i^T g_j = \hat{k}(x_i, x_j)$ and as $n, m \rightarrow \infty$, we have

$$\chi_i = \frac{1}{n} \sum_{j=1}^n (g_i^T g_j)^2 = \frac{1}{n} \sum_{j=1}^n \hat{k}^2(x_i, x_j) \rightarrow \mathbb{E}_x [k^2(x_i, x)]. \quad (12)$$

Therefore, in the limit, by our assumption $\chi_i \geq \chi := \inf_{x'} \mathbb{E}_x [k^2(x', x)]$ for any $i \in [n]$. The remaining issue is to transfer this to a non-asymptotic one.

We first make the following assumption on the feature function.

Assumption 1. Let $\Psi = \{\varphi(x; \cdot) \mid x \in \mathcal{X}\}$. Assume $\sup_{x \in \mathcal{X}, w \in \Omega} |\varphi(x; w)| \leq b$ and $\widehat{\text{Rad}}_n(\Psi) \leq b/\sqrt{n}$.

Note that the assumption $\widehat{\text{Rad}}_n(\Psi) = O(n^{-1/2})$ is satisfied by the popular feature function $\varphi(x; w) = \sigma(w^T x)$ with $\sigma : \mathbb{R} \mapsto \mathbb{R}$ being a Lipschitz activation function.

Proposition F.7. Suppose $m \geq n$. Let $\chi(x) = \mathbb{E}_{x' \sim \rho} [k^2(x, x')]$ and $\bar{\chi} = \mathbb{E}_{x \sim \rho} [\chi(x)]$. For any $\delta \in (0, 1/e)$, w.p. larger than $1 - \delta$ over the sampling of data and random features, we have $\mu_1(\theta) \geq \inf_x \frac{\chi(x)}{\bar{\chi}} - \epsilon_n$, $\mu_2(\theta) \leq \tau(\mathcal{K}) + \epsilon_n$, and $\mu(\theta) \geq \inf_x \frac{\chi(x)}{\bar{\chi}} - \tau(\mathcal{K}) - 2\epsilon_n$, where $\tau(\mathcal{K}) = \lambda_1^2(\mathcal{K}) / \sum_j \lambda_j^2(\mathcal{K})$ and $\epsilon_n = Cb^2 \sqrt{\frac{\log(1/\delta)}{n}}$.

By this proposition, ensuring the alignment property only needs $\inf_x \chi(x)/\bar{\chi} > 0$. This condition is very mild and satisfied by most popular kernels, e.g., the dot-product kernel, which appears naturally in the analysis of neural nets [18]. Next we proceed to prove this proposition.

We first need the following lemmas.

Lemma F.8. Let $\Phi_1 = \{k(\cdot, x) \mid x \in \mathcal{X}\}$. Then, $\widehat{\text{Rad}}_n(\Phi_1) \leq 1/\sqrt{n}$.

Proof. Denote by \mathcal{H}_k the reproducing kernel Hilbert space (RKHS). Then, by the Moore-Aronsjajn theorem [2],

$$\|k(\cdot, x)\|_{\mathcal{H}_k}^2 = \langle k(\cdot, x), k(\cdot, x) \rangle_{\mathcal{H}_k} = k(x, x) \leq 1.$$

Therefore, Φ_1 is a subset of the unit ball of \mathcal{H}_k , for which it is well-known that the Rademacher complexity is bounded by $\sqrt{\sum_{i=1}^n k(x_i, x_i)/n} \leq \sqrt{1/n}$. \square

Lemma F.9. For any $\delta \in (0, 1)$, w.p. at least $1 - \delta$, we have

$$\sup_{x, x' \in \mathcal{X}} |k(x, x') - \hat{k}(x, x')| \leq b^2 \sqrt{\frac{\log(1/\delta)}{m}}.$$

Proof. Let $\Psi_2 = \{\varphi(x; \cdot) \varphi(x'; \cdot) \mid x, x' \in \mathcal{X}\}$. Notice that

$$\varphi(x; w) \varphi(x'; w) = \frac{(\varphi(x; w) + \varphi(x'; w))^2}{4} - \frac{(\varphi(x; w) - \varphi(x'; w))^2}{4}.$$

Then by the same argument in the proof of Lemma F.6, we can obtain $\widehat{\text{Rad}}_n(\Psi_2) \lesssim b \widehat{\text{Rad}}_n(\Psi) \lesssim b^2/\sqrt{m}$, where the last step follows from Lemma 1. Then, applying Theorem F.2, we complete the proof. \square

Prove the first part of Proposition F.7 Notice that

$$\begin{aligned} |\chi(x_i) - \chi_i| &= |\mathbb{E}_x[k^2(x, x_i)] - \frac{1}{n} \sum_{j=1}^n k^2(x_j, x_i)| + |\frac{1}{n} \sum_{j=1}^n k^2(x_j, x_i) - \frac{1}{n} \sum_{j=1}^n \hat{k}(x_j, x_i)| \\ &\leq \sup_{x'} |\mathbb{E}_x[k^2(x, x')] - \frac{1}{n} \sum_{j=1}^n k^2(x_j, x')| + \sup_{x, x'} |k^2(x, x') - \hat{k}^2(x, x')| \end{aligned} \quad (13)$$

Let $\Phi_2 = \{k^2(x, \cdot) \mid x \in \mathcal{X}\}$. Since $\sup_{x, x'} |k(x, x')| \leq 1$, the Ledoux-Talagrand inequality (Lemma F.3) implies that $\widehat{\text{Rad}}_n(\Phi_2) \leq 0.5 \widehat{\text{Rad}}_n(\Phi_1) \leq 1/\sqrt{n}$, where the last inequality follows from Lemma F.8. Then, applying Theorem F.2, for any $\delta \in (0, 1/e)$, we have w.p. $1 - \delta$ that

$$\begin{aligned} \sup_{x'} |\mathbb{E}_x[k^2(x, x')] - \frac{1}{n} \sum_{j=1}^n k^2(x_j, x')| &\lesssim 2 \widehat{\text{Rad}}_n(\Phi_2) + \sqrt{\frac{\log(2/\delta)}{n}} \\ &\leq \frac{2}{\sqrt{n}} + \sqrt{\frac{\log(2/\delta)}{n}} \lesssim \sqrt{\frac{\log(2/\delta)}{n}}. \end{aligned} \quad (14)$$

Let $\Delta = \sup_{x, x'} |k(x, x') - \hat{k}(x, x')|$. Therefore,

$$|k^2(x, x') - \hat{k}^2(x, x')| = |k(x, x') + \hat{k}(x, x')| |k(x, x') - \hat{k}(x, x')| \leq (2k(x, x') + \Delta) \Delta \lesssim \Delta.$$

Applying Lemma F.9 and together with (14), we have

$$|\chi(x_i) - \chi_i| \lesssim \sqrt{\frac{\log(1/\delta)}{n}} + b \sqrt{\frac{\log(1/\delta)}{m}} =: \epsilon_n.$$

On the other hand, by the Hoeffding's inequality, w.p. at least $1 - \delta$, we have

$$\tilde{\chi} = \frac{1}{n} \sum_{i=1}^n \chi_i \leq \frac{1}{n} \sum_{i=1}^n \chi(x_i) + C\epsilon_n \leq \mathbb{E}_x[\chi(x)] + \sqrt{\frac{\log(1/\delta)}{n}} + C\epsilon_n \leq \bar{\chi} + C\epsilon_n.$$

By Proposition 2.2, we have

$$\mu_1(\theta) \geq \frac{\inf_i \chi_i}{\tilde{\chi}} \geq \frac{\inf \chi(x) - C\epsilon_n}{\bar{\chi} + C\epsilon_n} \geq \frac{\inf \chi(x)}{\bar{\chi}} - C\epsilon_n.$$

\square

The above proves the first part of Proposition F.7. We now turn to the second part. We will frequently use the following McDiarmid's inequality.

Theorem F.10 (McDiarmid's inequality). *Let X_1, \dots, X_n are i.i.d. random variables. Assume for all $i \in [n]$ and $x_1, \dots, x_n, \tilde{x}_i \in \mathcal{X}$ that*

$$|f(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_{i-1}, \tilde{x}_i, x_{i+1}, \dots, x_n)| \leq D_i.$$

Let $\sigma^2 := \frac{1}{4} \sum_{i=1}^n D_i^2$. Then, for any $\delta \in (0, 1)$, w.p. at least $1 - \delta$ over the sampling of X_1, \dots, X_n , we have

$$|f(X_1, \dots, X_n) - \mathbb{E}[f]| \leq \sqrt{2 \log(2/\delta) \sigma}.$$

We define two matrices

- The kernel matrix: $K = (K_{i,j}) \in \mathbb{R}^{n \times n}$ with $K_{i,j} = \frac{1}{n} k(x_i, x_j)$;
- The approximate kernel matrix: $\hat{K} = (\hat{K}_{i,j}) \in \mathbb{R}^{n \times n}$ with

$$\hat{K}_{i,j} = \frac{1}{nm} \sum_{s=1}^m \varphi(x_i; w_s) \varphi(x_j; w_s) =: \frac{1}{n} \hat{k}(x_i, x_j).$$

Note that the approximate kernel matrix is the normalized Fisher matrix, i.e., $\hat{K} = G/n$. We will frequently use the following inequality:

$$\mathbb{E}_w |\hat{k}(x, x') - k(x, x')|^2 \leq \frac{\mathbb{E}_w [\varphi^2(x; w) \varphi^2(x'; w)]}{m} \leq \frac{b^4}{m}. \quad (15)$$

Bounding the largest eigenvalue.

Lemma F.11. *There exists a constant $c > 0$ such that for any $\delta \in (0, 1)$, w.p. at least $1 - \delta$ over the sampling of random features, we have $\lambda_1(\hat{K}) \leq \lambda_1(K) + Cb^2 \sqrt{\log(2/\delta)/m}$.*

Proof. Let $\Delta(w_1, \dots, w_m) = \|K - \hat{K}\|_F$. By Jensen's inequality, we have

$$\mathbb{E}[\Delta] \leq \sqrt{\mathbb{E}[\Delta^2]} \leq \sqrt{\mathbb{E}\left[\frac{1}{n^2} \sum_{i,j=1}^n |k(x_i, x_j) - \hat{k}(x_i, x_j)|^2\right]} \lesssim \sqrt{\frac{b^4}{m}},$$

where the last inequality follows from (15). Denote by \hat{K}' the approximate kernel matrix associated with $(w_1, \dots, w'_s, \dots, w_m)$. Then,

$$\begin{aligned} D_s &= |\Delta(w_1, \dots, w_j, \dots, w_m) - \Delta(w_1, \dots, w'_j, \dots, w_m)| \leq \|\hat{K} - \hat{K}'\|_F \\ &= \sqrt{\frac{1}{n^2} \sum_{i,j=1}^n \frac{1}{m^2} (\varphi(x_i; w_s) \varphi(x_j; w_s) - \varphi(x_i; w'_s) \varphi(x_j; w'_s))^2} \lesssim \frac{b^2}{m}. \end{aligned}$$

Therefore $\sigma^2 = \frac{1}{4} \sum_{s=1}^m D_s^2 \lesssim b^4/m$. By McDiarmid's inequality (Theorem F.10), for any $\delta \in (0, 1)$, w.p. at least $1 - \delta$ over the sampling of random features, we have

$$\Delta \lesssim \mathbb{E}[\Delta] + d \sqrt{\frac{\log(2/\delta)}{m}} \leq b^2 \sqrt{\frac{\log(2/\delta)}{m}}.$$

Plugging the above estimate to the Weyl's inequality: $\lambda_1(\hat{K}) \leq \lambda_1(K) + \|K - \hat{K}\|_2$, we complete the proof. \square

Lemma F.12. *Suppose $k(\cdot, \cdot)$ to be positive semi-definite kernel and let $\phi : \mathcal{X} \mapsto \mathcal{H}$ be a feature map satisfying $k(x, y) = \langle \phi_x, \phi_y \rangle_{\mathcal{H}}$. Then,*

$$\lambda_1(\mathcal{K}) = \sup_{\|h\|_{\mathcal{H}}=1} \mathbb{E}_x [\langle h, \phi_x \rangle_{\mathcal{H}}^2]. \quad (16)$$

Proof. By the variational principle of eigenvalues, we have

$$\begin{aligned} \lambda_1(\mathcal{K}) &= \sup_{\|u\|_{L_2(\rho)}=1} \mathbb{E}_{x,y} [k(x, y) u(x) u(y)] = \sup_{\|u\|_{L_2(\rho)}=1} \mathbb{E}_{x,y} [\langle \phi_x, \phi_y \rangle_{\mathcal{H}} u(x) u(y)] \\ &= \sup_{\|u\|_{L_2(\rho)}=1} \|\mathbb{E}_x [u(x) \phi_x]\|_{\mathcal{H}}^2 = \sup_{\|u\|_{L_2(\rho)}=1} \sup_{\|h\|_{\mathcal{H}}=1} \langle h, \mathbb{E}_x [u(x) \phi_x] \rangle_{\mathcal{H}}^2 \\ &= \sup_{\|h\|_{\mathcal{H}}=1} \sup_{\|u\|_{L_2(\rho)}=1} \mathbb{E}_x [u(x) \langle h, \phi_x \rangle_{\mathcal{H}}]^2 = \sup_{\|h\|_{\mathcal{H}}=1} \mathbb{E}_x [\langle h, \phi_x \rangle_{\mathcal{H}}^2]. \end{aligned}$$

\square

Lemma F.13. For any $\delta \in (0, 1/e)$, w.p. at least $1 - \delta$ over the sampling of data, we have

$$\lambda_1(K) \leq \lambda_1(\mathcal{K}) + C \sqrt{\frac{\log(2/\delta)}{n}}.$$

Proof. By Lemma F.12, we have

$$\lambda_1(\mathcal{K}) = \sup_{\|h\|_{\mathcal{H}}=1} \mathbb{E}_x[\langle \phi_x, h \rangle_{\mathcal{H}}^2], \quad \lambda_1(K) = \sup_{\|h\|_{\mathcal{H}}=1} \hat{\mathbb{E}}_x[\langle \phi_x, h \rangle_{\mathcal{H}}^2]. \quad (17)$$

Let $\mathcal{F}_1 = \{f(x) := \langle h, \sigma_x \rangle_{\mathcal{H}} \mid \|h\|_{\mathcal{H}} \leq 1\}$ and $\mathcal{F}_2 = \{f^2 \mid f \in \mathcal{H}_1\}$. Then, we have: (1) for any $f \in \mathcal{H}_1$, $|f(x)| \leq \|h\|_{\mathcal{H}} \|\sigma_x\|_{\mathcal{H}} = \sqrt{k(x, x)} \lesssim 1$; (2) $\widehat{\text{Rad}}_n(\mathcal{F}_1) = \sqrt{\sum_{i=1}^n k(x_i, x_i)/n} \lesssim 1/\sqrt{n}$ by Lemma F.5. Using contraction inequality (Lemma F.3), we have

$$\widehat{\text{Rad}}_n(\mathcal{F}_2) \lesssim \widehat{\text{Rad}}_n(\mathcal{F}_1) \lesssim 1/\sqrt{n}.$$

Using Theorem F.2, for any $\delta \in (0, 1/e)$, w.p. larger than $1 - \delta$ over the sampling of data, we have

$$\sup_{\|h\|_{\mathcal{H}} \leq 1} |\hat{\mathbb{E}}[\langle h, \sigma_{x_i} \rangle^2] - \mathbb{E}[\langle h, \sigma_x \rangle^2]| \lesssim \widehat{\text{Rad}}_n(\mathcal{F}_2) + \sqrt{\frac{\log(1/\delta)}{n}} \leq \sqrt{\frac{\log(1/\delta)}{n}}. \quad (18)$$

For any $\varepsilon > 0$, let $\hat{h} \in \mathcal{H}$ such that $\lambda_1(K) \leq \hat{\mathbb{E}}[\langle \phi_x, \hat{h} \rangle_{\mathcal{H}}^2] + \varepsilon$. Then, using (18), we have

$$\begin{aligned} \lambda_1(K) &\leq \hat{\mathbb{E}}[\langle \phi_x, \hat{h} \rangle_{\mathcal{H}}^2] + \varepsilon = \mathbb{E}[\langle \phi_x, \hat{h} \rangle_{\mathcal{H}}^2] + (\hat{\mathbb{E}}[\langle \phi_x, \hat{h} \rangle_{\mathcal{H}}^2] - \mathbb{E}[\langle \phi_x, \hat{h} \rangle_{\mathcal{H}}^2]) + \varepsilon \\ &\lesssim \lambda_1(\mathcal{K}) + \sqrt{\frac{\log(1/\delta)}{n}} + \varepsilon. \end{aligned}$$

Taking $\varepsilon \rightarrow 0$ completes the proof. \square

Lemma F.14. For any $\delta \in (0, 1)$, w.p. at least $1 - \delta$ over the sampling of data and random features, we have

$$\lambda_1(\hat{K}) \leq \lambda_1(\mathcal{K}) + C \left(\sqrt{\frac{\log(1/\delta)}{n}} + b^2 \sqrt{\frac{\log(2/\delta)}{m}} \right).$$

Proof. The conclusion directly follows from the combination of Lemma F.13 and F.11. \square

Bounding the Frobenius norm. The following lemma provide a lower bound of the Frobenius norm of the approximate kernel matrix.

Lemma F.15. For any $\delta \in (0, 1/e)$, w.p. at least $1 - \delta$ over the sampling of data and random features, we have

$$\|\hat{K}\|_F^2 \geq \sum_i \lambda_i^2(\mathcal{K}) - C \left(b^2 \sqrt{\frac{\log(1/\delta)}{m}} + \sqrt{\frac{\log(1/\delta)}{n}} \right).$$

Proof. Denote by $\{(\lambda_s, v_s)\}_{s \geq 1}$ the eigen-pairs of the kernel $k(\cdot, \cdot)$, and thus $k(x, x') = \sum_s \lambda_s v_s(x) v_s(x')$. Let x, x' be independently drawn from ρ . Then,

$$\begin{aligned} \mathbb{E}[k^2(x, x')] &= \sum_{s_1, s_2} \lambda_{s_1} \lambda_{s_2} \mathbb{E}[v_{s_1}(x) v_{s_2}(x) v_{s_1}(x') v_{s_2}(x')] \\ &= \sum_{s_1, s_2} \lambda_{s_1} \lambda_{s_2} \mathbb{E}[v_{s_1}(x) v_{s_2}(x)] \mathbb{E}[v_{s_1}(x') v_{s_2}(x')] \\ &= \sum_{s_1, s_2} \lambda_{s_1} \lambda_{s_2} \delta_{s_1, s_2} = \sum_s \lambda_s^2. \end{aligned} \quad (19)$$

Using the above equality, we have

$$\begin{aligned} \mathbb{E}[\|K\|_F^2] &= \frac{1}{n^2} \sum_{i, j=1}^n \mathbb{E}[k(x_i, x_j)^2] = \frac{1}{n} \mathbb{E}[k^2(x, x)] + \frac{1}{n^2} \sum_{i \neq j} \mathbb{E}[k^2(x_i, x_j)] \\ &= \sum_s \lambda_s^2 + \frac{\mathbb{E}_x[k^2(x, x)] - \mathbb{E}_{x, x'}[k^2(x, x')]}{n} \geq \sum_s \lambda_s^2. \end{aligned} \quad (20)$$

Here, the last inequality is due to $\mathbb{E}[k^2(x, x)] \geq \mathbb{E}[k^2(x, x')]$ as explained as follows. Notice that for any $x, x' \in \mathcal{X}$, $k^2(x, x') \leq k(x, x)k(x', x')$ (See, e.g., [34, Lemma 35]). Therefore, $\mathbb{E}[k^2(x, x')] \leq \mathbb{E}[k(x, x)]\mathbb{E}[k(x', x')] = (\mathbb{E}[k(x, x)])^2 \leq \mathbb{E}[k^2(x, x)]$. The last inequality follows from that $(\mathbb{E}[X])^2 \leq \mathbb{E}[X^2]$ holds for any random variable X .

Let K' be the kernel matrix corresponding to the $(x_1, \dots, \tilde{x}_i, \dots, x_n)$. Then,

$$\begin{aligned} D_i &= |||K||_F^2 - |||K'||_F^2| \\ &= \left| \frac{1}{n^2} \sum_{j \neq i} (k^2(x_i, x_j) - k^2(\tilde{x}_i, x_j)) + \frac{1}{n^2} (k^2(x_i, x_i) - k^2(\tilde{x}_i, \tilde{x}_i)) \right| \lesssim \frac{1}{n}. \end{aligned}$$

Therefore, $\sigma^2 = \frac{1}{4} \sum_{i=1}^n |D_i|^2 \lesssim 1/n$. Using Theorem F.10 and Eq. (20), we have w.p. at least $1 - \delta$ over the sampling of data that

$$|||K||_F^2 \geq \mathbb{E}[|||K||_F^2] - C \sqrt{\frac{\log(1/\delta)}{n}} \geq \sum_s \lambda_s^2 - C \sqrt{\frac{\log(1/\delta)}{n}}. \quad (21)$$

In addition, following the proof of Lemma F.11, we have for any $\delta \in (0, 1/e)$, w.p. larger than $1 - \delta$ over the sampling of random features that $||\hat{K} - K||_F \leq d\sqrt{\log(2/\delta)/m}$. Thus,

$$|||\hat{K}||_F^2 - |||K||_F^2| \leq (||\hat{K}||_F + ||K||_F)(||\hat{K}||_F - ||K||_F) \lesssim b^2 ||\hat{K} - K||_F \leq b^2 \sqrt{\log(2/\delta)/m}.$$

Combining with (21), we complete the proof. \square

Prove the second part of Proposition F.7. Recall $\epsilon_n = \sqrt{\log(1/\delta)/n} + b^2 \sqrt{\log(1/\delta)/m}$. Combining Lemma F.14 and F.15, we have

$$\mu_2(\theta) = \frac{\lambda_1^2(\hat{K})}{||\hat{K}||_F^2} \leq \frac{(\lambda_1(\mathcal{K}) + C\epsilon_n)^2}{\sum_i \lambda_i^2(\mathcal{K}) - C\epsilon_n} \leq \frac{\lambda_1^2(\mathcal{K})}{\sum_i \lambda_i^2(\mathcal{K})} + C\epsilon_n. \quad (22)$$

F.3 Proof of Proposition 2.5

Denote by τ_{d-1} the uniform distribution over the unit sphere $\mathbb{S}^{d-1} = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}$. According to [5], the eigenfunctions of \mathcal{K} are the spherical harmonics. In particular, the eigenfunction corresponding to the largest eigenvalue is the first spherical harmonics: $Y_1(x) \equiv 1$. Therefore,

$$\lambda_1(\mathcal{K}) = \lambda_1(\mathcal{K})Y_1(x) = \mathbb{E}_{x' \sim \tau_{d-1}}[\kappa(x^T x')Y_1(x')] = \mathbb{E}_{x' \sim \tau_{d-1}}[\kappa(x'_1)],$$

where the last equality uses the rotational symmetry of τ_{d-1} . Moreover, by (19), $\sum_i \lambda_i^2(\mathcal{K}) = \mathbb{E}_{x, x'}[\kappa^2(x^T x')] = \mathbb{E}_x[\kappa^2(x_1)]$.

For the ReLU activation function, [7] shows that $k(x, x') = \kappa(x^T x')$ with

$$\kappa(z) = \frac{\sqrt{1 - z^2} + (\pi - \arccos(z))z}{2\pi}.$$

The eigenvalues of this kernel has been derived in [5, 45]. Specifically, we have

$$\lambda_1(\mathcal{K}) \sim 1, \quad \sum_{i=2}^{\infty} \lambda_i^2(\mathcal{K}) \sim \frac{1}{d}.$$

Hence,

$$\frac{\lambda_1^2(\mathcal{K})}{\sum_{i=1}^{\infty} \lambda_i^2(\mathcal{K})} \sim \frac{1}{1 + \frac{1}{d}} \leq 1 - \frac{C}{d}.$$

Then, applying Proposition F.7, we complete the proof. \square

G Proofs of Section 3

We first need the following technical lemma.

Lemma G.1. For $a, b > 0$, let $g(a, b, \theta) = -a \frac{\theta^T H^2 \theta}{\theta^T H \theta} + b \frac{\theta^T H^3 \theta}{\theta^T H \theta}$. Then, $\inf_{\theta} g(a, b, \theta) \geq -a^2/(4b)$.

Proof. Let $u = H^{1/2} \theta / \|H^{1/2} \theta\|$ and $H = \sum_j \lambda_j e_j e_j^T$ the eigen-decomposition of H . Suppose $u = \sum_j s_j e_j$. Then $\sum_j s_j^2 = 1$ and

$$\begin{aligned} -a \frac{\theta^T H^2 \theta}{\theta^T H \theta} + b \frac{\theta^T H^3 \theta}{\theta^T H \theta} &= -a u^T H u + b u^T H^2 u = \sum_j (b \lambda_j^2 - a \lambda_j) s_j^2 \\ &\geq \inf_{\lambda \geq 0} (b \lambda^2 - a \lambda) = \inf_{\lambda \geq 0} \left(b \left(\lambda - \frac{a}{2b} \right)^2 - \frac{a^2}{4b} \right) \geq -\frac{a^2}{4b}. \end{aligned}$$

□

We first consider GD, for which the stability only imposes the flatness constraint: $\lambda_1(H) \leq 2/\eta$. This means that the flatness seen by GD is only the largest eigenvalue of Hessian.

Lemma G.2. (1) $\inf_{\theta} r(\theta) \geq 0$; (2) $\sup_{\theta} r(\theta) \leq 1$ if $\eta \leq 2/\lambda_1(H)$.

Proof. Recall that $r(\theta) = 1 - 2\eta \frac{\theta^T H^2 \theta}{\theta^T H \theta} + \eta^2 \frac{\theta^T H^3 \theta}{\theta^T H \theta}$. Let $u = H^{1/2} \theta / \|H^{1/2} \theta\|$. Then we have

$$r(\theta) = 1 - 2\eta u^T H u + \eta^2 u^T H^2 u = \|\eta H u - u\|^2 \geq 0.$$

Let $u = \sum_j s_j e_j$ with $\{e_j\}_j$ being the eigenvectors of H . Then $\sum_j s_j^2 = 1$ due to $\|u\| = 1$. By the assumption, $(\eta \lambda_j - 1)^2 \leq 1$ for all $j \in [n]$. Then,

$$r(\theta) = \left\| \sum_j \eta \lambda_j s_j e_j - \sum_j s_j e_j \right\|^2 = \sum_j (\eta \lambda_j - 1)^2 s_j^2 \leq \sum_j s_j^2 = 1. \quad (23)$$

□

Proof of Proposition 3.4 Using $\mu_1(\theta) \geq \mu_1$ and $\Sigma(\theta) = \Sigma_1(\theta) - \Sigma_2(\theta)$, we have

$$\begin{aligned} \nu(\theta) &= \frac{1}{2B} \text{Tr}(H \Sigma(\theta)) = \frac{1}{2B} \text{Tr}(H \Sigma_1(\theta)) - \frac{1}{2B} \text{Tr}(H \Sigma_2(\theta)) \\ &\geq \frac{\mu_1 \|H\|_F^2}{B} L(\theta) - \frac{1}{2B} \nabla L(\theta)^T H \nabla L(\theta) = L(\theta) \left(\frac{\mu_1 \|H\|_F^2}{B} - \frac{\theta^T H^3 \theta}{B \theta^T H \theta} \right), \end{aligned}$$

where the contribution of $\Sigma_2(\theta)$ is $-\theta^T H^3 \theta / (B \theta^T H \theta)$. Notice that there is a similar term in the contribution of mean gradient $r(\theta)$. Together we have $\mathbb{E}[L(\theta_{t+1})] \geq \mathbb{E}[L(\theta_t) \gamma(\theta)]$ with

$$\begin{aligned} \gamma(\theta) &\geq 1 - 2\eta \frac{\theta^T H^2 \theta}{\theta^T H \theta} + \eta^2 \frac{\theta^T H^3 \theta}{\theta^T H \theta} + \frac{\mu_1 \eta^2}{B} \|H\|_F^2 - \frac{\eta^2 \theta^T H^3 \theta}{B \theta^T H \theta} \\ &= 1 - 2\eta \frac{\theta^T H^2 \theta}{\theta^T H \theta} + \eta^2 \left(1 - \frac{1}{B} \right) \frac{\theta^T H^3 \theta}{\theta^T H \theta} + \frac{\mu_1 \eta^2}{B} \|H\|_F^2 \end{aligned} \quad (24)$$

By Lemma G.1, we have

$$\gamma(\theta) \geq 1 - \frac{4\eta^2}{4\eta^2(1-B^{-1})} + \frac{\mu_1 \eta^2}{B} \|H\|_F^2 =: \gamma_0.$$

Let $\gamma_0 \leq 1$ yields $\|H\|_F^2 \leq B/(\eta \sqrt{(B-1)\mu_1})$. This bound is trivial for the case $B = 1$, where

$$\gamma(\theta) \geq 1 - 2\eta \frac{\theta^T H^2 \theta}{\theta^T H \theta} + \frac{\mu_1 \eta^2}{B} \|H\|_F^2 \geq 1 - 2\eta \lambda_1(H) + \frac{\mu_1 \eta^2}{B} \|H\|_F^2.$$

Noticing that $\|H\|_F^2 = \sum_j \lambda_j^2(H)$, the stability condition needs $1 - 2\eta \lambda_1(H) + \frac{\mu_1 \eta^2}{B} \|H\|_F^2 \leq 1$, leading to

$$\sum_j \lambda_j^2(H) \leq \frac{2B}{\mu_1 \eta} u^T H u \leq \frac{2B}{\mu_1 \eta} \lambda_1(H).$$

Thus, $\lambda_1(H) \leq 2B/(\mu_1 \eta)$. Consequently, $\|H\|_F^2 \leq \frac{2B}{\mu_1 \eta} \lambda_1(H) \leq (\frac{2B}{\mu_1 \eta})^2$.

Combing them, we complete the proof. □