# The Anisotropic Noise in Stochastic Gradient Descent: Its Behavior of Escaping from Sharp Minima and Regularization Effects

Zhanxing Zhu [* 1 2 3]   Jingfeng Wu [* 1]   Bing Yu [1]   Lei Wu [1]   Jinwen Ma [1]

## Abstract

Understanding the behavior of stochastic gradient descent (SGD) in the context of deep neural networks has raised lots of concerns recently. Along this line, we study a general form of gradient based optimization dynamics with unbiased noise, which unifies SGD and standard Langevin dynamics. Through investigating this general optimization dynamics, we analyze the behavior of SGD on escaping from minima and its regularization effects. A novel indicator is derived to characterize the efficiency of escaping from minima through measuring the alignment of noise covariance and the curvature of loss function. Based on this indicator, two conditions are established to show which type of noise structure is superior to isotropic noise in term of escaping efficiency. We further show that the anisotropic noise in SGD satisfies the two conditions, and thus helps to escape from sharp and poor minima effectively, towards more stable and flat minima that typically generalize well. We systematically design various experiments to verify the benefits of the anisotropic noise, compared with full gradient descent plus isotropic diffusion (i.e. Langevin dynamics).

## 1. Introduction

As a successful learning algorithm, stochastic gradient descent (SGD) was originally adopted for dealing with the computational bottleneck of training neural networks with large-scale datasets (Bottou, 1991). Its empirical efficiency and effectiveness have attracted lots of attention. Besides the aspect of empirical efficiency, recently, researchers started to

---
[*]Equal contribution  [1]School of Mathematical Sciences, Peking University, Beijing, China [2]Center for Data Science, Peking University, Beijing, China [3]Beijing Institute of Big Data Research, Beijing, China. Correspondence to: Zhanxing Zhu <zhanxing.zhu@pku.edu.cn>, Jingfeng Wu <pkuwjf@pku.edu.cn>.
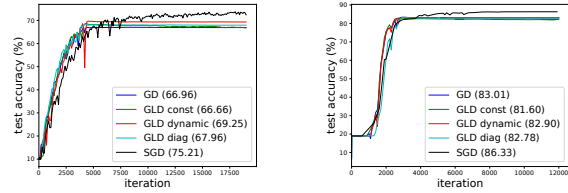
*Figure 1.* The generalization performance of dynamics in Table 1. The noise magnitude of SGD, GLD dynamic and GLD diag is tuned to be the same for fair comparison. The noise of GLD constant is tunded to the best. **Left**: SVHN. We only use $2,5000$ examples for training to compromise with the computational burden; **Right**: CIFAR-10. The model is VGG-11 since it achieves decent performance without using batch normalization, which causes uncontrollable affects for analyzing SGD.

analyze the optimization behaviors of SGD and its impacts on generalization.

The optimization properties of SGD have been studied from various perspectives. The convergence behaviors of SGD for simple one hidden layer neural networks were investigated in (Li & Yuan, 2017; Brutzkus et al., 2017). In non-convex settings, the characterization of how SGD escapes from stationary points, including saddle points and local minima, was analyzed in (Daneshmand et al., 2018; Jin et al., 2017; Hu et al., 2017). On the other hand, in the context of deep learning, researchers realized that the noise introduced by SGD impacts the generalization, thanks to the research on the phenomenon that training with a large batch could cause a significant drop of test accuracy (Keskar et al., 2017). Particularly, several works attempted to investigate how the magnitude of the noise influences the generalization during the process of SGD optimization, including the batch size and learning rate (Hoffer et al., 2017; Goyal et al., 2017; Chaudhari & Soatto, 2017; Jastrzkebski et al., 2017). Another line of research interpreted SGD from a Bayesian perspective. In (Mandt et al., 2017; Chaudhari & Soatto, 2017), SGD was interpreted as performing variational inference, where certain entropic regularization involves to prevent overfitting. And the work (Smith & Le, 2018) attempted to provide an understanding based on model evidence. These explanations are compatible with the flat/sharp minima argument (Hochreiter & Schmidhuber, 1997; Keskar et al., 2017), since Bayesian inference tends to targeting the re-

gion with large probability mass, corresponding to the flat minima.

When analyzing SGD, most of existing works assume the noise covariance of SGD is constant or upper bounded by some constant, and what role the noise structure of stochastic gradient plays in optimization and generalization was rarely studied in literature. On the other hand, experiments (Figure 1) show that the isotropic approximation of SGD like gradient Langevin dynamic (GLD) cannot fully explain the mystery of the good generalization performance of SGD, even they are tuned to have the same noise magnitude. Thus the analysis over the structure of SGD noise is on demand for fully understanding SGD.

In this work, we take the first step studying the anisotropic noise of SGD and its superiority over its isotropic equivalence. Specifically, we study a general form of gradient-based optimization dynamics with unbiased noise, which unifies SGD and standard Langevin dynamics. By investigating the general dynamics, we analyze how the noise structure of SGD influences the escaping behavior from sharp minima and its regularization effects. Several novel analysis and empirical justifications are made as follow.

(1) We derive a key indicator to characterize the efficiency of escaping from minima through measuring the alignment of noise covariance and the curvature of loss function. Based on this indicator, two conditions are established to show which type of noise structure is superior to isotropic noise in term of escaping efficiency;

(2) We further justify that SGD in the context of neural networks satisfies these two conditions, and thus provide a plausible explanation why SGD can escape from sharp minima more efficiently, converging to flat minima with a higher probability. Moreover, these flat minima typically generalize well according to various works (Hochreiter & Schmidhuber, 1997; Keskar et al., 2017; Neyshabur et al., 2017; Wu & Zhu, 2017). We also show that Langevin dynamics with well tuned isotropic noise cannot beat SGD, which further confirms the importance of noise structure of SGD;

(3) A large number of experiments are designed systematically to justify our understanding on the behavior of the anisotropic diffusion of SGD. We compare SGD with full gradient descent with different types of diffusion noise, including isotropic and position-dependent/independent noise. All these comparisons demonstrate the effectiveness of anisotropic diffusion for good generalization in training neural networks.

## 2. Background

In general, supervised learning usually involves an optimization process of minimizing an empirical loss over training data,

$$L(\theta) := \frac{1}{N} \sum_{i=1}^{N} \ell(x_i; \theta), \tag{1}$$

where $\{x_i | i = 1, \dots, N\}$ denotes the training set with $N$ *i.i.d.* training samples, the model is parameterized by $\theta \in \mathbb{R}^D$ and $\ell$ denotes the combination of the loss and the model for simplicity, e.g. deep networks with cross entropy loss. Under many circumstances, including deep networks, there could exist multiple global minima for Eq. (1), exhibiting diverse generalization performance. We call those solutions generalizing well good solutions or minima, and vice versa.

**Gradient descent and its stochastic variants** A typical approach to minimize Eq. (1) is gradient descent (GD),

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta L(\theta_t), \tag{2}$$

where $\eta$ denotes the learning rate and we assume it to be a small constant for the convenience of analysis, similarly hereinafter.

In practice, a more useful kind of gradient based optimizers act like GD with an unbiased noise, including gradient Langevin dynamics (GLD),

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta L(\theta_t) + \eta \epsilon_t, \epsilon_t \sim \mathcal{N}\left(0, \sigma_t^2 I\right); \tag{3}$$

and stochastic gradient descent (SGD),

$$\theta_{t+1} = \theta_t - \eta \tilde{g}(\theta_t), \tag{4}$$

where $\tilde{g}(\theta_t) = \frac{1}{m} \sum_{x \in B_t} \nabla_\theta \ell(x; \theta_t)$ is an unbiased estimator of the full gradient $\nabla_\theta L(\theta_t)$, with $B_t$ being a randomly selected minibatch of size $m$. Assume the size of minibatch $m$ is large enough for the central limit theorem to hold, thus $\tilde{g}(\theta_t)$ follows a Gaussian distribution (Chen et al., 2014; Ahn et al., 2012; Shang et al., 2015; Mandt et al., 2017),

$$\tilde{g}(\theta_t) \sim \mathcal{N}\left(\nabla L(\theta_t), \Sigma^{\text{sgd}}(\theta_t)\right), \Sigma^{\text{sgd}}(\theta_t) \approx$$

$$\frac{1}{m} \left[ \frac{1}{N} \sum_{i=1}^{N} \nabla \ell(x_i; \theta_t) \nabla \ell(x_i; \theta_t)^T - \nabla L(\theta_t) \nabla L(\theta_t)^T \right]. \tag{5}$$

Therefore we can rewrite Eq. (4) as,

$$\theta_{t+1} = \theta_t - \eta \nabla L(\theta_t) + \eta \epsilon_t, \quad \epsilon_t \sim \mathcal{N}\left(0, \Sigma^{\text{sgd}}(\theta_t)\right). \tag{6}$$

Inspired by the dynamics of GLD (Eq. (3)) and SGD (Eq. (6)), more generally, we study the dynamics of *gradient descent with unbiased noise*,

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta L(\theta_t) + \eta \epsilon_t, \quad \epsilon_t \sim \mathcal{N}\left(0, \Sigma_t\right). \tag{7}$$

*Table 1.* Compared dynamics defined in Eq. (7). The parameter $\sigma_t$ is adjusted to force the noise share the same expected norm as that of SGD noise, to meet constraint Eq. (14) for fair comparison.

| Dynamics | Noise $\epsilon_t$ | Remarks |
|---|---|---|
| **SGD** | $\epsilon_t \sim \mathcal{N}\left(0, \Sigma_t^{\mathrm{sgd}}\right)$ | $\Sigma_t^{\mathrm{sgd}}$ is defined as in Eq. (5). |
| **GLD constant** | $\epsilon_t \sim \mathcal{N}\left(0, \varrho_t^2 I\right)$ | $\varrho_t$ is a tunable constant. |
| **GLD dynamic** | $\epsilon_t \sim \mathcal{N}\left(0, \sigma_t^2 I\right)$ | $\sigma_t$ is adjusted to force the noise share the same magnitude with SGD noise, similarly hereinafter. |
| **GLD diagonal** | $\epsilon_t \sim \mathcal{N}\left(0, \mathrm{diag}(\Sigma_t^{\mathrm{sgd}})\right)$ | $\mathrm{diag}(\Sigma_t^{\mathrm{sgd}})$ is the diagonal of the covariance of SGD noise $\Sigma_t^{\mathrm{sgd}}$. |
| **GLD leading** | $\epsilon_t \sim \mathcal{N}\left(0, \sigma_t \tilde{\Sigma}_t\right)$ | $\tilde{\Sigma}_t$ is a low rank approximation of $\Sigma_t^{\mathrm{sgd}}$, i.e., $\tilde{\Sigma}_t = \sum_{i=1}^k \gamma_i v_i v_i^T$, where $\gamma_i, v_i$ are the first $k$ leading eigenvalues and corresponding unit eigenvectors of $\Sigma_t^{\mathrm{sgd}}$. |
| **GLD Hessian** | $\epsilon_t \sim \mathcal{N}\left(0, \sigma_t \tilde{H}_t\right)$ | $\tilde{H}_t$ is a low rank approximation of the Hessian matrix of loss $L(\theta)$ by its the first $k$ leading eigenvalues and corresponding eigenvalues. |
| **GLD 1st eigven**($H$) | $\epsilon_t \sim \mathcal{N}\left(0, \sigma_t \lambda_1 u_1 u_1^T\right)$ | $\lambda_1, u_1$ are the maximal eigenvalue and its corresponding unit eigenvector of the Hessian matrix of loss $L(\theta_t)$. |

For small enough constant learning rate $\eta$, Eq. (7) can be treated as the numerical discretization of the following stochastic differential equation (SDE) (Li et al., 2017; Jastrzkebski et al., 2017; Chaudhari & Soatto, 2017),

$$\mathrm{d}\theta_t = -\nabla_\theta L(\theta_t)\,\mathrm{d}t + \sqrt{\eta \Sigma_t}\,\mathrm{d}W_t, \qquad (8)$$

where $W_t$ is a standard Brownian motion in $\mathbb{R}^D$.

Let $\Sigma_t = \Sigma^{\mathrm{sgd}}(\theta_t)$ and $\sqrt{\eta \Sigma^{\mathrm{sgd}}(\theta_t)}$ be the coefficient of the the noise term, Hoffer et al. (2017) and Jastrzkebski et al. (2017) studied the generalization influence of the magnitude of the SGD noise, which is controlled by the quotient of learning rate and batch size, $\frac{\eta}{m}$.

Different from previous works either assuming the noise of SGD is constant or upper bounded by some constant, we are the first to study SGD from the perspective of its noise structure. In the following sections, we first show that for dynamics Eq. (8), the structure of $\Sigma_t$ indeed affects the escaping from minima, especially for the sharp ones containing rich curvature information; and then we demonstrate that for neural networks, the noise of SGD is closely related to the Hessian of loss surface. Hence we conclude that SGD can escape from sharp minima much faster than its isotropic equivalence, and converge to flatter minima which tend to generalize better. Finally we verify our understanding by numerous experiments.

## 3. The behaviors of escaping from minima

To ease the notation, we absorb $\eta$ into $\Sigma_t$ in Eq. (8),

$$\mathrm{d}\theta_t = -\nabla_\theta L(\theta_t)\,\mathrm{d}t + \Sigma_t^{\frac{1}{2}}\,\mathrm{d}W_t. \qquad (9)$$

We now analyze the escaping behaviors of dynamics Eq. (9) with different choices of noise structures, i.e., $\Sigma_t$.

### 3.1. The escaping efficiency

We define the *escaping efficiency* as the expected increase of the potential or the loss.

**Definition 1** (Escaping efficiency). *Suppose we start the dynamics of Eq. (9) from the minimum $\theta_0$, then for a fixed time $t$ small enough (such that $L(\theta_t) - L(\theta_0) \geq 0$), we call*

$$\mathbb{E}_{\theta_t}[L(\theta_t) - L(\theta_0)] \qquad (10)$$

*the* escaping efficiency.

There are two remarks about the definition of escaping efficiency. Firstly it characterizes the ability of the dynamic escaping from the minimum $\theta_0$. Secondly because $L(\theta_t) - L(\theta_0) \geq 0$, for any $\delta > 0$, the escaping probability $P(L(\theta_t) - L(\theta_0) \geq \delta)$ can be upper bounded by the expectation $\mathbb{E}[L(\theta_t) - L(\theta_0)]$, given the Markov's inequality, $P(L(\theta_t) - L(\theta_0) \geq \delta) \leq \frac{\mathbb{E}[L(\theta_t) - L(\theta_0)]}{\delta}$.

Now we calculate the escape efficiency of dynamics Eq. (9). Provided that the mild smoothness assumptions for Ito's lemma holds, we have

$$\mathbb{E}[L(\theta_t) - L(\theta_0)] = -\int_0^t \mathbb{E}\left[\nabla L^T \nabla L\right] + \int_0^t \frac{1}{2}\mathbb{E}\mathrm{Tr}(H_t \Sigma_t)\,\mathrm{d}t, \qquad (11)$$

where $H_t := \nabla_\theta^2 L(\theta_t)$ is the Hessian of $L(\theta_t)$. The derivation of Eq. (11) is provided in Supplementary Materials.

Generally, the escaping efficiency characterized by Eq. (11) is hard to analyze due to the intractableness of the integral. Nonetheless, focusing on the locally escaping process, we take the second-order approximation near the minima $\theta_0$, where $L(\theta) \approx L(\theta_0) + \frac{1}{2}(\theta - \theta_0)^T H(\theta - \theta_0)$. Without loss of generality, let $\theta_0 = 0$. Further, assume $H$ is a positive definite matrix and the diffusion covariance matrix $\Sigma_t = \Sigma$ is constant for $t$. Then Eq. (9) becomes an Ornstein-

Uhlenbeck process,

$$d\theta_t = -H\theta_t \, dt + \Sigma^{\frac{1}{2}} \, dW_t, \quad \theta_0 = 0. \qquad (12)$$

The escaping efficiency of Eq. (12) could be explicitly obtained as

$$\mathbb{E}[L(\theta_t) - L(\theta_0)] = \frac{1}{4}\text{Tr}\left(\left(I - e^{-2Ht}\right)\Sigma\right) \approx \frac{t}{2}\text{Tr}(H\Sigma). \tag{13}$$

We defer the derivation to Supplementary Materials.

Eq. (11) and Eq. (13) characterize the escaping efficiency of general process and Ornstein-Uhlenbeck process respectively, and they clearly show that the indicator $\text{Tr}(H_t\Sigma_t)$ plays an crucial role for stochastic processes escaping from minima. Since we only care about the locally escaping behavior near the minima, we could directly analyze this key indicator $\text{Tr}(H_t\Sigma_t)$, in order to understand the importance of noise structure $\Sigma_t$ for escaping.

### 3.2. Anisotropic noise helps escape from sharp minima

Now we study what factors affect the locally escaping behaviors by analyzing the indicator $\text{Tr}(H_t\Sigma_t)$.

**The magnitude of noise** Clearly, the magnitude of the noise affects the escaping efficiency and larger magnitude leads to faster escape. Along this line, Hoffer et al. (2017) and Jastrzkebski et al. (2017) studied the generalization influence of the magnitude of the SGD noise, which is controlled by the quotient of learning rate and batch size.

Hence to explore the role of the noise structure, we must eliminate the impact of noise magnitude for fair comparison. One reasonable evaluation of the noise magnitude is the expected squared norm of the noise vector (Li et al., 2017): suppose $\epsilon_t \sim \mathcal{N}(0, \Sigma_t), z \sim \mathcal{N}(0, I)$ and the eigen decomposition of $\Sigma_t$ is $\Sigma_t = V\Gamma V^T$, then

$$\|\epsilon_t\|_{\text{trace}} := \mathbb{E}[\epsilon_t^T \epsilon_t] = \mathbb{E}[(V\sqrt{\Gamma}z)^T(V\sqrt{\Gamma}z)] = \mathbb{E}[z^T\Gamma z]$$
$$= \mathbb{E}\text{Tr}(\Gamma zz^T) = \text{Tr}\mathbb{E}[\Gamma zz^T] = \text{Tr}(\Sigma_t).$$

Based on such measure of magnitude, we introduce the following important trace constraint,

$$\textbf{given time } t, \text{Tr}(\Sigma_t) \textbf{ is constant}. \tag{14}$$

From the statistical physics point of view, $\text{Tr}(\Sigma_t)$ characterizes the kinetic energy (Gardiner, 2018), thus it is natural to force the energy to be unchanging, otherwise it is trivial that the higher the energy is, the less stable the system is.

**The ill-conditioning of minima** Consider the isotropic minima where the Hessian is $H_t = \lambda I$, our escaping indicator becomes $\text{Tr}(H_t\Sigma_t) = \lambda\text{Tr}\Sigma_t$, which is invariant under

constraint Eq. (14). Thus the noise structure has no impact on escaping from isotropic minima. However, for the minima where the Hessian is highly ill-conditioned, which is the typical case in practical over-parameterized neural networks (Sagun et al., 2017), the noise structure could cause huge difference on escaping behaviors, as analyzed below.

**The structure of noise** For semi-positive definite $H_t, \Sigma_t$ and assuming $H_t$ has distinguished top eigenvalues, to achieve the maximum of $\text{Tr}(H_t\Sigma_t)$ under constraint Eq.(14), $\Sigma_t$ should be $\Sigma_t^* = (\text{Tr}\Sigma_t) \cdot u_1 u_1^T$, where $u_1$ is the first unit eigenvector of $H_t$. Note this rank-1 matrix $\Sigma_t^*$ is highly anisotropic. More generally, the following Proposition 1 characterizes one kind of anisotropic noise significantly outperforming its isotropic equivalence, given $H$ is ill-conditioned.

**Proposition 1** (The benefits of anisotropic noise). *Assume $H_{D \times D}$ and $\Sigma_{D \times D}$ are semi-positive definite. If*

*(1) H is ill-conditioned. Let $\lambda_1 \geq \lambda_2 \geq \ldots, \geq \lambda_D \geq 0$ be the eigenvalues of H in descent order, and for some constant $k \ll D$ and $d > \frac{1}{2}$,*

$$\lambda_1 > 0, \qquad \lambda_{k+1}, \lambda_{k+2}, \ldots, \lambda_D < \lambda_1 D^{-d}; \tag{15}$$

*(2) $\Sigma$ is "aligned" with H. Let $u_i$ be the corresponding unit eigenvector of eigenvalue $\lambda_i$, for some projection coefficient $a > 0$,*

$$u_1^T \Sigma u_1 \geq a\lambda_1 \frac{Tr\Sigma}{TrH}; \tag{16}$$

*then for such $\Sigma$ and its isotropic equivalence $\bar{\Sigma} = \frac{Tr\Sigma}{D}I$ under constraint Eq. (14), we have the follow ratio describing their difference in term of escaping efficiency,*

$$\frac{Tr(H\Sigma)}{Tr(H\bar{\Sigma})} = \mathcal{O}\left(aD^{(2d-1)}\right), \quad d > \frac{1}{2}. \tag{17}$$

The first condition Eq. (15) characterizes the illness of $H$. To give some geometric intuitions on the second condition Eq. (16), let the maximal eigenvalue and its corresponding unit eigenvector of $\Sigma$ be $\gamma_1, v_1$, then $u_1^T\Sigma u_1 \geq u_1^T v_1 \gamma_1 v_1^T u_1 = \gamma_1 \langle u_1, v_1 \rangle^2$. Thus if the maximal eigenvalues of $H$ and $\Sigma$ are aligned in proportion, $\gamma_1/\text{Tr}\Sigma \geq a_1\lambda_1/\text{Tr}H$, and the angle of their corresponding unit eigenvectors is close enough such that $\langle u_1, v_1 \rangle \geq a_2$, the second condition Eq. (16) holds for $a = a_1 a_2$.

Typically, in the scenario of modern neural networks, due to the over-parameterization, Hessian and the gradient covariance are usually ill-conditioned and anisotropic near minima (Sagun et al., 2017; Chaudhari & Soatto, 2017). Thus the first condition in Proposition 1 usually holds for neural networks, and we further justify it by experiments in Section 5.3. In the next section, we turn to discuss how

the covariance of SGD noise meets the second condition of Proposition 1 in the context of neural networks. Hence this explains the superiority of the anisotropic noise of SGD over the isotropic one, such as gradient Langevin dynamics.

## 4. The relationship between the noise of SGD and the curvature of loss surface

In this section we investigate the anisotropic structure of gradient covariance in SGD, and explore its connection with the Hessian of loss surface.

**Around the true parameter.** According to the classic statistical theory (Pawitan, 2001, Chap. 8), for population loss $L(\theta) = \mathbb{E}_x \ell(x; \theta)$, with $\ell$ being the negative log likelihood, when evaluating at the true parameter $\theta^*$, there is the exact equivalence between the Hessian $H$ of the population loss and *Fisher information matrix F* at $\theta^*$,

$$F(\theta^*) := \mathbb{E}_x[\nabla_\theta \ell(x; \theta^*)\nabla_\theta \ell(x; \theta^*)^T] = \mathbb{E}_x[\nabla_\theta^2 \ell(x; \theta^*)]$$
$$= \nabla_\theta^2 L(\theta^*) =: H(\theta^*).$$

In practice, with the assumptions that the sample size $N$ is large enough (i.e. indicating asymptotic behavior) and suitable smoothness conditions, when the current parameter $\theta_t$ is not far from the ground truth, Fisher is close to Hessian. Thus we can obtain the following approximate equality between gradient covariance and Hessian,

$$\hat{\Sigma}(\theta_t) = \hat{F}(\theta_t) - \nabla_{\theta_t}\hat{L}(\theta_t)\nabla_\theta \hat{L}^T(\theta_t) \approx \hat{F}(\theta_t) \approx \hat{F}(\theta^*)$$
$$\approx F(\theta^*) = H(\theta^*) \approx \hat{H}(\theta^*) \approx \hat{H}(\theta_t).$$
(18)

The first approximation is due to the dominance of noise over the mean of gradient in the later stage of SGD optimization (Shwartz-Ziv & Tishby, 2017), in which a similar experiment was conducted to demonstrate this observation, shown in Supplementary Materials due to the limit of space.

**One hidden layer network with fixed output layer.** In the following we provide theoretical characterization about the alignment between $\Sigma$ and $H$ in the context of one hidden layer neural networks with fixed output layer. We first show the connection of Fisher and Hessian in this specific case.

**Proposition 2** (The connection between Fisher and Hessian in one hidden layer network). *Consider a binary classification problem with data $\{(x_i, y_i)\}_{i \in I}, y \in \{0, 1\}$, and mean square loss (either population or empirical),*

$$L(\theta) = \mathbb{E}_{(x,y)}\|\phi \circ f(x; \theta) - y\|^2.$$
(19)

*Here f denotes the network and $\phi$ is a threshold activation function controlling the output of the model,*

$$\phi(f) = \min\{\max\{f, \delta\}, 1 - \delta\} \subset [\delta, 1 - \delta],$$
(20)

*where $\delta$ is a small positive constant.*

*Suppose the network f satisfies: (1) it has one hidden layer and piece-wise linear activation; (2) the parameters of its output layer are fixed during training (Brutzkus et al., 2017).*

*Then for Fisher F and Hessian H (either population or empirical), we have*

*(1) $F(\theta) \succeq \delta^2 H(\theta)$, almost everywhere; (2) $F(\theta) \preceq (\delta + \epsilon)^2 H(\theta)$, almost everywhere around the minima, $\{\theta : \|\phi \circ f(x; \theta) - y\| \le \delta + \epsilon, \forall(x,y)\}$. $A \preceq B$ means that $(B - A)$ is semi-positive definite.*

There are two remarks on Proposition 2. Firstly, the considered neural networks in Proposition 2 are non-convex and have multiple minima, and one example to show this is provided in Supplementary Materials. Thus it is non-trivial to consider the escaping from minima. Secondly, Proposition 2 holds in both population and empirical sense, since the proof does not distinguish the two circumstances.

Based on Proposition 2, we could show that this neural network meets the second condition in Proposition 1.

**Proposition 3** (The connection between gradient covariance and Hessian in one hidden layer network). *Assume the conditions in Proposition 2 hold, then there is a constant $a > 0$, for $\theta$ close enough to minima $\theta^*$ (local or global), we have*

$$u(\theta)^T\Sigma(\theta)u(\theta) \ge a\lambda(\theta)\frac{Tr\Sigma(\theta)}{TrH(\theta)}$$
(21)

*holds almost everywhere, for $\lambda(\theta)$ and $u(\theta)$ being the maximal eigenvalue and its corresponding eigenvector of Hessian $H(\theta)$.*

Therefore, based on the discussion on population loss around the true parameters and one hidden layer neural network with fixed output layer parameters, given the ill-conditioning of $H$ due to the over-parameterization of modern neural networks, according to Proposition 1, we can conclude the noise structure of SGD helps to escape from sharp minima significantly faster than the dynamics with isotropic noise. Hence SGD tends to converge to flatter solutions, which typically generalize well (Hochreiter & Schmidhuber, 1997; Keskar et al., 2017; Neyshabur et al., 2017; Wu & Zhu, 2017). Thus, the anisotropic noise of SGD might explain its better generalization performance comparing to GD, GLD and other dynamics with isotropic noise.

In the following, we conduct a series of experiments systematically to verify our understanding on the behavior of escaping from minima and its regularization effects for different optimization dynamics.

## 5. Experiments

For better understanding the difference between the anisotropic noise and the isotropic one, we introduce dy-
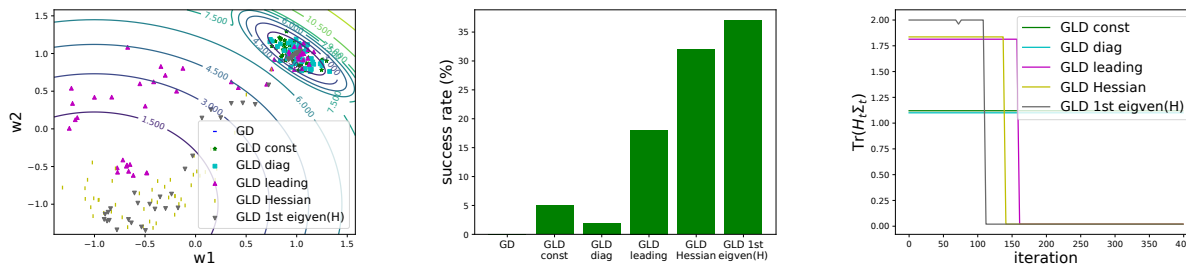
*Figure 2.* 2-D toy example. Compared dynamics are defined in Table 1, $k = 2$, $\sigma_t^2$ is tuned to keep noise of all dynamics sharing same expected squared norm, 0.01. All dynamics are run by 500 iterations with learning rate 0.005. **Left**: The trajectory of each compared dynamics for escaping from the sharp minimum in one run. **Middle**: Success rate of arriving the flat solution in 100 repeated runs. **Right**: $\text{Tr}(H_t \Sigma_t)$ of compared dynamics in one run.
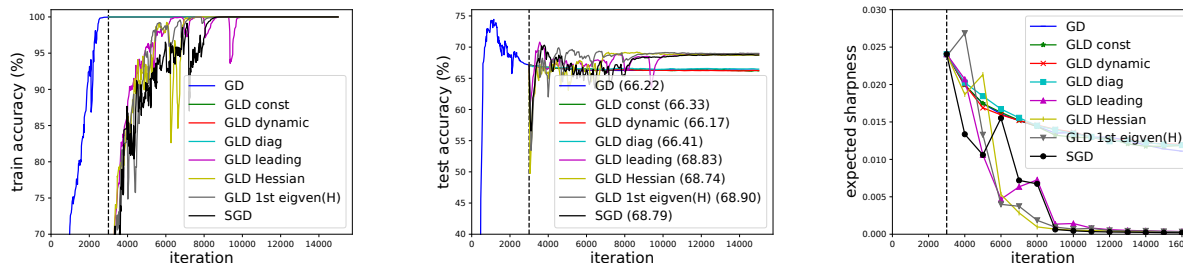


*Figure 3.* FashionMNIST experiments. Compared dynamics are initialized at $\theta_{GD}^*$ found by GD, marked by the vertical dashed line in iteration 3000. The learning rate is same for all the compared methods, $\eta_t = 0.07$, and batch size $m = 20$. **Left**: Training accuracy versus iteration. **Middle**: Test accuracy versus iteration. The final accuracy is noted within the parentheses. **Right**: Expected sharpness versus iteration. Expected sharpness is measured as $\mathbb{E}_{\nu \sim \mathcal{N}(0, \delta^2 I)} \left[ L(\theta + \nu) \right] - L(\theta)$, and $\delta = 0.01$, the expectation is computed by average on 1000 times sampling.

namics with various kinds of noise structure to empirical study with, as shown in Table 1.

### 5.1. Two-dimensional toy example

We design a 2-D toy example $L(w_1, w_2)$ with two basins, a small one and a large one, corresponding to a sharp and flat minima, $(1, 1)$ and $(-1, -1)$, respectively, both of which are global minima, see Supplementary Materials for more details. We initialize the dynamics of interest with the sharp minimum $(w_1, w_2) = (1, 1)$, and run them to study their behaviors escaping from this sharp minimum.

To explicitly control the noise magnitude, we only conduct experiments on GD, GLD const, GLD diag, GLD leading (with $k = 2 = D$ in Table 1, which is also the exactly covariance of SGD noise), GLD Hessian ($k = 2$) and GLD 1st eigen($H$). And we adjust $\sigma_t$ in each dynamics to force their noise to share the same expected squared norm the meet the constraint Eq. (14). Figure 2 (Left) shows the trajectories of the dynamics escaping from the sharp minimum $(1, 1)$ towards the flat one $(-1, -1)$, while Figure 2 (Middle) presents the success rate of escaping for each dynamic during 100 repeated experiments. Figure 2 (Right) demonstrates our derived indicator $\text{Tr}(H_t \Sigma_t)$ in one run.

As shown in Figure 2, GLD 1st eigvec($H$) achieves the highest success rate, indicating the fastest escaping speed from

the sharp minimum. The dynamics with anisotropic noise aligned with Hessian well, including GLD 1st eigvec($H$), GLD Hessian and GLD leading, greatly outperform GD, GLD const with isotropic noise, and GLD diag with noise poorly aligned with Hessian. These experiments are consistent with our theoretical analysis on OU process shown in Eq. (13) and Proposition 1, demonstrating the benefits of anisotropic noise for escaping from sharp minima.

### 5.2. One hidden layer network with fixed output layer

To verify the conclusion of Proposition 1 in neural network cases, three networks are trained to binary classify $1,000$ linearly separable two-dimensional points to show the benefits of anisotropic noise of SGD. The activations are all ReLU and $\delta$ (in Proposition 2) is set to be 0.001. The number of hidden nodes for each network varies in $\{32, 128, 512\}$. We plot the empirical indicator $\text{Tr}(H\Sigma)$ in Figure 4. We can easily observe that as the increase of the number of hidden nodes, the ratio $\text{Tr}(H\Sigma)/\text{Tr}(H\bar{\Sigma})$ is enlarged significantly, which is consistent with the Eq. (17) described in Proposition 1.

### 5.3. FashionMNIST with corrupted labels

We conduct a series of experiments in real deep learning scenarios to study the importance of SGD's noise covariance
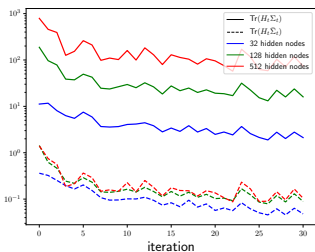
*Figure 4.* One hidden layer neural networks. The solid and the dotted lines represent the value of $\text{Tr}(H\Sigma)$ and $\text{Tr}(H\bar{\Sigma})$, respectively. The number of hidden nodes varies in $\{32, 128, 512\}$.



*Figure 5.* FashionMNIST experiments. **Left**: The first $400$ eigenvalues of Hessian at $\theta^*_{GD}$, the sharp minima found by GD after $3000$ iterations. **Middle**: The projection coefficient estimation $\hat{a} = \frac{u_1^T \Sigma u_1 \text{Tr}H}{\lambda_1 \text{Tr}\Sigma}$, as shown in Proposition 1. **Right**: $\text{Tr}(H_t\Sigma_t)$ versus $\text{Tr}(H_t\bar{\Sigma}_t)$ during SGD optimization initialized from $\theta^*_{GD}$, $\bar{\Sigma}_t = \frac{\text{Tr}\Sigma_t}{D}I$ denotes the isotropic noise with same expected squared norm as SGD noise.

structure and its implicit regularization effects. We construct a noisy training set based on FashionMNIST dataset. Concretely, the training set consist of $1000$ images with correct labels, and another $200$ images with random labels. A small LeNet-like network with $11,330$ parameters is utilized such that the spectral decomposition over $\Sigma$ and $H$ are computationally feasible.

We firstly run the full gradient decent for $3,000$ iterations to arrive at the parameters $\theta^*_{GD}$ near the global minima with nearly zero training loss and $100\%$ training accuracy, which are typically sharp minima that generalize poorly (Neyshabur et al., 2017). And then all other compared methods are initialized with $\theta^*_{GD}$ and run with the same learning rate $\eta_t = 0.07$ and same batch size $m = 20$ (if needed) for fair comparison.

**Behaviors of different dynamics escaping from minima and its generalization effects.** To compare the different dynamics on escaping behaviors and generalization performance, we run dynamics initialized from the sharp minima $\theta^*_{GD}$ found by GD. The settings for each compared method are as follows. The hyperparameter $\sigma^2$ for GLD const has already been tuned as optimal ($\sigma = 0.001$) by grid search. For GLD leading, we set $k = 20$ for comprising the computational cost and approximation accuracy. As for GLD Hessian, to reduce the expensive evaluation of such a huge Hessian in each iteration, we set $k = 20$ and update the Hessian every 10 iterations. We adjust $\sigma_t$ in GLD dynamic, GLD Hessian and GLD 1st eigvec($H$) to guarantee that they share the same expected squred noise norm defined in Eq. (14) as that of SGD. And we measure the expected sharpness of different minima as $\mathbb{E}_{\nu \sim \mathcal{N}(0, \delta^2 I)} \left[ L(\theta + \nu) \right] - L(\theta)$, as defined in ((Neyshabur et al., 2017), Eq.(7)).

As shown in Figure 3, SGD, GLD 1st eigvec($H$), GLD leading and GLD Hessian successfully escape from the sharp minima found by GD, while GLD, GLD dynamic and GLD diag are trapped in the minima. This demonstrates that the methods with anisotropic noise "aligned" with loss curvature can help to find flatter minima that generalize well.

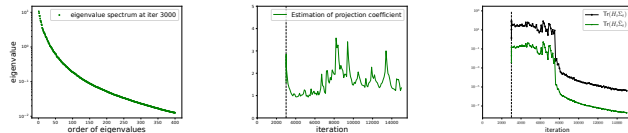**Verification of the conditions in Proposition 1.**

To check whether the noise of SGD in deep neural networks satisfies the two conditions in Proposition 1, we run SGD initialized from $\theta^*_{GD}$, i.e. the sharp minima found by GD.

Figure 5(Left) shows the first $400$ eigenvalues of Hessian at $\theta^*_{GD}$, from which we see that the $140$th eigenvalue has already decayed to about $1\%$ of the first eigenvalue. Note that Hessian $H \in \mathbb{R}^{D \times D}$, $D = 11330$, thus $H$ around $\theta^*_{GD}$ approximately meets the ill-conditioning requirement in Proposition 1. Figure 5(Middle) shows the projection coefficient estimated by $\hat{a} = \frac{u_1^T \Sigma u_1 \text{Tr}H}{\lambda_1 \text{Tr}\Sigma}$ along the trajectory of SGD. The plot indicates that the projection coefficient is in a descent scale comparing to $D^{2d-1}$, thus satisfying the second condition in Proposition 1. Therefore, Proposition 1 ensures that SGD would escape from minima $\theta^*_{GD}$ faster than GLD in order of $\mathcal{O}(D^{2d-1})$, as shown in Figure 5(Right). An interesting observation is that in the later stage of SGD optimization, $\text{Tr}(H\Sigma)$ becomes significantly ($10^7$ times) smaller than in the beginning stage, implying that SGD has already converged to minima being almost impossible to escape from. This phenomenon demonstrates the reasonability to employ $\text{Tr}(H\Sigma)$ as an empirical indicator for escaping efficiency.

### 5.4. SVHN and CIFAR-10

We also provide experiments on SVHN and CIFAR-10 datasets with VGG11 in Figure 1 and Figure 6. For CIFAR-10 we use the original datasets while we only use $2,5000$ training examples for SVHN to compromise with the computational burden of gradient descent. We choose VGG over ResNet since it achieves decent performance without using batch normalization, which causes extra affects on analyzing the noise of SGD. We re-estimate the noise structure of GLD dynamic and GLD diag every 10 iterations to ease the computational burden. Also, we only run GD, GLD const, GLD dynamic, GLD diag and SGD since the computational costs of these dynamics are relatively acceptable to our hardware.

From Figure 1 we can see the generalization gap between SGD and other dynamics, which demonstrates that the mag-
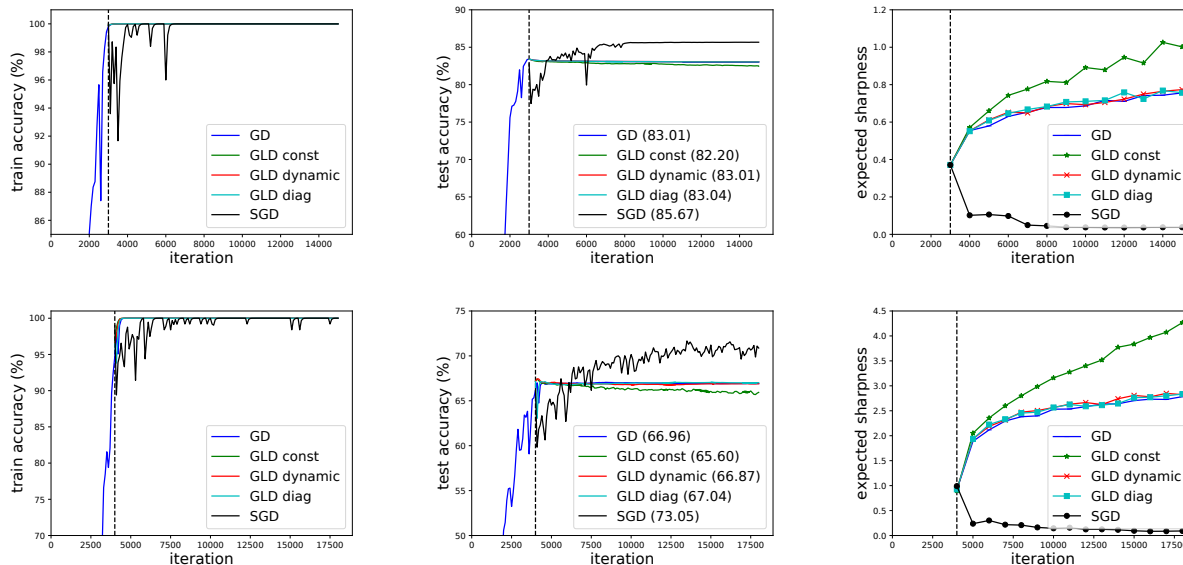
*Figure 6.* SVHN and CIFAR-10 experiments. **Top**: SVHN experiments; **Bottom**: CIFAR-10 experiments. Compared dynamics are initialized at $\theta^*_{GD}$ found by GD, marked by the vertical dashed line (in iteration $3,000$ for SVHN and iteration $4000$ for CIFAR-10). The learning rate is same for all the compared methods, $\eta_t = 0.05$, and batch size $m = 100$. **Left**: Training accuracy versus iteration. **Middle**: Test accuracy versus iteration. The final accuracy is noted within the parentheses. **Right**: Expected sharpness versus iteration. Expected sharpness is measured as $\mathbb{E}_{\nu \sim \mathcal{N}(0, \delta^2 I)}\left[L(\theta + \nu)\right] - L(\theta)$, and $\delta = 0.01$, the expectation is computed by average on 100 times sampling.

nitude of SGD cannot fully explain the performance of SGD. Figure 6 shows the escaping behavior of SGD and other dynamics, the results are consistent with experiments on FashionMNIST.

# 6. Discussions

**Benefits of considering covariance structure.** Previous works on SGD for deep learning typically ignores the co-variance structure, as we have shown in this work, which has significant effects on its dynamics behaviors and generalization performance as well. The key observation on connecting gradient noise structure with curvature of the loss landscape, especially near the minima, provides a new perspective for understanding why SGD can achieve good generalization in practice. Our work is an initial attempt to reveal the non-negligible benefits of SGD's covariance structure. More theoretical explorations are needed along this direction.

**Effects of learning rate and batch size.** As seen from the SGD dynamics in Eq. (6), when the learning rate is too small or batch size is overly large, the magnitude of gradient noise will become small, and thus effects of covariance structure is not obvious as before. In these cases, SGD often needs long time for diffusion towards flat minima to obtain better solutions, as shown in existing research (Keskar et al., 2017; Hoffer et al., 2017; Jastrzkebski et al., 2017).

**Designing optimizers that help to generalize better.** The derived indicator also sheds some light on designing the

optimizers that might generalize better than SGD by adding the noise along the direction of the maximum eigenvector of Hessian. We leave the exploration regarding this as future work.

# 7. Conclusion

We theoretically investigate a general optimization dynamics with unbiased noise, which unifies various existing optimization methods, including SGD. We provide some novel results on the behaviors of escaping from minima and its regularization effects. A novel indicator is derived for characterizing the escaping efficiency. Based on this indicator, two conditions are constructed for showing what type of noise structure is superior to isotropic noise in term of escaping. We then analyze the noise structure of SGD in neural networks and find that it indeed satisfies the two conditions, thus explaining the widely known observation that SGD can escape from sharp minima efficiently toward flat ones that generalize well. Various experimental evidence supports our arguments on the behavior of SGD and its effects on generalization. Our study also shows that isotropic noise helps little for escaping from sharp minima, due to the highly anisotropic nature of landscape. This indicates that it is not sufficient to analyze SGD by treating it as an isotropic diffusion over landscape (Zhang et al., 2017; Mou et al., 2017). A better understanding of this out-of-equilibrium behavior (Chaudhari & Soatto, 2017) is on demand.

## Acknowledgement

## References

Ahn, S., Korattikara, A., and Welling, M. Bayesian posterior sampling via stochastic gradient fisher scoring. In *Proceedings of the 29th International Coference on International Conference on Machine Learning*, pp. 1771–1778. Omnipress, 2012.

Bottou, L. Stochastic gradient learning in neural networks. *Proceedings of Neuro-Nımes*, 91(8), 1991.

Brutzkus, A., Globerson, A., Malach, E., and Shalev-Shwartz, S. Sgd learns over-parameterized networks that provably generalize on linearly separable data. *arXiv preprint arXiv:1710.10174*, 2017.

Chaudhari, P. and Soatto, S. Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. *arXiv preprint arXiv:1710.11029*, 2017.

Chen, T., Fox, E., and Guestrin, C. Stochastic gradient Hamiltonian Monte Carlo. In *International Conference on Machine Learning*, pp. 1683–1691, 2014.

Daneshmand, H., Kohler, J., Lucchi, A., and Hofmann, T. Escaping saddles with stochastic gradients. *arXiv preprint arXiv:1803.05999*, 2018.

Gardiner, C. *Stochastic Methods: A Handbook for the Natural and Social Sciences*. Springer Series in Synergetics. Springer Berlin Heidelberg, 2018. ISBN 9783540707127.

Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch sgd: training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.

Hochreiter, S. and Schmidhuber, J. Flat minima. *Neural Computation*, 9(1):1–42, 1997.

Hoffer, E., Hubara, I., and Soudry, D. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In *Advances in Neural Information Processing Systems 30*, pp. 1731–1741. 2017.

Hu, W., Junchi Li, C., Li, L., and Liu, J.-G. On the diffusion approximation of nonconvex stochastic gradient descent. *arXiv preprint arXiv:1705.07562*, 2017.

Jastrzkebski, S., Kenton, Z., Arpit, D., Ballas, N., Fischer, A., Bengio, Y., and Storkey, A. Three factors influencing minima in sgd. *arXiv preprint arXiv:1711.04623*, 2017.

Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. How to escape saddle points efficiently. *arXiv preprint arXiv:1703.00887*, 2017.

Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M., and Tang, P. T. P. On large-batch training for deep learning: Generalization gap and sharp minima. In *In International Conference on Learning Representations (ICLR)*, 2017.

Li, Q., Tai, C., and Weinan, E. Stochastic modified equations and adaptive stochastic gradient algorithms. In *International Conference on Machine Learning*, pp. 2101–2110, 2017.

Li, Y. and Yuan, Y. Convergence analysis of two-layer neural networks with relu activation. In *Advances in Neural Information Processing Systems*, pp. 597–607, 2017.

Mandt, S., Hoffman, M. D., and Blei, D. M. Stochastic gradient descent as approximate bayesian inference. *The Journal of Machine Learning Research*, 18(1):4873–4907, 2017.

Mou, W., Wang, L., Zhai, X., and Zheng, K. Generalization bounds of sgld for non-convex learning: Two theoretical viewpoints. *arXiv preprint arXiv:1707.05947*, 2017.

Neyshabur, B., Bhojanapalli, S., Mcallester, D., and Srebro, N. Exploring generalization in deep learning. In *Advances in Neural Information Processing Systems 30*, pp. 5947–5956. 2017.

Øksendal, B. Stochastic differential equations. In *Stochastic differential equations*, pp. 65–84. Springer, 2003.

Pawitan, Y. *In all likelihood: statistical modelling and inference using likelihood*. Oxford University Press, 2001.

Sagun, L., Evci, U., Guney, V. U., Dauphin, Y., and Bottou, L. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.

Shang, X., Zhu, Z., Leimkuhler, B., and Storkey, A. J. Covariance-controlled adaptive langevin thermostat for large-scale bayesian sampling. In *Advances in Neural Information Processing Systems*, pp. 37–45, 2015.

Shwartz-Ziv, R. and Tishby, N. Opening the black box of deep neural networks via information bottleneck. *arXiv preprint arXiv:1703.00810*, 2017.

Smith, S. L. and Le, Q. V. A bayesian perspective on generalization and stochastic gradient descent. *International Conference on Learning Representations*, 2018.

Wu, L. and Zhu, Z. Towards understanding generalization of deep learning: Perspective of loss landscapes. In *International Conference of Machine Learning Workshop*, 2017.

Zhang, Y., Liang, P., and Charikar, M. A hitting time analysis of stochastic gradient langevin dynamics. In *Conference on Learning Theory*, pp. 1980–2022, 2017.

## A. Derivations and Proofs for Main Paper

### A.1. Derivation of Eq. (11) in main paper

*Proof.* The "mild smoothness assumptions" refers that $L(\theta_t) \in C^2$. Then the Ito's lemma holds (Øksendal, 2003). Thus,

$$dL(\theta_t) \tag{22}$$

$$= \left( -\nabla L^T \nabla L + \frac{1}{2}\mathrm{Tr}\left( \Sigma_t^{\frac{1}{2}} H_t \Sigma_t^{\frac{1}{2}} \right) \right) dt + \nabla L^T \Sigma_t^{\frac{1}{2}} dW_t \tag{23}$$

$$= \left( -\nabla L^T \nabla L + \frac{1}{2}\mathrm{Tr}\left( H_t \Sigma_t \right) \right) dt + \nabla L^T \Sigma_t^{\frac{1}{2}} dW_t. \tag{24}$$

Taking expectation with respect to the distribution of $\theta_t$, we have

$$d\mathbb{E}_{\theta_t} L(\theta_t) = \mathbb{E}\left( -\nabla L^T \nabla L + \frac{1}{2}\mathrm{Tr}(H_t \Sigma_t) \right) dt, \tag{25}$$

since the expectation of Brownian motion is zero.

Thus the solution of $\mathbb{E}_{\theta_t} L(\theta_t)$ is,

$$\mathbb{E}L(\theta_t) = L(\theta_0) - \int_0^t \mathbb{E}\left( \nabla L^T \nabla L \right) + \int_0^t \frac{1}{2}\mathbb{E}\mathrm{Tr}(H_t \Sigma_t)\, dt. \tag{26}$$

$\square$

### A.2. Derivation of Eq. (13) in main paper

*Proof.* Without loss of generality, we assume that $L(\theta_0) = 0$.

For multivariate Ornstein-Uhlenbeck process, when $\theta_0 = 0$ is an constant, $\theta_t$ follows a multivariate Gaussian distribution (Øksendal, 2003).

For symmetric matrix $A$, let

$$e^A := U^T \mathrm{diag}(e^{\lambda_1}, \dots, e^{\lambda_n})U, \tag{27}$$

where $\lambda_1, \dots, \lambda_n$ and $U$ are the eigenvalues and eigenvector matrix of $A$.

Consider change of variables $\theta \to \phi(\theta, t) = e^{Ht}\theta_t$. Note that,

$$\frac{de^{Ht}}{dt} = He^{Ht}. \tag{28}$$

Thus by applying Ito's lemma, we have

$$d\phi(\theta_t, t) = e^{Ht}\Sigma^{\frac{1}{2}}\, dW_t, \tag{29}$$

which we can integrate form 0 to $t$ to obtain

$$\theta_t = 0 + \int_0^t e^{H(s-t)}\Sigma^{\frac{1}{2}}\, dW_s. \tag{30}$$

The expectation of $\theta_t$ is zero. And by Ito's isometry (Øksendal, 2003), the covariance of $\theta_t$ is,

$$\mathbb{E}\theta_t \theta_t^T \tag{31}$$

$$= \mathbb{E}\left[ \int_0^t e^{H(s-t)}\Sigma^{\frac{1}{2}}\, dW_s \left( \int_0^t e^{H(r-t)}\Sigma^{\frac{1}{2}}\, dW_r \right)^T \right] \tag{32}$$

$$= \mathbb{E}\left[ \int_0^t e^{H(s-t)}\Sigma^{\frac{1}{2}}\Sigma^{\frac{1}{2}}e^{H(s-t)}\, ds \right] \tag{33}$$

$$= \mathbb{E}\left[ \int_0^t e^{H(s-t)}\Sigma e^{H(s-t)}\, ds \right] \tag{34}$$

$$= \int_0^t e^{H(s-t)}\Sigma e^{H(s-t)}\, ds. \tag{35}$$

The last equation is because $H$ and $\Sigma$ are both constant.

Therefore

$$\mathbb{E}L(\theta_t) = \frac{1}{2}\mathbb{E}\mathrm{Tr}\left( \theta_t^T H \theta_t \right) \tag{36}$$

$$= \frac{1}{2}\mathrm{Tr}\left( H\mathbb{E}\theta_t \theta_t^T \right) \tag{37}$$

$$= \frac{1}{2}\int_0^t \mathrm{Tr}\left( He^{H(s-t)}\Sigma e^{H(s-t)} \right) ds \tag{38}$$

$$= \frac{1}{2}\int_0^t \mathrm{Tr}\left( e^{H(s-t)}H\Sigma e^{H(s-t)} \right) ds \tag{39}$$

$$= \frac{1}{2}\int_0^t \mathrm{Tr}\left( e^{2H(s-t)}H\Sigma \right) ds \tag{40}$$

$$= \frac{1}{2}\mathrm{Tr}\left( \frac{1}{2}H^{-1}\left( I - e^{-2Ht} \right)H\Sigma \right) \tag{41}$$

$$= \frac{1}{4}\mathrm{Tr}\left( \left( I - e^{-2Ht} \right)\Sigma \right). \tag{42}$$

Eq. (39) holds since $H$ is symmetric. Further, by Taylor's expansion we have

$$\mathbb{E}L(\theta_t) = \frac{1}{4}\mathrm{Tr}\left( \left( I - e^{-2Ht} \right)\Sigma \right) = \frac{t}{2}\mathrm{Tr}(H\Sigma). \tag{43}$$

$\square$

## A.3. Proof of Proposition 1

*Proof.* $\text{Tr}(H\Sigma)$ can be decomposed as

$$\text{Tr}(H\Sigma) = \sum_{i=1}^{D} \lambda_i u_i^T \Sigma u_i. \tag{44}$$

Thus by the conditions of Proposition 1, we can bound $\text{Tr}(H\Sigma)$ as

$$\text{Tr}(H\Sigma) \geq u_1^T \Sigma u_1 \geq a\lambda_1 \frac{\text{Tr}\Sigma}{\text{Tr}H}. \tag{45}$$

On the other hand,

$$\text{Tr}(H\bar{\Sigma}) = \frac{\text{Tr}\Sigma}{D}\text{Tr}H. \tag{46}$$

Thus,

$$\frac{\text{Tr}(H\Sigma)}{\text{Tr}(H\bar{\Sigma})} \geq \frac{a\lambda_1 D}{(\text{Tr}H)^2} \geq \frac{a\lambda_1 D}{\left(k\lambda_1 + (D-k)D^{-d}\lambda_1\right)^2} \tag{47}$$
$$= \mathcal{O}\left(aD^{2d-1}\right).$$

$\square$

## A.4. Proof of Proposition 2 in main paper

*Proof.* For simplicity, we define

$$\bar{f}(x;\theta) := \phi \circ f(x;\theta) \in [\delta, 1-\delta], \tag{48}$$

and

$$\ell(x,y;\theta) = \frac{1}{2}(\bar{f}(x;\theta) - y)^2. \tag{49}$$

Then the loss function becomes $L(\theta) = \mathbb{E}_{(x,y)}\ell(x,y;\theta)$.

Since both $f$ and $\phi$ are piecewise linear, $\bar{f}(x;\theta)$ is also piece-wise linear with respect to $\theta$. Thus the Hessian of $\bar{f}$ is zero almost everywhere.

We calculate the gradient and the Hessian of the loss:

$$\nabla_\theta L(\theta) = \mathbb{E}(\bar{f}(x;\theta) - y)\nabla_\theta \bar{f}(x;\theta); \tag{50}$$

$$H(\theta) = \nabla_\theta^2 L(\theta) \tag{51}$$
$$= \mathbb{E}\nabla_\theta \bar{f}(x;\theta) \cdot \nabla_\theta \bar{f}(x;\theta)^T + \mathbb{E}(\bar{f}(x;\theta) - y)\nabla_\theta^2 \bar{f}(x;\theta) \tag{52}$$
$$= \mathbb{E}\nabla_\theta \bar{f}(x;\theta) \cdot \nabla_\theta \bar{f}(x;\theta)^T. \quad \text{almost everywhere.} \tag{53}$$

The last equation holds almost everywhere, since $\bar{f}(x;\theta)$ is piece-wise linear and its Hessian is zero almost everywhere.

On the other hand, the Fisher is

$$F(\theta) = \mathbb{E}\nabla_\theta \ell(x,y;\theta) \cdot \nabla_\theta \ell(x,y;\theta)^T \tag{54}$$
$$= \mathbb{E}(\bar{f}(x;\theta) - y)^2 \nabla_\theta \bar{f}(x;\theta) \cdot \nabla_\theta \bar{f}(x;\theta)^2. \tag{55}$$

(1) Note that $\bar{f} \in [\delta, 1-\delta]$ and $y \in \{0, 1\}$, thus

$$(\bar{f}(x;\theta) - y)^2 \geq \delta^2. \tag{56}$$

Therefore

$$F(\theta) \succeq \mathbb{E}\delta^2 \nabla_\theta \bar{f}(x;\theta) \cdot \nabla_\theta \bar{f}(x;\theta)^2 = \delta^2 H(\theta), \tag{57}$$

holds almost everywhere.

(2) Around the minima where $\theta \in \{\theta : \|f(x;\theta) - y\| \leq \delta + \epsilon, \forall(x,y)\}$, we have

$$(\bar{f}(x;\theta) - y)^2 \leq (\delta + \epsilon)^2. \tag{58}$$

Therefore

$$F(\theta) \preceq \mathbb{E}(\delta+\epsilon)^2 \nabla_\theta \bar{f}(x;\theta) \cdot \nabla_\theta \bar{f}(x;\theta)^2 = (\delta+\epsilon)^2 H(\theta), \tag{59}$$

holds almost everywhere around the minima. $\square$

## A.5. Proof of Proposition 3 in main paper

*Proof.* We only consider $\theta$ around the minima $\theta^*$ such that $\{\theta : \|\phi \circ f(x;\theta) - y\| \leq \delta + \epsilon, \forall(x,y)\}$. On the other hand by construction $\|\phi \circ f(x;\theta) - y\| \geq \delta$. Thus according to Proposition 2,

$$\delta^2 H(\theta) \preceq F(\theta) \preceq (\delta + \epsilon)^2 H(\theta) \tag{60}$$

holds almost everywhere.

Thus let $\lambda(\theta)$ and $u(\theta)$ being the maximal eigenvalue and its corresponding eigenvector of $H(\theta)$,

$$u(\theta)^T F(\theta)u(\theta) \geq \delta^2 u(\theta)^T H(\theta)u(\theta) = \delta^2\lambda(\theta). \tag{61}$$

Since at the minimal $\theta^*$ the Hessian is not zero, thus there is a positive value $\lambda^* > 0$ such that $\lambda(\theta^*) > \lambda^* > 0$. Therefore by the continuity of $H(\theta)$, there are $\epsilon_1, \delta_1$, such that,

$$\lambda(\theta) > \lambda^* - \epsilon_1 > 0, \quad \forall\|\theta - \theta^*\| \leq \delta_1. \tag{62}$$

By Taylor's expansion,

$$\nabla L(\theta) = \nabla L(\theta^*) + H(\theta^*)(\theta - \theta^*) + o(\theta - \theta^*) \tag{63}$$
$$= H(\theta^*)(\theta - \theta^*) + o(\theta - \theta^*).$$

Hence,

$$\|\nabla L(\theta)\|_2^2 \leq \|H(\theta^*)\|_2^2\|\theta - \theta^*\|_2^2 + o\left(\|\theta - \theta^*\|_2^2\right). \tag{64}$$

Therefore, for all $\theta$ such that

$$\|\theta - \theta^*\|_2 \leq \frac{\sqrt{\delta^2\delta_2(\lambda^* - \epsilon_1)}}{\|H(\theta^*)\|_2} \tag{65}$$

$$\leq \frac{\sqrt{\delta^2\delta_2\lambda(\theta)}}{\|H(\theta^*)\|_2} \tag{66}$$

$$\leq \frac{\sqrt{\delta_2 u(\theta)^T F(\theta)u(\theta)}}{\|H(\theta^*)\|_2}, \tag{67}$$

we have

$$\left\|\nabla L(\theta)\right\|_2^2 \le \delta_2 u(\theta)^T F(\theta)u(\theta) + o\left(\left|\delta_2 u(\theta)^T F(\theta)u(\theta)\right|\right). \tag{68}$$

On the other hand, by definition, the gradient covariance $\Sigma$ and Fisher $F$ has the following relationship,

$$
\begin{aligned}
\Sigma(\theta) &= \mathbb{E}(\nabla\ell(x,y;\theta) - \nabla L(\theta)) \cdot (\nabla\ell(x,y;\theta) - \nabla L(\theta))^T \\
&= \mathbb{E}\nabla\ell(x,y;\theta) \cdot \nabla\ell(x,y;\theta)^T - \nabla L(\theta)\nabla L(\theta)^T \\
&= F(\theta) - \nabla L(\theta)\nabla L(\theta)^T.
\end{aligned}
\tag{69}
$$

Thus,

$$\frac{u(\theta)^T \Sigma(\theta) u(\theta)}{\text{Tr}\Sigma(\theta)} \tag{70}$$

$$= \frac{u(\theta)^T F(\theta)u(\theta) - u(\theta)^T \nabla L(\theta)\nabla L(\theta)^T u(\theta)}{\text{Tr}F(\theta) - \text{Tr}(\nabla L(\theta)\nabla L(\theta)^T)} \tag{71}$$

$$= \frac{u(\theta)^T F(\theta)u(\theta) - \left\|\nabla L(\theta)\right\|_2^2}{\text{Tr}F(\theta) - \left\|\nabla L(\theta)\right\|_2^2} \tag{72}$$

$$\ge \frac{u(\theta)^T F(\theta)u(\theta) - \left\|\nabla L(\theta)\right\|_2^2}{\text{Tr}F(\theta)} \tag{73}$$

$$= \frac{u(\theta)^T F(\theta)u(\theta)}{\text{Tr}F(\theta)}\left(1 - \frac{\left\|\nabla L(\theta)\right\|_2^2}{u(\theta)^T F(\theta)u(\theta)}\right) \tag{74}$$

$$\ge \frac{u(\theta)^T F(\theta)u(\theta)}{\text{Tr}F(\theta)}\left(1 - \delta_2 - o\left(|\delta_2|\right)\right) \tag{75}$$

$$\ge \frac{u(\theta)^T F(\theta)u(\theta)}{\text{Tr}F(\theta)}\left(1 - 2\delta_2\right). \tag{76}$$

Note that Eq. (60) indicates that

$$\forall u, \quad u^T(F(\theta) - \delta^2 H(\theta))u \ge 0 \tag{77}$$
$$\text{and} \quad \text{Tr}((\delta + \epsilon)^2 H(\theta) - F(\theta)) \ge 0. \tag{78}$$

Thus

$$\frac{u(\theta)^T F(\theta)u(\theta)}{\text{Tr}F(\theta)} \ge \frac{\delta^2 u(\theta)^T H(\theta)u(\theta)}{(\delta + \epsilon)^2 \text{Tr}H(\theta)} \tag{79}$$

$$= \frac{\delta^2 \lambda(\theta)}{(\delta + \epsilon)^2 \text{Tr}H(\theta)}. \tag{80}$$

Therefore for all $\theta$ in the set of

$$
\begin{aligned}
&\left\{\left\|\phi \circ f(x;\theta) - y\right\| \le \delta + \epsilon, \forall(x,y)\right\} \\
&\cap \left\{\left\|\theta - \theta^*\right\| \le \delta^1\right\} \\
&\cap \left\{\left\|\theta - \theta^*\right\|_2 \le \frac{\sqrt{\delta^2 \delta_2(\lambda^* - \epsilon_1)}}{\left\|H(\theta^*)\right\|_2}\right\},
\end{aligned}
\tag{81}
$$

we have

$$\frac{u(\theta)^T \Sigma(\theta) u(\theta)}{\text{Tr}\Sigma(\theta)} \ge \frac{u(\theta)^T F(\theta)u(\theta)}{\text{Tr}F(\theta)}(1 - 2\delta_2) \tag{82}$$

$$\ge \frac{(1 - 2\delta_2)\delta^2}{(\delta + \epsilon)^2}\frac{\lambda(\theta)}{\text{Tr}H(\theta)}. \tag{83}$$

$\square$

## B. About the non-convexity of the model in Proposition 2 in main paper

Suppose we only have one training data $\{x = (1,1); y = 1\}$, and the threshold activation is

$$\phi(f) = \min\{\max\{f, 0.1\}, 0.9\}. \tag{84}$$

Thus the loss is

$$L(w_1, w_2) = (\phi(relu(w_1) - relu(w_2)) - 1)^2. \tag{85}$$

Hence

$$
\begin{aligned}
L(1, 0) &= 0.01 \\
L(0, 1) &= 0.81 \\
L(0.5, 0.5) &= 0.81.
\end{aligned}
\tag{86}
$$

Therefore

$$\frac{1}{2}L(1, 0) + \frac{1}{2}L(0, 1) < L(0.5, 0.5), \tag{87}$$

which means that $L$ is not convex.

It is also easy to see that $L$ has multiple minima.

## C. Additional experiments

### C.1. Dominance of noise over gradient

Figure 7 shows the comparison of gradient mean and the expected norm of noise during training using SGD. The dataset and model are same as the experiments of FashionMNIST in main paper, or as in Section D.3. From Figure 7, we see that in the later stage of SGD optimization, the magnitude of noise indeed dominates that of gradient.

These experiments are implemented by TensorFlow 1.5.0.

### C.2. The first 50 iterations of FashionMNIST experiments in main paper

Figure 8 shows the first 50 iterations of FashionMNIST experiments in main paper. We observe that SGD, GLD 1st eigvec($H$), GLD Hessian and GLD leading successfully escape from the sharp minima found by GD, while GLD diag, GLD dynamic, GLD const and GD do not.

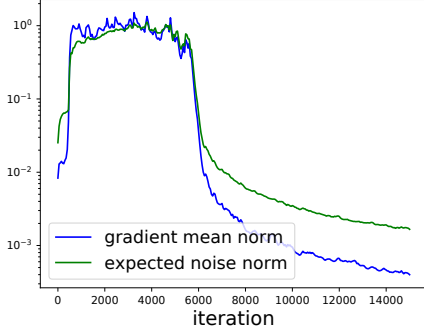These experiments are implemented by TensorFlow 1.5.0.

*Figure 7.* $L_2$ norm of gradient mean vs. the expected norm of noise during the training using SGD. The dataset and model are same as the experiments of FashionMNIST in main paper, or as in Section D.3
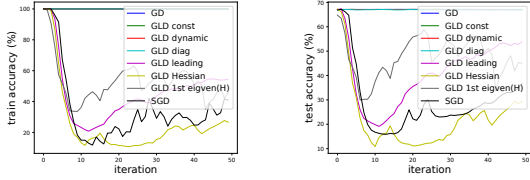


*Figure 8.* The fisrt 50 iterations of FashionMNIST experiments in main paper. Compared dynamics are initialized at $\theta_{GD}^*$ found by GD. The learning rate is same for all the compared methods, $\eta_t = 0.07$, and batch size $m = 20$. **Left**: Training accuracy versus iteration. **Right**: Test accuracy versus iteration.
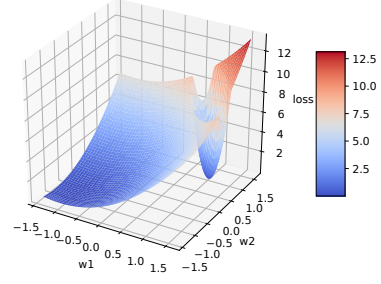


*Figure 9.* Constructed 2-dimensional surface in main paper.

# D. Detailed setups for experiments in main paper

### D.1. Two-dimensional toy example

**Loss Surface**    The loss surface $L(w_1, w_2)$ is constructed by,

$$s_1 = w_1 - 1 - x_1,$$
$$s_2 = w_2 - 1 - x_2,$$
$$\ell(w_1, w_2; x_1, x_2) = \min\{10(s_1 \cos \theta - s_2 \sin \theta)^2$$
$$+ 100(s_1 \cos \theta + s_2 \sin \theta)^2, (w_1 - x_1 + 1)^2 + (w_2 - x_2 + 1)^2\},$$
$$L(w_1, w_2) = \frac{1}{N} \sum_{k=1}^{N} \ell(w_1, w_2; x_1^k, x_2^k),$$

where

$$\theta = \frac{1}{4}\pi,$$
$$N = 100,$$
$$x^k \sim \mathcal{N}(0, \Sigma), \quad \Sigma = \begin{pmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{pmatrix}.$$

Note that $\Sigma$ is the inverse of the Hessian of the quadric form generalizeing the sharp minima. And the 3-dimensional plot of the loss surface is shown in Figure 9.

**Hyperparameters**    All learning rates are equal to 0.005. All dynamics concerned are tuned to share the same expected square norm, 0.01. The number of iteration during one run is 500.

These experiments are implemented by PyTorch 0.3.0.

### D.2. One hidden layer network

**Hyperparameters**    The $\delta$ is set to be 0.001. The learning rate is 0.001. The optimizer is Adam for fast convergence, which does not affect our point on studying $\text{Tr}(H\Sigma)$.

The code is implemented in TensorFlow 1.9.0.

## D.3. FashionMNIST with corrupted labels

**Dataset**   Our training set consists of $1,200$ examples randomly sampled from original FashionMNIST training set, and we further specify 200 of them with randomly wrong labels. The test set is same as the original FashionMNIST test set.

**Model**   Network architecture:

$$\text{input} \Rightarrow \text{conv1} \Rightarrow \text{max\_pool} \Rightarrow \text{ReLU} \Rightarrow \text{conv2}$$
$$\Rightarrow \text{max\_pool} \Rightarrow \text{ReLU} \Rightarrow \text{fc1} \Rightarrow \text{ReLU}$$
$$\Rightarrow \text{fc2} \Rightarrow \text{output}.$$

Both two convolutional layers use $5 \times 5$ kernels with $10$ channels and no padding. The number of hidden units between fully connected layers are $50$. The total number of parameters of this network are $11,330$.

**Training details**

- **GD**: Learning rate $\eta = 0.1$. We tuned the learning rate (in diffusion stage) in a wide range of $\{0.5, 0.2, 0.15, 0.1, 0.09, 0.08, \ldots, 0.01\}$ and no improvement on generalization.

- **GLD constant**: Learning rate $\eta = 0.07$, noise std $\sigma = 10^{-3}$. We tuned the noise std in range of $\{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$ and no improvement on generalization.

- **GLD dynamic**: Learning rate $\eta = 0.07$.

- **GLD diagnoal**: Learning rate $\eta = 0.07$.

- **GLD leading**: Learning rate $\eta = 0.07$, number of leading eigenvalues $k = 20$, batchsize $m = 20$. We first randomly divide the training set into 60 mini batches containing 20 examples, and then use those minibatches to estimate covariance matrix.

- **GLD Hessian**: Learning rate $\eta = 0.07$, number of leading eigenvalues $= 20$, update frequence $f = 10$. Do to the limit of computational resources, we only update Hessian matrix every 10 iterations. But add Hessian generated noise every iteration. And to the same reason, we simplify set the coefcent of Hessian noise to $\sqrt{\text{Tr}H/m\text{Tr}\Sigma}$, to avoid extensively tuning of hyperparameter.

- **GLD 1st eigvec**($H$): Learning rate $\eta = 0.07$, as for GLD Hessian, and we set the coefficient of noise to $\sqrt{\lambda_1/m\text{Tr}\Sigma}$, where $\lambda_1$ is the first eigenvalue of $H$.

- **SGD**: Learning rate $\eta = 0.07$, batchsize $m = 20$.

**Estimation of Sharpness**   The sharpness are estimated by

$$\frac{1}{M} \sum_{j=1}^{M} L(\theta + \nu_j) - L(\theta), \quad \nu_j \sim \mathcal{N}(0, \delta^2 I), \quad (88)$$

with $M = 1,000$ and $\delta = 0.01$.

These experiments are implemented by TensorFlow 1.5.0.

## D.4. SVHN and CIFAR-10

**Dataset**   For SVHN experiments, we use $2,5000$ examples for training and $7,5000$ examples for test, to compromise with the computational burden of gradient descent. And for CIFAR-10 experiments, we use standard CIFAR-10 datasets. We do not use data augmentation since it could cause uncontrollable affects on analyzing SGD noise.

**Model**   Standard VGG11 network without any regularizations including dropout, batch normalization, weight decay, etc. The total number of parameters of this network is $9,750,922$.

We choose VGG11 instead of ResNet because VGG11 achieves good generalization performance without using *Batch Normalization*, which has a subtle impact on SGD noise.

**Training details**   Learning rates $\eta_t = 0.05$ are fixed for all optimizers, which is tuned for the best generalization performance of GD. The batch size of SGD is $m = 100$. The noise std of GLD constant is $\sigma = 10^{-3}$, which is tuned to best. Due to computational limitation, we only conduct experiments on GD, GLD const, GLD dynamic, GLD diag and SGD.

**Estimation of Sharpness**   The sharpness are estimated by

$$\frac{1}{M} \sum_{j=1}^{M} L(\theta + \nu_j) - L(\theta), \quad \nu_j \sim \mathcal{N}(0, \delta^2 I), \quad (89)$$

with $M = 100$ and $\delta = 0.01$.

These experiments are implemented by PyTorch 1.0.0.