
Supplementary Material for Data-Dependent Stability of Stochastic Gradient Descent

Ilja Kuzborskij¹ Christoph H. Lampert²

1. Proofs

In this section we present proofs of all the statements.

Proof of Theorem 2. Indicate by $S = \{z_i\}_{i=1}^m$ and $S' = \{z'_i\}_{i=1}^m$ independent training sets sampled i.i.d. from \mathcal{D} , and let $S^{(i)} = \{z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_m\}$, such that $z'_i \stackrel{\text{iid}}{\sim} \mathcal{D}$. We relate expected empirical risk and expected risk by

$$\begin{aligned} \mathbb{E}_S \mathbb{E}_A \left[\widehat{R}_S(A_S) \right] &= \mathbb{E}_S \mathbb{E}_A \left[\frac{1}{m} \sum_{i=1}^m f(A_S, z_i) \right] \\ &= \mathbb{E}_{S, S'} \mathbb{E}_A \left[\frac{1}{m} \sum_{i=1}^m f(A_{S^{(i)}}, z'_i) \right] \\ &= \mathbb{E}_{S, S'} \mathbb{E}_A \left[\frac{1}{m} \sum_{i=1}^m f(A_S, z'_i) \right] - \delta \\ &= \mathbb{E}_S \mathbb{E}_A [R(A_S)] - \delta, \end{aligned}$$

where

$$\begin{aligned} \delta &= \mathbb{E}_{S, S'} \mathbb{E}_A \left[\frac{1}{m} \sum_{i=1}^m (f(A_S, z'_i) - f(A_{S^{(i)}}, z'_i)) \right] \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{S, z'_i} \mathbb{E}_A [f(A_S, z'_i) - f(A_{S^{(i)}}, z'_i)]. \end{aligned}$$

Renaming z'_i as z and taking sup over i we get that

$$\delta \leq \sup_{i \in [m]} \left\{ \mathbb{E}_{S, z} \mathbb{E}_A [f(A_S, z) - f(A_{S^{(i)}}, z)] \right\}.$$

This completes the proof. \square

1.1. Preliminaries

We say that the Stochastic Gradient Descent (SGD) gradient update rule is an operator $G_t : \mathcal{H} \mapsto \mathcal{H}$, such that

$$G_t(\mathbf{w}) := \mathbf{w} - \alpha_t \nabla f(\mathbf{w}, z_{i_t}),$$

¹University of Milan, Italy ²IST Austria. Correspondence to: Ilja Kuzborskij <ilja.kuzborskij@gmail.com>.

and it is also a function of the training set S and a random index set I . Then, $\mathbf{w}_{t+1} = G_t(\mathbf{w}_t)$, throughout $t = 1, \dots, T$. Recall the use of notation $\mathbf{w}_{S,t}$ to indicate the output of SGD ran on a training set S , at step t , and define

$$\delta_t(S, z) := \|\mathbf{w}_{S,t} - \mathbf{w}_{S^{(i)},t}\|.$$

Next, we summarize a few instrumental facts about G_t and few statements about the loss functions used in our proofs.

Definition 1 (Expansiveness). *A gradient update rule is η -expansive if for all \mathbf{w}, \mathbf{v} ,*

$$\|G_t(\mathbf{w}) - G_t(\mathbf{v})\| \leq \eta \|\mathbf{w} - \mathbf{v}\|.$$

The following lemma characterizes expansiveness for the gradient update rule under different assumptions on f .

Lemma 1 (Lemma 3.6 in (Hardt et al., 2016)). *Assume that f is β -smooth. Then, we have that:*

- 1) G_t is $(1 + \alpha_t \beta)$ -expansive,
- 2) If f in addition is convex, then, for any $\alpha_t \leq \frac{2}{\beta}$, the gradient update rule G_t is 1-expansive.

An important consequence of β -smoothness of f is self-boundedness (Shalev-Shwartz & Ben-David, 2014), which we will use on many occasions.

Lemma 2 (Self-boundedness). *For β -smooth non-negative function f we have that*

$$\|\nabla f(\mathbf{w}, z)\| \leq \sqrt{2\beta f(\mathbf{w}, z)}.$$

Self-boundedness in turn implies the following boundedness of a gradient update rule.

Corollary 1. *Assume that f is β -smooth and non-negative. Then,*

$$\begin{aligned} \|\mathbf{w} - G_t(\mathbf{w})\| &= \alpha_t \|\nabla f(\mathbf{w}, z_{j_t})\| \\ &\leq \alpha_t \min \left\{ \sqrt{2\beta f(\mathbf{w}, z_{j_t})}, L \right\}. \end{aligned}$$

Proof. By Lemma 2

$$\|\alpha_t \nabla f(\mathbf{w}, z_{j_t})\| \leq \alpha_t \sqrt{2\beta f(\mathbf{w}, z_{j_t})},$$

and also by Lipschitzness of f , $\|\alpha_t \nabla f(\mathbf{w}, z_{j_t})\| \leq \alpha_t L$. \square

Next we introduce a bound that relates the risk of the output at step t to the risk of the initialization point \mathbf{w}_1 through the variance of the gradient. Given an appropriate choice of step size, this bound will be crucial at stating stability bounds that depend on the risk at \mathbf{w}_1 . The proof idea is similar to the one of (Ghadimi & Lan, 2013). In particular, it does not require convexity of the loss function.

Lemma 3. *Suppose SGD is ran with step sizes $\alpha_1, \dots, \alpha_{t-1} \leq \frac{1}{\beta}$ w.r.t. the β -smooth loss f . Then we have that*

$$\begin{aligned} & \sum_{k=1}^{t-1} \left(\alpha_k - \frac{\alpha_k^2 \beta}{2} \right) \mathbb{E}_S [\|\nabla R(\mathbf{w}_{S,k})\|^2] \\ & \leq R(\mathbf{w}_1) - R(\mathbf{w}_t) \\ & + \frac{\beta}{2} \sum_{k=1}^{t-1} \alpha_k^2 \mathbb{E}_S [\|\nabla f(\mathbf{w}_{S,k}, z_{j_k}) - \nabla R(\mathbf{w}_{S,k})\|^2]. \end{aligned}$$

Proof. For brevity denote $f_k(\mathbf{w}) \equiv f(\mathbf{w}, z_{j_k})$. By β -smoothness of R and recalling that the SGD update rule $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_k(\mathbf{w}_k)$, we have

$$\begin{aligned} & R(\mathbf{w}_{k+1}) - R(\mathbf{w}_k) \\ & \leq \nabla R(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{\beta}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|^2 \\ & = -\alpha_k \nabla R(\mathbf{w}_k)^\top \nabla f_k(\mathbf{w}_k) + \frac{\beta \alpha_k^2}{2} \|\nabla f_k(\mathbf{w}_k)\|^2 \\ & = -\alpha_k \nabla R(\mathbf{w}_k)^\top \nabla f_k(\mathbf{w}_k) \\ & + \frac{\beta \alpha_k^2}{2} \|\nabla f_k(\mathbf{w}_k) - \nabla R(\mathbf{w}_k) + \nabla R(\mathbf{w}_k)\|^2 \\ & = -(\alpha_k + \alpha_k^2 \beta) \nabla R(\mathbf{w}_k)^\top \nabla f_k(\mathbf{w}_k) \\ & + \frac{3\alpha_k^2 \beta}{2} \|\nabla R(\mathbf{w}_k)\|^2 + \frac{\beta \alpha_k^2}{2} \|\nabla f_k(\mathbf{w}_k) - \nabla R(\mathbf{w}_k)\|^2. \end{aligned}$$

Taking expectation w.r.t. S on both sides, recalling that $\mathbb{E}_{z_k} [\nabla f_k(\mathbf{w}_k)] = \nabla R(\mathbf{w}_k)$ and rearranging terms we get

$$\begin{aligned} & \left(\alpha_k - \frac{\alpha_k^2 \beta}{2} \right) \mathbb{E} [\|\nabla R(\mathbf{w}_k)\|^2] \leq R(\mathbf{w}_k) - R(\mathbf{w}_{k+1}) \\ & + \frac{\beta \alpha_k^2}{2} \mathbb{E} [\|\nabla f_k(\mathbf{w}_k) - \nabla R(\mathbf{w}_k)\|^2], \end{aligned}$$

and summing above over $k = 1, \dots, t-1$ we get the statement. \square

Lemma 4. *Suppose SGD is ran with step sizes $\alpha_1, \dots, \alpha_{t-1} \leq \frac{1}{\beta}$ on the β -smooth loss f . Assume that the variance of stochastic gradients obeys*

$$\mathbb{E}_{S,z} [\|\nabla f(\mathbf{w}_{S,k}, z) - \nabla R(\mathbf{w}_{S,k})\|^2] \leq \sigma^2 \quad \forall k \in [T].$$

Then we have that

$$\begin{aligned} & \mathbb{E}_S \left[\sum_{k=1}^{t-1} \alpha_k \|\nabla f(\mathbf{w}_{S,k}, z_k)\| \right] \\ & \leq 2 \sqrt{\left(\sum_{k=1}^{t-1} \alpha_k \right) \left(R(\mathbf{w}_1) - R^* + \frac{\beta \sigma^2}{2} \sum_{k=1}^{t-1} \alpha_k^2 \right)} \\ & + \sigma \sum_{k=1}^{t-1} \alpha_k. \end{aligned}$$

Proof. First we perform the decomposition,

$$\begin{aligned} & \mathbb{E}_S \left[\sum_{k=1}^{t-1} \alpha_k \|\nabla f(\mathbf{w}_{S,k}, z_k)\| \right] \\ & = \sum_{k=1}^{t-1} \alpha_k \mathbb{E}_S [\|\nabla R(\mathbf{w}_{S,k})\|] \\ & + \sum_{k=1}^{t-1} \alpha_k \mathbb{E}_S [\|\nabla f(\mathbf{w}_{S,k}, z_k) - \nabla R(\mathbf{w}_{S,k})\|] \\ & \leq \sum_{k=1}^{t-1} \alpha_k \mathbb{E}_S [\|\nabla R(\mathbf{w}_{S,k})\|] + \sigma \sum_{k=1}^{t-1} \alpha_k. \end{aligned} \quad (1)$$

Introduce

$$Q_t := \sum_{k=1}^{t-1} \left(\alpha_k - \frac{\alpha_k^2 \beta}{2} \right).$$

Now we invoke the stationary-point argument to bound the first term above as

$$\begin{aligned} & \sum_{k=1}^{t-1} \alpha_k \mathbb{E}_S \left[\sqrt{\|\nabla R(\mathbf{w}_k)\|^2} \right] \\ & \leq \sum_{k=1}^{t-1} \frac{\left(1 - \frac{\alpha_k \beta}{2}\right)}{\left(1 - \frac{\alpha_k \beta}{2}\right)} \cdot \alpha_k \sqrt{\mathbb{E}_S [\|\nabla R(\mathbf{w}_k)\|^2]} \end{aligned} \quad (2)$$

(By Jensen's inequality)

$$\leq 2 \sum_{k=1}^{t-1} \left(\alpha_k - \frac{\alpha_k^2 \beta}{2} \right) \sqrt{\mathbb{E}_S [\|\nabla R(\mathbf{w}_k)\|^2]} \quad \text{(Assuming that } \alpha_k \leq \frac{1}{\beta} \text{)}$$

$$= \frac{2Q_t}{Q_t} \sum_{k=1}^{t-1} \left(\alpha_k - \frac{\alpha_k^2 \beta}{2} \right) \sqrt{\mathbb{E}_S [\|\nabla R(\mathbf{w}_k)\|^2]} \quad (3)$$

$$\leq 2\sqrt{Q_t} \sqrt{\sum_{k=1}^{t-1} \left(\alpha_k - \frac{\alpha_k^2 \beta}{2} \right) \mathbb{E}_S [\|\nabla R(\mathbf{w}_k)\|^2]}$$

(By Jensen's inequality)

$$\leq 2\sqrt{Q_t} \sqrt{R(\mathbf{w}_1) - R(\mathbf{w}_t) + \frac{\beta \sigma^2}{2} \sum_{k=1}^{t-1} \alpha_k^2}.$$

(By Lemma 3)

Combining this with (1) gives the statement and completes the proof. \square

The following lemma is similar to Lemma 3.11 of (Hardt et al., 2016), and is instrumental in bounding the stability of SGD. However, we make an adjustment and state it in expectation over the data. Note that it does not require convexity of the loss function.

Lemma 5. *Assume that the loss function $f(\cdot, z) \in [0, 1]$ is L -Lipschitz for all z . Then, for every $t_0 \in \{0, 1, 2, \dots, m\}$ we have that,*

$$\begin{aligned} & \mathbb{E}_{S,z} \mathbb{E}_A [f(\mathbf{w}_{S,T}, z) - f(\mathbf{w}_{S^{(i)},T}, z)] \\ & \leq L \mathbb{E}_{S,z} \left[\mathbb{E}_A [\delta_T(S, z) \mid \delta_{t_0}(S, z) = 0] \right] + \mathbb{E}_{S,A} [R(A_S)] \frac{t_0}{m}. \end{aligned} \quad (4)$$

Proof. We proceed with elementary decomposition, Lipschitzness of f , and using the fact that f is non-negative to have that

$$\begin{aligned} & f(\mathbf{w}_{S,T}, z) - f(\mathbf{w}_{S^{(i)},T}, z) \\ & = (f(\mathbf{w}_{S,T}, z) - f(\mathbf{w}_{S^{(i)},T}, z)) \mathbb{I} \{\delta_{t_0}(S, z) = 0\} \\ & \quad + (f(\mathbf{w}_{S,T}, z) - f(\mathbf{w}_{S^{(i)},T}, z)) \mathbb{I} \{\delta_{t_0}(S, z) \neq 0\} \\ & \leq L \delta_T(S, z) \mathbb{I} \{\delta_{t_0}(S, z) = 0\} \\ & \quad + f(\mathbf{w}_{S,T}, z) \mathbb{I} \{\delta_{t_0}(S, z) \neq 0\}. \end{aligned} \quad (7)$$

Taking expectation w.r.t. algorithm randomization, we get that

$$\begin{aligned} & \mathbb{E}_A [f(\mathbf{w}_{S,T}, z) - f(\mathbf{w}_{S^{(i)},T}, z)] \\ & \leq L \mathbb{E}_A [\delta_T(S, z) \mathbb{I} \{\delta_{t_0}(S, z) = 0\}] \\ & \quad + \mathbb{E}_A [f(\mathbf{w}_{S,T}, z) \mathbb{I} \{\delta_{t_0}(S, z) \neq 0\}]. \end{aligned} \quad (8)$$

Recall that $i \in [m]$ is the index where S and $S^{(i)}$ differ, and introduce a random variable τ_A taking on the index of the first time step where SGD uses the example z_i or a replacement z . Note also that τ_A does not depend on the data. When $\tau_A > t_0$, then it must be that $\delta_{t_0}(S, z) = 0$, because updates on both S and $S^{(i)}$ are identical until t_0 . A consequence of this is that $\mathbb{I} \{\delta_{t_0}(S, z) \neq 0\} \leq \mathbb{I} \{\tau_A \leq t_0\}$. Thus the rightmost term in (8) is bounded as

$$\begin{aligned} & \mathbb{E}_A [f(\mathbf{w}_{S,T}, z) \mathbb{I} \{\delta_{t_0}(S, z) \neq 0\}] \\ & \leq \mathbb{E}_A [f(\mathbf{w}_{S,T}, z) \mathbb{I} \{\tau_A \leq t_0\}]. \end{aligned}$$

Now, focus on the r.h.s. above. Recall that we assume randomization by sampling from the uniform distribution over $[m]$ without replacement, and denote a realization by $\{j_i\}_{i=1}^m$. Then, we can always express our randomization as permutation function $\pi_A(S) = \{z_{j_i}\}_{i=1}^m$. In addition,

introduce an algorithm $\text{GD} : \mathcal{Z}^m \mapsto \mathcal{H}$, which is identical to A , except that it passes over the training set S sequentially without randomization. That said, we have that

$$\mathbb{E}_A [f(\mathbf{w}_{S,T}, z) \mathbb{I} \{\tau_A \leq t_0\}] = \mathbb{E}_A [f(\text{GD}_{\pi_A(S)}, z) \mathbb{I} \{\tau_A \leq t_0\}],$$

and taking expectation over the data,

$$\begin{aligned} & \mathbb{E}_{S,z} \left[\mathbb{E}_A [f(\mathbf{w}_{S,T}, z) \mathbb{I} \{\tau_A \leq t_0\}] \right] \\ & = \mathbb{E}_A \left[\mathbb{E}_{S,z} [f(\text{GD}_{\pi_A(S)}, z) \mathbb{I} \{\tau_A \leq t_0\}] \right]. \end{aligned}$$

Now observe that for any realization of A , $\mathbb{E}_{S,z} [f(\text{GD}_{\pi_A(S)}, z)] = \mathbb{E}_A \mathbb{E}_{S,z} [f(A_S, z)]$ because expectation w.r.t. S and z does not change under our randomization¹. Thus, we have that

$$\begin{aligned} & \mathbb{E}_A \left[\mathbb{E}_{S,z} [f(\text{GD}_{\pi_A(S)}, z) \mathbb{I} \{\tau_A \leq t_0\}] \right] \\ & = \mathbb{E}_{S,A} [R(A_S)] \mathbb{P}(\tau_A \leq t_0). \end{aligned}$$

Now assuming that τ_A is uniformly distributed over $[m]$ we have that

$$\mathbb{P}(\tau_A \leq t_0) = \frac{t_0}{m}.$$

Putting this together with (6) and (7), we finally get that

$$\begin{aligned} & \mathbb{E}_{S,z} \mathbb{E}_A [f(\mathbf{w}_{S,T}, z) - f(\mathbf{w}_{S^{(i)},T}, z)] \\ & \leq L \mathbb{E}_{S,z} \left[\mathbb{E}_A [\delta_T(S, z) \mathbb{I} \{\delta_{t_0}(S, z) = 0\}] \right] + \mathbb{E}_{S,A} [R(A_S)] \frac{t_0}{m} \\ & \leq L \mathbb{E}_{S,z} \left[\mathbb{E}_A [\delta_T(S, z) \mid \delta_{t_0}(S, z) = 0] \right] + \mathbb{E}_{S,A} [R(A_S)] \frac{t_0}{m}. \end{aligned}$$

This completes the proof. \square

We spend a moment to highlight the role of conditional expectation in (5). Observe that we could naively bound (4) by the Lipschitzness of f , but Lemma 5 follows a more careful argument. First note that t_0 is a free parameter. The expected distance in (5) between SGD outputs $\mathbf{w}_{S,t}$ and $\mathbf{w}_{S^{(i)},t}$ is conditioned on the fact that at step t_0 outputs of SGD are still the same. This means that the perturbed point is encountered after t_0 . Then, the conditional expectation should be a decreasing function of t_0 : the later the perturbation occurs, the smaller deviation between $\mathbf{w}_{S,t}$ and $\mathbf{w}_{S^{(i)},t}$ we should expect. Later we use this fact to minimize the bound (5) over t_0 .

¹Strictly speaking we could omit $\mathbb{E}_A[\cdot]$ and consider *any* randomization by reshuffling, but we keep expectation for the sake of clarity.

1.2. Convex Losses

In this section we prove on-average stability for loss functions that are non-negative, β -smooth, and convex.

Theorem 1. *Assume that f is convex, and that SGD's is ran with step sizes $\{\alpha_t\}_{t=1}^T$. Then, for every $t_0 \in \{0, 1, 2, \dots, m\}$, SGD is $\epsilon(\mathcal{D}, \mathbf{w}_1)$ -on-average stable with*

$$\begin{aligned} & \epsilon(\mathcal{D}, \mathbf{w}_1) \\ & \leq \frac{2}{m} \sum_{t=t_0+1}^T \alpha_t \mathbb{E}_{S,z} [\|\nabla f(\mathbf{w}_t, z_{j_t})\|] + \mathbb{E}_{S,A} [R(A_S)] \frac{t_0}{m}. \end{aligned}$$

Proof. For brevity denote $\Delta_t(S, z) := \mathbb{E}_A [\delta_t(S, z) \mid \delta_{t_0}(S, z) = 0]$. We start by applying Lemma 5:

$$\begin{aligned} & \mathbb{E}_{S,z,A} [f(\mathbf{w}_{S,T}, z) - f(\mathbf{w}_{S^{(i)},T}, z)] \\ & \leq L \mathbb{E}_{S,z} [\Delta_T(S, z)] + \mathbb{E}_{S,A} [R(A_S)] \frac{t_0}{m}. \quad (9) \end{aligned}$$

Our goal is to bound the first term on the r.h.s. as a decreasing function of t_0 , so that eventually we can minimize the bound w.r.t. t_0 . At this point we focus on the first term, and the proof partially follows the outline of the proof of Theorem 3.7 in (Hardt et al., 2016). The strategy will be to establish the bound on $\Delta_T(S, z)$ by using a recursive argument. In fact we will state the bound on $\Delta_{t+1}(S, z)$ in terms of $\Delta_t(S, z)$ and then unravel the recursion. Finally, we will take expectation w.r.t. the data after we obtain the bound by recursion.

To do so, we distinguish two cases: 1) SGD encounters a perturbed point at step t , that is $t = i$, and 2) the current point is the same in S and $S^{(i)}$, so $t \neq i$. For the first case, we will use data-dependent boundedness of the gradient update rule, Corollary 1, that is

$$\|G_t(\mathbf{w}_{S,t}) - G_t(\mathbf{w}_{S^{(i)},t})\| \leq \delta_t(S, z) + 2\alpha_t \|\nabla f(\mathbf{w}_{S,t}, z_{j_t})\|.$$

To handle the second case, we will use the expansiveness of the gradient update rule, Lemma 1, which states that for convex loss functions, the gradient update rule is 1-expansive, so $\delta_{t+1}(S, z) \leq \delta_t(S, z)$. Considering both cases of example selection, and noting that SGD encounters the perturbation w.p. $\frac{1}{m}$, we write \mathbb{E}_A for a step t as

$$\begin{aligned} \Delta_{t+1}(S, z) & \leq \left(1 - \frac{1}{m}\right) \Delta_t(S, z) \\ & \quad + \frac{1}{m} (\Delta_t(S, z) + 2\alpha_t \|\nabla f(\mathbf{w}_{S,t}, z_{j_t})\|) \\ & = \Delta_t(S, z) + \frac{2\alpha_t \|\nabla f(\mathbf{w}_{S,t}, z_{j_t})\|}{m}. \end{aligned}$$

Unraveling the recursion from T to t_0 and plugging the above into (9) yields

$$\begin{aligned} & \mathbb{E}_A \mathbb{E}_{S,z} [\delta_T(S, z)] \\ & \leq \frac{2}{m} \sum_{t=t_0+1}^T \alpha_t \mathbb{E}_{S,z} [\|\nabla f(\mathbf{w}_t, z_{j_t})\|] + \mathbb{E}_{S,A} [R(A_S)] \frac{t_0}{m}. \end{aligned}$$

This completes the proof. \square

Next statement is a simple consequence of Theorem 1 and Lemma 4.

Proof of Theorem 3. Consider Theorem 1 and set $t_0 = 0$.

$$\epsilon(\mathcal{D}, \mathbf{w}_1) \leq \frac{2}{m} \sum_{t=1}^T \alpha_t \mathbb{E}_{S,z} [\|\nabla f(\mathbf{w}_{S,t}, z_{j_t})\|]. \quad (10)$$

Bounding the sum using Lemma 4 recalling that $\alpha_t = c/\sqrt{t}$, we get

$$\begin{aligned} & \mathbb{E}_S \left[\sum_{t=1}^T \alpha_t \|\nabla f(\mathbf{w}_t, z_{j_t})\| \right] \\ & \leq 2\sqrt{\left(\sum_{t=1}^T \alpha_t \right)} \left(R(\mathbf{w}_1) - R^* + \frac{\beta\sigma^2}{2} \sum_{t=1}^T \alpha_t^2 \right) \\ & \quad + \sigma \sum_{t=1}^T \alpha_t \\ & \leq 2\sqrt{2c} \cdot \sqrt[4]{T} \cdot \sqrt{R(\mathbf{w}_1) - R^*} \\ & \quad + 2c\sigma \left(\sqrt[4]{T} \sqrt{\frac{\beta}{2}} + \sqrt{T} \right). \end{aligned}$$

Combining above with (10) completes the proof. \square

1.3. Non-convex Losses

Our proof of a stability bound for non-convex loss functions, Theorem 4 (in the submission file), follows a general outline of (Hardt et al., 2016, Theorem 3.8). Namely, the outputs of SGD run on a training set S and its perturbed version $S^{(i)}$ will not differ too much, because by the time a perturbation is encountered, the step size has already decayed enough. So, on the one hand, stabilization is enforced by the diminishing the step size, and on the other hand, by how much updates expand the distance between the gradients after the perturbation. Since (Hardt et al., 2016) work with uniform stability, they capture the expansiveness of post-perturbation update by the Lipschitzness of the gradient. In combination with a recursive argument, their bound has exponential dependency on the Lipschitz constant of the gradient. We argue that the Lipschitz continuity of the gradient can be too pessimistic in general. Instead, we rely on a local

data-driven argument: considering that we initialize SGD at point \mathbf{w}_1 , how much do updates expand the gradient under the distribution of interest? The following crucial lemma characterizes such behavior in terms of the curvature at \mathbf{w}_1 .

Lemma 6. *Assume that the loss function $f(\cdot, z)$ is β -smooth and that its Hessian is ρ -Lipschitz. Then,*

$$\|G_t(\mathbf{w}_{S,t}) - G_t(\mathbf{w}_{S^{(i)},t})\| \leq (1 + \alpha_t \xi_t(S, z)) \delta_t(S, z) \quad (11)$$

where

$$\begin{aligned} \xi_t(S, z) &:= \|\nabla^2 f(\mathbf{w}_1, z_t)\|_2 \\ &+ \frac{\rho}{2} \left\| \sum_{k=1}^{t-1} \alpha_k \nabla f(\mathbf{w}_{S,k}, z_k) \right\| \\ &+ \frac{\rho}{2} \left\| \sum_{k=1}^{t-1} \alpha_k \nabla f(\mathbf{w}_{S^{(i)},k}, z_{k'}) \right\|. \end{aligned}$$

Furthermore, for any $t \in [T]$,

$$\begin{aligned} \mathbb{E}_{S,z} [\xi_t(S, z)] &\leq \mathbb{E}_{S,z} [\|\nabla^2 f(\mathbf{w}_1, z_t)\|_2] \\ &+ 2\rho\sqrt{(R(\mathbf{w}_1) - R^*)c(1 + \ln(T))} \\ &+ \rho\sigma\left(\sqrt{2c\beta} + c(1 + \ln(T))\right). \end{aligned}$$

Proof. Recall that the randomness of the algorithm is realized through sampling without replacement from the uniform distribution over $[m]$. Apart from that we will not be concerned with the randomness of the algorithm, and given the set of random variables $\{j_i\}_{i=1}^m$, for brevity we will use indexing notation z_1, z_2, \dots, z_m to indicate $z_{j_1}, z_{j_2}, \dots, z_{j_m}$. Next, let $S^{(i)} = \{z'_i\}_{i=1}^m$, and introduce a shorthand notation $f_k(\mathbf{w}) = f(\mathbf{w}, z_k)$ and $f_{k'}(\mathbf{w}) = f(\mathbf{w}, z'_k)$. We start by applying triangle inequality to get

$$\begin{aligned} \|G_t(\mathbf{w}_{S,t}) - G_t(\mathbf{w}_{S^{(i)},t})\| &\leq \|\mathbf{w}_{S,t} - \mathbf{w}_{S^{(i)},t}\| \\ &+ \alpha_t \|\nabla f_t(\mathbf{w}_{S,t}) - \nabla f_t(\mathbf{w}_{S^{(i)},t})\|. \end{aligned}$$

In the following we will focus on the second term of r.h.s. above. Given SGD outputs $\mathbf{w}_{S,t}$ and $\mathbf{w}_{S^{(i)},t}$ with $t > i$, our goal here is to establish how much do gradients grow apart with every new update. This behavior can be characterized assuming that gradient is Lipschitz continuous, however, we conduct a local analysis. Specifically, we observe how much do updates expand gradients, given that we start at some point \mathbf{w}_1 under the data-generating distribution. So, instead of the Lipschitz constant, expansiveness rather depends on the curvature around \mathbf{w}_1 . On the other hand, we are dealing with outputs at an arbitrary time step t , and therefore we first have to relate them to the initialization point \mathbf{w}_1 . We do so by using the gradient update rule and telescopic sums, and conclude that this relationship is controlled by the sum of gradient norms along the update path. We further

establish that this sum is controlled by the risk of \mathbf{w}_1 up to the noise of stochastic gradients, through stationary-point result of Lemma 4. Thus, the proof consists of two parts: 1) Decomposition into curvature and gradients along the update path, and 2) bounding those gradients.

1) Decomposition. Introduce $\delta_t := \mathbf{w}_{S^{(i)},t} - \mathbf{w}_{S,t}$. By Taylor theorem we get that

$$\begin{aligned} \nabla f_t(\mathbf{w}_{S,t}) - \nabla f_t(\mathbf{w}_{S^{(i)},t}) &= \nabla^2 f_t(\mathbf{w}_1) \delta_t \\ &+ \int_0^1 \left(\nabla^2 f_t(\mathbf{w}_{S,t} + \tau \delta_t) - \nabla^2 f_t(\mathbf{w}_1) \right) d\tau \delta_t. \end{aligned}$$

Taking norm on both sides, applying triangle inequality, Cauchy-Schwartz inequality, and assuming that Hessians are ρ -Lipschitz we obtain

$$\begin{aligned} \|\nabla f_t(\mathbf{w}_{S,t}) - \nabla f_t(\mathbf{w}_{S^{(i)},t})\| &\leq \rho \int_0^1 \|\mathbf{w}_{S,t} - \mathbf{w}_1 + \tau \delta_t\| d\tau \|\delta_t\| + \|\nabla^2 f_t(\mathbf{w}_1)\| \|\delta_t\|. \end{aligned} \quad (12)$$

2) Bounding gradients. Using telescoping sums and SGD update rule we get that

$$\begin{aligned} \mathbf{w}_{S,t} - \mathbf{w}_1 + \tau \delta_t &= \mathbf{w}_{S,t} - \mathbf{w}_1 + \tau (\mathbf{w}_{S^{(i)},t} - \mathbf{w}_1 + \mathbf{w}_1 - \mathbf{w}_{S,t}) \\ &= \sum_{k=1}^{t-1} (\mathbf{w}_{S,k+1} - \mathbf{w}_{S,k}) + \tau \sum_{k=1}^{t-1} (\mathbf{w}_{S^{(i)},k+1} - \mathbf{w}_{S^{(i)},k}) \\ &\quad - \tau \sum_{k=1}^{t-1} (\mathbf{w}_{S,k+1} - \mathbf{w}_{S,k}) \\ &= (\tau - 1) \sum_{k=1}^{t-1} \alpha_k \nabla f_k(\mathbf{w}_{S,k}) - \tau \sum_{k=1}^{t-1} \alpha_k \nabla f_{k'}(\mathbf{w}_{S^{(i)},k}). \end{aligned}$$

Plugging above into the integral of (12) we have

$$\begin{aligned} &\int_0^1 \left\| \sum_{k=1}^{t-1} \alpha_k ((\tau - 1) \nabla f_k(\mathbf{w}_{S,k}) - \tau \nabla f_{k'}(\mathbf{w}_{S^{(i)},k})) \right\| d\tau \\ &\leq \frac{1}{2} \left\| \sum_{k=1}^{t-1} \alpha_k \nabla f_k(\mathbf{w}_{S,k}) \right\| + \frac{1}{2} \left\| \sum_{k=1}^{t-1} \alpha_k \nabla f_{k'}(\mathbf{w}_{S^{(i)},k}) \right\| \\ &\leq \frac{1}{2} \sum_{k=1}^{t-1} \alpha_k \|\nabla f_k(\mathbf{w}_{S,k})\| + \frac{1}{2} \sum_{k=1}^{t-1} \alpha_k \|\nabla f_{k'}(\mathbf{w}_{S^{(i)},k})\|. \end{aligned}$$

Plugging this result back into (12) completes the proof of the first statement. The second statement comes from Lemma 4 with $\alpha_t = c/t$. \square

Next, we need the following statement to prove our stability bound.

Proposition 1 (Bernstein-type inequality). *Let Z be a zero-mean real-valued r.v., such that $|Z| \leq b$ and $\mathbb{E}[Z^2] \leq \sigma^2$. Then for all $|c| \leq \frac{1}{2b}$, we have that $\mathbb{E}[e^{cZ}] \leq e^{c^2\sigma^2}$.*

Proof. Stated inequality is a consequence of a Bernstein-type inequality for moment generating functions, Theorem 2.10 in (Boucheron et al., 2013). Observe that zero-centered r.v. Z bounded by b satisfies Bernstein's condition, that is

$$|\mathbb{E}[(Z - \mathbb{E}[Z])^q]| \leq \frac{q!}{2} \sigma^2 b^{q-2} \quad \text{for all integers } q \geq 3.$$

This in turn satisfies condition for Bernstein-type inequality stating that

$$\mathbb{E}[\exp(c(Z - \mathbb{E}[Z]))] \leq \exp\left(\frac{c^2\sigma^2/2}{1 - b|c|}\right).$$

Choosing $|c| \leq \frac{1}{2b}$ verifies the statement. \square

Now we are ready to prove Theorem 4, which bounds the $\epsilon(\mathcal{D}, \mathbf{w}_1)$ -on-average stability of SGD.

Proof of Theorem 4. For brevity denote

$$r := \mathbb{E}_{S,A} [R(A_S)]$$

and

$$\Delta_t(S, z) := \mathbb{E}_A [\delta_t(S, z) \mid \delta_{t_0}(S, z) = 0].$$

By Lemma 5, for all $t_0 \in [m]$,

$$\mathbb{E}_{S,z} \mathbb{E}_A [f(\mathbf{w}_{S,T}, z) - f(\mathbf{w}_{S^{(i)},T}, z)] \quad (13)$$

$$\leq L \mathbb{E}_{S,z} [\Delta_T(S, z)] + r \frac{t_0}{m}. \quad (14)$$

Most of the proof is dedicated to bounding the first term in (14). We deal with this similarly as in (Hardt et al., 2016). Specifically, we state the bound on $\Delta_T(S, z)$ by using a recursion. In our case, however, we also have an expectation w.r.t. the data, and to avoid complications with dependencies, we first unroll the recursion for the random quantities, and only then take the expectation. At this point the proof crucially relies on the product of exponentials arising from the recursion, and all relevant random quantities end up inside of them. We alleviate this by Proposition 1. Finally, we conclude by minimizing (14) w.r.t. t_0 . Thus we have three steps: 1) recursion, 2) bounding $\mathbb{E}[\exp(\dots)]$, and 3) tuning of t_0 .

1) Recursion. We begin by stating the bound on $\Delta_T(S, z)$ by recursion. Thus we will first state the bound on $\Delta_{t+1}(S, z)$ in terms of $\Delta_t(S, z)$, and other relevant quantities and then unravel the recursion. As in the convex case, we distinguish two cases: 1) SGD encounters the perturbed point at step t , that is $t = i$, and 2) the current point is the same in S and $S^{(i)}$, so $t \neq i$. For the first case, we will use worst-case boundedness of G_t , Corollary 1, that is, $\|G_t(\mathbf{w}_{S,t}) - G_t(\mathbf{w}_{S^{(i)},t})\| \leq \delta_t(S, z) + 2\alpha_t L$. To handle the second case we will use Lemma 6, namely,

$$\|G_t(\mathbf{w}_{S,t}) - G_t(\mathbf{w}_{S^{(i)},t})\| \leq (1 + \alpha_t \xi_t(S, z)) \delta_t(S, z).$$

In addition, as a safety measure we will also take into account that the gradient update rule is at most $(1 + \alpha_t \beta)$ -expansive by Lemma 1. So we will work with the function $\psi_t(S, z) := \min\{\xi_t(S, z), \beta\}$ instead of $\xi_t(S, z)$. and decompose the expectation w.r.t. A for a step t . Noting that SGD encounters the perturbed example with probability $\frac{1}{m}$,

$$\begin{aligned} \Delta_{t+1}(S, z) &\leq \left(1 - \frac{1}{m}\right) (1 + \alpha_t \psi_t(S, z)) \Delta_t(S, z) \\ &\quad + \frac{1}{m} (2\alpha_t L + \Delta_t(S, z)) \\ &= \left(1 + \left(1 - \frac{1}{m}\right) \alpha_t \psi_t(S, z)\right) \Delta_t(S, z) + \frac{2\alpha_t L}{m} \\ &\leq \exp(\alpha_t \psi_t(S, z)) \Delta_t(S, z) + \frac{2\alpha_t L}{m}, \end{aligned} \quad (15)$$

where the last inequality follows from $1 + x \leq \exp(x)$. This inequality is not overly loose for $x \in [0, 1]$, and, in our case it becomes instrumental in handling the recursion.

Now, observe that relation $x_{t+1} \leq a_t x_t + b_t$ with $x_{t_0} = 0$ unwinds from T to t_0 as $x_T \leq \sum_{t=t_0+1}^T b_t \prod_{k=t+1}^T a_k$. Consequently, having $\Delta_{t_0}(S, z) = 0$, we unwind (15) to get

$$\begin{aligned} \Delta_T(S, z) &\leq \sum_{t=t_0+1}^T \left(\prod_{k=t+1}^T \exp\left(\frac{c\psi_k(S, z)}{k}\right) \right) \frac{2cL}{mt} \\ &= \sum_{t=t_0+1}^T \exp\left(c \sum_{k=t+1}^T \frac{\psi_k(S, z)}{k}\right) \frac{2cL}{mt}. \end{aligned} \quad (16)$$

2) Bounding $\mathbb{E}[\exp(\dots)]$. We take expectation w.r.t. S and z on both sides and focus on the expectation of the exponential in (16). First, introduce $\mu_k := \mathbb{E}_{S,z}[\psi_k(S, z)]$, and proceed as

$$\begin{aligned} &\mathbb{E}_{S,z} \left[\exp\left(c \sum_{k=t+1}^T \frac{\psi_k(S, z)}{k}\right) \right] \\ &= \mathbb{E}_{S,z} \left[\exp\left(c \sum_{k=t+1}^T \frac{\psi_k(S, z) - \mu_k}{k}\right) \right] \exp\left(c \sum_{k=t+1}^T \frac{\mu_k}{k}\right). \end{aligned} \quad (17)$$

Observe that zero-mean version of $\psi_k(S, z)$ is bounded as

$$\sum_{k=t+1}^T \frac{|\psi_k(S, z) - \mu_k|}{k} \leq 2\beta \ln(T),$$

and assume the setting of c as $c \leq \frac{1}{2(2\beta \ln(T))^2}$. By Proposition 1, we have

$$\begin{aligned} & \mathbb{E} \left[\exp \left(c \sum_{k=t+1}^T \frac{\psi_k(S, z) - \mu_k}{k} \right) \right] \\ & \leq \exp \left(c^2 \mathbb{E} \left[\left(\sum_{k=t+1}^T \frac{\psi_k(S, z) - \mu_k}{k} \right)^2 \right] \right) \\ & = \exp \left(\frac{c}{2} \mathbb{E} \left[\left(\frac{1}{2\beta \ln(T)} \sum_{k=t+1}^T \frac{\psi_k(S, z) - \mu_k}{k} \right)^2 \right] \right) \\ & \leq \exp \left(\frac{c}{2} \mathbb{E} \left[\left| \sum_{k=t+1}^T \frac{\psi_k(S, z) - \mu_k}{k} \right| \right] \right) \\ & \leq \exp \left(\frac{c}{2} \sum_{k=t+1}^T \frac{\mathbb{E} [|\psi_k(S, z) - \mu_k|]}{k} \right) \\ & \leq \exp \left(c \sum_{k=t+1}^T \frac{\mu_k}{k} \right). \end{aligned}$$

Getting back to (17) we conclude that

$$\mathbb{E}_{S,z} \left[\exp \left(c \sum_{k=t+1}^T \frac{\psi_k(S, z)}{k} \right) \right] \leq \exp \left(c \sum_{k=t+1}^T \frac{2\mu_k}{k} \right). \quad (18)$$

Next, we give an upper-bound on μ_k , that is $\mu_k \leq \min \{ \beta, \mathbb{E}_{S,z} [\xi_k(S, z)] \}$. Finally, we bound $\mathbb{E}_{S,z} [\xi_k(S, z)]$ using the second result of Lemma 6, which holds for any $k \in [T]$, to get that $\mu_k \leq \gamma$, with γ defined in the statement of the theorem.

3) Tuning of t_0 . Now we turn our attention back to (16). Considering that we took an expectation w.r.t. the data, we use (18) and the fact that $\mu_k \leq \gamma$ to get that

$$\begin{aligned} \mathbb{E}_{S,z} [\Delta_T(S, z)] & \leq \sum_{t=t_0+1}^T \exp \left(2c\gamma \sum_{k=t+1}^T \frac{1}{k} \right) \frac{2cL}{mt} \\ & \leq \sum_{t=t_0+1}^T \exp \left(2c\gamma \ln \left(\frac{T}{t} \right) \right) \frac{2cL}{mt} \\ & = \frac{2cL}{m} (T^{2c\gamma}) \sum_{t=t_0+1}^T t^{-2c\gamma-1} \\ & \leq \frac{1}{2c\gamma} \frac{2cL}{m} \left(\frac{T}{t_0} \right)^{2c\gamma}. \end{aligned}$$

Plug the above into (14) to get

$$\begin{aligned} & \mathbb{E}_{S,z} \mathbb{E}_A [f(\mathbf{w}_{S,T}, z) - f(\mathbf{w}_{S^{(i)},T}, z)] \\ & \leq \frac{L^2}{\gamma m} \left(\frac{T}{t_0} \right)^{2c\gamma} + r \frac{t_0}{m}. \end{aligned} \quad (19)$$

Let $q = 2c\gamma$. Then, setting

$$t_0 = \left(\frac{2cL^2}{r} \right)^{\frac{1}{1+q}} T^{\frac{q}{1+q}}$$

minimizes (19). Plugging t_0 back we get that (19) equals to

$$\frac{1 + \frac{1}{q}}{m} (2cL^2)^{\frac{1}{1+q}} (rT)^{\frac{q}{1+q}}.$$

This completes the proof. \square

1.3.1. OPTIMISTIC RATES FOR LEARNING WITH NON-CONVEX LOSS FUNCTIONS

Next we will prove an optimistic bound based on Theorem 4, in other words, the bound that demonstrates fast convergence rate subject to the vanishing empirical risk. First we will need the following technical statement.

Lemma 7. (Cucker & Zhou, 2007, Lemma 7.2) *Let $c_1, c_2, \dots, c_l > 0$ and $s > q_1 > q_2 > \dots > q_{l-1} > 0$. Then the equation*

$$x^s - c_1 x^{q_1} - c_2 x^{q_2} - \dots - c_{l-1} x^{q_{l-1}} - c_l = 0$$

has a unique positive solution x^* . In addition,

$$x^* \leq \max \left\{ (lc_1)^{\frac{1}{s-q_1}}, (lc_2)^{\frac{1}{s-q_2}}, \dots, (lc_{l-1})^{\frac{1}{s-q_{l-1}}}, (lc_l)^{\frac{1}{s}} \right\}.$$

Next we prove a useful technical lemma similarly as in (Orabona, 2014, Lemma 7).

Lemma 8. *Let $a, c > 0$ and $0 < \alpha < 1$. Then the inequality*

$$x - ax^\alpha - c \leq 0$$

implies

$$x \leq \max \left\{ 2^{\frac{\alpha}{1-\alpha}} a^{\frac{1}{1-\alpha}}, (2c)^\alpha a \right\} + c.$$

Proof. Consider a function $h(x) = x - ax^\alpha - c$. Applying Lemma 7 with $s = 1$, $l = 2$, $c_1 = a$, $c_2 = c$, and $q_1 = \alpha$ we get that $h(x) = 0$ has a unique positive solution x^* and

$$x^* \leq \max \left\{ (2a)^{\frac{1}{1-\alpha}}, 2c \right\}. \quad (20)$$

Moreover, the inequality $h(x) \leq 0$ is verified for $x = 0$, and $\lim_{x \rightarrow +\infty} h(x) = +\infty$, so we have that $h(x) \leq 0$ implies

$x \leq x^*$. Now, using this fact and the fact that $h(x^*) = 0$, we have that

$$x \leq x^* = a(x^*)^\alpha + c,$$

and upper-bounding x^* by (20) we finally have

$$x \leq a \max \left\{ (2a)^{\frac{1-\alpha}{1+\alpha}}, (2c)^\alpha \right\} + c,$$

which completes the proof. \square

Proof of Corollary 2. Consider Theorem 4 and observe that it verifies condition of Lemma 8 with $x = \mathbb{E}_{S,A} [R(A_S)]$, $c = \mathbb{E}_{S,A} [\widehat{R}_S(A_S)]$, $\alpha = \frac{c\gamma}{1+c\gamma}$, and

$$a = \frac{1 + \frac{1}{c\gamma}}{m} (2cL^2)^{\frac{1}{1+c\gamma}} T^{\frac{c\gamma}{1+c\gamma}}.$$

Note that $\alpha/(1-\alpha) = c\gamma$ and $1/(1-\alpha) = 1+c\gamma$. Then, we obtain that

$$\begin{aligned} & \mathbb{E}_{S,A} [R(A_S) - \widehat{R}_S(A_S)] \\ & \leq \max \left\{ 2^{c\gamma} \left(\frac{1 + \frac{1}{c\gamma}}{m} \right)^{1+c\gamma} (2cL^2)^{c\gamma} T^{c\gamma}, \right. \\ & \left. \left(2 \mathbb{E}_{S,A} [\widehat{R}_S(A_S)] \right)^{\frac{c\gamma}{1+c\gamma}} \left(\frac{1 + \frac{1}{c\gamma}}{m} (2cL^2)^{\frac{1}{1+c\gamma}} T^{\frac{c\gamma}{1+c\gamma}} \right) \right\} \\ & = \max \left\{ \left(2 + \frac{2}{c\gamma} \right)^{1+c\gamma} (cL^2)^{c\gamma} \left(\frac{T^{c\gamma}}{m^{1+c\gamma}} \right), \right. \\ & \left. \frac{1 + \frac{1}{c\gamma}}{m} (2cL^2)^{\frac{1}{1+c\gamma}} \left(2 \mathbb{E}_{S,A} [\widehat{R}_S(A_S)] \cdot T \right)^{\frac{c\gamma}{1+c\gamma}} \right\}. \end{aligned}$$

This completes the proof. \square

Proof of Proposition 1. Consider minimizing the bound given by Corollary 1 (in the submission file) over a discrete set of source hypotheses $\{\mathbf{w}_k^{\text{src}}\}_{k=1}^K$,

$$\begin{aligned} & \min_{k \in [K]} \epsilon(\mathcal{D}, \mathbf{w}_k^{\text{src}}) \\ & \leq \min_{k \in [K]} \mathcal{O} \left(\frac{1 + \frac{1}{c\gamma_k}}{m} (R(\mathbf{w}_k^{\text{src}}) \cdot T)^{\frac{c\gamma_k}{1+c\gamma_k}} \right), \quad (21) \end{aligned}$$

and let

$$\begin{aligned} \gamma_k & = \mathcal{O} \left(\mathbb{E}_{z \sim \mathcal{D}} [\|\nabla^2 f(\mathbf{w}_k^{\text{src}}, z)\|_2] + \sqrt{R(\mathbf{w}_k^{\text{src}})} \right), \\ \widehat{\gamma}_k & = \frac{1}{m} \sum_{i=1}^m \|\nabla^2 f(\mathbf{w}_k^{\text{src}}, z_i)\|_2 + \sqrt{\widehat{R}_S(\mathbf{w}_k^{\text{src}})}. \end{aligned}$$

By Hoeffding inequality, with high probability, we have that $|\gamma_k - \widehat{\gamma}_k| \leq \mathcal{O} \left(\frac{1}{\sqrt{m}} \right)$. Now we further upper bound (21)

by upper bounding $R(\mathbf{w}_k^{\text{src}})$ and apply union bound to get

$$\begin{aligned} & \min_{k \in [K]} \epsilon(\mathcal{D}, \mathbf{w}_k^{\text{src}}) \\ & \leq \min_{k \in [K]} \mathcal{O} \left(\left(1 + \frac{1}{c\widehat{\gamma}_k} \right) \widehat{R}_S(\mathbf{w}_k^{\text{src}})^{\frac{c\widehat{\gamma}_k^+}{1+c\widehat{\gamma}_k^+}} \cdot \frac{\sqrt{\log(K)}}{m^{\frac{1}{1+c\widehat{\gamma}_k^+}}} \right), \end{aligned}$$

where $\widehat{\gamma}_k^\pm = \widehat{\gamma}_k \pm \frac{1}{\sqrt{m}}$. This completes the proof. \square

References

- Boucheron, S., Lugosi, G., and Massart, P. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, 2013.
- Cucker, F. and Zhou, D. X. *Learning theory: an approximation theory viewpoint*, volume 24. Cambridge University Press, 2007.
- Ghadimi, S. and Lan, G. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- Hardt, M., Recht, B., and Singer, Y. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning (ICML)*, 2016.
- Orabona, F. Simultaneous model selection and optimization through parameter-free stochastic learning. In *Conference on Neural Information Processing Systems (NIPS)*, pp. 1116–1124, 2014.
- Shalev-Shwartz, S. and Ben-David, S. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.