
Sharper Generalization Bounds for Pairwise Learning: Supplementary Material

A Proof of Theorem 1

To prove Theorem 1, we need to introduce some lemmas. The following lemma is attributed to [7], which provides far-reaching moment bounds for a summation of weakly dependent and mean-zero random functions with bounded increments under a change of any single coordinate. We denote $S \setminus \{z_i\}$ the set $\{z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n\}$. The L_p -norm of a random variable Z is denoted by $\|Z\|_p := (\mathbb{E}[|Z|^p])^{\frac{1}{p}}, p \geq 1$.

Lemma A.1 ([4]). *Let $S = \{z_1, \dots, z_n\}$ be a set of independent random variables each taking values in \mathcal{Z} and $M > 0$. Let g_1, \dots, g_n be some functions $g_i : \mathcal{Z}^n \mapsto \mathbb{R}$ such that the following holds for any $i \in [n]$*

- $|\mathbb{E}_{S \setminus \{z_i\}}[g_i(S)]| \leq M$ almost surely (a.s.),
- $\mathbb{E}_{z_i}[g_i(S)] = 0$ a.s.,
- for any $j \in [n]$ with $j \neq i$, and $z_j'' \in \mathcal{Z}$

$$|g_i(S) - g_i(z_1, \dots, z_{j-1}, z_j'', z_{j+1}, \dots, z_n)| \leq \beta. \quad (\text{A.1})$$

Then, for any $p \geq 2$

$$\left\| \sum_{i=1}^n g_i(S) \right\|_p \leq 12\sqrt{6}pn\beta[\log_2 n] + 3\sqrt{2}M\sqrt{pn}.$$

The bounds on moments of random variables can be used to establish concentration inequalities, as shown in the following lemma [4, 16].

Lemma A.2. *Let $a, b \in \mathbb{R}_+$ and $\delta \in (0, 1/e)$. Let Z be a random variable with $\|Z\|_p \leq \sqrt{pa} + pb$ for any $p \geq 2$. Then with probability at least $1 - \delta$*

$$|Z| \leq e \left(a\sqrt{\log(e/\delta)} + b \log(e/\delta) \right).$$

The following lemma controls the change on the output of stable algorithms if we perturb a training dataset by two examples.

Lemma A.3. *Let $A : \mathcal{Z}^n \mapsto \mathcal{W}$ be γ -uniformly stable. Then for any $S' = \{z'_1, \dots, z'_n\}$ and $i \neq j$, we have*

$$\sup_{z, \tilde{z} \in \mathcal{Z}} |\ell(A(S); z, \tilde{z}) - \ell(A(S_{i,j}); z, \tilde{z})| \leq 2\gamma,$$

where $S_{i,j}$ is defined in (3.4).

Proof. For any $i \in [n]$, introduce

$$S_i = \{z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n\}. \quad (\text{A.2})$$

Note that S, S_i differ only by a single example, and $S_i, S_{i,j}$ differ only by a single example. It then follows from the definition of uniform stability that

$$\begin{aligned} & \sup_{z, \tilde{z} \in \mathcal{Z}} |\ell(A(S); z, \tilde{z}) - \ell(A(S_{i,j}); z, \tilde{z})| \\ & \leq \sup_{z, \tilde{z} \in \mathcal{Z}} |\ell(A(S); z, \tilde{z}) - \ell(A(S_i); z, \tilde{z})| + \sup_{z, \tilde{z} \in \mathcal{Z}} |\ell(A(S_i); z, \tilde{z}) - \ell(A(S_{i,j}); z, \tilde{z})| \\ & \leq 2\gamma. \end{aligned}$$

The proof is complete. \square

With these lemmas, we can give the proof of Theorem 1 on high-probability bounds of the generalization gap. The concentration inequality established in Lemma A.1 applies to a summation of n random functions involving n independent random variables, which does not apply to the objective function in pairwise learning since it is a U -statistic. We introduce a novel decomposition to exploit the structure of pairwise learning problems. We abbreviate $\sum_{i,j \in [n]: i \neq j}$ as $\sum_{i \neq j}$.

Proof of Theorem 1. Let $p \geq 2$ be any number. We can decompose the generalization gap associated to $A(S)$ as follows

$$\begin{aligned} & n(n-1)\mathbb{E}_{Z, \tilde{Z}} [\ell(A(S); Z, \tilde{Z})] - \sum_{i \neq j} \ell(A(S); z_i, z_j) = \sum_{i \neq j} \mathbb{E}_{Z, \tilde{Z}} [\ell(A(S); Z, \tilde{Z}) - \mathbb{E}_{z'_i, z'_j} [\ell(A(S_{i,j}); Z, \tilde{Z})]] \\ & + \sum_{i \neq j} \mathbb{E}_{z'_i, z'_j} [\mathbb{E}_{Z, \tilde{Z}} [\ell(A(S_{i,j}); Z, \tilde{Z})] - \ell(A(S_{i,j}); z_i, z_j)] + \sum_{i \neq j} \mathbb{E}_{z'_i, z'_j} [\ell(A(S_{i,j}); z_i, z_j) - \ell(A(S); z_i, z_j)], \end{aligned}$$

where $S_{i,j}$ is defined in (3.4). According to Lemma A.3, we know

$$\left| \ell(A(S); Z, \tilde{Z}) - \mathbb{E}_{z'_i, z'_j} [\ell(A(S_{i,j}); Z, \tilde{Z})] \right| \leq 2\gamma$$

and

$$\left| \ell(A(S_{i,j}); z_i, z_j) - \ell(A(S); z_i, z_j) \right| \leq 2\gamma.$$

Therefore,

$$\left| n(n-1)\mathbb{E}_{Z, \tilde{Z}} [\ell(A(S); Z, \tilde{Z})] - \sum_{i \neq j} \ell(A(S); z_i, z_j) \right| \leq 4n(n-1)\gamma + \left| \sum_{i \neq j} g_{i,j}(S) \right|, \quad (\text{A.3})$$

where we introduce

$$g_{i,j}(S) = \mathbb{E}_{z'_i, z'_j} [\mathbb{E}_{Z, \tilde{Z}} [\ell(A(S_{i,j}); Z, \tilde{Z})] - \ell(A(S_{i,j}); z_i, z_j)], \quad \forall i, j \in [n].$$

For any $i, j \in [n]$, we can further decompose $g_{i,j}$ as $g_{i,j} = g_j^{(i)} + \tilde{g}_i^{(j)}$, where (we omit the argument S for brevity)

$$\begin{aligned} g_j^{(i)} &= \mathbb{E}_{z'_i, z'_j} [\mathbb{E}_{Z, \tilde{Z}} [\ell(A(S_{i,j}); Z, \tilde{Z})] - \mathbb{E}_Z [\ell(A(S_{i,j}); Z, z_j)]] \\ \tilde{g}_i^{(j)} &= \mathbb{E}_{z'_i, z'_j} [\mathbb{E}_Z [\ell(A(S_{i,j}); Z, z_j) - \ell(A(S_{i,j}); z_i, z_j)]]. \end{aligned}$$

Let us temporarily fix i , and consider $n-1$ random functions $g_1^{(i)}, \dots, g_{i-1}^{(i)}, g_{i+1}^{(i)}, \dots, g_n^{(i)}$. According to the assumption $|\mathbb{E}_S [\ell(A(S); z, \tilde{z})]| \leq M$ for all z, \tilde{z} , we know

$$|\mathbb{E}_{S \setminus \{z_j\}} [g_j^{(i)}(S)]| \leq 2M, \quad \forall j \in [n].$$

For any $j \neq i$, since z_j is independent of $S_{i,j}$ we know

$$\mathbb{E}_{z_j} [\mathbb{E}_{Z, \tilde{Z}} [\ell(A(S_{i,j}); Z, \tilde{Z})] - \mathbb{E}_Z [\ell(A(S_{i,j}); Z, z_j)]] = 0.$$

Therefore, $\mathbb{E}_{z_j} [g_j^{(i)}] = 0$. For any $k \neq j$ and any $z'_k \in \mathcal{Z}$, it is clear from the uniform stability of A that

$$\left| \mathbb{E}_{z'_i, z'_j} \mathbb{E}_{Z, \tilde{Z}} [\ell(A(S_{i,j}); Z, \tilde{Z})] - \mathbb{E}_{z'_i, z'_j} \mathbb{E}_{Z, \tilde{Z}} [\ell(A(S_{i,j}^{(k)}); Z, \tilde{Z})] \right| \leq \gamma,$$

where $S_{i,j}^{(k)}$ is the set derived by replacing the k -th element of $S_{i,j}$ with z_k'' . Similarly, one have

$$\left| \mathbb{E}_{z_i', z_j'} \mathbb{E}_Z [\ell(A(S_{i,j}); Z, z_j)] - \mathbb{E}_{z_i', z_j'} \mathbb{E}_Z [\ell(A(S_{i,j}^{(k)}); Z, z_j)] \right| \leq \gamma.$$

It then follows from the above two inequalities that $g_j^{(i)}$ satisfies the bounded increment condition (A.1) with $\beta = 2\gamma$ for all $k \neq j$, i.e.,

$$\left| \mathbb{E}_{z_i', z_j'} \left[\mathbb{E}_{Z, \tilde{Z}} [\ell(A(S_{i,j}); Z, \tilde{Z})] - \mathbb{E}_Z [\ell(A(S_{i,j}); Z, z_j)] \right] - \mathbb{E}_{z_i', z_j'} \left[\mathbb{E}_{Z, \tilde{Z}} [\ell(A(S_{i,j}^{(k)}); Z, \tilde{Z})] - \mathbb{E}_Z [\ell(A(S_{i,j}^{(k)}); Z, z_j)] \right] \right| \leq 2\gamma.$$

Therefore, all the assumptions of Lemma A.1 hold for the random functions $g_1^{(i)}, \dots, g_{i-1}^{(i)}, g_{i+1}^{(i)}, \dots, g_n^{(i)}$ with n there replaced by $n-1$ and $\beta = 2\gamma$. We can apply Lemma A.1 to derive

$$\left\| \sum_{j \in [n], j \neq i} g_j^{(i)} \right\|_p \leq 24\sqrt{6}p(n-1)\gamma \lceil \log_2(n-1) \rceil + 6\sqrt{2}M\sqrt{p(n-1)}, \quad \forall i \in [n].$$

Similarly, we can also show that

$$\left\| \sum_{i \in [n], i \neq j} \tilde{g}_i^{(j)} \right\|_p \leq 24\sqrt{6}p(n-1)\gamma \lceil \log_2(n-1) \rceil + 6\sqrt{2}M\sqrt{p(n-1)}, \quad \forall j \in [n].$$

It then follows from the subadditivity of $\|\cdot\|_p$ and the above two inequalities that

$$\begin{aligned} \left\| \sum_{i \neq j} g_{i,j} \right\|_p &\leq \left\| \sum_{i \neq j} g_j^{(i)} \right\|_p + \left\| \sum_{i \neq j} \tilde{g}_i^{(j)} \right\|_p \\ &\leq \sum_{i \in [n]} \left\| \sum_{j \in [n], j \neq i} g_j^{(i)} \right\|_p + \sum_{j \in [n]} \left\| \sum_{i \in [n], i \neq j} \tilde{g}_i^{(j)} \right\|_p \\ &\leq 48\sqrt{6}p(n-1)n\gamma \lceil \log_2(n-1) \rceil + 12\sqrt{2}M\sqrt{p(n-1)}n. \end{aligned}$$

We can combine the above p -norm and Lemma A.2 to derive the following inequality with probability at least $1 - \delta$

$$\left| \sum_{i \neq j} g_{i,j} \right| \leq e \left(12\sqrt{2}M\sqrt{(n-1)n\sqrt{\log(e/\delta)}} + 48\sqrt{6}(n-1)n\gamma \lceil \log_2(n-1) \rceil \log(e/\delta) \right).$$

Plugging the above inequality back into (A.3) and using the definition of R_S, R , we derive the following inequality with probability at least $1 - \delta$

$$\begin{aligned} |R_S(A(S)) - R(A(S))| &\leq 4\gamma + \frac{1}{n(n-1)} \left| \sum_{i \neq j} g_{i,j} \right| \\ &\leq 4\gamma + e \left(12\sqrt{2}M(n-1)^{-\frac{1}{2}} \sqrt{\log(e/\delta)} + 48\sqrt{6}\gamma \lceil \log_2(n-1) \rceil \log(e/\delta) \right). \end{aligned}$$

The proof is complete. \square

B Proof of Theorem 3

In this section, we prove Theorem 3 on high-probability bounds for learning with strongly convex objective functions. We first prove Lemma 2 on the norm of output model.

Proof of Lemma 2. Since $A(S)$ is the minimizer of F_S , we know there is a $F_S'(A(S)) = 0$ (F_S' is a subgradient of F_S at $A(S)$). This together with the definition of strong convexity implies

$$R_S(\mathbf{w}^*) + r(\mathbf{w}^*) - R_S(A(S)) - r(A(S)) \geq \frac{\sigma}{2} \|A(S) - \mathbf{w}^*\|^2. \quad (\text{B.1})$$

Analogous to (A.3), we know

$$n(n-1)\left(R(A(S)) - R_S(A(S))\right) \leq 4n(n-1)\gamma + \sum_{i,j \in [n]: i \neq j} g_{i,j},$$

where $g_{i,j}$ is defined in the proof of Theorem 1. In the proof of Theorem 1, we have shown $\mathbb{E}[g_{i,j}] = 0$. It then follows that

$$\mathbb{E}[R(A(S)) - R_S(A(S))] \leq 4\gamma.$$

We can plug the above inequality back into (B.1) to derive

$$\begin{aligned} \frac{\sigma}{2} \mathbb{E}[\|A(S) - \mathbf{w}^*\|^2] &\leq \mathbb{E}[R_S(\mathbf{w}^*) + r(\mathbf{w}^*) - R_S(A(S)) - r(A(S))] \\ &\leq \mathbb{E}[R_S(\mathbf{w}^*) + r(\mathbf{w}^*) - R(A(S)) - r(A(S))] + 4\gamma \\ &= \mathbb{E}[R(\mathbf{w}^*) + r(\mathbf{w}^*) - R(A(S)) - r(A(S))] + 4\gamma \leq 4\gamma, \end{aligned}$$

where the last inequality holds since \mathbf{w}^* minimizes $F = R + r$. The stated inequality then follows and finishes the proof. \square

To prove Theorem 3, we introduce some lemmas.

Lemma B.1. *For any $S \in \mathcal{Z}^n$, define A as $A(S) = \arg \min_{\mathbf{w} \in \mathcal{W}} F_S(\mathbf{w})$. For any $k \in [n]$, let S_k be defined by (A.2). Then*

$$\begin{aligned} F_S(A(S_k)) - F_S(A(S)) &\leq \\ \frac{1}{n(n-1)} \sum_{i \in [n]: i \neq k} &\left((\ell(A(S_k); z_i, z_k) - \ell(A(S); z_i, z_k)) + (\ell(A(S_k); z_k, z_i) - \ell(A(S); z_k, z_i)) \right. \\ &\left. + (\ell(A(S); z_i, z'_k) - \ell(A(S_k); z_i, z'_k)) + (\ell(A(S); z'_k, z_i) - \ell(A(S_k); z'_k, z_i)) \right). \end{aligned}$$

Proof. Without loss of generality, we can assume $k = n$. Since $A(S_n)$ is a minimizer of F_{S_n} , we know

$$\begin{aligned} F_S(A(S_n)) - F_S(A(S)) &= F_S(A(S_n)) - F_{S_n}(A(S_n)) + F_{S_n}(A(S_n)) - F_{S_n}(A(S)) + F_{S_n}(A(S)) - F_S(A(S)) \\ &\leq F_S(A(S_n)) - F_{S_n}(A(S_n)) + F_{S_n}(A(S)) - F_S(A(S)). \end{aligned} \tag{B.2}$$

By the definition of F_S and F_{S_n} , we know

$$\begin{aligned} n(n-1)(F_S(A(S_n)) - F_{S_n}(A(S_n))) &= \sum_{i,j \in [n]: i \neq j} f(A(S_n); z_i, z_j) \\ &- \left(\sum_{i,j \in [n-1]: i \neq j} f(A(S_n); z_i, z_j) + \sum_{i \in [n-1]} f(A(S_n); z_i, z'_n) + \sum_{i \in [n-1]} f(A(S_n); z'_n, z_i) \right) \\ &= \sum_{i \in [n-1]} \left(f(A(S_n); z_i, z_n) + f(A(S_n); z_n, z_i) - f(A(S_n); z_i, z'_n) - f(A(S_n); z'_n, z_i) \right). \end{aligned}$$

Similarly, we know

$$\begin{aligned} n(n-1)(F_{S_n}(A(S)) - F_S(A(S))) &= \\ \sum_{i \in [n-1]} &\left(f(A(S); z_i, z'_n) + f(A(S); z'_n, z_i) - f(A(S); z_i, z_n) - f(A(S); z_n, z_i) \right). \end{aligned}$$

Therefore, we can combine the above two identities to derive

$$\begin{aligned} n(n-1)(F_S(A(S_n)) - F_{S_n}(A(S_n)) + F_{S_n}(A(S)) - F_S(A(S))) &= \\ \sum_{i \in [n-1]} &\left((f(A(S_n); z_i, z_n) - f(A(S); z_i, z_n)) + (f(A(S_n); z_n, z_i) - f(A(S); z_n, z_i)) + \right. \\ &\left. (f(A(S); z_i, z'_n) - f(A(S_n); z_i, z'_n)) + (f(A(S); z'_n, z_i) - f(A(S_n); z'_n, z_i)) \right). \end{aligned}$$

This together with the structure of f ($f = \ell + r$ with r depending only on \mathbf{w}) implies

$$\begin{aligned} n(n-1)(F_S(A(S_n)) - F_{S_n}(A(S_n)) + F_{S_n}(A(S)) - F_S(A(S))) = \\ \sum_{i \in [n-1]} \left((\ell(A(S_n); z_i, z_n) - \ell(A(S); z_i, z_n)) + (\ell(A(S_n); z_n, z_i) - \ell(A(S); z_n, z_i)) + \right. \\ \left. (\ell(A(S); z_i, z'_n) - \ell(A(S_n); z_i, z'_n)) + (\ell(A(S); z'_n, z_i) - \ell(A(S_n); z'_n, z_i)) \right). \end{aligned}$$

Plugging the above identity back into (B.2), we derive

$$\begin{aligned} F_S(A(S_n)) - F_S(A(S)) \\ \leq \frac{1}{n(n-1)} \sum_{i \in [n-1]} \left((\ell(A(S_n); z_i, z_n) - \ell(A(S); z_i, z_n)) + (\ell(A(S_n); z_n, z_i) - \ell(A(S); z_n, z_i)) \right. \\ \left. + (\ell(A(S); z_i, z'_n) - \ell(A(S_n); z_i, z'_n)) + (\ell(A(S); z'_n, z_i) - \ell(A(S_n); z'_n, z_i)) \right). \end{aligned}$$

The proof is complete. \square

The following lemma establishes the uniform stability of pairwise learning with strongly convex objectives.

Lemma B.2. *Define A as $A(S) = \arg \min_{\mathbf{w} \in \mathcal{W}} F_S(\mathbf{w})$. Suppose F_S is σ -strongly convex w.r.t. $\|\cdot\|$. Assume for all z, \tilde{z} we have (4.3). Then A is $\frac{8L^2}{n\sigma}$ -uniformly stable.*

Proof. Let S, S' be two sets that differ by a single example and let $\mathbf{w}_S = A(S)$ and $\mathbf{w}_{S'} = A(S')$. Without loss of generality, we can assume $S' = \{z_1, \dots, z_{n-1}, z'_n\}$, i.e., S and S' differ by the last example.

Since \mathbf{w}_S is a minimizer of F_S we know there is a subgradient $F'_S(\mathbf{w}_S) = 0$, which together with the σ -strong convexity of F_S , implies

$$F_S(\mathbf{w}_{S'}) - F_S(\mathbf{w}_S) \geq \frac{\sigma}{2} \|\mathbf{w}_{S'} - \mathbf{w}_S\|^2. \quad (\text{B.3})$$

According to (4.3) and Lemma B.1, we know

$$F_S(\mathbf{w}_{S'}) - F_S(\mathbf{w}_S) \leq \frac{4(n-1)L\|\mathbf{w}_S - \mathbf{w}_{S'}\|}{n(n-1)}.$$

which, together with (B.3), implies

$$\|\mathbf{w}_S - \mathbf{w}_{S'}\| \leq \frac{8L}{n\sigma}.$$

This further together with (4.3) implies the $\frac{8L^2}{n\sigma}$ -uniform stability of A . The proof is complete. \square

To obtain tight control on the term $R(\mathbf{w}^*) - R_S(\mathbf{w}^*)$, we will need a version of Bernstein's inequality for U-statistics. The following theorem is attributed to [10], and can be found in [5] (inequality A.1 on page 868), and in [12] (Theorem 2). A complete proof is provided in [13] (page 4).

Lemma B.3 (Bernstein's inequality for U-Statistic [10, 13]). *Let Z_1, \dots, Z_n be independent variables taking values in \mathcal{Z} and $q: \mathcal{Z} \times \mathcal{Z} \mapsto \mathbb{R}$. Let $b = \sup_{z, \tilde{z}} |q(z, \tilde{z})|$ and σ_0^2 be the variance of $q(Z, \tilde{Z})$. Then for any $\delta \in (0, 1)$ with probability at least $1 - \delta$*

$$\left| \frac{1}{n(n-1)} \sum_{i, j \in [n]: i \neq j} q(Z_i, Z_j) - \mathbb{E}_{Z, \tilde{Z}}[q(Z, \tilde{Z})] \right| \leq \frac{2b \log(1/\delta)}{3 \lfloor n/2 \rfloor} + \sqrt{\frac{2\sigma_0^2 \log(1/\delta)}{\lfloor n/2 \rfloor}}. \quad (\text{B.4})$$

We now give the proof of Theorem 3.

Proof of Theorem 3. According to Lemma B.2, we know that A is $\frac{8L^2}{n\sigma}$ -uniformly stable. Using this together with Lemma 2 we derive $\mathbb{E}_S[\|\mathbf{w}^* - A(S)\|^2] \leq \frac{64L^2}{n\sigma^2}$ and therefore

$$\mathbb{E}_S[\|\mathbf{w}^* - A(S)\|] \leq \left(\mathbb{E}_S[\|\mathbf{w}^* - A(S)\|^2]\right)^{\frac{1}{2}} \leq \frac{8L}{\sqrt{n\sigma}}. \quad (\text{B.5})$$

For any $\mathbf{w} \in \mathcal{W}$ and z, \tilde{z} , define

$$\tilde{\ell}(\mathbf{w}; z, \tilde{z}) = \ell(\mathbf{w}; z, \tilde{z}) - \ell(\mathbf{w}^*; z, \tilde{z}).$$

Then it is clear from Lemma B.2 that A is also $\frac{8L^2}{n\sigma}$ -uniformly stable when measured by the ‘‘loss’’ $\tilde{\ell}$, i.e., for any S, S' differing by one example

$$\begin{aligned} & \sup_{z, \tilde{z}} |\tilde{\ell}(A(S); z, \tilde{z}) - \tilde{\ell}(A(S'); z, \tilde{z})| \\ &= \sup_{z, \tilde{z}} |\ell(A(S); z, \tilde{z}) - \ell(A(S'); z, \tilde{z}) - \ell(\mathbf{w}^*; z, \tilde{z}) + \ell(\mathbf{w}^*; z, \tilde{z})| \\ &= \sup_{z, \tilde{z}} |\ell(A(S); z, \tilde{z}) - \ell(A(S'); z, \tilde{z})| \leq 8L^2/(n\sigma). \end{aligned}$$

Furthermore, by the Lipschitz continuity (4.3) and (B.5), we know the following inequality for all $z, \tilde{z} \in \mathcal{Z}$

$$\begin{aligned} \left| \mathbb{E}_S[\tilde{\ell}(A(S); z, \tilde{z})] \right| &= \left| \mathbb{E}_S[\ell(A(S); z, \tilde{z}) - \ell(\mathbf{w}^*; z, \tilde{z})] \right| \\ &\leq L \mathbb{E}_S[\|\mathbf{w}^* - A(S)\|] \leq \frac{8L^2}{\sqrt{n\sigma}}. \end{aligned}$$

We can now apply Theorem 1, with $\gamma = \frac{8L^2}{n\sigma}$, $M = 8L^2/(\sqrt{n\sigma})$ and ℓ replaced by $\tilde{\ell}$, and show the following inequality with probability $1 - \delta/2$

$$\begin{aligned} \left| \frac{1}{n(n-1)} \sum_{i \neq j} \tilde{\ell}(A(S); z_i, z_j) - \mathbb{E}_{z, \tilde{z}}[\tilde{\ell}(A(S); z, \tilde{z})] \right| &\leq \frac{32L^2}{n\sigma} \\ &+ e \left(\frac{96\sqrt{2}L^2\sqrt{\log(2e/\delta)}}{\sqrt{n(n-1)}\sigma} \right) + \frac{384\sqrt{6}L^2[\log_2 n] \log(2e/\delta)}{n\sigma}, \end{aligned}$$

from which we derive the following inequality with probability $1 - \delta/2$

$$\begin{aligned} |R_S(A(S)) - R(A(S))| &\leq \left| \frac{1}{n(n-1)} \sum_{i \neq j} \ell(\mathbf{w}^*; z_i, z_j) - \mathbb{E}_{z, \tilde{z}}[\ell(\mathbf{w}^*; z, \tilde{z})] \right| \\ &+ \frac{32L^2}{n\sigma} \left(1 + 3\sqrt{\frac{2n \log(2e/\delta)}{n-1}} + 12\sqrt{6}[\log_2 n] \log(2e/\delta) \right). \quad (\text{B.6}) \end{aligned}$$

By the definition of \mathbf{w}^* ($R'(\mathbf{w}^*) + r'(\mathbf{w}^*) = 0$), the σ -strong convexity and Assumption 1 ($0 \leq \ell(0; z, \tilde{z})$), we know

$$\frac{\sigma\|\mathbf{w}^*\|^2}{2} \leq R(0) + r(0) - R(\mathbf{w}^*) - r(\mathbf{w}^*) \implies \|\mathbf{w}^*\| \leq \sqrt{\frac{2(R(0) + r(0))}{\sigma}}.$$

It then follows from the Lipschitz continuity (4.3) that

$$\begin{aligned} |\ell(\mathbf{w}^*; z, \tilde{z})| &= |\ell(\mathbf{w}^*; z, \tilde{z}) - \ell(0; z, \tilde{z}) + \ell(0; z, \tilde{z})| \leq L\|\mathbf{w}^*\| + \sup_{z, \tilde{z}} \ell(0; z, \tilde{z}) \\ &\leq L\sqrt{\frac{2(R(0) + r(0))}{\sigma}} + \sup_{z, \tilde{z}} \ell(0; z, \tilde{z}). \end{aligned}$$

According to Bernstein’s inequality (B.4), we derive the following inequality with probability $1 - \delta/2$ that

$$\begin{aligned} \left| \frac{1}{n(n-1)} \sum_{i \neq j} \ell(\mathbf{w}^*; z_i, z_j) - \mathbb{E}_{z, \tilde{z}}[\ell(\mathbf{w}^*; z, \tilde{z})] \right| &\leq \\ &\frac{2(L\sqrt{2(R(0) + r(0))/\sigma} + b) \log(2/\delta)}{3\lfloor n/2 \rfloor} + \sqrt{\frac{2\sigma_0^2 \log(2/\delta)}{\lfloor n/2 \rfloor}}. \quad (\text{B.7}) \end{aligned}$$

Plugging the above inequality back into (B.6), we derive the following inequality with probability at least $1 - \delta$

$$\begin{aligned} |R_S(A(S)) - R(A(S))| &\leq \frac{2(L\sqrt{2(R(0) + r(0))/\sigma + b}) \log(2/\delta)}{3\lfloor n/2 \rfloor} + \sqrt{\frac{2\sigma_0^2 \log(2/\delta)}{\lfloor n/2 \rfloor}} \\ &\quad + \frac{32L^2}{n\sigma} \left(1 + 3\sqrt{\frac{2n \log(2e/\delta)}{n-1}} + 12\sqrt{6} \lceil \log_2 n \rceil \log(2e/\delta)\right). \end{aligned}$$

The above inequality can be written as the stated bound (4.4).

We now turn to (4.5). According to the definition of R and F , we can decompose the excess risk $R(A(S)) - R(\mathbf{w}_R^*)$ as follows

$$\begin{aligned} &R(A(S)) - R(\mathbf{w}_R^*) \\ &= R(A(S)) - R_S(A(S)) + R_S(\mathbf{w}_R^*) - R(\mathbf{w}_R^*) + R_S(A(S)) - R_S(\mathbf{w}_R^*) \\ &= R(A(S)) - R_S(A(S)) + R_S(\mathbf{w}_R^*) - R(\mathbf{w}_R^*) + F_S(A(S)) - F_S(\mathbf{w}_R^*) + r(\mathbf{w}_R^*) - r(A(S)) \\ &\leq R(A(S)) - R_S(A(S)) + R_S(\mathbf{w}_R^*) - R(\mathbf{w}_R^*) + O(\sigma \|\mathbf{w}_R^*\|^2) - r(A(S)), \end{aligned} \quad (\text{B.8})$$

where we have used the inequality $F_S(A(S)) \leq F_S(\mathbf{w}_R^*)$ due to the definition of $A(S)$ and the assumption $r(\mathbf{w}) = O(\sigma \|\mathbf{w}\|^2)$ in the last step. Analogous to (B.7), one can use Bernstein's inequality (Lemma B.3) to show with probability at least $1 - \delta/2$ that (under a very mild assumption $\sup_{z, z'} \ell(\mathbf{w}_R^*; z, z') = O(\sqrt{n})$)

$$R_S(\mathbf{w}_R^*) - R(\mathbf{w}_R^*) = O\left(\frac{\log(1/\delta)}{\sqrt{n}} + \sqrt{\frac{\sigma_0^2 \log(1/\delta)}{n}}\right). \quad (\text{B.9})$$

Plugging the above inequality and (4.4) back into (B.8) shows the following inequality with probability at least $1 - \delta$

$$R(A(S)) - R(\mathbf{w}_R^*) = O\left((n\sigma)^{-1} \log n \log(1/\delta) + n^{-\frac{1}{2}} \log(1/\delta)\right) + O(\sigma \|\mathbf{w}_R^*\|^2).$$

The stated bound (4.5) follows with $\sigma \asymp n^{-1/2}$. The proof is complete. \square

Remark B.1. We show here that the existing stability bound (eq. (4.2) with $\gamma = O(1/(n\sigma))$) [1, 6, 11, 17]

$$|R_S(A(S)) - R(A(S))| = O(\sigma^{-1} n^{-\frac{1}{2}}) \quad (\text{B.10})$$

yields at best the excess risk bound $R(A(S)) - R(\mathbf{w}_R^*) = O(n^{-\frac{1}{4}})$. Indeed, plugging (B.10) and (B.9) back into (B.8), we derive the following inequality with high probability

$$R(A(S)) - R(\mathbf{w}_R^*) = O(\sigma^{-1} n^{-\frac{1}{2}}) + O(\sigma).$$

We can balance the above two terms by taking $\sigma \asymp n^{-\frac{1}{4}}$ and get

$$R(A(S)) - R(\mathbf{w}_R^*) = O(n^{-\frac{1}{4}}).$$

C Proof of Theorem 4

To prove Theorem 4, we first introduce some lemmas. Lemma C.1 shows the non-expansiveness of the gradient-update operator, which plays a key role in establishing the stability of SGD. Lemma C.2 is a Chernoff's bound for a summation of independent Bernoulli random variables [2]. In this section, we let $\|\cdot\|_2$ be the Euclidean norm.

Lemma C.1 ([8]). *Assume for all $z \in \mathcal{Z}$, the function $\mathbf{w} \mapsto \ell(\mathbf{w}; z, z')$ is convex and α -smooth. Then for all $\eta \leq 2/\alpha$ and $z, z' \in \mathcal{Z}$ there holds*

$$\|\mathbf{w} - \eta \ell'(\mathbf{w}; z, z') - \mathbf{w}' + \eta \ell'(\mathbf{w}'; z, z')\|_2 \leq \|\mathbf{w} - \mathbf{w}'\|_2.$$

Lemma C.2 (Chernoff's Bound). *Let X_1, \dots, X_T be independent random variables taking values in $\{0, 1\}$. Let $X = \sum_{t=1}^T X_t$ and $\mu = \mathbb{E}[X]$. Then for any $\epsilon \in (0, 1)$ with probability at least $1 - \exp(-\mu\epsilon^2/3)$ we have $X \leq (1 + \epsilon)\mu$.*

We now establish the uniform stability of SGD. The randomness of SGD can be characterized by $\{\{(i_t, j_t)_t\} : i_t, j_t \in [n], i_t \neq j_t\}$. Therefore, SGD can be considered as a deterministic algorithm if $\{\{(i_t, j_t)_t\} : i_t, j_t \in [n], i_t \neq j_t\}$ is fixed. For simplicity, we consider two datasets that differ by the last example. However, our discussion directly extends to the general case where two datasets differ by a single example. Notice that the Lipschitz continuity (4.3) implies $\|\ell'(\mathbf{w}; z, z')\|_2 \leq L$.

Lemma C.3. *Consider fixed $\{\{(i_t, j_t)_t\} : i_t, j_t \in [n], i_t \neq j_t\}$. Let $S = \{z_1, \dots, z_n\}$ and $S' = \{z'_1, \dots, z'_n\}$ be two datasets that differ only by the last example, i.e., $z_i = z'_i$ if $i \in [n-1]$. Suppose for all $z, z' \in \mathcal{Z}$ the function $\mathbf{w} \mapsto \ell(\mathbf{w}; z, z')$ is convex, α -smooth and L -Lipschitz w.r.t. $\|\cdot\|_2$. Let $\{\mathbf{w}_t\}, \{\mathbf{w}'_t\}$ be produced by SGD on S and S' respectively with $\eta_t \leq 2/\alpha$, i.e., (3.3) with $r(\mathbf{w}) = 0$. Then SGD with t iterations is γ -uniformly stable with*

$$\gamma \leq 2L^2 \sum_{k=1}^t \eta_k \mathbb{I}[i_k = n \text{ or } j_k = n].$$

Proof. Let us consider two cases. We first consider the case $i_t \in [n-1]$ and $j_t \in [n-1]$. In this case, according to the SGD update (3.3) with $r(\mathbf{w}) = 0$ we know

$$\begin{aligned} \mathbf{w}_{t+1} - \mathbf{w}'_{t+1} &= \mathbf{w}_t - \eta_t \ell'(\mathbf{w}_t; z_{i_t}, z_{j_t}) - \mathbf{w}'_t + \eta_t \ell'(\mathbf{w}'_t; z'_{i_t}, z'_{j_t}) \\ &= \mathbf{w}_t - \eta_t \ell'(\mathbf{w}_t; z_{i_t}, z_{j_t}) - \mathbf{w}'_t + \eta_t \ell'(\mathbf{w}'_t; z_{i_t}, z_{j_t}). \end{aligned}$$

It then follows from Lemma C.1 that

$$\|\mathbf{w}_{t+1} - \mathbf{w}'_{t+1}\|_2 \leq \|\mathbf{w}_t - \mathbf{w}'_t\|_2.$$

We now consider the case that either $i_t = n$ or $j_t = n$. In this case, we know

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{w}'_{t+1}\|_2 &= \|\mathbf{w}_t - \eta_t \ell'(\mathbf{w}_t; z_{i_t}, z_{j_t}) - \mathbf{w}'_t + \eta_t \ell'(\mathbf{w}'_t; z'_{i_t}, z'_{j_t})\|_2 \\ &\leq \|\mathbf{w}_t - \mathbf{w}'_t\|_2 + \|\eta_t \ell'(\mathbf{w}'_t; z'_{i_t}, z'_{j_t}) - \eta_t \ell'(\mathbf{w}_t; z_{i_t}, z_{j_t})\|_2 \\ &\leq \|\mathbf{w}_t - \mathbf{w}'_t\|_2 + 2\eta_t L, \end{aligned}$$

where we have used $\|\ell'(\mathbf{w}; z, z')\|_2 \leq L$ due to the L -Lipschitzness. As a combination of the above two cases, we derive

$$\|\mathbf{w}_{t+1} - \mathbf{w}'_{t+1}\|_2 \leq \|\mathbf{w}_t - \mathbf{w}'_t\|_2 + 2\eta_t L \mathbb{I}[i_t = n \text{ or } j_t = n],$$

where $\mathbb{I}[\cdot]$ is the indicator function taking 1 if the argument holds and 0 otherwise. Taking a summation of the above inequality gives $(\mathbf{w}_1 = \mathbf{w}'_1)$

$$\|\mathbf{w}_{t+1} - \mathbf{w}'_{t+1}\|_2 \leq 2L \sum_{k=1}^t \eta_k \mathbb{I}[i_k = n \text{ or } j_k = n].$$

This together with the Lipschitz continuity of ℓ implies the following inequality for all $z, z' \in \mathcal{Z}$

$$\begin{aligned} |\ell(\mathbf{w}_{t+1}; z, z') - \ell(\mathbf{w}'_{t+1}; z, z')| &\leq L \|\mathbf{w}_{t+1} - \mathbf{w}'_{t+1}\|_2 \\ &\leq 2L^2 \sum_{k=1}^t \eta_k \mathbb{I}[i_k = n \text{ or } j_k = n]. \end{aligned}$$

The proof is complete. \square

We now apply the above uniform stability bounds and Theorem 1 to prove Theorem 4.

Proof of Theorem 4. We can apply Theorem 1 with $A(S) = \mathbf{w}_T$ and the uniform stability bounds in Lemma C.3 to show with probability at least $1 - \delta/2$ that

$$|R_S(\mathbf{w}_T) - R(\mathbf{w}_T)| = O\left((\log n \log(1/\delta)) \sum_{t=1}^T \eta \mathbb{I}[i_t = n \text{ or } j_t = n]\right) + O(n^{-\frac{1}{2}} \sqrt{\log(1/\delta)}), \quad (\text{C.1})$$

where $\eta = c/\sqrt{T}$. Let $X_t = \mathbb{I}[i_t = n \text{ or } j_t = n]$. It is clear that

$$\mathbb{E}[X_t] = \Pr\{i_t = n \text{ or } j_t = n\} \leq \Pr\{i_t = n\} + \Pr\{j_t = n\} = 2/n.$$

Applying Lemma C.2 with $X_t = \mathbb{I}[i_t = n \text{ or } j_t = n]$ then gives with probability $1 - \delta/2$ that

$$\sum_{t=1}^T X_t \leq \left(1 + \sqrt{3\mu^{-1} \log(1/\delta)}\right) \mu,$$

where $\mu = \sum_{t=1}^T \mathbb{E}[X_t] \leq 2T/n$. It then follows with probability $1 - \delta/2$ that

$$\sum_{t=1}^T X_t \leq \frac{2T}{n} \left(1 + \sqrt{2nT^{-1} \log(1/\delta)}\right). \quad (\text{C.2})$$

Combining (C.1) and (C.2) together, we derive the following inequality with probability $1 - \delta$

$$\begin{aligned} |R_S(\mathbf{w}_T) - R(\mathbf{w}_T)| &= O(n^{-\frac{1}{2}} \sqrt{\log(1/\delta)} + T\eta(\log n \log(1/\delta))n^{-1} \\ &\quad + \eta \log n \log(1/\delta) \sqrt{n^{-1}T \log(1/\delta)}). \end{aligned}$$

The proof is complete with $\eta = c/\sqrt{T}$. \square

Remark C.1. We now give details on deriving excess risk bounds based on the estimation error bounds in Theorem 4. We can decompose the excess risk into optimization errors and estimation errors as follows (we omit $\log(1/\delta)$) [3]

$$\begin{aligned} R(\mathbf{w}_T) - R(\mathbf{w}_R^*) &= R(\mathbf{w}_T) - R_S(\mathbf{w}_T) + R_S(\mathbf{w}_T) - R_S(\mathbf{w}_R^*) + R_S(\mathbf{w}_R^*) - R(\mathbf{w}_R^*) \\ &= (R(\mathbf{w}_T) - R_S(\mathbf{w}_T)) + (R_S(\mathbf{w}_T) - R_S(\mathbf{w}_R^*)) + O(n^{-\frac{1}{2}}), \end{aligned} \quad (\text{C.3})$$

where we have used (B.9). The first term is the estimation error and comes from the approximation of testing errors by training errors. The second term is the optimization error which comes since the optimization algorithm may not output the exact minimizer. Then Theorem 4 actually presents estimation error bounds. If we further assume $\|\mathbf{w}_t\| \leq B$ for some $B > 0$ and all t , then it was shown with high probability that [9]

$$R_S(\mathbf{w}_T) - R_S(\mathbf{w}_R^*) = O(T^{-\frac{1}{2}} \log T). \quad (\text{C.4})$$

We can plug the above optimization error bounds and the estimation error bounds in Theorem 4 into (C.3), and get with high probability

$$R(\mathbf{w}_T) - R(\mathbf{w}_R^*) = O\left(\log n \sqrt{T}/n + n^{-\frac{1}{2}} \log n\right) + O(T^{-\frac{1}{2}} \log T).$$

One can take an optimal $T \asymp n$ to trade-off the optimization and estimation errors, and get

$$R(\mathbf{w}_T) - R(\mathbf{w}_R^*) = O(n^{-\frac{1}{2}} \log n).$$

Remark C.2. If we plug the uniform stability bounds in Lemma C.3 into the existing connection between stability and generalization established in (4.2), we get with high probability that

$$|R_S(\mathbf{w}_T) - R(\mathbf{w}_T)| = O\left(\sqrt{n} \sum_{t=1}^T \eta_t \mathbb{I}[i_t = n \text{ or } j_t = n] + n^{-\frac{1}{2}}\right).$$

This together with (C.2) shows the following inequality with high probability ($\eta_t = \eta = O(1/\sqrt{T})$)

$$|R_S(\mathbf{w}_T) - R(\mathbf{w}_T)| = O\left(\frac{T\eta}{\sqrt{n}}(1 + \sqrt{n/T}) + n^{-\frac{1}{2}}\right) = O\left(\frac{\sqrt{T}}{\sqrt{n}}(1 + \sqrt{n/T}) + n^{-\frac{1}{2}}\right).$$

We can plug the above estimation error bound, the optimization error bound (C.4) back into (C.3), and derive the following excess risk bound with high probability

$$R(\mathbf{w}_T) - R(\mathbf{w}_R^*) = O\left(\frac{\log T}{\sqrt{T}} + \frac{\sqrt{T}}{\sqrt{n}}(1 + \sqrt{n/T}) + n^{-\frac{1}{2}}\right) = O(1).$$

D Proofs on Optimistic Bounds

In this section, we prove optimistic bounds in Theorem 6 by using the smoothness of loss functions. We first prove Theorem 5 on the connection between generalization and on-average stability.

Proof of Theorem 5. For all $i, j \in [n]$, let $S_{i,j}$ be defined by (3.4). Due to the symmetry, we know $\mathbb{E}[R(A(S))] = \mathbb{E}[R(A(S_{i,j}))]$ for all $i, j \in [n]$ with $i \neq j$ and therefore

$$\begin{aligned} \mathbb{E}[R(A(S)) - R_S(A(S))] &= \frac{1}{n(n-1)} \sum_{i,j \in [n]: i \neq j} \mathbb{E}[R(A(S_{i,j})) - R_S(A(S))] \\ &= \frac{1}{n(n-1)} \sum_{i,j \in [n]: i \neq j} \mathbb{E}[\ell(A(S_{i,j}); z_i, z_j) - \ell(A(S); z_i, z_j)] \leq \gamma, \end{aligned}$$

where the second identity holds since $A(S_{i,j})$ is independent of z_i and z_j . The proof is complete. \square

We then introduce some basic properties of smooth functions. For a α -smooth and non-negative function g , we have the following self-bounding property [14]

$$\|g'(\mathbf{w})\|^2 \leq 2\alpha g(\mathbf{w}), \quad \forall \mathbf{w} \in \mathcal{W} \quad (\text{D.1})$$

and the following elementary inequality

$$g(\mathbf{w}) \leq g(\mathbf{w}') + \langle g'(\mathbf{w}'), \mathbf{w} - \mathbf{w}' \rangle + \frac{\alpha \|\mathbf{w} - \mathbf{w}'\|^2}{2}, \quad \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}. \quad (\text{D.2})$$

We then present a useful lemma.

Lemma D.1. *Let S, S' be defined in Definition 2. Assume for all z, z' , $\ell(\cdot, z, z')$ is α -smooth w.r.t. a norm. For all $i \in [n]$, let S_i be defined as (A.2) and $\epsilon > 0$. Then*

$$\mathbb{E}[R(A(S)) - R_S(A(S))] \leq \frac{\alpha \mathbb{E}[R_S(A(S))]}{\epsilon} + \frac{2(\epsilon + \alpha)}{n} \sum_{i \in [n]} \mathbb{E}[\|A(S_i) - A(S)\|^2].$$

Proof. For all $i, j \in [n]$, let $S_{i,j}$ be defined by (3.4). According to (D.2), the Cauchy-Schwartz inequality and (D.1), for all $i, j \in [n]$ we know

$$\begin{aligned} \ell(A(S_{i,j}); z_i, z_j) - \ell(A(S); z_i, z_j) &\leq \langle \ell'(A(S); z_i, z_j), A(S_{i,j}) - A(S) \rangle + \frac{\alpha}{2} \|A(S_{i,j}) - A(S)\|^2 \\ &\leq \|\ell'(A(S); z_i, z_j)\| \|A(S_{i,j}) - A(S)\| + \frac{\alpha}{2} \|A(S_{i,j}) - A(S)\|^2 \\ &\leq \frac{\|\ell'(A(S); z_i, z_j)\|^2}{2\epsilon} + \frac{\epsilon + \alpha}{2} \|A(S_{i,j}) - A(S)\|^2 \\ &\leq \frac{\alpha \ell(A(S); z_i, z_j)}{\epsilon} + \frac{\epsilon + \alpha}{2} \|A(S_{i,j}) - A(S)\|^2. \end{aligned}$$

We can plug the above inequality into Theorem 5 to derive

$$\begin{aligned} &\mathbb{E}[R(A(S)) - R_S(A(S))] \\ &\leq \frac{\alpha}{\epsilon n(n-1)} \sum_{i \neq j} \mathbb{E}[\ell(A(S); z_i, z_j)] + \frac{\epsilon + \alpha}{2n(n-1)} \sum_{i \neq j} \mathbb{E}[\|A(S_{i,j}) - A(S)\|^2] \\ &= \frac{\alpha \mathbb{E}[R_S(A(S))]}{\epsilon} + \frac{\epsilon + \alpha}{2n(n-1)} \sum_{i \neq j} \mathbb{E}[\|A(S_{i,j}) - A(S)\|^2]. \end{aligned} \quad (\text{D.3})$$

By the elementary inequality $(a+b)^2 \leq 2a^2 + 2b^2$, we get the following inequality for all $i \neq j$

$$\begin{aligned} \mathbb{E}[\|A(S_{i,j}) - A(S)\|^2] &\leq 2\mathbb{E}[\|A(S_{i,j}) - A(S_i)\|^2] + 2\mathbb{E}[\|A(S_i) - A(S)\|^2] \\ &= 2\mathbb{E}[\|A(S_i) - A(S)\|^2] + 2\mathbb{E}[\|A(S_j) - A(S)\|^2], \end{aligned}$$

where we have used the following identity due to the symmetry between z_i and z'_i

$$\mathbb{E}[\|A(S_{i,j}) - A(S_i)\|^2] = \mathbb{E}[\|A(S_j) - A(S)\|^2].$$

Plugging the above inequality back into (D.3), we know

$$\begin{aligned} \mathbb{E}[R(A(S)) - R_S(A(S))] &\leq \frac{\alpha \mathbb{E}[R_S(A(S))]}{\epsilon} \\ &\quad + \frac{\epsilon + \alpha}{n(n-1)} \sum_{i \neq j} \left(\mathbb{E}[\|A(S_i) - A(S)\|^2] + \mathbb{E}[\|A(S_j) - A(S)\|^2] \right). \end{aligned}$$

This yields the stated inequality and finishes the proof. \square

Proof of Theorem 6. According to Lemma B.1 and the α -smoothness of ℓ , we know the following inequality for any k

$$\begin{aligned} n(n-1)(F_S(A(S_k)) - F_S(A(S))) &\leq \sum_{i \in [n]: i \neq k} \left(\left\langle \ell'(A(S); z_i, z_k) + \ell'(A(S); z_k, z_i) \right. \right. \\ &\quad \left. \left. - \ell'(A(S_k); z_i, z'_k) - \ell'(A(S_k); z'_k, z_i); A(S_k) - A(S) \right\rangle + \frac{4\alpha \|A(S_k) - A(S)\|^2}{2} \right). \end{aligned}$$

It then follows from the Cauchy-Schwartz inequality that

$$\begin{aligned} n(n-1)(F_S(A(S_k)) - F_S(A(S))) &\leq \sum_{i \in [n]: i \neq k} \left(\|\ell'(A(S); z_i, z_k)\| + \|\ell'(A(S); z_k, z_i)\| \right. \\ &\quad \left. + \|\ell'(A(S_k); z_i, z'_k)\| + \|\ell'(A(S_k); z'_k, z_i)\| \right) \|A(S_k) - A(S)\| + 2\alpha(n-1)\|A(S_k) - A(S)\|^2. \end{aligned}$$

This, together with (D.1) and (B.3), implies

$$\begin{aligned} \frac{\sigma n(n-1)\|A(S_k) - A(S)\|^2}{2} &\leq \sqrt{2\alpha} \sum_{i \in [n]: i \neq k} \left(\sqrt{\ell(A(S); z_i, z_k)} + \sqrt{\ell(A(S); z_k, z_i)} \right. \\ &\quad \left. + \sqrt{\ell(A(S_k); z_i, z'_k)} + \sqrt{\ell(A(S_k); z'_k, z_i)} \right) \|A(S_k) - A(S)\| + 2\alpha(n-1)\|A(S_k) - A(S)\|^2 \end{aligned}$$

and further

$$\begin{aligned} \frac{\sigma n(n-1)\|A(S_k) - A(S)\|}{2} &\leq \sqrt{2\alpha} \sum_{i \in [n]: i \neq k} \left(\sqrt{\ell(A(S); z_i, z_k)} + \sqrt{\ell(A(S); z_k, z_i)} \right. \\ &\quad \left. + \sqrt{\ell(A(S_k); z_i, z'_k)} + \sqrt{\ell(A(S_k); z'_k, z_i)} \right) + 2\alpha(n-1)\|A(S_k) - A(S)\|. \end{aligned}$$

Since $2\alpha \leq \sigma n/4$, we further get

$$\begin{aligned} \frac{\sigma n(n-1)\|A(S_k) - A(S)\|}{4} &\leq \sqrt{2\alpha} \sum_{i \in [n]: i \neq k} \left(\sqrt{\ell(A(S); z_i, z_k)} + \sqrt{\ell(A(S); z_k, z_i)} \right. \\ &\quad \left. + \sqrt{\ell(A(S_k); z_i, z'_k)} + \sqrt{\ell(A(S_k); z'_k, z_i)} \right). \end{aligned}$$

Taking a square over both sides and using the standard inequality $(\sum_{i=1}^{n-1} a_i)^2 \leq (n-1) \sum_{i=1}^{n-1} a_i^2$, we derive

$$\begin{aligned} \frac{\sigma^2 n^2 (n-1)^2 \|A(S_k) - A(S)\|^2}{16} &\leq 2\alpha(n-1) \sum_{i \in [n]: i \neq k} \left(\sqrt{\ell(A(S); z_i, z_k)} + \sqrt{\ell(A(S); z_k, z_i)} \right. \\ &\quad \left. + \sqrt{\ell(A(S_k); z_i, z'_k)} + \sqrt{\ell(A(S_k); z'_k, z_i)} \right)^2. \end{aligned}$$

This, further together with the inequality $(a + b + c + d)^2 \leq 4(a^2 + b^2 + c^2 + d^2)$, implies

$$\begin{aligned} \sigma^2 n^2 (n-1) \|A(S_k) - A(S)\|^2 &\leq 128\alpha \sum_{i \in [n]: i \neq k} \left(\ell(A(S); z_i, z_k) + \ell(A(S); z_k, z_i) \right. \\ &\quad \left. + \ell(A(S_k); z_i, z'_k) + \ell(A(S_k); z'_k, z_i) \right). \end{aligned}$$

Taking a summation of the above inequality from $k = 1$ to n , we get

$$\begin{aligned} \sigma^2 n^2 (n-1) \sum_{k=1}^n \|A(S_k) - A(S)\|^2 &\leq 128\alpha \sum_{i, k \in [n]: i \neq k} \left(\ell(A(S); z_i, z_k) + \ell(A(S); z_k, z_i) \right. \\ &\quad \left. + \ell(A(S_k); z_i, z'_k) + \ell(A(S_k); z'_k, z_i) \right). \quad (\text{D.4}) \end{aligned}$$

Due to the symmetry, we know

$$\mathbb{E}[\ell(A(S_k); z_i, z'_k)] = \mathbb{E}[\ell(A(S); z_i, z_k)], \quad \forall i \neq k.$$

It then follows that

$$\begin{aligned} &\sum_{i, k \in [n]: i \neq k} \mathbb{E} \left[\ell(A(S); z_i, z_k) + \ell(A(S); z_k, z_i) + \ell(A(S_k); z_i, z'_k) + \ell(A(S_k); z'_k, z_i) \right] \\ &= \sum_{i, k \in [n]: i \neq k} \mathbb{E} \left[\ell(A(S); z_i, z_k) + \ell(A(S); z_k, z_i) + \ell(A(S); z_i, z_k) + \ell(A(S); z_k, z_i) \right] \\ &= 4n(n-1) \mathbb{E}[R_S(A(S))]. \end{aligned}$$

We can plug the above inequality back into (D.4) and derive that

$$\sigma^2 n \sum_{k=1}^n \mathbb{E}[\|A(S_k) - A(S)\|^2] \leq 512\alpha \mathbb{E}[R_S(A(S))].$$

We now plug the above inequality back into Lemma D.1 and derive that the following inequality for all $\epsilon > 0$

$$\mathbb{E}[R(A(S)) - R_S(A(S))] \leq \frac{\alpha \mathbb{E}[R_S(A(S))]}{\epsilon} + \frac{1024(\epsilon + \alpha)\alpha}{n^2 \sigma^2} \mathbb{E}[R_S(A(S))].$$

We can take $\epsilon = \frac{n\sigma}{32}$ to derive

$$\mathbb{E}[R(A(S)) - R_S(A(S))] \leq \left(\frac{1024\alpha^2}{n^2 \sigma^2} + \frac{64\alpha}{n\sigma} \right) \mathbb{E}[R_S(A(S))].$$

Furthermore, according to the definition of $A(S)$ we know (\mathbf{w}^* is independent of S)

$$\begin{aligned} \mathbb{E}[F(A(S))] - F(\mathbf{w}^*) &= \mathbb{E}[F(A(S)) - F_S(A(S))] + \mathbb{E}[F_S(A(S)) - F_S(\mathbf{w}^*)] \\ &\leq \mathbb{E}[F(A(S)) - F_S(A(S))] = \mathbb{E}[R(A(S)) - R_S(A(S))]. \end{aligned}$$

This finishes the proof of (4.8).

We now turn to the bound of $\mathbb{E}[R(A(S))] - R(\mathbf{w}_R^*)$. Analogously to (B.8), we know

$$\begin{aligned} \mathbb{E}[R(A(S)) - R(\mathbf{w}_R^*)] &\leq \mathbb{E}[R(A(S)) - R_S(A(S))] + O(\sigma \|\mathbf{w}_R^*\|^2) \\ &= O\left(\frac{1}{n\sigma}\right) \mathbb{E}[R_S(A(S))] + O(\sigma \|\mathbf{w}_R^*\|^2), \quad (\text{D.5}) \end{aligned}$$

where we have used (4.8) in the last step. According to the definition of $A(S)$, we further know

$$R_S(A(S)) + r(A(S)) \leq R_S(\mathbf{w}_R^*) + r(\mathbf{w}_R^*) = R_S(\mathbf{w}_R^*) + O(\sigma \|\mathbf{w}_R^*\|^2).$$

Since \mathbf{w}_R^* is independent of S , we can take expectation to derive

$$\mathbb{E}[R_S(A(S))] = R(\mathbf{w}_R^*) + O(\sigma \|\mathbf{w}_R^*\|^2).$$

We can plug the above inequality back into (D.5), and derive

$$\mathbb{E}[R(A(S)) - R(\mathbf{w}_R^*)] = O\left(\frac{R(\mathbf{w}_R^*)}{n\sigma} + O(n^{-1} + \sigma)\|\mathbf{w}_R^*\|^2\right).$$

We can take

$$\sigma = \max\left\{\frac{8\alpha}{n}, \sqrt{\frac{R(\mathbf{w}_R^*)}{n\|\mathbf{w}_R^*\|^2}}\right\}$$

and derive

$$\mathbb{E}[R(A(S)) - R(\mathbf{w}_R^*)] = O\left(\frac{\sqrt{R(\mathbf{w}_R^*)\|\mathbf{w}_R^*\|}}{\sqrt{n}} + \frac{\|\mathbf{w}_R^*\|^2}{n}\right).$$

This establishes (4.9) and finishes the proof. \square

E Proofs on Applications

In this section, we present proofs for applications of our general results to metric learning.

Proof of Corollary 10. It is well known that F_S is 2λ -strongly convex w.r.t. $\|\cdot\|$. To apply Theorem 3, we require to check (4.3). For all $\mathbf{w}, \mathbf{w}', z, z'$, we know

$$\begin{aligned} & |\ell^\psi(\mathbf{w}; z, z') - \ell^\psi(\mathbf{w}'; z, z')| \\ &= \left| \max\{0, 1 - \tau(y, y')(1 - h_{\mathbf{w}}(x, x'))\} - \max\{0, 1 - \tau(y, y')(1 - h_{\mathbf{w}'}(x, x'))\} \right| \\ &\leq |\tau(y, y')| |h_{\mathbf{w}}(x, x') - h_{\mathbf{w}'}(x, x')| \leq |\langle \mathbf{w} - \mathbf{w}', (x - x')(x - x')^\top \rangle| \\ &\leq 4B^2 \|\mathbf{w} - \mathbf{w}'\|. \end{aligned}$$

Therefore, (4.3) holds with $L = 4B^2$. The proof then completes by applying Theorem 3. \square

Proof of Corollary 11. To apply Theorem 4, it suffices to show the smoothness of the loss function. The gradient of ℓ^ψ w.r.t. \mathbf{w} can be calculated by

$$\nabla \ell^\psi(\mathbf{w}; z, z') = -\psi'(\tau(y, y')(1 - h_{\mathbf{w}}(x, x')))\tau(y, y')(x - x')(x - x')^\top.$$

Then, for any \mathbf{w} and $\mathbf{w}' \in \mathcal{W}$ we have

$$\begin{aligned} & \|\nabla \ell^\psi(\mathbf{w}; z, z') - \nabla \ell^\psi(\mathbf{w}'; z, z')\|_K \\ &\leq \|\tau(y, y')(x - x')(x - x')^\top\| \|\psi'(\tau(y, y')(1 - h_{\mathbf{w}}(x, x')) - \psi'(\tau(y, y')(1 - h_{\mathbf{w}'}(x, x')))\| \\ &\leq 4B^2 |\psi'(\tau(y, y')(1 - h_{\mathbf{w}}(x, x')) - \psi'(\tau(y, y')(1 - h_{\mathbf{w}'}(x, x')))| \\ &\leq 4B^2 |\tau(y, y')| |(1 - h_{\mathbf{w}}(x, x')) - (1 - h_{\mathbf{w}'}(x, x'))| \\ &= 4B^2 |\langle \mathbf{w} - \mathbf{w}', (x - x')(x - x')^\top \rangle| \\ &\leq 16B^4 \|\mathbf{w} - \mathbf{w}'\|, \end{aligned}$$

where we have used the 1-smoothness of the logistic loss in the third step. That is, ℓ^ψ is $(16B^4)$ -smooth w.r.t. the Frobenius norm. The stated bound then follows from Theorem 4. \square

F Minimax Optimal Excess Risk Bounds for Pairwise Learning

Here we explain that the bound $O(n^{-\frac{1}{2}})$ is minimax optimal for the excess risks in pairwise learning. To see this, we consider pairwise loss functions which do not depend on the second example, i.e., $\ell(\mathbf{w}; z, z') = \ell(\mathbf{w}; z, \tilde{z}')$ for all $z', \tilde{z}' \in \mathcal{Z}$. Then it is clear that R_S defined in (3.1) becomes

$$R_S(\mathbf{w}) = \frac{1}{n} \sum_{i \in [n]} \frac{1}{n-1} \sum_{j \in [n]: j \neq i} \ell(\mathbf{w}; z_i, z_j) = \frac{1}{n} \sum_{i \in [n]} \ell(\mathbf{w}; z_i, z_0) := \tilde{R}_S(\mathbf{w}),$$

where z_0 is any fixed point in \mathcal{Z} . This is actually an objective function for pointwise learning. We know that for any estimator we can find a pointwise learning problem such that this estimator has the excess risk bound $O(n^{-\frac{1}{2}})$ [15]. Then, for any estimator we can build a pairwise learning problem such that this estimator has at best the excess risk bound $O(n^{-\frac{1}{2}})$. Furthermore, we can construct such a pairwise learning problem with the loss function independent of the second example.

References

- [1] S. Agarwal and P. Niyogi. Generalization bounds for ranking algorithms via algorithmic stability. *Journal of Machine Learning Research*, 10(Feb):441–474, 2009.
- [2] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford university press, 2013.
- [3] O. Bousquet and L. Bottou. The tradeoffs of large scale learning. In *Advances in Neural Information Processing Systems*, pages 161–168, 2008.
- [4] O. Bousquet, Y. Klochkov, and N. Zhivotovskiy. Sharper bounds for uniformly stable algorithms. In *Conference on Learning Theory*, pages 610–626, 2020.
- [5] S. Cléménçon, G. Lugosi, and N. Vayatis. Ranking and empirical minimization of U-statistics. *The Annals of Statistics*, pages 844–874, 2008.
- [6] C. Cortes, M. Mohri, and A. Rastogi. Magnitude-preserving ranking algorithms. In *Proceedings of the 24th international conference on Machine learning*, pages 169–176, 2007.
- [7] V. Feldman and J. Vondrak. High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. In *Conference on Learning Theory*, pages 1270–1279, 2019.
- [8] M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning*, pages 1225–1234, 2016.
- [9] N. J. Harvey, C. Liaw, Y. Plan, and S. Randhawa. Tight analyses for non-smooth stochastic gradient descent. In *Conference on Learning Theory*, pages 1579–1613, 2019.
- [10] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- [11] R. Jin, S. Wang, and Y. Zhou. Regularized distance metric learning: Theory and algorithm. In *Advances in Neural Information Processing Systems*, pages 862–870, 2009.
- [12] T. Peel, S. Anthoine, and L. Ralaivola. Empirical bernstein inequalities for u-statistics. In *Advances in Neural Information Processing Systems*, pages 1903–1911, 2010.
- [13] Y. Pitcan. A note on concentration inequalities for u-statistics, 2017.
- [14] N. Srebro, K. Sridharan, and A. Tewari. Smoothness, low noise and fast rates. In *Advances in Neural Information Processing Systems*, pages 2199–2207, 2010.
- [15] A. B. Tsybakov. *Introduction to nonparametric estimation*. Springer Science & Business Media, 2008.
- [16] R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [17] B. Wang, H. Zhang, P. Liu, Z. Shen, and J. Pineau. Multitask metric learning: Theory and algorithm. In *International Conference on Artificial Intelligence and Statistics*, pages 3362–3371, 2019.