

A. Detailed Choices of Reserved Constants

The absolute constants c_0, c_1 and c_2 are specified in Lemma 24, and c_3 and c_4 are specified in Lemma 25. c_5 and c_6 are clarified in Section 3.1.1. The definition of c_7 and c_8 can be found in Lemma 26 and Lemma 27 respectively. The absolute constant C_1 acts as an upper bound of all b_k 's, and by our choice in Section 3.1.1, $C_1 = \bar{c}/16$. The absolute constant C_2 is defined in Lemma 3. Other absolute constants, such as C_3, C_4 are not quite crucial to our analysis or algorithmic design. Therefore, we do not track their definitions. The subscript variants of K , e.g. K_1 and K_2 , are also absolute constants but their values may change from appearance to appearance. We remark that the value of all these constants does not depend on the underlying distribution D chosen by the adversary, but rather depends on the knowledge that D is a member of the family of isotropic log-concave distributions.

B. Omitted Proofs in Section 2

We will frequently use the well-known Chernoff bound in our analysis. For convenience, we record it below.

Lemma 17 (Chernoff bound). *Let Z_1, Z_2, \dots, Z_n be n independent random variables that take value in $\{0, 1\}$. Let $Z = \sum_{i=1}^n Z_i$. For each Z_i , suppose that $\Pr(Z_i = 1) \leq \eta$. Then for any $\alpha \in [0, 1]$*

$$\Pr(Z \geq (1 + \alpha)\eta n) \leq e^{-\frac{\alpha^2 \eta n}{3}}.$$

When $\Pr(Z_i = 1) \geq \eta$, for any $\alpha \in [0, 1]$

$$\Pr(Z \leq (1 - \alpha)\eta n) \leq e^{-\frac{\alpha^2 \eta n}{2}}.$$

B.1. Proof of Lemma 4

Proof. We note that $(\rho^+ - \rho^-)^2 \leq 4(\rho^+)^2$. In addition, this inequality is almost tight up to a constant factor since ρ^- can be as small as 0. To see this, observe that $u \in W$ and x is such that $|u \cdot x| \leq b$.

Thus, it remains to upper bound ρ^+ . Due to localized sampling, for any $w \in W$ we have

$$|w \cdot x| \leq |(w - u) \cdot x| + |u \cdot x| \leq \|w - u\|_2 \cdot \|x\|_2 + b \leq r \cdot c_7 \sqrt{d} \log \frac{1}{b\delta} + b, \quad (5)$$

where the first step follows from the triangle inequality, the second step uses Cauchy-Schwarz inequality and the fact $x \sim D_{u,b}$, and the last step applies Lemma 26. The lemma follows by noting that $r = \Theta(b)$. \square

B.2. Proof of Lemma 7

Proof. For any unit vector v , observe that $w := rv + u$ is such that $\|w - u\|_2 \leq r$. Hence,

$$\begin{aligned} \mathbb{E}[(v \cdot x)^2] &= \frac{1}{r^2} \mathbb{E}[(r \cdot v \cdot x)^2] \\ &\leq \frac{2}{r^2} \mathbb{E}[(r \cdot v + u) \cdot x]^2 + \frac{2}{r^2} \mathbb{E}[(u \cdot x)^2] \\ &\leq \frac{2}{r^2} \cdot C_2(b^2 + r^2) + \frac{2}{r^2} \cdot b^2 \\ &\leq \frac{4C_2(b^2 + r^2)}{r^2}, \end{aligned}$$

where in the second step we use the basic inequality $a_1^2 \leq 2(a_1 - a_2)^2 + 2a_2^2$, and in the third step we apply Lemma 3. This proves the first desired inequality.

Next, by Lemma 26 we have with probability $1 - \delta$, $\|x\|_2 \leq c_7 \sqrt{d} \log \frac{1}{b\delta}$. Then for any unit vector v , we have

$$(v \cdot x)^2 \leq \|v\|_2^2 \cdot \|x\|_2^2 \leq c_7^2 \cdot d \log^2 \frac{1}{b\delta},$$

which implies the second desired inequality. \square

B.3. Proof of Proposition 8

Proof. In Lemma 6, we set $\alpha = 1$, $M_i = x_i x_i^\top$ where x_i is the i -th instance in the set T_C . Lemma 7 implies that $\mu_{\max} \leq \frac{4C_2(b^2+r^2)}{r^2} |T_C| \leq K \cdot |T_C|$ for some constant $K > 0$ since $r = \Theta(b)$, and with probability $1 - \delta$, $\Lambda \leq K_1 \cdot d \log^2 \frac{|T_C|}{b\delta}$ by union bound. By conditioning on these events and putting all pieces together, Lemma 6 asserts that with probability

$$1 - d \cdot \left(\frac{\epsilon}{4}\right)^{\frac{K}{K_1} \cdot \frac{|T_C|}{d \log^2 \frac{|T_C|}{b\delta}}}, \quad \lambda_{\max} \left(\sum_{x \in S} x x^\top \right) \leq 2K \cdot |T_C|. \quad (6)$$

Equivalently, the above holds with probability $1 - \delta$ as long as $|T_C| \geq K_2 d \log^2 \frac{|T_C|}{b\delta} \cdot \log \frac{d}{\delta}$ for some constant $K_2 > 0$. \square

B.4. Proof of Lemma 9

Proof. By Lemma 27

$$\Pr_{x \sim D}(x \in X) \geq c_8 b.$$

This implies that

$$\begin{aligned} & \Pr_{x \sim \text{EX}_\eta^x(D, w^*)}(x \in X_{u,b} \text{ and } x \text{ is clean}) \\ &= \Pr_{x \sim \text{EX}_\eta^x(D, w^*)}(x \in X_{u,b} \mid x \text{ is clean}) \cdot \Pr_{x \sim \text{EX}_\eta^x(D, w^*)}(x \text{ is clean}) \geq c_8 b(1 - \eta). \end{aligned}$$

We want to ensure that by drawing N instances from $\text{EX}_\eta^x(D, w^*)$, with probability at least $1 - \delta$, n out of them fall into the band $X_{u,b}$. We apply the second inequality of Lemma 17 by letting $Z_i = \mathbf{1}_{\{x_i \in X_{u,b} \text{ and } x_i \text{ is clean}\}}$ and $\alpha = 1/2$, and obtain

$$\Pr \left(|T_C| \leq \frac{c_8 b(1 - \eta)}{2} N \right) \leq \exp \left(-\frac{c_8 b(1 - \eta) N}{8} \right),$$

where the probability is taken over the event that we make a number of N calls to $\text{EX}_\eta^x(D, w^*)$. Thus, when $N \geq \frac{8}{c_8 b(1 - \eta)} \left(n + \ln \frac{1}{\delta} \right)$, we are guaranteed that at least n samples from $\text{EX}_\eta^x(D, w^*)$ fall into the band $X_{u,b}$ with probability $1 - \delta$. The lemma follows by observing $\eta < \frac{1}{2}$. \square

B.5. Proof of Lemma 10

This is a simplified version of Lemma 30 of Shen & Zhang (2021).

Proof. We calculate the noise rate within the band $X_k := \{x : |w_{k-1} \cdot x| \leq b_k\}$ by Lemma 18:

$$\Pr_{x \sim \text{EX}_\eta^x(D, w^*)}(x \text{ is dirty} \mid x \in X_{u,b}) \leq \frac{2\eta}{c_8 b} \leq \frac{2\eta}{c_8 \epsilon} \leq \frac{2c_5}{c_8} \leq \frac{1}{8},$$

where the second inequality applies the setting $b \geq \epsilon$, the third inequality is due to the condition $\eta \leq c_5 \epsilon$, and the last inequality is due to the condition that c_5 is assumed to be a sufficiently small constant. Now we apply the first inequality of Lemma 17 by specifying $Z_i = \mathbf{1}_{\{x_i \text{ is dirty}\}}$, $\alpha = 1$ therein, which gives

$$\Pr \left(|T_D| \geq \frac{1}{4} |T| \right) \leq \exp \left(-\frac{|T|}{24} \right),$$

where the probability is taken over the draw of T . The lemma follows by setting the right-hand side to δ and noting that $|T_C| = |T| - |T_D|$. \square

Lemma 18. Assume $\eta < \frac{1}{2}$. We have

$$\Pr_{x \sim \text{EX}_\eta^x(D, w^*)}(x \text{ is dirty} \mid x \in X_{u,b}) \leq \frac{2\eta}{c_8 b}$$

where c_8 was defined in Lemma 27.

Proof. For an instance x , we use $\text{tag}_x = 1$ to denote that x is drawn from D , and use $\text{tag}_x = -1$ to denote that x is adversarially generated.

We first calculate the probability that an instance returned by $\text{EX}_\eta^x(D, w^*)$ falls into the band $X_{u,b}$ as follows:

$$\begin{aligned}
 & \Pr_{x \sim \text{EX}_\eta^x(D, w^*)} (x \in X_{u,b}) \\
 &= \Pr_{x \sim \text{EX}_\eta^x(D, w^*)} (x \in X_{u,b} \text{ and } \text{tag}_x = 1) + \Pr_{x \sim \text{EX}_\eta^x(D, w^*)} (x \in X_{u,b} \text{ and } \text{tag}_x = -1) \\
 &\geq \Pr_{x \sim \text{EX}_\eta^x(D, w^*)} (x \in X_{u,b} \text{ and } \text{tag}_x = 1) \\
 &= \Pr_{x \sim \text{EX}_\eta^x(D, w^*)} (x \in X_{u,b} \mid \text{tag}_x = 1) \cdot \Pr_{x \sim \text{EX}_\eta^x(D, w^*)} (\text{tag}_x = 1) \\
 &= \Pr_{x \sim D} (x \in X_{u,b}) \cdot \Pr_{x \sim \text{EX}_\eta^x(D, w^*)} (\text{tag}_x = 1) \\
 &\stackrel{\zeta}{\geq} c_8 b \cdot (1 - \eta) \\
 &\geq \frac{1}{2} c_8 b,
 \end{aligned}$$

where in the inequality ζ we applied Part 1 of Lemma 27. It is thus easy to see that

$$\Pr_{x \sim \text{EX}_\eta^x(D, w^*)} (\text{tag}_x = -1 \mid x \in X_{u,b}) \leq \frac{\Pr_{x \sim \text{EX}_\eta^x(D, w^*)} (\text{tag}_x = -1)}{\Pr_{x \sim \text{EX}_\eta^x(D, w^*)} (x \in X_{u,b})} \leq \frac{2\eta}{c_8 b},$$

which is the desired result. \square

B.6. Rademacher analysis leads to suboptimal sample complexity for quadratic functions

To see why a general Rademacher analysis may not suffice, we can, for example, think of the quadratic function $(w \cdot x)^2$ as a composition of the functions $\phi(f) = f^2$ and $f_w(x) = w \cdot x$. Recall that we showed with high probability that $|w \cdot x| \leq O(b\sqrt{d})$ (omitting logarithmic factors for convenience). Now, the gradient of $\phi(\cdot)$ is $2w \cdot x$ which is upper bounded by $O(b\sqrt{d})$ and the function value of $\phi(\cdot)$ is upper bounded by $O(b^2 d)$. For the Rademacher complexity $\mathcal{R}_{\mathcal{F}}$ of the class of linear functions $\mathcal{F} := \{f_w(x) = w \cdot x : w \in W\}$ on $T_C = \{x_1, \dots, x_n\}$, let $V = \{v \in \mathbb{R}^d : \|v\|_2 \leq 1\}$ and note that for any $w \in W$, $w = u + rv$. We have by definition

$$\begin{aligned}
 \mathcal{R}_{\mathcal{F}} &= \frac{1}{n} \mathbb{E} \sup_{w \in W} \sum_{i=1}^n \sigma_i (w \cdot x_i) \\
 &= \frac{1}{n} \mathbb{E} \sup_{w \in W} w \cdot \sum_{i=1}^n \sigma_i x_i \\
 &\leq \frac{r}{n} \mathbb{E} \sup_{v \in V} v \cdot \sum_{i=1}^n \sigma_i x_i + \frac{1}{n} \mathbb{E} u \cdot \sum_{i=1}^n \sigma_i x_i \\
 &\leq \frac{r}{n} \mathbb{E} \left\| \sum_{i=1}^n \sigma_i x_i \right\|_2 \\
 &\leq \frac{r}{n} \cdot \sqrt{n} \max_{1 \leq i \leq n} \|x_i\|_2,
 \end{aligned}$$

where the expectation is taken over the i.i.d. Rademacher variables $\sigma_1, \dots, \sigma_n$. By Lemma 26, $\mathcal{R}_{\mathcal{F}} \leq \frac{r}{\sqrt{n}} \sqrt{d}$ with high probability. By the contraction lemma, the Rademacher complexity of the class of quadratic functions is $O(\frac{brd}{\sqrt{n}})$, and thus uniform concentration through Rademacher analysis requires $O(d^2)$ samples.

Similarly, a straightforward application of local Rademacher analysis (Bartlett et al., 2005) may not suffice as well. However, our discussion here does not rule out the possibility that a more sophisticated exploration of these techniques would lead to the desired sample complexity bound; we leave it as an open problem.

C. Omitted Proofs in Section 3

We present a full proof of the results in Section 3. Observe that the malicious noise is a special case of the nasty noise; hence this section can also be thought of as providing a complete proof for the results in Section 2.

To improve the transparency, we collect useful notations in Table 1.

Table 1. Summary of useful notations associated with the working set T at each phase k for learning with nasty noise.

\hat{A}'	labeled clean instance set obtained by drawing N instances from D and labeling them by w^*
A'	(clean) instance set obtained by hiding all the labels in \hat{A}'
\hat{A}	labeled corrupted instance set obtained by replacing ηN samples in \hat{A}'
A	(corrupted) instance set obtained by hiding all the labels in \hat{A}
A_C	set of clean instances in A
A_D	set of dirty instances in A , i.e. $A \setminus A_C$
A_E	set of clean instances erased from A' by the adversary
T	set of instances in A that satisfy $ w_{k-1} \cdot x \leq b_k$
T_C	set of clean instances in T
T_D	set of dirty instances in T , i.e. $T \setminus T_C$
\hat{T}_C	unrevealed labeled set of T_C
\hat{T}_E	unrevealed labeled set of T_E

C.1. Proof of Lemma 11

Proof. Since $\eta \leq c_5 \epsilon$ and $b \geq \epsilon$, we have $\eta \leq c_5 b \leq \frac{1}{2} c_8 \xi b$ where the second inequality follows from the fact that c_5 is a small constant and $\xi \geq \Omega(1)$. Thus $|A_D| = \eta N \leq \frac{1}{2} c_8 \xi b N$ and $|A_C| = N - |A_D| \geq (1 - \frac{1}{2} c_8 \xi b) N$. \square

C.2. Proof of Lemma 12

Proof. We first show that the following two events hold simultaneously with probability $1 - \frac{\delta_k}{24}$:

$$E_1 : |A_C| \geq \left(1 - \frac{1}{2} c_8 \xi b\right) N \text{ and } |A_D| \leq \frac{1}{2} c_8 \xi b N,$$

$$E_2 : |T_C| \geq \frac{1}{2} c_8 (1 - \xi) b N \text{ and } |T_E| \leq \frac{1}{2} c_8 \xi b N.$$

Observe that E_1 holds with certainty due to Lemma 11.

To see why E_2 holds with high probability, we recall that Part 1 of Lemma 27 shows that $\Pr_{x \sim D}(x \in X_{u,b}) \geq c_8 b$. For each $x_i \in A_C \cup A_E$, define $Z_i = \mathbf{1}_{\{x_i \in X_{u,b}\}}$. Since $A_C \cup A_E$ are i.i.d. draws from D , by applying the second part of Lemma 17 with $\alpha = 1/2$, we have

$$\Pr \left(\sum_{i=1}^N Z_i \leq \frac{1}{2} c_8 b N \right) \leq \exp \left(-\frac{c_8 b N}{8} \right).$$

This shows that

$$|T_C| + |T_E| \geq \frac{1}{2} c_8 b N$$

with probability $1 - \delta$ provided that $N \geq \frac{8}{c_8 b} \ln \frac{1}{\delta}$. On the other side, we have $|T_E| \leq |A_E| = |A_D| \leq \frac{1}{2} c_8 \xi b N$. Thus it follows that $|T_C| \geq \frac{1}{2} c_8 (1 - \xi) b N$.

For Part 1, we have

$$\frac{|T_C|}{|T_D|} \geq \frac{1 - \xi}{\xi}, \quad (7)$$

where the inequality follows from E_2 and the fact $|T_D| = |T_E|$. Therefore,

$$\frac{|T_D|}{|T|} = \frac{1}{1 + |T_C|/|T_D|} \leq \xi. \quad (8)$$

Part 2 of the lemma simply follows E_2 . \square

C.3. Proof of Proposition 13

Proof. Since $N \geq \frac{d}{b} \cdot \text{polylog}((d, \frac{1}{\delta}))$, we have by Part 2 that $|T_C \cup T_E| \geq |T_C| \geq d \cdot \text{polylog}(d, \frac{1}{\delta})$. Therefore, we can directly apply Proposition 8 by thinking of T_C therein as $T_C \cup T_E$ in the current proposition. \square

C.4. Proof of Theorem 14

Proof. We first show the existence of a feasible function $q(x)$ to Algorithm 2. Consider the specific function $q : T \rightarrow [0, 1]$ as follows: $q(x) = 1$ for all $x \in T_C$ and $q(x) = 0$ otherwise. We have

$$\frac{1}{|T|} \sum_{x \in T} q(x) = \frac{|T_C|}{|T|} = 1 - \frac{|T_D|}{|T|} \geq 1 - \xi,$$

in view of Part 1 of Lemma 12.

To show Part 3, we note that $T_C \cup T_E$ are i.i.d. draws from $D_{u,b}$ and Lemma 12 shows that $|T_C \cup T_E| \geq \Omega(bN)$. Therefore, as far as $N \geq \frac{d}{b} \cdot \text{polylog}(d)$, Theorem 5 implies that

$$\frac{1}{|T_C| + |T_E|} \sum_{x \in T_C \cup T_E} (w \cdot x)^2 \leq \frac{c}{2} (b^2 + r^2).$$

Since $(w \cdot x)^2$ is always non-negative, we have

$$\frac{1}{|T_C|} \sum_{x \in T_C} (w \cdot x)^2 \leq \frac{|T_C| + |T_E|}{|T_C|} \cdot \frac{1}{|T_C| + |T_E|} \sum_{x \in T_C \cup T_E} (w \cdot x)^2 \leq \frac{|T_C| + |T_E|}{|T_C|} \cdot \frac{c}{2} (b^2 + r^2).$$

Part 2 of Lemma 12 shows that $|T_E|/|T_C| \leq \frac{\xi}{1-\xi} \leq 1$ since $\xi \leq \frac{1}{2}$. Plugging this upper bound into the above inequality, we obtain

$$\frac{1}{|T_C|} \sum_{x \in T_C} (w \cdot x)^2 \leq c(b^2 + r^2).$$

In a nutshell, our construction of $q(x)$ ensures the feasibility to all constraints in Algorithm 2. By ellipsoid method we are able to find a feasible solution in polynomial time. \square

C.5. Proof of Proposition 15

Let $z = \sqrt{b^2 + r^2}$. We will in fact prove a stronger result, i.e.,

$$\ell_\tau(w; \hat{T}_C \cup \hat{T}_E) \leq \ell_\tau(w; p \circ \hat{T}) + 2\xi \left(2 + \sqrt{2K_2} \cdot \frac{z}{\tau} \right) + \sqrt{2K_2\xi} \cdot \frac{z}{\tau}, \quad (9)$$

$$\ell_\tau(w; p \circ \hat{T}) \leq \ell_\tau(w; \hat{T}_C \cup \hat{T}_E) + 2\xi + \sqrt{4K_2\xi} \cdot \frac{z}{\tau}. \quad (10)$$

The claim in the proposition immediately follows since $z/\tau = \Theta(1)$ and ξ can be chosen as an arbitrarily small constant.

Let $\{q(x)\}_{x \in T}$ be the output of Algorithm 2 under the nasty noise model. We extend the domain of $q(x)$ from T to $T \cup T_E$ as follows: for any $x \in T$, the value $q(x)$ remains unchanged; for any $x \in T_E$, we set $q(x) = 0$. With this in mind, we can, for the purpose of analysis, think of the probability mass function $\{p(x)\}_{x \in T}$ obtained in Algorithm 1 as over $T \cup T_E$, with the value $p(x)$ stays unchanged for $x \in T$ and $p(x) = 0$ for all $x \in T_E$.

Now with the *extended* probability mass function $\{p(x)\}_{x \in T \cup T_E}$, we can prove the proposition.

Proof. Let \hat{T}_C and \hat{T}_E be the labeled set of T_C and T_E that is correctly annotated by w^* respectively. For any x in the instance space, let y_x be the label that the adversary is committed to. Recall that the empirical distribution $\{p(x)\}_{x \in T \cup T_E}$ was defined as follows: $p(x) = \frac{q(x)}{\sum_{x \in T} q(x)}$ for $x \in T$ and $p(x) = 0$ for $x \in T_E$. The reweighted hinge loss on $T \cup T_E$ using $p(x)$ is given by

$$\ell_\tau(w; p \circ \hat{T}) = \frac{1}{|T \cup T_E|} \sum_{x \in T \cup T_E} p(x) \cdot \max \left\{ 0, 1 - \frac{1}{\tau} y_x w \cdot x \right\}. \quad (11)$$

The choice of N guarantees that Proposition 13, Lemma 12, and Theorem 14 hold simultaneously with probability $1 - \delta$. We thus have for all $w \in W$

$$\frac{1}{|T_C \cup T_E|} \sum_{x \in T_C \cup T_E} (w \cdot x)^2 \leq K_1 z^2, \quad (12)$$

$$\frac{|T_D|}{|T|} \leq \xi, \quad (13)$$

$$\frac{1}{|T|} \sum_{x \in T} q(x) (w \cdot x)^2 \leq K_2 z^2. \quad (14)$$

We now expand T to $T \cup T_E$ for the last two inequalities. Indeed, from (13), it is easy to show that

$$\frac{|T_D|}{|T \cup T_E|} \leq \frac{|T_D|}{|T|} \leq \xi. \quad (15)$$

Next, since we defined $q(x) = 0$ for all $x \in T_E$, (14) implies that

$$\frac{1}{|T \cup T_E|} \sum_{x \in T \cup T_E} q(x) (w \cdot x)^2 = \frac{1}{|T \cup T_E|} \sum_{x \in T} q(x) (w \cdot x)^2 \leq \frac{1}{|T|} \sum_{x \in T} q(x) (w \cdot x)^2 \leq K_2 z^2. \quad (16)$$

The remaining steps are exactly same as Proposition 33 of Shen & Zhang (2021) since all the analyses therein rely only on the conditions (12), (15) and (16). For completeness, we present the full proof here.

It follows from Eq. (15) and $\xi \leq 1/2$ that

$$\frac{|T \cup T_E|}{|T_C \cup T_E|} \leq \frac{|T \cup T_E|}{|T_C|} = \frac{|T \cup T_E|}{|T \cup T_E| - |T_D|} = \frac{1}{1 - |T_D|/|T \cup T_E|} \leq \frac{1}{1 - \xi} \leq 2. \quad (17)$$

In the following, we condition on the event that all these inequalities are satisfied.

Step 1. First we upper bound $\ell_\tau(w; \hat{T}_C \cup \hat{T}_E)$ by $\ell_\tau(w; p \circ \hat{T})$.

$$\begin{aligned} |T_C \cup T_E| \cdot \ell_\tau(w; \hat{T}_C \cup \hat{T}_E) &= \sum_{x \in T_C \cup T_E} \ell(w; x, y_x) \\ &= \sum_{x \in T \cup T_E} \left[q(x) \ell(w; x, y_x) + (\mathbf{1}_{\{x \in T_C \cup T_E\}} - q(x)) \ell(w; x, y_x) \right] \\ &\stackrel{\zeta_1}{\leq} \sum_{x \in T \cup T_E} q(x) \ell(w; x, y_x) + \sum_{x \in T_C \cup T_E} (1 - q(x)) \ell(w; x, y_x) \\ &\stackrel{\zeta_2}{\leq} \sum_{x \in T \cup T_E} q(x) \ell(w; x, y_x) + \sum_{x \in T_C \cup T_E} (1 - q(x)) \left(1 + \frac{|w \cdot x|}{\tau} \right) \\ &\stackrel{\zeta_3}{\leq} \sum_{x \in T \cup T_E} q(x) \ell(w; x, y_x) + \xi |T \cup T_E| + \frac{1}{\tau} \sum_{x \in T_C \cup T_E} (1 - q(x)) |w \cdot x| \\ &\stackrel{\zeta_4}{\leq} \sum_{x \in T \cup T_E} q(x) \ell(w; x, y_x) + \xi |T \cup T_E| + \frac{1}{\tau} \sqrt{\sum_{x \in T_C \cup T_E} (1 - q(x))^2} \cdot \sqrt{\sum_{x \in T_C \cup T_E} (w \cdot x)^2} \\ &\stackrel{\zeta_5}{\leq} \sum_{x \in T \cup T_E} q(x) \ell(w; x, y_x) + \xi |T \cup T_E| + \frac{1}{\tau} \sqrt{\xi |T \cup T_E|} \cdot \sqrt{K_1 |T_C \cup T_E|} \cdot z, \end{aligned} \quad (18)$$

where ζ_1 follows from the simple fact that

$$\begin{aligned} \sum_{x \in T \cup T_E} (\mathbf{1}_{\{x \in T_C \cup T_E\}} - q(x)) \ell(w; x, y_x) &= \sum_{x \in T_C \cup T_E} (1 - q(x)) \ell(w; x, y_x) + \sum_{x \in T_D} (-q(x)) \ell(w; x, y_x) \\ &\leq \sum_{x \in T_C \cup T_E} (1 - q(x)) \ell(w; x, y_x), \end{aligned}$$

ζ_2 explores the fact that the hinge loss is always upper bounded by $1 + \frac{|w \cdot x|}{\tau}$ and that $1 - q(x) \geq 0$, ζ_3 follows from Part 2 of Theorem 14, ζ_4 applies Cauchy-Schwarz inequality, and ζ_5 uses Eq. (12).

In view of Eq. (17), we have $\frac{|T \cup T_E|}{|T_C \cup T_E|} \leq 2$. Continuing Eq. (18), we obtain

$$\begin{aligned} \ell_\tau(w; \hat{T}_C \cup \hat{T}_E) &\leq \frac{1}{|T_C \cup T_E|} \sum_{x \in T \cup T_E} q(x) \ell(w; x, y_x) + 2\xi + \sqrt{2K_1} \xi \cdot \frac{z}{\tau} \\ &= \frac{\sum_{x \in T \cup T_E} q(x)}{|T_C \cup T_E|} \sum_{x \in T \cup T_E} p(x) \ell(w; x, y_x) + 2\xi + \sqrt{2K_1} \xi \cdot \frac{z}{\tau} \\ &= \ell_\tau(w; p \circ \hat{T}) + \left(\frac{\sum_{x \in T \cup T_E} q(x)}{|T_C \cup T_E|} - 1 \right) \sum_{x \in T \cup T_E} p(x) \ell(w; x, y_x) + 2\xi + \sqrt{2K_1} \xi \cdot \frac{z}{\tau} \\ &\leq \ell_\tau(w; p \circ \hat{T}) + \left(\frac{|T \cup T_E|}{|T_C \cup T_E|} - 1 \right) \sum_{x \in T \cup T_E} p(x) \ell(w; x, y_x) + 2\xi + \sqrt{2K_1} \xi \cdot \frac{z}{\tau} \\ &\leq \ell_\tau(w; p \circ \hat{T}) + 2\xi \sum_{x \in T \cup T_E} p(x) \ell(w; x, y_x) + 2\xi + \sqrt{2K_1} \xi \cdot \frac{z}{\tau}, \end{aligned} \tag{19}$$

where in the last inequality we use the fact that $|T_E| = |T_D|$ and $T \cap T_E = \emptyset$, and thus

$$\frac{|T \cup T_E|}{|T_C \cup T_E|} - 1 = \frac{|T| + |T_D|}{|T|} - 1 = \frac{|T_D|}{|T|} \leq \xi.$$

On the other hand, we have the following result which will be proved later on.

Claim 19. $\sum_{x \in T \cup T_E} p(x) \ell(w; x, y_x) \leq 1 + \sqrt{2K_2} \cdot \frac{z}{\tau}$.

Therefore, continuing Eq. (19) we have

$$\ell_\tau(w; \hat{T}_C \cup \hat{T}_E) \leq \ell_\tau(w; p \circ \hat{T}) + 2\xi \left(2 + \sqrt{2K_2} \cdot \frac{z}{\tau} \right) + \sqrt{2K_2} \xi \cdot \frac{z}{\tau}$$

which proves the first inequality of the proposition.

Step 2. We move on to prove the second inequality of the theorem, i.e. using $\ell_\tau(w; \hat{T}_C \cup \hat{T}_E)$ to upper bound $\ell_\tau(w; p \circ \hat{T})$. Let us denote by $p_D = \sum_{x \in T_D} p(x)$ the probability mass on dirty instances. Then

$$p_D = \frac{\sum_{x \in T_D} q(x)}{\sum_{x \in T} q(x)} \leq \frac{|T_D|}{(1 - \xi)|T|} \leq \frac{\xi}{1 - \xi} \leq 2\xi, \tag{20}$$

where the first inequality follows from $q(x) \leq 1$ and Part 2 of Theorem 14, the second inequality follows from (13), and the last inequality is by our choice $\xi \leq 1/2$.

Note that by Part 2 of Theorem 14 and the choice $\xi \leq 1/2$, we have

$$\sum_{x \in T} q(x) \geq (1 - \xi)|T| \geq |T|/2$$

Hence

$$\begin{aligned}
 \sum_{x \in T} p(x)(w \cdot x)^2 &= \frac{1}{\sum_{x \in T} q(x)} \sum_{x \in T} q(x)(w \cdot x)^2 \\
 &\leq \frac{2}{|T|} \sum_{x \in T} q(x)(w \cdot x)^2 \\
 &\leq 2 \cdot K_2 z^2
 \end{aligned} \tag{21}$$

where the last inequality holds because of (14). Thus,

$$\begin{aligned}
 \sum_{x \in T_D} p(x)\ell(w; x, y_x) &\leq \sum_{x \in T_D} p(x) \left(1 + \frac{|w \cdot x|}{\tau} \right) \\
 &= p_D + \frac{1}{\tau} \sum_{x \in T_D} p(x)|w \cdot x| \\
 &= p_D + \frac{1}{\tau} \sum_{x \in T} \left(\mathbf{1}_{\{x \in T_D\}} \sqrt{p(x)} \right) \cdot \left(\sqrt{p(x)} |w \cdot x| \right) \\
 &\leq p_D + \frac{1}{\tau} \sqrt{\sum_{x \in T} \mathbf{1}_{\{x \in T_D\}} p(x)} \cdot \sqrt{\sum_{x \in T} p(x)(w \cdot x)^2} \\
 &\stackrel{(21)}{\leq} p_D + \sqrt{p_D} \cdot \sqrt{2K_2} \cdot \frac{z}{\tau}.
 \end{aligned}$$

With the result on hand, we bound $\ell_\tau(w; p \circ \hat{T})$ as follows:

$$\begin{aligned}
 \ell_\tau(w; p \circ \hat{T}) &= \sum_{x \in T_C \cup T_E} p(x)\ell(w; x, y_x) + \sum_{x \in T_D} p(x)\ell(w; x, y_x) \\
 &\leq \sum_{x \in T_C \cup T_E} \ell(w; x, y_x) + \sum_{x \in T_D} p(x)\ell(w; x, y_x) \\
 &= \ell_\tau(w; \hat{T}_C \cup \hat{T}_E) + \sum_{x \in T_D} p(x)\ell(w; x, y_x) \\
 &\leq \ell_\tau(w; \hat{T}_C \cup \hat{T}_E) + p_D + \sqrt{p_D} \cdot \sqrt{2K_2} \cdot \frac{z}{\tau} \\
 &\stackrel{(20)}{\leq} \ell_\tau(w; \hat{T}_C \cup \hat{T}_E) + 2\xi + \sqrt{4K_2\xi} \cdot \frac{z}{\tau},
 \end{aligned}$$

which proves the second inequality of the proposition.

This completes the proof. \square

Proof of Claim 19. Since $\ell(w; x, y_x) \leq 1 + \frac{|w \cdot x|}{\tau}$, it follows that

$$\begin{aligned}
 \sum_{x \in T \cup T_E} p(x)\ell(w; x, y_x) &\leq \sum_{x \in T \cup T_E} p(x) \left(1 + \frac{|w \cdot x|}{\tau} \right) \\
 &= 1 + \frac{1}{\tau} \sum_{x \in T \cup T_E} p(x)|w \cdot x| \\
 &\leq 1 + \frac{1}{\tau} \sqrt{\sum_{x \in T \cup T_E} p(x)(w \cdot x)^2} \\
 &\stackrel{(21)}{\leq} 1 + \sqrt{2K_2} \cdot \frac{z}{\tau},
 \end{aligned}$$

which completes the proof of Claim 19. \square

C.6. Proof of Lemma 16

For any phase k , let $L_{\tau_k}(w) = \mathbb{E}_{x \sim D_{w_{k-1}, b_k}} [\ell_{\tau_k}(w; x, \text{sign}(w^* \cdot x))]$.

Proof. Proposition 35 of Shen & Zhang (2021) showed that if $|T_C \cup T_E| \geq d \cdot \text{polylog}(d, \frac{1}{\delta}, \frac{1}{\delta})$, then by Rademacher complexity of the hinge loss we have that with probability $1 - \frac{\delta}{2}$

$$\sup_{w \in W} \left| \ell_{\tau}(w; \hat{T}_C \cup \hat{T}_E) - \mathbb{E}_{x \sim D_{u, b}} [\ell_{\tau}(w; x, \text{sign}(w^* \cdot x))] \right| \leq \kappa. \quad (22)$$

Combining the above with Proposition 15 gives that with probability $1 - \delta$,

$$\sup_{w \in W} \left| \ell_{\tau}(w; p \circ \hat{T}) - \mathbb{E}_{x \sim D_{u, b}} [\ell_{\tau}(w; x, \text{sign}(w^* \cdot x))] \right| \leq 2\kappa.$$

Namely, in any phase $k \leq K$, if $|T_C \cup T_E| \geq d \cdot \text{polylog}(d, \frac{1}{\delta_k}, \frac{1}{\delta_k})$, then with probability $1 - \delta_k$,

$$\sup_{w \in W_k} \left| \ell_{\tau_k}(w; p) - L_{\tau_k}(w) \right| \leq 2\kappa. \quad (23)$$

On the other hand, since the (rescaled) hinge loss is always an upper bound of the error rate, we have

$$\text{err}_{D_{w_{k-1}, b_k}}(v_k) \leq L_k(v_k) \stackrel{\zeta_1}{\leq} \ell_{\tau_k}(v_k; p) + 2\kappa \stackrel{\zeta_2}{\leq} \min_{w \in W_k} \ell_{\tau_k}(w; p) + 3\kappa \leq \ell_{\tau_k}(w^*; p) + 3\kappa \stackrel{\zeta_3}{\leq} L_k(w^*) + 5\kappa \stackrel{\zeta_4}{\leq} 6\kappa \leq 8\kappa,$$

where we use the fact that $v_k \in W_k$ in ζ_1 , use the optimality condition of v_k in ζ_2 , use $w^* \in W_k$ in ζ_3 , and use Lemma 20 in ζ_4 . \square

Lemma 20 (Lemma 3.7 in Awasthi et al. (2017)). *Suppose Assumption 1 is satisfied. Then*

$$L_{\tau_k}(w^*) \leq \frac{\tau_k}{c_0 \min\{b_k, 1/9\}}.$$

In particular, by our choice of τ_k , it holds that

$$L_{\tau_k}(w^*) \leq \kappa.$$

Lemma 21. *For any $1 \leq k \leq K$, if $w^* \in W_k$, then with probability $1 - \delta_k$, $\theta(v_k, w^*) \leq 2^{-k-8}\pi$.*

Proof. For $k = 1$, by Lemma 16 with the facts that we actually sample from D and $w^* \in \mathbb{R}^d =: W_1$, we immediately have

$$\Pr_{x \sim D} (\text{sign}(v_1 \cdot x) \neq \text{sign}(w^* \cdot x)) \leq 8\kappa.$$

Hence Part 4 of Lemma 24 indicates that

$$\theta(v_1, w^*) \leq 8c_2\kappa = 16c_2\kappa \cdot 2^{-1}. \quad (24)$$

Now we consider $2 \leq k \leq K$. Denote $X_k = \{x : |w_{k-1} \cdot x| \leq b_k\}$, and $\bar{X}_k = \{x : |w_{k-1} \cdot x| > b_k\}$. We will show that the error of v_k on both X_k and \bar{X}_k is small, hence v_k is a good approximation to w^* .

First, we consider the error on X_k , which is given by

$$\begin{aligned} & \Pr_{x \sim D} (\text{sign}(v_k \cdot x) \neq \text{sign}(w^* \cdot x), x \in X_k) \\ &= \Pr_{x \sim D} (\text{sign}(v_k \cdot x) \neq \text{sign}(w^* \cdot x) \mid x \in X_k) \cdot \Pr_{x \sim D}(x \in X_k) \\ &= \text{err}_{D_{w_{k-1}, b_k}}(v_k) \cdot \Pr_{x \sim D}(x \in X_k) \\ &\leq 8\kappa \cdot 2b_k = 16\kappa b_k, \end{aligned} \quad (25)$$

where the inequality is due to Lemma 16 and Lemma 24. Note that the inequality holds with probability $1 - \delta_k$ in view of Lemma 16.

Next we derive the error on \bar{X}_k . Note that Lemma 10 of Zhang (2018) states for any unit vector u , and any general vector v , $\theta(v, u) \leq \pi \|v - u\|_2$. Hence,

$$\theta(v_k, w^*) \leq \pi \|v_k - w^*\|_2 \leq \pi(\|v_k - w_{k-1}\|_2 + \|w^* - w_{k-1}\|_2) \leq 2\pi r_k,$$

where we use the condition that both v_k and w^* are in W_k .

Recall that we set $r_k = 2^{-k-6} < 1/4$ in our algorithm and choose $b_k = \bar{c} \cdot r_k$ where $\bar{c} \geq 8\pi/c_4$, which allows us to apply Lemma 25 and obtain

$$\begin{aligned} \Pr_{x \sim D} (\text{sign}(v_k \cdot x) \neq \text{sign}(w^* \cdot x), x \notin X_k) &\leq c_3 \cdot 2\pi r_k \cdot \exp\left(-\frac{c_4 \bar{c} \cdot r_k}{2 \cdot 2\pi r_k}\right) \\ &= 2^{-k} \cdot \frac{c_3 \pi}{4} \exp\left(-\frac{c_4 \bar{c}}{4\pi}\right). \end{aligned}$$

This in allusion to (25) gives

$$\text{err}_D(v_k) \leq 16\kappa \cdot \bar{c} \cdot r_k + 2^{-k} \cdot \frac{c_3 \pi}{4} \exp\left(-\frac{c_4 \bar{c}}{4\pi}\right) = \left(2\kappa \bar{c} + \frac{c_3 \pi}{4} \exp\left(-\frac{c_4 \bar{c}}{4\pi}\right)\right) \cdot 2^{-k}.$$

Recall that we set $\kappa = \exp(-\bar{c})$. For convenience denote by $f(\bar{c})$ the coefficient of 2^{-k} in the above expression. By Part 4 of Lemma 24

$$\theta(v_k, w^*) \leq c_2 \text{err}_D(v_k) \leq c_2 f(\bar{c}) \cdot 2^{-k}. \quad (26)$$

Now let $g(\bar{c}) = c_2 f(\bar{c}) + 16c_2 \exp(-\bar{c})$. By our choice of \bar{c} , $g(\bar{c}) \leq 2^{-8}\pi$. This ensures that for both (24) and (26), $\theta(v_k, w^*) \leq 2^{-k-8}\pi$ for any $k \geq 1$. \square

Lemma 22. For any $1 \leq k \leq K$, if $\theta(v_k, w^*) \leq 2^{-k-8}\pi$, then $w^* \in W_{k+1}$.

Proof. We only need to show that $\|w_k - w^*\|_2 \leq r_{k+1}$. Let $\hat{v}_k = v_k / \|v_k\|_2$. By algebra $\|\hat{v}_k - w^*\|_2 = 2 \sin \frac{\theta(v_k, w^*)}{2} \leq \theta(v_k, w^*) \leq 2^{-k-8}\pi \leq 2^{-k-6}$. Now we have

$$\|w_k - w^*\|_2 = \|\hat{v}_k - w^*\|_2 \leq 2^{-k-6} = r_{k+1}.$$

The proof is complete. \square

C.7. Proof of Theorem 2

Proof. We will prove the theorem with the following claim.

Claim 23. For any $1 \leq k \leq K$, with probability at least $1 - \sum_{i=1}^k \delta_i$, w^* is in W_{k+1} .

Based on the claim, we immediately have that with probability at least $1 - \sum_{k=1}^K \delta_k \geq 1 - \delta$, w^* is in W_{K+1} . By our construction of W_{K+1} , we have

$$\|w^* - w_K\|_2 \leq 2^{-K-5}.$$

This, together with Part 4 of Lemma 24 and the fact that $\theta(w^*, w_K) \leq \pi \|w^* - w_K\|_2$ (see Lemma 10 of Zhang (2018)), implies

$$\text{err}_D(w_K) \leq \frac{\pi}{c_1} \cdot 2^{-K-5} = \epsilon.$$

The sample complexity of the algorithm is given by

$$N := \sum_{k=1}^K N_k = \sum_{k=1}^K \frac{d}{b_k} \cdot \text{polylog}\left(d, \frac{1}{b_k}, \frac{1}{\delta_k}\right) \leq \frac{d}{\epsilon} \cdot \text{polylog}\left(d, \frac{1}{\epsilon}, \frac{1}{\delta}\right),$$

where we use the fact that $b_k \geq K_1 \epsilon$ for some constant $K_1 > 0$ and $K = O(\log \frac{1}{\epsilon})$.

For each phase $k \leq K$, the number of calls to EX^y equals the size of T . For the size of T_C , by Lemma 24 we know that the probability mass of the band $X_k = \{x : |w_{k-1} \cdot x| \leq b_k\}$ is at most $2b_k$, implying that $|T_C| \leq O(b_k N_k)$ with high probability in view of Chernoff bound. On the other hand, by Part 2 of Lemma 12 we have $|T_D| = |T_E| \leq O(b_k N_k)$ since $\xi_k = \Theta(1)$ as indicated in Section 3.1.1. Therefore, $|T| \leq O(b_k N_k)$ and the label complexity m of the algorithm is given by

$$m \leq \sum_{k=1}^K b_k N_k = d \cdot \text{polylog} \left(d, \frac{1}{\epsilon}, \frac{1}{\delta} \right).$$

It remains to prove Claim 23 by induction. First, for $k = 1$, $W_1 = \{w : \|w\|_2 \leq 1\}$. Therefore, $w^* \in W_1$ with probability 1. Now suppose that Claim 23 holds for some $k \geq 2$, that is, there is an event E_{k-1} that happens with probability $1 - \sum_{i=1}^{k-1} \delta_i$, and on this event $w^* \in W_k$. By Lemma 21 we know that there is an event F_k that happens with probability $1 - \delta_k$, on which $\theta(v_k, w^*) \leq 2^{-k-8}\pi$. This further implies that $w^* \in W_{k+1}$ in view of Lemma 22. Therefore, consider the event $E_{k-1} \cap F_k$, on which $w^* \in W_{k+1}$ with probability $\Pr(E_{k-1}) \cdot \Pr(F_k | E_{k-1}) = (1 - \sum_{i=1}^{k-1} \delta_i)(1 - \delta_k) \geq 1 - \sum_{i=1}^k \delta_i$. \square

D. Properties of Isotropic Log-Concave Distributions

We record some useful properties of isotropic log-concave distributions.

Lemma 24. *There are absolute constants $c_0, c_1, c_2 > 0$, such that the following holds for all isotropic log-concave distributions $D \in \mathcal{D}$. Let f_D be the density function. We have*

1. *Orthogonal projections of D onto subspaces of \mathbb{R}^d are isotropic log-concave;*
2. *If $d = 1$, then $\Pr_{x \sim D}(a \leq x \leq b) \leq |b - a|$;*
3. *If $d = 1$, then $f_D(x) \geq c_0$ for all $x \in [-1/9, 1/9]$;*
4. *For any two vectors $u, v \in \mathbb{R}^d$,*

$$c_1 \cdot \Pr_{x \sim D}(\text{sign}(u \cdot x) \neq \text{sign}(v \cdot x)) \leq \theta(u, v) \leq c_2 \cdot \Pr_{x \sim D}(\text{sign}(u \cdot x) \neq \text{sign}(v \cdot x));$$

5. *$\Pr_{x \sim D}(\|x\|_2 \geq t\sqrt{d}) \leq \exp(-t + 1)$.*

We remark that Parts 1, 2, 3, and 5 are due to Lovász & Vempala (2007), and Part 4 is from Vempala (2010); Balcan & Long (2013).

The following lemma is implied by the proof of Theorem 21 of Balcan & Long (2013), which shows that if we choose a proper band width $b > 0$, the error outside the band will be small. This observation is crucial for controlling the error over the distribution D , and has been broadly recognized in the literature (Awasthi et al., 2017; Zhang, 2018).

Lemma 25 (Theorem 21 of Balcan & Long (2013)). *There are absolute constants $c_3, c_4 > 0$ such that the following holds for all isotropic log-concave distributions $D \in \mathcal{D}$. Let u and v be two unit vectors in \mathbb{R}^d and assume that $\theta(u, v) = \alpha < \pi/2$. Then for any $b \geq \frac{4}{c_4}\alpha$, we have*

$$\Pr_{x \sim D}(\text{sign}(u \cdot x) \neq \text{sign}(v \cdot x) \text{ and } |v \cdot x| \geq b) \leq c_3 \alpha \exp\left(-\frac{c_4 b}{2\alpha}\right).$$

Lemma 26. *Suppose x is randomly drawn from $D_{u,b}$. Then with probability $1 - \delta$, $\|x\|_2 \leq c_7 \sqrt{d} \log \frac{1}{\delta}$ for some constant $c_7 > 0$.*

Proof. Using Part 2 of Lemma 27, we have

$$\Pr_{x \sim D_{u,b}}(\|x\|_2 \geq \alpha) \leq \frac{1}{c_8 b} \Pr_{x \sim D}(\|x\|_2 \geq \alpha) \leq \frac{e}{c_8 b} \exp(-\alpha/\sqrt{d}),$$

where we applied Part 5 of Lemma 24 in the last inequality. The lemma follows by setting the right-hand side to δ . \square

Lemma 27. Let $c_8 = \min \left\{ 2c_0, \frac{2c_0}{9C_1}, \frac{1}{C_1} \right\}$. Then for all isotropic log-concave distributions $D \in \mathcal{D}$,

1. $\Pr_{x \sim D} (|u \cdot x| \leq b) \geq c_8 \cdot b$;
2. $\Pr_{x \sim D_{u,b}}(E) \leq \frac{1}{c_8 b} \Pr_{x \sim D}(E)$ for any event E .

Proof. We first consider the case that u is a unit vector.

For the lower bound, Part 3 of Lemma 24 shows that the density function of the random variable $u \cdot x$ is lower bounded by c_0 when $|u \cdot x| \leq 1/9$. Thus

$$\Pr_{x \sim D} (|u \cdot x| \leq b) \geq \Pr_{x \sim D} (|u \cdot x| \leq \min\{b, 1/9\}) \geq 2c_0 \min\{b, 1/9\} \geq 2c_0 \min \left\{ 1, \frac{1}{9C_1} \right\} \cdot b$$

where in the last inequality we use the condition $b \leq C_1$.

For any event E , we always have

$$\Pr_{x \sim D_{u,b}}(E) \leq \frac{\Pr_{x \sim D}(E)}{\Pr_{x \sim D}(|u \cdot x| \leq b)} \leq \frac{1}{c_8 b} \Pr_{x \sim D}(E).$$

Now we consider the case that u is the zero vector and $b = C_1$. Then $\Pr_{x \sim D} (|u \cdot x| \leq b) = 1 \geq c_8 \cdot b$ in view of the choice c_8 . Thus Part 2 still follows. The proof is complete. \square