
Appendix of “Binary Classification with Confidence Difference”

Anonymous Author(s)

Affiliation

Address

email

1 A Proof of Theorem 1

2 Before giving the proof of Theorem 1, we begin with the following lemmas:

3 **Lemma 2.** *The confidence difference $c(\mathbf{x}, \mathbf{x}')$ can be equivalently expressed as*

$$c(\mathbf{x}, \mathbf{x}') = \frac{\pi_+ p(\mathbf{x}) p_+(\mathbf{x}') - \pi_+ p_+(\mathbf{x}) p(\mathbf{x}')}{p(\mathbf{x}) p(\mathbf{x}')} \quad (1)$$

$$= \frac{\pi_- p_-(\mathbf{x}) p(\mathbf{x}') - \pi_- p(\mathbf{x}) p_-(\mathbf{x}')}{p(\mathbf{x}) p(\mathbf{x}')} \quad (2)$$

4 *Proof.* On one hand,

$$\begin{aligned} c(\mathbf{x}, \mathbf{x}') &= p(y' = 1 | \mathbf{x}') - p(y = 1 | \mathbf{x}) \\ &= \frac{p(\mathbf{x}', y' = 1)}{p(\mathbf{x}')} - \frac{p(\mathbf{x}, y = 1)}{p(\mathbf{x})} \\ &= \frac{\pi_+ p_+(\mathbf{x}')}{p(\mathbf{x}')} - \frac{\pi_+ p_+(\mathbf{x})}{p(\mathbf{x})} \\ &= \frac{\pi_+ p(\mathbf{x}) p_+(\mathbf{x}') - \pi_+ p_+(\mathbf{x}) p(\mathbf{x}')}{p(\mathbf{x}) p(\mathbf{x}')} \end{aligned}$$

5 On the other hand,

$$\begin{aligned} c(\mathbf{x}, \mathbf{x}') &= p(y' = 1 | \mathbf{x}') - p(y = 1 | \mathbf{x}) \\ &= (1 - p(y' = 0 | \mathbf{x}')) - (1 - p(y = 0 | \mathbf{x})) \\ &= p(y = 0 | \mathbf{x}) - p(y' = 0 | \mathbf{x}') \\ &= \frac{p(\mathbf{x}, y = 0)}{p(\mathbf{x})} - \frac{p(\mathbf{x}', y = 0)}{p(\mathbf{x}')} \\ &= \frac{\pi_- p_-(\mathbf{x})}{p(\mathbf{x})} - \frac{\pi_- p_-(\mathbf{x}')}{p(\mathbf{x}')} \\ &= \frac{\pi_- p_-(\mathbf{x}) p(\mathbf{x}') - \pi_- p(\mathbf{x}) p_-(\mathbf{x}')}{p(\mathbf{x}) p(\mathbf{x}')} \end{aligned}$$

6 which concludes the proof. □

7 **Lemma 3.** *The following equations hold:*

$$\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\pi_+ - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), +1)] = \pi_+ \mathbb{E}_{p_+(\mathbf{x})}[\ell(g(\mathbf{x}), +1)], \quad (3)$$

$$\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\pi_- + c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), -1)] = \pi_- \mathbb{E}_{p_-(\mathbf{x})}[\ell(g(\mathbf{x}), -1)], \quad (4)$$

$$\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\pi_+ + c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}'), +1)] = \pi_+ \mathbb{E}_{p_+(\mathbf{x}')}[\ell(g(\mathbf{x}'), +1)], \quad (5)$$

$$\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\pi_- - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}'), -1)] = \pi_- \mathbb{E}_{p_-(\mathbf{x}')}[\ell(g(\mathbf{x}'), -1)]. \quad (6)$$

8 *Proof.* Firstly, the proof of Eq. (3) is given:

$$\begin{aligned}
& \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\pi_+ - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), +1)] \\
&= \int \int \frac{\pi_+ p(\mathbf{x}) p(\mathbf{x}') - \pi_+ p(\mathbf{x}) p_+(\mathbf{x}') + \pi_+ p_+(\mathbf{x}) p(\mathbf{x}')}{p(\mathbf{x}) p(\mathbf{x}')} \ell(g(\mathbf{x}), +1) p(\mathbf{x}, \mathbf{x}') d\mathbf{x} d\mathbf{x}' \\
&= \int \int (\pi_+ p(\mathbf{x}) p(\mathbf{x}') - \pi_+ p(\mathbf{x}) p_+(\mathbf{x}') + \pi_+ p_+(\mathbf{x}) p(\mathbf{x}')) \ell(g(\mathbf{x}), +1) d\mathbf{x} d\mathbf{x}' \\
&= \int \pi_+ p(\mathbf{x}) \ell(g(\mathbf{x}), +1) d\mathbf{x} \int p(\mathbf{x}') d\mathbf{x}' - \int \pi_+ p(\mathbf{x}) \ell(g(\mathbf{x}), +1) d\mathbf{x} \int p_+(\mathbf{x}') d\mathbf{x}' \\
&\quad + \int \pi_+ p_+(\mathbf{x}) \ell(g(\mathbf{x}), +1) d\mathbf{x} \int p(\mathbf{x}') d\mathbf{x}' \\
&= \int \pi_+ p(\mathbf{x}) \ell(g(\mathbf{x}), +1) d\mathbf{x} - \int \pi_+ p(\mathbf{x}) \ell(g(\mathbf{x}), +1) d\mathbf{x} + \int \pi_+ p_+(\mathbf{x}) \ell(g(\mathbf{x}), +1) d\mathbf{x} \\
&= \int \pi_+ p_+(\mathbf{x}) \ell(g(\mathbf{x}), +1) d\mathbf{x} \\
&= \pi_+ \mathbb{E}_{p_+(\mathbf{x})}[\ell(g(\mathbf{x}), +1)].
\end{aligned}$$

9 After that, the proof of Eq. (4) is given:

$$\begin{aligned}
& \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\pi_- + c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), -1)] \\
&= \int \int \frac{\pi_- p(\mathbf{x}) p(\mathbf{x}') + \pi_- p_-(\mathbf{x}) p(\mathbf{x}') - \pi_- p(\mathbf{x}) p_-(\mathbf{x}')}{p(\mathbf{x}) p(\mathbf{x}')} \ell(g(\mathbf{x}), -1) p(\mathbf{x}, \mathbf{x}') d\mathbf{x} d\mathbf{x}' \\
&= \int \int (\pi_- p(\mathbf{x}) p(\mathbf{x}') + \pi_- p_-(\mathbf{x}) p(\mathbf{x}') - \pi_- p(\mathbf{x}) p_-(\mathbf{x}')) \ell(g(\mathbf{x}), -1) d\mathbf{x} d\mathbf{x}' \\
&= \int \pi_- p(\mathbf{x}) \ell(g(\mathbf{x}), -1) d\mathbf{x} \int p(\mathbf{x}') d\mathbf{x}' + \int \pi_- p_-(\mathbf{x}) \ell(g(\mathbf{x}), -1) d\mathbf{x} \int p(\mathbf{x}') d\mathbf{x}' \\
&\quad - \int \pi_- p(\mathbf{x}) \ell(g(\mathbf{x}), -1) d\mathbf{x} \int p_-(\mathbf{x}') d\mathbf{x}' \\
&= \int \pi_- p(\mathbf{x}) \ell(g(\mathbf{x}), -1) d\mathbf{x} + \int \pi_- p_-(\mathbf{x}) \ell(g(\mathbf{x}), -1) d\mathbf{x} - \int \pi_- p(\mathbf{x}) \ell(g(\mathbf{x}), -1) d\mathbf{x} \\
&= \int \pi_- p_-(\mathbf{x}) \ell(g(\mathbf{x}), -1) d\mathbf{x} \\
&= \pi_- \mathbb{E}_{p_-(\mathbf{x})}[\ell(g(\mathbf{x}), -1)].
\end{aligned}$$

10 It can be noticed that $c(\mathbf{x}, \mathbf{x}') = -c(\mathbf{x}', \mathbf{x})$ and $p(\mathbf{x}, \mathbf{x}') = p(\mathbf{x}', \mathbf{x})$. Therefore, it can be deduced
11 naturally that $\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\pi_+ - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), +1)] = \mathbb{E}_{p(\mathbf{x}', \mathbf{x})}[(\pi_+ + c(\mathbf{x}', \mathbf{x}))\ell(g(\mathbf{x}), +1)]$. Be-
12 cause \mathbf{x} and \mathbf{x}' are symmetric, we can swap them and deduce Eq. (5). Eq. (6) can be deduced in the
13 same manner, which concludes the proof. \square

14 Based on Lemma 3, the proof of Theorem 1 is given.

15 *Proof of Theorem 1.* To begin with, it can be noticed that $\mathbb{E}_{p_+(\mathbf{x})}[\ell(g(\mathbf{x}), +1)] =$
16 $\mathbb{E}_{p_+(\mathbf{x}')}[\ell(g(\mathbf{x}'), +1)]$ and $\mathbb{E}_{p_-(\mathbf{x})}[\ell(g(\mathbf{x}), -1)] = \mathbb{E}_{p_-(\mathbf{x}')}[\ell(g(\mathbf{x}'), -1)]$. Then, by summing up all
17 the equations from Eq. (3) to Eq. (6), we can get the following equation:

$$\begin{aligned}
& \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[\mathcal{L}_+(g(\mathbf{x}), g(\mathbf{x}')) + \mathcal{L}_-(g(\mathbf{x}), g(\mathbf{x}'))] \\
&= 2\pi_+ \mathbb{E}_{p_+(\mathbf{x})}[\ell(g(\mathbf{x}), +1)] + 2\pi_- \mathbb{E}_{p_-(\mathbf{x})}[\ell(g(\mathbf{x}), -1)]
\end{aligned}$$

18 After dividing each side of the equation above by 2, we can obtain Theorem 1. \square

19 B Analysis on Variance of Risk Estimator

20 B.1 Proof of Lemma 1 in the Main Paper

21 Based on Lemma 3, it can be observed that

$$\begin{aligned}\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[\mathcal{L}(\mathbf{x}, \mathbf{x}')] &= \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\pi_+ - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), +1) + (\pi_- - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}'), -1)] \\ &= \pi_+ \mathbb{E}_{p_+(\mathbf{x})}[\ell(g(\mathbf{x}), +1)] + \pi_- \mathbb{E}_{p_-(\mathbf{x}')}[\ell(g(\mathbf{x}'), -1)] \\ &= \pi_+ \mathbb{E}_{p_+(\mathbf{x})}[\ell(g(\mathbf{x}), +1)] + \pi_- \mathbb{E}_{p_-(\mathbf{x})}[\ell(g(\mathbf{x}), -1)] \\ &= R(g)\end{aligned}$$

22 and

$$\begin{aligned}\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[\mathcal{L}(\mathbf{x}', \mathbf{x})] &= \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\pi_+ + c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}'), +1) + (\pi_- + c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), -1)] \\ &= \pi_- \mathbb{E}_{p_-(\mathbf{x})}[\ell(g(\mathbf{x}), -1)] + \pi_+ \mathbb{E}_{p_+(\mathbf{x}')}[\ell(g(\mathbf{x}'), +1)] \\ &= \pi_- \mathbb{E}_{p_-(\mathbf{x})}[\ell(g(\mathbf{x}), -1)] + \pi_+ \mathbb{E}_{p_+(\mathbf{x})}[\ell(g(\mathbf{x}), +1)] \\ &= R(g).\end{aligned}$$

23 Therefore, for an arbitrary weight $\alpha \in [0, 1]$,

$$\begin{aligned}R(g) &= \alpha R(g) + (1 - \alpha)R(g) \\ &= \alpha \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[\mathcal{L}(\mathbf{x}, \mathbf{x}')] + (1 - \alpha) \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[\mathcal{L}(\mathbf{x}', \mathbf{x})],\end{aligned}$$

24 which indicates that

$$\frac{1}{n} \sum_{i=1}^n (\alpha \mathcal{L}(\mathbf{x}_i, \mathbf{x}'_i) + (1 - \alpha) \mathcal{L}(\mathbf{x}'_i, \mathbf{x}_i))$$

25 is also an unbiased risk estimator and concludes the proof. \square

26 B.2 Proof of Theorem 2

27 In this subsection, we show that Eq. (8) in the main paper achieves the minimum variance of

$$S(g; \alpha) = \frac{1}{n} \sum_{i=1}^n (\alpha \mathcal{L}(\mathbf{x}_i, \mathbf{x}'_i) + (1 - \alpha) \mathcal{L}(\mathbf{x}'_i, \mathbf{x}_i))$$

28 w.r.t. any $\alpha \in [0, 1]$. To begin with, we introduce the following notations:

$$\begin{aligned}\mu_1 &\triangleq \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \mathbf{x}'_i))^2] = \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}'_i, \mathbf{x}_i))^2], \\ \mu_2 &\triangleq \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[\frac{1}{n^2} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \mathbf{x}'_i) \sum_{i=1}^n \mathcal{L}(\mathbf{x}'_i, \mathbf{x}_i)].\end{aligned}\tag{7}$$

29 Furthermore, according to Lemma 1 in the main paper, we have

$$\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[S(g; \alpha)] = R(g).$$

30 Then, we provide the proof of Theorem 2 as follows.

Proof of Theorem 2.

$$\begin{aligned}\text{Var}(S(g; \alpha)) &= \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(S(g; \alpha) - R(g))^2] \\ &= \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[S(g; \alpha)^2] - R(g)^2 \\ &= \alpha^2 \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \mathbf{x}'_i))^2] + (1 - \alpha)^2 \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}'_i, \mathbf{x}_i))^2] \\ &\quad + 2\alpha(1 - \alpha) \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[\frac{1}{n^2} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \mathbf{x}'_i) \sum_{i=1}^n \mathcal{L}(\mathbf{x}'_i, \mathbf{x}_i)] - R(g)^2 \\ &= \mu_1 \alpha^2 + \mu_1 (1 - \alpha)^2 + 2\mu_2 \alpha(1 - \alpha) - R(g)^2 \\ &= (2\mu_1 - 2\mu_2)(\alpha - \frac{1}{2})^2 + \frac{1}{2}(\mu_1 + \mu_2) - R(g)^2.\end{aligned}$$

31 Besides, it can be observed that

$$2\mu_1 - 2\mu_2 = \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[\left(\frac{1}{n} \sum_{i=1}^n (\mathcal{L}(\mathbf{x}_i, \mathbf{x}'_i) - \mathcal{L}(\mathbf{x}'_i, \mathbf{x}_i))^2 \right) \right] \geq 0.$$

32 Therefore, $\text{Var}(S(g; \alpha))$ achieves the minimum value when $\alpha = 1/2$, which concludes the proof. \square

33 C Proof of Theorem 3

34 To begin with, we give the definition of Rademacher complexity.

35 **Definition 1** (Rademacher complexity). Let $\mathcal{X}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ denote n i.i.d. random variables
 36 drawn from a probability distribution with density $p(\mathbf{x})$, $\mathcal{G} = \{g : \mathcal{X} \mapsto \mathbb{R}\}$ denote a class of
 37 measurable functions, and $\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \dots, \sigma_n)$ denote Rademacher variables taking values from
 38 $\{+1, -1\}$ uniformly. Then, the (expected) Rademacher complexity of \mathcal{G} is defined as

$$\mathfrak{R}_n(\mathcal{G}) = \mathbb{E}_{\mathcal{X}_n} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{x}_i) \right]. \quad (8)$$

39 Let $\mathcal{D}_n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, \mathbf{x}')$ denote n pairs of ConfDiff data and $\mathcal{L}_{\text{CD}}(g; \mathbf{x}_i, \mathbf{x}'_i) = (\mathcal{L}(\mathbf{x}, \mathbf{x}') + \mathcal{L}(\mathbf{x}', \mathbf{x}))/2$,
 40 then we introduce the following lemma.

Lemma 4.

$$\bar{\mathfrak{R}}_n(\mathcal{L}_{\text{CD}} \circ \mathcal{G}) \leq 2L_\ell \mathfrak{R}_n(\mathcal{G}),$$

41 where $\mathcal{L}_{\text{CD}} \circ \mathcal{G} = \{\mathcal{L}_{\text{CD}} \circ g | g \in \mathcal{G}\}$ and $\bar{\mathfrak{R}}_n(\cdot)$ is the Rademacher complexity over ConfDiff data
 42 pairs \mathcal{D}_n of size n .

Proof.

$$\begin{aligned} \bar{\mathfrak{R}}_n(\mathcal{L}_{\text{CD}} \circ \mathcal{G}) &= \mathbb{E}_{\mathcal{D}_n} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathcal{L}_{\text{CD}}(g; \mathbf{x}_i, \mathbf{x}'_i) \right] \\ &= \mathbb{E}_{\mathcal{D}_n} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{g \in \mathcal{G}} \frac{1}{2n} \sum_{i=1}^n \sigma_i ((\pi_+ - c_i)\ell(g(\mathbf{x}_i), +1) + (\pi_- - c_i)\ell(g(\mathbf{x}'_i), -1) \right. \\ &\quad \left. + (\pi_+ + c_i)\ell(g(\mathbf{x}'_i), +1) + (\pi_- + c_i)\ell(g(\mathbf{x}_i), -1)) \right]. \end{aligned}$$

43 Then, we can induce that

$$\begin{aligned} &\|\nabla \mathcal{L}_{\text{CD}}(g; \mathbf{x}_i, \mathbf{x}'_i)\|_2 \\ &= \left\| \nabla \left(\frac{(\pi_+ - c_i)\ell(g(\mathbf{x}_i), +1) + (\pi_- - c_i)\ell(g(\mathbf{x}'_i), -1)}{2} \right. \right. \\ &\quad \left. \left. + \frac{(\pi_+ + c_i)\ell(g(\mathbf{x}'_i), +1) + (\pi_- + c_i)\ell(g(\mathbf{x}_i), -1)}{2} \right) \right\|_2 \\ &\leq \left\| \nabla \left(\frac{(\pi_+ - c_i)\ell(g(\mathbf{x}_i), +1)}{2} \right) \right\|_2 + \left\| \nabla \left(\frac{(\pi_- - c_i)\ell(g(\mathbf{x}'_i), -1)}{2} \right) \right\|_2 \\ &\quad + \left\| \nabla \left(\frac{(\pi_+ + c_i)\ell(g(\mathbf{x}'_i), +1)}{2} \right) \right\|_2 + \left\| \nabla \left(\frac{(\pi_- + c_i)\ell(g(\mathbf{x}_i), -1)}{2} \right) \right\|_2 \\ &\leq \frac{|\pi_+ - c_i|L_\ell}{2} + \frac{|\pi_- - c_i|L_\ell}{2} + \frac{|\pi_+ + c_i|L_\ell}{2} + \frac{|\pi_- + c_i|L_\ell}{2}. \end{aligned} \quad (9)$$

44 Suppose $\pi_+ \geq \pi_-$, the value of RHS of Eq. (9) can be determined as follows: when $c_i \in [-1, -\pi_+)$,
 45 the value is $-2c_iL_\ell$; when $c_i \in [-\pi_+, -\pi_-)$, the value is $(\pi_+ - c_i)L_\ell$; when $c_i \in [-\pi_-, \pi_-)$, the
 46 value is L_ℓ ; when $c_i \in [\pi_-, \pi_+)$, the value is $(\pi_+ + c_i)L_\ell$; when $c_i \in [\pi_+, 1]$, the value is $2c_iL_\ell$.
 47 To sum up, when $\pi_+ \geq \pi_-$, the value of RHS of Eq. (9) is less than $2L_\ell$. When $\pi_+ \leq \pi_-$, we can

deduce that the value of RHS of Eq. (9) is less than $2L_\ell$ in the same way. Therefore,

$$\begin{aligned}\bar{\mathfrak{R}}_n(\mathcal{L}_{\text{CD}} \circ \mathcal{G}) &\leq 2L_\ell \mathbb{E}_{\mathcal{D}_n} \mathbb{E}_{\sigma} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{x}_i) \right] \\ &= 2L_\ell \mathbb{E}_{\mathcal{X}_n} \mathbb{E}_{\sigma} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{x}_i) \right] \\ &= 2L_\ell \mathfrak{R}_n(\mathcal{G}),\end{aligned}$$

which concludes the proof. \square

After that, we introduce the following lemma.

Lemma 5. *The inequality below hold with probability at least $1 - \delta$:*

$$\sup_{g \in \mathcal{G}} |R(g) - \hat{R}_{\text{CD}}(g)| \leq 4L_\ell \mathfrak{R}_n(\mathcal{G}) + 2C_\ell \sqrt{\frac{\ln 2/\delta}{2n}}.$$

Proof. To begin with, we introduce $\Phi = \sup_{g \in \mathcal{G}} (R(g) - \hat{R}_{\text{CD}}(g))$ and $\bar{\Phi} = \sup_{g \in \mathcal{G}} (R(g) - \hat{\bar{R}}_{\text{CD}}(g))$, where $\hat{R}_{\text{CD}}(g)$ and $\hat{\bar{R}}_{\text{CD}}(g)$ denote the empirical risk over two sets of training examples with exactly one different point $\{(\mathbf{x}_i, \mathbf{x}'_i), c_i\}$ and $\{(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i), c(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i)\}$ respectively. Then we have

$$\begin{aligned}\bar{\Phi} - \Phi &\leq \sup_{g \in \mathcal{G}} (\hat{R}_{\text{CD}}(g) - \hat{\bar{R}}_{\text{CD}}(g)) \\ &\leq \sup_{g \in \mathcal{G}} \left(\frac{\mathcal{L}_{\text{CD}}(g; \mathbf{x}_i, \mathbf{x}'_i) - \mathcal{L}_{\text{CD}}(g; \bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i)}{n} \right) \\ &\leq \frac{2C_\ell}{n}.\end{aligned}$$

Accordingly, $\Phi - \bar{\Phi}$ can be bounded in the same way. The following inequalities holds with probability at least $1 - \delta/2$ by applying McDiarmid's inequality:

$$\sup_{g \in \mathcal{G}} (R(g) - \hat{R}_{\text{CD}}(g)) \leq \mathbb{E}_{\mathcal{D}_n} [\sup_{g \in \mathcal{G}} (R(g) - \hat{R}_{\text{CD}}(g))] + 2C_\ell \sqrt{\frac{\ln 2/\delta}{2n}},$$

Furthermore, we can bound $\mathbb{E}_{\mathcal{D}_n} [\sup_{g \in \mathcal{G}} (R(g) - \hat{R}_{\text{CD}}(g))]$ with Rademacher complexity. It is a routine work to show by symmetrization [16] that

$$\mathbb{E}_{\mathcal{D}_n} [\sup_{g \in \mathcal{G}} (R(g) - \hat{R}_{\text{CD}}(g))] \leq 2\bar{\mathfrak{R}}_n(\mathcal{L}_{\text{CD}} \circ \mathcal{G}) \leq 4L_\ell \mathfrak{R}_n(\mathcal{G}),$$

where the second inequality is from Lemma 4. Accordingly, $\sup_{g \in \mathcal{G}} (\hat{R}_{\text{CD}}(g) - R(g))$ has the same bound. By using the union bound, the following inequality holds with probability at least $1 - \delta$:

$$\sup_{g \in \mathcal{G}} |R(g) - \hat{R}_{\text{CD}}(g)| \leq 4L_\ell \mathfrak{R}_n(\mathcal{G}) + 2C_\ell \sqrt{\frac{\ln 2/\delta}{2n}},$$

which concludes the proof. \square

Finally, the proof of Theorem 3 is provided.

Proof of Theorem 3.

$$\begin{aligned}R(\hat{g}_{\text{CD}}) - R(g^*) &= (R(\hat{g}_{\text{CD}}) - \hat{R}_{\text{CD}}(\hat{g}_{\text{CD}})) + (\hat{R}_{\text{CD}}(\hat{g}_{\text{CD}}) - \hat{R}_{\text{CD}}(g^*)) + (\hat{R}_{\text{CD}}(g^*) - R(g^*)) \\ &\leq (R(\hat{g}_{\text{CD}}) - \hat{R}_{\text{CD}}(\hat{g}_{\text{CD}})) + (\hat{R}_{\text{CD}}(g^*) - R(g^*)) \\ &\leq |R(\hat{g}_{\text{CD}}) - \hat{R}_{\text{CD}}(\hat{g}_{\text{CD}})| + |\hat{R}_{\text{CD}}(g^*) - R(g^*)| \\ &\leq 2 \sup_{g \in \mathcal{G}} |R(g) - \hat{R}_{\text{CD}}(g)| \\ &\leq 8L_\ell \mathfrak{R}_n(\mathcal{G}) + 4C_\ell \sqrt{\frac{\ln 2/\delta}{2n}}.\end{aligned}$$

63 The first inequality is derived because \widehat{g}_{CD} is the minimizer of $\widehat{R}_{\text{CD}}(g)$. The last inequality is derived
 64 according to Lemma 5, which concludes the proof. \square

65 D Proof of Theorem 4

66 To begin with, we provide the following inequality:

$$\begin{aligned}
 & \sup_{g \in \mathcal{G}} |\bar{R}_{\text{CD}}(g) - \widehat{R}_{\text{CD}}(g)| \\
 &= \frac{1}{2n} \left| \sum_{i=1}^n ((\bar{\pi}_+ - \pi_+ + c_i - \bar{c}_i)\ell(g(\mathbf{x}_i), +1) + (\bar{\pi}_- - \pi_- + c_i - \bar{c}_i)\ell(g(\mathbf{x}'_i), -1)) \right. \\
 & \quad \left. + (\bar{\pi}_+ - \pi_+ + \bar{c}_i - c_i)\ell(g(\mathbf{x}'_i), +1) + (\bar{\pi}_- - \pi_- + \bar{c}_i - c_i)\ell(g(\mathbf{x}_i), -1) \right| \\
 &\leq \frac{1}{2n} \sum_{i=1}^n (|(\bar{\pi}_+ - \pi_+ + c_i - \bar{c}_i)\ell(g(\mathbf{x}_i), +1)| + |(\bar{\pi}_- - \pi_- + c_i - \bar{c}_i)\ell(g(\mathbf{x}'_i), -1)| \\
 & \quad + |(\bar{\pi}_+ - \pi_+ + \bar{c}_i - c_i)\ell(g(\mathbf{x}'_i), +1)| + |(\bar{\pi}_- - \pi_- + \bar{c}_i - c_i)\ell(g(\mathbf{x}_i), -1)|) \\
 &= \frac{1}{2n} \sum_{i=1}^n (|\bar{\pi}_+ - \pi_+ + c_i - \bar{c}_i|\ell(g(\mathbf{x}_i), +1) + |\bar{\pi}_- - \pi_- + c_i - \bar{c}_i|\ell(g(\mathbf{x}'_i), -1) \\
 & \quad + |\bar{\pi}_+ - \pi_+ + \bar{c}_i - c_i|\ell(g(\mathbf{x}'_i), +1) + |\bar{\pi}_- - \pi_- + \bar{c}_i - c_i|\ell(g(\mathbf{x}_i), -1)) \\
 &\leq \frac{1}{2n} \sum_{i=1}^n ((|\bar{\pi}_+ - \pi_+| + |c_i - \bar{c}_i|)\ell(g(\mathbf{x}_i), +1) + (|\bar{\pi}_- - \pi_-| + |c_i - \bar{c}_i|)\ell(g(\mathbf{x}'_i), -1) \\
 & \quad + (|\bar{\pi}_+ - \pi_+| + |\bar{c}_i - c_i|)\ell(g(\mathbf{x}'_i), +1) + (|\bar{\pi}_- - \pi_-| + |\bar{c}_i - c_i|)\ell(g(\mathbf{x}_i), -1)) \\
 &= \frac{1}{2n} \sum_{i=1}^n ((|\bar{\pi}_+ - \pi_+| + |c_i - \bar{c}_i|)\ell(g(\mathbf{x}_i), +1) + (|\pi_+ - \bar{\pi}_+| + |c_i - \bar{c}_i|)\ell(g(\mathbf{x}'_i), -1) \\
 & \quad + (|\bar{\pi}_+ - \pi_+| + |\bar{c}_i - c_i|)\ell(g(\mathbf{x}'_i), +1) + (|\pi_- - \bar{\pi}_-| + |\bar{c}_i - c_i|)\ell(g(\mathbf{x}_i), -1)) \\
 &\leq \frac{2C_\ell \sum_{i=1}^n |\bar{c}_i - c_i|}{n} + 2C_\ell |\bar{\pi}_+ - \pi_+|.
 \end{aligned}$$

67 Then, we deduce the following inequality:

$$\begin{aligned}
 R(\bar{g}_{\text{CD}}) - R(g^*) &= (R(\bar{g}_{\text{CD}}) - \widehat{R}_{\text{CD}}(\bar{g}_{\text{CD}})) + (\widehat{R}_{\text{CD}}(\bar{g}_{\text{CD}}) - \bar{R}_{\text{CD}}(\bar{g}_{\text{CD}})) + (\bar{R}_{\text{CD}}(\bar{g}_{\text{CD}}) - \bar{R}_{\text{CD}}(\widehat{g}_{\text{CD}})) \\
 & \quad + (\bar{R}_{\text{CD}}(\widehat{g}_{\text{CD}}) - \widehat{R}_{\text{CD}}(\widehat{g}_{\text{CD}})) + (\widehat{R}_{\text{CD}}(\widehat{g}_{\text{CD}}) - R(\widehat{g}_{\text{CD}})) + (R(\widehat{g}_{\text{CD}}) - R(g^*)) \\
 &\leq 2 \sup_{g \in \mathcal{G}} |R(g) - \widehat{R}_{\text{CD}}(g)| + 2 \sup_{g \in \mathcal{G}} |\bar{R}_{\text{CD}}(g) - \widehat{R}_{\text{CD}}(g)| + (R(\widehat{g}_{\text{CD}}) - R(g^*)) \\
 &\leq 4 \sup_{g \in \mathcal{G}} |R(g) - \widehat{R}_{\text{CD}}(g)| + 2 \sup_{g \in \mathcal{G}} |\bar{R}_{\text{CD}}(g) - \widehat{R}_{\text{CD}}(g)| \\
 &\leq 16L_\ell \mathfrak{R}_n(\mathcal{G}) + 8C_\ell \sqrt{\frac{\ln 2/\delta}{2n}} + \frac{4C_\ell \sum_{i=1}^n |\bar{c}_i - c_i|}{n} + 4C_\ell |\bar{\pi}_+ - \pi_+|.
 \end{aligned}$$

68 The first inequality is derived because \bar{g}_{CD} is the minimizer of $\bar{R}(g)$. The second and third inequality
 69 are derived according to the proof of Theorem 3 and Lemma 5 respectively. \square

70 E Proof of Theorem 5

71 To begin with, let $\mathfrak{D}_n^+(g) = \{\mathcal{D}_n | \widehat{A}(g) \geq 0 \cap \widehat{B}(g) \geq 0 \cap \widehat{C}(g) \geq 0 \cap \widehat{D}(g) \geq 0\}$ and $\mathfrak{D}_n^-(g) =$
 72 $\{\mathcal{D}_n | \widehat{A}(g) \leq 0 \cup \widehat{B}(g) \leq 0 \cup \widehat{C}(g) \leq 0 \cup \widehat{D}(g) \leq 0\}$. Before giving the proof of Theorem 5, we
 73 give the following lemma based on the assumptions in Section 3.

74 **Lemma 6.** *The probability measure of $\mathfrak{D}_n^-(g)$ can be bounded as follows:*

$$\mathbb{P}(\mathfrak{D}_n^-(g)) \leq \exp\left(\frac{-2a^2n}{C_\ell^2}\right) + \exp\left(\frac{-2b^2n}{C_\ell^2}\right) + \exp\left(\frac{-2c^2n}{C_\ell^2}\right) + \exp\left(\frac{-2d^2n}{C_\ell^2}\right). \quad (10)$$

75 *Proof.* It can be observed that

$$\begin{aligned} p(\mathcal{D}_n) &= p(\mathbf{x}_1, \mathbf{x}'_1) \cdots p(\mathbf{x}_n, \mathbf{x}'_n) \\ &= p(\mathbf{x}_1) \cdots p(\mathbf{x}'_n) p(\mathbf{x}_1) \cdots p(\mathbf{x}'_n). \end{aligned}$$

76 Therefore, the probability measure $\mathbb{P}(\mathfrak{D}_n^-(g))$ can be defined as follows:

$$\begin{aligned} \mathbb{P}(\mathfrak{D}_n^-(g)) &= \int_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} p(\mathcal{D}_n) d\mathcal{D}_n \\ &= \int_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} p(\mathcal{D}_n) d\mathbf{x}_1 \cdots d\mathbf{x}_n d\mathbf{x}'_1 \cdots d\mathbf{x}'_n. \end{aligned}$$

77 When exactly one ConfDiff data pair in S_n is replaced, the change of $\hat{A}(g)$, $\hat{B}(g)$, $\hat{C}(g)$ and $\hat{D}(g)$
78 will be no more than C_ℓ/n . By applying McDiarmid's inequality, we can obtain the following
79 inequalities:

$$\begin{aligned} \mathbb{P}(\mathbb{E}[\hat{A}(g)] - \hat{A}(g) \geq a) &\leq \exp\left(-\frac{2a^2n}{C_\ell^2}\right), \\ \mathbb{P}(\mathbb{E}[\hat{B}(g)] - \hat{B}(g) \geq b) &\leq \exp\left(-\frac{2b^2n}{C_\ell^2}\right), \\ \mathbb{P}(\mathbb{E}[\hat{C}(g)] - \hat{C}(g) \geq c) &\leq \exp\left(-\frac{2c^2n}{C_\ell^2}\right), \\ \mathbb{P}(\mathbb{E}[\hat{D}(g)] - \hat{D}(g) \geq d) &\leq \exp\left(-\frac{2d^2n}{C_\ell^2}\right). \end{aligned}$$

80 Furthermore,

$$\begin{aligned} \mathbb{P}(\mathfrak{D}_n^-(g)) &\leq \mathbb{P}(\hat{A}(g) \leq 0) + \mathbb{P}(\hat{B}(g) \leq 0) + \mathbb{P}(\hat{C}(g) \leq 0) + \mathbb{P}(\hat{D}(g) \leq 0) \\ &\leq \mathbb{P}(\hat{A}(g) \leq \mathbb{E}[\hat{A}(g)] - a) + \mathbb{P}(\hat{B}(g) \leq \mathbb{E}[\hat{B}(g)] - b) \\ &\quad + \mathbb{P}(\hat{C}(g) \leq \mathbb{E}[\hat{C}(g)] - c) + \mathbb{P}(\hat{D}(g) \leq \mathbb{E}[\hat{D}(g)] - d) \\ &= \mathbb{P}(\mathbb{E}[\hat{A}(g)] - \hat{A}(g) \geq a) + \mathbb{P}(\mathbb{E}[\hat{B}(g)] - \hat{B}(g) \geq b) \\ &\quad + \mathbb{P}(\mathbb{E}[\hat{C}(g)] - \hat{C}(g) \geq c) + \mathbb{P}(\mathbb{E}[\hat{D}(g)] - \hat{D}(g) \geq d) \\ &\leq \exp\left(-\frac{2a^2n}{C_\ell^2}\right) + \exp\left(-\frac{2b^2n}{C_\ell^2}\right) + \exp\left(-\frac{2c^2n}{C_\ell^2}\right) + \exp\left(-\frac{2d^2n}{C_\ell^2}\right), \end{aligned}$$

81 which concludes the proof. □

82 Then, the proof of Theorem 5 is given.

83 *Proof of Theorem 5.* To begin with, we prove the first inequality in Theorem 5.

$$\begin{aligned} &\mathbb{E}[\tilde{R}_{\text{CD}}(g)] - R(g) \\ &= \mathbb{E}[\tilde{R}_{\text{CD}}(g) - \hat{R}_{\text{CD}}(g)] \\ &= \int_{\mathcal{D}_n \in \mathfrak{D}_n^+(g)} (\tilde{R}_{\text{CD}}(g) - \hat{R}_{\text{CD}}(g)) p(\mathcal{D}_n) d\mathcal{D}_n \\ &\quad + \int_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} (\tilde{R}_{\text{CD}}(g) - \hat{R}_{\text{CD}}(g)) p(\mathcal{D}_n) d\mathcal{D}_n \\ &= \int_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} (\tilde{R}_{\text{CD}}(g) - \hat{R}_{\text{CD}}(g)) p(\mathcal{D}_n) d\mathcal{D}_n \geq 0, \end{aligned}$$

84 where the last inequality is derived because $\tilde{R}_{\text{CD}}(g)$ is an upper bound of $\hat{R}_{\text{CD}}(g)$. Furthermore,

$$\begin{aligned}
& \mathbb{E}[\tilde{R}_{\text{CD}}(g)] - R(g) \\
&= \int_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} (\tilde{R}_{\text{CD}}(g) - \hat{R}_{\text{CD}}(g)) p(\mathcal{D}_n) d\mathcal{D}_n \\
&\leq \sup_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} (\tilde{R}_{\text{CD}}(g) - \hat{R}_{\text{CD}}(g)) \int_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} p(\mathcal{D}_n) d\mathcal{D}_n \\
&= \sup_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} (\tilde{R}_{\text{CD}}(g) - \hat{R}_{\text{CD}}(g)) \mathbb{P}(\mathfrak{D}_n^-(g)) \\
&= \sup_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} (f(\hat{A}(g)) + f(\hat{B}(g)) + f(\hat{C}(g)) + f(\hat{D}(g)) \\
&\quad - \hat{A}(g) - \hat{B}(g) - \hat{C}(g) - \hat{D}(g)) \mathbb{P}(\mathfrak{D}_n^-(g)) \\
&\leq \sup_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} (L_f |\hat{A}(g)| + L_f |\hat{B}(g)| + L_f |\hat{C}(g)| + L_f |\hat{D}(g)| \\
&\quad + |\hat{A}(g)| + |\hat{B}(g)| + |\hat{C}(g)| + |\hat{D}(g)|) \mathbb{P}(\mathfrak{D}_n^-(g)) \\
&= \sup_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} \frac{L_f + 1}{2n} (|\sum_{i=1}^n (\pi_+ - c_i) \ell(g(\mathbf{x}_i), +1)| + |\sum_{i=1}^n (\pi_- - c_i) \ell(g(\mathbf{x}'_i), -1)| \\
&\quad + |\sum_{i=1}^n (\pi_+ + c_i) \ell(g(\mathbf{x}'_i), +1)| + |\sum_{i=1}^n (\pi_- + c_i) \ell(g(\mathbf{x}_i), -1)|) \mathbb{P}(\mathfrak{D}_n^-(g)) \\
&\leq \sup_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} \frac{L_f + 1}{2n} (\sum_{i=1}^n |(\pi_+ - c_i) \ell(g(\mathbf{x}_i), +1)| + \sum_{i=1}^n |(\pi_- - c_i) \ell(g(\mathbf{x}'_i), -1)| \\
&\quad + \sum_{i=1}^n |(\pi_+ + c_i) \ell(g(\mathbf{x}'_i), +1)| + \sum_{i=1}^n |(\pi_- + c_i) \ell(g(\mathbf{x}_i), -1)|) \mathbb{P}(\mathfrak{D}_n^-(g)) \\
&= \sup_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} \frac{L_f + 1}{2n} \sum_{i=1}^n (|(\pi_+ - c_i) \ell(g(\mathbf{x}_i), +1)| + |(\pi_- - c_i) \ell(g(\mathbf{x}'_i), -1)| \\
&\quad + |(\pi_+ + c_i) \ell(g(\mathbf{x}'_i), +1)| + |(\pi_- + c_i) \ell(g(\mathbf{x}_i), -1)|) \mathbb{P}(\mathfrak{D}_n^-(g)) \\
&\leq \sup_{\mathcal{D}_n \in \mathfrak{D}_n^-(g)} \frac{(L_f + 1)C_\ell}{2n} \sum_{i=1}^n (|\pi_+ - c_i| + |\pi_- - c_i| + |\pi_+ + c_i| + |\pi_- + c_i|) \mathbb{P}(\mathfrak{D}_n^-(g)).
\end{aligned}$$

85 Similar to the proof of Theorem 3, we can obtain

$$|\pi_+ - c_i| + |\pi_- - c_i| + |\pi_+ + c_i| + |\pi_- + c_i| \leq 4.$$

86 Therefore, we have

$$\mathbb{E}[\tilde{R}_{\text{CD}}(g)] - R(g) \leq 2(L_f + 1)C_\ell \Delta,$$

87 which concludes the proof of the first inequality in Theorem 5. Before giving the proof of the second
88 inequality, we give the upper bound of $|\tilde{R}_{\text{CD}}(g) - \mathbb{E}[\tilde{R}_{\text{CD}}(g)]|$. When exactly one ConfDiff data
89 pair in \mathcal{D}_n is replaced, the change of $\tilde{R}_{\text{CD}}(g)$ is no more than $2C_\ell L_f/n$. By applying McDiarmid's
90 inequality, we have the following inequalities with probability at least $1 - \delta/2$:

$$\begin{aligned}
\tilde{R}_{\text{CD}}(g) - \mathbb{E}[\tilde{R}_{\text{CD}}(g)] &\leq 2C_\ell L_f \sqrt{\frac{\ln 2/\delta}{2n}}, \\
\mathbb{E}[\tilde{R}_{\text{CD}}(g)] - \tilde{R}_{\text{CD}}(g) &\leq 2C_\ell L_f \sqrt{\frac{\ln 2/\delta}{2n}}.
\end{aligned}$$

91 Therefore, with probability at least $1 - \delta$, we have

$$|\tilde{R}_{\text{CD}}(g) - \mathbb{E}[\tilde{R}_{\text{CD}}(g)]| \leq 2C_\ell L_f \sqrt{\frac{\ln 2/\delta}{2n}}.$$

Table 1: Characteristics of experimental data sets.

Data Set	# Train	# Test	# Features	# Class Labels	Model
MNIST	60,000	10,000	784	10	MLP
Kuzushiji	60,000	10,000	784	10	MLP
Fashion	60,000	10,000	784	10	MLP
CIFAR-10	50,000	10,000	3,072	10	ResNet-34
Optdigits	4,495	1,125	62	10	MLP
USPS	7,437	1,861	256	10	MLP
Pendigits	8,793	2,199	16	10	MLP
Letter	16,000	4,000	16	26	MLP

Finally, we have

$$\begin{aligned}
|\tilde{R}_{\text{CD}}(g) - R(g)| &= |\tilde{R}_{\text{CD}}(g) - \mathbb{E}[\tilde{R}_{\text{CD}}(g)] + \mathbb{E}[\tilde{R}_{\text{CD}}(g)] - R(g)| \\
&\leq |\tilde{R}_{\text{CD}}(g) - \mathbb{E}[\tilde{R}_{\text{CD}}(g)]| + |\mathbb{E}[\tilde{R}_{\text{CD}}(g)] - R(g)| \\
&= |\tilde{R}_{\text{CD}}(g) - \mathbb{E}[\tilde{R}_{\text{CD}}(g)]| + \mathbb{E}[\tilde{R}_{\text{CD}}(g)] - R(g) \\
&\leq 2C_\ell L_f \sqrt{\frac{\ln 2/\delta}{2n}} + 2(L_f + 1)C_\ell \Delta,
\end{aligned} \tag{11}$$

with probability at least $1 - \delta$, which concludes the proof. \square

F Proof of Theorem 6

With probability at least $1 - \delta$, we have

$$\begin{aligned}
R(\tilde{g}_{\text{CD}}) - R(g^*) &= (R(\tilde{g}_{\text{CD}}) - \tilde{R}_{\text{CD}}(\tilde{g}_{\text{CD}})) + (\tilde{R}_{\text{CD}}(\tilde{g}_{\text{CD}}) - \tilde{R}_{\text{CD}}(\hat{g}_{\text{CD}})) \\
&\quad + (\tilde{R}_{\text{CD}}(\hat{g}_{\text{CD}}) - R(\hat{g}_{\text{CD}})) + (R(\hat{g}_{\text{CD}}) - R(g^*)) \\
&\leq |R(\tilde{g}_{\text{CD}}) - \tilde{R}_{\text{CD}}(\tilde{g}_{\text{CD}})| + |\tilde{R}_{\text{CD}}(\hat{g}_{\text{CD}}) - R(\hat{g}_{\text{CD}})| + (R(\hat{g}_{\text{CD}}) - R(g^*)) \\
&\leq 4C_\ell(L_f + 1)\sqrt{\frac{\ln 2/\delta}{2n}} + 4(L_f + 1)C_\ell \Delta + 8L_\ell \mathfrak{R}_n(\mathcal{G}).
\end{aligned}$$

The first inequality is derived because \tilde{g}_{CD} is the minimizer of $\tilde{R}_{\text{CD}}(g)$. The second inequality is derived from Theorem 5 and Theorem 3. The proof is completed. \square

G Related Work

Learning with pairwise comparisons has been investigated pervasively in the community [1, 2, 10, 11, 17, 19, 23], with applications in information retrieval [15], computer vision [6], regression [21, 22], crowdsourcing [3, 24], and graph learning [9]. It is noteworthy that there exist distinct differences between our work and previous works on learning with pairwise comparisons. Previous works have mainly tried to learn a ranking function that can rank candidate examples according to relevance or preference. In this paper, we try to learn a *pointwise binary classifier* by conducting empirical risk minimization under the binary classification setting.

H Limitations and Potential Negative Social Impacts

H.1 Limitations

This work focuses on binary classification problems. To generalize it to multi-class problems, we need to convert multi-class classification to a set of binary classification problems via the one-versus-rest or the one-versus-one strategies. In the future, developing methods directly handling multi-class classification problems is promising.

H.2 Potential Negative Social Impacts

This work is within the scope of weakly supervised learning, which aims to achieve comparable performance while reducing labeling costs. Therefore, when this technique is very effective and prevalent in society, the demand for data annotations may be reduced, leading to the increasing unemployment rate of data annotation workers.

I Additional Information about Experiments

In this section, the details of experimental data sets and hyperparameters are provided.

I.1 Details of Experimental Data Sets

The detailed statistics and corresponding model architectures are summarized in Table 1 while the basic information, sources and data split details are elaborated in this subsection.

For the four benchmark data sets,

- MNIST [14]: It is a grayscale handwritten digits recognition data set. It is composed of 60,000 training examples and 10,000 test examples. The original feature dimension is 28×28 , and the label space is 0-9. The even digits are regarded as the positive class while the odd digits are regarded as the negative class. We sampled 15,000 unlabeled data pairs as training data. The data set can be downloaded from <http://yann.lecun.com/exdb/mnist/>.
- Kuzushiji-MNIST [4]: It is a grayscale Japanese character recognition data set. It is composed of 60,000 training examples and 10,000 test examples. The original feature dimension is 28×28 , and the label space is {'o', 'su', 'na', 'ma', 're', 'ki', 'tsu', 'ha', 'ya', 'wo'}. The positive class is composed of 'o', 'su', 'na', 'ma', and 're' while the negative class is composed of 'ki', 'tsu', 'ha', 'ya', and 'wo'. We sampled 15,000 unlabeled data pairs as training data. The data set can be downloaded from <https://github.com/rois-codh/kmnist>.
- Fashion-MNIST [20]: It is a grayscale fashion item recognition data set. It is composed of 60,000 training examples and 10,000 test examples. The original feature dimension is 28×28 , and the label space is {'T-shirt', 'trouser', 'pullover', 'dress', 'sandal', 'coat', 'shirt', 'sneaker', 'bag', 'ankle boot'}. The positive class is composed of 'T-shirt', 'pullover', 'coat', 'shirt', and 'bag' while the negative class is composed of 'trouser', 'dress', 'sandal', 'sneaker', and 'ankle boot'. We sampled 15,000 unlabeled data pairs as training data. The data set can be downloaded from <https://github.com/zalandoresearch/fashion-mnist>.
- CIFAR-10 [13]: It is a colorful object recognition data set. It is composed of 50,000 training examples and 10,000 test examples. The original feature dimension is $32 \times 32 \times 3$, and the label space is {'airplane', 'bird', 'automobile', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck'}. The positive class is composed of 'bird', 'deer', 'dog', 'frog', 'cat', and 'horse' while the negative class is composed of 'airplane', 'automobile', 'ship', and 'truck'. We sampled 10,000 unlabeled data pairs as training data. The data set can be downloaded from <https://www.cs.toronto.edu/~kriz/cifar.html>.

For the four UCI data sets, they can be downloaded from [5].

- Optdigits, USPS, Pendigits [5]: They are handwritten digit recognition data set. The train-test split can be found in Table 1. The feature dimensions are 62, 256, and 16 respectively and the label space is 0-9. The even digits are regarded as the positive class while the odd digits are regarded as the negative class. We sampled 1,200, 2,000, and 2,500 unlabeled data pairs for training respectively.
- Letter [5]: It is a letter recognition data set. It is composed of 16,000 training examples and 4,000 test examples. The feature dimension is 16 and the label space is the 26 capital letters in the English alphabet. The positive class is composed of the top 13 letters while the negative class is composed of the latter 13 letters. We sampled 4,000 unlabeled data pairs for training.

I.2 Details of Experiments on the KuaiRec Data Set

We used the small matrix of the KuaiRec [7] data set since it has dense confidence scores. It has 1,411 users and 3,327 items. We clipped the watching ratio above 2 and regarded the examples with watching ratio greater than 2 as positive examples. Following the experimental protocol of [8],

we regarded the latest positive example for each user as the positive testing data, and sampled 49 negative testing data to form the testing set for each user. The HR and NDCG were calculated at top 10. The learning rate was set to $1e-3$ and the dropout rate was set to 0.5. The number of epochs was set to 50 and the batch size was set to 256. The number of MLP layers was 2 and the embedding dimension was 128. The hyperparameters were the same for all the approaches for a fair comparison.

I.3 Details of Hyperparameters

All the experiments were conducted on NVIDIA GeForce RTX 3090. The number of training epochs was set to 200 and we obtained the testing accuracy by averaging the results in the last 10 epochs. All the methods were implemented in Pytorch [18]. We used the Adam optimizer [12]. To ensure fair comparisons, we set the same hyperparameter values for all the compared approaches.

For MNIST, Kuzushiji-MNIST and Fashion-MNIST, the learning rate was set to $1e-3$ and the weight decay was set to $1e-5$. The batch size was set to 256 data pairs. For training the probabilistic classifier to generate confidence, the batch size was set to 256 and the epoch number was set to 10.

For CIFAR10, the learning rate was set to $5e-4$ and the weight decay was set to $1e-5$. The batch size was set to 128 data pairs. For training the probabilistic classifier to generate confidence, the batch size was set to 128 and the epoch number was set to 10.

For all the UCI data sets, the learning rate was set to $1e-3$ and the weight decay was set to $1e-5$. The batch size was set to 128 data pairs. For training the probabilistic classifier to generate confidence, the batch size was set to 128 and the epoch number was set to 10.

The learning rate and weight decay for training the probabilistic classifier were the same as the setting for each data set correspondingly.

J More Experimental Results with Fewer Training Data

Figure 1 shows extra experimental results with fewer training data on other data sets with different class priors.

References

- [1] Christopher J. C. Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Gregory N. Hullender. Learning to rank using gradient descent. In *Proceedings of the 22nd International Conference on Machine Learning*, pages 89–96, 2005.
- [2] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: From pairwise approach to listwise approach. In *Proceedings of the 24th International Conference on Machine Learning*, pages 129–136, 2007.
- [3] Xi Chen, Paul N. Bennett, Kevyn Collins-Thompson, and Eric Horvitz. Pairwise ranking aggregation in a crowdsourced setting. In *Proceedings of the 6th ACM International Conference on Web Search and Data Mining*, pages 193–202, 2013.
- [4] Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha. Deep learning for classical Japanese literature. *CoRR*, abs/1812.01718, 2018.
- [5] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [6] Yanwei Fu, Timothy M. Hospedales, Tao Xiang, Jiechao Xiong, Shaogang Gong, Yizhou Wang, and Yuan Yao. Robust subjective visual property prediction from crowdsourced pairwise labels. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(3):563–577, 2015.
- [7] Chongming Gao, Shijun Li, Wenqiang Lei, Jiawei Chen, Biao Li, Peng Jiang, Xiangnan He, Jiaxin Mao, and Tat-Seng Chua. KuaiRec: A fully-observed dataset and insights for evaluating recommender systems. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pages 540–550, 2022.

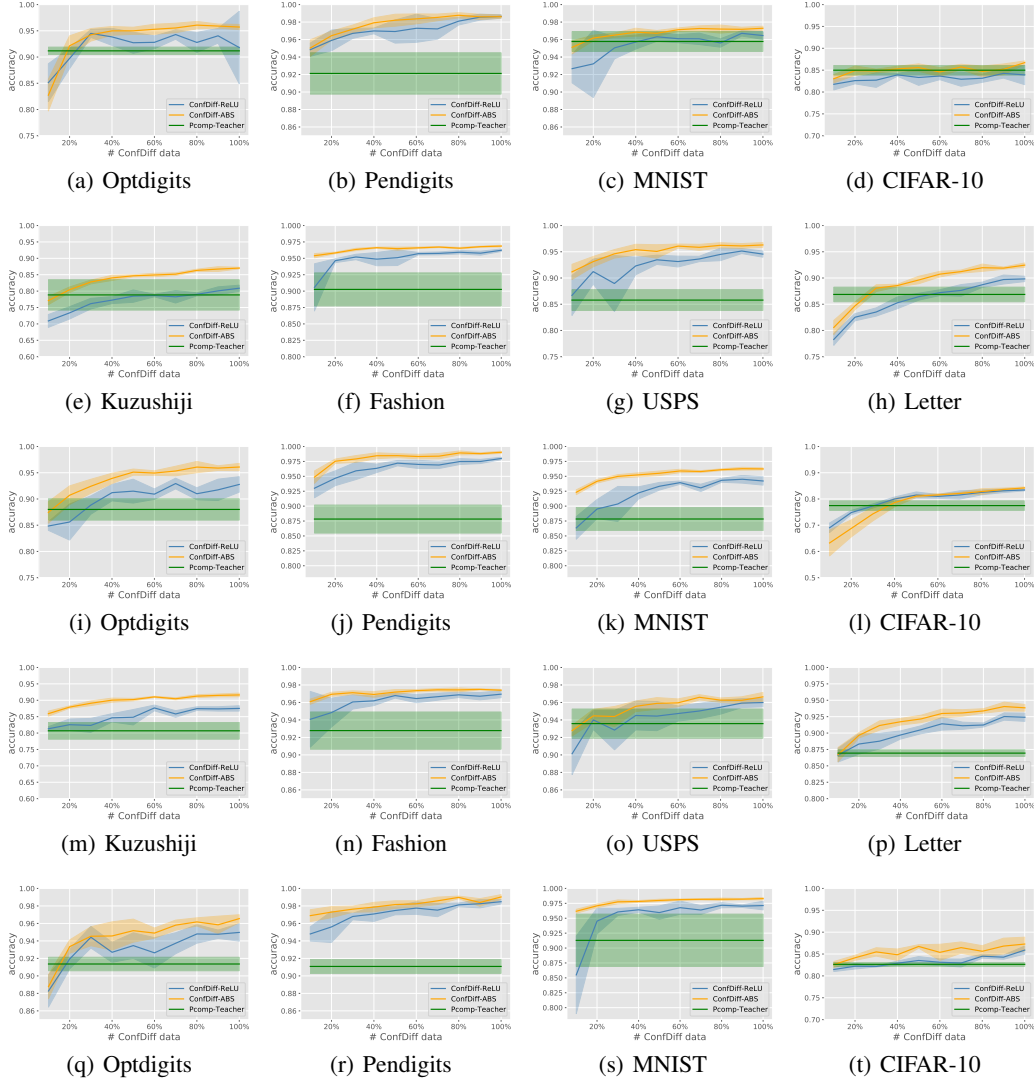


Figure 1: Classification performance of ConfDiff-ReLU and ConfDiff-ABS given a fraction of training data as well as Pcomp-Teacher given 100% of training data with different prior settings ($\pi_+ = 0.2$ for the first row, $\pi_+ = 0.5$ for the second and the third row, and $\pi_+ = 0.8$ for the fourth and the fifth row).

- 205 [8] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural
206 collaborative filtering. In *Proceedings of the 26th International Conference on World Wide Web*,
207 pages 173–182, 2017.
- 208 [9] Yixuan He, Quan Gan, David Wipf, Gesine D. Reinert, Junchi Yan, and Mihai Cucuringu.
209 GNNRank: Learning global rankings from pairwise comparisons via directed graph neural
210 networks. In *Proceedings of the 39th International Conference on Machine Learning*, pages
211 8581–8612, 2022.
- 212 [10] Kevin G. Jamieson and Robert D. Nowak. Active ranking using pairwise comparisons. In
213 *Advances in Neural Information Processing Systems 24*, pages 2240–2248, 2011.
- 214 [11] Daniel M. Kane, Shachar Lovett, Shay Moran, and Jiapeng Zhang. Active classification with
215 comparison queries. In *2017 IEEE 58th Annual Symposium on Foundations of Computer
216 Science*, pages 355–366, 2017.

- 217 [12] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings*
218 *of the 3rd International Conference on Learning Representations*, 2015.
- 219 [13] Alex Krizhevsky and Geoffrey E. Hinton. Learning multiple layers of features from tiny images.
220 Technical report, University of Toronto, 2009.
- 221 [14] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning
222 applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 223 [15] Tie-Yan Liu. *Learning to Rank for Information Retrieval*. Springer, 2011.
- 224 [16] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*.
225 The MIT Press, 2012.
- 226 [17] Dohyung Park, Joe Neeman, Jin Zhang, Sujay Sanghavi, and Inderjit S. Dhillon. Preference
227 completion: Large-scale collaborative ranking from pairwise comparisons. In *Proceedings of*
228 *the 32nd International Conference on Machine Learning*, pages 1907–1916, 2015.
- 229 [18] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
230 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative
231 style, high-performance deep learning library. In *Advances in Neural Information Processing*
232 *Systems 32*, pages 8026–8037, 2019.
- 233 [19] Nihar B. Shah, Sivaraman Balakrishnan, and Martin J. Wainwright. Feeling the bern: Adap-
234 tive estimators for bernoulli probabilities of pairwise comparisons. *IEEE Transactions on*
235 *Information Theory*, 65(8):4854–4874, 2019.
- 236 [20] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: A novel image dataset for
237 benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017.
- 238 [21] Liyuan Xu, Junya Honda, Gang Niu, and Masashi Sugiyama. Uncoupled regression from
239 pairwise comparison data. In *Advances in Neural Information Processing Systems 32*, pages
240 3992–4002, 2019.
- 241 [22] Yichong Xu, Sivaraman Balakrishnan, Aarti Singh, and Artur Dubrawski. Regression with
242 comparisons: Escaping the curse of dimensionality with ordinal information. *Journal of*
243 *Machine Learning Research*, 21(1):6480–6533, 2020.
- 244 [23] Yichong Xu, Hongyang Zhang, Kyle Miller, Aarti Singh, and Artur Dubrawski. Noise-tolerant
245 interactive learning using pairwise comparisons. In *Advances in Neural Information Processing*
246 *Systems 30*, pages 2428–2437, 2017.
- 247 [24] Shiwei Zeng and Jie Shen. Efficient pac learning from the crowd with pairwise comparisons. In
248 *Proceedings of the 39th International Conference on Machine Learning*, pages 25973–25993,
249 2022.