

Appendix A. Proofs of Loss Function Properties

As stated in Section 4, if we define a loss function \mathcal{L} in an additive form by (12), there exist constants ξ_1 and ξ_2 which satisfy (10) and (11). We first describe proof of (10). The following formulation holds in accordance with $\ell(z) + \ell(-z) = a$.

$$\begin{aligned}
& \sum_{y \in Y} \mathcal{L}_{\text{OVA}}(f(x), y) \\
&= \sum_{y \in Y} \left(\ell(g_y(x)) + \frac{1}{K-1} \sum_{y' \neq y} \ell(-g_{y'}(x)) \right) \\
&= \sum_{y \in Y} \left(\frac{K}{K-1} \ell(g_y(x)) - \frac{1}{K-1} \ell(g_y(x)) + \frac{1}{K-1} \sum_{y' \in \mathcal{Y}} \ell(-g_{y'}(x)) - \frac{1}{K-1} \ell(-g_y(x)) \right) \\
&= \frac{1}{K-1} \sum_{y \in Y} \left(K \ell(g_y(x)) + \sum_{y' \in \mathcal{Y}} \ell(-g_{y'}(x)) - a \right) \\
&= \frac{1}{K-1} \left(K \sum_{y \in Y} \ell(g_y(x)) + \sum_{y \in Y} \sum_{y' \in \mathcal{Y}} \ell(-g_{y'}(x)) - \sum_{y \in Y} a \right) \\
&= \frac{1}{K-1} \left((K-N+N) \sum_{y \in Y} \ell(g_y(x)) + N \sum_{y \in Y} \ell(-g_y(x)) + N \sum_{y' \notin Y} \ell(-g_{y'}(x)) - aN \right) \\
&= \frac{1}{K-1} \left((K-N) \sum_{y \in Y} \ell(g_y(x)) + N \sum_{y' \notin Y} \ell(-g_{y'}(x)) + aN(N-1) \right) \\
&= \frac{K-N}{K-1} \sum_{y \in Y} \ell(g_y(x)) + \frac{N}{K-1} \sum_{y' \notin Y} \ell(-g_{y'}(x)) + \frac{aN(N-1)}{K-1} \\
&= \mathcal{L}_{\text{OVA}}(f(x), Y) + \frac{aN(N-1)}{K-1}
\end{aligned}$$

■

Next, we describe proof of (11). Similar to the above formulation, in accordance with $\ell(z) + \ell(-z) = a$,

$$\begin{aligned}
& \sum_{y \in Y} \mathcal{L}_{\text{PC}}(f(x), y) \\
&= \sum_{y \in Y} \sum_{y' \neq y} \ell(g_y(x) - g_{y'}(x)) \\
&= \sum_{y \in Y} \sum_{y' \in \mathcal{Y}} \ell(g_y(x) - g_{y'}(x)) - \sum_{y \in Y} \ell(g_y(x) - g_y(x))
\end{aligned}$$

$$\begin{aligned}
&= \sum_{y \in Y} \sum_{y' \notin Y} \ell(g_y(x) - g_{y'}(x)) + \sum_{\substack{y, y' \in Y \\ y=y'}} \ell(g_y(x) - g_{y'}(x)) + \sum_{\substack{y, y' \in Y \\ y \neq y'}} \ell(g_y(x) - g_{y'}(x)) - \frac{aN}{2} \\
&= \sum_{y \in Y} \sum_{y' \notin Y} \ell(g_y(x) - g_{y'}(x)) + a \cdot NC_2 \\
&= \mathcal{L}_{\text{PC}}(f(x), Y) + a \cdot NC_2
\end{aligned}$$

The third and fourth equality hold due to $\ell(0) = a/2$. ■

Appendix B. Property of the Data-Generation Probability Model

In this section, we describe proof of (15). To begin with, we introduce the following lemma.

Lemma 5 *Let any finite sets be X with size of K , and any elements of the N -size power set $\mathfrak{P}_N(X)$ be A . Then, the following equation holds for any function f and g over X .*

$$\sum_{A \in \mathfrak{P}_N(X)} \sum_{a_1, a_2 \in A} f(a_1)g(a_2) = {}_{K-2}C_{N-2} \sum_{\substack{x_1, x_2 \in X \\ x_1 \neq x_2}} f(x_1)g(x_2) + {}_{K-1}C_{N-1} \sum_{\substack{x_1, x_2 \in X \\ x_1 = x_2}} f(x_1)g(x_2) \quad (23)$$

Proof The left-hand side can be formulated as:

$$\sum_{A \in \mathfrak{P}_N(X)} \sum_{a_1, a_2 \in A} f(a_1)g(a_2) = \sum_{A \in \mathfrak{P}_N(X)} \sum_{\substack{a_1, a_2 \in A \\ a_1 \neq a_2}} f(a_1)g(a_2) + \sum_{A \in \mathfrak{P}_N(X)} \sum_{\substack{a_1, a_2 \in A \\ a_1 = a_2}} f(a_1)g(a_2)$$

For the first term, any A can be chosen from $\mathfrak{P}_N(X)$ in ${}_K C_N$ patterns, and any a_1, a_2 ($a_1 \neq a_2$) can be chosen from A in ${}_N P_2$ patterns. Because $f(a_1)g(a_2)$ with any a_1, a_2 ($a_1 \neq a_2$) can be chosen from X in ${}_K P_2$ patterns, the first term in (23) holds due to ${}_{K-2}C_{N-2} = \frac{{}_K C_N \cdot {}_N P_2}{{}_K P_2}$. Similar for the second term. ■

Here, we prove (15). Let $P(\alpha \in Y|x)$ be the probability where a set of candidate labels Y includes a label $\alpha \in \mathcal{Y}$. Denoting I as an indicator function, then

$$\begin{aligned}
P(\alpha \in Y) &= \sum_{Y \in \mathfrak{P}_N(\mathcal{Y})} I(\alpha \in Y) P_N(Y|x) \\
&= \frac{1}{{}_{K-1}C_{N-1}} \sum_{Y \in \mathfrak{P}_N(\mathcal{Y})} I(\alpha \in Y) \sum_{y \in Y} P(y|x) \\
&= \frac{1}{{}_{K-1}C_{N-1}} \sum_{Y \in \mathfrak{P}_N(\mathcal{Y})} \sum_{y_1, y_2 \in Y} I(y_1 = \alpha) P(y_2|x) \\
&= \frac{1}{{}_{K-1}C_{N-1}} \left({}_{K-1}C_{N-1} \sum_{y \in \mathcal{Y}} I(y = \alpha) P(y|x) + {}_{K-2}C_{N-2} \sum_{\substack{y_1, y_2 \in \mathcal{Y} \\ y_1 \neq y_2}} I(y_1 = \alpha) P(y_2|x) \right)
\end{aligned}$$

$$\begin{aligned}
&= P(\alpha|x) + \frac{N-1}{K-1} \sum_{y \neq \alpha} P(y|x) \\
&= \frac{K-N}{K-1} P(\alpha|x) + \frac{N-1}{K-1}
\end{aligned}$$

The fourth equality holds due to lemma 5. In the fifth equality, the second term holds because $y_2 \neq \alpha$ holds if $y_1 = \alpha$. Therefore,

$$\begin{aligned}
Q(y_o = \alpha|x) &= Q(y_o = \alpha|x, \alpha \in Y)Q(\alpha \in Y|x) + Q(y_o = \alpha|x, \alpha \notin Y)Q(\alpha \notin Y|x) \\
&= \frac{1}{N} \frac{K-N}{N-1} P(\alpha|x) + \frac{1}{N} \frac{N-1}{K-1} \\
&= \beta P(\alpha|x) + (1-\beta) \frac{1}{K}
\end{aligned}$$

The second equality holds due to the assumption where $Q(y_o = \alpha|x, \alpha \in Y) = 1/N$ and $Q(y_o = \alpha|x, \alpha \notin Y) = 0$ hold. \blacksquare

Appendix C. Proof of Theorem 1

We first derive the expectation of the sum of loss for all the ordinary labels.

$$\begin{aligned}
&\mathbb{E}_{P_N(Y|x)} \left[\sum_{y \in Y} \mathcal{L}(f(x), y) \right] \\
&= \sum_{Y \in \mathfrak{P}_N(\mathcal{Y})} \sum_{y \in Y} \mathcal{L}(f(x), y) P_N(Y|x) \\
&= \frac{1}{K-1} \sum_{Y \in \mathfrak{P}_N(\mathcal{Y})} \sum_{y \in Y} \sum_{y' \in Y} \mathcal{L}(f(x), y') P_N(y|x) \\
&= \frac{1}{K-1} \sum_{Y \in \mathfrak{P}_N(\mathcal{Y})} \left\{ \sum_{y \in Y} \sum_{y' \in Y, y' \neq y} \mathcal{L}(f(x), y') P_N(y|x) + \sum_{y \in Y} \mathcal{L}(f(x), y) P_N(y|x) \right\} \\
&= \frac{N-1}{K-1} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}, y' \neq y} \mathcal{L}(f(x), y') P_N(y|x) + \sum_{y \in \mathcal{Y}} \mathcal{L}(f(x), y) P_N(y|x) \\
&= \frac{N-1}{K-1} \mathbb{E}_{P(y|x)} \left[\sum_{y' \in \mathcal{Y}} \mathcal{L}(f(x), y') - \mathcal{L}(f(x), y) \right] + \mathbb{E}_{P(y|x)} [\mathcal{L}(f(x), y)] \\
&= \frac{N-1}{K-1} \sum_{y \in \mathcal{Y}} \mathcal{L}(f(x), y) + \frac{K-N}{K-1} \mathbb{E}_{P(y|x)} [\mathcal{L}(f(x), y)]
\end{aligned}$$

The second equality holds due to the definition of $P_N(Y|x)$. The third equality holds due to Lemma 5. Thus, the following formulation holds due to (12).

$$\mathbb{E}_{P_N(x,Y)} [\mathcal{L}(f(x), Y)] = \xi_1 \mathbb{E}_{P_N(x,Y)} \left[\sum_{y \in Y} \mathcal{L}(f(x), y) \right] + \xi_2$$

$$= \frac{\xi_1(K-N)}{K-1} \mathbb{E}_{P(x,y)} [\mathcal{L}(f(x), y)] + \frac{\xi_1(N-1)}{K-1} \sum_{y \in \mathcal{Y}} \mathcal{L}(f(x), y) + \xi_2$$

■

Appendix D. Proof of Lemma 2

According to the duality described in (12) and (13), loss function \mathcal{L} can be formulated by $\mathcal{L}(f(x), y)$ for N candidate labels $y \in Y$, or $\tilde{\mathcal{L}}(f(x), \bar{y})$ for $K-N$ complementary labels $\bar{y} \in \bar{Y}$. Thus, we can redefine loss function \mathcal{L} as the following:

$$\mathcal{L}(f(x), Y) = \sum_{y \in \tilde{Y}} \tilde{\mathcal{L}}(f(x), y) = \begin{cases} \sum_{y \in Y} \mathcal{L}(f(x), y), & \text{if } N \leq \frac{K}{2} \\ \sum_{\bar{y} \in \bar{Y}} \tilde{\mathcal{L}}(f(x), \bar{y}), & \text{otherwise} \end{cases}$$

where $\tilde{\mathcal{L}}$ and \tilde{Y} denote \mathcal{L} and Y respectively if $N \leq K/2$, otherwise $\tilde{\mathcal{L}}, \bar{Y}$ respectively. Therefore, given $\tilde{N} := |\tilde{Y}|$ it always satisfies $\tilde{N} \leq K/2$. Note that $\xi_2 = \tilde{\xi}_2 = 0$ due to the assumption of $a = 0$, as discussed in Section 4.1. Similarly, $\tilde{\ell}(z)$ denotes $\tilde{\ell} : z \mapsto \ell(z)$ if $N \leq K/2$ otherwise $\tilde{\ell} : z \mapsto \ell(-z)$. For the rest of this work, we prove Lemma 2 according to those definitions.

First we describe proof for one-versus-all classification. Under the assumption of $a = 0$, the following formulation holds.

$$\sum_{y \in \tilde{Y}} \tilde{\mathcal{L}}_{\text{OVA}}(f(x), y) = \sum_{y \in \tilde{Y}} \left(\frac{K}{K-1} \tilde{\ell}(g_y(x)) + \frac{1}{K-1} \sum_{y' \in \mathcal{Y}} \tilde{\ell}(-g_{y'}(x)) \right)$$

Thus, the following formulation holds for any $\tilde{Y}, \tilde{Y}' \in \mathfrak{P}_N(\mathcal{Y})$.

$$\begin{aligned} \|\mathcal{L}_{\text{OVA}}\|_{\infty} &= \sup_{g_1, \dots, g_K \in \mathcal{G}} \left(\sum_{y \in \tilde{Y}} \tilde{\mathcal{L}}_{\text{OVA}}(f(x), y) - \sum_{y' \in \tilde{Y}'} \tilde{\mathcal{L}}_{\text{OVA}}(f(x), y') \right) \\ &= \sup_{g_1, \dots, g_K \in \mathcal{G}} \left\{ \frac{K}{K-1} \left(\sum_{y \in \tilde{Y}} \tilde{\ell}(g_y(x)) - \sum_{y' \in \tilde{Y}'} \tilde{\ell}(g_{y'}(x)) \right) \right. \\ &\quad \left. + \frac{1}{K-1} \left(\sum_{y \in \mathcal{Y}} \tilde{\ell}(-g_y(x)) - \sum_{y \in \mathcal{Y}} \tilde{\ell}(-g_y(x)) \right) \right\} \\ &\leq \frac{K}{K-1} \left(\sup_{g_1, \dots, g_K \in \mathcal{G}} \sum_{y \in \tilde{Y}} \tilde{\ell}(g_y(x)) - \inf_{g_1, \dots, g_K} \sum_{y' \in \tilde{Y}'} \tilde{\ell}(g_{y'}(x)) \right) \\ &\leq \frac{K}{K-1} \left(\frac{\tilde{N}}{2} + \frac{\tilde{N}}{2} \right) \end{aligned}$$

$$= \frac{K\tilde{N}}{K-1}$$

The second inequality holds because supremum and infimum of ℓ are $1/2$ and $-1/2$ respectively. \blacksquare

We further describe proof for pairwise classification. Under the assumption of $a = 0$, the following formulation holds.

$$\sum_{y \in \tilde{Y}} \tilde{\mathcal{L}}_{\text{PC}}(f(x), y) = \sum_{y \in \tilde{Y}} \sum_{y' \notin \tilde{Y}} \tilde{\ell}(g_y(x) - g_{y'}(x))$$

Thus,

$$\begin{aligned} \|\mathcal{L}_{\text{PC}}\|_{\infty} &= \sup_{g_1, \dots, g_K \in \mathcal{G}} \left(\sum_{y' \in \tilde{Y}} \tilde{\mathcal{L}}_{\text{PC}}(f(x), y') - \sum_{y \in \tilde{Y}'} \tilde{\mathcal{L}}_{\text{PC}}(f(x), y) \right) \\ &= \sup_{g_1, \dots, g_K \in \mathcal{G}} \sum_{y \in \tilde{Y}} \sum_{y' \notin \tilde{Y}} \tilde{\ell}(g_y(x) - g_{y'}(x)) - \inf_{g_1, \dots, g_K \in \mathcal{G}} \sum_{y' \in \tilde{Y}'} \sum_{y \notin \tilde{Y}'} \tilde{\ell}(g_{y'}(x) - g_y(x)) \\ &\leq \frac{\tilde{N}(K - \tilde{N})}{2} - \left(-\frac{\tilde{N}(K - \tilde{N})}{2} \right) \\ &= \tilde{N}(K - \tilde{N}) \end{aligned}$$

\blacksquare

Appendix E. Proof of Lemma 3

We describe proof for one-versus-all classification. Under the assumption that $h \in \mathcal{H}_{\text{OVA}}$ is equivalent to \mathcal{L}_{OVA} , the following formulation holds due to the definition of \mathcal{H}_{OVA} .

$$\begin{aligned} \mathfrak{R}_n(\mathcal{H}_{\text{OVA}}) &= \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{h \in \mathcal{H}_{\text{OVA}}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(x_i, Y_i) \right] \\ &= \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \tilde{Y}_i} \tilde{\mathcal{L}}_{\text{OVA}}(f(x_i), y) \right] \\ &= \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \tilde{Y}_i} \left\{ \frac{K}{K-1} \tilde{\ell}(g_y(x_i)) + \frac{1}{K-1} \sum_{y' \in \mathcal{Y}} \tilde{\ell}(-g_{y'}(x_i)) \right\} \right] \\ &\leq \frac{K}{K-1} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \tilde{Y}_i} \tilde{\ell}(g_y(x_i)) \right] \end{aligned}$$

$$+ \frac{1}{K-1} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \tilde{Y}_i} \sum_{y' \in \mathcal{Y}} \tilde{\ell}(-g_{y'}(x_i)) \right]$$

Let $I(y \in \tilde{Y}_i)$ be an indicator function and define $\alpha_i := 2I(y \in \tilde{Y}_i) - 1$, then for the first term,

$$\begin{aligned} & \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \tilde{Y}_i} \tilde{\ell}(g_y(x_i)) \right] \\ &= \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \mathcal{Y}} \tilde{\ell}(g_y(x_i)) I(y \in \tilde{Y}_i) \right] \\ &= \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{2n} \sum_{i=1}^n \sigma_i \sum_{y \in \mathcal{Y}} \tilde{\ell}(g_y(x_i)) (\alpha_i + 1) \right] \\ &\leq \sum_{y \in \mathcal{Y}} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{2n} \sum_{i=1}^n \sigma_i \tilde{\ell}(g_y(x_i)) (\alpha_i + 1) \right] \\ &\leq \sum_{y \in \mathcal{Y}} \left\{ \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{2n} \sum_{i=1}^n \alpha_i \sigma_i \sum_{y \in \mathcal{Y}} \tilde{\ell}(g_y(x_i)) \right] + \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{2n} \sum_{i=1}^n \sigma_i \sum_{y \in \mathcal{Y}} \tilde{\ell}(g_y(x_i)) \right] \right\} \\ &= K \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \mathcal{Y}} \tilde{\ell}(g(x_i)) \right] \\ &= K \mathfrak{R}_n(\tilde{\ell} \circ \mathcal{G}) \end{aligned}$$

The second equality from the last holds because σ_i and $\alpha_i \sigma_i$ are drawn from the same probabilistic distribution. For the second term,

$$\begin{aligned} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \tilde{Y}_i} \sum_{y' \in \mathcal{Y}} \tilde{\ell}(-g_{y'}(x_i)) \right] &\leq \sum_{y \in \mathcal{Y}} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_y \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{N} \tilde{\ell}(-g_y(x_i)) \right] \\ &= K \tilde{N} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{\ell}(-g(x_i)) \right] \\ &= K \tilde{N} \mathfrak{R}_n(\tilde{\ell} \circ \mathcal{G}) \end{aligned}$$

Thus,

$$\begin{aligned} \mathfrak{R}_n(\mathcal{H}_{\text{OVA}}) &\leq \frac{K^2 + K \tilde{N}}{K-1} \mathfrak{R}_n(\tilde{\ell} \circ \mathcal{G}) \\ &\leq \frac{K(K + \tilde{N})}{K-1} L_{\ell} \mathfrak{R}_n(\mathcal{G}) \end{aligned}$$

The second inequality holds due to $\mathfrak{R}_n(\tilde{\ell} \circ \mathcal{G}) \leq L_\ell \mathfrak{R}_n(\mathcal{G})$ according to Talagrand's contraction lemma (Ledoux and Talagrand, 2013). \blacksquare

Note that $\ell(z) - \ell(-z) = a$ is incorrectly assumed in (Ishida et al., 2017), which causes miscalculation in proof of Lemma 3.

We further describe proof for pairwise classification. Under the assumption that h and \mathcal{L}_{PC} are equivalent,

$$\begin{aligned}
& \mathfrak{R}_n(\mathcal{H}_{\text{PC}}) \\
&= \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{h \in \mathcal{H}_{\text{PC}}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(x_i, Y_i) \right] \\
&= \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \tilde{Y}_i} \tilde{\mathcal{L}}_{\text{PC}}(f(x_i), y) \right] \\
&= \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \tilde{Y}_i} \sum_{y' \neq y} \tilde{\ell}(g_y(x_i) - g_{y'}(x_i)) \right] \\
&= \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{y \in \mathcal{Y}} \sum_{y' \neq y} \tilde{\ell}(g_y(x_i) - g_{y'}(x_i)) I(y \in \tilde{Y}_i) \right] \\
&\leq \sum_{y \in \mathcal{Y}} \sum_{y' \neq y} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_y, g_{y'} \in \mathcal{G}} \frac{1}{2n} \sum_{i=1}^n \sigma_i \tilde{\ell}(g_y(x_i) - g_{y'}(x_i)) (\alpha_i + 1) \right] \\
&\leq \sum_{y \in \mathcal{Y}} \sum_{y' \neq y} \left\{ \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_y, g_{y'} \in \mathcal{G}} \frac{1}{2n} \sum_{i=1}^n \alpha_i \sigma_i \tilde{\ell}(g_y(x_i) - g_{y'}(x_i)) \right] \right. \\
&\quad \left. + \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_y, g_{y'} \in \mathcal{G}} \frac{1}{2n} \sum_{i=1}^n \sigma_i \tilde{\ell}(g_y(x_i) - g_{y'}(x_i)) \right] \right\} \\
&\leq \sum_{y \in \mathcal{Y}} \sum_{y' \neq y} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_y, g_{y'} \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{\ell}(g_y(x_i) - g_{y'}(x_i)) \right]
\end{aligned}$$

Let we define $\mathcal{G}_{g_y, g_{y'}} := \{x \mapsto g_y(x) - g_{y'}(x) | g_y, g_{y'} \in \mathcal{G}\}$, then:

$$\begin{aligned}
& \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_y, g_{y'} \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{\ell}(g_y(x_i) - g_{y'}(x_i)) \right] \\
&= \mathfrak{R}_n(\tilde{\ell} \circ \mathcal{G}_{g_y, g_{y'}}) \\
&\leq L_\ell \mathfrak{R}_n(\mathcal{G}_{g_y, g_{y'}}) \\
&= L_\ell \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_y, g_{y'} \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i (g_y(x_i) - g_{y'}(x_i)) \right]
\end{aligned}$$

$$\begin{aligned}
&\leq L_\ell \mathbb{E}_S \mathbb{E}_\sigma \left[\sup_{g_y \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g_y(x_i) \right] + L_\ell \mathbb{E}_S \mathbb{E}_\sigma \left[\sup_{g_{y'} \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n (-\sigma_i) g_{y'}(x_i) \right] \\
&= 2L_\ell \mathfrak{R}_n(\mathcal{G})
\end{aligned}$$

The third equality holds because σ_i and $-\sigma_i$ are drawn from the same probabilistic distribution. Then,

$$\mathfrak{R}_n(\mathcal{H}_{\text{PC}}) \leq 2K(K-1)L_\ell \mathfrak{R}_n(\mathcal{G})$$

■

Appendix F. Proof of Theorem 4

We only describe proof for one-versus-all classification; proof for pairwise classification is similar. We substitute the j th data (x_j, Y_j) in \mathcal{S} with any data (x'_j, Y'_j) , and define the data set as \mathcal{S}' . Let a set of empirical discrimination functions and empirical risk for \mathcal{S}' be $\mathcal{G}' := \{g'\}$ and $\hat{R}'(f)$ respectively. Then the following formulation holds due to Lemma 2.

$$\begin{aligned}
&\sup_{g_1, \dots, g_K \in \mathcal{G}} \left(\hat{R}(f) - R(f) \right) - \sup_{g'_1, \dots, g'_K \in \mathcal{G}'} \left(\hat{R}'(f) - R(f) \right) \\
&\leq \frac{K-1}{n(K-N)} \sup_{g_1, \dots, g_K \in \mathcal{G}} \inf_{g'_1, \dots, g'_K \in \mathcal{G}'} \left\{ \left(\mathcal{L}_{\text{OVA}}(f(x_j), Y_j) - \mathcal{L}_{\text{OVA}}(f(x'_j), Y'_j) \right) \right\} \\
&\leq \frac{K-1}{n(K-N)} \|\mathcal{L}_{\text{OVA}}\|_\infty \\
&\leq \frac{K\tilde{N}}{n(K-N)}
\end{aligned}$$

According to McDiarmid's inequality (McDiarmid, 1989), for any integer $\delta > 0$ the following formulation holds with a probability at least $1 - \delta/2$.

$$\sup_{g_1, \dots, g_K} \left(\hat{R}(f) - R(f) \right) - \mathbb{E} \left[\sup_{g_1, \dots, g_K} \left(\hat{R}(f) - R(f) \right) \right] \leq \frac{K\tilde{N}}{2(K-N)} \sqrt{\frac{2 \ln(2/\delta)}{n}}$$

Let $\mathcal{S}' := \{(x'_i, Y'_i)\}_{i=1}^n$ be any dataset where each data is drawn from the data-generation probability model $P_N(x, Y)$. Due to $R(f) = \mathbb{E} \left[\hat{R}(f) \right]$,

$$\begin{aligned}
&\mathbb{E} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \left(\hat{R}(f) - R(f) \right) \right] \\
&= \frac{K-1}{K-N} \mathbb{E}_S \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \left(\frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{OVA}}(f(x_i), Y_i) - \mathbb{E}_{S'} \left[\frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{OVA}}(f(x'_i), Y'_i) \right] \right) \right] \\
&\leq \frac{K-1}{K-N} \mathbb{E}_S \mathbb{E}_{S'} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \left(\frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{OVA}}(f(x_i), Y_i) - \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{OVA}}(f(x'_i), Y'_i) \right) \right]
\end{aligned}$$

$$\begin{aligned}
&= \frac{K-1}{K-N} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\mathcal{S}'} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \left(\frac{1}{n} \sum_{i=1}^n \sigma_i \mathcal{L}_{\text{OVA}}(f(x_i), Y_i) + \frac{1}{n} \sum_{i=1}^n (-\sigma_i) \mathcal{L}_{\text{OVA}}(f(x'_i), Y'_i) \right) \right] \\
&\leq \frac{2(K-1)}{K-N} \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\sigma} \left[\sup_{g_1, \dots, g_K \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathcal{L}_{\text{OVA}}(f(x_i), Y_i) \right] \\
&= \frac{2(K-1)}{K-N} \mathfrak{R}_n(\mathcal{H}_{\text{OVA}}) \\
&\leq \frac{2K(K+\tilde{N})}{K-N} L_{\ell} \mathfrak{R}_n(\mathcal{G})
\end{aligned}$$

The second equality holds because $\mathcal{L}(f(x_i, Y_i))$ and $\sigma_i \mathcal{L}(f(x_i, Y_i))$ are drawn from the same probabilistic distribution; similar for $\mathcal{L}(f(x'_i, Y'_i))$. The last inequality holds due to Lemma 3. $\hat{R}(\hat{f}) \leq \hat{R}(f^*)$ holds according to (20), thus,

$$\begin{aligned}
\mathcal{E}_N &= \left(\hat{R}(\hat{f}) - \hat{R}(f^*) \right) + \left(R(\hat{f}) - \hat{R}(\hat{f}) \right) + \left(\hat{R}(f^*) - R(f^*) \right) \\
&\leq 2 \sup_{g_1, \dots, g_K \in \mathcal{G}} \left| \hat{R}(f) - R(f) \right| \\
&\leq \frac{4K(K+\tilde{N})}{K-N} L_{\ell} \mathfrak{R}_n(\mathcal{G}) + \frac{K\tilde{N}}{K-N} \sqrt{\frac{2 \ln(2/\delta)}{n}}
\end{aligned}$$

■

Appendix G. Comparison with Existing Work

Our major contribution is the finding that loss functions reflecting the properties of label space naturally bridge ordinary-label learning and complementary-label learning. However, unlike (Ishida et al., 2019; Feng et al., 2020), it is also true that satisfaction of additivity and duality is a strict limitation. In this section we present an experiment with a loss function which does not satisfy the assumption to evaluate the limitation of our work.

Here we introduce a loss function \mathcal{L}_{CE} as follows:

$$\mathcal{L}_{\text{CE}}(f(x), Y) = - \sum_{y \in Y} \log \left(\frac{\exp(g_y(x))}{\sum_{y' \in \mathcal{Y}} \exp(g_{y'}(x))} \right)$$

\mathcal{L}_{CE} is represented as an additive form of cross entropy loss, the most commonly used loss function. Unlike \mathcal{L}_{OVA} or \mathcal{L}_{PC} , learning with \mathcal{L}_{CE} is outside the scope of our framework because additivity and duality are not satisfied.

We performed an experiment with \mathcal{L}_{CE} on the same dataset described in Section 6.1 to evaluate the classification accuracy and error. We set the batch size as 64 and the number of epochs as 300. MLP was adopted for MNIST, Fashion-MNIST, and Kuzushiji-MNIST datasets; *Adam* was used for optimization with weight decay of 10^{-4} and learning rate of 5×10^{-5} . DenseNet was adopted for CIFAR-10 dataset; *stochastic gradient descent* was used for optimization with weight decay of 5×10^{-4} and initial learning rate of 5×10^{-2}

Table 4: Experimental classification accuracies for 10 and 5 class classification (%). The experiments were performed 5 times for each case; the mean accuracy and standard deviation are presented by the upper and lower values, respectively. The highest accuracy is boldfaced.

N	K = 10									K = 5			
	1	2	3	4	5	6	7	8	9	1	2	3	4
MNIST (CE)	94.69 (±0.13)	93.70 (±0.19)	93.1 (±0.25)	92.17 (±0.50)	91.11 (±0.35)	89.60 (±0.73)	86.42 (±0.79)	80.99 (±0.81)	71.33 (±0.78)	98.05 (±0.10)	97.05 (±0.14)	96.02 (±0.52)	94.28 (±0.49)
Fashion (CE)	86.81 (±0.43)	85.34 (±0.23)	84.97 (±0.36)	83.75 (±0.36)	82.93 (±0.30)	80.60 (±0.66)	77.89 (±0.35)	73.83 (±0.63)	63.05 (±1.91)	87.21 (±0.27)	85.91 (±0.65)	84.16 (±0.45)	81.63 (±0.66)
Kuzushiji (CE)	77.19 (±0.53)	71.38 (±0.60)	69.22 (±0.53)	66.43 (±0.76)	63.05 (±0.36)	57.31 (±0.36)	51.62 (±1.13)	42.25 (±1.22)	30.99 (±0.51)	82.01 (±0.48)	78.23 (±0.77)	73.84 (±1.12)	65.58 (±2.91)
CIFAR-10 (CE)	73.23 (±0.20)	56.84 (±0.52)	43.00 (±0.55)	36.03 (±0.77)	32.57 (±1.25)	35.47 (±1.76)	40.18 (±0.45)	40.88 (±0.89)	29.70 (±0.53)	74.81 (±0.36)	64.69 (±0.93)	58.54 (±0.80)	49.43 (±1.97)

Table 5: Experimental classification errors for 10 and 5 class classification ($\times 10^{-3}$). The experiments were performed 5 times for each case; the mean error and standard deviation are presented by the upper and lower values, respectively. The highest error is boldfaced.

N	K = 10									K = 5			
	1	2	3	4	5	6	7	8	9	1	2	3	4
MNIST (CE)	0.07 (±0.01)	0.40 (±0.16)	0.49 (±0.22)	0.49 (±0.23)	0.61 (±0.18)	0.68 (±0.11)	0.51 (±0.27)	0.46 (±0.23)	0.86 (±0.26)	0.02 (±0.02)	0.43 (±0.12)	0.51 (±0.15)	0.33 (±0.20)
Fashion (CE)	0.08 (±0.06)	0.39 (±0.22)	0.39 (±0.13)	0.41 (±0.23)	0.34 (±0.19)	0.52 (±0.22)	0.42 (±0.16)	0.38 (±0.29)	0.75 (±0.49)	0.05 (±0.02)	0.22 (±0.14)	0.30 (±0.04)	0.48 (±0.24)
Kuzushiji (CE)	0.33 (±0.13)	0.36 (±0.07)	0.27 (±0.14)	0.23 (±0.14)	0.28 (±0.13)	0.28 (±0.09)	0.23 (±0.10)	0.65 (±0.30)	0.92 (±0.36)	0.18 (±0.12)	0.26 (±0.15)	0.27 (±0.16)	0.46 (±0.30)
CIFAR-10 (CE)	1.31 (±0.21)	1.34 (±0.17)	1.49 (±0.15)	1.55 (±0.14)	1.72 (±0.08)	1.01 (±0.10)	0.28 (±0.08)	0.08 (±0.03)	0.07 (±0.02)	0.04 (±0.03)	0.25 (±0.07)	0.22 (±0.06)	0.13 (±0.02)

which was halved every 30 epochs. The results of accuracy and error are listed in Table 4 and 5 correspondingly.

Compared with the results in Table 2 and 3, the classifier with \mathcal{L}_{CE} tends to perform better than those of \mathcal{L}_{OVA} and \mathcal{L}_{PC} . This shows that although the assumption of additivity and duality naturally reflects the relationship between ordinary-label and complementary-label, it could be a limitation on performance in a context of practical use.