
\mathcal{H} -Consistency Bounds for Pairwise Misranking Loss Surrogates

Anqi Mao¹ Mehryar Mohri² Yutao Zhong¹

Abstract

We present a detailed study of \mathcal{H} -consistency bounds for score-based ranking. These are upper bounds on the target loss estimation error of a predictor in a hypothesis set \mathcal{H} , expressed in terms of the surrogate loss estimation error of that predictor. We will show that both in the *general pairwise ranking* scenario and in the *bipartite ranking* scenario, there are no meaningful \mathcal{H} -consistency bounds for most hypothesis sets used in practice including the family of linear models and that of the neural networks, which satisfy the equicontinuous property with respect to the input. To come up with ranking surrogate losses with theoretical guarantees, we show that a natural solution consists of resorting to a *pairwise abstention loss* in the general pairwise ranking scenario, and similarly, a *bipartite abstention loss* in the bipartite ranking scenario, to abstain from making predictions at some limited cost c . For surrogate losses of these abstention loss functions, we give a series of \mathcal{H} -consistency bounds for both the family of linear functions and that of neural networks with one hidden-layer. Our experimental results illustrate the effectiveness of ranking with abstention.

1. Introduction

In many applications, ranking is a more appropriate formulation of the learning task than classification, given the crucial significance of the ordering of the items. As an example, for movie recommendation systems, an ordered list of movies is preferable to a comprehensive list of recommended titles, since users are more likely to watch those ranked highest.

The problem of learning to rank has been studied in a large number of publications. Ailon & Mohri (2008; 2010) distin-

guish two general formulations of the problem: the *score-based setting* and the *preference-based setting*. In the score-based setting, a real-valued function over the input space is learned, whose values determine a total ordering of all input points. In the preference-based setting, a pairwise preference function is first learned, typically by training a classifier over a sample of labeled pairs; next, that function is used to derive an ordering, potentially randomized, of any subset of points.

This paper deals with the score-based ranking formulation both in the general ranking setting, where items are not assigned any specific category, and the bipartite setting, where they are labeled with one of two classes. The evaluation of a ranking solution in this context is based on the average pairwise misranking metric. In the bipartite setting, this metric is directly related to the AUC (Area Under the ROC Curve), which coincides with the average correct pairwise ranking (Hanley & McNeil, 1982; Cortes & Mohri, 2003), also known as the Wilcoxon-Mann-Whitney statistic.

For most hypothesis sets, directly optimizing the pairwise misranking loss is intractable. Instead, ranking algorithms resort to a surrogate loss. As an example, the surrogate loss for RankBoost (Freund et al., 2003; Rudin et al., 2005) is based on the exponential function and that of SVM ranking (Joachims, 2002) on the hinge loss. But, what guarantees can we rely on when minimizing a surrogate loss instead of the original pairwise misranking loss?

The property often invoked in this context is *Bayes consistency*, which has been extensively studied for classification (Zhang, 2004; Bartlett et al., 2006; Tewari & Bartlett, 2007). The Bayes consistency of ranking surrogate losses has been studied in the special case of bipartite ranking: in particular, Uematsu & Lee (2017) proved the inconsistency of the pairwise ranking loss based on the hinge loss and Gao & Zhou (2015) gave excess loss bounds for pairwise ranking losses based on the exponential or the logistic loss (see also (Menon & Williamson, 2014)). A related but distinct consistency question has been studied in several publications (Agarwal et al., 2005; Kotlowski et al., 2011; Agarwal, 2014). It is one with respect to binary classification, that is whether a near minimizer of the surrogate loss of the binary classification loss is a near minimizer of the bipartite misranking loss (Cortes & Mohri, 2003).

¹Courant Institute of Mathematical Sciences, New York, NY;
²Google Research, New York, NY. Correspondence to: Anqi Mao <aqmao@cims.nyu.edu>, Mehryar Mohri <mohri@google.com>, Yutao Zhong <yutao@cims.nyu.edu>.

However, as recently argued by Awasthi, Mao, Mohri, and Zhong (2022a), Bayes consistency is not a sufficiently informative notion since it only applies to the entire class of measurable functions and does not hold for specific subsets, such as sub-families of linear functions or neural networks. Furthermore, Bayes consistency is solely an asymptotic concept and does not offer insights into the performance of predictors trained on finite samples. In response, the authors proposed an alternative concept called \mathcal{H} -consistency bounds, which provide non-asymptotic guarantees tailored to a given hypothesis set \mathcal{H} . They proceeded to establish such bounds within the context of classification both in binary and multi-class classification (Awasthi et al., 2022a;b), see also (Mao et al., 2023). These are stronger and more informative guarantees than Bayes consistency.

But, can we derive \mathcal{H} -consistency guarantees for ranking? In other words, can we find surrogate losses for pairwise misranking whose approximate estimation error minimization guarantees approximate minimization of the pairwise or bipartite misranking loss? We will show that, surprisingly, this is not possible for most hypothesis sets used in practice, including the family of constrained linear models or that of the constrained neural networks, or any family of equicontinuous functions with respect to the input. In fact, we will give a relatively simple example where the pairwise misranking error of the RankBoost algorithm remains significant, even after training with relatively large sample sizes. How can we then come up with ranking surrogate losses with theoretical guarantees?

We will show that a natural solution consists of resorting to a *pairwise abstention loss* in the general pairwise ranking scenario and, similarly, a *bipartite abstention loss* in the bipartite ranking scenario, to abstain from making predictions at some limited cost c . For surrogate losses of these abstention loss functions, we give a series of \mathcal{H} -consistency bounds for both the family of linear functions and that of neural networks with one hidden-layer. A key term appearing in these bounds is the *minimizability gap*, which measures the difference between the best-in-class expected loss and the expected infimum of the pointwise expected loss. This plays a crucial role in these bounds and we give a detailed analysis of these terms. We also present the results of experiments illustrating the effectiveness of ranking with abstention.

Comparison with previous work. Here, we briefly discuss the relationship of prior work on \mathcal{H} -consistency bounds (Awasthi et al., 2022a;b; Mao et al., 2023; Zheng et al., 2023) with ours. Awasthi et al. (2022a) studied \mathcal{H} -consistency bounds in the binary classification setting. They provided a series of positive results for common binary surrogate losses with the hypothesis sets of linear models and one-hidden-layer ReLU networks. Awasthi et al. (2022b); Mao et al.

(2023); Zheng et al. (2023) studied \mathcal{H} -consistency bounds in the context of multi-class classification. Awasthi et al. (2022b) provided an extensive analysis of \mathcal{H} -consistency bounds for multi-class *max losses* such as those of Cramer & Singer (2001), *sum losses* such as that of Weston & Watkins (1998), and *constrained losses*, such as the loss function adopted by Lee et al. (2004) in the analysis of multi-class SVM. They further gave the analysis of \mathcal{H} -consistency bounds for all these multi-class losses in the adversarial setting. More recently, Mao et al. (2023) presented a theoretical analysis of a broad family of loss functions, *comp-sum losses*, that includes *cross-entropy (or logistic loss)* (Verhulst, 1838; 1845; Berkson, 1944; 1951), *generalized cross-entropy* (Zhang & Sabuncu, 2018), the *mean absolute error* (Ghosh et al., 2017) and other cross-entropy-like loss functions. They gave tight \mathcal{H} -consistency bounds for these loss functions with any *complete* hypothesis set. They also introduced new *smooth adversarial comp-sum losses* in the adversarial setting and proved \mathcal{H} -consistency bounds guarantees for these loss functions. Zheng et al. (2023) also provided \mathcal{H} -consistency bounds for the logistic loss with the hypothesis set of linear models, under some distributional assumptions. They used these bounds to compare multi-class logistic regression and naive Bayes methods.

Our paper primarily concentrates on score-based ranking with a binary label space, setting it apart from the multi-class scenario (Awasthi et al., 2022b; Mao et al., 2023; Zheng et al., 2023). The primary technical differences and challenges between the ranking and binary classification settings (Awasthi et al., 2022a) stem from the fundamental distinction that ranking loss functions take as argument a pair of samples rather than a single one, as is the case for binary classification loss functions. This makes it more challenging to derive \mathcal{H} -consistency bounds, as upper bounding the calibration gap of the target loss by that of the surrogate loss becomes technically more difficult.

Additionally, this fundamental difference leads to a negative result for ranking, as \mathcal{H} -consistency bounds cannot be guaranteed for most commonly used hypothesis sets, including the family of constrained linear models and that of constrained neural networks, both of which satisfy the equicontinuity property concerning the input. As a result, a natural alternative involves using ranking with abstention, for which \mathcal{H} -consistency bounds can be proven. In the abstention setting, an extra challenge lies in carefully monitoring the effect of a threshold γ to relate the calibration gap of the target loss to that of the surrogate loss.

Furthermore, the bipartite ranking setting introduces an added layer of complexity, as each element of a pair of samples has an independent conditional distribution, which results in a more intricate calibration gap.

Structure of the paper. The remaining sections of this pa-

per are organized as follows. In Section 2, we study general pairwise ranking. We first prove several negative results in Section 2.2 showing that there exists no meaningful \mathcal{H} -consistency bound for general surrogate loss functions with an equicontinuous hypothesis set \mathcal{H} . We then present a series of positive results by considering a family of piecewise continuous functions (Section 2.3), which we will simply refer to as piecewise functions, and the family of all measurable functions (Section 2.4). In Section 3, we study general pairwise ranking with abstention. We provide a series of explicit \mathcal{H} -consistency bounds in the case of the pairwise abstention loss, with multiple choices of the surrogate loss and for both the family of linear functions (Section 3.1) and that of neural networks with one hidden-layer (Section 3.2). We also study bipartite ranking in Section 4. Here too, we provide both negative results with general hypothesis sets (Section 4.2) and positive results with the family of all measurable functions (Section 4.3). We then present \mathcal{H} -consistency bounds for bipartite ranking with abstention in Section 5, for linear hypothesis sets (Section 5.1) and the family of neural networks with one hidden-layer (Section 5.2). In Section 6, we report the results of experiments illustrating the effectiveness of ranking with abstention.

We give a detailed discussion of related work in Appendix A.

2. General Pairwise Ranking

In this section, we analyze the properties of surrogate losses for the general pairwise misranking loss. We begin by introducing the necessary definitions and concepts. Next, we present negative results demonstrating the impossibility of deriving non-trivial \mathcal{H} -consistency bounds for widely used surrogate losses and hypothesis sets. In contrast, we give positive results for a family of piecewise functions and that of all measurable functions.

2.1. Preliminaries

We study the learning scenario of score-based ranking in the *general pairwise ranking* scenario (e.g. see (Mohri et al., 2018)). Let \mathcal{X} denote the input space and $\mathcal{Y} = \{-1, +1\}$ the label space. We denote by \mathcal{H} a hypothesis set of functions mapping from \mathcal{X} to \mathbb{R} . The *general pairwise misranking loss* L_{0-1} is defined for all h in \mathcal{H} , x, x' in \mathcal{X} and y in \mathcal{Y} by

$$L_{0-1}(h, x, x', y) = \mathbb{1}_{y \neq \text{sign}(h(x') - h(x))}, \quad (1)$$

where $\text{sign}(u) = \mathbb{1}_{u \geq 0} - \mathbb{1}_{u < 0}$. Thus, h incurs a loss of one on the labeled pair (x, x', y) when it ranks the pair (x, x') opposite to the sign of y , where, by convention, x' is considered as ranked above x when $h(x') \geq h(x)$. Otherwise, the loss incurred is zero. Optimizing the pairwise misranking loss L_{0-1} is intractable for most hypothesis sets. Thus, general ranking algorithms rely on a surrogate loss function L instead of L_{0-1} . We will analyze the properties

of such surrogate loss functions.

Let \mathcal{D} denote a distribution over $\mathcal{X} \times \mathcal{X} \times \mathcal{Y}$. We denote by $\eta(x, x') = \mathcal{D}(Y = +1 \mid (X, X') = (x, x'))$ the conditional probability of $Y = +1$ given $(X, X') = (x, x')$. We also denote by $\mathcal{R}_L(h)$ the *expected L-loss* of a hypothesis h and by $\mathcal{R}_L^*(\mathcal{H})$ its infimum over \mathcal{H} :

$$\mathcal{R}_L(h) = \mathbb{E}_{(x, x', y) \sim \mathcal{D}} [L(h, x, x', y)] \quad \mathcal{R}_L^*(\mathcal{H}) = \inf_{h \in \mathcal{H}} \mathcal{R}_L(h)$$

\mathcal{H} -consistency bounds. We will analyze the *\mathcal{H} -consistency bounds* properties (Awasthi et al., 2022a) of such surrogate loss functions. An \mathcal{H} -consistency bound for a surrogate loss L is a guarantee of the form:

$$\forall h \in \mathcal{H}, \quad \mathcal{R}_{L_{0-1}}(h) - \mathcal{R}_{L_{0-1}}^*(\mathcal{H}) \leq f(\mathcal{R}_L(h) - \mathcal{R}_L^*(\mathcal{H})),$$

for some non-decreasing function $f: \mathbb{R}_+ \rightarrow \mathbb{R}_+$. This provides a quantitative relationship between the estimation loss of L_{0-1} and that of the surrogate loss L . The guarantee is stronger and more informative than Bayes consistency, or \mathcal{H} -consistency, \mathcal{H} -calibration or the excess error bounds (Zhang, 2004; Bartlett et al., 2006; Steinwart, 2007; Mohri et al., 2018) discussed in the literature.

A key quantity appearing in \mathcal{H} -consistency bounds is the *minimizability gap*, which is the difference between the best-in-class expected loss and the expected pointwise infimum of the loss:

$$\mathcal{M}_L(\mathcal{H}) = \mathcal{R}_L^*(\mathcal{H}) - \mathbb{E}_{(x, x')} \left[\inf_{h \in \mathcal{H}} \mathbb{E}_y [L(h, x, x', y) \mid (x, x')] \right].$$

By the super-additivity of the infimum, the minimizability gap is always non-negative.

We will specifically study the hypothesis set of all measurable functions, \mathcal{H}_{all} ; that of linear hypotheses, $\mathcal{H}_{\text{lin}} = \{x \mapsto w \cdot x + b \mid \|w\|_q \leq W, |b| \leq B\}$; and the hypothesis set of one-hidden-layer ReLU networks: $\mathcal{H}_{\text{NN}} = \{x \mapsto \sum_{j=1}^n u_j (w_j \cdot x + b_j)_+ \mid \|u\|_1 \leq \Lambda, \|w_j\|_q \leq W, |b_j| \leq B\}$, where $(\cdot)_+ = \max(\cdot, 0)$. A table of notation (Table 7) is presented in Appendix B. We will say that a hypothesis set is *regular for general pairwise ranking* if, for any $x \neq x' \in \mathcal{X}$, we have $\{\text{sign}(h(x') - h(x)) : h \in \mathcal{H}\} = \{-1, +1\}$. Hypothesis sets commonly used in practice all admit this property.

2.2. Negative Results

Here, we give a negative result for a broad family of surrogate losses and hypothesis sets.

The general pairwise ranking surrogate losses widely used in practice admit the following form:

$$L_\Phi(h, x, x', y) = \Phi(y(h(x') - h(x))), \quad (2)$$

where Φ is a non-increasing function that is continuous at 0 and upper bounding $u \mapsto \mathbb{1}_{u \leq 0}$ over \mathbb{R} . The following

result shows that these surrogate losses do not benefit from a non-trivial \mathcal{H} -consistency bound when the hypothesis set used is equicontinuous, which includes most hypothesis sets used in practice, in particular the family of linear hypotheses and that of neural networks.

Theorem 2.1 (Negative results). *Assume that \mathcal{X} contains an interior point x_0 and that \mathcal{H} is regular for general pairwise ranking, contains 0 and is equicontinuous at x_0 . If for some function f that is non-decreasing and continuous at 0, the following bound holds for all $h \in \mathcal{H}$ and any distribution,*

$$\mathcal{R}_{L_{0-1}}(h) - \mathcal{R}_{L_{0-1}}^*(\mathcal{H}) \leq f(\mathcal{R}_{L_\Phi}(h) - \mathcal{R}_{L_\Phi}^*(\mathcal{H})),$$

then, $f(t) \geq 1$ for any $t \geq 0$.

Theorem 2.1 shows that for equicontinuous hypothesis sets, any \mathcal{H} -consistency bound is vacuous, assuming that f is a non-decreasing function continuous at zero. This is because for any such bound, a small L_Φ -estimation loss does not guarantee a small L_{0-1} -estimation loss, as the right-hand side remains lower-bounded by one.

The proof is given in Appendix D, where we give a simple example on pairs whose distance is relatively small for which the standard surrogate losses including the RankBoost algorithm (L_{exp}) fail (see also Section 6). It is straightforward to see that the assumptions of Theorem 2.1 hold for the case $\mathcal{H} = \mathcal{H}_{\text{lin}}$ or $\mathcal{H} = \mathcal{H}_{\text{NN}}$. Indeed, we can take $x_0 = 0$ as the interior point and thus for any $h \in \mathcal{H}_{\text{lin}}$, $|h(x) - h(x_0)| = |w \cdot x| < \epsilon$ for any $x \in \{x \in \mathcal{X} : \|x\|_p < \frac{\epsilon}{W}\}$, which implies that \mathcal{H}_{lin} is equicontinuous at x_0 . As with the linear hypothesis set, for any $h \in \mathcal{H}_{\text{NN}}$, $|h(x) - h(x_0)| = |\sum_{j=1}^n u_j (w_j \cdot x + b_j)_+ - \sum_{j=1}^n u_j (b_j)_+| = |\sum_{j=1}^n u_j [(w_j \cdot x + b_j)_+ - (b_j)_+]| \leq \Lambda W \|x\|_p < \epsilon$, for any $x \in \{x \in \mathcal{X} : \|x\|_p < \frac{\epsilon}{\Lambda W}\}$, which implies that \mathcal{H}_{NN} is equicontinuous at x_0 . In fact, Theorem 2.1 holds for any family of Lipschitz constrained neural networks, since a family of functions that share the same Lipschitz constant is equicontinuous.

It is straightforward to verify that the proof of Theorem 2.1 also holds in the deterministic case where $\eta(x, x')$ equals 0 or 1 for any $x \neq x'$, which yields the following corollary.

Corollary 2.2 (Negative results in the deterministic case). *In the deterministic case where $\eta(x, x')$ equals 0 or 1 for any $x \neq x'$, the negative result of Theorem 2.1 still holds.*

2.3. Positive Results: A Family of Piecewise Functions

In this section, we seek alternative positive results. In light of the negative results just presented, we need to consider hypothesis sets that are not equicontinuous. A natural choice

is a family of piecewise functions. For any fixed parameter $\tau > 0$, we consider the family of piecewise functions $\mathcal{H}_{\text{pw}} = \{u \mapsto \alpha(\mathbb{1}_{u \neq x \wedge \|u\| > \tau} - \mathbb{1}_{u \neq x \wedge \|u\| \leq \tau}) \mid x \in \mathcal{X}, \alpha \in \mathbb{R}\}$. Then, we have the following positive result in the deterministic setting.

Theorem 2.3 (Positive results for piecewise functions). *Assume that Φ satisfies $\lim_{u \rightarrow +\infty} \Phi(u) = 0$. Then, for all $h \in \mathcal{H}_{\text{pw}}$ and any deterministic distribution,*

$$\begin{aligned} & \mathcal{R}_{L_{0-1}}(h) - \mathcal{R}_{L_{0-1}}^*(\mathcal{H}_{\text{pw}}) \\ & \leq \mathcal{R}_{L_\Phi}(h) - \mathcal{R}_{L_\Phi}^*(\mathcal{H}_{\text{pw}}) + \mathcal{M}_{L_\Phi}(\mathcal{H}_{\text{pw}}) - \mathcal{M}_{L_{0-1}}(\mathcal{H}_{\text{pw}}). \end{aligned}$$

The proof is included in Appendix E. Theorem 2.3 provides a meaningful \mathcal{H} -consistency bound for the hypothesis set of piecewise functions: modulo the minimizability gaps, which are zero when the best-in-class error coincides with the Bayes error or can be small in some other cases, reducing the surrogate estimation loss appearing on the right-hand side guarantees a small target estimation loss (left-hand side).

One example of the corresponding hypotheses in \mathcal{H}_{pw} is $u \mapsto \mathbb{1}_{\|u\| > \tau} - \mathbb{1}_{0 < \|u\| \leq \tau}$, where $\alpha = 1$ and $x = 0$. This is a family of piecewise functions based on the magnitude of u in two distinct ranges. Note that such a hypothesis set is not typical and in particular does not admit equicontinuity, a necessary condition according to our negative results. Nevertheless, it provides an example of a hypothesis set supported by \mathcal{H} -consistency bounds in the general ranking setting. We will further show in section 3 that, for a standard hypothesis set such as an equicontinuous one, we need to resort to ranking with abstention. This approach will be the primary focus of the positive results presented in our paper, as it allows us to leverage hypothesis sets that are more commonly used in real-world applications.

2.4. Positive Results: \mathcal{H}_{all} -Consistency Bounds

In this section, we also present a series of positive results in the case of the family of all measurable functions, \mathcal{H}_{all} .

We prove \mathcal{H}_{all} -consistency bounds for the surrogate loss L_Φ when using as auxiliary function Φ the hinge-loss, the ρ -margin loss, the exponential loss, the logistic loss, the squared hinge-loss, or the sigmoid loss defined in Table 1. Table 8 (Appendix F) gives the full expression of our \mathcal{H}_{all} -consistency upper bounds, with detailed proofs given in Appendix F. Table 1 gives the expression of the corresponding minimizability gaps. As an example, we have

$$\begin{aligned} & \mathcal{M}_{L_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{all}}) \\ & = \mathcal{R}_{L_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[2\sqrt{\eta(x, x')(1 - \eta(x, x'))} \right]. \end{aligned}$$

In Appendix G, we show that these minimizability gaps are in general not null for $\mathcal{H} = \mathcal{H}_{\text{all}}$ in general pairwise ranking,

Table 1: Auxiliary functions and the minimizability gaps of their pairwise ranking and bipartite ranking surrogates.

Auxiliary Functions	Definitions	$\mathcal{M}_{L_\Phi}(\mathcal{H}_{\text{all}})$	$\mathcal{M}_{L_\Phi}(\mathcal{H}_{\text{all}})$
Hinge	$\Phi_{\text{hinge}}(t) = \max\{0, 1 - t\}$	(12)	(25)
ρ -Margin	$\Phi_\rho(t) = \min\left\{1, \max\left\{0, 1 - \frac{t}{\rho}\right\}\right\}, \rho > 0$	(14)	(27)
Exponential	$\Phi_{\text{exp}}(t) = e^{-t}$	(16)	(29)
Logistic	$\Phi_{\text{log}}(t) = \log_2(1 + e^{-t})$	(18)	(31)
Squared hinge	$\Phi_{\text{sq}}(t) = (1 - t)^2 \mathbb{1}_{t \leq 1}$	(20)	(33)
Sigmoid	$\Phi_{\text{sig}}(t) = 1 - \tanh(kt), k > 0$	(22)	(35)

in contrast with binary classification where the minimizability gaps for \mathcal{H}_{all} are zero. This is because the *distribution order* for general pairwise ranking cannot always be induced by a real-valued function. We first introduce the definition of that order. Next, we characterize the distribution order that can be induced by a predictor h , which leads to the zero minimizability gap of pairwise misranking loss.

Definition 2.4. The *distribution order* is a homogeneous relation $\stackrel{\mathcal{D}}{\leq}$ over \mathcal{X} , defined as follows for all $x, x' \in \mathcal{X}$,

$$x \stackrel{\mathcal{D}}{\leq} x' \iff \eta(x, x') \geq \eta(x', x).$$

We say that a hypothesis h induces the distribution order if, for all $x, x' \in \mathcal{X}$, $(h(x) \leq h(x'))$ holds iff $(x \stackrel{\mathcal{D}}{\leq} x')$. We say that a subset $\tilde{\mathcal{X}} \subset \mathcal{X}$ is *dense countable* in \mathcal{X} with respect to the distribution order if $\tilde{\mathcal{X}}$ is countable and, for all $x, x' \in \mathcal{X}$ satisfying $x \stackrel{\mathcal{D}}{\leq} x'$ and not $x' \stackrel{\mathcal{D}}{\leq} x$, there exists $\bar{x} \in \tilde{\mathcal{X}}$ such that $x \stackrel{\mathcal{D}}{\leq} \bar{x} \stackrel{\mathcal{D}}{\leq} x'$.

The following result characterizes the distribution order induced by a hypothesis h .

Theorem 2.5 (Characterization of distribution order). *The distribution order is transitive and there exists a dense countable subset $\tilde{\mathcal{X}} \subset \mathcal{X}$ with respect to the distribution order if and only if there exists $h \in \mathcal{H}_{\text{all}}$ inducing the distribution order.*

A special case of Theorem 2.5 is when the distribution order is a total order and $\eta(x, x')$ is continuous.

Theorem 2.6. *Assume that the distribution order is a total order and $\eta(x, x')$ is continuous on $\mathcal{X} \times \mathcal{X}$. Then, there exists $h \in \mathcal{H}_{\text{all}}$ inducing the distribution order.*

Theorem 2.5 and Theorem 2.6 characterize cases where the distribution order can be induced by a hypothesis. These results actually characterize the case where the minimizability gap of the pairwise misranking loss is null, since we will see immediately that, when the distribution order is induced by a hypothesis $h \in \mathcal{H}$, the minimizability gaps of the pairwise misranking loss will be zero.

Theorem 2.7. *Assume that for all $x, x' \in \mathcal{X}$, $\eta(x, x') + \eta(x', x) = 1$. Then, for any hypothesis set \mathcal{H} , if there exists*

$h \in \mathcal{H}$ inducing the distribution order, the minimizability gap of the pairwise misranking loss is null, $\mathcal{M}_{L_{0-1}}(\mathcal{H}) = 0$.

The proofs of Theorems 2.5, 2.6 and 2.7 are presented in Appendix H. These results provide a detailed analysis of the distribution order and the minimizability gaps appearing in the \mathcal{H}_{all} -consistency bounds in general pairwise misranking. However, learning algorithms are not based on \mathcal{H}_{all} . Instead, they rely on a restricted hypothesis set such as \mathcal{H}_{lin} or \mathcal{H}_{NN} . To come up with ranking surrogate losses with theoretical guarantees for such hypothesis sets, a natural solution consists of resorting to a pairwise abstention loss.

3. General Pairwise Ranking with Abstention

The negative results of the previous section suggest that general pairwise ranking with theoretical guarantees is difficult with common hypothesis sets. The inherent issue for pairwise ranking is that for equicontinuous hypotheses, when x and x' are arbitrarily close, the confidence value $|h(x) - h(x')|$ can be arbitrary close to zero. This motivates us to study the learning scenario of *general pairwise ranking with abstention*.

In this scenario, the learner abstains from making a prediction on input pair (x, x') if the distance between x' and x is relatively small, in which case a cost c is incurred. Let $\|\cdot\|$ denote the norm adopted, which is typically an ℓ_p -norm, $p \in [1, +\infty]$. The *pairwise abstention loss* is defined as follows for any $h \in \mathcal{H}$ and $(x, x', y) \in \mathcal{X} \times \mathcal{X} \times \mathcal{Y}$:

$$\begin{aligned} L_{0-1}^{\text{abs}}(h, x, x', y) &= \mathbb{1}_{y \neq \text{sign}(h(x') - h(x))} \mathbb{1}_{\|x - x'\| > \gamma} + c \mathbb{1}_{\|x - x'\| \leq \gamma}, \end{aligned} \quad (3)$$

where γ is a given threshold value. For $\gamma = 0$, L_{0-1}^{abs} reduces to the pairwise misranking loss L_{0-1} without abstention. Let $p, q \in [1, +\infty]$ be conjugate numbers, that is $\frac{1}{p} + \frac{1}{q} = 1$. Without loss of generality, we consider $\mathcal{X} = B_p^d(1)$ and $\|\cdot\|$ in (3) to be the ℓ_p norm. The corresponding conjugate ℓ_q norm is adopted in the hypothesis sets \mathcal{H}_{lin} and \mathcal{H}_{NN} . In the following, we will prove \mathcal{H} -consistency bounds for L_Φ when using as an auxiliary function Φ the hinge-loss, the ρ -margin loss, the exponential loss, the logistic loss, the squared hinge-loss or the sigmoid loss with respect to

Table 2: \mathcal{H}_{lin} -consistency upper bounds for general pairwise abstention.

Loss function	\mathcal{H}_{lin} -consistency upper bound
$\mathcal{L}_{\Phi_{\text{hinge}}}$	$\frac{\mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}^*}(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{lin}})}{\min\{W\gamma, 1\}} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$
$\mathcal{L}_{\Phi_{\rho}}$	$\frac{\rho(\mathcal{R}_{\mathcal{L}_{\Phi_{\rho}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\rho}}}(\mathcal{H}_{\text{lin}}))}{\min\{W\gamma, \rho\}} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$
$\mathcal{L}_{\Phi_{\text{exp}}}$	$\Gamma_{\text{exp}}(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}^*}(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{lin}})) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$ where $\Gamma_{\text{exp}}(t) = \max\left\{\sqrt{2t}, 2\left(\frac{e^{2W\gamma} + 1}{e^{2W\gamma} - 1}\right)t\right\}$
$\mathcal{L}_{\Phi_{\log}}$	$\Gamma_{\log}(\mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\log}}}(\mathcal{H}_{\text{lin}})) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$ where $\Gamma_{\log}(t) = \max\left\{\sqrt{2t}, 2\left(\frac{e^{W\gamma} + 1}{e^{W\gamma} - 1}\right)t\right\}$
$\mathcal{L}_{\Phi_{\text{sq}}}$	$\Gamma_{\text{sq}}(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}^*}(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{lin}})) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$ where $\Gamma_{\text{sq}}(t) = \max\left\{\sqrt{t}, \frac{t}{2W\gamma} + \frac{W\gamma}{2}\right\}$
$\mathcal{L}_{\Phi_{\text{sig}}}$	$\frac{\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}^*}(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{lin}})}{\tanh(kW\gamma)} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$

 Table 3: \mathcal{H}_{NN} -consistency upper bounds for general pairwise abstention.

Loss function	\mathcal{H}_{NN} -consistency upper bound
$\mathcal{L}_{\Phi_{\text{hinge}}}$	$\frac{\mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}^*}(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{NN}})}{\min\{\Lambda W\gamma, 1\}} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$
$\mathcal{L}_{\Phi_{\rho}}$	$\frac{\rho(\mathcal{R}_{\mathcal{L}_{\Phi_{\rho}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\rho}}}(\mathcal{H}_{\text{NN}}))}{\min\{\Lambda W\gamma, \rho\}} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$
$\mathcal{L}_{\Phi_{\text{exp}}}$	$\Gamma_{\text{exp}}(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}^*}(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{NN}})) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$ where $\Gamma_{\text{exp}}(t) = \max\left\{\sqrt{2t}, 2\left(\frac{e^{2\Lambda W\gamma} + 1}{e^{2\Lambda W\gamma} - 1}\right)t\right\}$
$\mathcal{L}_{\Phi_{\log}}$	$\Gamma_{\log}(\mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\log}}}(\mathcal{H}_{\text{NN}})) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$ where $\Gamma_{\log}(t) = \max\left\{\sqrt{2t}, 2\left(\frac{e^{\Lambda W\gamma} + 1}{e^{\Lambda W\gamma} - 1}\right)t\right\}$
$\mathcal{L}_{\Phi_{\text{sq}}}$	$\Gamma_{\text{sq}}(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}^*}(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{NN}})) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$ where $\Gamma_{\text{sq}}(t) = \max\left\{\sqrt{t}, \frac{t}{2\Lambda W\gamma} + \frac{\Lambda W\gamma}{2}\right\}$
$\mathcal{L}_{\Phi_{\text{sig}}}$	$\frac{\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}^*}(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{NN}})}{\tanh(k\Lambda W\gamma)} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$

$\mathcal{L}_{0-1}^{\text{abs}}$, in the case of the linear hypothesis set \mathcal{H}_{lin} or that of one-hidden-layer ReLU networks \mathcal{H}_{NN} .

3.1. Linear Hypotheses

Table 2 supplies the \mathcal{H}_{lin} -consistency upper bounds for \mathcal{L}_{Φ} when using as Φ the auxiliary functions in Table 1. The bounds of Table 2 depend directly on the threshold value γ , the parameter W in the linear models and parameters of the loss function (e.g., k in sigmoid loss).

As an example, when using as Φ the exponential loss function, modulo the minimizability gaps (which are zero when the best-in-class error coincides with the Bayes error or can be small in some other cases), the bound implies that if the surrogate estimation loss $\mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}^*}(\mathcal{H}_{\text{lin}})$ is reduced to ϵ , then, the target estimation loss $\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}})$ is upper bounded by $\Gamma_{\text{exp}}(\epsilon)$. For sufficiently small values of ϵ , the dependence of Γ_{exp} on ϵ exhibits a square root relationship. However, if this is not the case, the dependence becomes linear, subject to a constant factor depending on the threshold value γ and the parameter W in the linear models.

The proofs consist of analyzing calibration gaps of the target loss and that of each surrogate loss and seeking a tight lower bound of the surrogate calibration gap in terms of the target

one. As an example, for $\Phi = \Phi_{\text{exp}}$, we have the tight lower bound $\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \geq \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}(h_0, x, x') = \Psi_{\text{exp}}\left(\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}}}(h, x, x')\right)$, where h_0 can be the null hypothesis when $\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}}}(h, x, x') \neq 0$ and Ψ_{exp} is an increasing and piecewise convex function on $[0, 1]$ defined by $\Psi_{\text{exp}}(t) = \begin{cases} 1 - \sqrt{1 - t^2}, & t \leq \frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \\ 1 - \frac{t+1}{2}e^{-W\gamma} - \frac{1-t}{2}e^{W\gamma}, & t > \frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \end{cases}$. The detailed derivation and the expression of the corresponding minimizability gaps are included in Appendix K.1.

3.2. One-hidden-layer ReLU Neural Networks

Table 3 gives \mathcal{H}_{NN} -consistency upper bounds for \mathcal{L}_{Φ} when using as Φ the auxiliary functions in Table 1. Different from the bounds in the linear case, all the bounds in Table 3 not only depend on W , but also depend on Λ , which is a parameter appearing in \mathcal{H}_{NN} . The proof is similar to that of the linear case. The detailed derivation and the expression of the corresponding minimizability gaps are given in Appendix K.2.

As with the linear case, taking the exponential loss function Φ_{exp} as an example, modulo the minimizability gaps, the bound implies that if the surrogate estimation loss $\mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}^*}(\mathcal{H}_{\text{NN}})$ is reduced to ϵ , then, the target estimation loss $\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}})$ is upper bounded

by $\Gamma_{\text{exp}}(\epsilon)$. For sufficiently small values of ϵ , the dependence of Γ_{exp} on ϵ exhibits a square root relationship. However, if this is not the case, the dependence becomes linear, subject to a constant factor depending on the threshold value γ , the parameter W , and an additional parameter Λ in the one-hidden-layer ReLU networks.

4. Bipartite Ranking

In this section, we analyze the properties of surrogate loss functions in the bipartite ranking setting. We first introduce the relevant definitions and concepts. Next, as in the general pairwise ranking setting, we present general negative results, as well as positive results for the family of all measurable functions.

4.1. Preliminaries

In the bipartite setting, each point x admits a label $y \in \{-1, +1\}$. The *bipartite misranking loss* $\tilde{\mathcal{L}}_{0-1}$ is defined for all h in \mathcal{H} , and $(x, y), (x', y')$ in $(\mathcal{X} \times \mathcal{Y})$ by

$$\tilde{\mathcal{L}}_{0-1}(h, x, x', y, y') = \mathbb{1}_{(y-y')(h(x)-h(x')) < 0} + \frac{1}{2} \mathbb{1}_{(h(x)=h(x')) \wedge (y \neq y')}. \quad (4)$$

Optimizing the bipartite misranking loss $\tilde{\mathcal{L}}_{0-1}$ is intractable for most hypothesis sets and bipartite ranking algorithms rely instead on a surrogate loss $\tilde{\mathcal{L}}$. We will analyze the properties of such surrogate losses.

Let \mathcal{D} be a distribution over $\mathcal{X} \times \mathcal{Y}$. We denote by $\eta(x) = \mathcal{D}(Y = +1 \mid X = x)$ the conditional probability of $Y = +1$ given $X = x$. We will use a definition and notation for the expected $\tilde{\mathcal{L}}$ -loss of $h \in \mathcal{H}$, its infimum, and the minimizability gaps similar to what we used in the general pairwise misranking setting:

$$\begin{aligned} \mathcal{R}_{\tilde{\mathcal{L}}}(h) &= \mathbb{E}_{(x, x', y) \sim \mathcal{D}} [\tilde{\mathcal{L}}(h, x, x', y)] & \mathcal{R}_{\tilde{\mathcal{L}}}^*(\mathcal{H}) &= \inf_{h \in \mathcal{H}} \mathcal{R}_{\tilde{\mathcal{L}}}(h) \\ \mathcal{M}_{\tilde{\mathcal{L}}}(\mathcal{H}) & & & \\ &= \mathcal{R}_{\tilde{\mathcal{L}}}^*(\mathcal{H}) - \mathbb{E}_{(x, x')} \left[\inf_{h \in \mathcal{H}} \mathbb{E}_{(y, y')} [\tilde{\mathcal{L}}(h, x, x', y, y') \mid (x, x')] \right]. \end{aligned}$$

We say that a hypothesis set is *regular for bipartite ranking* if, for any $x \neq x' \in \mathcal{X}$, there exists $h_+ \in \mathcal{H}$ such that $h_+(x) < h_+(x')$ and $h_- \in \mathcal{H}$ such that $h_-(x) > h_-(x')$. Hypothesis sets commonly used in practice all admit this property.

4.2. Negative Results

Here, as in general pairwise misranking scenario, we present a negative result for a broad family of surrogate losses and hypothesis sets in the bipartite setting.

The bipartite ranking surrogate losses widely used in prac-

tice, admit the following form:

$$\tilde{\mathcal{L}}_{\Phi}(h, x, x', y, y') = \Phi\left(\frac{(y-y')(h(x)-h(x'))}{2}\right) \mathbb{1}_{y \neq y'}, \quad (5)$$

where Φ is a non-increasing function that is continuous at 0 upper bounding $u \mapsto \mathbb{1}_{u \leq 0}$ over \mathbb{R} . As with the general pairwise ranking, we show that these surrogate losses do not benefit from \mathcal{H} -consistency bounds when \mathcal{H} is an equicontinuous family.

Theorem 4.1 (Negative results for bipartite ranking). *Assume that \mathcal{X} contains an interior point x_0 and that \mathcal{H} is regular for bipartite ranking, contains 0 and is equicontinuous at x_0 . If for some function f that is non-decreasing and continuous at 0, the following bound holds for all $h \in \mathcal{H}$ and any distribution,*

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}^*(\mathcal{H}) \leq f\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi}}^*(\mathcal{H})\right),$$

then, $f(t) \geq \frac{1}{2}$ for any $t \geq 0$.

As with Theorem 2.1, Theorem 4.1 shows that in the bipartite ranking setting, any \mathcal{H} -consistency bound with an equicontinuous hypothesis sets is vacuous, assuming a non-decreasing function f continuous at zero.

The proof is given in Appendix I. It is straightforward to verify that the proof holds in the deterministic case where $\eta(x)$ equals 0 or 1 for any $x \in \mathcal{X}$, which yields the following corollary.

Corollary 4.2 (Negative results in the bipartite deterministic case). *In the bipartite deterministic case where $\eta(x)$ equals 0 or 1 for any $x \in \mathcal{X}$, the same negative result as in Theorem 4.1 holds.*

4.3. Positive Results: \mathcal{H}_{all} -Consistency Bounds

In this section, we present a series of positive results by proving \mathcal{H}_{all} -consistency bounds for $\tilde{\mathcal{L}}_{0-1}$ and the surrogate loss $\tilde{\mathcal{L}}$ when using as an auxiliary function Φ the hinge-loss, the ρ -margin loss, the exponential loss, the logistic loss, the squared hinge-loss and the sigmoid loss, as summarized in Table 9 of Appendix J, where the corresponding proofs are also provided. The expression of the corresponding minimizability gaps are summarized in Table 1. In bipartite ranking with $\mathcal{H} = \mathcal{H}_{\text{all}}$, the minimizability gaps are zero for $\Phi_{\rho}, \Phi_{\text{exp}}, \Phi_{\text{log}}, \Phi_{\text{sig}}$, while they are non-zero for Φ_{hinge} and Φ_{sq} in general.

5. Bipartite Ranking with Abstention

As with the general pairwise ranking case, the negative results shown in Section 4.2 motivate us to study *bipartite ranking with abstention*, where the learner can abstain from making prediction on a pair (x, x') with x an x' relatively

Table 4: \mathcal{H}_{lin} -consistency upper bounds for bipartite abstention loss.

Loss function	\mathcal{H}_{lin} -consistency upper bound
$\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}$	$\frac{\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}^*}(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{lin}})}{\min\{W\gamma, 1\}} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$
$\tilde{\mathcal{L}}_{\Phi_{\rho}}$	$\frac{\rho\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}\right)}{\min\{W\gamma, \rho\}} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$
$\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}$	$\Gamma_{\text{exp}}\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}^*}(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{lin}})\right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$ where $\Gamma_{\text{exp}}(t) = \max\left\{\sqrt{t}, \left(\frac{e^{2W\gamma+1}}{e^{2W\gamma}-1}\right)t\right\}$
$\tilde{\mathcal{L}}_{\Phi_{\log}}$	$\Gamma_{\log}\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}\right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$ where $\Gamma_{\log}(t) = \max\left\{\sqrt{t}, \left(\frac{e^{W\gamma+1}}{e^{W\gamma}-1}\right)t\right\}$
$\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}$	$\Gamma_{\text{sq}}\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}^*}(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{lin}})\right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$ where $\Gamma_{\text{sq}}(t) = \max\left\{\sqrt{t}, \frac{t}{2W\gamma} + \frac{W\gamma}{2}\right\}$
$\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}$	$\frac{\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}^*}(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{lin}})}{\tanh(kW\gamma)} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$

 Table 5: \mathcal{H}_{NN} -consistency upper bounds for bipartite abstention loss.

Loss function	\mathcal{H}_{NN} -consistency upper bound
$\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}$	$\frac{\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}^*}(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{NN}})}{\min\{\Lambda W\gamma, 1\}} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$
$\tilde{\mathcal{L}}_{\Phi_{\rho}}$	$\frac{\rho\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}\right)}{\min\{\Lambda W\gamma, \rho\}} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$
$\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}$	$\Gamma_{\text{exp}}\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}^*}(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{NN}})\right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$ where $\Gamma_{\text{exp}}(t) = \max\left\{\sqrt{t}, \left(\frac{e^{2\Lambda W\gamma+1}}{e^{2\Lambda W\gamma}-1}\right)t\right\}$
$\tilde{\mathcal{L}}_{\Phi_{\log}}$	$\Gamma_{\log}\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}\right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$ where $\Gamma_{\log}(t) = \max\left\{\sqrt{t}, \left(\frac{e^{\Lambda W\gamma+1}}{e^{\Lambda W\gamma}-1}\right)t\right\}$
$\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}$	$\Gamma_{\text{sq}}\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}^*}(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{NN}})\right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$ where $\Gamma_{\text{sq}}(t) = \max\left\{\sqrt{t}, \frac{t}{2\Lambda W\gamma} + \frac{\Lambda W\gamma}{2}\right\}$
$\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}$	$\frac{\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}^*}(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{NN}})}{\tanh(k\Lambda W\gamma)} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})$

close. The *bipartite abstention loss* is defined as follows for any $h \in \mathcal{H}$ and $(x, y), (x', y') \in \mathcal{X} \times \mathcal{Y}$:

$$\begin{aligned} \tilde{\mathcal{L}}_{0-1}^{\text{abs}}(h, x, x', y, y') \\ = \tilde{\mathcal{L}}_{0-1}(h, x, x', y, y') \mathbb{1}_{\|x-x'\| > \gamma} + c \mathbb{1}_{\|x-x'\| \leq \gamma}, \end{aligned} \quad (6)$$

where γ is a given threshold value. When $\gamma = 0$, $\tilde{\mathcal{L}}_{0-1}^{\text{abs}}$ reduces to bipartite misranking loss $\tilde{\mathcal{L}}_{0-1}$ without abstention.

5.1. Linear Hypotheses

Table 4 presents a series of \mathcal{H}_{lin} -consistency upper bounds for $\tilde{\mathcal{L}}_{\Phi}$ when using as Φ the auxiliary functions in Table 1. The bounds of Table 4 depend directly on the threshold value γ , the parameter W in the linear models and parameters of the loss function (e.g., k in sigmoid loss).

As an example, when adopting the exponential loss function as Φ , modulo the minimizability gaps (which are zero when the best-in-class error coincides with the Bayes error or can be small in some other cases), the bound implies that if the surrogate estimation loss $\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}^*}(\mathcal{H}_{\text{lin}})$ is reduced to ϵ , then, the target estimation loss $\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) -$

$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}})$ is upper bounded by $\Gamma_{\text{exp}}(\epsilon)$. For sufficiently small values of ϵ , the dependence of Γ_{exp} on ϵ exhibits a square root relationship. However, if this is not the case, the dependence becomes linear, subject to a constant factor depending on the threshold value γ and the parameter W in the linear models.

As with the general pairwise ranking setting, the proofs consist of analyzing calibration gaps of the target loss and that of each surrogate loss and seeking a tight lower bound of the surrogate calibration gap in terms of the target one. Additionally, the bipartite ranking setting introduces an added layer of complexity, as x and x' in a pair have independent conditional distributions $\eta(x)$ and $\eta(x')$, which results in a more intricate calibration gap that is harder to address.

As an example, for $\Phi = \Phi_{\text{exp}}$ the exponential loss function, we have the lower bound $\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \geq \Psi_{\text{exp}}\left(\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}}}(h, x, x')\right)$, where Ψ_{exp} is an increasing and piece-wise convex function on $[0, 2]$ defined by $\Psi_{\text{exp}}(t) = \min\left\{t^2, \left(\frac{e^{2W\gamma+1}}{e^{2W\gamma}-1}\right)t\right\}$. The detailed derivation and the expression of the corresponding minimizability gaps are included in Appendix L.1.

Pairwise Misranking Loss Surrogates

Table 6: General pairwise abstention loss for the Rankboost loss on CIFAR-10; mean \pm standard deviation over three runs for various γ and cost c .

γ	0	0.3	0.5	0.7	0.9
Cost 0.1	8.33% \pm 0.15%	8.33% \pm 0.15%	8.33% \pm 0.15%	8.25% \pm 0.07%	8.54% \pm 0.07%
Cost 0.3	8.33% \pm 0.15%	8.33% \pm 0.15%	8.35% \pm 0.15%	9.73% \pm 0.11%	20.41% \pm 0.06%
Cost 0.5	8.33% \pm 0.15%	8.33% \pm 0.15%	8.36% \pm 0.14%	11.20% \pm 0.14%	32.28% \pm 0.07%

5.2. One-hidden-layer ReLU Neural Networks.

Table 5 presents the \mathcal{H}_{NN} -consistency upper bounds for $\tilde{\mathcal{L}}_{\Phi}$ when using as Φ the auxiliary functions in Table 1. Different from the bounds in the linear case, all the bounds in Table 5 not only depend on W , but also depend on Λ , a parameter in \mathcal{H}_{NN} . The proof is similar to that of the linear case. The detailed derivation and the expression of the corresponding minimizability gaps are given in Appendix L.2.

As with the linear case, taking the exponential loss function Φ_{exp} as an example, modulo the minimizability gaps, the bound implies that if the surrogate estimation loss $\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{NN}})$ is reduced to ϵ , then, the target estimation loss $\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}})$ is upper bounded by $\Gamma_{\text{exp}}(\epsilon)$. For sufficiently small values of ϵ , the dependence of Γ_{exp} on ϵ exhibits a square root relationship. However, if this is not the case, the dependence becomes linear, subject to a constant factor depending on the threshold value γ , the parameter W , and an additional parameter Λ in the one-hidden-layer ReLU networks.

6. Experiments

In this section, we provide empirical results for general pairwise ranking with abstention on the CIFAR-10 dataset (Krizhevsky, 2009).

We used ResNet-34 with ReLU activations (He et al., 2016). Here, ResNet- n denotes a residual network with n convolutional layers. Standard data augmentations, 4-pixel padding with 32×32 random crops and random horizontal flips are applied for CIFAR-10. For training, we used Stochastic Gradient Descent (SGD) with Nesterov momentum (Nesterov, 1983). We set the batch size, weight decay, and initial learning rate to 1,024, 1×10^{-4} and 0.1 respectively. We adopted the cosine decay learning rate schedule (Loshchilov & Hutter, 2016) for a total of 200 epochs. The pairs (x, x', y) are randomly sampled from CIFAR-10 during training, with $y = \pm 1$ indicating if x is ranked above or below x' per the natural ordering of labels of x and x' .

We evaluated the models based on their averaged pairwise abstention loss (3) with γ selected from $\{0.0, 0.3, 0.5, 0.7, 0.9\}$ and the cost c selected from $\{0.1, 0.3, 0.5\}$. We randomly sampled 10,000 pairs (x, x') from the test data for evaluation. The ℓ_{∞} distance is adopted in the algorithm. We averaged losses over three runs and

report the standard deviation as well.

We used the surrogate loss (2) with $\Phi(t) = \exp(-t)$ the exponential loss, $\mathcal{L}_{\Phi_{\text{exp}}}$, which coincides with the loss function of RankBoost. Table 6 shows that when γ is as small as 0.3, no abstention takes place and the abstention loss coincides with the standard misranking loss ($\gamma = 0$) for any cost c . As γ increases, there are more samples that are abstained. When using a minimal cost c of 0.1 (as demonstrated in the first row of Table 6), abstaining on pairs with a relatively small distance ($\gamma = 0.7$) results in a lower target abstention loss compared to the scenario without abstention ($\gamma = 0$). Conversely, abstaining on pairs with larger distances ($\gamma = 0.9$) led to a higher abstention loss. This can be attributed to the fact that rejected samples at $\gamma = 0.7$ had lower accuracy compared to those at $\gamma = 0.9$. This empirically verifies that the surrogate loss $\mathcal{L}_{\Phi_{\text{exp}}}$ is not favorable on pairs whose distance is relatively small, for equicontinuous hypotheses. When the cost c is larger, the abstention loss, in general, increases with γ , since the number of samples rejected increases with γ .

Overall, the experiment shows that, in practice, for small γ , abstention actually does not take place. Thus, the abstention loss coincides with the standard pairwise misranking loss in those cases, and the surrogate loss is consistent with respect to both of them. Our results also indicate that the surrogate loss $\mathcal{L}_{\Phi_{\text{exp}}}$, a commonly used loss function, for example for RankBoost, is not optimal for pairs with a relatively small distance. Instead, rejecting these pairs at a minimal cost proves to be a more effective strategy.

7. Conclusion

We presented a series of theoretical \mathcal{H} -consistency guarantees for surrogate losses in pairwise misranking. Our proposed abstention methods are important when using common equicontinuous hypothesis sets in practice. It will be useful to explore alternative non-equicontinuous hypothesis sets that may be of practical use, and to further study the choice of the parameter γ for abstention in practice. We have also initiated the study of randomized ranking solutions with theoretical guarantees without resorting to abstention.

Acknowledgments

We thank reviewers for their comments on the presentation.

References

- Agarwal, S. Surrogate regret bounds for bipartite ranking via strongly proper losses. *The Journal of Machine Learning Research*, 15(1):1653–1674, 2014.
- Agarwal, S., Graepel, T., Herbrich, R., Har-Peled, S., Roth, D., and Jordan, M. I. Generalization bounds for the area under the ROC curve. *Journal of Machine Learning Research*, 6(4), 2005.
- Ailon, N. and Mohri, M. An efficient reduction of ranking to classification. In *Conference on Learning Theory*, 2008.
- Ailon, N. and Mohri, M. Preference-based learning to rank. *Machine Learning*, 80(2-3):189–211, 2010.
- Awasthi, P., Frank, N., Mao, A., Mohri, M., and Zhong, Y. Calibration and consistency of adversarial surrogate losses. In *Advances in Neural Information Processing Systems*, 2021a.
- Awasthi, P., Mao, A., Mohri, M., and Zhong, Y. A finer calibration analysis for adversarial robustness. *arXiv preprint arXiv:2105.01550*, 2021b.
- Awasthi, P., Mao, A., Mohri, M., and Zhong, Y. H-consistency bounds for surrogate loss minimizers. In *International Conference on Machine Learning*, 2022a.
- Awasthi, P., Mao, A., Mohri, M., and Zhong, Y. Multi-class \mathcal{H} -consistency bounds. In *Advances in neural information processing systems*, 2022b.
- Awasthi, P., Mao, A., Mohri, M., and Zhong, Y. DC-programming for neural network optimizations. *Journal of Global Optimization*, 2023a.
- Awasthi, P., Mao, A., Mohri, M., and Zhong, Y. Theoretically grounded loss functions and algorithms for adversarial robustness. In *International Conference on Artificial Intelligence and Statistics*, pp. 10077–10094, 2023b.
- Bartlett, P. L., Jordan, M. I., and McAuliffe, J. D. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- Berkson, J. Application of the logistic function to bio-assay. *Journal of the American Statistical Association*, 39:357–365, 1944.
- Berkson, J. Why I prefer logits to probits. *Biometrics*, 7(4): 327–339, 1951.
- Buffoni, D., Calauzenes, C., Gallinari, P., and Usunier, N. Learning scoring functions with order-preserving losses and standardized supervision. In *International Conference on Machine Learning*, pp. 825–832, 2011.
- Calauzenes, C., Usunier, N., and Gallinari, P. On the (non-) existence of convex, calibrated surrogate losses for ranking. In *Advances in Neural Information Processing Systems*, 2012.
- Carlini, N. and Wagner, D. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy (SP)*, pp. 39–57, 2017.
- Clemençon, S., Lugosi, G., and Vayatis, N. Ranking and empirical minimization of U-statistics. *The Annals of Statistics*, 36(2):844–874, 2008.
- Cohen, W. W., Schapire, R. E., and Singer, Y. Learning to order things. *Advances in neural information processing systems*, 10, 1997.
- Cortes, C. and Mohri, M. AUC optimization vs. error rate minimization. *Advances in neural information processing systems*, 16, 2003.
- Cossock, D. and Zhang, T. Statistical analysis of bayes optimal subset ranking. *IEEE Transactions on Information Theory*, 54(11):5140–5154, 2008.
- Crammer, K. and Singer, Y. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of machine learning research*, 2(Dec):265–292, 2001.
- Duchi, J. C., Mackey, L. W., and Jordan, M. I. On the consistency of ranking algorithms. In *International conference on Machine learning*, pp. 327–334, 2010.
- Freund, Y., Iyer, R., Schapire, R. E., and Singer, Y. An efficient boosting algorithm for combining preferences. *Journal of machine learning research*, 4(Nov):933–969, 2003.
- Gao, W. and Zhou, Z.-H. On the consistency of multi-label learning. In *Conference on learning theory*, pp. 341–358, 2011.
- Gao, W. and Zhou, Z.-H. On the consistency of AUC pairwise optimization. In *International Joint Conference on Artificial Intelligence*, 2015.
- Gao, W., Jin, R., Zhu, S., and Zhou, Z.-H. One-pass auc optimization. In *International conference on machine learning*, pp. 906–914, 2013.
- Ghosh, A., Kumar, H., and Sastry, P. S. Robust loss functions under label noise for deep neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, 2017.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.

- Hanley, J. A. and McNeil, B. J. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 1982.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Joachims, T. Optimizing search engines using clickthrough data. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 133–142, 2002.
- Kotlowski, W., Dembczynski, K. J., and Huellermeier, E. Bipartite ranking through minimization of univariate loss. In *International Conference on Machine Learning*, pp. 1113–1120, 2011.
- Krizhevsky, A. Learning multiple layers of features from tiny images. Technical report, Toronto University, 2009.
- Kuznetsov, V., Mohri, M., and Syed, U. Multi-class deep boosting. In *Advances in Neural Information Processing Systems*, pp. 2501–2509, 2014.
- Lan, Y., Guo, J., Cheng, X., and Liu, T.-Y. Statistical consistency of ranking methods in a rank-differentiable probability space. In *Advances in Neural Information Processing Systems*, 2012.
- Lee, Y., Lin, Y., and Wahba, G. Multicategory support vector machines: Theory and application to the classification of microarray data and satellite radiance data. *Journal of the American Statistical Association*, 99(465):67–81, 2004.
- Long, P. and Servedio, R. Consistency versus realizable H-consistency for multiclass classification. In *International Conference on Machine Learning*, pp. 801–809, 2013.
- Loshchilov, I. and Hutter, F. SGDR: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Mao, A., Mohri, M., and Zhong, Y. Cross-entropy loss functions: Theoretical analysis and applications. *arXiv preprint arXiv:2304.07288*, 2023.
- Menon, A. K. and Williamson, R. C. Bayes-optimal scorers for bipartite ranking. In *Conference on Learning Theory*, pp. 68–106, 2014.
- Mohri, M., Rostamizadeh, A., and Talwalkar, A. *Foundations of Machine Learning*. MIT Press, second edition, 2018.
- Nesterov, Y. E. A method for solving the convex programming problem with convergence rate $o(1/k^2)$. *Dokl. akad. nauk Sssr*, 269:543–547, 1983.
- Ramaswamy, H. G. and Agarwal, S. Classification calibration dimension for general multiclass losses. In *Advances in Neural Information Processing Systems*, 2012.
- Ramaswamy, H. G., Agarwal, S., and Tewari, A. Convex calibrated surrogates for low-rank loss matrices with applications to subset ranking losses. In *Advances in Neural Information Processing Systems*, 2013.
- Ramaswamy, H. G., Babu, B. S., Agarwal, S., and Williamson, R. C. On the consistency of output code based learning algorithms for multiclass learning problems. In *Conference on Learning Theory*, pp. 885–902, 2014.
- Ravikumar, P., Tewari, A., and Yang, E. On ndcg consistency of listwise ranking methods. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pp. 618–626, 2011.
- Rudin, C., Cortes, C., Mohri, M., and Schapire, R. E. Margin-based ranking meets boosting in the middle. In *Conference on Learning Theory*, pp. 63–78, 2005.
- Steinwart, I. How to compare different loss functions and their risks. *Constructive Approximation*, 26(2):225–287, 2007.
- Tewari, A. and Bartlett, P. L. On the consistency of multi-class classification methods. *Journal of Machine Learning Research*, 8(36):1007–1025, 2007.
- Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., and Madry, A. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- Uematsu, K. and Lee, Y. On theoretically optimal ranking functions in bipartite ranking. *Journal of the American Statistical Association*, 112(519):1311–1322, 2017.
- Verhulst, P. F. Notice sur la loi que la population suit dans son accroissement. *Correspondance mathématique et physique*, 10:113—121, 1838.
- Verhulst, P. F. Recherches mathématiques sur la loi d’accroissement de la population. *Nouveaux Mémoires de l’Académie Royale des Sciences et Belles-Lettres de Bruxelles*, 18:1—42, 1845.
- Weston, J. and Watkins, C. Multi-class support vector machines. Technical report, Citeseer, 1998.

- Xia, F., Liu, T.-Y., Wang, J., Zhang, W., and Li, H. Listwise approach to learning to rank: theory and algorithm. In *International conference on Machine learning*, pp. 1192–1199, 2008.
- Zhang, M. and Agarwal, S. Bayes consistency vs. H-consistency: The interplay between surrogate loss functions and the scoring function class. In *Advances in Neural Information Processing Systems*, 2020.
- Zhang, M., Ramaswamy, H. G., and Agarwal, S. Convex calibrated surrogates for the multi-label f-measure. In *International Conference on Machine Learning*, pp. 11246–11255, 2020.
- Zhang, T. Statistical behavior and consistency of classification methods based on convex risk minimization. *The Annals of Statistics*, 32(1):56–85, 2004.
- Zhang, Z. and Sabuncu, M. Generalized cross entropy loss for training deep neural networks with noisy labels. In *Advances in neural information processing systems*, 2018.
- Zheng, C., Wu, G., Bao, F., Cao, Y., Li, C., and Zhu, J. Revisiting discriminative vs. generative classifiers: Theory and implications. *arXiv preprint arXiv:2302.02334*, 2023.

Contents of Appendix

A	Related work	15
B	Notation	15
C	General tools	15
D	Negative results for general pairwise ranking (Proof of Theorem 2.1)	18
E	Positive results for piecewise functions in general pairwise ranking (Proof of Theorem 2.3)	18
F	\mathcal{H}_{all} - consistency bounds for pairwise misranking losses	19
F.1	Derivation for $L_{\Phi_{\text{hinge}}}$	20
F.2	Derivation for $L_{\Phi_{\rho}}$	20
F.3	Derivation for $L_{\Phi_{\text{exp}}}$	21
F.4	Derivation for $L_{\Phi_{\text{log}}}$	21
F.5	Derivation for $L_{\Phi_{\text{sq}}}$	22
F.6	Derivation for $L_{\Phi_{\text{sig}}}$	23
G	Minimizability gaps can be non-zero for $\mathcal{H} = \mathcal{H}_{\text{all}}$ in the general pairwise misranking case.	23
H	Characterization of distribution order and minimizability gap (Proof of Theorem 2.5, Theorem 2.6 and Theorem 2.7)	24
I	Negative results for bipartite ranking (Proof of Theorem 4.1)	26
J	\mathcal{H}_{all} - consistency bounds for bipartite misranking losses	27
J.1	Derivation for $\tilde{L}_{\Phi_{\text{hinge}}}$	28
J.2	Derivation for $\tilde{L}_{\Phi_{\rho}}$	29
J.3	Derivation for $\tilde{L}_{\Phi_{\text{exp}}}$	29
J.4	Derivation for $\tilde{L}_{\Phi_{\text{log}}}$	30
J.5	Derivation for $\tilde{L}_{\Phi_{\text{sq}}}$	31
J.6	Derivation for $\tilde{L}_{\Phi_{\text{sig}}}$	32
K	\mathcal{H} - consistency bounds for pairwise abstention loss	33
K.1	Linear Hypotheses	34
K.1.1	Derivation for $L_{\Phi_{\text{hinge}}}$	34
K.1.2	Derivation for $L_{\Phi_{\rho}}$	35
K.1.3	Derivation for $L_{\Phi_{\text{exp}}}$	35
K.1.4	Derivation for $L_{\Phi_{\text{log}}}$	37

K.1.5	Derivation for $L_{\Phi_{sq}}$	39
K.1.6	Derivation for $L_{\Phi_{sig}}$	40
K.2	One-Hidden-Layer ReLU Neural Networks	40
K.2.1	Derivation for $L_{\Phi_{hinge}}$	41
K.2.2	Derivation for $L_{\Phi_{\rho}}$	41
K.2.3	Derivation for $L_{\Phi_{exp}}$	42
K.2.4	Derivation for $L_{\Phi_{log}}$	44
K.2.5	Derivation for $L_{\Phi_{sq}}$	45
K.2.6	Derivation for $L_{\Phi_{sig}}$	46
L	\mathcal{H} - consistency bounds for bipartite abstention losses	47
L.1	Linear Hypotheses	48
L.1.1	Derivation for $\tilde{L}_{\Phi_{hinge}}$	48
L.1.2	Derivation for $\tilde{L}_{\Phi_{\rho}}$	49
L.1.3	Derivation for $\tilde{L}_{\Phi_{exp}}$	50
L.1.4	Derivation for $\tilde{L}_{\Phi_{log}}$	51
L.1.5	Derivation for $\tilde{L}_{\Phi_{sq}}$	53
L.1.6	Derivation for $\tilde{L}_{\Phi_{sig}}$	54
L.2	One-Hidden-Layer ReLU Neural Networks	54
L.2.1	Derivation for $\tilde{L}_{\Phi_{hinge}}$	54
L.2.2	Derivation for $\tilde{L}_{\Phi_{\rho}}$	55
L.2.3	Derivation for $\tilde{L}_{\Phi_{exp}}$	56
L.2.4	Derivation for $\tilde{L}_{\Phi_{log}}$	57
L.2.5	Derivation for $\tilde{L}_{\Phi_{sq}}$	59
L.2.6	Derivation for $\tilde{L}_{\Phi_{sig}}$	60

A. Related work

The notions of Bayes consistency (also known as consistency) and calibration have been extensively studied for classification (Zhang, 2004; Bartlett et al., 2006; Tewari & Bartlett, 2007). The Bayes consistency of ranking surrogate losses has been studied in the special case of bipartite score-based ranking: in particular, Uematsu & Lee (2017) proved the inconsistency of the pairwise ranking loss based on the hinge loss and Gao & Zhou (2015) gave excess loss bounds for pairwise ranking losses based on the exponential or the logistic loss. Later, these results were further generalized by Menon & Williamson (2014). A related but distinct consistency question has been studied in several publications (Agarwal et al., 2005; Kotlowski et al., 2011; Agarwal, 2014). It is one with respect to binary classification, that is whether a near minimizer of the surrogate loss of the binary classification loss is a near minimizer of the bipartite misranking loss (Cortes & Mohri, 2003).

Considerable attention has been devoted to the study of the learning to rank algorithms and their related problems: including one-pass AUC pairwise optimization (Gao et al., 2013), preference-based ranking (Cohen et al., 1997; Clemençon et al., 2008), subset ranking with Discounted Cumulative Gain (DCG) (Cosssock & Zhang, 2008; Buffoni et al., 2011), listwise ranking (Xia et al., 2008), subset ranking based on Pairwise Disagreement (PD) (Duchi et al., 2010; Lan et al., 2012), subset ranking using Normalized Discounted Cumulative Gain (NDCG) (Ravikumar et al., 2011), subset ranking with Average Precision (AP) (Calauzenes et al., 2012; Ramaswamy et al., 2013), general multi-class problems (Ramaswamy & Agarwal, 2012; Ramaswamy et al., 2014) and multi-label problems (Gao & Zhou, 2011; Zhang et al., 2020).

Bayes consistency only holds for the full family of measurable functions, which of course is distinct from the more restricted hypothesis set used by a learning algorithm. Therefore, a hypothesis set-dependent notion of \mathcal{H} -consistency has been proposed by Long & Servedio (2013) in the realizable setting, which was used by Zhang & Agarwal (2020) for linear models, and generalized by Kuznetsov et al. (2014) to the structured prediction case. Long & Servedio (2013) showed that there exists a case where a Bayes-consistent loss is not \mathcal{H} -consistent while inconsistent loss functions can be \mathcal{H} -consistent. Zhang & Agarwal (2020) further investigated the phenomenon in (Long & Servedio, 2013) and showed that the situation of loss functions that are not \mathcal{H} -consistent with linear models can be remedied by carefully choosing a larger piecewise linear hypothesis set. Kuznetsov et al. (2014) proved positive results for the \mathcal{H} -consistency of several multi-class ensemble algorithms, as an extension of \mathcal{H} -consistency results in (Long & Servedio, 2013).

Recently, Awasthi et al. (2022a) presented a series of results providing \mathcal{H} -consistency bounds in binary classification. These guarantees are significantly stronger than the \mathcal{H} -calibration or \mathcal{H} -consistency properties studied by Awasthi et al. (2021a;b). Awasthi et al. (2022b) and Mao et al. (2023) (see also (Zheng et al., 2023)) generalized \mathcal{H} -consistency bounds to the scenario of multi-class classification. Awasthi et al. (2023b) proposed a family of loss functions that benefit from such \mathcal{H} -consistency bounds guarantees for adversarial robustness (Goodfellow et al., 2014; Madry et al., 2017; Tsipras et al., 2018; Carlini & Wagner, 2017; Awasthi et al., 2023a). \mathcal{H} -consistency bounds are also more informative than similar excess error bounds derived in the literature, which correspond to the special case where \mathcal{H} is the family of all measurable functions (Zhang, 2004; Bartlett et al., 2006; Mohri et al., 2018). Our work significantly generalizes the results of Awasthi et al. (2022a) to the score-based ranking setting, including both the general pairwise ranking and bipartite ranking scenarios.

B. Notation

We provide a table of notation in Table 7.

C. General tools

To begin with the proof, we first introduce some notation. In general pairwise ranking scenario, we denote by \mathcal{D} a distribution over $\mathcal{X} \times \mathcal{X} \times \mathcal{Y}$ and by \mathcal{P} a set of such distributions. We further denote by $\eta(x, x') = \mathcal{D}(Y = 1 | (X, X') = (x, x'))$ the conditional probability of $Y = 1$ given $(X, X') = (x, x')$. Without loss of generality, we assume that $\eta(x, x) = 1/2$. The generalization error for a surrogate loss L can be rewritten as $\mathcal{R}_L(h) = \mathbb{E}_X[\mathcal{C}_L(h, x, x')]$, where $\mathcal{C}_L(h, x, x')$ is the conditional L-risk, defined by

$$\mathcal{C}_L(h, x, x') = \eta(x, x')L(h, x, x', +1) + (1 - \eta(x, x'))L(h, x, x', -1).$$

We denote by $\mathcal{C}_L^*(\mathcal{H}, x, x') = \inf_{h \in \mathcal{H}} \mathcal{C}_L(h, x, x')$ the minimal conditional L-risk. Then, the minimizability gap can be rewritten as follows:

$$\mathcal{M}_L(\mathcal{H}) = \mathcal{R}_L^*(\mathcal{H}) - \mathbb{E}_X[\mathcal{C}_L^*(\mathcal{H}, x)].$$

Table 7: Summary of notation.

\mathcal{X}	Input space
\mathcal{Y}	Label space
\mathcal{H}	A hypothesis set of functions mapping from \mathcal{X} to \mathbb{R}
\mathcal{D}	A distribution over $\mathcal{X} \times \mathcal{X} \times \mathcal{Y}$ or $\mathcal{X} \times \mathcal{Y}$
L_{0-1}	General pairwise misranking loss
$\mathcal{R}_{L_{0-1}}$	Expected general pairwise misranking loss
$\text{sign}(u)$	$\mathbb{1}_{u \geq 0} - \mathbb{1}_{u < 0}$
$\eta(x, x')$	The conditional probability of $Y = +1$ given $(X, X') = (x, x')$
\tilde{L}_{0-1}	Bipartite misranking loss
$\mathcal{R}_{\tilde{L}_{0-1}}$	Expected bipartite misranking loss
$\eta(x)$	The conditional probability of $Y = +1$ given $X = x$
L	A surrogate loss for L_{0-1}
\tilde{L}	A surrogate loss for \tilde{L}_{0-1}
$\mathcal{R}_L^*(\mathcal{H})$ ($\mathcal{R}_{\tilde{L}}^*(\mathcal{H})$)	The minimal generalization error
$\mathcal{M}_L(\mathcal{H})$ ($\mathcal{M}_{\tilde{L}}(\mathcal{H})$)	The minimizability gap
\mathcal{H}_{all}	The hypothesis set of all measurable functions
\mathcal{H}_{lin}	Linear hypothesis set
\mathcal{H}_{NN}	The hypothesis set of one-hidden-layer ReLU networks
$(\cdot)_+$	$\max(\cdot, 0)$
\mathcal{H}_{pw}	The hypothesis set of piecewise functions
\mathcal{D}	
\leq	The distribution order
L_{0-1}^{abs}	The pairwise abstention loss
γ	A given threshold value
c	Cost
$\tilde{L}_{0-1}^{\text{abs}}$	The bipartite abstention loss
\mathcal{C}_L ($\mathcal{C}_{\tilde{L}}$)	The conditional L-risk (\tilde{L} -risk)
$\mathcal{C}_L^*(\mathcal{H}, x, x')$ ($\mathcal{C}_{\tilde{L}}^*(\mathcal{H}, x, x')$)	The minimal conditional L-risk (\tilde{L} -risk)
$\Delta \mathcal{C}_{L, \mathcal{H}}$ ($\Delta \mathcal{C}_{\tilde{L}, \mathcal{H}}$)	The calibration gap
$\langle t \rangle_\epsilon$	The ϵ -truncation of t

We further refer to $\mathcal{C}_L(h, x, x') - \mathcal{C}_L^*(\mathcal{H}, x, x')$ as the calibration gap and denote it by $\Delta \mathcal{C}_{L, \mathcal{H}}(h, x, x')$.

In bipartite ranking scenario, we denote by \mathcal{D} a distribution over $\mathcal{X} \times \mathcal{Y}$ and by \mathcal{P} a set of such distributions. We further denote by $\eta(x) = \mathcal{D}(Y = 1 | X = x)$ the conditional probability of $Y = 1$ given $X = x$. The generalization error for a surrogate loss L can be rewritten as $\mathcal{R}_L(h) = \mathbb{E}_X[\mathcal{C}_L(h, x, x')]$, where $\mathcal{C}_L(h, x, x')$ is the conditional L -risk, defined by

$$\mathcal{C}_L(h, x, x') = \eta(x)(1 - \eta(x'))\tilde{L}(h, x, x', +1, -1) + \eta(x')(1 - \eta(x))\tilde{L}(h, x, x', -1, +1).$$

We denote by $\mathcal{C}_L^*(\mathcal{H}, x, x') = \inf_{h \in \mathcal{H}} \mathcal{C}_L(h, x, x')$ the minimal conditional L -risk. Then, the minimizability gap can be rewritten as follows:

$$\mathcal{M}_L(\mathcal{H}) = \mathcal{R}_L^*(\mathcal{H}) - \mathbb{E}_X[\mathcal{C}_L^*(\mathcal{H}, x, x')].$$

We further refer to $\mathcal{C}_L(h, x, x') - \mathcal{C}_L^*(\mathcal{H}, x, x')$ as the calibration gap and denote it by $\Delta \mathcal{C}_{L, \mathcal{H}}(h, x, x')$. For any $\epsilon > 0$, we will denote by $\langle t \rangle_\epsilon$ the ϵ -truncation of $t \in \mathbb{R}$ defined by $t\mathbb{1}_{t > \epsilon}$.

We first prove two general results, which provide bounds between any loss functions L_1 and L_2 in both general pairwise ranking scenario and bipartite ranking scenario.

Theorem C.1. *Assume that there exists a convex function $\Psi: \mathbb{R}_+ \rightarrow \mathbb{R}$ with $\Psi(0) \geq 0$ and $\epsilon \geq 0$ such that the following holds for all $h \in \mathcal{H}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $\mathcal{D} \in \mathcal{P}$:*

$$\Psi(\langle \Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \rangle_\epsilon) \leq \langle \Delta \mathcal{C}_{L_1, \mathcal{H}}(h, x, x') \rangle_\epsilon. \quad (7)$$

Then, the following inequality holds for any $h \in \mathcal{H}$ and $\mathcal{D} \in \mathcal{P}$:

$$\Psi(\mathcal{R}_{L_2}(h) - \mathcal{R}_{L_2}^*(\mathcal{H}) + \mathcal{M}_{L_2}(\mathcal{H})) \leq \mathcal{R}_{L_1}(h) - \mathcal{R}_{L_1}^*(\mathcal{H}) + \mathcal{M}_{L_1}(\mathcal{H}) + \max\{\Psi(0), \Psi(\epsilon)\}. \quad (8)$$

Proof. By the definition of the generalization error and the minimizability gap, for any $h \in \mathcal{H}$ and $\mathcal{D} \in \mathcal{P}$, we can write the left hand side of (8) as

$$\Psi(\mathcal{R}_{L_2}(h) - \mathcal{R}_{L_2}^*(\mathcal{H}) + \mathcal{M}_{L_2}(\mathcal{H})) = \Psi(\mathcal{R}_{L_2}(h) - \mathbb{E}_{(X, X')}[\mathcal{C}_{L_2}^*(\mathcal{H}, x, x')]) = \Psi(\mathbb{E}_{(X, X')}[\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x')]).$$

Since Ψ is convex, by Jensen's inequality, it can be upper bounded by $\mathbb{E}_{(X, X')}[\Psi(\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x'))]$. Due to the decomposition

$$\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') = \langle \Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \rangle_\epsilon + \Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \mathbb{1}_{\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \leq \epsilon},$$

and the assumption $\Psi(0) \geq 0$, we have the following inequality:

$$\mathbb{E}_{(X, X')}[\Psi(\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x'))] \leq \mathbb{E}_{(X, X')}[\Psi(\langle \Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \rangle_\epsilon)] + \mathbb{E}_{(X, X')}[\Psi(\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \mathbb{1}_{\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \leq \epsilon})].$$

By assumption (7), the first term can be bounded as follows:

$$\mathbb{E}_{(X, X')}[\Psi(\langle \Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \rangle_\epsilon)] \leq \mathbb{E}_{(X, X')}[\Delta \mathcal{C}_{L_1, \mathcal{H}}(h, x, x')] = \mathcal{R}_{L_1}(h) - \mathcal{R}_{L_1}^*(\mathcal{H}) + \mathcal{M}_{L_1}(\mathcal{H}).$$

Since $\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \mathbb{1}_{\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \leq \epsilon} \in [0, \epsilon]$, we can bound $\mathbb{E}_{(X, X')}[\Psi(\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \mathbb{1}_{\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \leq \epsilon})]$ by $\sup_{t \in [0, \epsilon]} \Psi(t)$, which equals $\max\{\Psi(0), \Psi(\epsilon)\}$ due to the convexity of Ψ . \square

Theorem C.2. Assume that there exists a non-decreasing concave function $\Gamma: \mathbb{R}_+ \rightarrow \mathbb{R}$ and $\epsilon \geq 0$ such that the following holds for all $h \in \mathcal{H}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $\mathcal{D} \in \mathcal{P}$:

$$\langle \Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \rangle_\epsilon \leq \Gamma(\langle \Delta \mathcal{C}_{L_1, \mathcal{H}}(h, x, x') \rangle_\epsilon). \quad (9)$$

Then, the following inequality holds for any $h \in \mathcal{H}$ and $\mathcal{D} \in \mathcal{P}$:

$$\mathcal{R}_{L_2}(h) - \mathcal{R}_{L_2}^*(\mathcal{H}) \leq \Gamma(\mathcal{R}_{L_1}(h) - \mathcal{R}_{L_1}^*(\mathcal{H}) + \mathcal{M}_{L_1}(\mathcal{H})) - \mathcal{M}_{L_2}(\mathcal{H}) + \epsilon. \quad (10)$$

Proof. By the definition of the generalization error and the minimizability gap, for any $h \in \mathcal{H}$ and $\mathcal{D} \in \mathcal{P}$, we can write the left hand side of (10) as

$$\begin{aligned} & \mathcal{R}_{L_2}(h) - \mathcal{R}_{L_2}^*(\mathcal{H}) \\ &= \mathbb{E}_{(X, X')}[\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x')] - \mathcal{M}_{L_2}(\mathcal{H}) \\ &= \mathbb{E}_{(X, X')}[\langle \Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \rangle_\epsilon] + \mathbb{E}_{(X, X')}[\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \mathbb{1}_{\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \leq \epsilon}] - \mathcal{M}_{L_2}(\mathcal{H}) \end{aligned}$$

By assumption (9) and that Γ is non-decreasing, the following inequality holds:

$$\mathbb{E}_{(X, X')}[\langle \Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \rangle_\epsilon] \leq \mathbb{E}_{(X, X')}[\Gamma(\Delta \mathcal{C}_{L_1, \mathcal{H}}(h, x, x'))].$$

Since Γ is concave, by Jensen's inequality,

$$\mathbb{E}_{(X, X')}[\Gamma(\Delta \mathcal{C}_{L_1, \mathcal{H}}(h, x, x'))] \leq \Gamma(\mathbb{E}_{(X, X')}[\Delta \mathcal{C}_{L_1, \mathcal{H}}(h, x, x')]) = \Gamma(\mathcal{R}_{L_1}(h) - \mathcal{R}_{L_1}^*(\mathcal{H}) + \mathcal{M}_{L_1}(\mathcal{H})).$$

We complete the proof by noting that $\mathbb{E}_{(X, X')}[\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \mathbb{1}_{\Delta \mathcal{C}_{L_2, \mathcal{H}}(h, x, x') \leq \epsilon}] \leq \epsilon$. \square

D. Negative results for general pairwise ranking (Proof of Theorem 2.1)

Theorem 2.1 (Negative results). *Assume that \mathcal{X} contains an interior point x_0 and that \mathcal{H} is regular for general pairwise ranking, contains 0 and is equicontinuous at x_0 . If for some function f that is non-decreasing and continuous at 0, the following bound holds for all $h \in \mathcal{H}$ and any distribution,*

$$\mathcal{R}_{L_{0-1}}(h) - \mathcal{R}_{L_{0-1}}^*(\mathcal{H}) \leq f(\mathcal{R}_{L_\Phi}(h) - \mathcal{R}_{L_\Phi}^*(\mathcal{H})),$$

then, $f(t) \geq 1$ for any $t \geq 0$.

Proof. Assume $x_0 \in \mathcal{X}$ is an interior point and $h_0 = 0 \in \mathcal{H}$. By the assumption that x_0 is an interior point and \mathcal{H} is equicontinuous at x_0 , for any $\epsilon > 0$, we are able to take $x' \neq x_0 \in \mathcal{X}$ such that $|h(x') - h(x_0)| < \epsilon$ for all $h \in \mathcal{H}$. Consider the distribution that supports on $\{(x_0, x')\}$ with $\eta(x_0, x') = 0$. Then, for any $h \in \mathcal{H}$,

$$\mathcal{R}_{L_{0-1}}(h) = \mathcal{C}_{L_{0-1}}(h, x_0, x') = \mathbb{1}_{h(x') \geq h(x_0)} \geq 0,$$

where the equality can be achieved for some $h \in \mathcal{H}$ since \mathcal{H} is regular for general pairwise ranking. Therefore,

$$\mathcal{R}_{L_{0-1}}^*(\mathcal{H}) = \mathcal{C}_{L_{0-1}}^*(\mathcal{H}, x_0, x') = \inf_{h \in \mathcal{H}} \mathcal{C}_{L_{0-1}}(h, x_0, x') = 0.$$

Note $\mathcal{R}_{L_{0-1}}(h_0) = 1$. For the surrogate loss L_Φ , for any $h \in \mathcal{H}$,

$$\mathcal{R}_{L_\Phi}(h) = \mathcal{C}_{L_\Phi}(h, x_0, x') = \Phi(h(x_0) - h(x')) \in [\Phi(\epsilon), \Phi(-\epsilon)]$$

since $|h(x') - h(x_0)| < \epsilon$ and Φ is non-increasing. Therefore,

$$\mathcal{R}_{L_\Phi}^*(\mathcal{H}) = \mathcal{C}_{L_\Phi}^*(\mathcal{H}, x_0, x') \geq \Phi(\epsilon).$$

Note $\mathcal{R}_{L_\Phi}(h_0) = \Phi(0)$. If for some function f that is non-decreasing and continuous at 0, the bound holds, then, we obtain for any $h \in \mathcal{H}$ and $\epsilon > 0$,

$$\mathcal{R}_{L_{0-1}}(h) - 0 \leq f(\mathcal{R}_{L_\Phi}(h) - \mathcal{R}_{L_\Phi}^*(\mathcal{H})) \leq f(\mathcal{R}_{L_\Phi}(h) - \Phi(\epsilon)).$$

Let $h = h_0$, then $f(\Phi(0) - \Phi(\epsilon)) \geq 1$ for any $\epsilon > 0$. Take $\epsilon \rightarrow 0$, we obtain $f(0) \geq 1$ using the fact that Φ and f are both continuous at 0. Since f is non-decreasing, for any $t \in [0, 1]$, $f(t) \geq 1$. \square

E. Positive results for piecewise functions in general pairwise ranking (Proof of Theorem 2.3)

Theorem 2.3 (Positive results for piecewise functions). *Assume that Φ satisfies $\lim_{u \rightarrow +\infty} \Phi(u) = 0$. Then, for all $h \in \mathcal{H}_{\text{pw}}$ and any deterministic distribution,*

$$\mathcal{R}_{L_{0-1}}(h) - \mathcal{R}_{L_{0-1}}^*(\mathcal{H}_{\text{pw}}) \leq \mathcal{R}_{L_\Phi}(h) - \mathcal{R}_{L_\Phi}^*(\mathcal{H}_{\text{pw}}) + \mathcal{M}_{L_\Phi}(\mathcal{H}_{\text{pw}}) - \mathcal{M}_{L_{0-1}}(\mathcal{H}_{\text{pw}}).$$

Proof. For any $x \neq x' \in \mathcal{X}$,

$$\begin{aligned} \mathcal{C}_{L_{0-1}}(h, x, x') &= \begin{cases} \mathbb{1}_{h(x') < h(x)} & \eta(x, x') = 1 \\ \mathbb{1}_{h(x') \geq h(x)} & \eta(x, x') = 0 \end{cases} \\ &\geq 0 \\ \mathcal{C}_{L_\Phi}(h, x, x') &= \begin{cases} \Phi(h(x') - h(x)) & \eta(x, x') = 1 \\ \Phi(h(x) - h(x')) & \eta(x, x') = 0 \end{cases} \\ &\geq 0, \end{aligned}$$

where both equities can be achieved by $h(u) = \pm\alpha(\mathbb{1}_{u \neq x \wedge \|u\| > c} - \mathbb{1}_{u \neq x \wedge \|u\| \leq c})$ or $h(u) = \pm\alpha(\mathbb{1}_{u \neq x' \wedge \|u\| > c} - \mathbb{1}_{u \neq x' \wedge \|u\| \leq c})$ for $\alpha \rightarrow +\infty$. Therefore, we have $\mathcal{C}_{L_{0-1}}^*(\mathcal{H}, x, x') = \mathcal{C}_{L_\Phi}^*(\mathcal{H}, x, x') = 0$ and

$$\Delta \mathcal{C}_{L, \mathcal{H}}(h, x, x') = \mathcal{C}_{L_\Phi}(h, x, x') \geq \mathcal{C}_{L_{0-1}}(h, x, x') = \Delta \mathcal{C}_{L_{0-1}, \mathcal{H}}(h, x, x').$$

By Theorem C.1 or Theorem C.2, we complete the proof. \square

Table 8: \mathcal{H}_{all} -consistency upper bounds for general pairwise abstention.

Loss function	\mathcal{H}_{all} -consistency upper bound
$\mathcal{L}_{\Phi_{\text{hinge}}}$	$\mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{all}}) - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}})$
$\mathcal{L}_{\Phi_{\rho}}$	$\mathcal{R}_{\mathcal{L}_{\Phi_{\rho}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\rho}}}(\mathcal{H}_{\text{all}}) - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}})$
$\mathcal{L}_{\Phi_{\text{exp}}}$	$\sqrt{2} \left(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}} - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}})$
$\mathcal{L}_{\Phi_{\log}}$	$\sqrt{2} \left(\mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\log}}}(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}} - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}})$
$\mathcal{L}_{\Phi_{\text{sq}}}$	$\left(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}} - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}})$
$\mathcal{L}_{\Phi_{\text{sig}}}$	$\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{all}}) - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}})$

F. \mathcal{H}_{all} - consistency bounds for pairwise misranking losses

We first characterize the minimal conditional \mathcal{L}_{0-1} -risk and the calibration gap of the pairwise misranking loss for a broad class of hypothesis sets. We let $\overline{\mathcal{H}}(x, x') = \{h \in \mathcal{H}: \text{sign}(h(x') - h(x))(2\eta(x, x') - 1) \leq 0\}$ for convenience.

Lemma F.1. *Assume that \mathcal{H} is regular for general pairwise ranking. Then, the minimal conditional \mathcal{L}_{0-1} -risk is*

$$\mathcal{C}_{\mathcal{L}_{0-1}}^*(\mathcal{H}, x, x') = \min\{\eta(x, x'), 1 - \eta(x, x')\}.$$

The calibration gap of \mathcal{L}_{0-1} can be characterized as

$$\Delta \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}}(h, x, x') = |2\eta(x, x') - 1| \mathbb{1}_{h \in \overline{\mathcal{H}}(x, x')}.$$

Proof. By the definition, the conditional \mathcal{L}_{0-1} -risk is

$$\mathcal{C}_{\mathcal{L}_{0-1}}(h, x, x') = \eta(x, x') \mathbb{1}_{h(x') < h(x)} + (1 - \eta(x, x')) \mathbb{1}_{h(x') \geq h(x)}.$$

For any $x \in \mathcal{X}$, $\mathcal{C}_{\mathcal{L}_{0-1}}(h, x, x) = \mathcal{C}_{\mathcal{L}_{0-1}}^*(\mathcal{H}, x, x) = 1 - \eta(x, x) = 1/2$. For any $x \neq x' \in \mathcal{X}$, by the assumption, there exists $h^* \in \mathcal{H}$ such that $\text{sign}(h^*(x') - h^*(x)) = \text{sign}(\Delta\eta(x, x'))$. Therefore, the optimal conditional \mathcal{L}_{0-1} -risk can be characterized as for any $x, x' \in \mathcal{X}$,

$$\mathcal{C}_{\mathcal{L}_{0-1}}^*(\mathcal{H}, x, x') = \mathcal{C}_{\mathcal{L}_{0-1}}(h^*, x, x') = \min\{\eta(x, x'), 1 - \eta(x, x')\}$$

which proves the first part of lemma. By the definition,

$$\begin{aligned} \Delta \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}}(h, x, x') &= \mathcal{C}_{\mathcal{L}_{0-1}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{0-1}}^*(\mathcal{H}, x, x') \\ &= \eta(x, x') \mathbb{1}_{h(x') < h(x)} + (1 - \eta(x, x')) \mathbb{1}_{h(x') \geq h(x)} - \min\{\eta(x, x'), 1 - \eta(x, x')\} \\ &= \begin{cases} |2\eta(x, x') - 1|, & h \in \overline{\mathcal{H}}(x, x'), \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

This leads to

$$\langle \Delta \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}}(h, x, x') \rangle_{\epsilon} = \langle |2\eta(x, x') - 1| \rangle_{\epsilon} \mathbb{1}_{h \in \overline{\mathcal{H}}(x, x')}.$$

□

By Lemma F.1, the $(\mathcal{L}_{0-1}, \mathcal{H}_{\text{all}})$ -minimizability gap is

$$\mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}}) = \mathcal{R}_{\mathcal{L}_{0-1}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')}[\min\{\eta(x, x'), 1 - \eta(x, x')\}]. \quad (11)$$

F.1. Derivation for $\mathcal{L}_{\Phi_{\text{hinge}}}$.

For the hinge loss function $\Phi_{\text{hinge}}(u) := \max\{0, 1 - u\}$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\text{hinge}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{hinge}}}(h(x) - h(x')) \\ &= \eta(x, x') \max\{0, 1 - h(x') + h(x)\} + (1 - \eta(x, x')) \max\{0, 1 + h(x') - h(x)\}. \end{aligned}$$

Then,

$$\mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h, x, x') = 1 - |2\eta(x, x') - 1|.$$

The $(\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}})$ -minimizability gap is

$$\begin{aligned} \mathcal{M}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{all}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} [\mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}^*(x, x')] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} [1 - |2\eta(x, x') - 1|]. \end{aligned} \quad (12)$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}(h, x, x') &\geq \inf_{h \in \overline{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= \eta(x, x') \max\{0, 1 - 0\} + (1 - \eta(x, x')) \max\{0, 1 + 0\} - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= 1 - [1 - |2\eta(x, x') - 1|] \\ &= |2\eta(x, x') - 1|, \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}(h, x, x') \geq \langle |2\eta(x, x') - 1| \rangle_0 \mathbb{1}_{h \in \overline{\mathcal{H}}_{\text{all}}(x, x')} = \Delta \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}_{\text{all}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\mathcal{L}_{\Phi_{\text{hinge}}}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{all}}) - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}}). \quad (13)$$

F.2. Derivation for $\mathcal{L}_{\Phi_{\rho}}$.

For the ρ -margin loss function $\Phi_{\rho}(u) := \min\left\{1, \max\left\{0, 1 - \frac{u}{\rho}\right\}\right\}$, $\rho > 0$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\rho}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\rho}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\rho}}(h(x) - h(x')) \\ &= \eta(x, x') \min\left\{1, \max\left\{0, 1 - \frac{h(x') - h(x)}{\rho}\right\}\right\} + (1 - \eta(x, x')) \min\left\{1, \max\left\{0, 1 + \frac{h(x') - h(x)}{\rho}\right\}\right\} \\ &\geq \mathcal{C}_{\mathcal{L}_{0-1}}(h, x, x'). \end{aligned}$$

Then,

$$\mathcal{C}_{\mathcal{L}_{\Phi_{\rho}}, \mathcal{H}_{\text{all}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\rho}}}(h, x, x') = \min\{\eta(x, x'), 1 - \eta(x, x')\} = \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}_{\text{all}}}^*(x, x').$$

The $(\mathcal{L}_{\Phi_{\rho}}, \mathcal{H}_{\text{all}})$ -minimizability gap is

$$\begin{aligned} \mathcal{M}_{\mathcal{L}_{\Phi_{\rho}}}(\mathcal{H}_{\text{all}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} [\mathcal{C}_{\mathcal{L}_{\Phi_{\rho}}, \mathcal{H}_{\text{all}}}^*(x, x')] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} [\min\{\eta(x, x'), 1 - \eta(x, x')\}]. \end{aligned} \quad (14)$$

Therefore, for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{all}}}(h, x, x') \geq \langle |2\eta(x, x') - 1| \rangle_0 \mathbb{1}_{h \in \overline{\mathcal{H}}_{\text{all}}(x, x')} = \Delta \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}_{\text{all}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for \mathcal{L}_{Φ_ρ} , valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \mathcal{R}_{\mathcal{L}_{\Phi_\rho}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_\rho}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_\rho}}(\mathcal{H}_{\text{all}}) - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}}). \quad (15)$$

F.3. Derivation for $\mathcal{L}_{\Phi_{\text{exp}}}$.

For the exponential loss function $\Phi_{\text{exp}}(u) := e^{-u}$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\text{exp}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{exp}}}(h(x) - h(x')) \\ &= \eta(x, x') e^{-h(x') + h(x)} + (1 - \eta(x, x')) e^{h(x') - h(x)}. \end{aligned}$$

Then,

$$\mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{all}})(x, x') = \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h, x, x') = 2\sqrt{\eta(x, x')(1 - \eta(x, x'))}$$

. The $(\Phi_{\text{exp}}, \mathcal{H}_{\text{all}})$ -minimizability gap is:

$$\begin{aligned} \mathcal{M}_{\mathcal{L}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{all}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^*(x, x') \right] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[2\sqrt{\eta(x, x')(1 - \eta(x, x'))} \right]. \end{aligned} \quad (16)$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}(h, x, x') &\geq \inf_{h \in \overline{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= \eta(x, x') e^{-0} + (1 - \eta(x, x')) e^0 - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= 1 - 2\sqrt{\eta(x, x')(1 - \eta(x, x'))} \\ &= \left(\frac{2\eta(x, x') - 1}{\sqrt{\eta(x, x')} + \sqrt{1 - \eta(x, x')}} \right)^2 \\ &\geq \frac{(2\eta(x, x') - 1)^2}{2}, \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}(h, x, x') \geq \frac{(\Delta \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}_{\text{all}}}(h, x, x'))^2}{2}.$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\mathcal{L}_{\Phi_{\text{exp}}}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \sqrt{2} \left(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}} - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}}). \quad (17)$$

F.4. Derivation for $\mathcal{L}_{\Phi_{\text{log}}}$.

For the logistic loss function $\Phi_{\text{log}}(u) := \log_2(1 + e^{-u})$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{log}}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\text{log}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{log}}}(h(x) - h(x')) \\ &= \eta(x, x') \log_2 \left(1 + e^{-h(x') + h(x)} \right) + (1 - \eta(x, x')) \log_2 \left(1 + e^{h(x') - h(x)} \right). \end{aligned}$$

Then,

$$\begin{aligned}\mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{all}}}^*(x, x') &= \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\Phi_{\log}}(h, x, x') \\ &= -\eta(x, x') \log_2(\eta(x, x')) - (1 - \eta(x, x')) \log_2(1 - \eta(x, x')) \\ &\leq 2\sqrt{\eta(x, x')(1 - \eta(x, x'))} \quad (-a \log_2(a) - b \log_2(b) \leq 2\sqrt{ab}, a, b \in [0, 1])\end{aligned}$$

. The $(\Phi_{\log}, \mathcal{H}_{\text{all}})$ -minimizability gap is:

$$\begin{aligned}\mathcal{M}_{\mathcal{L}_{\Phi_{\log}}}(\mathcal{H}_{\text{all}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{all}}}^*(x, x') \right] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[-\eta(x, x') \log_2(\eta(x, x')) - (1 - \eta(x, x')) \log_2(1 - \eta(x, x')) \right].\end{aligned}\quad (18)$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned}\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{all}}}(h, x, x') &\geq \inf_{h \in \overline{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= \eta(x, x') \log_2(1 + e^{-0}) + (1 - \eta(x, x')) \log_2(1 + e^0) - \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &\geq 1 - 2\sqrt{\eta(x, x')(1 - \eta(x, x'))} \\ &= \left(\frac{2\eta(x, x') - 1}{\sqrt{\eta(x, x')} + \sqrt{1 - \eta(x, x')}} \right)^2 \\ &\geq \frac{(2\eta(x, x') - 1)^2}{2},\end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{all}}}(h, x, x') \geq \frac{(\Delta \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}_{\text{all}}}(h, x, x'))^2}{2}.$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\mathcal{L}_{\Phi_{\log}}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \sqrt{2} \left(\mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\log}}}(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}} - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}}).\quad (19)$$

F.5. Derivation for $\mathcal{L}_{\Phi_{\text{sq}}}$.

For the squared hinge loss function $\Phi_{\text{sq}}(u) := (1 - u)^2 \mathbb{1}_{u \leq 1}$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned}\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\text{sq}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{sq}}}(h(x) - h(x')) \\ &= \eta(x, x') (1 - h(x') + h(x))^2 \mathbb{1}_{h(x') - h(x) \leq 1} + (1 - \eta(x, x')) (1 + h(x') - h(x))^2 \mathbb{1}_{h(x') - h(x) \geq -1}.\end{aligned}$$

Then,

$$\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\Phi_{\text{sq}}}(h, x, x') = 4\eta(x, x')(1 - \eta(x, x')).$$

The $(\Phi_{\text{sq}}, \mathcal{H}_{\text{all}})$ -minimizability gap is:

$$\begin{aligned}\mathcal{M}_{\mathcal{L}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{all}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}^*(x, x') \right] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} [4\eta(x, x')(1 - \eta(x, x'))].\end{aligned}\quad (20)$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}(h, x, x') &\geq \inf_{h \in \overline{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= \eta(x, x') + (1 - \eta(x, x')) - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= 1 - 4\eta(x, x')(1 - \eta(x, x')) \\ &= (2\eta(x, x') - 1)^2, \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}(h, x, x') \geq (\Delta \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}_{\text{all}}}(h, x, x'))^2.$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\mathcal{L}_{\Phi_{\text{sq}}}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \left(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}} - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}}). \quad (21)$$

F.6. Derivation for $\mathcal{L}_{\Phi_{\text{sig}}}$.

For the sigmoid loss function $\Phi_{\text{sig}}(u) := 1 - \tanh(ku)$, $k > 0$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\text{sig}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{sig}}}(h(x) - h(x')) \\ &= \eta(x, x')(1 - \tanh(k[h(x') - h(x)])) + (1 - \eta(x, x'))(1 + \tanh(k[h(x) - h(x')])). \end{aligned}$$

Then,

$$\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h, x, x') = 1 - |1 - 2\eta(x, x')|.$$

. The $(\Phi_{\text{sig}}, \mathcal{H}_{\text{all}})$ -minimizability gap is:

$$\begin{aligned} \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{all}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') \right] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} [1 - |1 - 2\eta(x, x')|]. \end{aligned} \quad (22)$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}(h, x, x') &\geq \inf_{h \in \overline{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= 1 - |1 - 2\eta(x, x')| \tanh(0) - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= |1 - 2\eta(x, x')|, \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}(h, x, x') \geq \Delta \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}_{\text{all}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\mathcal{L}_{\Phi_{\text{sig}}}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{all}}) - \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}}). \quad (23)$$

G. Minimizability gaps can be non-zero for $\mathcal{H} = \mathcal{H}_{\text{all}}$ in the general pairwise misranking case.

Consider the uniform distribution that supports on three pairs $\{(x, x'), (x', x''), (x, x'')\}$. Let $\eta(x, x') = \eta(x', x'') = 1$ and $\eta(x, x'') = 0$. Note for any $h \in \mathcal{H}_{\text{all}}$, at least one of three difference $h(x') - h(x)$, $h(x'') - h(x')$, $h(x) - h(x'')$ is less than

or equal to 0 since the sum of all the difference is 0. Therefore,

$$\begin{aligned}
 \mathcal{R}_{\mathcal{L}_{0-1}}(h) &= \frac{1}{3}\mathcal{C}_{\mathcal{L}_{0-1}}(h, x, x') + \frac{1}{3}\mathcal{C}_{\mathcal{L}_{0-1}}(h, x', x'') + \frac{1}{3}\mathcal{C}_{\mathcal{L}_{0-1}}(h, x, x'') \\
 &= \frac{1}{3}\mathbb{1}_{h(x') < h(x)} + \frac{1}{3}\mathbb{1}_{h(x'') < h(x')} + \frac{1}{3}\mathbb{1}_{h(x'') \geq h(x)} \\
 &\geq \frac{1}{3} \\
 \mathcal{R}_{\mathcal{L}_{\Phi_\rho}}(h) &= \frac{1}{3}\mathcal{C}_{\mathcal{L}_{\Phi_\rho}}(h, x, x') + \frac{1}{3}\mathcal{C}_{\mathcal{L}_{0-1}}(h, x', x'') + \frac{1}{3}\mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h, x, x'') \\
 &= \frac{1}{3}\Phi_\rho(h(x') - h(x)) + \frac{1}{3}\Phi_\rho(h(x'') - h(x')) + \frac{1}{3}\Phi_\rho(h(x) - h(x'')) \\
 &\geq \frac{1}{3} \tag*{(\Phi_\rho(u) = 1, u \leq 0)} \\
 \mathcal{R}_{\mathcal{L}_{\Phi_{\text{convex}}}}(h) &= \frac{1}{3}\mathcal{C}_{\mathcal{L}_{\Phi_{\text{convex}}}}(h, x, x') + \frac{1}{3}\mathcal{C}_{\mathcal{L}_{0-1}}(h, x', x'') + \frac{1}{3}\mathcal{C}_{\mathcal{L}_{\Phi_{\text{convex}}}}(h, x, x'') \\
 &= \frac{1}{3}\Phi_{\text{convex}}(h(x') - h(x)) + \frac{1}{3}\Phi_{\text{convex}}(h(x'') - h(x')) + \frac{1}{3}\Phi_{\text{convex}}(h(x) - h(x'')) \\
 &\geq \Phi_{\text{convex}}\left(\frac{1}{3}[h(x') - h(x) + h(x'') - h(x') + h(x) - h(x'')]\right) \tag*{(convexity)} \\
 &= \Phi_{\text{convex}}(0) \\
 &\geq 1
 \end{aligned}$$

where all equality can be achieved by $h = 0$. Thus, using the fact that $\eta(x, x') = \eta(x', x'') = 1$ and $\eta(x, x'') = 0$, we obtain

$$\begin{aligned}
 \mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}_{\text{all}}) &= \mathcal{R}_{\mathcal{L}_{0-1}}^*(\mathcal{H}_{\text{all}}) = \frac{1}{3} \neq 0 \\
 \mathcal{M}_{\mathcal{L}_{\Phi_\rho}}(\mathcal{H}_{\text{all}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_\rho}}^*(\mathcal{H}_{\text{all}}) = \frac{1}{3} \neq 0 \\
 \mathcal{M}_{\mathcal{L}_{\Phi_{\text{convex}}}}(\mathcal{H}_{\text{all}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{convex}}}}^*(\mathcal{H}_{\text{all}}) \geq 1 \neq 0.
 \end{aligned}$$

H. Characterization of distribution order and minimizability gap (Proof of Theorem 2.5, Theorem 2.6 and Theorem 2.7)

Theorem 2.5 (Characterization of distribution order). *The distribution order is transitive and there exists a dense countable subset $\tilde{\mathcal{X}} \subset \mathcal{X}$ with respect to the distribution order if and only if there exists $h \in \mathcal{H}_{\text{all}}$ inducing the distribution order.*

Proof. \Leftarrow : Assume there exists $h \in \mathcal{H}_{\text{all}}$ inducing the distribution order. For all $x, x', x'' \in \mathcal{X}$ such that $x \stackrel{\mathcal{D}}{\preceq} x'$ and $x' \stackrel{\mathcal{D}}{\preceq} x''$, by definition, we have $h(x) \leq h(x')$ and $h(x') \leq h(x'')$, which implies that $h(x) \leq h(x'')$. Then, by definition, we obtain $x \stackrel{\mathcal{D}}{\preceq} x''$ and conclude the distribution is transitive. Furthermore, we can construct the countable set $\tilde{\mathcal{X}}$ using the following procedure: for any open interval $(h(x), h(x')), x, x' \in \mathcal{X}$ such that $h^{-1}((h(x), h(x'))) = \emptyset$, we pick x, x' in $\tilde{\mathcal{X}}$. Those intervals are countable since any of them contain a rational number and any two of them are disjoint. For any open interval with rational endpoints m, n such that $h^{-1}((m, n)) \neq \emptyset$, pick any $x \in h^{-1}((m, n))$ in $\tilde{\mathcal{X}}$. Again, those open intervals are also countable since rational numbers are countable. Thus, $\tilde{\mathcal{X}}$ is countable. Next, we verify that $\tilde{\mathcal{X}}$ is dense. Indeed, for any $x, x' \in \mathcal{X}$ that satisfy $x \stackrel{\mathcal{D}}{\preceq} x'$ and not $x' \stackrel{\mathcal{D}}{\preceq} x$, we have $h(x) < h(x')$. If $h^{-1}((h(x), h(x'))) = \emptyset$, by the procedure, we know there exists $\bar{x} \in \tilde{\mathcal{X}}$ such that $x \stackrel{\mathcal{D}}{\preceq} \bar{x} \stackrel{\mathcal{D}}{\preceq} x'$; if $h^{-1}((h(x), h(x'))) \neq \emptyset$ and assume $x^* \in h^{-1}((h(x), h(x')))$, then, we can take rational numbers m, n such that $h(x) < m < h(x^*) < n < h(x')$, by the procedure, we know there exists $\bar{x} \in \tilde{\mathcal{X}}$ such that $x \stackrel{\mathcal{D}}{\preceq} \bar{x} \stackrel{\mathcal{D}}{\preceq} x'$. In conclusion, $\tilde{\mathcal{X}}$ is dense countable.

\Rightarrow : Let $\tilde{\mathcal{X}} = \{\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots\}$. For any $x \in \mathcal{X}$, we use the following notation for convenience: $\tilde{\mathcal{L}}(x) =$

$\left\{n \in \mathbb{N} : \bar{x}_n \stackrel{\mathcal{D}}{\leq} x \text{ and not } x \stackrel{\mathcal{D}}{\leq} \bar{x}_n\right\}$ and $\mathcal{R}(x) = \left\{n \in \mathbb{N} : x \stackrel{\mathcal{D}}{\leq} \bar{x}_n \text{ and not } \bar{x}_n \stackrel{\mathcal{D}}{\leq} x\right\}$. Then, take

$$h^*(x) = \sum_{n \in \mathcal{L}(x)} \frac{1}{2^n} - \sum_{n \in \mathcal{R}(x)} \frac{1}{2^n}.$$

Next, we verify that h^* induces the distribution order. Indeed, for any $x \stackrel{\mathcal{D}}{\leq} x'$, by the transitivity of distribution order, we have

$$\mathcal{L}(x) \subseteq \mathcal{L}(x'), \quad \mathcal{R}(x') \subseteq \mathcal{R}(x).$$

which implies that $h(x) \leq h(x')$. Also, for any $x' \stackrel{\mathcal{D}}{\leq} x$ and $x \stackrel{\mathcal{D}}{\leq} x'$, we have $h(x) = h(x')$. Moreover, for any $x' \stackrel{\mathcal{D}}{\leq} x$ and not $x \stackrel{\mathcal{D}}{\leq} x'$, there exists $\bar{x} \in \bar{\mathcal{X}}$ such that $x' \stackrel{\mathcal{D}}{\leq} \bar{x} \stackrel{\mathcal{D}}{\leq} x$. Therefore, \bar{x} belongs to at least one of $\mathcal{L}(x)$ and $\mathcal{R}(x')$. Since $\bar{x} \notin \mathcal{L}(x')$ and $\bar{x} \notin \mathcal{R}(x)$, we obtain

$$\text{either } \mathcal{L}(x') \subset \mathcal{L}(x) \text{ or } \mathcal{R}(x) \subset \mathcal{R}(x'),$$

which implies that $h(x) > h(x')$. Therefore, h^* induces the distribution order. \square

Theorem 2.6. *Assume that the distribution order is a total order and $\eta(x, x')$ is continuous on $\mathcal{X} \times \mathcal{X}$. Then, there exists $h \in \mathcal{H}_{\text{all}}$ inducing the distribution order.*

Proof. Define $\stackrel{\mathcal{D}}{<}$ to be the relation associated with $\stackrel{\mathcal{D}}{\leq}$ that defined as $x \stackrel{\mathcal{D}}{<} x'$ if $x \stackrel{\mathcal{D}}{\leq} x'$ and $x \neq x'$. Let $f(x, x') = \eta(x, x') - \eta(x', x)$. By the assumption, for all $x, x' \in \mathcal{X}$,

- $x \stackrel{\mathcal{D}}{<} x' \iff f(x, x') > 0$
- $f(x, x')$ is continuous on both x and $x' \implies \{\bar{x} \in \mathcal{X} \mid f(x, \bar{x}) > 0\} = \{\bar{x} \in \mathcal{X} \mid x \stackrel{\mathcal{D}}{<} \bar{x}\}$ is open in \mathcal{X} .
- $f(x, x') = -f(x', x)$. In particular, $f(x, x) = 0$.

Now let's assume $x \stackrel{\mathcal{D}}{<} x'$, i.e., $f(x, x') > 0$ so $f(x', x) < 0$. Consider the following continuous functions on $[0, 1]$,

$$\begin{aligned} g_x(t) &= f(x, tx + (1-t)x') \\ g_{x'}(t) &= f(x', tx + (1-t)x'). \end{aligned}$$

We know that

- $g_x(0) > 0, g_x(1) = 0 \implies g_x(t) > 0$ when $t \in (0, 1)$ (because $t = 1$ is the only zero point of $g_x(t)$)
- $g_{x'}(0) = 0, g_{x'}(1) < 0 \implies g_{x'}(t) < 0$ when $t \in (0, 1)$ (because $t = 0$ is the only zero point of $g_{x'}(t)$).

Therefore, $\{\bar{x} \in \mathcal{X} \mid x \stackrel{\mathcal{D}}{<} \bar{x} \stackrel{\mathcal{D}}{<} x'\} \neq \emptyset$. Note $\{\bar{x} \in \mathcal{X} \mid x \stackrel{\mathcal{D}}{<} \bar{x} \stackrel{\mathcal{D}}{<} x'\} = \{\bar{x} \in \mathcal{X} \mid \bar{x} \stackrel{\mathcal{D}}{<} x'\} \cup \{\bar{x} \in \mathcal{X} \mid x \stackrel{\mathcal{D}}{<} \bar{x}\}$, the intersection of two open subsets, which is also open. Any nonempty open set includes at least one rational point, pick such point in $\bar{\mathcal{X}}$ and we obtain that $\bar{\mathcal{X}}$ is dense countable. By Theorem 2.5, we conclude that there exists $h \in \mathcal{H}_{\text{all}}$ inducing the distribution order. \square

Theorem 2.7. *Assume that for all $x, x' \in \mathcal{X}$, $\eta(x, x') + \eta(x', x) = 1$. Then, for any hypothesis set \mathcal{H} , if there exists $h \in \mathcal{H}$ inducing the distribution order, the minimizability gap of the pairwise misranking loss is null, $\mathcal{M}_{\mathcal{L}_{0-1}}(\mathcal{H}) = 0$.*

Proof. Assume that $h^* \in \mathcal{H}$ induces the distribution order. Then,

$$\begin{aligned}
 \mathcal{R}_{L_{0-1}}(h^*) &= \mathbb{E}_{(x,x') \sim \mathcal{D}} [L_{0-1}(h, x, x', y)] \\
 &= \mathbb{E}_{(X, X')} [\eta(x, x') \mathbb{1}_{h^*(x') < h^*(x)} + (1 - \eta(x, x')) \mathbb{1}_{h^*(x') \geq h^*(x)}] \\
 &= \mathbb{E}_{(X, X')} [\eta(x, x') \mathbb{1}_{\eta(x', x) > \eta(x, x')} + (1 - \eta(x, x')) \mathbb{1}_{\eta(x, x') \geq \eta(x', x)}] \quad (h^* \in \mathcal{H} \text{ induces the distribution order.}) \\
 &= \mathbb{E}_{(X, X')} [\eta(x, x') \mathbb{1}_{1 - \eta(x, x') > \eta(x, x')} + (1 - \eta(x, x')) \mathbb{1}_{\eta(x, x') \geq 1 - \eta(x, x')}] \quad (\eta(x, x') + \eta(x', x) = 1) \\
 &= \mathbb{E}_{(X, X')} [\min\{\eta(x, x'), 1 - \eta(x, x')\}] \\
 &= \mathbb{E}_{(X, X')} [\mathcal{C}_{L_{0-1}}^*(\mathcal{H}, x, x')].
 \end{aligned}$$

Therefore, $\mathcal{M}_{L_{0-1}}(\mathcal{H}) = \mathcal{R}_{L_{0-1}}^*(\mathcal{H}) - \mathbb{E}_{(X, X')} [\mathcal{C}_{L_{0-1}}^*(\mathcal{H}, x, x')] = 0$. \square

I. Negative results for bipartite ranking (Proof of Theorem 4.1)

Theorem 4.1 (Negative results for bipartite ranking). *Assume that \mathcal{X} contains an interior point x_0 and that \mathcal{H} is regular for bipartite ranking, contains 0 and is equicontinuous at x_0 . If for some function f that is non-decreasing and continuous at 0, the following bound holds for all $h \in \mathcal{H}$ and any distribution,*

$$\mathcal{R}_{L_{0-1}}(h) - \mathcal{R}_{L_{0-1}}^*(\mathcal{H}) \leq f(\mathcal{R}_{L_{\Phi}}(h) - \mathcal{R}_{L_{\Phi}}^*(\mathcal{H})),$$

then, $f(t) \geq \frac{1}{2}$ for any $t \geq 0$.

Proof. Assume $x_0 \in \mathcal{X}$ is an interior point and $h_0 = 0 \in \mathcal{H}$. By the assumption that x_0 is an interior point and \mathcal{H} is equicontinuous at x_0 , for any $\epsilon > 0$, we are able to take $x' \neq x_0 \in \mathcal{X}$ such that $|h(x') - h(x_0)| < \epsilon$ for all $h \in \mathcal{H}$. Consider the distribution that supports on $\{x_0, x'\}$ with $\eta(x_0) = 1$ and $\eta(x') = 0$. Then, for any $h \in \mathcal{H}$,

$$\mathcal{R}_{L_{0-1}}(h) = \mathcal{C}_{L_{0-1}}(h, x_0, x') = \mathbb{1}_{h(x_0) < h(x')} + \frac{1}{2} \mathbb{1}_{h(x_0) = h(x')} \geq 0,$$

where the equality can be achieved for some $h \in \mathcal{H}$ since \mathcal{H} is regular for bipartite ranking. Therefore,

$$\mathcal{R}_{L_{0-1}}^*(\mathcal{H}) = \mathcal{C}_{L_{0-1}}^*(\mathcal{H}, x_0, x') = \inf_{h \in \mathcal{H}} \mathcal{C}_{L_{0-1}}(h, x_0, x') = 0.$$

Note $\mathcal{R}_{L_{0-1}}(h_0) = \frac{1}{2}$. For the surrogate loss L_{Φ} , for any $h \in \mathcal{H}$,

$$\mathcal{R}_{L_{\Phi}}(h) = \mathcal{C}_{L_{\Phi}}(h, x_0, x') = \Phi(h(x_0) - h(x')) \in [\Phi(\epsilon), \Phi(-\epsilon)]$$

since $|h(x') - h(x_0)| < \epsilon$ and Φ is non-increasing. Therefore,

$$\mathcal{R}_{L_{\Phi}}^*(\mathcal{H}) = \mathcal{C}_{L_{\Phi}}^*(\mathcal{H}, x_0, x') \geq \Phi(\epsilon).$$

Note $\mathcal{R}_{L_{\Phi}}(h_0) = \Phi(0)$. If for some function f that is non-decreasing and continuous at 0, the bound holds, then, we obtain for any $h \in \mathcal{H}$ and $\epsilon > 0$,

$$\mathcal{R}_{L_{0-1}}(h) - 0 \leq f(\mathcal{R}_{L_{\Phi}}(h) - \mathcal{R}_{L_{\Phi}}^*(\mathcal{H})) \leq f(\mathcal{R}_{L_{\Phi}}(h) - \Phi(\epsilon)).$$

Let $h = h_0$, then $f(\Phi(0) - \Phi(\epsilon)) \geq \frac{1}{2}$ for any $\epsilon > 0$. Take $\epsilon \rightarrow 0$, we obtain $f(0) \geq \frac{1}{2}$ using the fact that Φ and f are both continuous at 0. Since f is non-decreasing, for any $t \in [0, 1]$, $f(t) \geq \frac{1}{2}$. \square

J. \mathcal{H}_{all} - consistency bounds for bipartite misranking losses

We first characterize the minimal conditional $\tilde{\mathcal{L}}_{0-1}$ -risk and the calibration gap of the bipartite misranking loss for a broad class of hypothesis sets. We let $\tilde{\mathcal{H}}(x, x') = \{h \in \mathcal{H}: (h(x) - h(x'))(\eta(x) - \eta(x')) < 0\}$ and $\mathring{\mathcal{H}}(x, x') = \{h \in \mathcal{H}: h(x) = h(x')\}$ for convenience.

Lemma J.1. *Assume that \mathcal{H} is regular for bipartite ranking. Then, the minimal conditional $\tilde{\mathcal{L}}_{0-1}$ -risk is*

$$\mathcal{C}_{\tilde{\mathcal{L}}_{0-1}}^*(\mathcal{H}, x, x') = \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\}.$$

The calibration gap of $\tilde{\mathcal{L}}_{0-1}$ can be characterized as

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}}(h, x, x') = |\eta(x) - \eta(x')| \mathbb{1}_{h \in \tilde{\mathcal{H}}(x, x')} + \frac{1}{2} |\eta(x) - \eta(x')| \mathbb{1}_{h \in \mathring{\mathcal{H}}(x, x')}.$$

Proof. By the definition, the conditional $\tilde{\mathcal{L}}_{0-1}$ -risk is

$$\mathcal{C}_{\tilde{\mathcal{L}}_{0-1}}(h, x, x') = \eta(x)(1 - \eta(x')) \left[\mathbb{1}_{h(x) - h(x') < 0} + \frac{1}{2} \mathbb{1}_{h(x) = h(x')} \right] + \eta(x')(1 - \eta(x)) \left[\mathbb{1}_{h(x) - h(x') > 0} + \frac{1}{2} \mathbb{1}_{h(x) = h(x')} \right].$$

For any $x \neq x' \in \mathcal{X}$, by the assumption, there exists $h^* \in \mathcal{H}$ such that

$$(h^*(x) - h^*(x'))(\eta(x) - \eta(x')) \mathbb{1}_{\eta(x) \neq \eta(x')} > 0.$$

Therefore, the optimal conditional $\tilde{\mathcal{L}}_{0-1}$ -risk can be characterized as for any $x \neq x' \in \mathcal{X}$,

$$\mathcal{C}_{\tilde{\mathcal{L}}_{0-1}}^*(\mathcal{H}, x, x') = \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}}(h^*, x, x') = \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\}$$

which proves the first part of lemma. By the definition,

$$\begin{aligned} \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}}(h, x, x') &= \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}}^*(\mathcal{H}, x, x') \\ &= \eta(x)(1 - \eta(x')) \left[\mathbb{1}_{h(x) - h(x') < 0} + \frac{1}{2} \mathbb{1}_{h(x) = h(x')} \right] \\ &\quad + \eta(x')(1 - \eta(x)) \left[\mathbb{1}_{h(x) - h(x') > 0} + \frac{1}{2} \mathbb{1}_{h(x) = h(x')} \right] \\ &\quad - \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \\ &= \begin{cases} |\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))|, & h \in \tilde{\mathcal{H}}(x, x'), \\ \frac{1}{2} |\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))|, & h \in \mathring{\mathcal{H}}(x, x'), \\ 0, & \text{otherwise.} \end{cases} \\ &= \begin{cases} |\eta(x) - \eta(x')|, & h \in \tilde{\mathcal{H}}(x, x'), \\ \frac{1}{2} |\eta(x) - \eta(x')|, & h \in \mathring{\mathcal{H}}(x, x'), \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

Table 9: \mathcal{H}_{all} -consistency upper bounds for bipartite abstention losses.

Loss function	\mathcal{H}_{all} -consistency upper bound
$\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}$	$\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}^*}(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{all}})$
$\tilde{\mathcal{L}}_{\Phi_{\rho}}$	$\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}^*}(\mathcal{H}_{\text{all}})$
$\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}$	$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}^*}(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}}$
$\tilde{\mathcal{L}}_{\Phi_{\text{log}}}$	$\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{log}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{log}}}^*}(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}}$
$\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}$	$\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}^*}(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}}$
$\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}$	$\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}^*}(\mathcal{H}_{\text{all}})$

□

Note that for $h_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}_{\text{all}}}^{\text{Bayes}}(x) := \eta(x)$,

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}\left(h_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}_{\text{all}}}^{\text{Bayes}}\right) = \mathbb{E}_{(X, X')}[\min\{\eta(x, x'), 1 - \eta(x, x')\}].$$

Therefore, by Lemma J.1, the $(\tilde{\mathcal{L}}_{0-1}, \mathcal{H}_{\text{all}})$ -minimizability gap is

$$\mathcal{M}_{\tilde{\mathcal{L}}_{0-1}}(\mathcal{H}_{\text{all}}) = \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')}[\min\{\eta(x, x'), 1 - \eta(x, x')\}] = 0. \quad (24)$$

J.1. Derivation for $\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}$.

For the hinge loss function $\Phi_{\text{hinge}}(u) := \max\{0, 1 - u\}$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h, x, x') &= \eta(x)(1 - \eta(x'))\Phi_{\text{hinge}}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\text{hinge}}(h(x') - h(x)) \\ &= \eta(x)(1 - \eta(x')) \max\{0, 1 - h(x) + h(x')\} + \eta(x')(1 - \eta(x)) \max\{0, 1 + h(x) - h(x')\}. \end{aligned}$$

Then,

$$\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h, x, x') = 2 \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\}.$$

The $(\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}})$ -minimizability gap is

$$\begin{aligned} \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{all}}) &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}^*(x, x') \right] \\ &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} [2 \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\}]. \end{aligned} \quad (25)$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}(h, x, x') &\geq \inf_{h \in \tilde{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= \eta(x)(1 - \eta(x')) \max\{0, 1 - 0\} + \eta(x')(1 - \eta(x)) \max\{0, 1 + 0\} - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 2 \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \\ &= |\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))|. \end{aligned}$$

Similarly, $\forall h \in \mathring{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}(h, x, x') &\geq \inf_{h \in \mathring{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= \eta(x)(1 - \eta(x')) \max\{0, 1 - 0\} + \eta(x')(1 - \eta(x)) \max\{0, 1 + 0\} - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ &= \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 2 \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \\ &= |\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))|, \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{all}}}(h, x, x') \geq \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{all}}). \quad (26)$$

J.2. Derivation for $\tilde{\mathcal{L}}_{\Phi_\rho}$.

For the ρ -margin loss function $\Phi_\rho(u) := \min\left\{1, \max\left\{0, 1 - \frac{u}{\rho}\right\}\right\}$, $\rho > 0$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(h, x, x') &= \eta(x)(1 - \eta(x'))\Phi_\rho(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_\rho(h(x') - h(x)) \\ &= \eta(x)(1 - \eta(x')) \min\left\{1, \max\left\{0, 1 - \frac{h(x) - h(x')}{\rho}\right\}\right\} \\ &\quad + \eta(x')(1 - \eta(x)) \min\left\{1, \max\left\{0, 1 + \frac{h(x) - h(x')}{\rho}\right\}\right\} \end{aligned}$$

Then,

$$\begin{aligned} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{all}}}^*(x, x') &= \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(h, x, x') \\ &= \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \\ &= \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}_{\text{all}}}^*(x, x'). \end{aligned}$$

Note that for $h_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{all}}}^\alpha(x) := \alpha \eta(x)$, $\alpha > 0$, by the Lebesgue's dominated convergence theorem,

$$\liminf_{\alpha \rightarrow +\infty} \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(h_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{all}}}^\alpha) = \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{all}}}^*(x, x') \right] \geq \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}^*(\mathcal{H}_{\text{all}}) \geq \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{all}}}^*(x, x') \right].$$

Therefore, the $(\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{all}})$ -minimizability gap is

$$\mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(\mathcal{H}_{\text{all}}) = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{all}}}^*(x, x') \right] = 0. \quad (27)$$

Furthermore, for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{all}}}(h, x, x') \geq \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_\rho}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}^*(\mathcal{H}_{\text{all}}) \leq \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}^*(\mathcal{H}_{\text{all}}). \quad (28)$$

J.3. Derivation for $\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}$.

For the exponential loss function $\Phi_{\text{exp}}(u) := e^{-u}$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h, x, x') &= \eta(x)(1 - \eta(x'))\Phi_{\text{exp}}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\text{exp}}(h(x') - h(x)) \\ &= \eta(x)(1 - \eta(x'))e^{-h(x)+h(x')} + \eta(x')(1 - \eta(x))e^{h(x)-h(x')}. \end{aligned}$$

Then,

$$\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h, x, x') = 2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))}.$$

Note that for $h_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^{\text{Bayes}}(x) := \frac{1}{2} \log \frac{\eta(x)(1 - \eta(x^*))}{\eta(x^*)(1 - \eta(x))}$ with fixed $x^* \in \mathcal{X}$,

$$\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}\left(h_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^{\text{Bayes}}\right) = \mathbb{E}_{(X, X')} \left[2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))} \right].$$

Therefore, the $(\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}})$ -minimizability gap is

$$\mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{all}}) = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^*(x, x') \right] = 0. \quad (29)$$

Furthermore, $\forall h \in \tilde{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}(h, x, x') \\ & \geq \inf_{h \in \tilde{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x'))e^{-0} + \eta(x')(1 - \eta(x))e^0 - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x')))} \\ & = \left(\frac{\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))}{\sqrt{\eta(x)(1 - \eta(x'))} + \sqrt{\eta(x')(1 - \eta(x))}} \right)^2 \\ & \geq (\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x)))^2 \end{aligned}$$

Similarly, $\forall h \in \mathring{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}(h, x, x') \\ & \geq \inf_{h \in \mathring{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x'))e^{-0} + \eta(x')(1 - \eta(x))e^0 - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x')))} \\ & \geq (\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x)))^2, \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{all}}}(h, x, x') \geq (\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}}(h, x, x'))^2.$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}}. \quad (30)$$

J.4. Derivation for $\tilde{\mathcal{L}}_{\Phi_{\log}}$.

For the logistic loss function $\Phi_{\log}(u) := \log_2(1 + e^{-u})$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(h, x, x') \\ & = \eta(x)(1 - \eta(x'))\Phi_{\log}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\log}(h(x') - h(x)) \\ & = \eta(x)(1 - \eta(x'))\log_2\left(1 + e^{-h(x)+h(x')}\right) + \eta(x')(1 - \eta(x))\log_2\left(1 + e^{h(x)-h(x')}\right). \end{aligned}$$

Then,

$$\begin{aligned} & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(h, x, x') \\ & = -\eta(x)(1 - \eta(x'))\log_2(\eta(x)(1 - \eta(x'))) - \eta(x')(1 - \eta(x))\log_2(\eta(x')(1 - \eta(x))) \\ & \leq 2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))} \quad (-a \log_2(a) - b \log_2(b) \leq 2\sqrt{ab}, a, b \in [0, 1]) \end{aligned}$$

Note that for $h_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}^{\text{Bayes}}(x) := \log \frac{\eta(x)(1-\eta(x^*))}{\eta(x^*)(1-\eta(x))}$ with fixed $x^* \in \mathcal{X}$,

$$\mathcal{R}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}}} \left(h_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}^{\text{Bayes}} \right) = \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}^*(x, x') \right].$$

Therefore, the $(\widetilde{\mathcal{L}}_{\Phi_{1\log}}, \mathcal{H}_{\text{all}})$ -minimizability gap is

$$\mathcal{M}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}}}(\mathcal{H}_{\text{all}}) = \mathcal{R}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}^*(x, x') \right] = 0. \quad (31)$$

Furthermore, $\forall h \in \widetilde{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}(h, x, x') \\ & \geq \inf_{h \in \widetilde{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}}}(h, x, x') - \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}^*(x, x') \\ & = \eta(x)(1-\eta(x')) \log_2(1+e^{-0}) + \eta(x')(1-\eta(x)) \log_2(1+e^0) - \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}^*(x, x') \\ & \geq \eta(x)(1-\eta(x')) + \eta(x')(1-\eta(x)) - 2\sqrt{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x')))} \\ & = \left(\frac{\eta(x)(1-\eta(x')) - \eta(x')(1-\eta(x))}{\sqrt{\eta(x)(1-\eta(x'))} + \sqrt{\eta(x')(1-\eta(x))}} \right)^2 \\ & \geq (\eta(x)(1-\eta(x')) - \eta(x')(1-\eta(x)))^2 \end{aligned}$$

Similarly, $\forall h \in \widetilde{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}(h, x, x') \\ & \geq \inf_{h \in \widetilde{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}}}(h, x, x') - \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}^*(x, x') \\ & = \eta(x)(1-\eta(x'))e^{-0} + \eta(x')(1-\eta(x))e^0 - \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}^*(x, x') \\ & \geq \eta(x)(1-\eta(x')) + \eta(x')(1-\eta(x)) - 2\sqrt{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x')))} \\ & \geq (\eta(x)(1-\eta(x')) - \eta(x')(1-\eta(x)))^2, \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}, \mathcal{H}_{\text{all}}}}(h, x, x') \geq (\Delta \mathcal{C}_{\widetilde{\mathcal{L}}_{0-1}, \mathcal{H}_{\mathcal{C}}}(h, x, x'))^2.$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\widetilde{\mathcal{L}}_{\Phi_{1\log}}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\widetilde{\mathcal{L}}_{0-1}}(h) - \mathcal{R}_{\widetilde{\mathcal{L}}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \left(\mathcal{R}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}}}(h) - \mathcal{R}_{\widetilde{\mathcal{L}}_{\Phi_{1\log}}}^*(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}}. \quad (32)$$

J.5. Derivation for $\widetilde{\mathcal{L}}_{\Phi_{\text{sq}}}$.

For the squared hinge loss function $\Phi_{\text{sq}}(u) := (1-u)^2 \mathbb{1}_{u \leq 1}$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} & \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h, x, x') \\ & = \eta(x)(1-\eta(x'))\Phi_{\text{sq}}(h(x) - h(x')) + \eta(x')(1-\eta(x))\Phi_{\text{sq}}(h(x') - h(x)) \\ & = \eta(x)(1-\eta(x'))(1-h(x) + h(x'))^2 \mathbb{1}_{h(x)-h(x') \leq 1} + \eta(x')(1-\eta(x))(1+h(x) - h(x'))^2 \mathbb{1}_{h(x)-h(x') \geq -1}. \end{aligned}$$

Then,

$$\mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{\text{sq}}, \mathcal{H}_{\text{all}}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\widetilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h, x, x') = 4 \frac{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x'))}{\eta(x)(1-\eta(x')) + \eta(x')(1-\eta(x))}.$$

The $(\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}})$ -minimizability gap is

$$\mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{all}}) = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}^*(x, x') \right]. \quad (33)$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}(h, x, x') \\ & \geq \inf_{h \in \tilde{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 4 \frac{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))}{\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x))} \\ & = \frac{(\eta(x)(1 - \eta(x)) - \eta(x')(1 - \eta(x')))^2}{\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x))} \\ & \geq (\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x)))^2 \quad (a + b - 2ab \leq 1, a, b \in [0, 1]) \end{aligned}$$

Similarly, $\forall h \in \mathcal{H}_{\text{all}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}(h, x, x') \\ & \geq \inf_{h \in \mathcal{H}_{\text{all}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 4 \frac{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))}{\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x))} \\ & \geq (\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x)))^2, \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{all}}}(h, x, x') \geq (\Delta \mathcal{C}_{\mathcal{L}_{0-1}, \mathcal{H}_{\text{all}}}(h, x, x'))^2.$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{all}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{all}}) \right)^{\frac{1}{2}}. \quad (34)$$

J.6. Derivation for $\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}$.

For the sigmoid loss function $\Phi_{\text{sig}}(u) := 1 - \tanh(ku)$, $k > 0$, for all $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$, $x' \in \mathcal{X}$ and $x \neq x'$:

$$\begin{aligned} & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h, x, x') \\ & = \eta(x)(1 - \eta(x'))\Phi_{\text{sig}}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\text{sig}}(h(x') - h(x)) \\ & = \eta(x)(1 - \eta(x'))(1 - \tanh(k[h(x) - h(x')])) + \eta(x')(1 - \eta(x))(1 + \tanh(k[h(x) - h(x')])) \end{aligned}$$

Then,

$$\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{all}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h, x, x') = 2 \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\}.$$

Note that for $h_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^\alpha(x) := \alpha \eta(x)$, $\alpha > 0$, by the Lebesgue's dominated convergence theorem,

$$\liminf_{\alpha \rightarrow +\infty} \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^\alpha) = \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') \right] \geq \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{all}}) \geq \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') \right].$$

Therefore, the $(\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}})$ -minimizability gap is

$$\mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{all}}) = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{all}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') \right] = 0. \quad (35)$$

Furthermore, $\forall h \in \tilde{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}(h, x, x') \\ & \geq \inf_{h \in \tilde{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 2 \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \\ & = |\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))|. \end{aligned}$$

Similarly, $\forall h \in \mathring{\mathcal{H}}_{\text{all}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}(h, x, x') \\ & \geq \inf_{h \in \mathring{\mathcal{H}}_{\text{all}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}^*(x, x') \\ & = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 2 \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \\ & = |\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))|, \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{all}}$, $x \in \mathcal{X}$ and $x' \in \mathcal{X}$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{all}}}(h, x, x') \geq \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}_{\mathcal{C}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{all} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}$, valid for all $h \in \mathcal{H}_{\text{all}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}}^*(\mathcal{H}_{\text{all}}) \leq \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{all}}). \quad (36)$$

K. \mathcal{H} - consistency bounds for pairwise abstention loss

We first characterize the minimal conditional $\mathbb{L}_{0-1}^{\text{abs}}$ -risk and the calibration gap of $\mathbb{L}_{0-1}^{\text{abs}}$ for a broad class of hypothesis sets. We let $\tilde{\mathcal{H}}(x, x') = \{h \in \mathcal{H} : \text{sign}(h(x') - h(x))(2\eta(x, x') - 1) \leq 0\}$ for convenience.

Lemma K.1. *Assume that \mathcal{H} is regular for general pairwise ranking. Then, the minimal conditional $\mathbb{L}_{0-1}^{\text{abs}}$ -risk is*

$$\mathcal{C}_{\mathbb{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}, x, x') = \min\{\eta(x, x'), 1 - \eta(x, x')\} \mathbb{1}_{\|x-x'\| > \gamma} + c \mathbb{1}_{\|x-x'\| \leq \gamma}.$$

The calibration gap of $\mathbb{L}_{0-1}^{\text{abs}}$ can be characterized as

$$\Delta \mathcal{C}_{\mathbb{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\mathcal{C}}}(h, x, x') = |2\eta(x, x') - 1| \mathbb{1}_{h \in \tilde{\mathcal{H}}(x, x')} \mathbb{1}_{\|x-x'\| > \gamma}.$$

Proof. By the definition, the conditional $\mathbb{L}_{0-1}^{\text{abs}}$ -risk is

$$\mathcal{C}_{\mathbb{L}_{0-1}^{\text{abs}}}(h, x, x') = (\eta(x, x') \mathbb{1}_{h(x') < h(x)} + (1 - \eta(x, x')) \mathbb{1}_{h(x') \geq h(x)}) \mathbb{1}_{\|x-x'\| > \gamma} + c \mathbb{1}_{\|x-x'\| \leq \gamma}.$$

For any (x, x') such that $\|x - x'\| \leq \gamma$ and $h \in \mathcal{H}$, $\mathcal{C}_{\mathbb{L}_{0-1}^{\text{abs}}}(h, x, x) = \mathcal{C}_{\mathbb{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}, x, x) = c$. For any (x, x') such that $\|x - x'\| > \gamma$, by the assumption, there exists $h^* \in \mathcal{H}$ such that $\text{sign}(h^*(x') - h^*(x)) = \text{sign}(2\eta(x, x') - 1)$. Therefore, the optimal conditional $\mathbb{L}_{0-1}^{\text{abs}}$ -risk can be characterized as for any $x, x' \in \mathcal{X}$,

$$\mathcal{C}_{\mathbb{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}, x, x') = \mathcal{C}_{\mathbb{L}_{0-1}^{\text{abs}}}(h^*, x, x') = \min\{\eta(x, x'), 1 - \eta(x, x')\} \mathbb{1}_{\|x-x'\| > \gamma} + c \mathbb{1}_{\|x-x'\| \leq \gamma}.$$

which proves the first part of lemma. By the definition, for any (x, x') such that $\|x - x'\| \leq \gamma$ and $h \in \mathcal{H}$, $\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x') = \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}, x, x') = 0$. For any (x, x') such that $\|x - x'\| > \gamma$ and $h \in \mathcal{H}$,

$$\begin{aligned} \Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x') &= \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}, x, x') \\ &= \eta(x, x') \mathbb{1}_{h(x') < h(x)} + (1 - \eta(x, x')) \mathbb{1}_{h(x') \geq h(x)} - \min\{\eta(x, x'), 1 - \eta(x, x')\} \\ &= \begin{cases} |2\eta(x, x') - 1|, & h \in \overline{\mathcal{H}}(x, x'), \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

This leads to

$$\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x') = |2\eta(x, x') - 1| \mathbb{1}_{h \in \overline{\mathcal{H}}(x, x')} \mathbb{1}_{\|x - x'\| > \gamma}.$$

□

K.1. Linear Hypotheses

Since \mathcal{H}_{lin} satisfies the condition of Lemma K.1, by Lemma K.1 the $(\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}})$ -minimizability gap can be expressed as follows:

$$\mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}) = \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')}[\min\{\eta(x, x'), 1 - \eta(x, x')\} \mathbb{1}_{\|x - x'\| > \gamma} + c \mathbb{1}_{\|x - x'\| \leq \gamma}]. \quad (37)$$

By the definition of \mathcal{H}_{lin} , for any $(x, x') \in \mathcal{X} \times \mathcal{X}$, $\{h(x') - h(x) \mid h \in \mathcal{H}_{\text{lin}}\} = [-W\|x - x'\|_p, W\|x - x'\|_p]$.

K.1.1. DERIVATION FOR $\mathcal{L}_{\Phi_{\text{hinge}}}$.

For the hinge loss function $\Phi_{\text{hinge}}(u) := \max\{0, 1 - u\}$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\text{hinge}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{hinge}}}(h(x) - h(x')) \\ &= \eta(x, x') \max\{0, 1 - h(x') + h(x)\} + (1 - \eta(x, x')) \max\{0, 1 + h(x') - h(x)\}. \end{aligned}$$

Then,

$$\mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}, \mathcal{H}_{\text{lin}}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h, x, x') = 1 - |2\eta(x, x') - 1| \min\{W\|x - x'\|_p, 1\}.$$

The $(\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is

$$\begin{aligned} \mathcal{M}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{lin}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')}[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}, \mathcal{H}_{\text{lin}}}}^*(x, x')] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')}[1 - |2\eta(x, x') - 1| \min\{W\|x - x'\|_p, 1\}]. \end{aligned} \quad (38)$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{lin}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}, \mathcal{H}_{\text{lin}}}}(h, x, x') &\geq \inf_{h \in \overline{\mathcal{H}}_{\text{lin}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}, \mathcal{H}_{\text{lin}}}}^*(x, x') \\ &= \eta(x, x') \max\{0, 1 - 0\} + (1 - \eta(x, x')) \max\{0, 1 + 0\} - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}, \mathcal{H}_{\text{lin}}}}^*(x, x') \\ &= 1 - [1 - |2\eta(x, x') - 1| \min\{W\|x - x'\|_p, 1\}] \\ &= |2\eta(x, x') - 1| \min\{W\|x - x'\|_p, 1\} \\ &\geq |2\eta(x, x') - 1| \min\{W\gamma, 1\} \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}, \mathcal{H}_{\text{lin}}}}(h, x, x') \geq \min\{W\gamma, 1\} (|2\eta(x, x') - 1|) \mathbb{1}_{h \in \overline{\mathcal{H}}_{\text{lin}}(x, x')} = \Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\mathcal{L}_{\Phi_{\text{hinge}}}$, valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \frac{\mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{lin}})}{\min\{W\gamma, 1\}} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \quad (39)$$

K.1.2. DERIVATION FOR \mathcal{L}_{Φ_ρ} .

For the ρ -margin loss function $\Phi_\rho(u) := \min\left\{1, \max\left\{0, 1 - \frac{u}{\rho}\right\}\right\}$, $\rho > 0$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_\rho}}(h, x, x') &= \eta(x, x')\mathcal{L}_{\Phi_\rho}(h(x') - h(x)) + (1 - \eta(x, x'))\mathcal{L}_{\Phi_\rho}(h(x) - h(x')) \\ &= \eta(x, x') \min\left\{1, \max\left\{0, 1 - \frac{h(x') - h(x)}{\rho}\right\}\right\} + (1 - \eta(x, x')) \min\left\{1, \max\left\{0, 1 + \frac{h(x') - h(x)}{\rho}\right\}\right\}. \end{aligned}$$

Then,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{lin}}}^*(x, x') &= \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\mathcal{L}_{\Phi_\rho}}(h, x, x') \\ &= \min\{\eta(x, x'), 1 - \eta(x, x')\} + \max\{\eta(x, x'), 1 - \eta(x, x')\} \left(1 - \frac{\min\{W\|x - x'\|_p, \rho\}}{\rho}\right) \end{aligned}$$

The $(\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{lin}})$ -minimizability gap is

$$\begin{aligned} \mathcal{M}_{\mathcal{L}_{\Phi_\rho}}(\mathcal{H}_{\text{lin}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_\rho}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{lin}}}^*(x, x') \right] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_\rho}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[\min\{\eta(x, x'), 1 - \eta(x, x')\} + \max\{\eta(x, x'), 1 - \eta(x, x')\} \left(1 - \frac{\min\{W\|x - x'\|_p, \rho\}}{\rho}\right) \right]. \end{aligned} \quad (40)$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{lin}}(x, x')$,

$$\begin{aligned} \Delta_{\mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{lin}}}}(h, x, x') &\geq \inf_{h \in \overline{\mathcal{H}}_{\text{lin}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_\rho}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= \max\{\eta(x, x'), 1 - \eta(x, x')\} + \min\{\eta(x, x'), 1 - \eta(x, x')\} \left(1 - \frac{\min\{W\|x - x'\|_p, \rho\}}{\rho}\right) - \mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= |2\eta(x, x') - 1| \frac{\min\{W\|x - x'\|_p, \rho\}}{\rho} \\ &\geq |2\eta(x, x') - 1| \frac{\min\{W\gamma, \rho\}}{\rho} \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta_{\mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{lin}}}}(h, x, x') \geq \frac{\min\{W\gamma, \rho\}}{\rho} (|2\eta(x, x') - 1|)_0 \mathbb{1}_{h \in \overline{\mathcal{H}}_{\text{lin}}(x, x')} = \Delta_{\mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for \mathcal{L}_{Φ_ρ} , valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \frac{\rho \left(\mathcal{R}_{\mathcal{L}_{\Phi_\rho}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_\rho}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_\rho}}(\mathcal{H}_{\text{lin}}) \right)}{\min\{W\gamma, \rho\}} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \quad (41)$$

 K.1.3. DERIVATION FOR $\mathcal{L}_{\Phi_{\text{exp}}}$.

For the exponential loss function $\Phi_{\text{exp}}(u) := e^{-u}$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h, x, x') &= \eta(x, x')\mathcal{L}_{\Phi_{\text{exp}}}(h(x') - h(x)) + (1 - \eta(x, x'))\mathcal{L}_{\Phi_{\text{exp}}}(h(x) - h(x')) \\ &= \eta(x, x')e^{-h(x') + h(x)} + (1 - \eta(x, x'))e^{h(x') - h(x)}. \end{aligned}$$

Then,

$$\begin{aligned}
 & \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 &= \inf_{h \in \overline{\mathcal{H}_{\text{lin}}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h, x, x') \\
 &= \begin{cases} 2\sqrt{\eta(x, x')(1 - \eta(x, x'))} & \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq W \|x - x'\|_p \\ \max\{\eta(x, x'), 1 - \eta(x, x')\} e^{-W \|x - x'\|_p} + \min\{\eta(x, x'), 1 - \eta(x, x')\} e^{W \|x - x'\|_p} & \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > W \|x - x'\|_p. \end{cases}
 \end{aligned}$$

The $(\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is:

$$\begin{aligned}
 \mathcal{M}_{\mathcal{L}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{lin}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \right] \\
 &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[2\sqrt{\eta(x, x')(1 - \eta(x, x'))} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq W \|x - x'\|_p} \right. \\
 &\quad \left. - \mathbb{E}_{(X, X')} \left[\max\{\eta(x, x'), 1 - \eta(x, x')\} e^{-W \|x - x'\|_p} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > W \|x - x'\|_p} \right] \right. \\
 &\quad \left. - \mathbb{E}_{(X, X')} \left[\min\{\eta(x, x'), 1 - \eta(x, x')\} e^{W \|x - x'\|_p} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > W \|x - x'\|_p} \right] \right]. \tag{42}
 \end{aligned}$$

Therefore, $\forall h \in \overline{\mathcal{H}_{\text{lin}}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \\
 & \geq \inf_{h \in \overline{\mathcal{H}_{\text{lin}}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 &= \eta(x, x') e^{-0} + (1 - \eta(x, x')) e^0 - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 &= \begin{cases} 1 - 2\sqrt{\eta(x, x')(1 - \eta(x, x'))} & \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq W \|x - x'\|_p \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\} e^{-W \|x - x'\|_p} - \min\{\eta(x, x'), 1 - \eta(x, x')\} e^{W \|x - x'\|_p} & \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > W \|x - x'\|_p \end{cases} \\
 & \geq \begin{cases} 1 - 2\sqrt{\eta(x, x')(1 - \eta(x, x'))} & \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq W \gamma \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\} e^{-W \gamma} - \min\{\eta(x, x'), 1 - \eta(x, x')\} e^{W \gamma} & \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > W \gamma \end{cases} \\
 &= \Psi_{\text{exp}}(|2\eta(x, x') - 1|),
 \end{aligned}$$

where Ψ_{exp} is the increasing and convex function on $[0, 1]$ defined by

$$\forall t \in [0, 1], \quad \Psi_{\text{exp}}(t) = \begin{cases} 1 - \sqrt{1 - t^2}, & t \leq \frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \\ 1 - \frac{t+1}{2} e^{-W\gamma} - \frac{1-t}{2} e^{W\gamma}, & t > \frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \end{cases}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \geq \Psi_{\text{exp}}\left(\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}}}(h, x, x')\right).$$

To simplify the expression, using the fact that

$$\begin{aligned}
 & 1 - \sqrt{1 - t^2} \geq \frac{t^2}{2}, \\
 & 1 - \frac{t+1}{2} e^{-W\gamma} - \frac{1-t}{2} e^{W\gamma} = 1 - \frac{e^{W\gamma}}{2} - \frac{e^{-W\gamma}}{2} + \frac{e^{W\gamma} - e^{-W\gamma}}{2} t,
 \end{aligned}$$

Ψ_{exp} can be lower bounded by

$$\tilde{\Psi}_{\text{exp}}(t) = \begin{cases} \frac{t^2}{2}, & t \leq \frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \\ \frac{1}{2} \left(\frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \right) t, & t > \frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1}. \end{cases}$$

Thus, we adopt an upper bound of Ψ^{-1} as follows:

$$\begin{aligned}\Gamma_{\text{exp}}(t) &= \tilde{\Psi}_{\text{exp}}^{-1}(t) = \begin{cases} \sqrt{2t}, & t \leq \frac{1}{2} \left(\frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \right)^2 \\ 2 \left(\frac{e^{2W\gamma} + 1}{e^{2W\gamma} - 1} \right) t, & t > \frac{1}{2} \left(\frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \right)^2 \end{cases} \\ &= \max \left\{ \sqrt{2t}, 2 \left(\frac{e^{2W\gamma} + 1}{e^{2W\gamma} - 1} \right) t \right\}.\end{aligned}$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\mathcal{L}_{\Phi_{\text{exp}}}$, valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \Gamma_{\text{exp}} \left(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{lin}}) \right) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \quad (43)$$

where $\Gamma_{\text{exp}}(t) = \max \left\{ \sqrt{2t}, 2 \left(\frac{e^{2W\gamma} + 1}{e^{2W\gamma} - 1} \right) t \right\}$.

K.1.4. DERIVATION FOR $\mathcal{L}_{\Phi_{\log}}$.

For the logistic loss function $\Phi_{\log}(u) := \log_2(1 + e^{-u})$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned}\mathcal{C}_{\mathcal{L}_{\Phi_{\log}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\log}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\log}}(h(x) - h(x')) \\ &= \eta(x, x') \log_2(1 + e^{-h(x') + h(x)}) + (1 - \eta(x, x')) \log_2(1 + e^{h(x') - h(x)}).\end{aligned}$$

Then,

$$\begin{aligned}\mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}^*(x, x') &= \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}}(h, x, x') \\ &= \begin{cases} -\eta(x, x') \log_2(\eta(x, x')) - (1 - \eta(x, x')) \log_2(1 - \eta(x, x')) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq W \|x - x'\|_p \\ \max\{\eta(x, x'), 1 - \eta(x, x')\} \log_2(1 + e^{-W \|x - x'\|_p}) + \min\{\eta(x, x'), 1 - \eta(x, x')\} \log_2(1 + e^{W \|x - x'\|_p}) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > W \|x - x'\|_p. \end{cases}\end{aligned}$$

The $(\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is:

$$\begin{aligned}\mathcal{M}_{\mathcal{L}_{\Phi_{\log}}}(\mathcal{H}_{\text{lin}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}^*(x, x') \right] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[-\eta(x, x') \log_2(\eta(x, x')) - (1 - \eta(x, x')) \log_2(1 - \eta(x, x')) \mathbb{1}_{\left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq W \|x - x'\|_p} \right. \\ &\quad \left. - \mathbb{E}_{(X, X')} \left[\max\{\eta(x, x'), 1 - \eta(x, x')\} \log_2(1 + e^{-W \|x - x'\|_p}) \mathbb{1}_{\left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > W \|x - x'\|_p} \right] \right. \\ &\quad \left. - \mathbb{E}_{(X, X')} \left[\min\{\eta(x, x'), 1 - \eta(x, x')\} \log_2(1 + e^{W \|x - x'\|_p}) \mathbb{1}_{\left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > W \|x - x'\|_p} \right] \right]. \quad (44)\end{aligned}$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{lin}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}(h, x, x') \\
 & \geq \inf_{h \in \overline{\mathcal{H}}_{\text{lin}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 & = \eta(x, x') \log_2(1 + e^{-0}) + (1 - \eta(x, x')) \log_2(1 + e^0) - \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 & = \begin{cases} 1 + \eta(x, x') \log_2(\eta(x, x')) + (1 - \eta(x, x')) \log_2(1 - \eta(x, x')) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq W \|x - x'\|_p \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\} \log_2(1 + e^{-W \|x - x'\|_p}) - \min\{\eta(x, x'), 1 - \eta(x, x')\} \log_2(1 + e^{W \|x - x'\|_p}) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > W \|x - x'\|_p \end{cases} \\
 & \geq \begin{cases} 1 + \eta(x, x') \log_2(\eta(x, x')) + (1 - \eta(x, x')) \log_2(1 - \eta(x, x')) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq W\gamma \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\} \log_2(1 + e^{-W\gamma}) - \min\{\eta(x, x'), 1 - \eta(x, x')\} \log_2(1 + e^{W\gamma}) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > W\gamma \end{cases} \\
 & = \Psi_{\log}(|2\eta(x, x') - 1|)
 \end{aligned}$$

where Ψ_{\log} is the increasing and convex function on $[0, 1]$ defined by

$$\forall t \in [0, 1], \quad \mathcal{J}(t) = \begin{cases} \frac{t+1}{2} \log_2(t+1) + \frac{1-t}{2} \log_2(1-t), & t \leq \frac{e^{W\gamma}-1}{e^{W\gamma}+1} \\ 1 - \frac{t+1}{2} \log_2(1 + e^{-W\gamma}) - \frac{1-t}{2} \log_2(1 + e^{W\gamma}), & t > \frac{e^{W\gamma}-1}{e^{W\gamma}+1} \end{cases}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}(h, x, x') \geq \Psi_{\log}(\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}}}(h, x, x')).$$

To simplify the expression, using the fact that

$$\begin{aligned}
 \frac{t+1}{2} \log_2(t+1) + \frac{1-t}{2} \log_2(1-t) & = 1 - \left(-\frac{t+1}{2} \log_2\left(\frac{t+1}{2}\right) - \frac{1-t}{2} \log_2\left(\frac{1-t}{2}\right) \right) \\
 & \geq 1 - \sqrt{4 \frac{1-t}{2} \frac{t+1}{2}} \\
 & = 1 - \sqrt{1-t^2} \\
 & \geq \frac{t^2}{2}, \\
 1 - \frac{t+1}{2} \log_2(1 + e^{-W\gamma}) - \frac{1-t}{2} \log_2(1 + e^{W\gamma}) & = \frac{1}{2} \log_2\left(\frac{4}{2 + e^{-W\gamma} + e^{W\gamma}}\right) + \frac{1}{2} \log_2\left(\frac{1 + e^{W\gamma}}{1 + e^{-W\gamma}}\right) t,
 \end{aligned}$$

Ψ_{\log} can be lower bounded by

$$\tilde{\Psi}_{\log}(t) = \begin{cases} \frac{t^2}{2}, & t \leq \frac{e^{W\gamma}-1}{e^{W\gamma}+1} \\ \frac{1}{2} \left(\frac{e^{W\gamma}-1}{e^{W\gamma}+1} \right) t, & t > \frac{e^{W\gamma}-1}{e^{W\gamma}+1} \end{cases}$$

Thus, we adopt an upper bound of Ψ_{\log}^{-1} as follows:

$$\begin{aligned}
 \Gamma_{\log}(t) = \tilde{\Psi}_{\log}^{-1}(t) & = \begin{cases} \sqrt{2t}, & t \leq \frac{1}{2} \left(\frac{e^{W\gamma}-1}{e^{W\gamma}+1} \right)^2 \\ 2 \left(\frac{e^{W\gamma}+1}{e^{W\gamma}-1} \right) t, & t > \frac{1}{2} \left(\frac{e^{W\gamma}-1}{e^{W\gamma}+1} \right)^2 \end{cases} \\
 & = \max \left\{ \sqrt{2t}, 2 \left(\frac{e^{W\gamma}+1}{e^{W\gamma}-1} \right) t \right\}.
 \end{aligned}$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\mathcal{L}_{\Phi_{\log}}$, valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\mathcal{L}_{\Phi_{\log}}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \Gamma_{\log} \left(\mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\log}}}(\mathcal{H}_{\text{lin}}) \right) - \mathcal{M}_{\mathcal{L}_{\Phi_{\log}}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \quad (45)$$

where $\Gamma_{\log}(t) = \max\left\{\sqrt{2t}, 2\left(\frac{e^{W\gamma} + 1}{e^{W\gamma} - 1}\right)t\right\}$.

K.1.5. DERIVATION FOR $\mathcal{L}_{\Phi_{\text{sq}}}$.

For the squared hinge loss function $\Phi_{\text{sq}}(u) := (1 - u)^2 \mathbb{1}_{u \leq 1}$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} & \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h, x, x') \\ &= \eta(x, x') \mathcal{L}_{\Phi_{\text{sq}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{sq}}}(h(x) - h(x')) \\ &= \eta(x, x') (1 - h(x') + h(x))^2 \mathbb{1}_{h(x') - h(x) \leq 1} + (1 - \eta(x, x')) (1 + h(x') - h(x))^2 \mathbb{1}_{h(x') - h(x) \geq -1}. \end{aligned}$$

Then,

$$\begin{aligned} & \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h, x, x') \\ &= \begin{cases} 4\eta(x, x')(1 - \eta(x, x')) & \text{if } |2\eta(x, x') - 1| \leq W\|x - x'\|_p \\ \max\{\eta(x, x'), 1 - \eta(x, x')\} (1 - W\|x - x'\|_p)^2 + \min\{\eta(x, x'), 1 - \eta(x, x')\} (1 + W\|x - x'\|_p)^2 & \text{if } |2\eta(x, x') - 1| > W\|x - x'\|_p. \end{cases} \end{aligned}$$

The $(\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is:

$$\begin{aligned} \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{lin}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \right] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[4\eta(x, x')(1 - \eta(x, x')) \mathbb{1}_{|2\eta(x, x') - 1| \leq W\|x - x'\|_p} \right] \\ &\quad - \mathbb{E}_{(X, X')} \left[\max\{\eta(x, x'), 1 - \eta(x, x')\} (1 - W\|x - x'\|_p)^2 \mathbb{1}_{|2\eta(x, x') - 1| > W\|x - x'\|_p} \right] \\ &\quad - \mathbb{E}_{(X, X')} \left[\min\{\eta(x, x'), 1 - \eta(x, x')\} (1 + W\|x - x'\|_p)^2 \mathbb{1}_{|2\eta(x, x') - 1| > W\|x - x'\|_p} \right]. \end{aligned} \quad (46)$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{lin}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \\ & \geq \inf_{h \in \overline{\mathcal{H}}_{\text{lin}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= \eta(x, x') + (1 - \eta(x, x')) - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= \begin{cases} 1 - 4\eta(x, x')(1 - \eta(x, x')) & \text{if } |2\eta(x, x') - 1| \leq W\|x - x'\|_p \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\} (1 - W\|x - x'\|_p)^2 - \min\{\eta(x, x'), 1 - \eta(x, x')\} (1 + W\|x - x'\|_p)^2 & \text{if } |2\eta(x, x') - 1| > W\|x - x'\|_p \end{cases} \\ & \geq \begin{cases} 1 - 4\eta(x, x')(1 - \eta(x, x')) & |2\eta(x, x') - 1| \leq W\gamma \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\} (1 - W\gamma)^2 - \min\{\eta(x, x'), 1 - \eta(x, x')\} (1 + W\gamma)^2 & |2\eta(x, x') - 1| > W\gamma \end{cases} \\ &= \Psi_{\text{sq}}(|2\eta(x, x') - 1|), \end{aligned}$$

where Ψ_{sq} is the increasing and convex function on $[0, 1]$ defined by

$$\forall t \in [0, 1], \quad \Psi_{\text{sq}}(t) = \begin{cases} t^2 & t \leq W\gamma \\ 2W\gamma t - (W\gamma)^2 & t > W\gamma \end{cases}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \geq \Psi_{\text{sq}}(\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}}}(h, x, x')).$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\mathcal{L}_{\Phi_{\text{sq}}}$, valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \Gamma_{\text{sq}}(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{lin}})) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}) \quad (47)$$

where $\Gamma_{\text{sq}}(t) = \Phi_{\text{sq}}^{-1}(t) = \begin{cases} \sqrt{t}, & t \leq (W\gamma)^2 \\ \frac{t}{2W\gamma} + \frac{W\gamma}{2}, & t > (W\gamma)^2 \end{cases} = \max\left\{\sqrt{t}, \frac{t}{2W\gamma} + \frac{W\gamma}{2}\right\}.$

K.1.6. DERIVATION FOR $\mathcal{L}_{\Phi_{\text{sig}}}$.

For the sigmoid loss function $\Phi_{\text{sig}}(u) := 1 - \tanh(ku)$, $k > 0$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\text{sig}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{sig}}}(h(x) - h(x')) \\ &= \eta(x, x')(1 - \tanh(k[h(x') - h(x)])) + (1 - \eta(x, x'))(1 + \tanh(k[h(x') - h(x)])). \end{aligned}$$

Then,

$$\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{lin}})(x, x') = \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h, x, x') = 1 - |1 - 2\eta(x, x')| \tanh(kW\|x - x'\|_p).$$

The $(\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is:

$$\begin{aligned} \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{lin}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')}[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}}}^*(x, x')] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')}[1 - |1 - 2\eta(x, x')| \tanh(kW\|x - x'\|_p)]. \end{aligned} \quad (48)$$

Therefore, $\forall h \in \overline{\mathcal{H}_{\text{lin}}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}}}(h, x, x') &\geq \inf_{h \in \overline{\mathcal{H}_{\text{lin}}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= 1 - |1 - 2\eta(x, x')| \tanh(0) - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= |1 - 2\eta(x, x')| \tanh(kW\|x - x'\|_p) \\ &\geq |1 - 2\eta(x, x')| \tanh(kW\gamma) \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \geq \tanh(kW\gamma) \Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\mathcal{L}_{\Phi_{\text{sig}}}$, valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \frac{\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{lin}})}{\tanh(kW\gamma)} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \quad (49)$$

K.2. One-Hidden-Layer ReLU Neural Networks

Since \mathcal{H}_{NN} satisfies the condition of Lemma K.1, by Lemma K.1 the $(\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{NN}})$ -minimizability gap can be expressed as follows:

$$\mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}) = \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')}[\min\{\eta(x, x'), 1 - \eta(x, x')\} \mathbb{1}_{\|x - x'\| > \gamma} + c \mathbb{1}_{|x - x'| \leq \gamma}]. \quad (50)$$

By the definition of \mathcal{H}_{NN} , for any $(x, x') \in \mathcal{X} \times \mathcal{X}$, $\{h(x') - h(x) \mid h \in \mathcal{H}_{\text{NN}}\} = [-\Lambda W\|x - x'\|_p, \Lambda W\|x - x'\|_p]$.

K.2.1. DERIVATION FOR $\mathcal{L}_{\Phi_{\text{hinge}}}$.

For the hinge loss function $\Phi_{\text{hinge}}(u) := \max\{0, 1 - u\}$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\text{hinge}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{hinge}}}(h(x) - h(x')) \\ &= \eta(x, x') \max\{0, 1 - h(x') + h(x)\} + (1 - \eta(x, x')) \max\{0, 1 + h(x') - h(x)\}. \end{aligned}$$

Then,

$$\mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h, x, x') = 1 - |2\eta(x, x') - 1| \min\{\Lambda W \|x - x'\|_p, 1\}.$$

The $(\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}})$ -minimizability gap is

$$\begin{aligned} \mathcal{M}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{NN}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[1 - |2\eta(x, x') - 1| \min\{\Lambda W \|x - x'\|_p, 1\} \right]. \end{aligned} \quad (51)$$

Therefore, $\forall h \in \overline{\mathcal{H}_{\text{NN}}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}(h, x, x') &\geq \inf_{h \in \overline{\mathcal{H}_{\text{NN}}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= \eta(x, x') \max\{0, 1 - 0\} + (1 - \eta(x, x')) \max\{0, 1 + 0\} - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= 1 - \left[1 - |2\eta(x, x') - 1| \min\{\Lambda W \|x - x'\|_p, 1\} \right] \\ &= |2\eta(x, x') - 1| \min\{\Lambda W \|x - x'\|_p, 1\} \\ &\geq |2\eta(x, x') - 1| \min\{\Lambda W \gamma, 1\} \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \min\{\Lambda W \gamma, 1\} (|2\eta(x, x') - 1|) \mathbb{1}_{h \in \overline{\mathcal{H}_{\text{NN}}}(x, x')} = \Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{NN}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\mathcal{L}_{\Phi_{\text{hinge}}}$, valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \frac{\mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{NN}})}{\min\{\Lambda W \gamma, 1\}} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \quad (52)$$

K.2.2. DERIVATION FOR $\mathcal{L}_{\Phi_{\rho}}$.

For the ρ -margin loss function $\Phi_{\rho}(u) := \min\left\{1, \max\left\{0, 1 - \frac{u}{\rho}\right\}\right\}$, $\rho > 0$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\rho}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\rho}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\rho}}(h(x) - h(x')) \\ &= \eta(x, x') \min\left\{1, \max\left\{0, 1 - \frac{h(x') - h(x)}{\rho}\right\}\right\} + (1 - \eta(x, x')) \min\left\{1, \max\left\{0, 1 + \frac{h(x') - h(x)}{\rho}\right\}\right\}. \end{aligned}$$

Then,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\rho}}, \mathcal{H}_{\text{NN}}}^*(x, x') &= \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\rho}}}(h, x, x') \\ &= \min\{\eta(x, x'), 1 - \eta(x, x')\} + \max\{\eta(x, x'), 1 - \eta(x, x')\} \left(1 - \frac{\min\{\Lambda W \|x - x'\|_p, \rho\}}{\rho}\right) \end{aligned}$$

The $(\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{NN}})$ -minimizability gap is

$$\begin{aligned}
 & \mathcal{M}_{\mathcal{L}_{\Phi_\rho}}(\mathcal{H}_{\text{NN}}) \\
 &= \mathcal{R}_{\mathcal{L}_{\Phi_\rho}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\
 &= \mathcal{R}_{\mathcal{L}_{\Phi_\rho}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\min\{\eta(x, x'), 1 - \eta(x, x')\} + \max\{\eta(x, x'), 1 - \eta(x, x')\} \left(1 - \frac{\min\{\Lambda W \|x - x'\|_p, \rho\}}{\rho} \right) \right]. \tag{53}
 \end{aligned}$$

Therefore, $\forall h \in \overline{\mathcal{H}_{\text{NN}}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{NN}}}(h, x, x') \\
 & \geq \inf_{h \in \overline{\mathcal{H}_{\text{NN}}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_\rho}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 &= \max\{\eta(x, x'), 1 - \eta(x, x')\} + \min\{\eta(x, x'), 1 - \eta(x, x')\} \left(1 - \frac{\min\{\Lambda W \|x - x'\|_p, \rho\}}{\rho} \right) - \mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 &= |2\eta(x, x') - 1| \frac{\min\{\Lambda W \|x - x'\|_p, \rho\}}{\rho} \\
 & \geq |2\eta(x, x') - 1| \frac{\min\{\Lambda W \gamma, \rho\}}{\rho}
 \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_\rho}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \frac{\min\{\Lambda W \gamma, \rho\}}{\rho} (|2\eta(x, x') - 1|)_0 \mathbb{1}_{h \in \overline{\mathcal{H}_{\text{NN}}}(x, x')} = \Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{NN}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for \mathcal{L}_{Φ_ρ} , valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \frac{\rho \left(\mathcal{R}_{\mathcal{L}_{\Phi_\rho}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_\rho}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_\rho}}(\mathcal{H}_{\text{NN}}) \right)}{\min\{\Lambda W \gamma, \rho\}} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \tag{54}$$

K.2.3. DERIVATION FOR $\mathcal{L}_{\Phi_{\text{exp}}}$.

For the exponential loss function $\Phi_{\text{exp}}(u) = e^{-u}$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned}
 \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\text{exp}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{exp}}}(h(x) - h(x')) \\
 &= \eta(x, x') e^{-h(x') + h(x)} + (1 - \eta(x, x')) e^{h(x') - h(x)}.
 \end{aligned}$$

Then,

$$\begin{aligned}
 & \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 &= \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h, x, x') \\
 &= \begin{cases} 2\sqrt{\eta(x, x')(1 - \eta(x, x'))} & \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq \Lambda W \|x - x'\|_p \\ \max\{\eta(x, x'), 1 - \eta(x, x')\} e^{-\Lambda W \|x - x'\|_p} + \min\{\eta(x, x'), 1 - \eta(x, x')\} e^{\Lambda W \|x - x'\|_p} & \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > \Lambda W \|x - x'\|_p \end{cases}
 \end{aligned}$$

The $(\Phi_{\text{exp}}, \mathcal{H}_{\text{NN}})$ -minimizability gap is:

$$\begin{aligned}
 \mathcal{M}_{\mathcal{L}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{NN}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\
 &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[2\sqrt{\eta(x, x')(1 - \eta(x, x'))} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq \Lambda W \|x - x'\|_p} \right. \\
 & \quad \left. - \mathbb{E}_{(X, X')} \left[\max\{\eta(x, x'), 1 - \eta(x, x')\} e^{-\Lambda W \|x - x'\|_p} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > \Lambda W \|x - x'\|_p} \right] \right. \\
 & \quad \left. - \mathbb{E}_{(X, X')} \left[\min\{\eta(x, x'), 1 - \eta(x, x')\} e^{\Lambda W \|x - x'\|_p} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > \Lambda W \|x - x'\|_p} \right] \right] \tag{55}
 \end{aligned}$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{NN}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \\
 & \geq \inf_{h \in \overline{\mathcal{H}}_{\text{NN}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 & = \eta(x, x')e^{-0} + (1 - \eta(x, x'))e^0 - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 & = \begin{cases} 1 - 2\sqrt{\eta(x, x')(1 - \eta(x, x'))} \\ \quad \text{if } \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq \Lambda W \|x - x'\|_p \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\} e^{-\Lambda W \|x - x'\|_p} - \min\{\eta(x, x'), 1 - \eta(x, x')\} e^{\Lambda W \|x - x'\|_p} \\ \quad \text{if } \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > \Lambda W \|x - x'\|_p \end{cases} \\
 & \geq \begin{cases} 1 - 2\sqrt{\eta(x, x')(1 - \eta(x, x'))} & \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq \Lambda W \gamma \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\} e^{-\Lambda W \gamma} - \min\{\eta(x, x'), 1 - \eta(x, x')\} e^{\Lambda W \gamma} & \frac{1}{2} \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > \Lambda W \gamma \end{cases} \\
 & = \Psi_{\text{exp}}(|2\eta(x, x') - 1|),
 \end{aligned}$$

where Ψ_{exp} is the increasing and convex function on $[0, 1]$ defined by

$$\forall t \in [0, 1], \quad \Psi_{\text{exp}}(t) = \begin{cases} 1 - \sqrt{1 - t^2}, & t \leq \frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \\ 1 - \frac{t+1}{2} e^{-\Lambda W \gamma} - \frac{1-t}{2} e^{\Lambda W \gamma}, & t > \frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \end{cases}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \Psi_{\text{exp}}(\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{NN}}}(h, x, x')).$$

To simplify the expression, using the fact that

$$\begin{aligned}
 1 - \sqrt{1 - t^2} & \geq \frac{t^2}{2}, \\
 1 - \frac{t+1}{2} e^{-\Lambda W \gamma} - \frac{1-t}{2} e^{\Lambda W \gamma} & = 1 - \frac{e^{\Lambda W \gamma}}{2} - \frac{e^{-\Lambda W \gamma}}{2} + \frac{e^{\Lambda W \gamma} - e^{-\Lambda W \gamma}}{2} t,
 \end{aligned}$$

Ψ_{exp} can be lower bounded by

$$\tilde{\Psi}_{\text{exp}}(t) = \begin{cases} \frac{t^2}{2}, & t \leq \frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \\ \frac{1}{2} \left(\frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \right) t, & t > \frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1}. \end{cases}$$

Thus, we adopt an upper bound of Ψ^{-1} as follows:

$$\begin{aligned}
 \Gamma_{\text{exp}}(t) & = \tilde{\Psi}_{\text{exp}}^{-1}(t) = \begin{cases} \sqrt{2t}, & t \leq \frac{1}{2} \left(\frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \right)^2 \\ 2 \left(\frac{e^{2W\gamma} + 1}{e^{2W\gamma} - 1} \right) t, & t > \frac{1}{2} \left(\frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \right)^2 \end{cases} \\
 & = \max \left\{ \sqrt{2t}, 2 \left(\frac{e^{2W\gamma} + 1}{e^{2W\gamma} - 1} \right) t \right\}.
 \end{aligned}$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\mathcal{L}_{\Phi_{\text{exp}}}$, valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \Gamma_{\text{exp}} \left(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{NN}}) \right) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \quad (56)$$

where $\Gamma_{\text{exp}}(t) = \max \left\{ \sqrt{2t}, 2 \left(\frac{e^{2W\gamma} + 1}{e^{2W\gamma} - 1} \right) t \right\}$.

K.2.4. DERIVATION FOR $\mathcal{L}_{\Phi_{\log}}$.

For the logistic loss function $\Phi_{\log}(u) := \log_2(1 + e^{-u})$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}}(h, x, x') &= \eta(x, x')\mathcal{L}_{\Phi_{\log}}(h(x') - h(x)) + (1 - \eta(x, x'))\mathcal{L}_{\Phi_{\log}}(h(x) - h(x')) \\ &= \eta(x, x')\log_2\left(1 + e^{-h(x') + h(x)}\right) + (1 - \eta(x, x'))\log_2\left(1 + e^{h(x') - h(x)}\right). \end{aligned}$$

Then,

$$\begin{aligned} &\mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}}(h, x, x') \\ &= \begin{cases} -\eta(x, x')\log_2(\eta(x, x')) - (1 - \eta(x, x'))\log_2(1 - \eta(x, x')) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq \Lambda W \|x - x'\|_p \\ \max\{\eta(x, x'), 1 - \eta(x, x')\}\log_2\left(1 + e^{-\Lambda W \|x - x'\|_p}\right) + \min\{\eta(x, x'), 1 - \eta(x, x')\}\log_2\left(1 + e^{\Lambda W \|x - x'\|_p}\right) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > \Lambda W \|x - x'\|_p \end{cases} \end{aligned}$$

The $(\Phi_{\log}, \mathcal{H}_{\text{NN}})$ -minimizability gap is:

$$\begin{aligned} &\mathcal{M}_{\mathcal{L}_{\Phi_{\log}}}(\mathcal{H}_{\text{NN}}) \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[-\eta(x, x')\log_2(\eta(x, x')) - (1 - \eta(x, x'))\log_2(1 - \eta(x, x')) \mathbb{1}_{\left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq \Lambda W \|x - x'\|_p} \right] \quad (57) \\ &\quad - \mathbb{E}_{(X, X')} \left[\max\{\eta(x, x'), 1 - \eta(x, x')\}\log_2\left(1 + e^{-\Lambda W \|x - x'\|_p}\right) \mathbb{1}_{\left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > \Lambda W \|x - x'\|_p} \right] \\ &\quad - \mathbb{E}_{(X, X')} \left[\min\{\eta(x, x'), 1 - \eta(x, x')\}\log_2\left(1 + e^{\Lambda W \|x - x'\|_p}\right) \mathbb{1}_{\left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > \Lambda W \|x - x'\|_p} \right]. \end{aligned}$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{NN}}(x, x')$,

$$\begin{aligned} &\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}(h, x, x') \\ &\geq \inf_{h \in \overline{\mathcal{H}}_{\text{NN}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= \eta(x, x')\log_2(1 + e^{-0}) + (1 - \eta(x, x'))\log_2(1 + e^0) - \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= \begin{cases} 1 + \eta(x, x')\log_2(\eta(x, x')) + (1 - \eta(x, x'))\log_2(1 - \eta(x, x')) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq \Lambda W \|x - x'\|_p \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\}\log_2\left(1 + e^{-\Lambda W \|x - x'\|_p}\right) - \min\{\eta(x, x'), 1 - \eta(x, x')\}\log_2\left(1 + e^{\Lambda W \|x - x'\|_p}\right) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > \Lambda W \|x - x'\|_p \end{cases} \\ &\geq \begin{cases} 1 + \eta(x, x')\log_2(\eta(x, x')) + (1 - \eta(x, x'))\log_2(1 - \eta(x, x')) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| \leq \Lambda W \gamma \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\}\log_2(1 + e^{-\Lambda W \gamma}) - \min\{\eta(x, x'), 1 - \eta(x, x')\}\log_2(1 + e^{\Lambda W \gamma}) \\ \quad \text{if } \left| \log \frac{\eta(x, x')}{1 - \eta(x, x')} \right| > \Lambda W \gamma \end{cases} \\ &= \Psi_{\log}(|2\eta(x, x') - 1|) \end{aligned}$$

where Ψ_{\log} is the increasing and convex function on $[0, 1]$ defined by

$$\forall t \in [0, 1], \quad \mathcal{J}(t) = \begin{cases} \frac{t+1}{2} \log_2(t+1) + \frac{1-t}{2} \log_2(1-t), & t \leq \frac{e^{\Lambda W \gamma} - 1}{e^{\Lambda W \gamma} + 1} \\ 1 - \frac{t+1}{2} \log_2(1 + e^{-\Lambda W \gamma}) - \frac{1-t}{2} \log_2(1 + e^{\Lambda W \gamma}), & t > \frac{e^{\Lambda W \gamma} - 1}{e^{\Lambda W \gamma} + 1} \end{cases}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \Psi_{\log} \left(\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{NN}}}(h, x, x') \right).$$

To simplify the expression, using the fact that

$$\begin{aligned} \frac{t+1}{2} \log_2(t+1) + \frac{1-t}{2} \log_2(1-t) &= 1 - \left(-\frac{t+1}{2} \log_2\left(\frac{t+1}{2}\right) - \frac{1-t}{2} \log_2\left(\frac{1-t}{2}\right) \right) \\ &\geq 1 - \sqrt{4 \frac{1-t}{2} \frac{t+1}{2}} \\ &= 1 - \sqrt{1-t^2} \\ &\geq \frac{t^2}{2}, \\ 1 - \frac{t+1}{2} \log_2(1+e^{-\Lambda W \gamma}) - \frac{1-t}{2} \log_2(1+e^{\Lambda W \gamma}) &= \frac{1}{2} \log_2\left(\frac{4}{2+e^{-\Lambda W \gamma}+e^{\Lambda W \gamma}}\right) + \frac{1}{2} \log_2\left(\frac{1+e^{\Lambda W \gamma}}{1+e^{-\Lambda W \gamma}}\right) t, \end{aligned}$$

Ψ_{\log} can be lower bounded by

$$\tilde{\Psi}_{\log}(t) = \begin{cases} \frac{t^2}{2}, & t \leq \frac{e^{\Lambda W \gamma} - 1}{e^{\Lambda W \gamma} + 1} \\ \frac{1}{2} \left(\frac{e^{\Lambda W \gamma} - 1}{e^{\Lambda W \gamma} + 1} \right) t, & t > \frac{e^{\Lambda W \gamma} - 1}{e^{\Lambda W \gamma} + 1} \end{cases}$$

Thus, we adopt an upper bound of Ψ_{\log}^{-1} as follows:

$$\begin{aligned} \Gamma_{\log}(t) = \tilde{\Psi}_{\log}^{-1}(t) &= \begin{cases} \sqrt{2t}, & t \leq \frac{1}{2} \left(\frac{e^{\Lambda W \gamma} - 1}{e^{\Lambda W \gamma} + 1} \right)^2 \\ 2 \left(\frac{e^{\Lambda W \gamma} + 1}{e^{\Lambda W \gamma} - 1} \right) t, & t > \frac{1}{2} \left(\frac{e^{\Lambda W \gamma} - 1}{e^{\Lambda W \gamma} + 1} \right)^2 \end{cases} \\ &= \max \left\{ \sqrt{2t}, 2 \left(\frac{e^{\Lambda W \gamma} + 1}{e^{\Lambda W \gamma} - 1} \right) t \right\}. \end{aligned}$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\mathcal{L}_{\Phi_{\log}}$, valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \Gamma_{\log} \left(\mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\log}}}(\mathcal{H}_{\text{NN}}) \right) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \quad (58)$$

where $\Gamma_{\log}(t) = \max \left\{ \sqrt{2t}, 2 \left(\frac{e^{\Lambda W \gamma} + 1}{e^{\Lambda W \gamma} - 1} \right) t \right\}$.

K.2.5. DERIVATION FOR $\mathcal{L}_{\Phi_{\text{sq}}}$.

For the squared hinge loss function $\Phi_{\text{sq}}(u) := (1-u)^2 \mathbb{1}_{u \leq 1}$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h, x, x') &= \eta(x, x') \mathcal{L}_{\Phi_{\text{sq}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{sq}}}(h(x) - h(x')) \\ &= \eta(x, x') (1 - h(x') + h(x))^2 \mathbb{1}_{h(x') - h(x) \leq 1} + (1 - \eta(x, x')) (1 + h(x') - h(x))^2 \mathbb{1}_{h(x') - h(x) \geq -1}. \end{aligned}$$

Then,

$$\begin{aligned} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}^*(x, x') &= \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h, x, x') \\ &= \begin{cases} 4\eta(x, x')(1 - \eta(x, x')) & \text{if } |2\eta(x, x') - 1| \leq \Lambda W \|x - x'\|_p \\ \max\{\eta(x, x'), 1 - \eta(x, x')\} (1 - \Lambda W \|x - x'\|_p)^2 + \min\{\eta(x, x'), 1 - \eta(x, x')\} (1 + \Lambda W \|x - x'\|_p)^2 & \text{if } |2\eta(x, x') - 1| > \Lambda W \|x - x'\|_p. \end{cases} \end{aligned}$$

The $(\Phi_{\text{sq}}, \mathcal{H}_{\text{NN}})$ -minimizability gap is:

$$\begin{aligned}
 \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{NN}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\
 &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[4\eta(x, x')(1 - \eta(x, x')) \mathbb{1}_{|2\eta(x, x') - 1| \leq \Lambda W \|x - x'\|_p} \right] \\
 &\quad - \mathbb{E}_{(X, X')} \left[\max\{\eta(x, x'), 1 - \eta(x, x')\} (1 - \Lambda W \|x - x'\|_p)^2 \mathbb{1}_{|2\eta(x, x') - 1| > \Lambda W \|x - x'\|_p} \right] \\
 &\quad - \mathbb{E}_{(X, X')} \left[\min\{\eta(x, x'), 1 - \eta(x, x')\} (1 + \Lambda W \|x - x'\|_p)^2 \mathbb{1}_{|2\eta(x, x') - 1| > \Lambda W \|x - x'\|_p} \right].
 \end{aligned} \tag{59}$$

Therefore, $\forall h \in \overline{\mathcal{H}}_{\text{NN}}(x, x')$,

$$\begin{aligned}
 &\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \\
 &\geq \inf_{h \in \overline{\mathcal{H}}_{\text{NN}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 &= \eta(x, x') + (1 - \eta(x, x')) - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 &= \begin{cases} 1 - 4\eta(x, x')(1 - \eta(x, x')) & \text{if } |2\eta(x, x') - 1| \leq \Lambda W \|x - x'\|_p \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\} (1 - \Lambda W \|x - x'\|_p)^2 - \min\{\eta(x, x'), 1 - \eta(x, x')\} (1 + \Lambda W \|x - x'\|_p)^2 & \text{if } |2\eta(x, x') - 1| > \Lambda W \|x - x'\|_p \end{cases} \\
 &\geq \begin{cases} 1 - 4\eta(x, x')(1 - \eta(x, x')) & |2\eta(x, x') - 1| \leq \Lambda W \gamma \\ 1 - \max\{\eta(x, x'), 1 - \eta(x, x')\} (1 - \Lambda W \gamma)^2 - \min\{\eta(x, x'), 1 - \eta(x, x')\} (1 + \Lambda W \gamma)^2 & |2\eta(x, x') - 1| > \Lambda W \gamma \end{cases} \\
 &= \Psi_{\text{sq}}(|2\eta(x, x') - 1|),
 \end{aligned}$$

where Ψ_{sq} is the increasing and convex function on $[0, 1]$ defined by

$$\forall t \in [0, 1], \quad \Psi_{\text{sq}}(t) = \begin{cases} t^2 & t \leq \Lambda W \gamma \\ 2\Lambda W \gamma t - (\Lambda W \gamma)^2 & t > \Lambda W \gamma \end{cases}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \Psi_{\text{sq}}(\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{NN}}}(h, x, x')).$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\mathcal{L}_{\Phi_{\text{sq}}}$, valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \Gamma_{\text{sq}} \left(\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{NN}}) \right) - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}) \tag{60}$$

where $\Gamma_{\text{sq}}(t) = \Phi_{\text{sq}}^{-1}(t) = \begin{cases} \sqrt{t}, & t \leq (\Lambda W \gamma)^2 \\ \frac{t}{2\Lambda W \gamma} + \frac{\Lambda W \gamma}{2}, & t > (\Lambda W \gamma)^2 \end{cases} = \max\left\{ \sqrt{t}, \frac{t}{2\Lambda W \gamma} + \frac{\Lambda W \gamma}{2} \right\}.$

K.2.6. DERIVATION FOR $\mathcal{L}_{\Phi_{\text{sig}}}$.

For the sigmoid loss function $\Phi_{\text{sig}}(u) := 1 - \tanh(ku)$, $k > 0$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned}
 &\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h, x, x') \\
 &= \eta(x, x') \mathcal{L}_{\Phi_{\text{sig}}}(h(x') - h(x)) + (1 - \eta(x, x')) \mathcal{L}_{\Phi_{\text{sig}}}(h(x) - h(x')) \\
 &= \eta(x, x') (1 - \tanh(k[h(x') - h(x)])) + (1 - \eta(x, x')) (1 + \tanh(k[h(x) - h(x')])).
 \end{aligned}$$

Then,

$$\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h, x, x') = 1 - |1 - 2\eta(x, x')| \tanh(k\Lambda W \|x - x'\|_p).$$

. The $(\Phi_{\text{sig}}, \mathcal{H}_{\text{NN}})$ -minimizability gap is:

$$\begin{aligned} \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{NN}}) &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\ &= \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[1 - |1 - 2\eta(x, x')| \tanh(k\Lambda W \|x - x'\|_p) \right]. \end{aligned} \quad (61)$$

Therefore, $\forall h \in \overline{\mathcal{H}_{\text{NN}}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}(h, x, x') &\geq \inf_{h \in \overline{\mathcal{H}_{\text{NN}}}(x, x')} \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h, x, x') - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= 1 - |1 - 2\eta(x, x')| \tanh(0) - \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= |1 - 2\eta(x, x')| \tanh(k\Lambda W \|x - x'\|_p) \\ &\geq |1 - 2\eta(x, x')| \tanh(k\Lambda W \gamma) \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\mathcal{L}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \tanh(k\Lambda W \gamma) \Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}})(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\mathcal{L}_{\Phi_{\text{sig}}}$, valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \frac{\mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\mathcal{L}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\mathcal{L}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{NN}})}{\tanh(k\Lambda W \gamma)} - \mathcal{M}_{\mathcal{L}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \quad (62)$$

L. \mathcal{H} - consistency bounds for bipartite abstention losses

We first characterize the minimal conditional $\widetilde{\mathcal{L}}_{0-1}^{\text{abs}}$ -risk and the calibration gap of $\widetilde{\mathcal{L}}_{0-1}^{\text{abs}}$ for a broad class of hypothesis sets. We let $\widetilde{\mathcal{H}}(x, x') = \{h \in \mathcal{H}: (h(x) - h(x'))(\eta(x) - \eta(x')) < 0\}$ and $\mathring{\mathcal{H}}(x, x') = \{h \in \mathcal{H}: h(x) = h(x')\}$ for convenience.

Lemma L.1. *Assume that \mathcal{H} is regular for bipartite ranking. Then, the minimal conditional $\widetilde{\mathcal{L}}_{0-1}^{\text{abs}}$ -risk is*

$$\mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}, x, x') = \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \mathbb{1}_{\|x - x'\| > \gamma} + c \mathbb{1}_{\|x - x'\| \leq \gamma}.$$

The calibration gap of $\widetilde{\mathcal{L}}_{0-1}^{\text{abs}}$ can be characterized as

$$\Delta \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x') = |\eta(x) - \eta(x')| \mathbb{1}_{h \in \widetilde{\mathcal{H}}(x, x')} \mathbb{1}_{\|x - x'\| > \gamma} + \frac{1}{2} |\eta(x) - \eta(x')| \mathbb{1}_{h \in \mathring{\mathcal{H}}(x, x')} \mathbb{1}_{\|x - x'\| > \gamma}.$$

Proof. By the definition, the conditional $\widetilde{\mathcal{L}}_{0-1}^{\text{abs}}$ -risk is

$$\begin{aligned} &\mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}(h, x, x') \\ &= \left(\eta(x)(1 - \eta(x')) \left[\mathbb{1}_{h(x) - h(x') < 0} + \frac{1}{2} \mathbb{1}_{h(x) = h(x')} \right] + \eta(x')(1 - \eta(x)) \left[\mathbb{1}_{h(x) - h(x') > 0} + \frac{1}{2} \mathbb{1}_{h(x) = h(x')} \right] \right) \mathbb{1}_{\|x - x'\| > \gamma} + c \mathbb{1}_{\|x - x'\| \leq \gamma}. \end{aligned}$$

For any (x, x') such that $\|x - x'\| \leq \gamma$ and $h \in \mathcal{H}$, $\mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}(h, x, x) = \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}, x, x) = c$. For any (x, x') such that $\|x - x'\| > \gamma$, by the assumption, there exists $h^* \in \mathcal{H}$ such that

$$(h^*(x) - h^*(x'))(\eta(x) - \eta(x')) \mathbb{1}_{\eta(x) \neq \eta(x')} > 0.$$

Therefore, the optimal conditional $\widetilde{\mathcal{L}}_{0-1}^{\text{abs}}$ -risk can be characterized as for any $x, x' \in \mathcal{X}$,

$$\mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}^*(\mathcal{H}, x, x') = \mathcal{C}_{\mathcal{L}_{0-1}^{\text{abs}}}(h^*, x, x') = \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \mathbb{1}_{\|x - x'\| > \gamma} + c \mathbb{1}_{\|x - x'\| \leq \gamma}.$$

which proves the first part of lemma. By the definition, for any (x, x') such that $\|x - x'\| \leq \gamma$ and $h \in \mathcal{H}$, $\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x') = \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}, x, x') = 0$. For any (x, x') such that $\|x - x'\| > \gamma$ and $h \in \mathcal{H}$,

$$\begin{aligned} \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x') &= \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}, x, x') \\ &= \eta(x)(1 - \eta(x')) \left[\mathbb{1}_{h(x) - h(x') < 0} + \frac{1}{2} \mathbb{1}_{h(x) = h(x')} \right] \\ &\quad + \eta(x')(1 - \eta(x)) \left[\mathbb{1}_{h(x) - h(x') > 0} + \frac{1}{2} \mathbb{1}_{h(x) = h(x')} \right] \\ &\quad - \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \\ &= \begin{cases} |\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))|, & h \in \tilde{\mathcal{H}}(x, x'), \\ \frac{1}{2} |\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))|, & h \in \mathring{\mathcal{H}}(x, x'), \\ 0, & \text{otherwise.} \end{cases} \\ &= \begin{cases} |\eta(x) - \eta(x')|, & h \in \tilde{\mathcal{H}}(x, x'), \\ \frac{1}{2} |\eta(x) - \eta(x')|, & h \in \mathring{\mathcal{H}}(x, x'), \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

This leads to

$$\left\langle \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x') \right\rangle_{\epsilon} = \langle |\eta(x) - \eta(x')| \rangle_{\epsilon} \mathbb{1}_{h \in \tilde{\mathcal{H}}(x, x')} \mathbb{1}_{\|x - x'\| > \gamma} + \left\langle \frac{1}{2} |\eta(x) - \eta(x')| \right\rangle_{\epsilon} \mathbb{1}_{h \in \mathring{\mathcal{H}}(x, x')} \mathbb{1}_{\|x - x'\| > \gamma}.$$

□

L.1. Linear Hypotheses

Since \mathcal{H}_{lin} satisfies the condition of Lemma L.1, by Lemma L.1 the $(\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{lin}})$ -minimizability gap can be expressed as follows:

$$\mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}) = \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} [\min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \mathbb{1}_{\|x - x'\| > \gamma} + c \mathbb{1}_{\|x - x'\| \leq \gamma}]. \quad (63)$$

By the definition of \mathcal{H}_{lin} , for any $(x, x') \in \mathcal{X} \times \mathcal{X}$, $\{h(x') - h(x) \mid h \in \mathcal{H}_{\text{lin}}\} = [-W\|x - x'\|_p, W\|x - x'\|_p]$.

L.1.1. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}$.

For the hinge loss function $\Phi_{\text{hinge}}(u) := \max\{0, 1 - u\}$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h, x, x') &= \eta(x)(1 - \eta(x')) \Phi_{\text{hinge}}(h(x) - h(x')) + \eta(x')(1 - \eta(x)) \Phi_{\text{hinge}}(h(x') - h(x)) \\ &= \eta(x)(1 - \eta(x')) \max\{0, 1 - h(x) + h(x')\} + \eta(x')(1 - \eta(x)) \max\{0, 1 + h(x) - h(x')\}. \end{aligned}$$

Then,

$$\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{lin}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h, x, x') = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - |\eta(x) - \eta(x')| \min\{W\|x - x'\|_p, 1\}.$$

The $(\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is

$$\begin{aligned} \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{lin}}) &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} [\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{lin}}}^*(x, x')] \\ &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} [\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - |\eta(x) - \eta(x')| \min\{W\|x - x'\|_p, 1\}]. \end{aligned} \quad (64)$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \hat{\mathcal{H}}_{\text{lin}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}, \mathcal{H}_{\text{lin}}}}(h, x, x') \\
 & \geq \inf_{h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \hat{\mathcal{H}}_{\text{lin}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}, \mathcal{H}_{\text{lin}}}}^*(x, x') \\
 & = \eta(x)(1 - \eta(x')) \max\{0, 1 - 0\} + \eta(x')(1 - \eta(x)) \max\{0, 1 + 0\} - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}, \mathcal{H}_{\text{lin}}}}^*(x, x') \\
 & = |\eta(x) - \eta(x')| \min\{W\|x - x'\|_p, 1\} \\
 & \geq |\eta(x) - \eta(x')| \min\{W\gamma, 1\}
 \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}, \mathcal{H}_{\text{lin}}}}(h, x, x') \geq \min\{W\gamma, 1\} \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{C}}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}$, valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \frac{\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{lin}})}{\min\{W\gamma, 1\}} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \quad (65)$$

L.1.2. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\rho}}$.

For the ρ -margin loss function $\Phi_{\rho}(u) := \min\left\{1, \max\left\{0, 1 - \frac{u}{\rho}\right\}\right\}$, $\rho > 0$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned}
 & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}(h, x, x') \\
 & = \eta(x)(1 - \eta(x'))\Phi_{\rho}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\rho}(h(x') - h(x)) \\
 & = \eta(x)(1 - \eta(x')) \min\left\{1, \max\left\{0, 1 - \frac{h(x) - h(x')}{\rho}\right\}\right\} \\
 & \quad + \eta(x')(1 - \eta(x)) \min\left\{1, \max\left\{0, 1 + \frac{h(x) - h(x')}{\rho}\right\}\right\}
 \end{aligned}$$

Then,

$$\begin{aligned}
 & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 & = \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}(h, x, x') \\
 & = [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] + [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \left(1 - \frac{\min\{W\|x - x'\|_p, \rho\}}{\rho}\right).
 \end{aligned}$$

The $(\tilde{\mathcal{L}}_{\Phi_{\rho}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is

$$\begin{aligned}
 & \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}(\mathcal{H}_{\text{lin}}) \\
 & = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}, \mathcal{H}_{\text{lin}}}^*(x, x') \right] \\
 & = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[[\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \right. \\
 & \quad \left. + [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \left(1 - \frac{\min\{W\|x - x'\|_p, \rho\}}{\rho}\right) \right]. \quad (66)
 \end{aligned}$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \mathring{\mathcal{H}}_{\text{lin}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{lin}}}(h, x, x') \\
 & \geq \inf_{\tilde{\mathcal{H}}_{\text{lin}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 & = [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \left(1 - \frac{\min\{W\|x - x'\|_p, \rho\}}{\rho}\right) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 & = |\eta(x) - \eta(x')| \frac{\min\{W\|x - x'\|_p, \rho\}}{\rho} \\
 & \geq |\eta(x) - \eta(x')| \frac{\min\{W\gamma, \rho\}}{\rho}
 \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{lin}}}(h, x, x') \geq \frac{\min\{W\gamma, \rho\}}{\rho} \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}^{\text{abs}}, \mathcal{H}_{\text{lin}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_\rho}$, valid for $\text{lin } h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \frac{\rho \left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(\mathcal{H}_{\text{lin}}) \right)}{\min\{W\gamma, \rho\}} - \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_\rho}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \quad (67)$$

L.1.3. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}$.

For the exponential loss function $\Phi_{\text{exp}}(u) = e^{-u}$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned}
 & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h, x, x') \\
 & = \eta(x)(1 - \eta(x'))\Phi_{\text{exp}}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\text{exp}}(h(x') - h(x)) \\
 & = \eta(x)(1 - \eta(x'))e^{-h(x)+h(x')} + \eta(x')(1 - \eta(x))e^{h(x)-h(x')}.
 \end{aligned}$$

Then,

$$\begin{aligned}
 & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 & = \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h, x, x') \\
 & = \begin{cases} 2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))} \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| \leq W\|x - x'\|_p \\ \max\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} e^{-W\|x - x'\|_p} + \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} e^{W\|x - x'\|_p} \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| > W\|x - x'\|_p. \end{cases}
 \end{aligned}$$

The $(\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is:

$$\begin{aligned}
 \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{lin}}) & = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \right] \\
 & = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| \leq W\|x - x'\|_p} \right. \\
 & \quad \left. - \mathbb{E}_{(X, X')} \left[[\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] e^{-W\|x - x'\|_p} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| > W\|x - x'\|_p} \right] \right. \\
 & \quad \left. - \mathbb{E}_{(X, X')} \left[[\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] e^{W\|x - x'\|_p} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| > W\|x - x'\|_p} \right] \right]. \quad (68)
 \end{aligned}$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \hat{\mathcal{H}}_{\text{lin}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \\
 & \geq \inf_{h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \hat{\mathcal{H}}_{\text{lin}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 & = \eta(x)(1 - \eta(x'))e^{-0} + \eta(x')(1 - \eta(x))e^0 - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 & = \begin{cases} \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))} \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| \leq W \|x - x'\|_p \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - e^{-W \|x - x'\|_p}) + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - e^{W \|x - x'\|_p}) \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| > W \|x - x'\|_p \end{cases} \\
 & \geq \begin{cases} \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))} \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| \leq W\gamma \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - e^{-W\gamma}) + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - e^{W\gamma}) \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| > W\gamma \end{cases} \\
 & = \begin{cases} \left(\frac{\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))}{\sqrt{\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x))}} \right)^2 & \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| \leq W\gamma \\ \frac{\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x))}{2} (2 - e^{-W\gamma} - e^{W\gamma}) + \frac{1}{2} |\eta(x) - \eta(x')| (e^{W\gamma} - e^{-W\gamma}) & \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| > W\gamma \end{cases} \\
 & \geq \min \left\{ (\eta(x) - \eta(x'))^2, \left(\frac{e^{2W\gamma} + 1}{e^{2W\gamma} - 1} \right) |\eta(x) - \eta(x')| \right\}
 \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \geq \Psi_{\text{exp}} \left(\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x') \right).$$

where Ψ_{exp} is the increasing function on $[0, 2]$ defined by

$$\forall t \in [0, 1], \quad \Psi_{\text{exp}}(t) = \min \left\{ t^2, \left(\frac{e^{2W\gamma} + 1}{e^{2W\gamma} - 1} \right) t \right\}.$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}$, valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \Gamma_{\text{exp}} \left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}^*}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{lin}}) \right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \quad (69)$$

where $\Gamma_{\text{exp}}(t) = \max \left\{ \sqrt{t}, \left(\frac{e^{2W\gamma} - 1}{e^{2W\gamma} + 1} \right) t \right\}$.

L.1.4. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\log}}$.

For the logistic loss function $\Phi_{\log}(u) := \log_2(1 + e^{-u})$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned}
 & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(h, x, x') \\
 & = \eta(x)(1 - \eta(x'))\Phi_{\log}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\log}(h(x') - h(x)) \\
 & = \eta(x)(1 - \eta(x')) \log_2 \left(1 + e^{-h(x) + h(x')} \right) + \eta(x')(1 - \eta(x)) \log_2 \left(1 + e^{h(x) - h(x')} \right).
 \end{aligned}$$

Then,

$$\begin{aligned}
 & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 &= \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}^*(h, x, x') \\
 &= \begin{cases} -\eta(x)(1-\eta(x')) \log_2(\eta(x)(1-\eta(x'))) - \eta(x')(1-\eta(x)) \log_2(\eta(x')(1-\eta(x))) \\ \text{if } \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| \leq W \|x-x'\|_p \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \log_2(1 + e^{-W\|x-x'\|_p}) + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \log_2(1 + e^{W\|x-x'\|_p}) \\ \text{if } \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > W \|x-x'\|_p \end{cases}
 \end{aligned}$$

The $(\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is

$$\begin{aligned}
 & \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(\mathcal{H}_{\text{lin}}) \\
 &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}^*(x, x') \right] \\
 &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[-\eta(x)(1-\eta(x')) \log_2(\eta(x)(1-\eta(x'))) \right. \\
 &\quad \left. - \eta(x')(1-\eta(x)) \log_2(\eta(x')(1-\eta(x))) \mathbb{1}_{\left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| \leq W \|x-x'\|_p} \right] \\
 &\quad - \mathbb{E}_{(X, X')} \left[[\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \log_2(1 + e^{-W\|x-x'\|_p}) \mathbb{1}_{\left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > W \|x-x'\|_p} \right] \\
 &\quad - \mathbb{E}_{(X, X')} \left[[\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \log_2(1 + e^{W\|x-x'\|_p}) \mathbb{1}_{\left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > W \|x-x'\|_p} \right].
 \end{aligned} \tag{70}$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \hat{\mathcal{H}}_{\text{lin}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}^*(h, x, x') \\
 &\geq \inf_{h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \hat{\mathcal{H}}_{\text{lin}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}^*(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 &= \eta(x)(1-\eta(x')) \log_2(1 + e^{-0}) + \eta(x')(1-\eta(x)) \log_2(1 + e^0) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\
 &\geq \begin{cases} \eta(x)(1-\eta(x')) [1 - \log_2(\eta(x)(1-\eta(x')))] + \eta(x')(1-\eta(x)) [1 - \log_2(\eta(x')(1-\eta(x)))] \\ \text{if } \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| \leq W\gamma \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - \log_2(1 + e^{-W\gamma})) + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - \log_2(1 + e^{W\gamma})) \\ \text{if } \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > W\gamma \end{cases} \\
 &\geq \min \left\{ (\eta(x) - \eta(x'))^2, \left(\frac{e^{W\gamma} + 1}{e^{W\gamma} - 1} \right) |\eta(x) - \eta(x')| \right\}
 \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x-x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{lin}}}^*(h, x, x') \geq \Psi_{\log} \left(\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}}^*(h, x, x') \right).$$

where Ψ_{\log} is the increasing function on $[0, 2]$ defined by

$$\forall t \in [0, 1], \quad \Psi_{\log}(t) = \min \left\{ t^2, \left(\frac{e^{W\gamma} + 1}{e^{W\gamma} - 1} \right) t \right\}.$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\log}}$, valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \Gamma_{\log} \left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(\mathcal{H}_{\text{lin}}) \right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \tag{71}$$

where $\Gamma_{\log}(t) = \max \left\{ \sqrt{t}, \left(\frac{e^{W\gamma} - 1}{e^{W\gamma} + 1} \right) t \right\}$.

L.1.5. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}$.

For the squared hinge loss function $\Phi_{\text{sq}}(u) := (1-u)^2 \mathbb{1}_{u \leq 1}$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h, x, x') \\ &= \eta(x)(1 - \eta(x'))\Phi_{\text{sq}}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\text{sq}}(h(x') - h(x)) \\ &= \eta(x)(1 - \eta(x'))(1 - h(x) + h(x'))^2 \mathbb{1}_{h(x) - h(x') \leq 1} + \eta(x')(1 - \eta(x))(1 + h(x) - h(x'))^2 \mathbb{1}_{h(x) - h(x') \geq -1}. \end{aligned}$$

Then,

$$\begin{aligned} & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h, x, x') \\ &= \begin{cases} 4 \frac{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x'))}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \\ \text{if } \frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \leq W\|x-x'\|_p \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')](1 - W\|x-x'\|_p)^2 + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')](1 + W\|x-x'\|_p)^2 \\ \text{if } \frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} > W\|x-x'\|_p. \end{cases} \end{aligned}$$

The $(\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is

$$\begin{aligned} & \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{lin}}) \\ &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \right] \\ &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} \left[4 \frac{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x'))}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \mathbb{1}_{\frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \leq W\|x-x'\|_p} \right] \\ & \quad - \mathbb{E}_{(X, X')} \left[[\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')](1 - W\|x-x'\|_p)^2 \mathbb{1}_{\frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} > W\|x-x'\|_p} \right] \\ & \quad - \mathbb{E}_{(X, X')} \left[[\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')](1 + W\|x-x'\|_p)^2 \mathbb{1}_{\frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} > W\|x-x'\|_p} \right]. \end{aligned} \tag{72}$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \mathring{\mathcal{H}}_{\text{lin}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \\ & \geq \inf_{h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \mathring{\mathcal{H}}_{\text{lin}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ & \geq \begin{cases} \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 4 \frac{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x'))}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \\ \text{if } \frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \leq W\gamma \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')][1 - (1 - W\gamma)^2] + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')][1 - (1 + W\gamma)^2] \\ \text{if } \frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} > W\gamma. \end{cases} \\ & \geq \min\left\{(\eta(x) - \eta(x'))^2, 2W\gamma|\eta(x) - \eta(x')| - (W\gamma)^2\right\} \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \geq \Psi_{\text{sq}}\left(\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x')\right).$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}$, valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \Gamma_{\text{sq}}\left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{lin}})\right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \tag{73}$$

where $\Gamma_{\text{sq}} = \max\left\{\sqrt{t}, \frac{t}{2W\gamma} + \frac{W\gamma}{2}\right\}$.

L.1.6. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}$.

For the sigmoid loss function $\Phi_{\text{sig}}(u) := 1 - \tanh(ku)$, $k > 0$, for all $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h, x, x') &= \eta(x)(1 - \eta(x'))\Phi_{\text{sig}}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\text{sig}}(h(x') - h(x)) \\ &= \eta(x)(1 - \eta(x'))(1 - \tanh(k[h(x) - h(x')])) + \eta(x')(1 - \eta(x))(1 + \tanh(k[h(x) - h(x')])) \end{aligned}$$

Then,

$$\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{lin}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h, x, x') = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - |\eta(x) - \eta(x')| \tanh(kW\|x - x'\|_p).$$

The $(\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}})$ -minimizability gap is

$$\mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{lin}}) = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{lin}}) - \mathbb{E}_{(X, X')} [\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - |\eta(x) - \eta(x')| \tanh(kW\|x - x'\|_p)]. \quad (74)$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \mathring{\mathcal{H}}_{\text{lin}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}}}(h, x, x') &\geq \inf_{h \in \tilde{\mathcal{H}}_{\text{lin}}(x, x') \cup \mathring{\mathcal{H}}_{\text{lin}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}}}^*(x, x') \\ &= |\eta(x) - \eta(x')| \tanh(kW\|x - x'\|_p) \\ &\geq |\eta(x) - \eta(x')| \tanh(kW\gamma) \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{lin}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{lin}}}(h, x, x') \geq \tanh(kW\gamma) \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}_{\mathcal{C}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{lin} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}$, valid for all $h \in \mathcal{H}_{\text{lin}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{lin}}) \leq \frac{\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{lin}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{lin}})}{\tanh(kW\gamma)} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{lin}}). \quad (75)$$

L.2. One-Hidden-Layer ReLU Neural Networks

Since \mathcal{H}_{NN} satisfies the condition of Lemma L.1, by Lemma L.1 the $(\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}_{\text{NN}})$ -minimizability gap can be expressed as follows:

$$\mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}) = \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} [\min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \mathbb{1}_{\|x - x'\|_p > \gamma} + c \mathbb{1}_{\|x - x'\|_p \leq \gamma}]. \quad (76)$$

By the definition of \mathcal{H}_{NN} , for any $(x, x') \in \mathcal{X} \times \mathcal{X}$, $\{h(x') - h(x) \mid h \in \mathcal{H}_{\text{NN}}\} = [-\Lambda W\|x - x'\|_p, \Lambda W\|x - x'\|_p]$.

 L.2.1. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}$.

For the hinge loss function $\Phi_{\text{hinge}}(u) := \max\{0, 1 - u\}$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h, x, x') &= \eta(x)(1 - \eta(x'))\Phi_{\text{hinge}}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\text{hinge}}(h(x') - h(x)) \\ &= \eta(x)(1 - \eta(x')) \max\{0, 1 - h(x) + h(x')\} + \eta(x')(1 - \eta(x)) \max\{0, 1 + h(x) - h(x')\}. \end{aligned}$$

Then,

$$\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}^*(x, x') = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - |\eta(x) - \eta(x')| \min\{\Lambda W\|x - x'\|_p, 1\}.$$

The $(\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}})$ -minimizability gap is

$$\begin{aligned} & \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{NN}}) \\ &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\ &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - |\eta(x) - \eta(x')| \min\{\Lambda W \|x - x'\|_p, 1\} \right]. \end{aligned} \quad (77)$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \hat{\mathcal{H}}_{\text{NN}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \\ & \geq \inf_{h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \hat{\mathcal{H}}_{\text{NN}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= \eta(x)(1 - \eta(x')) \max\{0, 1 - 0\} + \eta(x')(1 - \eta(x)) \max\{0, 1 + 0\} - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= |\eta(x) - \eta(x')| \min\{\Lambda W \|x - x'\|_p, 1\} \\ & \geq |\eta(x) - \eta(x')| \min\{\Lambda W \gamma, 1\} \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \min\{\Lambda W \gamma, 1\} \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}$, valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \frac{\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{hinge}}}}(\mathcal{H}_{\text{NN}})}{\min\{\Lambda W \gamma, 1\}} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \quad (78)$$

L.2.2. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\rho}}$.

For the ρ -margin loss function $\Phi_{\rho}(u) := \min\{1, \max\{0, 1 - \frac{u}{\rho}\}\}$, $\rho > 0$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}(h, x, x') \\ &= \eta(x)(1 - \eta(x')) \Phi_{\rho}(h(x) - h(x')) + \eta(x')(1 - \eta(x)) \Phi_{\rho}(h(x') - h(x)) \\ &= \eta(x)(1 - \eta(x')) \min\left\{1, \max\left\{0, 1 - \frac{h(x) - h(x')}{\rho}\right\}\right\} \\ & \quad + \eta(x')(1 - \eta(x)) \min\left\{1, \max\left\{0, 1 + \frac{h(x) - h(x')}{\rho}\right\}\right\} \end{aligned}$$

Then,

$$\begin{aligned} & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}(h, x, x') \\ &= \min\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} + \max\{\eta(x)(1 - \eta(x')), \eta(x')(1 - \eta(x))\} \left(1 - \frac{\min\{\Lambda W \|x - x'\|_p, \rho\}}{\rho}\right). \end{aligned}$$

The $(\tilde{\mathcal{L}}_{\Phi_{\rho}}, \mathcal{H}_{\text{NN}})$ -minimizability gap is

$$\begin{aligned} & \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}(\mathcal{H}_{\text{NN}}) \\ &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\ &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\rho}}}^*(\mathcal{H}_{\text{NN}}) \\ & \quad - \mathbb{E}_{(X, X')} \left[[\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] + [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \left(1 - \frac{\min\{\Lambda W \|x - x'\|_p, \rho\}}{\rho}\right) \right]. \end{aligned} \quad (79)$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \mathcal{H}_{\text{NN}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{NN}}}(h, x, x') \\
 & \geq \inf_{\tilde{\mathcal{H}}_{\text{NN}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 & = [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \left(1 - \frac{\min\{\Lambda W \|x - x'\|_p, \rho\}}{\rho}\right) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 & = |\eta(x) - \eta(x')| \frac{\min\{\Lambda W \|x - x'\|_p, \rho\}}{\rho} \\
 & \geq |\eta(x) - \eta(x')| \frac{\min\{\Lambda W \gamma, \rho\}}{\rho}
 \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_\rho}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \frac{\min\{\Lambda W \gamma, \rho\}}{\rho} \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}_{\mathcal{C}}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_\rho}$, valid for NN $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \frac{\rho \left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_\rho}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_\rho}}(\mathcal{H}_{\text{NN}}) \right)}{\min\{\Lambda W \gamma, \rho\}} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \quad (80)$$

L.2.3. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}$.

For the exponential loss function $\Phi_{\text{exp}}(u) := e^{-u}$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned}
 & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h, x, x') \\
 & = \eta(x)(1 - \eta(x'))\Phi_{\text{exp}}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\text{exp}}(h(x') - h(x)) \\
 & = \eta(x)(1 - \eta(x'))e^{-h(x)+h(x')} + \eta(x')(1 - \eta(x))e^{h(x)-h(x')}.
 \end{aligned}$$

Then,

$$\begin{aligned}
 & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 & = \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h, x, x') \\
 & = \begin{cases} 2\sqrt{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x'))} \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| \leq \Lambda W \|x - x'\|_p \\ \max\{\eta(x)(1-\eta(x')), \eta(x')(1-\eta(x))\} e^{-\Lambda W \|x-x'\|_p} + \min\{\eta(x)(1-\eta(x')), \eta(x')(1-\eta(x))\} e^{\Lambda W \|x-x'\|_p} \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > \Lambda W \|x - x'\|_p. \end{cases}
 \end{aligned}$$

The $(\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}})$ -minimizability gap is:

$$\begin{aligned}
 \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{NN}}) & = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}^*}(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\
 & = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}^*}(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[2\sqrt{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x'))} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| \leq \Lambda W \|x-x'\|_p} \right. \\
 & \quad \left. - \mathbb{E}_{(X, X')} \left[[\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] e^{-\Lambda W \|x-x'\|_p} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > \Lambda W \|x-x'\|_p} \right] \right. \\
 & \quad \left. - \mathbb{E}_{(X, X')} \left[[\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] e^{\Lambda W \|x-x'\|_p} \mathbb{1}_{\frac{1}{2} \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > \Lambda W \|x-x'\|_p} \right] \right]. \quad (81)
 \end{aligned}$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \mathring{\mathcal{H}}_{\text{NN}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \\
 & \geq \inf_{h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \mathring{\mathcal{H}}_{\text{NN}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 & = \eta(x)(1 - \eta(x'))e^{-0} + \eta(x')(1 - \eta(x))e^0 - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 & = \begin{cases} \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))} \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| \leq \Lambda W \|x - x'\|_p \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - e^{-\Lambda W \|x - x'\|_p}) + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - e^{\Lambda W \|x - x'\|_p}) \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| > \Lambda W \|x - x'\|_p \end{cases} \\
 & \geq \begin{cases} \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 2\sqrt{\eta(x)\eta(x')(1 - \eta(x))(1 - \eta(x'))} \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| \leq \Lambda W \gamma \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - e^{-\Lambda W \gamma}) + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - e^{\Lambda W \gamma}) \\ \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| > \Lambda W \gamma \end{cases} \\
 & = \begin{cases} \left(\frac{\eta(x)(1 - \eta(x')) - \eta(x')(1 - \eta(x))}{\sqrt{\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x))}} \right)^2 & \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| \leq \Lambda W \gamma \\ \frac{\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x))}{2} (2 - e^{-\Lambda W \gamma} - e^{\Lambda W \gamma}) + \frac{1}{2} |\eta(x) - \eta(x')| (e^{\Lambda W \gamma} - e^{-\Lambda W \gamma}) & \text{if } \frac{1}{2} \left| \log \frac{\eta(x)(1 - \eta(x'))}{\eta(x')(1 - \eta(x))} \right| > \Lambda W \gamma \end{cases} \\
 & \geq \min \left\{ (\eta(x) - \eta(x'))^2, \left(\frac{e^{2\Lambda W \gamma} + 1}{e^{2\Lambda W \gamma} - 1} \right) |\eta(x) - \eta(x')| \right\}
 \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \Psi_{\text{exp}} \left(\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}_{\mathcal{C}}}(h, x, x') \right).$$

where Ψ_{exp} is the increasing function on $[0, 2]$ defined by

$$\forall t \in [0, 1], \quad \Psi_{\text{exp}}(t) = \min \left\{ t^2, \left(\frac{e^{2\Lambda W \gamma} + 1}{e^{2\Lambda W \gamma} - 1} \right) t \right\}.$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}$, valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \Gamma_{\text{exp}} \left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{exp}}}}(\mathcal{H}_{\text{NN}}) \right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \quad (82)$$

where $\Gamma_{\text{exp}}(t) = \max \left\{ \sqrt{t}, \left(\frac{e^{2\Lambda W \gamma} - 1}{e^{2\Lambda W \gamma} + 1} \right) t \right\}$.

L.2.4. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\log}}$.

For the logistic loss function $\Phi_{\log}(u) := \log_2(1 + e^{-u})$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned}
 & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(h, x, x') \\
 & = \eta(x)(1 - \eta(x'))\Phi_{\log}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\log}(h(x') - h(x)) \\
 & = \eta(x)(1 - \eta(x')) \log_2 \left(1 + e^{-h(x) + h(x')} \right) + \eta(x')(1 - \eta(x)) \log_2 \left(1 + e^{h(x) - h(x')} \right).
 \end{aligned}$$

Then,

$$\begin{aligned}
 & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 &= \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(h, x, x') \\
 &= \begin{cases} -\eta(x)(1-\eta(x')) \log_2(\eta(x)(1-\eta(x'))) - \eta(x')(1-\eta(x)) \log_2(\eta(x')(1-\eta(x))) \\ \text{if } \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| \leq \Lambda W \|x-x'\|_p \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \log_2(1 + e^{-\Lambda W \|x-x'\|_p}) + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \log_2(1 + e^{\Lambda W \|x-x'\|_p}) \\ \text{if } \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > \Lambda W \|x-x'\|_p \end{cases}
 \end{aligned}$$

The $(\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}})$ -minimizability gap is

$$\begin{aligned}
 & \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(\mathcal{H}_{\text{NN}}) \\
 &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\
 &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[-\eta(x)(1-\eta(x')) \log_2(\eta(x)(1-\eta(x'))) \right. \\
 &\quad \left. - \eta(x')(1-\eta(x)) \log_2(\eta(x')(1-\eta(x))) \mathbb{1}_{\left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| \leq \Lambda W \|x-x'\|_p} \right] \\
 &\quad - \mathbb{E}_{(X, X')} \left[[\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \log_2(1 + e^{-\Lambda W \|x-x'\|_p}) \mathbb{1}_{\left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > \Lambda W \|x-x'\|_p} \right] \\
 &\quad - \mathbb{E}_{(X, X')} \left[[\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] \log_2(1 + e^{\Lambda W \|x-x'\|_p}) \mathbb{1}_{\left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > \Lambda W \|x-x'\|_p} \right].
 \end{aligned} \tag{83}$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \hat{\mathcal{H}}_{\text{NN}}(x, x')$,

$$\begin{aligned}
 & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}(h, x, x') \\
 &\geq \inf_{h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \hat{\mathcal{H}}_{\text{NN}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 &= \eta(x)(1-\eta(x')) \log_2(1 + e^{-0}) + \eta(x')(1-\eta(x)) \log_2(1 + e^0) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\
 &\geq \begin{cases} \eta(x)(1-\eta(x'))[1 - \log_2(\eta(x)(1-\eta(x')))] + \eta(x')(1-\eta(x))[1 - \log_2(\eta(x')(1-\eta(x)))] \\ \text{if } \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| \leq \Lambda W \gamma \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - \log_2(1 + e^{-\Lambda W \gamma})) + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')] (1 - \log_2(1 + e^{\Lambda W \gamma})) \\ \text{if } \left| \log \frac{\eta(x)(1-\eta(x'))}{\eta(x')(1-\eta(x))} \right| > \Lambda W \gamma \end{cases} \\
 &\geq \min \left\{ (\eta(x) - \eta(x'))^2, \left(\frac{e^{\Lambda W \gamma} + 1}{e^{\Lambda W \gamma} - 1} \right) |\eta(x) - \eta(x')| \right\}
 \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\log}}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \Psi_{\log} \left(\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}, \mathcal{H}}(h, x, x') \right).$$

where Ψ_{\log} is the increasing function on $[0, 2]$ defined by

$$\forall t \in [0, 1], \quad \Psi_{\log}(t) = \min \left\{ t^2, \left(\frac{e^{\Lambda W \gamma} + 1}{e^{\Lambda W \gamma} - 1} \right) t \right\}.$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\log}}$, valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \Gamma_{\log} \left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\log}}}(\mathcal{H}_{\text{NN}}) \right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \tag{84}$$

where $\Gamma_{\log}(t) = \max \left\{ \sqrt{t}, \left(\frac{e^{\Lambda W \gamma} - 1}{e^{\Lambda W \gamma} + 1} \right) t \right\}$.

L.2.5. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}$.

For the squared hinge loss function $\Phi_{\text{sq}}(u) := (1-u)^2 \mathbb{1}_{u \leq 1}$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h, x, x') \\ &= \eta(x)(1 - \eta(x'))\Phi_{\text{sq}}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\text{sq}}(h(x') - h(x)) \\ &= \eta(x)(1 - \eta(x'))(1 - h(x) + h(x'))^2 \mathbb{1}_{h(x) - h(x') \leq 1} + \eta(x')(1 - \eta(x))(1 + h(x) - h(x'))^2 \mathbb{1}_{h(x) - h(x') \geq -1}. \end{aligned}$$

Then,

$$\begin{aligned} & \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h, x, x') \\ &= \begin{cases} 4 \frac{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x'))}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \\ \text{if } \frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \leq \Lambda W \|x-x'\|_p \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')](1 - \Lambda W \|x-x'\|_p)^2 + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')](1 + \Lambda W \|x-x'\|_p)^2 \\ \text{if } \frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} > \Lambda W \|x-x'\|_p. \end{cases} \end{aligned}$$

The $(\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}})$ -minimizability gap is

$$\begin{aligned} & \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{NN}}) \\ &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[\mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \right] \\ &= \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} \left[4 \frac{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x'))}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \mathbb{1}_{\frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \leq \Lambda W \|x-x'\|_p} \right. \\ & \quad \left. - \mathbb{E}_{(X, X')} \left[[\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')](1 - \Lambda W \|x-x'\|_p)^2 \mathbb{1}_{\frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} > \Lambda W \|x-x'\|_p} \right] \right. \\ & \quad \left. - \mathbb{E}_{(X, X')} \left[[\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')](1 + \Lambda W \|x-x'\|_p)^2 \mathbb{1}_{\frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} > \Lambda W \|x-x'\|_p} \right] \right]. \end{aligned} \quad (85)$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \mathring{\mathcal{H}}_{\text{NN}}(x, x')$,

$$\begin{aligned} & \Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \\ & \geq \inf_{h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \mathring{\mathcal{H}}_{\text{NN}}(x, x')} \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h, x, x') - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ & \geq \begin{cases} \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - 4 \frac{\eta(x)\eta(x')(1-\eta(x))(1-\eta(x'))}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \\ \text{if } \frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} \leq \Lambda W \gamma \\ [\max\{\eta(x), \eta(x')\} - \eta(x)\eta(x')][1 - (1 - \Lambda W \gamma)^2] + [\min\{\eta(x), \eta(x')\} - \eta(x)\eta(x')][1 - (1 + \Lambda W \gamma)^2] \\ \text{if } \frac{|\eta(x)-\eta(x')|}{\eta(x)(1-\eta(x'))+\eta(x')(1-\eta(x))} > \Lambda W \gamma. \end{cases} \\ & \geq \min\left\{(\eta(x) - \eta(x'))^2, 2\Lambda W \gamma |\eta(x) - \eta(x')| - (\Lambda W \gamma)^2\right\} \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \Psi_{\text{sq}}\left(\Delta \mathcal{C}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x')\right).$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}$, valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \Gamma_{\text{sq}} \left(\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sq}}}}(\mathcal{H}_{\text{NN}}) \right) - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \quad (86)$$

where $\Gamma_{\text{sq}} = \max\left\{\sqrt{t}, \frac{t}{2\Lambda W \gamma} + \frac{\Lambda W \gamma}{2}\right\}$.

L.2.6. DERIVATION FOR $\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}$.

For the sigmoid loss function $\Phi_{\text{sig}}(u) := 1 - \tanh(ku)$, $k > 0$, for all $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\begin{aligned} \mathcal{E}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h, x, x') &= \eta(x)(1 - \eta(x'))\Phi_{\text{sig}}(h(x) - h(x')) + \eta(x')(1 - \eta(x))\Phi_{\text{sig}}(h(x') - h(x)) \\ &= \eta(x)(1 - \eta(x'))(1 - \tanh(k[h(x) - h(x')])) + \eta(x')(1 - \eta(x))(1 + \tanh(k[h(x) - h(x')])) \end{aligned}$$

Then,

$$\mathcal{E}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}^*(x, x') = \inf_{h \in \mathcal{H}_{\text{NN}}} \mathcal{E}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h, x, x') = \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - |\eta(x) - \eta(x')| \tanh(k\Lambda W \|x - x'\|_p).$$

The $(\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}})$ -minimizability gap is

$$\mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{NN}}) = \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{NN}}) - \mathbb{E}_{(X, X')} [\eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - |\eta(x) - \eta(x')| \tanh(k\Lambda W \|x - x'\|_p)]. \quad (87)$$

Therefore, $\forall h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \hat{\mathcal{H}}_{\text{NN}}(x, x')$,

$$\begin{aligned} \Delta \mathcal{E}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}(h, x, x') &\geq \inf_{h \in \tilde{\mathcal{H}}_{\text{NN}}(x, x') \cup \hat{\mathcal{H}}_{\text{NN}}(x, x')} \mathcal{E}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h, x, x') - \mathcal{E}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= \eta(x)(1 - \eta(x')) + \eta(x')(1 - \eta(x)) - \mathcal{E}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}^*(x, x') \\ &= |\eta(x) - \eta(x')| \tanh(k\Lambda W \|x - x'\|_p) \\ &\geq |\eta(x) - \eta(x')| \tanh(k\Lambda W \gamma) \end{aligned}$$

which implies that for any $h \in \mathcal{H}_{\text{NN}}$ and (x, x') such that $\|x - x'\|_p > \gamma$,

$$\Delta \mathcal{E}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}, \mathcal{H}_{\text{NN}}}(h, x, x') \geq \tanh(k\Lambda W \gamma) \Delta \mathcal{E}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}, \mathcal{H}}(h, x, x').$$

Thus, by Theorem C.1 or Theorem C.2, setting $\epsilon = 0$ yields the \mathcal{H}_{NN} -consistency bound for $\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}$, valid for all $h \in \mathcal{H}_{\text{NN}}$:

$$\mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}^*(\mathcal{H}_{\text{NN}}) \leq \frac{\mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(h) - \mathcal{R}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}^*(\mathcal{H}_{\text{NN}}) + \mathcal{M}_{\tilde{\mathcal{L}}_{\Phi_{\text{sig}}}}(\mathcal{H}_{\text{NN}})}{\tanh(k\Lambda W \gamma)} - \mathcal{M}_{\tilde{\mathcal{L}}_{0-1}^{\text{abs}}}(\mathcal{H}_{\text{NN}}). \quad (88)$$