

A. Proof of Theorem 1

Given θ (the parameter of the prediction function), we obtain the adversarial risk by the following optimization:

$$\max_{r \in \mathcal{U}_f} \mathbb{E}_{p(x,y)} [r(x,y) \ell(g_\theta(x), y)] \quad (24)$$

$$\text{where } \mathcal{U}_f \equiv \{r(x,y) \mid \mathbb{E}_{p(x,y)} [f(r(x,y))] \leq \delta, \mathbb{E}_{p(x,y)} [r(x,y)] = 1, r(x,y) \geq 0 \ (\forall (x,y) \in \mathcal{X} \times \mathcal{Y})\}. \quad (25)$$

Here, we are considering that $\ell(\cdot, \cdot)$ is the 0-1 loss. Let $\Omega_\theta^{(0)} \equiv \{(x,y) \mid \ell(g_\theta(x), y) = 0\} \subseteq \mathcal{X} \times \mathcal{Y}$. Then, we have $\Omega_\theta^{(1)} \equiv \{(x,y) \mid \ell(g_\theta(x), y) = 1\} = \mathcal{X} \times \mathcal{Y} \setminus \Omega_\theta^{(0)}$. Let $r^*(\cdot, \cdot)$ be the optimal solution of Eq. (24).

Note that Eq. (24) is a convex optimization problem. This is because the objective is linear in $r(\cdot, \cdot)$ and the uncertainty set is convex, which follows from the fact that $f(\cdot)$ is convex. Therefore, any local maximum of Eq. (24) is the global maximum. Nonetheless, there can be multiple solutions that attain the same global maxima. Among those solutions, we now show that there exists $r^*(\cdot, \cdot)$ such that it takes the same values within $\Omega_\theta^{(0)}$ and $\Omega_\theta^{(1)}$, respectively, i.e.,

$$r^*(x,y) = r_0^* \text{ for } \forall (x,y) \in \Omega_\theta^{(0)}, \quad r^*(x,y) = r_1^* \text{ for } \forall (x,y) \in \Omega_\theta^{(1)}, \quad (26)$$

where r_0^* and r_1^* are some constant values. This is because for any given optimal solution $r^{i*}(\cdot, \cdot)$ of Eq. (24), we can always obtain optimal solution $r^*(\cdot, \cdot)$ that satisfies Eq. (26) by the following transformation:

$$r_0^* = \mathbb{E}_{p(x,y)} [r^{i*}(x,y) \cdot \mathbf{1}\{(x,y) \in \Omega_\theta^{(0)}\}] / \mathbb{E}_{p(x,y)} [\mathbf{1}\{(x,y) \in \Omega_\theta^{(0)}\}], \quad (27)$$

$$r_1^* = \mathbb{E}_{p(x,y)} [r^{i*}(x,y) \cdot \mathbf{1}\{(x,y) \in \Omega_\theta^{(1)}\}] / \mathbb{E}_{p(x,y)} [\mathbf{1}\{(x,y) \in \Omega_\theta^{(1)}\}], \quad (28)$$

where $\mathbf{1}\{\cdot\}$ is an indicator function. Eqs. (27) and (28) are simple average operations of $r^{i*}(\cdot, \cdot)$ on regions $\Omega_\theta^{(0)}$ and $\Omega_\theta^{(1)}$, respectively. Utilizing the convexity of $f(\cdot)$, it is straightforward to see that $r^*(\cdot, \cdot)$ constructed in this way is still in the feasible region of Eq. (25). It also attains exactly the same objective of Eq. (24) as the original solution $r^{i*}(\cdot, \cdot)$ does. This concludes our proof that there exists an optimal solution of Eq. (24) such that it takes the same values within $\Omega_\theta^{(0)}$ and $\Omega_\theta^{(1)}$, respectively.

Let $p_{\Omega_\theta^{(i)}} = \mathbb{E}_{p(x,y)} [\mathbf{1}\{(x,y) \in \Omega_\theta^{(i)}\}]$ for $i = 0, 1$, which is the proportion of data that is correctly and incorrectly classified, respectively. We note that $p_{\Omega_\theta^{(0)}} + p_{\Omega_\theta^{(1)}} = 1$. Also, we see that $p_{\Omega_\theta^{(1)}}$ is by definition the misclassification rate; thus, it is equal to the ordinary risk, i.e., $\mathbb{E}_{p(x,y)} [\ell(g_\theta(x), y)]$. By using Eq. (26), we can simplify Eq. (24) as

$$\sup_{(r_0, r_1) \in \mathcal{U}'_f} p_{\Omega_\theta^{(1)}} r_1, \quad (29)$$

$$\text{where } \mathcal{U}'_f \equiv \{(r_0, r_1) \mid p_{\Omega_\theta^{(0)}} f(r_0) + p_{\Omega_\theta^{(1)}} f(r_1) \leq \delta, p_{\Omega_\theta^{(0)}} r_0 + p_{\Omega_\theta^{(1)}} r_1 = 1, r_0 \geq 0, r_1 \geq 0\}. \quad (30)$$

In the following, we show that Eq. (29) has monotonic relationship with $p_{\Omega_\theta^{(1)}}$. With a fixed value of $p_{\Omega_\theta^{(1)}}$, we can obtain the optimal (r_0, r_1) by solving Eq. (29). Let $(r_0^*(p_1), r_1^*(p_1))$ be the solution of Eq. (29) when $p_{\Omega_\theta^{(1)}}$ is fixed to p_1 with $0 \leq p_1 \leq 1$.

First, we note that the first inequality constraint in Eq. (30) is a convex set and includes $(1, 1)$ in its relative interior because $p_{\Omega_\theta^{(0)}} f(1) + p_{\Omega_\theta^{(1)}} f(1) = 0 < \delta$ and $p_{\Omega_\theta^{(0)}} \cdot 1 + p_{\Omega_\theta^{(1)}} \cdot 1 = 1$. Note that in a two dimensional space, the number of intersections between a line and a boundary of a convex set is at most two. For $\delta > 0$, there are always exactly two different points that satisfies both $p_{\Omega_\theta^{(0)}} f(r_0) + p_{\Omega_\theta^{(1)}} f(r_1) = \delta$ and $p_{\Omega_\theta^{(0)}} r_0 + p_{\Omega_\theta^{(1)}} r_1 = 1$. We further see that the optimal solution of r_1 is always greater than 1 because the objective in Eq. (29) is an increasing function of r_1 . Taking these facts into account, we can see that the optimal solution, $(r_0^*(p_1), r_1^*(p_1))$, satisfies either of the following two cases depending on whether the inequality constraint $r_0 \geq 0$ in Eq. (30) is active or not.

$$\textbf{Case 1:} \quad p_0 \cdot f(r_0^*(p_1)) + p_1 \cdot f(r_1^*(p_1)) = \delta, \quad p_0 \cdot r_0^*(p_1) + p_1 \cdot r_1^*(p_1) = 1, \quad 0 < r_0^*(p_1) < 1 < r_1^*(p_1). \quad (31)$$

$$\textbf{Case 2:} \quad r_0^*(p_1) = 0, \quad r_1^*(p_1) = \frac{1}{p_1}, \quad (32)$$

where $p_0 = 1 - p_1$. We now show that there is the monotonic relation between $p_{\Omega_\theta^{(1)}}$ and Eq. (29) for both the cases. To this end, pick any p'_1 such that $p_1 < p'_1 \leq 1$, and let $(r_0^*(p'_1), r_1^*(p'_1))$ be the solution of Eq. (29) when $p_{\Omega_\theta^{(1)}}$ is fixed to p'_1 .

Regarding Case 2 in Eq. (32), the adversarial risk in Eq. (29) is $p_1 \cdot \frac{1}{p_1} = 1$. On the other hand, it is easy to see that the active equality constraint stays $r_0 \geq 0$ in Eq. (30) for $p_{\Omega_\theta^{(1)}}$ larger p_1 . Hence, we can show that $r_0^*(p'_1) = 0, r_1^*(p'_1) = \frac{1}{p'_1}$. Therefore, the adversarial risk in Eq. (29) stays $p'_1 \cdot \frac{1}{p'_1} = 1$. This concludes our proof for Eq. (10) in Theorem 1.

Regarding Case 1 in Eq. (31), we note that both the ordinary risk p_1 and the adversarial risk $p_1 \cdot r_1^*(p_1)$ are strictly less than 1. Our goal is to show Eq. (9) in Theorem 1, which is equivalent to showing

$$p_1 \cdot r_1^*(p_1) < p'_1 \cdot r_1^*(p'_1). \quad (33)$$

To do so, we further consider the following two sub-cases of Case 1 in Eq. (31):

$$\text{Case 1-a: } p'_1 < p_1 \cdot r_1^*(p_1) \quad (34)$$

$$\text{Case 1-b: } p'_1 \geq p_1 \cdot r_1^*(p_1). \quad (35)$$

In Case 1-b, we can straightforwardly show Eq. (33) as follows.

$$p_1 \cdot r_1^*(p_1) \leq p'_1 < p'_1 \cdot r_1^*(p'_1), \quad (36)$$

where the last inequality follows from $1 < r_1^*(p'_1)$.

Now, assume Cases 1 and 1-a in Eqs. (31) and (34). Our goal is to show that $r_1^*(p'_1)$ satisfies $r_1^*(p'_1) > \frac{p_1}{p'_1} r_1^*(p_1)$ because then Eq. (33) holds. To this end, we show that

$$r'_1 = \frac{p_1}{p'_1} \cdot r_1^*(p_1) \quad (37)$$

is contained in the relative interior (excluding the boundary) of Eq. (30) with $p_{\Omega_\theta^{(1)}}$ fixed to p'_1 . Then, because our objective in Eq. (29) is linear in r_1 , $r'_1 < r_1^*(p'_1)$ holds in our setting. Then, we arrive at $r_1^*(p'_1) > \frac{p_1}{p'_1} \cdot r_1^*(p_1)$. Formally, our goal is to show r'_1 in Eq. (37) satisfies

$$p'_0 \cdot f(r'_0) + p'_1 \cdot f(r'_1) < \delta, \quad p'_0 r'_0 + p'_1 r'_1 = 1, \quad r'_0 > 0, \quad r'_1 > 0, \quad (38)$$

where $p'_0 = 1 - p'_1$. By Eqs. (31), (37) and the second equality of Eq. (38), we have

$$r'_0 = \frac{1 - p'_1 r'_1}{p'_0} = \frac{1 - p_1 \cdot r_1^*(p_1)}{p'_0} = \frac{p_0}{p'_0} \cdot r_0^*(p_1). \quad (39)$$

The latter two inequalities of Eq. (38), i.e., $r'_0 > 0$ and $r'_1 > 0$ follow straightforwardly from the assumptions. Combining the assumption in Eq. (34) and the last inequality in Eq. (31), we have the following inequality.

$$0 < r_0^*(p_1) < r'_0 < 1 < r'_1 < r_1^*(p_1). \quad (40)$$

Thus, we can write r'_0 (resp. r'_1) as a linear interpolation of $r_0^*(p_1)$ and 1 (resp. 1 and $r_1^*(p_1)$) as follows.

$$r'_0 = \alpha \cdot r_0^*(p_1) + (1 - \alpha) \cdot 1, \quad r'_1 = \beta \cdot r_1^*(p_1) + (1 - \beta) \cdot 1, \quad (41)$$

where $0 < \alpha, \beta < 1$. Substituting Eqs. (37) and (39), we have

$$\alpha = \frac{1}{1 - r_0^*(p_1)} \cdot \frac{p'_0 - p_0 \cdot r_0^*(p_1)}{p'_0}, \quad (42)$$

$$\beta = \frac{1}{r_1^*(p_1) - 1} \cdot \frac{p_1 \cdot r_1^*(p_1) - p'_1}{p'_1} = \frac{1}{r_1^*(p_1) - 1} \cdot \frac{p'_0 - p_0 \cdot r_0^*(p_1)}{p'_1}. \quad (43)$$

Then, we have

$$\begin{aligned}
 p'_0 f(r'_0) + p'_1 f(r'_1) &= p'_0 f(\alpha \cdot r_0^*(p_1) + (1 - \alpha) \cdot 1) + p'_1 f(\beta \cdot r_1^*(p_1) + (1 - \beta) \cdot 1) \\
 &\leq p'_0 \cdot \{\alpha \cdot f(r_0^*(p_1)) + (1 - \alpha) \cdot f(1)\} + p'_1 \cdot \{\beta \cdot f(r_1^*(p_1)) + (1 - \beta) \cdot f(1)\} \quad (\because \text{convexity of } f(\cdot)) \\
 &= p'_0 \alpha \cdot f(r_0^*(p_1)) + p'_1 \beta \cdot f(r_1^*(p_1)) \quad (\because f(1) = 0) \\
 &= (p'_0 - p_0 \cdot r_0^*(p_1)) \left(\frac{1}{1 - r_0^*(p_1)} \cdot f(r_0^*(p_1)) + \frac{1}{r_1^*(p_1) - 1} \cdot f(r_1^*(p_1)) \right) \quad (\because \text{Eqs. (42) and (43)}) \\
 &= (p_1 \cdot r_1^*(p_1) - p'_1) \left(\frac{1}{1 - r_0^*(p_1)} \cdot f(r_0^*(p_1)) + \frac{1}{r_1^*(p_1) - 1} \cdot f(r_1^*(p_1)) \right) \\
 &< (p_1 \cdot r_1^*(p_1) - p_1) \left(\frac{1}{1 - r_0^*(p_1)} \cdot f(r_0^*(p_1)) + \frac{1}{r_1^*(p_1) - 1} \cdot f(r_1^*(p_1)) \right) \quad (\because p'_1 > p_1) \\
 &= p_0 \cdot f(r_0^*(p_1)) + p_1 \cdot f(r_1^*(p_1)) \\
 &= \delta. \quad (\because \text{the first equation of Eq. (31).})
 \end{aligned}$$

This concludes our proof that Eq. (29) has monotonic relationship with $p_{\Omega_\theta^{(1)}}$. Recall that $p_{\Omega_\theta^{(1)}}$ is by definition equal to the ordinary risk, $\mathcal{R}(\theta)$. Therefore, for any pair of parameters θ_1 and θ_2 , if $\mathcal{R}_{\text{adv}}(\theta_1) < 1$, we have

$$\mathcal{R}(\theta_1) < \mathcal{R}(\theta_2) \implies \mathcal{R}_{\text{adv}}(\theta_1) < \mathcal{R}_{\text{adv}}(\theta_2). \quad (44)$$

To show that the opposite direction of Eq. (44) holds, we need to show that any pair of parameters θ_1 and θ_2 , the following holds:

$$\mathcal{R}(\theta_1) = \mathcal{R}(\theta_2) \implies \mathcal{R}_{\text{adv}}(\theta_1) = \mathcal{R}_{\text{adv}}(\theta_2). \quad (45)$$

This is obvious from Eq. (29) because the adversarial risk depends on the parameter of the model *only through the risk of the model*. This concludes the proof of Theorem 1, in which the adversarial risk and ordinary risk are compared. For the case of empirical approximations, the same argument can be used by replacing the expectations with empirical averages. \square

B. Proof of Theorem 2

We prove by contradiction. Let Ω be the subset of $\mathbb{R}^{|\mathcal{Y}|}$. We consider the multi-class classification and assume that the loss $\ell(\cdot, \cdot) : \Omega \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$ is classification-calibrated. Although we will mainly focus on a multi-class classification scenario, our proof easily extends to a binary classification scenario, which we will discuss at the end of the proof.

Let $g(\cdot) : \mathcal{X} \rightarrow \mathbb{R}^K$ be prediction function, where K is the number of classes. Assume the prediction function, g , can take any measurable functions. Then, g^* that minimizes the ordinary risk using the classification-calibrated loss, $\ell(\cdot, \cdot)$, i.e.,

$$\begin{aligned}
 g^* &= \arg \min_g \mathbb{E}_{p(x,y)} [\ell(g(x), y)] \\
 &= \arg \min_g \mathbb{E}_{p(x)} \left[\sum_{y \in \mathcal{Y}} p(y|x) \ell(g(x), y) \right], \quad (46)
 \end{aligned}$$

is the Bayes optimal classifier¹⁰ (Bartlett et al., 2006; Tewari & Bartlett, 2007).

Our goal is to show that $g^{(\text{adv})}$ that minimizes the adversarial risk using classification-calibrated loss is also Bayes optimal w.r.t. $p(x, y)$. More specifically, we consider

$$g^{(\text{adv})} = \arg \min_g \sup_{q: D_f(q||p) \leq \delta} \mathbb{E}_{q(x,y)} [\ell(g(x), y)] \quad (47)$$

$$= \arg \min_g \sup_{r(\cdot, \cdot) \in \mathcal{U}_f} \mathbb{E}_{p(x,y)} [r(x, y) \ell(g(x), y)]. \quad (48)$$

¹⁰The classifier that minimizes the mis-classification rate for the training density $p(x, y)$ (the 0-1 loss is considered), i.e., the classifier whose prediction on x is equal to $\arg \max_{y \in \mathcal{Y}} p(y|x)$.

Recall that $q(x, y)$ in Eq. (47) and $r(x, y)$ in Eq. (48) are related by $r(x, y) \equiv q(x, y)/p(x, y)$. In the following, with a slight abuse of notation, we denote $r(x) \equiv q(x)/p(x)$ and $r(y|x) \equiv q(y|x)/p(y|x)$. Obviously, we have $r(x, y) = r(x)r(y|x)$.

Let $r^*(\cdot, \cdot)$ be the solution of the inner maximization of Eq. (48) with $g^{(\text{adv})}$, i.e.,

$$r^*(\cdot, \cdot) = \arg \max_{r(\cdot, \cdot) \in \mathcal{U}_f} \mathbb{E}_{p(x, y)} [r(x, y) \ell(g^{(\text{adv})}(x, y))]. \quad (49)$$

Then, by Danskin's theorem (Danskin, 1966), Eq. (48) can be rewritten as

$$\begin{aligned} g^{(\text{adv})} &= \arg \min_g \mathbb{E}_{p(x, y)} [r^*(x, y) \ell(g(x), y)] \\ &= \arg \min_g \mathbb{E}_{p(x)} [r^*(x) \mathbb{E}_{p(y|x)} [r^*(y|x) \ell(g(x), y)]] \\ &= \arg \min_g \mathbb{E}_{p(x)} \left[r^*(x) \sum_{y \in \mathcal{Y}} p(y|x) r^*(y|x) \ell(g(x), y) \right]. \end{aligned} \quad (50)$$

Now, suppose that $g^{(\text{adv})}$ is not Bayes optimal almost surely over $q^*(x) \equiv r^*(x)p(x)$. Then, we have

$$\int_{x \in \mathcal{S}} q^*(x) dx > 0, \quad (51)$$

where

$$\mathcal{S} \equiv \left\{ x \mid x \in \mathcal{X}, p(x) > 0, q^*(x) > 0, \arg \max_{y \in \mathcal{Y}} p(y|x) \neq \arg \max_{y \in \mathcal{Y}} g_y^{(\text{adv})}(x) \right\}. \quad (52)$$

In the following, we denote $x \in \mathcal{S}$ by x^\dagger , i.e., whenever we denote x^\dagger , we implicitly assume $x^\dagger \in \mathcal{S}$. We immediately have $r^*(x^\dagger) = q^*(x^\dagger)/p(x^\dagger) > 0$. We let $y^{(\text{max})}(x^\dagger) \equiv \arg \max_{y \in \mathcal{Y}} p(y|x^\dagger)$ and $y^{(\text{adv})}(x^\dagger) \equiv \arg \max_{y \in \mathcal{Y}} g_y^{(\text{adv})}(x^\dagger)$. Since $\ell(\cdot, \cdot)$ is classification-calibrated, from Eq. (50) and the definition of the classification-calibrated loss (Bartlett et al., 2006; Tewari & Bartlett, 2007), we have

$$y^{(\text{adv})}(x^\dagger) = \arg \max_{y \in \mathcal{Y}} p(y|x^\dagger) r^*(y|x^\dagger). \quad (53)$$

Because we have $x^\dagger \in \mathcal{S}$, $y^{(\text{adv})}(x^\dagger) \neq y^{(\text{max})}(x^\dagger)$ holds. Thus, we have

$$p(y^{(\text{adv})}(x^\dagger)|x^\dagger) r^*(y^{(\text{adv})}(x^\dagger)|x^\dagger) > p(y^{(\text{max})}(x^\dagger)|x^\dagger) r^*(y^{(\text{max})}(x^\dagger)|x^\dagger). \quad (54)$$

Combining this with $p(y^{(\text{max})}(x^\dagger)|x^\dagger) > p(y^{(\text{adv})}(x^\dagger)|x^\dagger)$, we have

$$r^*(y^{(\text{adv})}(x^\dagger)|x^\dagger) > r^*(y^{(\text{max})}(x^\dagger)|x^\dagger) > 0. \quad (55)$$

We construct a new ratio function, $r_{\text{new}}(\cdot, \cdot)$, by the following operations. We first set $r_{\text{new}}(x, y) \leftarrow r^*(x, y)$ for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$. Then, for $\forall x \in \mathcal{S}$, we let

$$r_{\text{new}}(y^{(\text{max})}(x)|x) \leftarrow r^*(y^{(\text{max})}(x)|x) + \epsilon \cdot p(y^{(\text{adv})}(x)|x), \quad (56)$$

$$r_{\text{new}}(y^{(\text{adv})}(x)|x) \leftarrow r^*(y^{(\text{adv})}(x)|x) - \epsilon \cdot p(y^{(\text{max})}(x)|x), \quad (57)$$

where $\epsilon > 0$ is a sufficiently small number. We show that such $r_{\text{new}}(\cdot, \cdot)$ is still in \mathcal{U}_f . As shown in Eqs. (56) and (57), the value of $r_{\text{new}}(\cdot, \cdot)$ changed from $r^*(\cdot, \cdot)$ only in \mathcal{S} . Therefore, given that $r^*(\cdot, \cdot) \in \mathcal{U}_f$, in order to show $r_{\text{new}}(\cdot, \cdot) \in \mathcal{U}_f$, it is sufficient to show the following three equality/inequalities:

$$\mathbb{E}_{p(x, y)} [f(r_{\text{new}}(x, y))] \leq \delta, \mathbb{E}_{p(x, y)} [r_{\text{new}}(x, y)] = 1, r_{\text{new}}(x, y) \geq 0 (\forall (x, y) \in \mathcal{X} \times \mathcal{Y}). \quad (58)$$

Because we know $r^*(\cdot, \cdot) \in \mathcal{U}_f$, it is sufficient to show

$$\begin{aligned} \mathbb{E}_{p(x, y)} [f(r_{\text{new}}(x, y))] &\leq \mathbb{E}_{p(x, y)} [f(r^*(x, y))], \\ \mathbb{E}_{p(x, y)} [r_{\text{new}}(x, y)] &= \mathbb{E}_{p(x, y)} [r^*(x, y)], \\ r_{\text{new}}(x, y) &\geq 0, \forall (x, y) \in \mathcal{X} \times \mathcal{Y}. \end{aligned} \quad (59)$$

Furthermore, $r_{\text{new}}(\cdot, \cdot)$ only differs from $r^*(\cdot, \cdot)$ in $(x, y^{(\max)}) \in \mathcal{S} \times \mathcal{Y}$ and $(x, y^{(\text{adv})}) \in \mathcal{S} \times \mathcal{Y}$. Therefore, to show Eq. (59), all we need to show is

$$\begin{aligned} & \int_{x \in \mathcal{S}} \left\{ p(x, y^{(\text{adv})}(x)) \cdot f(r_{\text{new}}(x, y^{(\text{adv})}(x))) + p(x^\dagger, y^{(\max)}(x)) \cdot f(r_{\text{new}}(x, y^{(\max)}(x))) \right\} dx \\ & \leq \int_{x \in \mathcal{S}} \left\{ p(x, y^{(\text{adv})}(x)) \cdot f(r^*(x, y^{(\text{adv})}(x))) + p(x, y^{(\max)}(x)) \cdot f(r^*(x, y^{(\max)}(x))) \right\} dx, \end{aligned} \quad (60)$$

$$\begin{aligned} & \int_{x \in \mathcal{S}} \left\{ p(x, y^{(\text{adv})}(x)) \cdot r_{\text{new}}(x, y^{(\text{adv})}(x)) + p(x, y^{(\max)}(x)) \cdot r_{\text{new}}(x, y^{(\max)}(x)) \right\} dx \\ & = \int_{x \in \mathcal{S}} \left\{ p(x, y^{(\text{adv})}(x)) \cdot r^*(x, y^{(\text{adv})}(x)) + p(x, y^{(\max)}(x)) \cdot r^*(x, y^{(\max)}(x)) \right\} dx, \end{aligned} \quad (61)$$

$$r_{\text{new}}(x, y^{(\text{adv})}(x)) \geq 0, \quad r_{\text{new}}(x, y^{(\max)}(x)) \geq 0, \quad \forall x \in \mathcal{S}. \quad (62)$$

First, since $r^*(y^{(\text{adv})}(x^\dagger)|x^\dagger) > 0$ holds from Eq. (54), $r_{\text{new}}(y^{(\text{adv})}(x^\dagger)|x^\dagger) = r^*(y^{(\text{adv})}(x^\dagger)|x^\dagger) - \epsilon \cdot p(y^{(\max)}(x^\dagger)|x^\dagger) \geq 0$ holds for sufficiently small ϵ . Thus, $r_{\text{new}}(x^\dagger, y^{(\text{adv})}(x^\dagger)) = r_{\text{new}}(y^{(\text{adv})}(x^\dagger)|x^\dagger)r(x^\dagger) \geq 0$. Hence, Eq. (62) holds for sufficiently small ϵ . Also, Eq. (61) follows because

Integrand of L.H.S. of Eq. (61)

$$\begin{aligned} & = p(x^\dagger, y^{(\text{adv})}(x^\dagger)) \cdot \left[r_{\text{new}}(x^\dagger) \cdot \left\{ r^*(y^{(\text{adv})}(x^\dagger)|x^\dagger) - \epsilon \cdot p(y^{(\max)}(x^\dagger)|x^\dagger) \right\} \right] \\ & \quad + p(x^\dagger, y^{(\max)}(x^\dagger)) \cdot \left[r_{\text{new}}(x^\dagger) \cdot \left\{ r^*(y^{(\max)}(x^\dagger)|x^\dagger) + \epsilon \cdot p(y^{(\text{adv})}(x^\dagger)|x^\dagger) \right\} \right] \end{aligned} \quad (63)$$

$$\begin{aligned} & = p(x^\dagger, y^{(\text{adv})}(x^\dagger)) \cdot \left[r^*(x^\dagger) \cdot \left\{ r^*(y^{(\text{adv})}(x^\dagger)|x^\dagger) - \epsilon \cdot p(y^{(\max)}(x^\dagger)|x^\dagger) \right\} \right] \\ & \quad + p(x^\dagger, y^{(\max)}(x^\dagger)) \cdot \left[r^*(x^\dagger) \cdot \left\{ r^*(y^{(\max)}(x^\dagger)|x^\dagger) + \epsilon \cdot p(y^{(\text{adv})}(x^\dagger)|x^\dagger) \right\} \right] \end{aligned} \quad (64)$$

$$\begin{aligned} & = p(x^\dagger, y^{(\text{adv})}(x^\dagger)) \cdot r^*(x^\dagger, y^{(\text{adv})}(x^\dagger)) + p(x^\dagger, y^{(\max)}(x^\dagger)) \cdot r^*(x^\dagger, y^{(\max)}(x^\dagger)) \\ & \quad + \underbrace{\epsilon \cdot r^*(x^\dagger)p(x^\dagger) \left\{ p(y^{(\text{adv})}(x^\dagger)|x^\dagger)p(y^{(\max)}(x^\dagger)|x^\dagger) - p(y^{(\max)}(x^\dagger)|x^\dagger)p(y^{(\text{adv})}(x^\dagger)|x^\dagger) \right\}}_{=0} \end{aligned} \quad (65)$$

$$= \text{Integrand of R.H.S. of Eq. (61)} \quad (66)$$

Finally, we show Eq. (60). Substituting Eqs. (56) and (57) into the L.H.S. of Eq. (60), we have

Integrand of L.H.S. of Eq. (60)

$$\begin{aligned} & = p(x^\dagger) \cdot \left\{ p(y^{(\text{adv})}(x^\dagger)|x^\dagger) \cdot f\left(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger)) - \epsilon \cdot r^*(x^\dagger)p(y^{(\max)}(x^\dagger)|x^\dagger)\right) \right. \\ & \quad \left. + p(y^{(\max)}(x^\dagger)|x^\dagger) \cdot f\left(r^*(x^\dagger, y^{(\max)}(x^\dagger)) + \epsilon \cdot r^*(x^\dagger)p(y^{(\text{adv})}(x^\dagger)|x^\dagger)\right) \right\}. \end{aligned} \quad (67)$$

Because $f(\cdot)$ is differentiable, we can apply the first order Taylor expansion to the two terms involving $f(\cdot)$ in Eq. (67) as

$$\begin{aligned} & f\left(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger)) - \epsilon \cdot r^*(x^\dagger)p(y^{(\max)}(x^\dagger)|x^\dagger)\right) \\ & = f\left(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))\right) - \epsilon \cdot f'\left(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))\right) \cdot r^*(x^\dagger)p(y^{(\max)}(x^\dagger)|x^\dagger) + \mathcal{O}(\epsilon^2), \\ & f\left(r^*(x^\dagger, y^{(\max)}(x^\dagger)) + \epsilon \cdot r^*(x^\dagger)p(y^{(\text{adv})}(x^\dagger)|x^\dagger)\right) \\ & = f\left(r^*(x^\dagger, y^{(\max)}(x^\dagger))\right) + \epsilon \cdot f'\left(r^*(x^\dagger, y^{(\max)}(x^\dagger))\right) \cdot r^*(x^\dagger)p(y^{(\text{adv})}(x^\dagger)|x^\dagger) + \mathcal{O}(\epsilon^2). \end{aligned} \quad (68)$$

Substituting Eq. (68) into Eq. (67), we have

Integrand of L.H.S. of Eq. (60)

$$\begin{aligned} & = p(x^\dagger, y^{(\text{adv})}(x^\dagger)) \cdot f\left(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))\right) + p(x^\dagger, y^{(\max)}(x^\dagger)) \cdot f\left(r^*(x^\dagger, y^{(\max)}(x^\dagger))\right) \\ & \quad + \epsilon \cdot r^*(x^\dagger)p(y^{(\max)}(x^\dagger)|x^\dagger)p(y^{(\text{adv})}(x^\dagger)|x^\dagger) \cdot \left\{ f'\left(r^*(x^\dagger, y^{(\max)}(x^\dagger))\right) - f'\left(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))\right) \right\} + \mathcal{O}(\epsilon^2), \end{aligned} \quad (69)$$

Since $f(\cdot)$ is convex, its derivative $f'(\cdot)$ is non-decreasing. Also $r^*(x^\dagger, y^{(\text{adv})}(x^\dagger)) > r^*(x^\dagger, y^{(\text{max})}(x^\dagger))$ holds because of Eq. (55). Therefore, we have $f'(r^*(x^\dagger, y^{(\text{max})}(x^\dagger))) - f'(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))) \leq 0$.

First, assume $f'(r^*(x^\dagger, y^{(\text{max})}(x^\dagger))) - f'(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))) = 0$. Then, $f(\cdot)$ is exactly linear in the interval of $[r^*(x^\dagger, y^{(\text{max})}(x^\dagger)), r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))]$; hence, $\mathcal{O}(\epsilon^2)$ term in Eq. (69) is exactly 0 for sufficiently small ϵ . Hence, we have

$$\begin{aligned} \text{Eq. (69)} &= p(x^\dagger, y^{(\text{adv})}(x^\dagger)) \cdot f(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))) + p(x^\dagger, y^{(\text{max})}(x^\dagger)) \cdot f(r^*(x^\dagger, y^{(\text{max})}(x^\dagger))) \\ &= \text{Integrand of R.H.S. of Eq. (60)}. \end{aligned} \quad (70)$$

Next, assume $f'(r^*(x^\dagger, y^{(\text{max})}(x^\dagger))) - f'(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))) < 0$. In this case, since the coefficient of ϵ in Eq. (71) is negative, there exists sufficiently small $\epsilon > 0$ such that

$$\epsilon \cdot \underbrace{r^*(x^\dagger)}_{>0} \underbrace{p(y^{(\text{max})}(x^\dagger)|x^\dagger)p(y^{(\text{adv})}(x^\dagger)|x^\dagger)}_{>0} \cdot \underbrace{\left\{ f'(r^*(x^\dagger, y^{(\text{max})}(x^\dagger))) - f'(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))) \right\}}_{<0} + \mathcal{O}(\epsilon^2) < 0 \quad (71)$$

Thus, we have

$$\begin{aligned} \text{Eq. (69)} &< p(x^\dagger, y^{(\text{adv})}(x^\dagger)) \cdot f(r^*(x^\dagger, y^{(\text{adv})}(x^\dagger))) + p(x^\dagger, y^{(\text{max})}(x^\dagger)) \cdot f(r^*(x^\dagger, y^{(\text{max})}(x^\dagger))) \\ &= \text{R.H.S. of Eq. (60)}. \end{aligned} \quad (72)$$

In both Eqs. (70) and (72), by taking an integral in \mathcal{S} , Eq. (60) holds.

In summary, since Eqs. (60), (61) and (62) all hold, the newly constructed, $r_{\text{new}}(\cdot, \cdot)$, is still in \mathcal{U}_f .

We now show that $r_{\text{new}}(\cdot, \cdot)$ actually gives larger objective of Eq. (49) than $r^*(\cdot, \cdot)$, which contradicts Eq. (49). Since the value of $r_{\text{new}}(\cdot, \cdot)$ is mostly the same as that of $r^*(\cdot, \cdot)$ except that we have Eqs. (56) and (57), we only need to consider the part they differ. Therefore, it is sufficient to show

$$\begin{aligned} &\int_{x \in \mathcal{S}} p(x)r^*(x) \underbrace{\left\{ p(y^{(\text{adv})}(x)|x)r_{\text{new}}(y^{(\text{adv})}(x)|x)\ell(g^{(\text{adv})}(x), y^{(\text{adv})}(x)) + p(y^{(\text{max})}(x)|x)r_{\text{new}}(y^{(\text{max})}(x)|x)\ell(g^{(\text{adv})}(x), y^{(\text{max})}(x)) \right\}}_{(a)} dx \\ &> \int_{x \in \mathcal{S}} p(x)r^*(x) \underbrace{\left\{ p(y^{(\text{adv})}(x)|x)r^*(y^{(\text{adv})}(x)|x)\ell(g^{(\text{adv})}(x), y^{(\text{adv})}(x)) + p(y^{(\text{max})}(x)|x)r^*(y^{(\text{max})}(x)|x)\ell(g^{(\text{adv})}(x), y^{(\text{max})}(x)) \right\}}_{(b)} dx. \end{aligned} \quad (73)$$

Subtracting (b) from (a) in Eq. (73), and using Eqs. (56) and (57), we have

$$\begin{aligned} &p(y^{(\text{adv})}(x^\dagger)|x^\dagger)\ell(g^{(\text{adv})}(x^\dagger), y^{(\text{adv})}(x^\dagger)) \left\{ r_{\text{new}}(y^{(\text{adv})}(x^\dagger)|x^\dagger) - r^*(y^{(\text{adv})}(x^\dagger)|x^\dagger) \right\} \\ &+ p(y^{(\text{max})}(x^\dagger)|x^\dagger)\ell(g^{(\text{adv})}(x^\dagger), y^{(\text{max})}(x^\dagger)) \left\{ r_{\text{new}}(y^{(\text{max})}(x^\dagger)|x^\dagger) - r^*(y^{(\text{max})}(x^\dagger)|x^\dagger) \right\} \\ &= \epsilon \cdot p(y^{(\text{adv})}(x^\dagger)|x^\dagger)p(y^{(\text{max})}(x^\dagger)|x^\dagger) \left\{ -\ell(g^{(\text{adv})}(x^\dagger), y^{(\text{adv})}(x^\dagger)) + \ell(g^{(\text{adv})}(x^\dagger), y^{(\text{max})}(x^\dagger)) \right\}. \end{aligned} \quad (74)$$

We now show $\ell(g^{(\text{adv})}(x^\dagger), y^{(\text{max})}(x^\dagger)) > \ell(g^{(\text{adv})}(x^\dagger), y^{(\text{adv})}(x^\dagger))$. Suppose $\ell(g^{(\text{adv})}(x^\dagger), y^{(\text{max})}(x^\dagger)) \leq \ell(g^{(\text{adv})}(x^\dagger), y^{(\text{adv})}(x^\dagger))$. Construct $g' \in \mathbb{R}^K$ by swapping the $y^{(\text{max})}$ -th and $y^{(\text{adv})}$ -th elements of $g^{(\text{adv})}(x^\dagger)$, while retaining other elements to be exactly the same. Then, because of the assumption that $\ell(\cdot, \cdot)$ is invariant to class permutation, we have $\ell(g', y^{(\text{max})}(x^\dagger)) = \ell(g^{(\text{adv})}(x^\dagger), y^{(\text{adv})}(x^\dagger))$ and $\ell(g', y^{(\text{adv})}(x^\dagger)) = \ell(g^{(\text{adv})}(x^\dagger), y^{(\text{max})}(x^\dagger))$. Combining this with Eq. (54), we have

$$\begin{aligned} &p(y^{(\text{adv})}(x^\dagger)|x^\dagger)r^*(y^{(\text{adv})}(x^\dagger)|x^\dagger) \cdot \ell(g', y^{(\text{adv})}(x^\dagger)) + p(y^{(\text{max})}(x^\dagger)|x^\dagger)r^*(y^{(\text{max})}(x^\dagger)|x^\dagger) \cdot \ell(g', y^{(\text{max})}(x^\dagger)) \\ &< p(y^{(\text{adv})}(x^\dagger)|x^\dagger)r^*(y^{(\text{adv})}(x^\dagger)|x^\dagger) \cdot \ell(g^{(\text{adv})}(x^\dagger), y^{(\text{adv})}(x^\dagger)) + p(y^{(\text{max})}(x^\dagger)|x^\dagger)r^*(y^{(\text{max})}(x^\dagger)|x^\dagger) \cdot \ell(g^{(\text{adv})}(x^\dagger), y^{(\text{max})}(x^\dagger)), \end{aligned} \quad (75)$$

which is in contradiction with the fact that $g^{(\text{adv})}$ achieves the minimal value of Eq. (50). Therefore, we have $\ell(g^{(\text{adv})}(x^\dagger), y^{(\text{max})}(x^\dagger)) > \ell(g^{(\text{adv})}(x^\dagger), y^{(\text{adv})}(x^\dagger))$. Hence, Eq. (74) is positive. Multiplying Eq. (74) by $p(x)r^*(x)$ and taking an integral in \mathcal{S} , it is still positive because of Eq. (51). Thus, we have Eq. (73).

In conclusion, using $r_{\text{new}}(\cdot, \cdot)$ gives larger objective of Eq. (49) than $r^*(\cdot, \cdot)$. This contradicts the fact that $r^*(\cdot, \cdot)$ is the solution of Eq. (49). Therefore, $g^{(\text{adv})}$, which is obtained via ARM in Eq. (47), is Bayes optimal and coincides with g^* that is obtained via ordinary RM in Eq. (46).

So far, we showed that $g^{(\text{adv})}$, which is any solution of ARM, coincides with the Bayes optimal classifier almost surely over $q^*(x)$. We now turn our focus to the region

$$\mathcal{S}' \equiv \{x \mid x \in \mathcal{X}, p(x) > 0, q^*(x) = 0\}. \quad (76)$$

Because $g^{(\text{adv})}$ is chosen from all measurable functions, its function value at x is obtained in a point-wise manner:

$$g^{(\text{adv})}(x) = \arg \min_{\hat{y} \in \mathbb{R}^K} \sum_{y \in \mathcal{Y}} p(y|x) r^*(y|x) \ell(\hat{y}, y). \quad (77)$$

For $x \in \mathcal{S}'$, we immediately have $r^*(x) = q^*(x)/p(x) = 0$. Then, for all $x \in \mathcal{S}'$, we have $r^*(x, y) = 0$ for any $r(y|x), y \in \mathcal{Y}$. Therefore, for $x \in \mathcal{S}'$, we can *virtually* set $r^*(y|x) = 1$ for all $y \in \mathcal{Y}$. Substituting this to Eq. (77), we have

$$g^{(\text{adv})}(x) = \arg \min_{\hat{y} \in \mathbb{R}^K} \sum_{y \in \mathcal{Y}} p(y|x) \ell(\hat{y}, y), \quad \text{for } x \in \mathcal{S}'. \quad (78)$$

It follows from Eq. (78) and the use of the classification-calibrated loss that

$$\arg \max_y g_y^{(\text{adv})}(x) = \arg \max_y p(y|x), \quad \text{for } x \in \mathcal{S}'. \quad (79)$$

This particular $g^{(\text{adv})}$ coincides with the Bayes optimal classifier for all $x \in \mathcal{S}'$. Define

$$\begin{aligned} \mathcal{S}_{\text{diff}} &\equiv \left\{ x \mid x \in \mathcal{X}, p(x) > 0, \arg \max_{y \in \mathcal{Y}} p(y|x) \neq \arg \max_{y \in \mathcal{Y}} g_y^{(\text{adv})}(x) \right\}. \\ &= \underbrace{\mathcal{S} \cup \left\{ x \mid x \in \mathcal{X}, p(x) > 0, q^*(x) = 0, \arg \max_{y \in \mathcal{Y}} p(y|x) \neq \arg \max_{y \in \mathcal{Y}} g_y^{(\text{adv})}(x) \right\}}_{\equiv \mathcal{S}_{\text{diff}}}. \end{aligned} \quad (80)$$

Then we have

$$\begin{aligned} \int_{x \in \mathcal{S}_{\text{diff}}} p(x) dx &= \int_{x \in \mathcal{S}} p(x) dx + \underbrace{\int_{x \in \mathcal{S}'_{\text{diff}}} p(x) dx}_{=0 \quad \because \text{Eq. (79)}} \\ &= \int_{x \in \mathcal{S}} \frac{q^*(x)}{r^*(x)} dx \\ &\leq \frac{1}{\min_{x \in \mathcal{S}} r^*(x)} \int_{x \in \mathcal{S}} q^*(x) dx \\ &= 0. \end{aligned} \quad (81)$$

Therefore, the particular $g^{(\text{adv})}$ coincides with the Bayes optimal classifier almost surely over $p(x)$.

Finally, we consider binary classification, where the key differences to multi-class classification are that the prediction function is $\mathcal{X} \rightarrow \mathbb{R}$, $\mathcal{Y} = \{+1, -1\}$, and the prediction is made based on the sign of the output of the prediction function. Therefore, we need to replace all ‘argmax’ in the above proof with ‘sign’. In addition, to show $\ell(g^{(\text{adv})}(x^\dagger), y^{(\text{max})}(x^\dagger)) > \ell(g^{(\text{adv})}(x^\dagger), y^{(\text{adv})}(x^\dagger))$, we construct g' in Eq. (75) by $-g^{(\text{adv})}(x^\dagger)$. With these two modifications, all the arguments for multi-class classification hold for binary classification. \square

Remark 2. Here, we show that *if the KL divergence is used* for the f -divergence, the prediction of *any* solution $g^{(\text{adv})}$ of ARM coincides with that of the Bayes optimal classifier almost surely over $p(x)$.

In the following, assume the KL divergence is used. Given prediction function g , the density ratio put by the adversary becomes

$$r^*(x, y) = \frac{1}{Z(\gamma)} \exp\left(\frac{\ell(g(x), y)}{\gamma}\right), \quad (82)$$

where

$$Z(\gamma) = \mathbb{E}_{p(x,y)} \left[\exp \left(\frac{\ell(g(x), y)}{\gamma} \right) \right], \quad (83)$$

and γ is chosen so that the following equality holds:

$$\mathbb{E}_{p(x,y)} [r^*(x, y) \log r^*(x, y)] = \delta. \quad (84)$$

From Eq. (82), we see that $r^*(\cdot, \cdot)$ is a positive function for any g as long as $\ell(\cdot, \cdot)$ is bounded. Thus, $q^*(x) \equiv \sum_{y \in \mathcal{Y}} r^*(x, y) p(x, y)$ is also positive for $x \in \mathcal{X}$ such that $p(x) > 0$. On the other hand, because we assume $q \ll p$, $p(x) = 0$ implies $q(x) = 0$. Thus, we have $q^*(x) > 0$ iff $p(x) > 0$.

Now, assume that $q^*(x) > 0$ iff $p(x) > 0$, and $g^{(\text{adv})}$ coincides with the Bayes optimal classifier almost surely over $q^*(x)$. Then, we have

$$\int_{x \in \mathcal{S}} q^*(x) dx = 0, \quad (85)$$

where

$$\mathcal{S}' \equiv \left\{ x \mid x \in \mathcal{X}, p(x) > 0, q^*(x) > 0, \arg \max_{y \in \mathcal{Y}} p(y|x) \neq \arg \max_{y \in \mathcal{Y}} g_y^{(\text{adv})}(x) \right\}. \quad (86)$$

Because $r^*(x) \equiv q^*(x)/p(x)$ is positive, $\epsilon \equiv \min_{x \in \mathcal{S}} r^*(x)$ is also positive. Then, we have

$$\begin{aligned} \int_{x \in \mathcal{S}} p(x) dx &= \int_{x \in \mathcal{S}} \frac{q^*(x)}{r^*(x)} dx \\ &\leq \frac{1}{\epsilon} \int_{x \in \mathcal{S}} q^*(x) dx \\ &= 0. \end{aligned} \quad (87)$$

Therefore, $g^{(\text{adv})}$ coincides with the Bayes optimal classifier almost surely over the training density $p(x)$.

C. Proof of Lemma 1

By assumption, $\ell(\hat{y}, y)$ is a convex margin loss. Thus, we can let $\ell(\hat{y}, y) = \phi(y\hat{y})$, where $\phi(\cdot)$ is a convex function. On the other hand, by Definition 1, for some non-constant, non-decreasing and non-negative function $h(\cdot)$, the steeper loss $\ell_{\text{steep}}(\hat{y}, y)$, satisfies

$$\begin{aligned} \frac{\partial \ell_{\text{steep}}(\hat{y}, y)}{\partial \hat{y}} &= h(\ell(\hat{y}, y)) \frac{\partial \ell(\hat{y}, y)}{\partial \hat{y}} \\ &= h(\phi(y\hat{y})) \frac{\partial \phi(y\hat{y})}{\partial \hat{y}}. \end{aligned} \quad (88)$$

From Eq. (88), it is easy to see that $\ell_{\text{steep}}(\hat{y}, y)$ is also a margin loss and can be written as $\ell_{\text{steep}}(\hat{y}, y) = \phi_{\text{steep}}(y\hat{y})$. Our first goal is to show that $\phi_{\text{steep}}(\cdot)$ is convex. To this end, it is sufficient to show that $\frac{\partial \phi_{\text{steep}}(y\hat{y})}{\partial \hat{y}} = h(\phi(y\hat{y})) \frac{\partial \phi(y\hat{y})}{\partial \hat{y}}$ is non-decreasing in \hat{y} , $\mathbb{Y} = \{+1, -1\}$.

Since $\phi(y\hat{y})$ is convex in \hat{y} , $\frac{\partial \phi(y\hat{y})}{\partial \hat{y}}$ is non-decreasing in \hat{y} . Let \hat{y}_α be the smallest \hat{y}_α such that $\frac{\partial \phi(y\hat{y}_\alpha)}{\partial \hat{y}_\alpha} = 0$, if such \hat{y}_α exists. In the following, we analyze $\phi(y\hat{y}) \frac{\partial \phi(y\hat{y})}{\partial \hat{y}}$, considering two cases: 1) $\hat{y} \leq \hat{y}_\alpha$ and 2) $\hat{y}_\alpha \leq \hat{y}$. Note that \hat{y}_α may not always exist because $\frac{\partial \phi(y\hat{y})}{\partial \hat{y}}$ can be negative for any finite \hat{y} , which is the case for the widely-used classification losses such as the exponential loss and the logistic loss. In such a case, we only consider the first case, letting \hat{y}_α arbitrarily large.

Case 1 $\hat{y} \leq \hat{y}_\alpha$: By convexity of $\phi(y\hat{y})$, for $\hat{y} \leq \hat{y}_\alpha$, $\frac{\partial \phi(y\hat{y})}{\partial \hat{y}} \leq 0$ holds and therefore, $\phi(y\hat{y})$ is non-increasing in \hat{y} . Since $h(\cdot)$ is a non-decreasing function, $h(\phi(y\hat{y}))$ is non-increasing in \hat{y} for $\hat{y} \leq \hat{y}_\alpha$. In summary, for $\hat{y} \leq \hat{y}_\alpha$, $\frac{\partial \phi_{\text{steep}}(y\hat{y})}{\partial \hat{y}}$ is a non-positive

non-decreasing function of \hat{y} , and $h(\phi(y\hat{y}))$ is a non-negative non-increasing function of \hat{y} . Thus, for $\hat{y} \leq \hat{y}_\alpha$, their product $h(\phi(y\hat{y})) \frac{\phi(y\hat{y})}{\partial \hat{y}}$ is a non-decreasing function of \hat{y} .

Case 2 $\hat{y}_\alpha \leq \hat{y}$: By convexity of $\phi(y\hat{y})$, for $\hat{y}_\alpha \leq \hat{y}$, $\frac{\phi(y\hat{y})}{\partial \hat{y}} \geq 0$ holds and therefore, $\phi(y\hat{y})$ is non-decreasing in \hat{y} . Since $h(\cdot)$ is a non-decreasing function, $h(\phi(y\hat{y}))$ is non-decreasing in \hat{y} for $\hat{y} \leq \hat{y}_\alpha$. In summary, for $\hat{y}_\alpha \leq \hat{y}$, $\frac{\phi(y\hat{y})}{\partial \hat{y}}$ is a non-negative non-decreasing function of \hat{y} , and $h(\phi(y\hat{y}))$ is a non-negative non-decreasing function of \hat{y} . Thus, for $\hat{y} \leq \hat{y}_\alpha$, their product $h(\phi(y\hat{y})) \frac{\phi(y\hat{y})}{\partial \hat{y}}$ is a non-decreasing function of \hat{y} .

Therefore, for any \hat{y} , $h(\phi(y\hat{y})) \frac{\phi(y\hat{y})}{\partial \hat{y}}$ is a non-decreasing function of \hat{y} , which directly indicates that the steeper loss, $\phi_{\text{steep}}(y\hat{y})$, is convex.

We now utilize the fact that a convex margin loss $\psi(y\hat{y})$ is classification calibrated iff $\psi'(0) < 0$ [Theorem 6 in (Bartlett et al., 2006)]. Using this fact, because $\phi(y\hat{y})$ is classification calibrated, we have $\phi'(0) < 0$. Furthermore, from the assumption, we have $h(\phi(0)) > 0$. Therefore, we have $\phi'_{\text{steep}}(0) = h(\phi(0))\phi'(0) < 0$. Using the fact again, we immediately have that $\phi_{\text{steep}}(y\hat{y})$ is classification calibrated. \square

Remark 3. In the proof, we need to assume $h(\phi(0)) > 0$. From Appendix D, we know that $h(\ell)$ corresponds to the weight put by the adversary to data points with a loss value of ℓ . We see from Eq. (22) that when the KL divergence is used, the adversary will only assign positive weights to data losses. Therefore, Lemma 1 always holds when the KL divergence is used.

D. Proof of Theorem 3

Let θ^* be the stationary point of Eq. (7). By using a chain rule and Danskin's theorem (Danskin, 1966), θ^* satisfies

$$\frac{1}{N} \sum_{i=1}^N r_i^* \left. \frac{\partial \ell(\hat{y}, y_i)}{\partial \hat{y}} \right|_{\hat{y}=g_{\theta^*}(x_i)} \cdot \nabla_{\theta} g_{\theta}(x_i) \Big|_{\theta=\theta^*} \in \mathbf{0}, \quad (89)$$

where r^* is the solution of inner maximization of Eq. (7) at the stationary point.

Now, we analyze r^* , which is the solution of Eq. (7) at the stationary point θ^* . For notational convenience, for $1 \leq i \leq N$, let us denote $\ell_i(\theta^*)$ by ℓ_i^* . Then, r^* is the solution of the following optimization problem.

$$\max_{\mathbf{r} \in \mathcal{U}_f} \frac{1}{N} \sum_{i=1}^N r_i \ell_i^*, \quad (90)$$

$$\hat{\mathcal{U}}_f = \left\{ \mathbf{r} \mid \frac{1}{N} \sum_{i=1}^N f(r_i) \leq \delta, \frac{1}{N} \sum_{i=1}^N r_i = 1, \mathbf{r} \geq 0 \right\}, \quad (91)$$

Note that Eq. (90) is a convex optimization problem because it has a linear objective with a convex constraint; thus, any local maximum is the global maximum. Nonetheless, there can be multiple solutions that attain the same global maxima. Among those solutions, we now show that there exists r^* such that elements of r^* has the monotonic relationship to the corresponding data losses, i.e., for any $1 \leq j \neq k \leq N$,

$$\ell_j^* < \ell_k^* \Rightarrow 0 \leq r_j^* \leq r_k^*, \quad (92)$$

$$\ell_j^* = \ell_k^* \Rightarrow 0 \leq r_j^* = r_k^*. \quad (93)$$

To prove this, we assume we obtain one of optimal solutions of Eq. (90), which we denote as r'^* . If this r'^* satisfies Eqs. (92) and (93) for any j and k , then we are done. In the following, we assume r'^* does not satisfy either Eqs. (92) or (93).

First, assume that r'^* does not satisfy Eq. (92). Then, there exist $1 \leq j \neq k \leq N$ such that $\ell_j^* < \ell_k^*$ but $r_j'^* > r_k'^*$. Define r''^* such that

$$r_i''^* = \begin{cases} r_i'^* & \text{if } i \neq j, k \\ r_j'^* & \text{if } i = k \\ r_k'^* & \text{if } i = j \end{cases} \quad \text{for } 1 \leq i \leq N. \quad (94)$$

Then, it is easy to see $\mathbf{r}''^* \in \widehat{\mathcal{U}}_f$, and the following holds:

$$\begin{aligned}
 \frac{1}{N} \sum_{i=1}^N r_i''^* \ell_i^* - \frac{1}{N} \sum_{i=1}^N r_i'^* \ell_i^* &= \frac{1}{N} (r_j''^* \ell_j^* + r_k''^* \ell_k^* - r_j'^* \ell_j^* - r_k'^* \ell_k^*) \\
 &= \frac{1}{N} (r_k''^* \ell_j^* + r_j''^* \ell_k^* - r_j'^* \ell_j^* - r_k'^* \ell_k^*) \\
 &= \frac{1}{N} (r_j''^* - r_k''^*) (\ell_k^* - \ell_j^*) \\
 &> 0.
 \end{aligned} \tag{95}$$

Therefore, the newly defined \mathbf{r}''^* attains the larger objective value of Eq. (90), which contradicts the assumption that \mathbf{r}^* is the optimal solution of Eq. (90). Thus, \mathbf{r}^* always satisfies Eq. (92).

Second, assume that \mathbf{r}^* does not satisfy Eq. (93). Then, there exist $1 \leq j \neq k \leq N$ such that $\ell_j^* = \ell_k^*$ but $r_j^* \neq r_k^*$. Define \mathbf{r}''^* such that

$$r_i''^* = \begin{cases} r_i'^* & \text{if } i \neq j, k \\ (r_i'^* + r_j'^*)/2 & \text{if } i = j, k \end{cases} \quad \text{for } 1 \leq i \leq N. \tag{96}$$

Then, it is easy to see $\mathbf{r}''^* \in \widehat{\mathcal{U}}_f$ because

$$\begin{aligned}
 \frac{1}{N} \sum_{i=1}^N f(r_i''^*) &= \frac{1}{N} \left\{ \left(\sum_{i \neq j, k} f(r_i''^*) \right) + f(r_j''^*) + f(r_k''^*) \right\} \\
 &= \frac{1}{N} \left\{ \left(\sum_{i \neq j, k} f(r_i'^*) \right) + f((r_j'^* + r_k'^*)/2) + f((r_j'^* + r_k'^*)/2) \right\} \\
 &\leq \frac{1}{N} \left\{ \left(\sum_{i \neq j, k} f(r_i'^*) \right) + f(r_j'^*) + f(r_k'^*) \right\} \quad (\because \text{convexity of } f(\cdot).) \\
 &= \frac{1}{N} \sum_{i=1}^N f(r_i'^*) \\
 &\leq \delta. \quad (\because \mathbf{r}^* \in \widehat{\mathcal{U}}_f.)
 \end{aligned} \tag{97}$$

Also, it is easy to see that \mathbf{r}''^* attain the same maximum value as \mathbf{r}^* ; thus, \mathbf{r}''^* is the optimal solution of Eq. (90), and notably, we have $r_j''^* = r_k''^* = r_k'^*$ for $\ell_j^* = \ell_k^*$. In general, we can start from any $\mathbf{r}^* \in \widehat{\mathcal{U}}_f$ and equally distribute the weights to the same losses to obtain \mathbf{r}''^* , which is still in $\widehat{\mathcal{U}}_f$ and attains exactly the same global optimal value in Eq. (90).

In the following, we assume we have \mathbf{r}^* that satisfies Eqs (92) and (93) for any $1 \leq j \neq k \leq N$. Then, there exists a non-decreasing non-negative function $r^*(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$, such that

$$r^*(\ell_i^*) = r_i^*, \quad \text{for } 1 \leq i \leq N. \tag{98}$$

Let us construct a new loss function $\ell_{\text{DRSL}}(\widehat{y}, y)$ by its derivative:

$$\frac{\partial \ell_{\text{DRSL}}(\widehat{y}, y)}{\partial \widehat{y}} \equiv r^*(\ell(\widehat{y}, y)) \frac{\partial \ell(\widehat{y}, y)}{\partial \widehat{y}}. \tag{99}$$

Then, from Eqs. (89), (98) and (99), we immediately have

$$\frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_{\text{DRSL}}(\widehat{y}, y_i)}{\partial \widehat{y}} \Big|_{\widehat{y}=g_{\theta^*}(x_i)} \cdot \nabla_{\theta} g_{\theta}(x_i) \Big|_{\theta=\theta^*} \in \mathbf{0}. \tag{100}$$

This readily implies that θ^* is a stationary point of Eq. (13), i.e., ERM using $\ell_{\text{DRSL}}(\widehat{y}, y)$. Furthermore, from Eq. (99) and the non-negativeness and non-decreasingness of $r^*(\cdot)$, we see that the newly constructed loss, $\ell_{\text{DRSL}}(\widehat{y}, y)$, is steeper than the original loss, $\ell(\widehat{y}, y)$ (see Definition 1 for the definition of the steeper loss). Here we see that $h(\cdot)$ in Definition 1 exactly corresponds to $r^*(\cdot)$ defined in Eq. (98). \square

E. Derivation of the Decomposition of the Adversarial Risk

Here, we derive Eq. (18) for the PE divergence.

$$\begin{aligned}
 \mathcal{R}_{\text{s-adv}}(\theta) - \mathcal{R}(\theta) &\equiv \sup_{w \in \mathcal{W}_{\text{PE}}} \mathbb{E}_{p(x,y,z)} [\{w(z) - 1\} \ell(g_\theta(x), y)] \\
 &= \sup_{w \in \mathcal{W}_{\text{PE}}} \mathbb{E}_{p(z)} [\{w(z) - 1\} \mathbb{E}_{p(x,y|z)} [\ell(g_\theta(x), y)]] \\
 &= \sup_{w \in \mathcal{W}_{\text{PE}}} \sum_{z \in \mathcal{Z}} p(z) \{w(z) - 1\} \mathcal{R}_z(\theta).
 \end{aligned} \tag{101}$$

It follows from Eq. (23) that for $z \in \mathcal{Z}$, we have the adversarial weight as¹¹

$$w^*(z) = \sqrt{\frac{\delta}{\sum_{z' \in \mathcal{Z}} p(z') (\mathcal{R}_{z'}(\theta) - \mathcal{R}(\theta))^2}} (\mathcal{R}_z(\theta) - \mathcal{R}(\theta)) + 1. \tag{102}$$

Hence, Eq. (101) becomes

$$\begin{aligned}
 \sum_{z \in \mathcal{Z}} p(z) \{w^*(z) - 1\} \mathcal{R}_z(\theta) &= \sqrt{\frac{\delta}{\sum_{z' \in \mathcal{Z}} p(z') (\mathcal{R}_{z'}(\theta) - \mathcal{R}(\theta))^2}} \sum_{z \in \mathcal{Z}} p(z) (\mathcal{R}_z(\theta) - \mathcal{R}(\theta)) \mathcal{R}_z(\theta) \\
 &= \sqrt{\frac{\delta}{\sum_{z' \in \mathcal{Z}} p(z') (\mathcal{R}_{z'}(\theta) - \mathcal{R}(\theta))^2}} \sum_{z \in \mathcal{Z}} p(z) (\mathcal{R}_z(\theta) - \mathcal{R}(\theta))^2 \\
 &= \sqrt{\delta} \cdot \sqrt{\sum_{z \in \mathcal{Z}} p(z) (\mathcal{R}_z(\theta) - \mathcal{R}(\theta))^2},
 \end{aligned} \tag{103}$$

which concludes our derivation.

F. Comparison between the Use of Different f -divergences

We qualitatively compare the use of different f -divergences. For $1 \leq x$, the f functions for the PE, KL divergences are $(x - 1)^2$, $x \log x$, respectively. The function f in Eq. (20) penalizes the deviation of the adversarial weights from the uniform weights, $\mathbf{1}_S$. With the quadratic penalty of the PE divergence, it is hard for the adversary to concentrate large weights onto a small portion of latent categories. In contrast, when the KL divergence is used, the adversary tends to put large weights to a small portion of latent categories. Hence, users can choose the appropriate divergence depending on their belief on how concentrated the distribution shift occurs.

G. Formal Statement of the Convergence Rate

Denote by $p_z = p(z)$ and $w_z = w(z)$ for $z \in \mathcal{Z}$ and define a set-valued function $\Phi : \mathbb{R}^S \rightarrow 2^{\mathbb{R}^S}$ as

$$\Phi(\mathbf{u}) = \{\mathbf{w} \in \mathbb{R}^S \mid \sum_s (p_s + u_s) f(w_s) \leq \delta, \sum_s (p_s + u_s) w_s = 1, w_s \geq 0\}.$$

Then, $\mathcal{W}_f = \Phi(\mathbf{0})$ and $\widehat{\mathcal{W}}_f = \Phi(\mathbf{u})$ where $u_s = n_s/N - p_s$ for $s = 1, \dots, S$. Similarly, denote by $l_z = \mathbb{E}_{p(x,y)} [p(z \mid x, y) \ell(g_\theta(x), y)]$ and define a function $R_\theta : \mathbb{R}^S \rightarrow \mathbb{R}$ indexed by θ as

$$R_\theta(\mathbf{u}') = \sum_s w_s (l_s + u'_s).$$

Then, $\mathbb{E}_{p(x,y,z)} [w_z \ell(g_\theta(x), y)] = R_\theta(\mathbf{0})$ and $\widehat{\mathcal{R}}(\mathbf{w}, \theta) = R_\theta(\mathbf{u}')$ where $u'_s = n_s \bar{\ell}_s(\theta)/N - l_s$ for $s = 1, \dots, S$. Finally, the perturbed objective function can be defined by

$$J(\theta, \mathbf{u}, \mathbf{u}') = \sup_{\mathbf{w} \in \Phi(\mathbf{u})} R_\theta(\mathbf{u}') + \lambda(\mathbf{u}, \mathbf{u}') \Omega(\theta),$$

¹¹Here, we need to assume that δ is not so large. Then, we can validly drop the non-negativity inequality constraint of $\widehat{\mathcal{W}}_f$, which is needed to obtain the analytic solution in Eq. (23).

where the function $\lambda(\mathbf{u}, \mathbf{u}') \geq 0$ serves as the regularization parameter, so that the truly optimal θ^* is the minimizer of $J(\theta, \mathbf{0}, \mathbf{0})$ and the empirically optimal $\hat{\theta}$ is the minimizer of $J(\theta, \mathbf{u}, \mathbf{u}')$ with the aforementioned perturbations \mathbf{u} and \mathbf{u}' .

According to the central limit theorem (Chung, 1968), $u_s = \mathcal{O}_p(1/\sqrt{N})$, and $u'_s = \mathcal{O}_p(1/\sqrt{N})$ if the loss ℓ is finite. Therefore, we only consider perturbations \mathbf{u} and \mathbf{u}' such that $\|\mathbf{u}\|_2 \leq \epsilon$ and $\|\mathbf{u}'\|_2 \leq \epsilon$ in our analysis, where $0 < \epsilon \leq \delta/(5\sqrt{S}|f'(1)|)$ is a sufficiently small constant.

We make the following assumptions:

- (a) $g_\theta(x)$ is linear in θ , and for all θ under consideration, $\|\nabla_\theta g_\theta\|_\infty = \sup_x \|\nabla_\theta g_\theta(x)\|_2 < \infty$, which implies $\|g_\theta\|_\infty = \sup_x |g_\theta(x)| < \infty$;¹²
- (b) $\partial \ell(t, y)/\partial t$ is bounded from below and above for all t such that $|t| \leq \|g_\theta\|_\infty$;¹³
- (c) $f(t)$ is twice differentiable, and this second derivative is bounded from below by a positive number for all t such that $0 \leq t \leq \sup_{\|\mathbf{u}\|_2 \leq \epsilon} \sup_{\mathbf{w} \in \Phi(\mathbf{u})} \max_s w_s$;¹⁴
- (d) $\Omega(\theta)$ is Lipschitz continuous, and $\lambda(\mathbf{u}, \mathbf{u}')$ converges to $\lambda(\mathbf{0}, \mathbf{0})$ in $\mathcal{O}(\|\mathbf{u}\|_2 + \|\mathbf{u}'\|_2)$.

We also assume either one of the two conditions holds:

- (e1) $\Omega(\theta)$ is strongly convex in θ and $\lambda(\mathbf{0}, \mathbf{0}) > 0$;
- (e2) $\ell(t, y)$ is twice differentiable w.r.t. t , and $\partial^2 \ell(t, y)/\partial t^2$ is lower bounded by a positive number for all t such that $|t| \leq \|g_{\theta^*}\|_\infty$. If t is vector-valued, $\partial^2 \ell(t, y)/\partial t_i^2$ is lower bounded for all dimensions of t such that $\|t\|_\infty \leq \sup_x \|g_{\theta^*}(x)\|_\infty$.¹⁵

Theorem 5 (Perturbation analysis). *Assume (a), (b), (c), (d), and (e1) or (e2). Let θ^* be the minimizer of $J(\theta, \mathbf{0}, \mathbf{0})$ and $\theta_{\mathbf{u}, \mathbf{u}'}$ be the minimizer of $J(\theta, \mathbf{u}, \mathbf{u}')$. Then, for all \mathbf{u} and \mathbf{u}' such that $\|\mathbf{u}\|_2 \leq \epsilon$ and $\|\mathbf{u}'\|_2 \leq \epsilon$,*

$$\begin{aligned} \|\theta_{\mathbf{u}, \mathbf{u}'} - \theta^*\|_2 &= \mathcal{O}(\|\mathbf{u}\|_2^{1/2} + \|\mathbf{u}'\|_2), \\ \|J(\theta_{\mathbf{u}, \mathbf{u}'}, \mathbf{0}, \mathbf{0}) - J(\theta^*, \mathbf{0}, \mathbf{0})\|_2 &= \mathcal{O}(\|\mathbf{u}\|_2^{1/2} + \|\mathbf{u}'\|_2). \end{aligned}$$

The convergence rate of the model parameter and the order of the estimation error are immediate corollaries of Theorem 5.

Theorem 6 (Convergence rate and estimation error). *Assume (a), (b), (c), (d), and (e1) or (e2). Let θ^* be the minimizer of the adversarial expected risk and $\hat{\theta}_N$ be the minimizer of the adversarial empirical risk given some training data of size N . Then, as $N \rightarrow \infty$,*

$$\|\hat{\theta}_N - \theta^*\|_2 = \mathcal{O}(N^{-1/4}),$$

and

$$\left\| \mathcal{R}_{\text{s-adv}}(\hat{\theta}_N) - \mathcal{R}_{\text{s-adv}}(\theta^*) \right\|_2 = \mathcal{O}(N^{-1/4})$$

H. Proof of the Convergence Rate

We begin with the growth condition of $J(\theta, \mathbf{0}, \mathbf{0})$ at $\theta = \theta^*$.

Lemma 2 (Second-order growth condition). *There exists a constant $C_{J''} > 0$ such that*

$$J(\theta, \mathbf{0}, \mathbf{0}) \geq J(\theta^*, \mathbf{0}, \mathbf{0}) + C_{J''} \|\theta - \theta^*\|_2^2.$$

Proof. First consider the assumption (e1). Let $C_{J''} = (1/2)\lambda(\mathbf{0}, \mathbf{0})$, so that $J(\theta, \mathbf{0}, \mathbf{0})$ is strongly convex with parameter $C_{J''}$, i.e.,

$$J(\theta, \mathbf{0}, \mathbf{0}) \geq J(\theta^*, \mathbf{0}, \mathbf{0}) + \nabla_\theta J(\theta^*, \mathbf{0}, \mathbf{0})^\top (\theta - \theta^*) + C_{J''} \|\theta - \theta^*\|_2^2.$$

¹²This makes $\ell(g_\theta(x), y)$ convex in θ for all $\ell(t, y)$ convex in t and $\nabla_\theta^2 \ell(g_\theta(x), y)$ easy to handle.

¹³This is a sufficient condition for the Lipschitz continuity of ℓ . In fact, it must be valid given (a) since ℓ is continuously differentiable w.r.t. t .

¹⁴This is for the Lipschitz continuity of $J(\theta, \mathbf{u}, \mathbf{u}') - J(\theta, \mathbf{0}, \mathbf{0})$. It is satisfied by the KL divergence since $f''(t) = 1/t$ and $\Phi(\mathbf{u})$ is bounded and then $\sup_{\|\mathbf{u}\|_2 \leq \epsilon} \sup_{\mathbf{w} \in \Phi(\mathbf{u})} \max_s w_s < \infty$, and by the PE divergence since $f''(t) = 2$

¹⁵This makes $J(\theta, \mathbf{0}, \mathbf{0})$ locally strongly convex in θ around θ^* . It is satisfied by the logistic loss with the lower bound as $1/(2 + \exp(\|g_{\theta^*}\|_\infty) + \exp(-\|g_{\theta^*}\|_\infty))$ and the softmax cross-entropy loss with the lower bound as $\min_y \min_{t_i = \pm \sup_x \|g_{\theta^*}\|_\infty} \exp(t_y) \sum_{i \neq y} \exp(t_i) / (\sum_i \exp(t_i))^2$.

The lemma follows from the optimality condition which says $\nabla_{\theta} J(\theta^*, \mathbf{0}, \mathbf{0}) = \mathbf{0}$.

Second consider the assumption (e2) if (e1) does not hold. Without loss of generality, assume that $\lambda(\mathbf{0}, \mathbf{0}) = 0$. Let $w^* = \arg \sup_{w \in \Phi(\mathbf{0})} R_{\theta^*}(\mathbf{0})$, then according to Danskin's theorem (Danskin, 1966),

$$\begin{aligned} \nabla_{\theta} J(\theta^*, \mathbf{0}, \mathbf{0}) &= \mathbb{E}_{p(x,y,z)} [w_z^* \nabla_{\theta} \ell(g_{\theta^*}(x), y)] \\ &= \mathbb{E}_{p(x,y,z)} [w_z^* \ell'(g_{\theta^*}(x), y) \nabla_{\theta} g_{\theta^*}(x)] \end{aligned}$$

where $\ell'(g_{\theta^*}(x), y)$ means $\partial \ell(t, y) / \partial t|_{t=g_{\theta^*}(x)}$. The assumption (a) guarantees that $\nabla_{\theta} g_{\theta^*}(x)$ is no longer a function of θ , and thus

$$\begin{aligned} \nabla_{\theta}^2 J(\theta^*, \mathbf{0}, \mathbf{0}) &= \mathbb{E}_{p(x,y,z)} [w_z^* \nabla_{\theta} \ell'(g_{\theta^*}(x), y) \nabla_{\theta} g_{\theta^*}(x)^{\top}] \\ &= \mathbb{E}_{p(x,y,z)} [w_z^* \ell''(g_{\theta^*}(x), y) \nabla_{\theta} g_{\theta^*}(x) \nabla_{\theta} g_{\theta^*}(x)^{\top}] \end{aligned}$$

where $\ell''(g_{\theta^*}(x), y)$ means $\partial^2 \ell(t, y) / \partial t^2|_{t=g_{\theta^*}(x)}$.

Let $C_{\ell''} = \inf_{|t| \leq \|g_{\theta^*}\|_{\infty}} \min_y \ell''(t, y)$, and by assumption $C_{\ell''} > 0$. Also let $C_{\lambda,z}$ be the smallest eigenvalue of $\mathbb{E}_{p(x|z)} [\nabla_{\theta} g_{\theta}(x) \nabla_{\theta} g_{\theta}(x)^{\top}]$ at $\theta = \theta^*$ for $z \in \mathcal{Z}$. Note that $p(x | z)$ generates infinite number of x , and $\mathbb{E}_{p(x|z)} [\nabla_{\theta} g_{\theta}(x) \nabla_{\theta} g_{\theta}(x)^{\top}]$ as an average of infinitely many independent positive semi-definite matrices $\nabla_{\theta} g_{\theta}(x) \nabla_{\theta} g_{\theta}(x)^{\top}$ (they are independent as long as $\nabla_{\theta} g_{\theta}(x)$ depends on x) is positive definite. Thus, $C_{\lambda,z} > 0$ for all $z \in \mathcal{Z}$, and subsequently,

$$\begin{aligned} &(\theta - \theta^*)^{\top} \nabla_{\theta}^2 J(\theta^*, \mathbf{0}, \mathbf{0}) (\theta - \theta^*) \\ &\geq \left(\inf_{|t| \leq \|g_{\theta^*}\|_{\infty}} \min_y \ell''(t, y) \right) \cdot (\theta - \theta^*)^{\top} \mathbb{E}_{p(x,y,z)} [w_z^* \nabla_{\theta} g_{\theta^*}(x) \nabla_{\theta} g_{\theta^*}(x)^{\top}] (\theta - \theta^*) \\ &= C_{\ell''} \cdot (\theta - \theta^*)^{\top} \mathbb{E}_{p(x,z)} [w_z^* \nabla_{\theta} g_{\theta^*}(x) \nabla_{\theta} g_{\theta^*}(x)^{\top}] (\theta - \theta^*) \\ &= C_{\ell''} \cdot (\theta - \theta^*)^{\top} \left(\sum_{s=1}^S p_s w_s^* \mathbb{E}_{p(x|z=s)} [\nabla_{\theta} g_{\theta^*}(x) \nabla_{\theta} g_{\theta^*}(x)^{\top}] \right) (\theta - \theta^*) \\ &= C_{\ell''} \left(\sum_{s=1}^S p_s w_s^* (\theta - \theta^*)^{\top} \mathbb{E}_{p(x|z=s)} [\nabla_{\theta} g_{\theta^*}(x) \nabla_{\theta} g_{\theta^*}(x)^{\top}] (\theta - \theta^*) \right) \\ &\geq C_{\ell''} \left(\sum_{s=1}^S p_s w_s^* C_{\lambda,s} \|\theta - \theta^*\|_2^2 \right) \\ &\geq C_{\ell''} \min_s C_{\lambda,s} \left(\sum_{s=1}^S p_s w_s^* \right) \|\theta - \theta^*\|_2^2 \\ &= C_{\ell''} \min_s C_{\lambda,s} \|\theta - \theta^*\|_2^2. \end{aligned}$$

This completes the proof by letting $C_{J''} = C_{\ell''} \min_s C_{\lambda,s}$. □

We then study the Lipschitz continuity of $J(\theta, \mathbf{u}, \mathbf{u}')$.

Lemma 3 (Lipschitz continuity of the perturbed objective). *For all \mathbf{u} and \mathbf{u}' such that $\|\mathbf{u}\|_2 \leq \epsilon$ and $\|\mathbf{u}'\|_2 \leq \epsilon$, $J(\theta, \mathbf{u}, \mathbf{u}')$ is Lipschitz continuous with a (not necessarily the best) Lipschitz constant independent of \mathbf{u} and \mathbf{u}' .*

Proof. Define $F(\theta, \mathbf{u}, \mathbf{u}') = \sup_{w \in \Phi(\mathbf{u})} R_{\theta}(\mathbf{u}')$ and let $w^* = \arg \sup_{w \in \Phi(\mathbf{u})} R_{\theta}(\mathbf{u}')$. According to Danskin's theorem (Danskin, 1966), $\nabla_{\theta} F(\theta, \mathbf{u}, \mathbf{u}') = \sum_s w_s^* \nabla_{\theta} l_s$ where

$$\nabla_{\theta} l_s = \mathbb{E}_{p(x,y)} [p(z=s | x, y) \ell'(g_{\theta}(x), y) \nabla_{\theta} g_{\theta}(x)].$$

The assumptions (a) and (b) say that $\|\nabla_{\theta} g_{\theta}\|_{\infty} < \infty$ and $|\ell'(g_{\theta}(x), y)| < \infty$ so that

$$\begin{aligned} \|\nabla_{\theta} l_s\|_2 &\leq \|\nabla_{\theta} g_{\theta}\|_{\infty} \left(\sup_{|t| \leq \|g_{\theta}\|_{\infty}} \max_y |\ell'(t, y)| \right) \mathbb{E}_{p(x, y)} [p(z = s \mid x, y)] \\ &= \|\nabla_{\theta} g_{\theta}\|_{\infty} \left(\sup_{|t| \leq \|g_{\theta}\|_{\infty}} \max_y |\ell'(t, y)| \right) p_s \\ &< \infty, \end{aligned}$$

and it is clear that $w_s^* < \infty$. Hence,

$$\|\nabla_{\theta} F(\theta, \mathbf{u}, \mathbf{u}')\|_2 \leq \sum_{s=1}^S w_s^* \|\nabla_{\theta} l_s\|_2 < \infty,$$

which means $F(\theta, \mathbf{u}, \mathbf{u}')$ is Lipschitz continuous with a Lipschitz constant independent of \mathbf{u} and \mathbf{u}' .

By the assumption (d), $\Omega(\theta)$ is Lipschitz continuous and there exists a constant $C_{\lambda} > 0$ such that

$$\begin{aligned} \lambda(\mathbf{u}, \mathbf{u}') &\leq \lambda(\mathbf{0}, \mathbf{0}) + C_{\lambda} (\|\mathbf{u}\|_2 + \|\mathbf{u}'\|_2) \\ &\leq \lambda(\mathbf{0}, \mathbf{0}) + 2C_{\lambda} \epsilon \\ &< \infty. \end{aligned}$$

As a result, $\lambda(\mathbf{u}, \mathbf{u}')\Omega(\theta)$ possesses a Lipschitz constant independent of \mathbf{u} and \mathbf{u}' as well. □

From now on, we investigate the Lipschitz continuity of the difference function

$$D(\theta) = J(\theta, \mathbf{u}, \mathbf{u}') - J(\theta, \mathbf{0}, \mathbf{0}),$$

which is the most challenging task in our perturbation analysis. Define

$$\begin{aligned} D_1(\theta) &= F(\theta, \mathbf{u}, \mathbf{u}') - F(\theta, \mathbf{u}, \mathbf{0}), \\ D_2(\theta) &= F(\theta, \mathbf{u}, \mathbf{0}) - F(\theta, \mathbf{0}, \mathbf{0}), \end{aligned}$$

where $F(\theta, \mathbf{u}, \mathbf{u}') = \sup_{\mathbf{w} \in \Phi(\mathbf{u})} R_{\theta}(\mathbf{u}')$ defined in Lemma 3, and then $D(\theta)$ can be decomposed as

$$D(\theta) = D_1(\theta) + D_2(\theta) + (\lambda(\mathbf{u}, \mathbf{u}') - \lambda(\mathbf{0}, \mathbf{0}))\Omega(\theta).$$

Given the assumption (d), the third function $(\lambda(\mathbf{u}, \mathbf{u}') - \lambda(\mathbf{0}, \mathbf{0}))\Omega(\theta)$ is Lipschitz continuous with a Lipschitz constant of order $\mathcal{O}(\|\mathbf{u}\|_2 + \|\mathbf{u}'\|_2)$. We are going to prove the same property for $D_1(\theta)$ and $D_2(\theta)$ using the assumptions (a), (b) and (c).

Lemma 4 (Lipschitz continuity of the difference function, I). *For any fixed \mathbf{u} and all \mathbf{u}' such that $\|\mathbf{u}'\|_2 \leq \epsilon$, $D_1(\theta)$ is Lipschitz continuous with a Lipschitz constant of order $\mathcal{O}(\|\mathbf{u}'\|_2)$.*

Proof. According to the chain rule in calculus,

$$\begin{aligned} \|\nabla_{\theta} D_1(\theta)\|_2 &= \left\| \sum_{s=1}^S \frac{\partial D_1(\theta)}{\partial l_s} \nabla_{\theta} l_s \right\|_2 \\ &\leq \left| \sum_{s=1}^S \frac{\partial D_1(\theta)}{\partial l_s} \right| \cdot \max_s \|\nabla_{\theta} l_s\|_2 \\ &= \mathcal{O} \left(\left| \sum_{s=1}^S \frac{\partial D_1(\theta)}{\partial l_s} \right| \right), \end{aligned}$$

since we have proven that $\|\nabla_{\theta} l_s\|_2 < \infty$ given the assumptions (a) and (b) in Lemma 3.

By definition,

$$\begin{aligned} D_1(\theta) &= \sup_{\mathbf{w} \in \Phi(\mathbf{u})} R_\theta(\mathbf{u}') - \sup_{\mathbf{w} \in \Phi(\mathbf{u})} R_\theta(\mathbf{0}) \\ &= \sup_{\mathbf{w} \in \Phi(\mathbf{u})} \sum_{s=1}^S w_s(l_s + u'_s) - \sup_{\mathbf{w} \in \Phi(\mathbf{u})} \sum_{s=1}^S w_s l_s. \end{aligned}$$

Let $\mathbf{w}^* = \arg \sup_{\mathbf{w} \in \Phi(\mathbf{u})} \sum_{s=1}^S w_s l_s$ and $\mathbf{v}^* = \arg \sup_{\mathbf{w} \in \Phi(\mathbf{u})} \sum_{s=1}^S w_s(l_s + u'_s)$, then according to Danskin's theorem (Danskin, 1966), $\partial D_1(\theta)/\partial l_s = v_s^* - w_s^*$ and

$$\begin{aligned} \left| \sum_{s=1}^S \frac{\partial D_1(\theta)}{\partial l_s} \right| &\leq \sum_{s=1}^S |v_s^* - w_s^*| \\ &\leq \sqrt{S} \|\mathbf{v}^* - \mathbf{w}^*\|_2, \end{aligned}$$

which means $\mathcal{O}(\|\nabla_\theta D_1(\theta)\|_2) = \mathcal{O}(\|\mathbf{v}^* - \mathbf{w}^*\|_2)$.

Consider the perturbation analysis of the following optimization problem

$$\min_{\mathbf{w}} - \sum_{s=1}^S w_s(l_s + u'_s) \quad \text{s.t. } \mathbf{w} \in \Phi(\mathbf{u}), \quad (104)$$

whose objective is perturbed and feasible region is unperturbed. Let

$$\begin{aligned} L(\mathbf{w}, \alpha, \alpha', \mathbf{u}') &= - \sum_{s=1}^S w_s(l_s + u'_s) + \alpha \left(\sum_{s=1}^S (p_s + u_s) f(w_s) - \delta \right) \\ &\quad + \alpha' \left(\sum_{s=1}^S (p_s + u_s) w_s - 1 \right) \end{aligned}$$

be the Lagrangian function, where $\alpha \geq 0$ and α' are Lagrange multipliers, and for simplicity the nonnegative constraints are omitted. Note that given the assumption (c), if $\alpha \neq 0$,

$$\frac{\partial^2}{\partial w_i \partial w_j} L(\mathbf{w}, \alpha, \alpha', \mathbf{u}') = \begin{cases} \alpha f''(w_i) > 0, & i = j, \\ 0, & i \neq j, \end{cases}$$

namely, $L(\mathbf{w}, \alpha, \alpha', \mathbf{u}')$ is locally strongly convex in \mathbf{w} . Thus,

- if $\alpha^* > 0$, the second-order sufficient condition (see Definition 6.2 in (Bonnans & Shapiro, 1998)) holds at \mathbf{w}^* that implies the corresponding second-order growth condition according to Theorem 6.3 in (Bonnans & Shapiro, 1998);
- if $\alpha^* = 0$, (104) is locally a standard linear programming around \mathbf{w}^* and it is fairly easy to see $\|\mathbf{v}^* - \mathbf{w}^*\|_2 = \mathcal{O}(\|\mathbf{u}'\|_2)$ according to Theorem 1 in (Robinson, 1977).

In the former case, it is obvious that for (104),

- the objective $-\sum_{s=1}^S w_s(l_s + u'_s)$ is Lipschitz continuous with a Lipschitz constant $\|\mathbf{l}\|_2 + \epsilon$ independent of \mathbf{u}' ;
- the difference function $-\sum_{s=1}^S w_s u'_s$ is Lipschitz continuous with a Lipschitz constant of order $\mathcal{O}(\|\mathbf{u}'\|_2)$.

Therefore, $\|\mathbf{v}^* - \mathbf{w}^*\|_2 = \mathcal{O}(\|\mathbf{u}'\|_2)$ by applying Proposition 6.1 in (Bonnans & Shapiro, 1998). \square

In order to prove the same property for $D_2(\theta)$, we need several lemmas.

Lemma 5. Denote by $f'(\mathbf{w}) = (f'(w_1), \dots, f'(w_S))^\top$. There exists a constant $C_{\cos} > 0$, such that $\cos(\mathbf{d} \circ f'(\mathbf{w}), \mathbf{d}) \leq 1 - C_{\cos}$ for all \mathbf{u} satisfying $\|\mathbf{u}\|_2 \leq \epsilon$, \mathbf{w} satisfying $\sum_{s=1}^S (p_s + u_s) f(w_s) = \delta$ and $\sum_{s=1}^S (p_s + u_s) w_s = 1$, and $\mathbf{d} > \mathbf{0}$.

Proof. Suppose the lemma is false, i.e., for any sufficiently large n , there exists some \mathbf{w}_n such that $\cos(\mathbf{d} \circ f'(\mathbf{w}_n), \mathbf{d}) =$

$1 - 1/(2n^2)$. Let $\zeta_n = \|\mathbf{d} \circ f'(\mathbf{w}_n)\|_2$ and $\eta_n = \zeta_n/\|\mathbf{d}\|_2$, then

$$\begin{aligned} \|\mathbf{d} \circ f'(\mathbf{w}_n) - \eta_n \mathbf{d}\|_2^2 &= \|\mathbf{d} \circ f'(\mathbf{w}_n)\|_2^2 + \eta_n^2 \|\mathbf{d}\|_2^2 - 2\eta_n (\mathbf{d} \circ f'(\mathbf{w}_n))^\top \mathbf{d} \\ &= 2\zeta_n^2 - 2\eta_n \cos(\mathbf{d} \circ f'(\mathbf{w}_n), \mathbf{d}) \|\mathbf{d} \circ f'(\mathbf{w}_n)\|_2 \|\mathbf{d}\|_2 \\ &= 2\zeta_n^2 - 2(1 - 1/(2n^2))\zeta_n^2 \\ &= \zeta_n^2/n^2. \end{aligned}$$

In other words, for $s = 1, \dots, S$,

$$\begin{aligned} |f'(w_{n,s}) - \eta_n| &= |d_s f'(w_{n,s}) - \eta_n d_s|/d_s \\ &\leq \|\mathbf{d} \circ f'(\mathbf{w}_n) - \eta_n \mathbf{d}\|_2/d_s \\ &= (\zeta_n/n)/d_s \\ &\leq \zeta'_n/n, \end{aligned}$$

where $\zeta'_n = \zeta_n/\min_s d_s$. Consequently, for any $1 \leq i, j \leq S$ and $i \neq j$,

$$\begin{aligned} |f'(w_{n,i}) - f'(w_{n,j})| &\leq |f'(w_{n,i}) - \eta_n| + |f'(w_{n,j}) - \eta_n| \\ &\leq 2\zeta'_n/n. \end{aligned}$$

Let $C_{f''} > 0$ be the lower bound of $f''(t)$ mentioned in the assumption (c). This assumption also guarantees that $f'(t)$ is continuous, and by the mean value theorem, there is some t between $w_{n,i}$ and $w_{n,j}$ such that

$$\begin{aligned} |w_{n,i} - w_{n,j}| &= \left| \frac{f'(w_{n,i}) - f'(w_{n,j})}{f''(t)} \right| \\ &\leq 2\zeta'_n/(C_{f''}n). \end{aligned}$$

Let $\eta'_n = \sum_s w_{n,s}/S$, then $|w_{n,s} - \eta'_n| \leq 2\zeta'_n/(C_{f''}n)$ for $s = 1, \dots, S$.

Recall that $\sum_s (p_s + u_s)w_{n,s} = 1$, then

$$\begin{aligned} \left(1 + \sum_{s=1}^S u_s\right) \eta'_n - 1 &= \sum_{s=1}^S (p_s + u_s) \eta'_n - \sum_{s=1}^S (p_s + u_s) w_{n,s} \\ &= \sum_{s=1}^S (p_s + u_s) (\eta'_n - w_{n,s}), \end{aligned}$$

and hence

$$\begin{aligned} \left| \left(1 + \sum_{s=1}^S u_s\right) \eta'_n - 1 \right| &\leq \frac{2\zeta'_n}{C_{f''}n} \sum_{s=1}^S (p_s + u_s) \\ &= \frac{2\zeta'_n}{C_{f''}n} \left(1 + \sum_{s=1}^S u_s\right). \end{aligned}$$

This ensures $|\eta'_n - 1/(1 + \sum_s u_s)| \leq 2\zeta'_n/(C_{f''}n)$ and implies $|w_{n,s} - 1/(1 + \sum_s u_s)| \leq 4\zeta'_n/(C_{f''}n)$ for $s = 1, \dots, S$. Since $f(t)$ is twice differentiable and $\|\mathbf{w}_n\|_2 < \infty$, we must have $\zeta'_n < \infty$ and then $\lim_{n \rightarrow \infty} w_{n,s} = 1/(1 + \sum_s u_s)$ for all $s = 1, \dots, S$.

The Taylor expansion of $f(t)$ at $t = 1$ is $f(t) = f'(1)(t - 1) + \mathcal{O}((t - 1)^2)$ since $f(1) = 0$, and if $t = 1/(1 + \sum_s u_s)$,

$$f\left(\frac{1}{1 + \sum_{s=1}^S u_s}\right) = -f'(1) \cdot \frac{\sum_{s=1}^S u_s}{1 + \sum_{s=1}^S u_s} + \mathcal{O}\left(\left(\frac{\sum_{s=1}^S u_s}{1 + \sum_{s=1}^S u_s}\right)^2\right).$$

When $\|\mathbf{u}\|_2 \leq \epsilon$, $|\sum_s u_s| \leq \|\mathbf{u}\|_1 \leq \sqrt{S}\epsilon$, and

$$\begin{aligned} f\left(\frac{1}{1 + \sum_{s=1}^S u_s}\right) &\leq |f'(1)| \frac{\sqrt{S}\epsilon}{1 - \sqrt{S}\epsilon} + \mathcal{O}(\epsilon^2) \\ &\leq 2\sqrt{S}|f'(1)|\epsilon. \end{aligned}$$

As a result,

$$\begin{aligned} \lim_{n \rightarrow \infty} \sum_{s=1}^S (p_s + u_s) f(w_{n,s}) &= \sum_{s=1}^S (p_s + u_s) f\left(\frac{1}{1 + \sum_{s=1}^S u_s}\right) \\ &= \left(1 + \sum_{s=1}^S u_s\right) f\left(\frac{1}{1 + \sum_{s=1}^S u_s}\right) \\ &\leq (1 + \sqrt{S}\epsilon) \cdot 2\sqrt{S}|f'(1)|\epsilon \\ &\leq 4\sqrt{S}|f'(1)|\epsilon. \end{aligned}$$

However, this is impossible since $\sum_s (p_s + u_s) f(w_{n,s}) = \delta \geq 5\sqrt{S}|f'(1)|\epsilon$. \square

Based on Lemma 5, we derive the convergence rate of $\Phi(\mathbf{u})$ to $\Phi(\mathbf{0})$.

Lemma 6. Let $d_{\mathcal{H}}(V, W)$ be the Hausdorff distance between two sets V and W :

$$d_{\mathcal{H}}(V, W) = \max \left\{ \sup_{\mathbf{v} \in V} \inf_{\mathbf{w} \in W} \|\mathbf{v} - \mathbf{w}\|_2, \sup_{\mathbf{w} \in W} \inf_{\mathbf{v} \in V} \|\mathbf{v} - \mathbf{w}\|_2 \right\}.$$

Then $d_{\mathcal{H}}(\Phi(\mathbf{u}), \Phi(\mathbf{0})) = \mathcal{O}(\|\mathbf{u}\|_2)$ for all \mathbf{u} satisfying $\|\mathbf{u}\|_2 \leq \epsilon$.

Proof. We are going to prove $\sup_{\mathbf{w} \in \Phi(\mathbf{0})} \inf_{\mathbf{v} \in \Phi(\mathbf{u})} \|\mathbf{v} - \mathbf{w}\|_2 = \mathcal{O}(\|\mathbf{u}\|_2)$, and the other direction can be proven similarly.

Pick an arbitrary $\mathbf{w}_0 \in \Phi(\mathbf{0})$. Let $\beta = \delta / (\delta + \|f(\mathbf{w}_0)\|_2 \|\mathbf{u}\|_2)$ and consider $\mathbf{v}_1 = \beta \mathbf{w}_0 + (1 - \beta) \mathbf{1}$,

$$\begin{aligned} \|\mathbf{v}_1 - \mathbf{w}_0\|_2 &= \|(\beta - 1)\mathbf{w}_0 + (1 - \beta)\mathbf{1}\|_2 \\ &= (1 - \beta) \|\mathbf{w}_0 - \mathbf{1}\|_2 \\ &\leq \|f(\mathbf{w}_0)\|_2 \|\mathbf{u}\|_2 \|\mathbf{w}_0 - \mathbf{1}\|_2 / \delta \\ &= \mathcal{O}(\|\mathbf{u}\|_2). \end{aligned}$$

Moreover,

$$\begin{aligned} \sum_{s=1}^S (p_s + u_s) f(v_{1,s}) &= \sum_{s=1}^S (p_s + u_s) f(\beta w_{0,s} + (1 - \beta)) \\ &\leq \sum_{s=1}^S (p_s + u_s) (\beta f(w_{0,s}) + (1 - \beta) f(1)) \\ &= \beta \left(\sum_{s=1}^S p_s f(w_{0,s}) + \sum_{s=1}^S u_s f(w_{0,s}) \right) \\ &\leq \beta (\delta + \|f(\mathbf{w}_0)\|_2 \|\mathbf{u}\|_2) \\ &= \delta, \end{aligned}$$

where the second line is due to the convexity of $f(t)$, the third line is because $f(1) = 0$, and the fourth line is according to Jensen's inequality. This means \mathbf{v}_1 belongs to the set $V_1 = \{\mathbf{w} \in \mathbb{R}^S \mid \sum_s (p_s + u_s) f(w_s) \leq \delta, w_s \geq 0\}$.

However, \mathbf{v}_1 does not belong to the set $V_2 = \{\mathbf{w} \in \mathbb{R}^S \mid \sum_s (p_s + u_s)w_s = 1, w_s \geq 0\}$. Since V_2 is a hyperplane, we can easily project \mathbf{v}_1 onto V_2 to obtain \mathbf{v}_2 , and

$$\begin{aligned}
 \|\mathbf{v}_2 - \mathbf{v}_1\|_2 &= \frac{1}{\|\mathbf{p} + \mathbf{u}\|_2} \left| \sum_{s=1}^S (p_s + u_s)v_{1,s} - 1 \right| \\
 &\leq \frac{2}{\|\mathbf{p}\|_2} \left| \sum_{s=1}^S (p_s + u_s)(\beta w_{0,s} + 1 - \beta) - 1 \right| \\
 &\leq \sqrt{4S} \left| \beta \sum_{s=1}^S p_s w_{0,s} + \beta \sum_{s=1}^S u_s w_{0,s} + (1 - \beta) \sum_{s=1}^S p_s + (1 - \beta) \sum_{s=1}^S u_s - 1 \right| \\
 &= \sqrt{4S} \left| \beta + \beta \sum_{s=1}^S u_s w_{0,s} + (1 - \beta) + (1 - \beta) \sum_{s=1}^S u_s - 1 \right| \\
 &\leq \sqrt{4S} \beta \left| \sum_{s=1}^S u_s w_{0,s} \right| + \mathcal{O}(\|\mathbf{u}\|_2^2) \\
 &= \mathcal{O}(\|\mathbf{u}\|_2).
 \end{aligned}$$

After this projection, $\mathbf{v}_2 \notin V_1$ again.

Let \mathbf{v}_3 be the projection of \mathbf{v}_2 onto $\Phi(\mathbf{u}) = V_1 \cap V_2$, $T_w V_1(\mathbf{v}_3)$ be the tangent hyperplane to V_1 at \mathbf{v}_3 of $S - 1$ dimensions, and $T_w(V_1 \cap V_2)(\mathbf{v}_3)$ be that to $V_1 \cap V_2$ at \mathbf{v}_3 of $S - 2$ dimensions. As a consequence, $\mathbf{v}_3 - \mathbf{v}_2 \in V_2$ is one of normal vectors to $T_w(V_1 \cap V_2)(\mathbf{v}_3)$ at \mathbf{v}_3 , which is also the projection of the normal vector to $T_w V_1(\mathbf{v}_3)$ at \mathbf{v}_3 onto V_2 . This means the normal vector to $T_w V_1(\mathbf{v}_3)$ at \mathbf{v}_3 belongs to the 2-dimensional plane determined by $\mathbf{v}_3 - \mathbf{v}_2$ and $\mathbf{v}_1 - \mathbf{v}_2$, since the latter is a normal vector to V_2 at \mathbf{v}_2 .

Consider the triangle $(\mathbf{v}_1 - \mathbf{v}_2, \mathbf{v}_3 - \mathbf{v}_2, \mathbf{v}_1 - \mathbf{v}_3)$. This is a right-angled triangle since $\mathbf{v}_1 - \mathbf{v}_2 \perp V_2$ and $\mathbf{v}_3 - \mathbf{v}_2 \in V_2$, so that

$$\|\mathbf{v}_1 - \mathbf{v}_3\|_2 = \frac{\|\mathbf{v}_1 - \mathbf{v}_2\|_2}{\sin(\mathbf{v}_2 - \mathbf{v}_3, \mathbf{v}_1 - \mathbf{v}_3)}.$$

Subsequently, let \mathbf{v}_4 be the intersection of $\mathbf{v}_1 - \mathbf{v}_2$ and $T_w V_1(\mathbf{v}_3)$, due to the convexity of V_1 ,

$$\begin{aligned}
 \sin^2(\mathbf{v}_2 - \mathbf{v}_3, \mathbf{v}_1 - \mathbf{v}_3) &\geq \sin^2(\mathbf{v}_2 - \mathbf{v}_3, \mathbf{v}_4 - \mathbf{v}_3) \\
 &= 1 - \cos^2(\mathbf{v}_2 - \mathbf{v}_3, \mathbf{v}_4 - \mathbf{v}_3) \\
 &= 1 - \cos^2(\mathbf{p} + \mathbf{u}, (\mathbf{p} + \mathbf{u}) \circ f'(\mathbf{v}_3)),
 \end{aligned}$$

where $\mathbf{p} + \mathbf{u}$ is a normal vector to V_2 containing $\mathbf{v}_2 - \mathbf{v}_3$, $(\mathbf{p} + \mathbf{u}) \circ f'(\mathbf{v}_3)$ is a normal vector to $T_w V_1(\mathbf{v}_3)$ containing $\mathbf{v}_4 - \mathbf{v}_3$, and both of them belong to the 2-dimensional plane determined by $\mathbf{v}_3 - \mathbf{v}_2$ and $\mathbf{v}_1 - \mathbf{v}_2$ as we have proven above. By definition, \mathbf{v}_3 is on the boundary of V_1 such that $\sum_s (p_s + u_s)f(v_{3,s}) = \delta$, and according to Lemma 5,

$$\begin{aligned}
 1 - \cos^2(\mathbf{p} + \mathbf{u}, (\mathbf{p} + \mathbf{u}) \circ f'(\mathbf{v}_3)) &\geq 1 - (1 - C_{\cos})^2 \\
 &= C_{\cos}(2 - C_{\cos}),
 \end{aligned}$$

which implies

$$\begin{aligned}
 \|\mathbf{v}_1 - \mathbf{v}_3\|_2 &\leq \frac{\|\mathbf{v}_1 - \mathbf{v}_2\|_2}{\sqrt{C_{\cos}(2 - C_{\cos})}} \\
 &= \mathcal{O}(\|\mathbf{v}_1 - \mathbf{v}_2\|_2) \\
 &= \mathcal{O}(\|\mathbf{u}\|_2).
 \end{aligned}$$

Combining $\|\mathbf{v}_1 - \mathbf{w}_0\|_2 = \mathcal{O}(\|\mathbf{u}\|_2)$ and $\|\mathbf{v}_1 - \mathbf{v}_3\|_2 = \mathcal{O}(\|\mathbf{u}\|_2)$ gives us $\|\mathbf{v}_3 - \mathbf{w}_0\|_2 = \mathcal{O}(\|\mathbf{u}\|_2)$, and thus $\inf_{\mathbf{v} \in \Phi(\mathbf{u})} \|\mathbf{v} - \mathbf{w}_0\|_2 \leq \|\mathbf{v}_3 - \mathbf{w}_0\|_2 = \mathcal{O}(\|\mathbf{u}\|_2)$. Since \mathbf{w}_0 is arbitrarily picked from $\Phi(\mathbf{0})$, the proof is completed. \square

Table 2: Summary of dataset statistics.

Dataset	#Points	#classes	Dimension
blood	748	2	4
adult	32561	2	123
fourclass	862	2	2
phishing	11055	2	68
20news	18040	20	50
satimage	4435	6	36
letter	20000	26	16
mnist	70000	10	50

Lemma 7 (Lipschitz continuity of the difference function, II). *For all \mathbf{u} such that $\|\mathbf{u}\|_2 \leq \epsilon$, $D_2(\theta)$ is Lipschitz continuous with a Lipschitz constant of order $\mathcal{O}(\|\mathbf{u}\|_2^{1/2})$.*

Proof. The proof goes along the same line with Lemma 4. Let $\mathbf{w}^* = \arg \sup_{\mathbf{w} \in \Phi(\mathbf{0})} \sum_s w_s l_s$ and $\mathbf{v}^* = \arg \sup_{\mathbf{w} \in \Phi(\mathbf{u})} \sum_s w_s l_s$, and consider the perturbation analysis of the following optimization problem

$$\min_{\mathbf{w}} - \sum_{s=1}^S w_s l_s \quad \text{s.t. } \mathbf{w} \in \Phi(\mathbf{u}),$$

whose objective is unperturbed and feasible region is perturbed. According to Lemma 6, we have $d_{\mathcal{H}}(\Phi(\mathbf{u}), \Phi(\mathbf{0})) = \mathcal{O}(\|\mathbf{u}\|_2)$, which ensures that the multifunction $\mathbf{u} \mapsto \Phi(\mathbf{u})$ is upper Lipschitz continuous and that $d_{\mathcal{H}}(\{\mathbf{w}^*\}, \Phi(\mathbf{u})) = \mathcal{O}(\|\mathbf{u}\|_2)$. Hence, $\|\mathbf{v}^* - \mathbf{w}^*\|_2 = \mathcal{O}(\|\mathbf{u}\|_2^{1/2})$ by applying Proposition 6.4 in (Bonnans & Shapiro, 1998). \square

Let us summarize what we have obtained so far:

- a second-order growth condition of $J(\theta, \mathbf{0}, \mathbf{0})$ at $\theta = \theta^*$;
- the Lipschitz continuity of $J(\theta, \mathbf{u}, \mathbf{u}')$ with a Lipschitz constant independent of \mathbf{u} and \mathbf{u}' ;
- the Lipschitz continuity of $D(\theta)$ with a Lipschitz constant of order $\mathcal{O}(\|\mathbf{u}\|_2^{1/2} + \|\mathbf{u}'\|_2)$.

Note that θ is unconstrained, by applying Proposition 6.1 in (Bonnans & Shapiro, 1998), we can obtain

$$\|\theta_{\mathbf{u}, \mathbf{u}'} - \theta^*\|_2 = \mathcal{O}(\|\mathbf{u}\|_2^{1/2} + \|\mathbf{u}'\|_2).$$

This immediately implies

$$\|J(\theta_{\mathbf{u}, \mathbf{u}'}, \mathbf{0}, \mathbf{0}) - J(\theta^*, \mathbf{0}, \mathbf{0})\|_2 = \mathcal{O}(\|\mathbf{u}\|_2^{1/2} + \|\mathbf{u}'\|_2),$$

due to the Lipschitz continuity of $J(\theta, \mathbf{0}, \mathbf{0})$. \square

I. Datasets

We obtained six classification datasets from the UCI repository¹⁶ and also obtained 20newsgroups¹⁷ and MNIST datasets. We used the raw features for the datasets from the UCI repository. For the 20newsgroups dataset, we removed stop words and retained the 2000 most frequent words. We then removed documents with fewer than 10 words. We extracted tf-idf features and applied Principle Component Analysis (PCA) to reduce the dimensionality to 50. For the MNIST dataset, we applied PCA on the raw features to reduce the dimensionality to 50. The dataset statistics are summarized in Table 2.

J. Details of the Subcategory Shift Scenario

In this section, we give details on how we converted the original multi-class classification problems into multi-class classification problems with fewer classes.

¹⁶<http://archive.ics.uci.edu/ml/index.html>

¹⁷<http://qwone.com/~jason/20NewsGroups/>

Table 3: Experimental comparisons of the three methods w.r.t. the estimated ordinary risk and the estimated structural adversarial risk *using the surrogate loss (the logistic loss)*. The lower these values are, the better the performance of the method is. The KL divergence is used and distribution shift is assumed to be (a) class prior change and (b) sub-category prior change. Mean and standard deviation over 50 random train-test splits were reported. The best method and comparable ones based on the t-test at the significance level 1% are highlighted in boldface.

(a) Class prior change.

Dataset	Estimated ordinary risk			Estimated adversarial risk			Estimated structural adversarial risk		
	ERM	AERM	Structural AERM	ERM	AERM	Structural AERM	ERM	AERM	Structural AERM
blood	0.52 (0.05)	0.69 (0.0)	0.62 (0.02)	1.04 (0.1)	0.69 (0.0)	0.97 (0.03)	0.86 (0.23)	0.69 (0.0)	0.63 (0.19)
adult	0.33 (0.0)	0.65 (0.03)	0.39 (0.0)	1.28 (0.02)	0.69 (0.01)	1.42 (0.01)	0.59 (0.3)	0.67 (0.01)	0.4 (0.38)
fourclass	0.51 (0.05)	0.69 (0.0)	0.54 (0.05)	0.91 (0.04)	0.69 (0.0)	0.88 (0.04)	0.65 (0.13)	0.69 (0.0)	0.56 (0.13)
phishing	0.15 (0.01)	0.41 (0.08)	0.15 (0.0)	0.86 (0.01)	0.59 (0.08)	0.85 (0.01)	0.18 (0.06)	0.41 (0.02)	0.16 (0.05)
20news	1.05 (0.01)	1.49 (0.04)	1.22 (0.02)	3.42 (0.02)	3.0 (0.04)	3.58 (0.03)	1.43 (0.1)	1.74 (0.19)	1.32 (0.13)
satimage	1.01 (0.01)	1.26 (0.02)	1.29 (0.01)	2.54 (0.02)	2.15 (0.02)	2.86 (0.02)	1.41 (0.05)	1.59 (0.01)	1.38 (0.04)
letter	0.37 (0.01)	0.47 (0.03)	0.51 (0.02)	1.65 (0.02)	1.19 (0.03)	2.27 (0.03)	0.77 (0.17)	0.77 (0.09)	0.58 (0.21)
mnist	0.35 (0.0)	0.59 (0.05)	0.45 (0.0)	1.96 (0.01)	1.38 (0.04)	1.85 (0.01)	0.49 (0.06)	0.73 (0.02)	0.47 (0.04)

(b) Sub-category prior change.

Dataset	Estimated ordinary risk			Estimated adversarial risk			Estimated structural adversarial risk		
	ERM	AERM	Structural AERM	ERM	AERM	Structural AERM	ERM	AERM	Structural AERM
20news	0.61 (0.01)	0.84 (0.06)	0.76 (0.05)	2.76 (0.02)	2.03 (0.05)	2.91 (0.03)	1.02 (0.14)	1.08 (0.29)	0.89 (0.21)
satimage	0.63 (0.0)	0.69 (0.0)	0.68 (0.0)	0.99 (0.01)	0.69 (0.0)	0.74 (0.0)	0.81 (0.02)	0.69 (0.0)	0.69 (0.02)
letter	0.47 (0.0)	0.69 (0.0)	0.64 (0.02)	1.05 (0.03)	0.69 (0.0)	0.86 (0.01)	0.93 (0.02)	0.69 (0.0)	0.68 (0.06)
mnist	0.31 (0.0)	0.68 (0.01)	0.39 (0.0)	1.28 (0.01)	0.69 (0.0)	1.25 (0.0)	0.49 (0.05)	0.68 (0.0)	0.42 (0.03)

For the datasets from the UCI repository, we systematically grouped the class labels into binary categories by the following procedure. First, class labels are sorted by the number of data points in the classes. Then, 1, 3, 5, . . . -th labels are assigned to a positive category and the others are assigned to a negative category. For MNIST, we considered a binary classification between odd and even numbers and set the original classes as subcategories. For 20newsgroups, we converted the original 20-class classification problem into a 7-class one with each class corresponding to a high-level topic: comp, rec, sci, misc, alt, soc, and talk. We then set the original classes as subcategories.

K. Experimental Results measured by Surrogate Loss

In this section, we report the experimental results measured by the surrogate loss (the logistic loss). We used the KL and the PE divergences, where we set $\delta = 0.5$. We used the same f -divergence and the same δ during training and testing. The experimental results using the KL and the PE divergences are reported in Tables 3 and 4, respectively. We empirically confirmed that *in terms of the surrogate loss*, each method indeed achieved the best performance in terms of the metric it optimizes for.

L. Experimental Results with the PE divergence

In this section, we report the experimental results using the PE divergence, where we set $\delta = 0.5$. The experimental results are reported in Table 5.

Table 4: Experimental comparisons of the three methods w.r.t. the estimated ordinary risk and the estimated structural adversarial risk *using the surrogate loss (the logistic loss)*. The lower these values are, the better the performance of the method is. The PE divergence is used and distribution shift is assumed to be (a) class prior change and (b) sub-category prior change. Mean and standard deviation over 50 random train-test splits were reported. The best method and comparable ones based on the t-test at the significance level 1% are highlighted in boldface.

(a) Class prior change.

Dataset	Estimated ordinary risk			Estimated adversarial risk			Estimated structural adversarial risk		
	ERM	AERM	Structural AERM	ERM	AERM	Structural AERM	ERM	AERM	Structural AERM
blood	0.52 (0.05)	0.67 (0.02)	0.61 (0.03)	0.77 (0.04)	0.69 (0.0)	0.81 (0.02)	0.71 (0.07)	0.69 (0.0)	0.62 (0.07)
adult	0.33 (0.0)	0.41 (0.02)	0.39 (0.01)	0.69 (0.01)	0.61 (0.01)	0.77 (0.01)	0.49 (0.02)	0.51 (0.0)	0.4 (0.03)
fourclass	0.52 (0.05)	0.66 (0.02)	0.53 (0.05)	0.73 (0.02)	0.69 (0.01)	0.77 (0.04)	0.6 (0.04)	0.67 (0.0)	0.54 (0.06)
phishing	0.15 (0.01)	0.2 (0.02)	0.15 (0.0)	0.43 (0.01)	0.4 (0.02)	0.43 (0.01)	0.17 (0.02)	0.21 (0.01)	0.15 (0.01)
20news	1.04 (0.01)	1.15 (0.11)	1.17 (0.02)	1.99 (0.01)	1.95 (0.1)	2.2 (0.02)	1.29 (0.03)	1.36 (0.06)	1.24 (0.04)
satimage	1.01 (0.01)	1.1 (0.01)	1.1 (0.01)	1.81 (0.01)	1.72 (0.02)	2.0 (0.01)	1.26 (0.02)	1.32 (0.01)	1.17 (0.02)
letter	0.36 (0.01)	0.4 (0.01)	0.42 (0.02)	0.89 (0.02)	0.81 (0.02)	1.0 (0.02)	0.6 (0.04)	0.59 (0.03)	0.53 (0.05)
mnist	0.35 (0.0)	0.41 (0.01)	0.43 (0.0)	0.96 (0.01)	0.87 (0.01)	1.04 (0.01)	0.45 (0.02)	0.5 (0.01)	0.44 (0.01)

(b) Sub-category prior change.

Dataset	Estimated ordinary risk			Estimated adversarial risk			Estimated structural adversarial risk		
	ERM	AERM	Structural AERM	ERM	AERM	Structural AERM	ERM	AERM	Structural AERM
20news	0.61 (0.01)	0.67 (0.01)	0.68 (0.02)	1.34 (0.01)	1.22 (0.01)	1.4 (0.01)	0.86 (0.04)	0.86 (0.03)	0.81 (0.05)
satimage	0.63 (0.0)	0.69 (0.0)	0.68 (0.0)	0.86 (0.01)	0.69 (0.0)	0.73 (0.0)	0.76 (0.01)	0.69 (0.0)	0.69 (0.01)
letter	0.47 (0.0)	0.67 (0.0)	0.57 (0.04)	0.79 (0.01)	0.69 (0.0)	0.73 (0.01)	0.74 (0.01)	0.69 (0.0)	0.66 (0.03)
mnist	0.31 (0.0)	0.4 (0.02)	0.36 (0.0)	0.72 (0.0)	0.6 (0.01)	0.72 (0.0)	0.44 (0.01)	0.48 (0.0)	0.41 (0.01)

Table 5: Experimental comparisons of the three methods w.r.t. the estimated ordinary risk and the estimated structural adversarial risk *using the 0-1 loss (%)*. The lower these values are, the better the performance of the method is. The PE divergence is used and distribution shift is assumed to be (a) class prior change and (b) sub-category prior change. Mean and standard deviation over 50 random train-test splits were reported. The best method and comparable ones based on the t-test at the significance level 1% are highlighted in boldface.

(a) Class prior change.

Dataset	Estimated ordinary risk			Estimated structural adversarial risk		
	ERM	AERM	Structural AERM	ERM	AERM	Structural AERM
blood	22.2 (0.6)	22.4 (0.5)	33.0 (2.0)	47.4 (1.6)	49.1 (1.1)	36.6 (2.4)
adult	15.3 (0.1)	15.3 (0.1)	18.7 (0.2)	24.9 (0.3)	24.8 (0.3)	19.1 (0.3)
fourclass	23.7 (1.2)	23.5 (1.2)	27.0 (1.3)	32.7 (1.7)	32.6 (1.8)	28.7 (1.7)
phishing	6.0 (0.2)	6.1 (0.2)	5.9 (0.2)	7.1 (0.3)	7.4 (0.3)	6.4 (0.3)
20news	28.8 (0.3)	29.7 (0.3)	33.7 (0.3)	37.9 (0.3)	38.4 (0.4)	37.5 (0.4)
satimage	25.1 (0.2)	26.7 (0.3)	28.1 (0.3)	33.7 (0.4)	35.6 (0.4)	32.0 (0.4)
letter	14.2 (0.5)	14.5 (0.5)	15.5 (0.5)	26.5 (0.9)	25.6 (0.9)	20.8 (0.7)
mnist	10.0 (0.1)	10.1 (0.1)	12.2 (0.1)	13.3 (0.1)	13.2 (0.1)	13.1 (0.1)

(b) Sub-category prior change.

Dataset	Estimated ordinary risk			Estimated structural adversarial risk		
	ERM	AERM	Structural AERM	ERM	AERM	Structural AERM
20news	19.0 (0.3)	19.6 (0.4)	20.8 (0.4)	29.2 (0.4)	29.7 (0.4)	27.3 (0.4)
satimage	36.4 (0.3)	41.1 (1.9)	39.6 (0.5)	53.9 (0.4)	56.7 (2.7)	47.0 (0.5)
letter	17.4 (0.4)	18.5 (0.5)	23.1 (3.3)	38.0 (0.5)	39.5 (0.6)	38.9 (1.0)
mnist	13.3 (0.1)	13.5 (0.1)	15.6 (0.2)	20.0 (0.2)	20.2 (0.2)	18.6 (0.2)