

## A. Proofs

In this section, we provide the proofs of our theoretical results in Section 3.

For ease of presentation, we introduce the following notations and an alternative definition of the  $\infty$ -Wasserstein distance. Let  $(\mathcal{X}, \mu)$  and  $(\mathcal{Y}, \nu)$  be two probability spaces. We say that  $T : \mathcal{X} \rightarrow \mathcal{Y}$  transports  $\mu \in \mathcal{P}(\mathcal{X})$  to  $\nu \in \mathcal{P}(\mathcal{Y})$ , and we call  $T$  a transport map, if  $\nu(B) = \mu(T^{-1}(B))$ , for all  $\nu$ -measurable sets  $B$ . In addition, for any measurable map  $T : \mathcal{X} \rightarrow \mathcal{Y}$ , we define the *pushforward* of  $\mu$  through  $T$  as  $(T_{\#}(\mu))$  given by

$$(T_{\#}(\mu))(B) = \mu(T^{-1}(B)), \quad \text{for any measurable } B \subseteq \mathcal{Y}.$$

**Alternative definition of  $\infty$ -Wasserstein distance.** From the perspective of transportation theory, given two probability measures  $\mu$  and  $\nu$  on  $(\mathcal{X}, \Delta)$ , any joint probability distribution  $\gamma \in \Gamma(\mu, \nu)$  corresponds to a specific transport map  $T : \mathcal{X} \rightarrow \mathcal{X}$  that moves  $\mu$  to  $\nu$ . Then, the  $p$ -th Wasserstein distance can be viewed as finding the optimal transport map to move from  $\mu$  to  $\nu$  that minimizes some cost functional depending on  $p$  (Kolouri et al., 2017). For the case where  $p = \infty$ , if we let  $T$  be the transport map induced by a given  $\gamma \in \Gamma(\mu, \nu)$ , then the cost functional can be informally understood as the maximum of all the transport distances  $\Delta(T(\mathbf{x}), \mathbf{x})$ . More rigorously, the  $\infty$ -Wasserstein distance can be alternatively defined as

$$\begin{aligned} W_{\infty}(\mu, \nu) &:= \inf_{\gamma \in \Gamma(\mu, \nu)} \gamma\text{-ess sup}_{(\mathbf{x}, \mathbf{x}') \in \mathcal{X}^2} \Delta(\mathbf{x}, \mathbf{x}') \\ &= \inf_{\gamma \in \Gamma(\mu, \nu)} \inf \{t \geq 0 : \gamma(\Delta(\mathbf{x}, \mathbf{x}') > t) = 0\}. \end{aligned}$$

A more detailed discussion of  $\infty$ -Wasserstein distance can be found in Champion et al. (2008).

### A.1. Proof of Theorem 3.2

Theorem 3.2 (restated here) connects the vulnerability of a given representation with the minimum adversarial gap of any classifier based on that representation.

**Theorem 3.2.** Let  $(\mathcal{X}, \|\cdot\|_p)$  be the input metric space and  $\mathcal{Y} = \{-1, 1\}$  be the label space. Assume the underlying data are generated according to (3.1). Consider the feature space  $\mathcal{Z} = \{-1, 1\}$  and the set of representations,

$$\mathcal{G}_{\text{bin}} = \{g : \mathbf{x} \mapsto \text{sgn}(\mathbf{w}^{\top} \mathbf{x}), \forall \mathbf{x} \in \mathcal{X} \mid \|\mathbf{w}\|_2 = 1\}.$$

Let  $\mathcal{H} = \{h : \mathcal{Z} \rightarrow \mathcal{Y}\}$  be the set of non-trivial downstream classifiers.<sup>6</sup> Given  $\epsilon \geq 0$ , for any  $g \in \mathcal{G}_{\text{bin}}$ , we have

$$\int_{\frac{1}{2} - \text{AG}_{\epsilon}(f^*)}^{\frac{1}{2}} H_2'(\theta) d\theta \leq \text{RV}_{\epsilon}(g) \leq \int_{\frac{1}{2} - \frac{1}{2} \text{AG}_{\epsilon}(f^*)}^{\frac{1}{2}} H_2'(\theta) d\theta,$$

where  $f^* = \text{argmin}_{h \in \mathcal{H}} \text{AdvRisk}_{\epsilon}(h \circ g)$  is the optimal classifier based on  $g$ ,  $H_2(\theta) = -\theta \log \theta - (1 - \theta) \log(1 - \theta)$  is the binary entropy function and  $H_2'$  denotes its derivative.

*Proof.* Let  $\mu_{XY}$  be the underlying joint probability distribution of the examples according to (3.1) and  $\mu_X$  be corresponding the marginal distribution of  $X$ . To begin with, we compute the explicit formulation for the defined representation vulnerability in Definition 3.1. Note that  $I(U; V) = H(U) - H(U|V) = H(V) - H(V|U)$ . Thus, for any  $g_{\mathbf{w}} \in \mathcal{G}$  we have

$$\begin{aligned} \text{RV}_{\epsilon}(g_{\mathbf{w}}) &= H(g_{\mathbf{w}}(X)) - \inf_{\mu_{X'} \in \mathcal{B}_{W_{\infty}}(\mu_X, \epsilon)} H(g_{\mathbf{w}}(X')) \\ &= 1 - \inf_{\mu_{X'} \in \mathcal{B}_{W_{\infty}}(\mu_X, \epsilon)} \left( - \Pr_{\mathbf{x}' \sim \mu_{X'}}(\mathbf{w}^{\top} \mathbf{x}' \geq 0) \cdot \log \left[ \Pr_{\mathbf{x}' \sim \mu_{X'}}(\mathbf{w}^{\top} \mathbf{x}' \geq 0) \right] \right. \\ &\quad \left. - \Pr_{\mathbf{x}' \sim \mu_{X'}}(\mathbf{w}^{\top} \mathbf{x}' < 0) \cdot \log \left[ \Pr_{\mathbf{x}' \sim \mu_{X'}}(\mathbf{w}^{\top} \mathbf{x}' < 0) \right] \right), \end{aligned}$$

<sup>6</sup>To be more specific, we do not consider the case where  $h$  is a constant function. Under our problem setting, there are two elements in  $\mathcal{H}$ , namely  $h_1(z) = z$ ,  $h_2(z) = -z$ , for any  $z \in \mathcal{Z}$ .

where the first equality holds because  $H(g_w(U) | U) = 0$  for any random variable  $U$ , and the second equality is due to the fact that the distribution of  $X$  is symmetric with respect to  $\mathbf{w}^\top \mathbf{x} = 0$ . Note that the binary entropy function  $H_2(\theta) = -\theta \log \theta - (1-\theta) \log(1-\theta)$  is monotonically increasing with respect to  $\theta$  in  $[0, 1/2)$  and monotonically decreasing in  $(1/2, 1]$ . Therefore, the optimal value of  $\text{RV}_\epsilon(g_w)$  is achieved when  $\mu_{X'}$  either minimizes  $\Pr_{\mathbf{x}' \sim \mu_{X'}}(\mathbf{w}^\top \mathbf{x}' \geq 0)$  or maximizes  $\Pr_{\mathbf{x}' \sim \mu_{X'}}(\mathbf{w}^\top \mathbf{x}' \geq 0)$ .

According to the Hölder's inequality, we have  $|\langle \mathbf{a}, \mathbf{b} \rangle| \leq \|\mathbf{a}\|_p \cdot \|\mathbf{b}\|_q$  for any  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ , where  $1/p + 1/q = 1$ . By the alternative definition of  $\infty$ -Wasserstein distance, for any  $\mu_{X'}$  that satisfies  $W_\infty(\mu_{X'}, \mu_X) \leq \epsilon$ , it induces a transport map  $T: \mathcal{X} \rightarrow \mathcal{X}$  such that  $\mu'_{X'} = T_\#(\mu_X)$  and  $\|\Delta(T(X), X)\|_p \leq \epsilon$  holds almost surely with respect to the randomness of  $X$  and  $T$ . Thus, we have

$$\Pr_{\mathbf{x} \sim \mu_X} \left[ -\epsilon \cdot \|\mathbf{w}\|_q \leq \mathbf{w}^\top (T(\mathbf{x}) - \mathbf{x}) \leq \epsilon \cdot \|\mathbf{w}\|_q \right] \geq \Pr_{\mathbf{x} \sim \mu_X} \left[ \|T(\mathbf{x}) - \mathbf{x}\|_p \leq \epsilon \right] = 1,$$

which implies

$$\Pr_{\mathbf{x} \sim \mu_X} (\mathbf{w}^\top \mathbf{x} - \epsilon \cdot \|\mathbf{w}\|_q \geq 0) \leq \Pr_{\mathbf{x}' \sim \mu_{X'}} (\mathbf{w}^\top \mathbf{x}' \geq 0) \leq \Pr_{\mathbf{x} \sim \mu_X} (\mathbf{w}^\top \mathbf{x} + \epsilon \cdot \|\mathbf{w}\|_q \geq 0).$$

We remark that the equality can be achieved when the transport map  $T$  is constructed by perturbing the  $i$ -th element of any sampled  $\mathbf{x} \sim \mu_X$  by  $\epsilon \cdot (w_i^q / \sum_i w_i^q)^{1/p}$ , for any  $i = 1, 2, \dots, d$ . In addition, according to the assumed Gaussian Mixture model (3.1), we have

$$\begin{aligned} \Pr_{\mathbf{x} \sim \mu_X} (\mathbf{w}^\top \mathbf{x} - \epsilon \cdot \|\mathbf{w}\|_q \geq 0) &= \frac{1}{2} \Pr_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\theta}^*, \boldsymbol{\Sigma}^*)} [\mathbf{w}^\top \mathbf{x} \geq \epsilon \cdot \|\mathbf{w}\|_q] + \frac{1}{2} \Pr_{\mathbf{x} \sim \mathcal{N}(-\boldsymbol{\theta}^*, \boldsymbol{\Sigma}^*)} [\mathbf{w}^\top \mathbf{x} \geq \epsilon \cdot \|\mathbf{w}\|_q] \\ &= \frac{1}{2} - \frac{1}{2} \Pr_{Z \sim \mathcal{N}(0,1)} \left[ \frac{-\epsilon \cdot \|\mathbf{w}\|_q + \mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \leq Z \leq \frac{\epsilon \cdot \|\mathbf{w}\|_q + \mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right]. \end{aligned}$$

Similarly, we have

$$\Pr_{\mathbf{x} \sim \mu_X} (\mathbf{w}^\top \mathbf{x} + \epsilon \cdot \|\mathbf{w}\|_q \geq 0) = \frac{1}{2} + \frac{1}{2} \Pr_{Z \sim \mathcal{N}(0,1)} \left[ \frac{-\epsilon \cdot \|\mathbf{w}\|_q + \mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \leq Z \leq \frac{\epsilon \cdot \|\mathbf{w}\|_q + \mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right].$$

Therefore, we derive the explicit formulation for  $\text{RV}_\epsilon(g_w)$

$$\text{RV}_\epsilon(g_w) = H_2\left(\frac{1}{2}\right) - H_2\left(\frac{1}{2} - \Pr_{Z \sim \mathcal{N}(0,1)} \left[ \frac{\mathbf{w}^\top \boldsymbol{\theta}^* - \epsilon \cdot \|\mathbf{w}\|_q}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \leq Z \leq \frac{\mathbf{w}^\top \boldsymbol{\theta}^* + \epsilon \cdot \|\mathbf{w}\|_q}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right]\right), \quad (\text{A.1})$$

where  $H_2(\cdot)$  is denotes binary entropy function.

Next, given  $g_w \in \mathcal{G}_{\text{bin}}$ , we are going to compute the adversarial gap of  $f \circ g_w$  for each  $h \in \mathcal{H}$ . To begin with, we consider the first case  $h_1(z) = z$  for any  $z \in \mathcal{Z}$ . According to the definition of adversarial risk, we have

$$\begin{aligned} \text{AdvRisk}_\epsilon(h_1 \circ g_w) &= \Pr_{(\mathbf{x}, y) \sim \mu_{XY}} [\exists \mathbf{x}' \in \mathcal{B}(\mathbf{x}, \epsilon) \text{ s.t. } \text{sgn}(\mathbf{w}^\top \mathbf{x}') \neq y] \\ &= \Pr_{(\mathbf{x}, y) \sim \mu_{XY}} \left[ \min_{\mathbf{x}' \in \mathcal{B}(\mathbf{x}, \epsilon)} y \cdot \mathbf{w}^\top \mathbf{x}' \leq 0 \right] \\ &= \Pr_{(\mathbf{x}, y) \sim \mu_{XY}} \left[ y \cdot \mathbf{w}^\top \mathbf{x} \leq - \min_{\boldsymbol{\Delta} \in \mathcal{B}(\mathbf{0}, \epsilon)} \mathbf{w}^\top \boldsymbol{\Delta} \right] \\ &= \Pr_{Z \sim \mathcal{N}(0,1)} \left[ Z \leq \frac{\epsilon \|\mathbf{w}\|_q - \mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right], \end{aligned}$$

where the equality is due to the fact that  $\mathcal{B}(\mathbf{0}, \epsilon)$  is symmetric with respect to  $\mathbf{0}$ , and the last equality holds because of the Hölder's inequality: for any  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ , it holds that  $\mathbf{a}^\top \mathbf{b} \geq -\|\mathbf{a}\|_p \cdot \|\mathbf{b}\|_q$  and the equality is achieved when  $(a_i / \|\mathbf{a}\|_p)^p = (b_i / \|\mathbf{b}\|_q)^q$  for any  $i \in \{1, 2, \dots, d\}$ .

Similarly, the standard risk can be computed as:

$$\text{Risk}(h_1 \circ g_w) = \Pr_{(\mathbf{x}, y) \sim \mathcal{D}} [\text{sgn}(\mathbf{w}^\top \mathbf{x}) \neq y] = \Pr_{(\mathbf{x}, y) \sim \mathcal{D}} [y \cdot \mathbf{w}^\top \mathbf{x} \leq 0] = \Pr_{Z \sim \mathcal{N}(0,1)} \left[ Z \leq \frac{-\mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right].$$

Thus, we derive the gap between standard and adversarial risk with respect to  $h_1 \circ g_w$ :

$$\text{AG}_\epsilon(h_1 \circ g_w) = \Pr_{Z \sim \mathcal{N}(0,1)} \left[ \frac{\mathbf{w}^\top \boldsymbol{\theta}^* - \epsilon \|\mathbf{w}\|_q}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \leq Z \leq \frac{\mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right].$$

For the other case where  $h_2(z) = -z$  for any  $z \in \mathcal{Z}$ , note that  $h_1 \circ g_w = h_2 \circ g_{-w}$  for any  $g_w \in \mathcal{G}_{\text{bin}}$ . Thus, a similar proof technique can be applied to compute the adversarial risk,

$$\text{AdvRisk}_\epsilon(h_2 \circ g_w) = \Pr_{Z \sim \mathcal{N}(0,1)} \left[ Z \leq \frac{\epsilon \|\mathbf{w}\|_q + \mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right],$$

and the adversarial gap,

$$\text{AG}_\epsilon(h_2 \circ g_w) = \Pr_{Z \sim \mathcal{N}(0,1)} \left[ \frac{\mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \leq Z \leq \frac{\mathbf{w}^\top \boldsymbol{\theta}^* + \epsilon \|\mathbf{w}\|_q}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right].$$

Note that  $f^* : \mathcal{X} \rightarrow \mathcal{Y}$  is the optimal classifier based on  $g_w$  that minimizes the adversarial risk of  $h \circ g_w$  for any  $h \in \mathcal{H}$ . By comparing the adversarial risk of  $h_1 \circ g_w$  and  $h_2 \circ g_w$ , we have  $f^* = h_1 \circ g_w$ , if  $\mathbf{w}^\top \boldsymbol{\theta}^* \geq 0$ ;  $f^* = h_2 \circ g_w$ , otherwise. Thus, we derive the adversarial gap with respect to  $f^*$  as follows

$$\text{AG}_\epsilon(f^*) = \begin{cases} \Pr_{Z \sim \mathcal{N}(0,1)} \left( \frac{\mathbf{w}^\top \boldsymbol{\theta}^* - \epsilon \|\mathbf{w}\|_q}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \leq Z \leq \frac{\mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right), & \text{if } \mathbf{w}^\top \boldsymbol{\theta}^* \geq 0; \\ \Pr_{Z \sim \mathcal{N}(0,1)} \left( \frac{\mathbf{w}^\top \boldsymbol{\theta}^*}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \leq Z \leq \frac{\mathbf{w}^\top \boldsymbol{\theta}^* + \epsilon \|\mathbf{w}\|_q}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right), & \text{otherwise.} \end{cases} \quad (\text{A.2})$$

Based on (A.2), we further obtain the following inequality

$$\text{AG}_\epsilon(f^*) \leq \Pr_{Z \sim \mathcal{N}(0,1)} \left[ \frac{\mathbf{w}^\top \boldsymbol{\theta}^* - \epsilon \cdot \|\mathbf{w}\|_q}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \leq Z \leq \frac{\mathbf{w}^\top \boldsymbol{\theta}^* + \epsilon \cdot \|\mathbf{w}\|_q}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}^* \mathbf{w}}} \right] \leq 2 \cdot \text{AG}_\epsilon(f^*).$$

Finally, according to the formulation of representation vulnerability (A.1), we have

$$\int_{\frac{1}{2} - \text{AG}_\epsilon(f^* \circ g_w)}^{\frac{1}{2}} \text{H}'_2(\theta) d\theta \leq \text{RV}_\epsilon(g_w) \leq \int_{\frac{1}{2} - \frac{1}{2} \text{AG}_\epsilon(f^* \circ g_w)}^{\frac{1}{2}} \text{H}'_2(\theta) d\theta,$$

which completes the proof.  $\square$

## A.2. Proof of Lemma 3.3

Lemma 3.3, restated below, connects adversarial risk and input distribution perturbations bounded in an  $\infty$ -Wasserstein ball.

**Lemma 3.3.** Let  $(\mathcal{X}, \Delta)$  be the input metric space and  $\mathcal{Y}$  be the set of labels. Assume all the examples are generated from a joint probability distribution  $(X, Y) \sim \mu_{XY}$ . Let  $\mu_X$  be the marginal distribution of  $X$ . Then, for any classifier  $f : \mathcal{X} \rightarrow \mathcal{Y}$  and  $\epsilon \geq 0$ , we have

$$\text{AdvRisk}_\epsilon(f) = \sup_{\mu_{X'} \in \mathcal{B}_{W_\infty}(\mu_X, \epsilon)} \Pr [f(X') \neq Y],$$

where  $X'$  denotes the random variable that follows  $\mu_{X'}$ .

*Proof.* Our proof proves the equality by proving  $\leq$  inequalities in both directions. First, we prove

$$\text{AdvRisk}_\epsilon(f) \leq \sup_{\mu_{X'} \in \mathcal{B}_{W_\infty}(\mu_X, \epsilon)} \Pr [f(X') \neq Y]. \quad (\text{A.3})$$

For any classifier  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , according to Definition 2.3, we have

$$\text{AdvRisk}_\epsilon(f) = \Pr_{(\mathbf{x}, y) \sim \mu_{XY}} [\exists \mathbf{x}' \in \mathcal{B}(\mathbf{x}, \epsilon) \text{ s.t. } f(\mathbf{x}') \neq y].$$

Since  $f$  is a given deterministic function, the optimal perturbation scheme that achieves  $\text{AdvRisk}_\epsilon(f)$  essentially defines a transport map  $T : \mathcal{X} \rightarrow \mathcal{X}$ . More specifically, let  $\mathcal{C}_y(f) = \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \neq y\}$ . Then, for any sampled pair  $(\mathbf{x}, y) \sim \mu_{XY}$ , we can construct  $T$  such that

$$T(\mathbf{x}) = \begin{cases} \operatorname{argmin}_{\mathbf{x}' \in \mathcal{C}_y(f)} \Delta(\mathbf{x}', \mathbf{x}), & \text{if } \mathcal{C}_y(f) \cap \mathcal{B}(\mathbf{x}, \epsilon) \neq \emptyset; \\ \mathbf{x}, & \text{otherwise.} \end{cases}$$

Let  $(X, Y)$  be the random variable that follows  $\mu_{XY}$ . By construction, it can be easily verified that  $T_{\#}(\mu_X) \in \mathcal{B}_{W_\infty}(\mu_X, \epsilon)$  and  $\text{AdvRisk}_\epsilon(f) = \Pr[f(T(X)) \neq Y]$ . Therefore, we have proven (A.3).

It remains to prove the other direction of the inequality:

$$\text{AdvRisk}_\epsilon(f) \geq \sup_{\mu_{X'} \in \mathcal{B}_{W_\infty}(\mu_X, \epsilon)} \Pr[f(X') \neq Y]. \quad (\text{A.4})$$

According to the alternative definition of  $\infty$ -Wasserstein distance, the optimal solution  $\mu_{X'}^*$  that achieves the supremum of the right hand side of (A.4) can be captured by a transport map  $T^* : \mathcal{X} \rightarrow \mathcal{X}$  such that  $\mu_{X'}^* = T_{\#}^*(\mu_X)$  and  $\Delta(T^*(X), X) \leq \epsilon$  holds almost surely with respect to the randomness of  $X$  and  $T^*$ . Thus, we have

$$\begin{aligned} \Pr[f(T^*(X)) \neq Y] &= \Pr_{(\mathbf{x}, y) \sim \mu_{XY}} [f(T^*(\mathbf{x})) \neq y] \\ &= \Pr_{(\mathbf{x}, y) \sim \mu_{XY}} [\Delta(T^*(\mathbf{x}), \mathbf{x}) \leq \epsilon \text{ and } f(T^*(\mathbf{x})) \neq y] \\ &\leq 1 - \Pr_{(\mathbf{x}, y) \sim \mu_{XY}} [\forall \mathbf{x}' \in \mathcal{B}(\mathbf{x}, \epsilon) \text{ s.t. } f(\mathbf{x}') = y] = \text{AdvRisk}_\epsilon(f). \end{aligned}$$

Therefore, we have proven the second direction and completed the proof.  $\square$

### A.3. Proof of Theorem 3.4

Theorem 3.4, restated below, gives a lower bound for the adversarial risk for any downstream classifier in terms of the worst-case mutual information between the representation's input and output distributions.

**Theorem 3.4.** Let  $(\mathcal{X}, \Delta)$  be the input metric space,  $\mathcal{Y}$  be the set of labels and  $\mu_{XY}$  be the underlying joint probability distribution. Assume the marginal distribution of labels  $\mu_Y$  is a uniform distribution over  $\mathcal{Y}$ . Consider the feature space  $\mathcal{Z}$  and the set of downstream classifiers  $\mathcal{H} = \{h : \mathcal{Z} \rightarrow \mathcal{Y}\}$ . Given  $\epsilon \geq 0$ , for any  $g : \mathcal{X} \rightarrow \mathcal{Z}$ , we have

$$\inf_{h \in \mathcal{H}} \text{AdvRisk}_\epsilon(h \circ g) \geq 1 - \frac{I(X; Z) - \text{RV}_\epsilon(g) + \log 2}{\log |\mathcal{Y}|},$$

where  $X$  is the random variable that follows the marginal distribution of inputs  $\mu_X$  and  $Z = g(X)$ .

Before starting the proof, we state two useful lemmas on Markov chains. A Markov chain is defined to be a collection of random variables  $\{X_t\}_{t \in \mathbb{Z}}$  with the property that given the present, the future is conditionally independent of the past. Namely,

$$\Pr(X_t = j | X_0 = i_0, X_1 = i_1, \dots, X_{(t-1)} = i_{(t-1)}) = \Pr(X_t = j | X_{(t-1)} = i_{(t-1)}).$$

**Lemma A.1** (Fano's Inequality). Let  $X$  be a random variable uniformly distributed over a finite set of outcomes  $\mathcal{X}$ . For any estimator  $\hat{X}$  such that  $X \rightarrow Y \rightarrow \hat{X}$  forms a Markov chain, we have

$$\Pr(\hat{X} \neq X) \geq 1 - \frac{I(X; \hat{X}) - \log 2}{\log |\mathcal{X}|}.$$

**Lemma A.2** (Data-Processing Inequality). For any Markov chain  $X \rightarrow Y \rightarrow Z$ , we have

$$I(X; Y) \geq I(X; Z) \quad \text{and} \quad I(Y; Z) \geq I(X; Z).$$

Chapter 2 in Cover & Thomas (2012) provides proofs of Lemmas A.1 and A.2.

*Proof of Theorem 3.4.* For any classifier  $h : \mathcal{Z} \rightarrow \mathcal{Y}$ , according to Lemma 3.3, we have

$$\text{AdvRisk}_\epsilon(h \circ g) = \sup_{\mu_{X'} \in \mathcal{B}_{W_\infty}(\mu_X, \epsilon)} \Pr [h(g(X')) \neq Y]. \quad (\text{A.5})$$

Let  $\mu_{X'} \in \mathcal{B}_{W_\infty}(\mu_X, \epsilon)$  be a probability measure over  $(\mathcal{X}, \Delta)$ . According to the alternative definition of  $\infty$ -Wasserstein distance using optimal transport,  $\mu_{X'}$  corresponds to a transport map  $T : \mathcal{X} \rightarrow \mathcal{X}$  such that  $\mu_{X'} = T_\#(\mu_X)$ . Thus, for any given  $\mu_{X'} \in \mathcal{B}_{W_\infty}(\mu_X, \epsilon)$  and  $h \in \mathcal{H}$ , we have the Markov chain

$$Y \rightarrow X \xrightarrow{T} X' \xrightarrow{g} g(X') \xrightarrow{h} (h \circ g)(X').$$

where  $X, Y$  are random variables for input and label distributions respectively. The first Markov chain  $Y \rightarrow X$  can be understood as a generative model for generating inputs according to the conditional probability distribution  $\mu_{X|Y}$ . Therefore, applying Lemmas A.1 and A.2, we obtain the inequality,

$$\Pr [h(g(X')) \neq Y] \geq 1 - \frac{\mathbb{I}(Y; (h \circ g)(X')) + \log 2}{\log |\mathcal{Y}|} \geq 1 - \frac{\mathbb{I}(X'; g(X')) + \log 2}{\log |\mathcal{Y}|}. \quad (\text{A.6})$$

Taking the supremum over the distribution of  $X'$  in  $\mathcal{B}_{W_\infty}(\mu_X, \epsilon)$  and infimum over  $h \in \mathcal{H}$  on both sides of (A.6) yields

$$\begin{aligned} \inf_{h \in \mathcal{H}} [\text{AdvRisk}_\epsilon(h \circ g)] &= \inf_{h \in \mathcal{H}} \sup_{\mu_{X'} \in \mathcal{B}_{W_\infty}(\mu_X, \epsilon)} \Pr [h(g(X')) \neq Y] \\ &\geq 1 - \frac{\inf_{\mu_{X'} \in \mathcal{B}_{W_\infty}(\mu_X, \epsilon)} \mathbb{I}(X'; g(X')) + \log 2}{\log |\mathcal{Y}|} \\ &= 1 - \frac{\mathbb{I}(X; g(X)) - \text{RV}_\epsilon(g) + \log 2}{\log |\mathcal{Y}|}, \end{aligned}$$

where the first equality is due to (A.5) and the inequality holds because of (A.6). Thus, we completed the proof.  $\square$

## B. Algorithm for Estimating the Worst-Case Mutual Information

This section presents the pseudocode of our heuristic algorithm for solving the empirical estimation problem (4.3). More specifically, given a training sample set  $\mathcal{S}_{\text{train}}$ , our algorithm alternatively optimizes for the worst-case input perturbations using projected gradient descent (Algorithm 1) and conducts gradient ascent for the network parameters  $\theta$  (training phase in Algorithm 2). Based on the best parameter  $\theta_{\text{opt}}$  selected from the training phase, our algorithm then estimates the worst-case mutual information with respect to the given representation  $g$  using a testing sample set  $\mathcal{S}_{\text{test}}$  (testing phase in Algorithm 2). Since we only have access to a finite set of data sampled from  $\mu_X$ , we use an additional testing phase in Algorithm 2 to minimize the overfitting effect of the training samples on the optimal network parameter  $\theta_{\text{opt}}$  for mutual information neural estimation (MINE).

Moreover, we adopt the negative sampling scheme (Hjelm et al., 2018) to estimate the expectation term with respect to  $\widehat{\mu}_X^{(m)} \otimes \widehat{\mu}_Z^{(m)}$  in mutual information neural estimation for better performance. Here, the pairing scheme defines a correspondence from each input to a set of inputs for a given sample set. To be more specific, given a set of samples  $\{\mathbf{x}_i\}_{i \in [B]}$ , a pairing scheme with negative sampling size  $N \leq B$  corresponds to a set of vectors  $\{\boldsymbol{\pi}_i\}_{i \in [B]}$  such that each  $\boldsymbol{\pi}_i$  is a randomly selected subset from  $\{1, 2, \dots, B\}$  with size  $N$ , and  $\pi_{ij}$  denotes the  $j$ -th element of  $\boldsymbol{\pi}_i$ . Compared with the algorithm in Hjelm et al. (2018) for estimating standard mutual information, Algorithm 2 requires additional  $B \cdot S$  steps of forward and backward propagations with respect to the input for finding the worst-case input perturbations in each iteration.

## C. Worst-case Mutual Information for Individual Neuron Features

The following tensorization inequality (Scarlett & Cevher, 2019) characterizes the connection between the mutual information of individual neuron features and that of the whole representation. According to Theorem 3.4, such connection suggests the necessity of enough worst-case mutual information for each individual neuron.

---

**Algorithm 1** Heuristic Search for Worst-Case Input Perturbations

---

**Input:** samples  $\{\mathbf{x}_i\}_{i \in [B]}$ , representation  $g$ , MINE estimator  $T_\theta$ , pairing scheme  $\{\pi_i\}_{i \in [B]}$ , perturbation strength  $\epsilon$  in  $\ell_p$

**Hyperparameters:** negative sampling size  $N$ , number of iterations  $S$ , step size  $\eta_a$

- 1: Initialize  $\{\mathbf{x}'_i\}_{i \in [B]} \leftarrow \{\mathbf{x}_i\}_{i \in [B]}$
- 2: **for**  $s = 1, 2, \dots, S$  **do**
- 3:  $J(\mathbf{x}'_1, \dots, \mathbf{x}'_B, \theta) \leftarrow \frac{1}{B} \sum_{i=1}^B T_\theta(\mathbf{x}'_i, g(\mathbf{x}'_i)) - \log \left( \frac{1}{BN} \sum_{i=1}^B \sum_{j=1}^N \exp [T_\theta(\mathbf{x}'_i, g(\mathbf{x}'_{\pi_{ij}}))] \right)$
- 4: **for**  $i = 1, 2, \dots, B$  **do**
- 5:  $\mathbf{x}'_i \leftarrow \mathcal{P}_{\mathcal{B}(\mathbf{x}_i, \epsilon)}[\mathbf{x}'_i - \eta_a \cdot \nabla_{\mathbf{x}'_i} J(\mathbf{x}'_1, \dots, \mathbf{x}'_B, \theta)]$  //  $\mathcal{P}_{\mathcal{B}(\mathbf{x}_i, \epsilon)}$  denotes the projection operator onto  $\mathcal{B}(\mathbf{x}_i, \epsilon)$
- 6: **end for**
- 7: **end for**
- 8:  $V_1 \leftarrow J(\mathbf{x}'_1, \dots, \mathbf{x}'_B, \theta)$
- 9:  $V_2 \leftarrow \nabla_\theta J(\mathbf{x}'_1, \dots, \mathbf{x}'_B, \theta)$

**Output:**  $\{V_1, V_2\}$

---

**Algorithm 2** Empirical Estimation of Worst-Case Mutual Information

---

**Input:** training and testing sample sets  $(\mathcal{S}_{\text{train}}, \mathcal{S}_{\text{test}})$  sampled from  $\mu_X$ , representation  $g$ , perturbation strength  $\epsilon$  in  $\ell_p$

**Hyperparameters:** number of training epochs  $T$ , step size  $\eta_e$ , number of testing mini-batches  $K$

- 1: // **Training Phase**
- 2:  $\theta_1 \leftarrow$  initialize network parameter for MINE estimator
- 3: **for**  $t = 1, 2, \dots, T$  **do**
- 4:  $\{\mathbf{x}_i\}_{i \in [B]}, \{\pi_i\}_{i \in [B]} \leftarrow$  randomly generate a batch of  $B$  training samples and a pairing scheme
- 5:  $\{V_1(t), V_2(t)\} \leftarrow$  Algorithm 1  $(\{\mathbf{x}_i\}_{i \in [B]}, g, T_{\theta_t}, \{\pi_i\}_{i \in [B]}, \epsilon)$
- 6:  $\theta_{t+1} \leftarrow \theta_t + \eta_e \cdot V_2(t)$
- 7: **end for**
- 8:  $\theta_{\text{opt}} \leftarrow \operatorname{argmax} \{t \in [T] : V_1(t)\}$  // choose the best parameter  $\theta_{\text{opt}}$  based on history
- 9: // **Testing Phase**
- 10: Randomly split the testing set  $\mathcal{S}_{\text{test}}$  into  $K$  mini-batches  $\{\mathcal{S}_1, \dots, \mathcal{S}_K\}$  with equal size
- 11: **for**  $k = 1, 2, \dots, K$  **do**
- 12:  $\{\pi(\mathbf{x})\}_{\mathbf{x} \in \mathcal{S}_k} \leftarrow$  randomly generate a pairing scheme with respect to  $\mathcal{S}_k$
- 13:  $\{V_1(k), V_2(k)\} \leftarrow$  Algorithm 1  $(\mathcal{S}_k, g, T_{\theta_{\text{opt}}}, \{\pi(\mathbf{x})\}_{\mathbf{x} \in \mathcal{S}_k}, \epsilon)$
- 14: **end for**
- 15:  $\hat{I}_{\text{worst}} \leftarrow \frac{1}{K} \sum_{k=1}^K V_1(k)$

**Output:**  $\hat{I}_{\text{worst}}$

---

**Lemma C.1** (Tensorization of Mutual Information). Let  $\mathbf{Z} = (Z_1, \dots, Z_n)$  be a product distributions over random variables. If  $Z_1, \dots, Z_n$  are mutually independent conditioned on  $X$ , then

$$I(X; \mathbf{Z}) \leq \sum_{i=1}^n I(X; Z_i)$$

Suppose neurons within a single layer have no interconnection, then each neuron's output is mutually independent conditioned on the model input. If a perturbation imposed on the input distribution makes the perturbed mutual information  $I(X'; Z'_i)$  relatively low for each neuron, then the perturbed mutual information with respect to the entire layer  $I(X'; \mathbf{Z}')$  will also be low, which further implies a low adversarial accuracy for any downstream classifier based on Theorem 3.4.

---

## D. Experiments

### D.1. Implementation Details

Here, we provide additional implementation details of our experiments presented in Section 6.

**Model architectures.** For all experiments, we follow Hjelm et al. (2018) in implementing the MINE estimator. We adopt the *encode-and-dot-product* model architecture in Hjelm et al. (2018) which maps  $x$  and  $z$  respectively to two high-dimensional vectors and then takes the dot-product to calculate the output. The basic modules used in our experiments are listed in Table 2. A slight difference in training the feature (encoder) is that Hjelm et al. (2018) shares parameters between parts of the mutual information estimator and the encoder, while we separate the two parts completely to be consistent with our mutual information estimation experiments.

Module	Structure
Feature Extractor	Conv(64, 4 × 4, 2) → Conv(128, 4 × 4, 2) → Conv(256, 4 × 4, 2) → FC(1024) → FC(64)
Top Classifier (MLP)	FC(200) → FC(10)
Top Classifier (Linear)	FC(10)
Baseline-H	Feature Extractor → Top Classifier (MLP)
Estimator Part 1	Conv(64, 4 × 4, 2) → Conv(128, 4 × 4, 2) → Conv(256, 4 × 4, 2)
Estimator Part 2	Conv(2048, 1 × 1, 1) → Conv(2048, 1 × 1, 1)
Estimator	( $x$ → Estimator Part 1 → Estimator Part 2) · ( $z$ → Estimator Part 2)

Table 2. Basic model structures used in our experiments. Batch-normalization and ReLU activation are used between layers (not including the output of each module). Shortcut-connection is omitted for Estimator Part 2. For scalar feature  $z$ , Estimator Part 2 is replaced by an identity mapping. Average operation is needed in the dot-product operation of Estimator. For more details, see Hjelm et al. (2018)

**Hyperparameters.** We use simple hyperparameter settings to control their effect on our various ablation experiments. We use  $l_\infty$  constrained perturbations and PGD attack (Madry et al., 2018) on all datasets. For CIFAR-10, we set the radius  $\epsilon = 8/255$  and use 7 attack steps with step size 0.01. For MNIST, we set the radius  $\epsilon = 0.3$  and use 10 attack steps with step size 0.1. For Fashion-MNIST, we set the radius  $\epsilon = 0.1$  and use 10 attack steps with step size 0.02. For SVHN, we set the radius  $\epsilon = 4/255$  and use 10 attack steps with step size 0.005. The batch size is set as 128 for both datasets, and our results are consistent with different batch sizes between 128 to 512 (we did not test other sizes). A total training epochs of 200 is set for VGG, ResNet, and DenseNet, with an initial learning rate of 0.1 which decays by a factor of 10 every 50 epochs. For the Baseline-H model and the similar mutual information estimator, we set the training epoch to 300 and use a fixed learning rate of 0.0001 as in Hjelm et al. (2018).

### D.2. Additional Results

**Results for MNIST, Fashion-MNIST, and SVHN.** We present the downstream classification results for MNIST, Fashion-MNIST, and SVHN in Table 3, 4, 5. These results support similar conclusions as those drawn from CIFAR-10 dataset in Table 1. That is, our training principle always produces representations that have significantly better adversarial accuracy for downstream adversarial classification. In many cases, our training principle also produces representations that have better natural accuracy, despite the worst-case situation that our training principle considers.

**Saliency maps of internal features.** In section 6.1, we evaluated the internal feature vulnerability of all the convolutional kernels in the second layer of Baseline-H. Here, we further visualize the saliency maps of the those internal features to evaluate the underlying correlations. As shown in Figure 5, features in robust model have less noisy saliency maps, which is consistent with the observations of lower representation vulnerability shown in Figure 3.

**Saliency maps of learned representations.** More comparison results of saliency maps are given in Figure 6. The saliency maps of representations learned using our unsupervised training method shows comparable interpretability results to the models learned using fully-supervised adversarial training. Saliency maps computed by different losses also show consistent interpretability results. This indicates that our training principle indeed produces adversarially robust representations.

Representation ( $g$ )	Classifier ( $h$ )	MLP $h$		Linear $h$	
		Natural	Adversarial	Natural	Adversarial
Hjelm et al. (2018)	Standard	<b>96.96 <math>\pm</math> 0.35</b>	0.00 $\pm$ 0.00	<b>84.88 <math>\pm</math> 1.11</b>	0.00 $\pm$ 0.00
Hjelm et al. (2018)	Robust	44.99 $\pm$ 14.49	16.70 $\pm$ 2.22	11.35 $\pm$ 0.00	11.35 $\pm$ 0.00
Ours	Standard	96.67 $\pm$ 0.12	9.97 $\pm$ 1.88	82.10 $\pm$ 0.37	3.69 $\pm$ 0.55
Ours	Standard (E.S.)	94.68 $\pm$ 0.93	12.79 $\pm$ 2.24	77.72 $\pm$ 2.47	4.73 $\pm$ 1.29
Ours	Robust	95.05 $\pm$ 0.19	<b>60.64 <math>\pm</math> 1.82</b>	73.99 $\pm$ 1.16	<b>30.55 <math>\pm</math> 1.34</b>
Fully-Supervised Standard		99.13 $\pm$ 0.23	0.45 $\pm$ 0.33	99.13 $\pm$ 0.04	0.00 $\pm$ 0.00
Fully-Supervised Robust		99.25 $\pm$ 0.05	95.73 $\pm$ 0.09	99.21 $\pm$ 0.06	95.29 $\pm$ 0.18

Table 3. Comparisons of different representation learning methods for downstream classification on MNIST. *E.S.* denotes early stopping under the criterion of the best adversarial accuracy. We present the mean accuracy and the standard deviation over 4 repeated trials.

Representation ( $g$ )	Classifier ( $h$ )	MLP $h$		Linear $h$	
		Natural	Adversarial	Natural	Adversarial
Hjelm et al. (2018)	Standard	89.58 $\pm$ 0.13	0.00 $\pm$ 0.00	85.93 $\pm$ 0.26	0.00 $\pm$ 0.00
Hjelm et al. (2018)	Robust	48.61 $\pm$ 4.96	14.95 $\pm$ 0.79	10.00 $\pm$ 0.00	10.00 $\pm$ 0.00
Ours	Standard	<b>90.45 <math>\pm</math> 0.19</b>	5.38 $\pm$ 1.00	<b>87.37 <math>\pm</math> 0.10</b>	18.20 $\pm$ 2.87
Ours	Standard (E.S.)	81.66 $\pm$ 0.18	29.71 $\pm$ 2.00	86.27 $\pm$ 0.64	23.40 $\pm$ 2.65
Ours	Robust	84.31 $\pm$ 0.29	<b>70.44 <math>\pm</math> 3.62</b>	81.05 $\pm$ 0.30	<b>61.33 <math>\pm</math> 0.49</b>
Fully-Supervised Standard		92.09 $\pm$ 0.23	0.00 $\pm$ 0.00	85.93 $\pm$ 0.26	0.00 $\pm$ 0.00
Fully-Supervised Robust		87.94 $\pm$ 0.18	77.59 $\pm$ 0.38	88.05 $\pm$ 0.46	77.15 $\pm$ 0.24

Table 4. Comparisons of different representation learning methods for downstream classification on Fashion-MNIST. *E.S.* denotes early stopping under the criterion of the best adversarial accuracy. We present the mean accuracy and the standard deviation over 4 repeated trials.

Representation ( $g$ )	Classifier ( $h$ )	MLP $h$		Linear $h$	
		Natural	Adversarial	Natural	Adversarial
Hjelm et al. (2018)	Standard	50.15 $\pm$ 0.89	0.00 $\pm$ 0.00	38.94 $\pm$ 1.52	0.00 $\pm$ 0.00
Hjelm et al. (2018)	Robust	19.59 $\pm$ 0.00	19.59 $\pm$ 0.00	19.59 $\pm$ 0.00	19.59 $\pm$ 0.00
Ours	Standard	<b>74.32 <math>\pm</math> 0.49</b>	26.29 $\pm$ 1.41	<b>58.37 <math>\pm</math> 0.54</b>	21.62 $\pm$ 0.91
Ours	Standard (E.S.)	71.85 $\pm$ 0.59	29.59 $\pm$ 0.83	54.76 $\pm$ 0.86	25.00 $\pm$ 0.53
Ours	Robust	68.25 $\pm$ 0.83	<b>40.23 <math>\pm</math> 0.83</b>	49.04 $\pm$ 0.79	<b>30.56 <math>\pm</math> 0.38</b>
Fully-Supervised Standard		91.97 $\pm$ 0.13	9.77 $\pm$ 1.58	91.33 $\pm$ 0.15	9.29 $\pm$ 1.73
Fully-Supervised Robust		90.14 $\pm$ 0.83	65.35 $\pm$ 0.44	89.60 $\pm$ 0.54	64.48 $\pm$ 1.06

Table 5. Comparisons of different representation learning methods for downstream classification on SVHN. *E.S.* denotes early stopping under the criterion of the best adversarial accuracy. We present the mean accuracy and the standard deviation over 4 repeated trials.



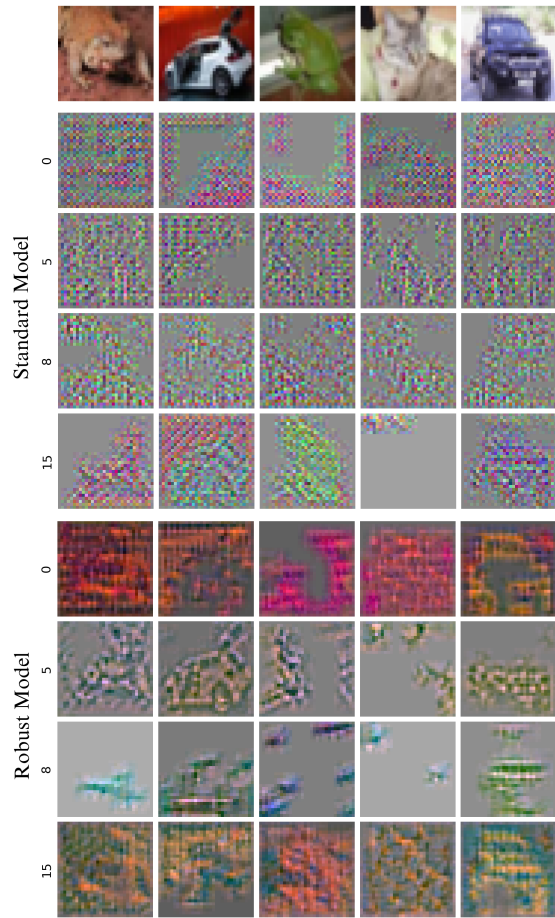


Figure 5. Saliency maps of four arbitrarily selected features in the second convolutional layer of Baseline-H. The feature saliency map is computed by the gradient of a kernel’s averaged activation over a input image. Each row presents saliency maps of a specific convolutional kernel.

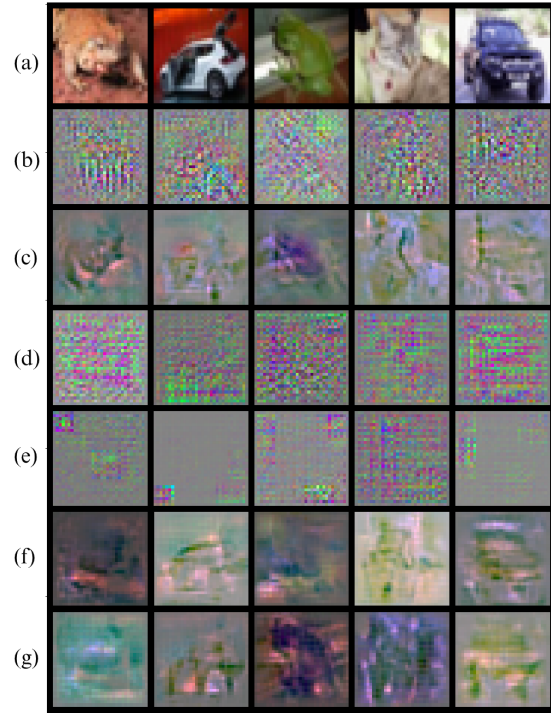


Figure 6. Saliency maps of different models on CIFAR-10: (a) original images (b) fully-supervised standard model (c) fully-supervised robust model (d) representations learned using Hjelm et al. (2018) with cross-entropy loss (e) representations learned using Hjelm et al. (2018) with mutual information maximization loss (f) representations learned using our method with cross-entropy loss (g) representations learned using our method with mutual information maximization loss.