

## A Proofs

### A.1 Proof and discussion of Theorem 1

**Theorem 1.** Given a random vector  $Y = (y_1, y_2, \dots, y_d)$  where  $y_i$  follows a positive-valued distribution, and two arbitrary vectors with the same dimension,  $X, X' \in \mathbb{R}^d$  that  $x_i, x'_i \geq 0$ , assume that there exists a sequence  $\mathcal{S} = \{X_i\}_{i=1}^N$  with  $X = X_0$  and  $X' = X_N$ , where the vectors satisfy the condition that  $\cos(X_i, Y) \geq \cos(X_{i+1}, Y)$ , and each  $X_{i+1}$  can be induced from its previous vector  $X_i$  through one of the following two operations,

- (i) arbitrarily exchanging two entries of  $X_i$
- (ii) multiplying one entry in  $X_i$  by  $\alpha \in (0, 1]$

Then Kendall's rank correlations of  $Y$  with  $X$  and  $X'$  have the property that  $\mathbb{E}[\tau(X, Y)] \geq \mathbb{E}[\tau(X', Y)]$ , where the expectation is taken over  $Y$  satisfying  $\cos(X_i, Y) \geq \cos(X_{i+1}, Y)$ .

*Proof.* To prove this theorem, we show that the property holds when  $N = 2$ , i.e.,  $\mathcal{S} = \{X, X'\}$ , which indicates that each one of the above operations on  $X$  would preserve the order of Kendall's rank correlation. The case when  $N \geq 3$  can be trivially generalized using mathematical induction.

Since  $X$  and  $X'$  are two arbitrary vectors, it is safe to fix  $X$  that  $X = (x_1, x_2, \dots, x_d)$ . To analyze the cosine similarities and Kendall's rank correlation,  $X$  can be assumed to be in descending order, i.e.,  $x_1 > x_2 > \dots > x_d$ , since the order of  $X'$  and  $Y$  can be changed correspondingly without affecting the cosine similarities and Kendall's rank correlation. Formally, we show in the following proof that for a random vector  $Y$  following exponential distribution and an arbitrary vector  $X$ , if the cosine similarities satisfy that  $\cos(X, Y) \geq \cos(X', Y)$ , then their corresponding Kendall's rank correlations have the property that  $\mathbb{E}[\tau(X, Y)] \geq \mathbb{E}[\tau(X', Y)]$ , where  $X'$  is generated from  $X$  by (1) exchanging two entries and (2) scalar multiplications.

**(1) Preservation under exchanging** Following the assumption, we consider a random vector  $Y = (y_1, y_2, \dots, y_d)$  where  $y_i$  positively distributed, and another vector  $X = (x_1, x_2, \dots, x_d)$ . We define the new vector  $X'$  that is produced by arbitrarily exchanging two entries in  $X$ . Suppose we exchange the  $p$ -th and  $q$ -th entry in  $X$ , where  $1 \leq p < q \leq d$ , then  $X' = (x_1, \dots, x_{p-1}, x_q, x_{p+1}, \dots, x_{q-1}, x_p, x_{q+1}, \dots, x_d)$ . We further assume both  $X$  and  $Y$  are normalized, i.e.,  $\|X\| = \|Y\| = \|X'\| = 1$ .

Now if we consider the cosine similarity and the assumption that  $\cos(X, Y) > \cos(X', Y)$ , we then have

$$\cos(X, Y) = \sum_{i=1}^d x_i y_i > \cos(X', Y) = \sum_{i \neq p, q}^d x_i y_i + x_p y_q + x_q y_p, \quad (8)$$

which can be simplified as

$$(x_p - x_q)(y_p - y_q) > 0 \quad (9)$$

For Kendall's rank correlation, we denote that  $\Omega(X, Y) = \sum_{i < j} \text{sign}(x_i - x_j) \text{sign}(y_i - y_j)$ , and  $\tau(X, Y) = \frac{2}{d(d-1)} \Omega(X, Y)$ . We notice that the difference between  $\Omega(X, Y)$  and  $\Omega(X', Y)$  only occurs when  $x_p$  and  $x_q$  are involved. We can write down the explicit expression of  $\Omega(X, Y) - \Omega(X', Y)$

$$\begin{aligned} \Omega(X, Y) - \Omega(X', Y) &= (\text{sign}(x_p - x_q) - \text{sign}(x_q - x_p)) \text{sign}(y_p - y_q) \\ &\quad + \sum_{p < i < q} (\text{sign}(x_p - x_i) - \text{sign}(x_q - x_i)) \text{sign}(y_p - y_i) \\ &\quad + \sum_{p < i < q} (\text{sign}(x_i - x_q) - \text{sign}(x_i - x_p)) \text{sign}(y_i - y_q) \end{aligned} \quad (10)$$

$$= 2 + 2 \sum_{p < i < q} (\text{sign}(y_p - y_i) + \text{sign}(y_i - y_q)) \quad (11)$$

Since

$$\mathbb{E} \left[ \sum_{p < i < q} \text{sign}(y_p - y_i) + \text{sign}(y_i - y_q) \middle| y_p - y_q > 0 \right] \geq 0, \quad (12)$$

we then have,

$$\mathbb{E} [\Omega(X, Y)] \geq \mathbb{E} [\Omega(X', Y)]. \quad (13)$$

**(2) Preservation under scalar multiplication** We use the assumptions mentioned before that  $y_i$  is positive-valued, and  $x_1 > x_2 > \dots > x_d > 0$ . Without loss of generality, we multiply  $x_1$  by a scalar  $\alpha \in [0, 1]$ , such that  $x_2 > \alpha x_1 > x_3$ , *i.e.*,  $X' = (\alpha x_1, x_2, \dots, x_d)$ .

To compare  $\tau(X, Y)$  and  $\tau(X', Y)$ , it is noticed that, under our assumptions, only the sign of  $y_1 - y_2$  is needed as other terms in  $\Omega$  are equal for  $\Omega(X, Y)$  and  $\Omega(X', Y)$ ,

$$\Omega(X, Y) - \Omega(X', Y) = \text{sign}(x_1 - x_2) \text{sign}(y_1 - y_2) - \text{sign}(\alpha x_1 - x_2) \text{sign}(y_1 - y_2) = 2 \text{sign}(y_1 - y_2). \quad (14)$$

Under the condition that the cosine similarity  $\cos(X, Y) > \cos(X', Y)$ , we have

$$x_1 y_1 + x_2 y_2 + \dots + x_d y_d \geq \frac{\alpha x_1 y_1 + x_2 y_2 + \dots + x_d y_d}{\sqrt{\alpha^2 x_1^2 + x_2^2 + \dots + x_d^2}}. \quad (15)$$

Note that  $\|X\| = \|Y\| = 1$ . For simplicity, we denote that  $A = \sqrt{\alpha^2 x_1^2 + x_2^2 + \dots + x_d^2}$ , and it is obvious that  $\alpha \leq A \leq 1$ , where  $A$  is close to 1 when  $d$  is large,

$$x_1 y_1 + x_2 y_2 + \left(1 - \frac{1}{A}\right) \left(\sum_{i=3}^d x_i y_i\right) \geq \frac{\alpha x_1 y_1 + x_2 y_2}{A}, \quad (16)$$

which can be relaxed as

$$x_1 y_1 + x_2 y_2 \geq \frac{\alpha x_1 y_1 + x_2 y_2}{A}, \quad (17)$$

*i.e.*,

$$y_1 \geq K y_2 \quad (18)$$

where  $K = \frac{(1-A)x_2}{(A-\alpha)x_1} > 0$ . Thus, we want to show that

$$\mathbb{E} [\text{sign}(y_1 - y_2) | y_1 \geq K y_2] = \mathbb{P}(\text{sign}(y_1 - y_2) \geq 0 | y_1 \geq K y_2) - \mathbb{P}(\text{sign}(y_1 - y_2) < 0 | y_1 \geq K y_2) \quad (19)$$

$$= 2\mathbb{P}(\text{sign}(y_1 - y_2) \geq 0 | y_1 \geq K y_2) - 1 \geq 0 \quad (20)$$

We consider two cases when  $K \geq 1$  and  $0 < K < 1$ . In the case that  $K \geq 1$ , it is obvious that

$$\mathbb{P}(\text{sign}(y_1 - y_2) \geq 0 | y_1 \geq K y_2) = 1. \quad (21)$$

If  $0 < K < 1$ ,

$$\mathbb{P}(\text{sign}(y_1 - y_2) \geq 0 | y_1 \geq K y_2) = \frac{\mathbb{P}(\text{sign}(y_1 - y_2) \geq 0 \cap y_1 \geq K y_2)}{\mathbb{P}(y_1 \geq K y_2)} \quad (22)$$

$$= \frac{\mathbb{P}(\text{sign}(y_1 - y_2) \geq 0)}{\mathbb{P}(y_1 \geq K y_2)} \geq \frac{1}{2} \quad (23)$$

Thus, after combining the above two cases, we have

$$\mathbb{P}(\text{sign}(y_1 - y_2) \geq 0 | y_1 \geq K y_2) \geq \frac{1}{2} \quad (24)$$

which concludes our proof.  $\square$

**Discussions** In practice, the attribution values are taken absolute values to emphasize the importance of features, regardless of whether the impact are positive or negative. Thus, without loss of generality,  $y_i$  is assumed to follow a positive-valued distribution in Theorem 1. We also consider the existence of the sequence  $\mathcal{S}$  as an assumption that assist the formulation of the theorem. Although searching for such sequence of every pair of attributions  $X$  and  $X'$  can be a combinatorial problem and is constrained by computation power, the numerical simulations of finding such sequences in lower dimensions still show a high success rate ( $\geq 0.8$  when  $d \leq 10$ ), and the number of possible sequences increases drastically when the dimension is higher.

## A.2 Proof of Proposition 1

**Proposition 1.** *Given a single-layer neural network with ReLU activation, and with the above parameterization, if, for all  $i$ ,  $\mathbf{W}_i$  and  $u_i$  are all independent and identically distributed random variables following Gaussian distributions, i.e.,  $\mathbf{W}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_w^2 I_d)$  and  $u_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_u^2)$ , and two input images that each has small variance,  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$ , then*

$$\cos(\text{IG}(\mathbf{x}), \text{IG}(\tilde{\mathbf{x}})) \approx \frac{\mathbb{P}(\mathbf{W}^\top \mathbf{x} > 0 \cap \mathbf{W}^\top \tilde{\mathbf{x}} > 0)}{\sqrt{\mathbb{P}(\mathbf{W}^\top \mathbf{x} > 0)\mathbb{P}(\mathbf{W}^\top \tilde{\mathbf{x}} > 0)}}. \quad (6)$$

*Proof.* Recall that  $\mathbf{x} \in \mathbb{R}^d$  is an input image, and the network function  $f$  is parameterized by  $(\mathbf{W}, \mathbf{u}, c) \in \mathbb{R}^{d \times m} \times \mathbb{R}^m \times \mathbb{R}$ , where  $\mathbf{W}_i$  is the column vector of  $\mathbf{W}$ ,  $w_{ij}$  is the  $ij$ -th entry of matrix  $\mathbf{W}$  and  $u_i$  is the  $i$ -th entry of vector  $\mathbf{u}$ , i.e.,  $f(\mathbf{x}) = \mathbf{u}^\top \text{ReLU}(\mathbf{W}^\top \mathbf{x}) + c$ .

Following the above notations, we first write the function as

$$f(\mathbf{x}) = \mathbf{u}^\top \text{ReLU}(\mathbf{W}^\top \mathbf{x}) + c = \sum_{i=1}^m u_i (\mathbf{W}_i^\top \mathbf{x}) \mathbb{1}_{\mathbf{W}_i^\top \mathbf{x} > 0} + c, \quad (25)$$

where  $\mathbb{1}_{\{\cdot\}}$  denotes the indicator function, and its gradient

$$\nabla_{x_k} f(\mathbf{x}) = (\nabla f(\mathbf{x}))_k = \sum_{i=1}^m u_i w_{ki} \mathbb{1}_{\mathbf{W}_i^\top \mathbf{x} > 0}$$

We assume the bias terms are zeros without loss of generality, i.e.,  $c = 0$ , and approximate the cosine similarity of IG using the small variance assumption that  $\frac{1}{n} \sum_i x_i^2 - (\frac{1}{n} \sum_i x_i)^2 \approx 0$  as

$$\begin{aligned} & \cos(\text{IG}(\mathbf{x}), \text{IG}(\tilde{\mathbf{x}})) \\ & \approx \frac{\sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{W}_i, \mathbf{W}_j \rangle u_i u_j \int_0^1 \mathbb{1}_{\mathbf{W}_i^\top (r\mathbf{x}) > 0} dr \int_0^1 \mathbb{1}_{\mathbf{W}_j^\top (r\tilde{\mathbf{x}}) > 0} dr}{\sqrt{\sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{W}_i, \mathbf{W}_j \rangle u_i u_j \int_0^1 \mathbb{1}_{\mathbf{W}_i^\top (r\mathbf{x}) > 0} dr} \sqrt{\sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{W}_i, \mathbf{W}_j \rangle u_i u_j \int_0^1 \mathbb{1}_{\mathbf{W}_i^\top (r\tilde{\mathbf{x}}) > 0} dr}} \quad (26) \end{aligned}$$

Since  $\langle \mathbf{W}_i, \mathbf{W}_j \rangle$  is close to 0 in high dimensional space when  $i \neq j$ , we approximate the above expression as

$$\begin{aligned} & \frac{\sum_{i=1}^m \|\mathbf{W}_i\|_2^2 u_i^2 \int_0^1 \mathbb{1}_{\mathbf{W}_i^\top (r\mathbf{x}) > 0} dr \int_0^1 \mathbb{1}_{\mathbf{W}_i^\top (r\tilde{\mathbf{x}}) > 0} dr}{\sqrt{\sum_{i=1}^m \|\mathbf{W}_i\|_2^2 u_i^2 \int_0^1 \mathbb{1}_{\mathbf{W}_i^\top (r\mathbf{x}) > 0} dr} \sqrt{\sum_{i=1}^m \|\mathbf{W}_i\|_2^2 u_i^2 \int_0^1 \mathbb{1}_{\mathbf{W}_i^\top (r\tilde{\mathbf{x}}) > 0} dr}} \quad (27) \end{aligned}$$

Notice that the indicator function is integrated from 0 to 1, which does not affect the sign of  $\mathbf{W}_i^\top (r\mathbf{x})$  and  $\mathbf{W}_i^\top (r\tilde{\mathbf{x}})$ , i.e., the activation states. This implies that the activation states is the same for all

samples from baseline to the corresponding image. Thus, we can write the cosine similarity as

$$\frac{\sum_{i=1}^m \|\mathbf{W}_i\|_2^2 u_i^2 \mathbb{1}_{\mathbf{W}_i^\top \mathbf{x} > 0} \mathbb{1}_{\mathbf{W}_i^\top \tilde{\mathbf{x}} > 0}}{\sqrt{\sum_{i=1}^m \|\mathbf{W}_i\|_2^2 u_i^2 \mathbb{1}_{\mathbf{W}_i^\top \mathbf{x} > 0}} \sqrt{\sum_{i=1}^m \|\mathbf{W}_i\|_2^2 u_i^2 \mathbb{1}_{\mathbf{W}_i^\top \tilde{\mathbf{x}} > 0}}} \quad (28)$$

Since  $\mathbf{W}_i$  and  $u_i$  are independent random variables following Gaussian distributions, *i.e.*,  $\mathbf{W}_i \sim \mathcal{N}(0, \sigma_w^2 I_d)$  and  $u_i \sim \mathcal{N}(0, \sigma_u^2)$ , when  $m$  is sufficiently large, we have

$$\frac{1}{m} \sum_{i=1}^m \|\mathbf{W}_i\|_2^2 u_i^2 \mathbb{1}_{\mathbf{W}_i^\top \mathbf{x} > 0} \mathbb{1}_{\mathbf{W}_i^\top \tilde{\mathbf{x}} > 0} = \mathbb{E}_{W,u} [\|W\|_2^2 u^2 \mathbb{1}_{W^\top \mathbf{x} > 0} \mathbb{1}_{W^\top \tilde{\mathbf{x}} > 0}] \quad (29)$$

The cosine similarity is then transformed into expectations

$$\cos(\text{IG}(\mathbf{x}), \text{IG}(\tilde{\mathbf{x}})) \approx \frac{\mathbb{E}_{W,u} [\|W\|_2^2 u^2 \mathbb{1}_{W^\top \mathbf{x} > 0} \mathbb{1}_{W^\top \tilde{\mathbf{x}} > 0}]}{\sqrt{\mathbb{E}_{W,u} [\|W\|_2^2 u^2 \mathbb{1}_{W^\top \mathbf{x} > 0}] \mathbb{E}_{W,u} [\|W\|_2^2 u^2 \mathbb{1}_{W^\top \tilde{\mathbf{x}} > 0}]}} \quad (30)$$

Based on the assumption on Gaussian distribution, we have  $\mathbb{E} [\|W\|_2^2] = \text{tr}(\sigma_w^2 I_d) = d\sigma_w^2$  and  $\mathbb{E} [u^2] = \sigma_u^2$ , and

$$\cos(\text{IG}(\mathbf{x}), \text{IG}(\tilde{\mathbf{x}})) \approx \frac{\sigma_u^2 \mathbb{E}_W [\|W\|_2^2 \mathbb{1}_{W^\top \mathbf{x} > 0} \mathbb{1}_{W^\top \tilde{\mathbf{x}} > 0}]}{\sigma_u \sqrt{\mathbb{E}_W [\|W\|_2^2 \mathbb{1}_{W^\top \mathbf{x} > 0}] \sigma_u \sqrt{\mathbb{E}_W [\|W\|_2^2 \mathbb{1}_{W^\top \tilde{\mathbf{x}} > 0}]}} \quad (31)$$

$$= \frac{\mathbb{E}_W [\|W\|_2^2 \mathbb{1}_{W^\top \mathbf{x} > 0} \mathbb{1}_{W^\top \tilde{\mathbf{x}} > 0}]}{\sqrt{\mathbb{E}_W [\|W\|_2^2 \mathbb{1}_{W^\top \mathbf{x} > 0}] \mathbb{E}_W [\|W\|_2^2 \mathbb{1}_{W^\top \tilde{\mathbf{x}} > 0}]}} \quad (32)$$

$$= \frac{d\sigma_w^2 \mathbb{P}(W^\top \mathbf{x} > 0 \cap W^\top \tilde{\mathbf{x}} > 0)}{d\sigma_w^2 \sqrt{\mathbb{P}(W^\top \mathbf{x} > 0) \mathbb{P}(W^\top \tilde{\mathbf{x}} > 0)}} \quad (33)$$

$$= \frac{\mathbb{P}(W^\top \mathbf{x} > 0 \cap W^\top \tilde{\mathbf{x}} > 0)}{\sqrt{\mathbb{P}(W^\top \mathbf{x} > 0) \mathbb{P}(W^\top \tilde{\mathbf{x}} > 0)}} \quad (34)$$

□

## B Unstable Pearson's correlation

In this section, we discuss the unstable Pearson's correlation for small variance inputs, *i.e.*,  $\frac{1}{n} \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i\right)^2 \approx 0$ . Consider the Pearson's correlation between  $\mathbf{x}$  and  $\mathbf{x} + \boldsymbol{\eta}$ , where  $\boldsymbol{\eta}$  is vector and bounded by  $\|\boldsymbol{\eta}\| \leq \epsilon$  for small  $\epsilon$ . Then the Pearson's correlation between  $\mathbf{x}$  and  $\mathbf{x} + \boldsymbol{\eta}$  can be written as

$$\rho(\mathbf{x}, \mathbf{x} + \boldsymbol{\eta}) = \frac{\frac{1}{n} \sum_{i=1}^n x_i(x_i + \eta_i) - \left(\frac{1}{n} \sum_{i=1}^n x_i\right) \left(\frac{1}{n} \sum_{i=1}^n (x_i + \eta_i)\right)}{\sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i\right)^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i + \eta_i)^2 - \left(\frac{1}{n} \sum_{i=1}^n (x_i + \eta_i)\right)^2}} \quad (35)$$

Consider the numerator of  $\rho(\mathbf{x}, \mathbf{x} + \boldsymbol{\eta})$  as

$$N_{\rho(\mathbf{x}, \mathbf{x} + \boldsymbol{\eta})} = \frac{1}{n} \sum_{i=1}^n x_i^2 + \frac{1}{n} \sum_{i=1}^n x_i \eta_i - \left(\frac{1}{n} \sum_{i=1}^n x_i\right)^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i\right) \left(\frac{1}{n} \sum_{i=1}^n \eta_i\right) \quad (36)$$

$$\approx \frac{1}{n} \sum_{i=1}^n x_i \eta_i - \left(\frac{1}{n} \sum_{i=1}^n x_i\right) \left(\frac{1}{n} \sum_{i=1}^n \eta_i\right) \quad (37)$$

Similarly, we can obtain

$$N_{\rho(\mathbf{x}, \mathbf{x} - \boldsymbol{\eta})} \approx -\frac{1}{n} \sum_{i=1}^n x_i \eta_i + \left(\frac{1}{n} \sum_{i=1}^n x_i\right) \left(\frac{1}{n} \sum_{i=1}^n \eta_i\right) \quad (38)$$

---

**Algorithm 1** Adversarial Training with IGR

---

**Input:** classifier  $f$ , data  $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^n$ , number of PGD attack  $n$ , PGD step-size  $\alpha$ , maximum allowable perturbation  $\varepsilon$ , scaling parameter of IGR  $\lambda$

**for**  $epoch \in \{1, 2, \dots\}$  **do**

  Compute  $\text{IG}(\mathbf{x})$

  Randomly initiate  $\tilde{\mathbf{x}} = \mathbf{x} + \mathcal{U}[-\varepsilon, \varepsilon]$

**for**  $i = 1$  **to**  $n$  **do**

$\tilde{\mathbf{x}} = \tilde{\mathbf{x}} + \alpha * \text{sign}(\nabla \mathcal{L}_{at}(\tilde{\mathbf{x}}, y))$

$\tilde{\mathbf{x}} = \text{Proj}_{\mathcal{B}_\varepsilon}(\tilde{\mathbf{x}})$

**end for**

  Compute  $\text{IG}(\tilde{\mathbf{x}})$

  Compute loss  $\mathcal{L}_{igr} = \mathcal{L}_{at}(\tilde{\mathbf{x}}, y) + \lambda(1 - \cos(\text{IG}(\mathbf{x}), \text{IG}(\tilde{\mathbf{x}})))$

  Update model parameters  $\theta$  using  $\mathcal{L}_{igr}$

**end for**

**Return**  $f$

---

Thus,  $N_{\rho(\mathbf{x}, \mathbf{x}+\boldsymbol{\eta})} \approx -N_{\rho(\mathbf{x}, \mathbf{x}-\boldsymbol{\eta})}$ . Since  $\frac{1}{n} \sum_{i=1}^n x_i^2 - (\frac{1}{n} \sum_{i=1}^n x_i)^2 \approx 0$ , the denominator of  $\rho(\mathbf{x}, \mathbf{x} + \boldsymbol{\eta})$  and  $\rho(\mathbf{x}, \mathbf{x} - \boldsymbol{\eta})$  are both small. Thus,  $\rho(\mathbf{x}, \mathbf{x} + \boldsymbol{\eta})$  and  $\rho(\mathbf{x}, \mathbf{x} - \boldsymbol{\eta})$  would be drastically different. However, since  $\|\boldsymbol{\eta}\|$  is very small, both  $\mathbf{x} + \boldsymbol{\eta}$  and  $\mathbf{x} - \boldsymbol{\eta}$  are in fact close to  $\mathbf{x}$ . Therefore, the Pearson's correlation can be unstable.

## C Additional experimental details and results

### C.1 Pseudo-code of IGR training

Algorithm 1 shows the pseudo-code for AT+IGR, where  $\tilde{\mathbf{x}}$  is generated from PGD in adversarial training. Similarly, for TRADES+IGR and MART+IGR,  $\tilde{\mathbf{x}}$  is obtained by replacing  $\mathcal{L}_{at}$  using  $\mathcal{L}_{trades}$  and  $\mathcal{L}_{mart}$ .

### C.2 Implementation details of baseline attribution robustness methods

The objective functions of the baseline attribution robustness methods are defined as follows.

#### IG-NORM [4]

$$\mathbb{E}_{\mathcal{D}} \left[ \mathcal{L}(\mathbf{x}, y; \theta) + \lambda \max_{\tilde{\mathbf{x}} \in \mathcal{B}_\varepsilon(\mathbf{x})} \|\text{IG}(\mathbf{x}, \tilde{\mathbf{x}})\|_1 \right] \quad (39)$$

#### IG-SUM-NORM [4]

$$\mathbb{E}_{\mathcal{D}} \left[ \max_{\tilde{\mathbf{x}} \in \mathcal{B}_\varepsilon(\mathbf{x})} \{ \mathcal{L}(\tilde{\mathbf{x}}, y; \theta) + \lambda \|\text{IG}(\mathbf{x}, \tilde{\mathbf{x}})\|_1 \} \right] \quad (40)$$

#### AdvAAT [13]

$$\mathbb{E}_{\mathcal{D}} \left[ \max_{\tilde{\mathbf{x}} \in \mathcal{B}_\varepsilon(\mathbf{x})} \{ \mathcal{L}(\tilde{\mathbf{x}}, y; \theta) + \lambda \text{PCL}(\text{IG}(\mathbf{x}), \text{IG}(\tilde{\mathbf{x}})) \} \right], \quad (41)$$

where  $\text{PCL}(\cdot) = 1 - [\text{PCC}(\cdot) + 1]/2$  is derived from *Pearson's Correlation Coefficient* ( $\text{PCC}(\cdot)$ ). Different from AT [21], AdvAAT adds a regularizer monitoring the attributions to the maximization problem. It generates perturbed samples that maximize both cross entropy and regularizer.

#### ART [28]

$$\mathbb{E}_{\mathcal{D}} [\mathcal{L}(\tilde{\mathbf{x}}, y; \theta) + \lambda \log(1 + \exp(-(d(g^*(\mathbf{x}), \mathbf{x}) - d(g^y(\mathbf{x}), \mathbf{x})))], \quad (42)$$

where

$$d(g^i(\mathbf{x}), \mathbf{x}) = 1 - \frac{g^i(\mathbf{x})^\top \mathbf{x}}{\|g^i(\mathbf{x})\|_2 \|\mathbf{x}\|_2}, i^* = \arg \max_{i \neq y} f(\mathbf{x})_i \quad (43)$$

and

$$\tilde{\mathbf{x}} = \arg \max_{\tilde{\mathbf{x}} \in \mathcal{B}_\varepsilon(\mathbf{x})} \log(1 + \exp(-(d(g^*(\mathbf{x}), \mathbf{x}) - d(g^y(\mathbf{x}), \mathbf{x}))) \quad (44)$$

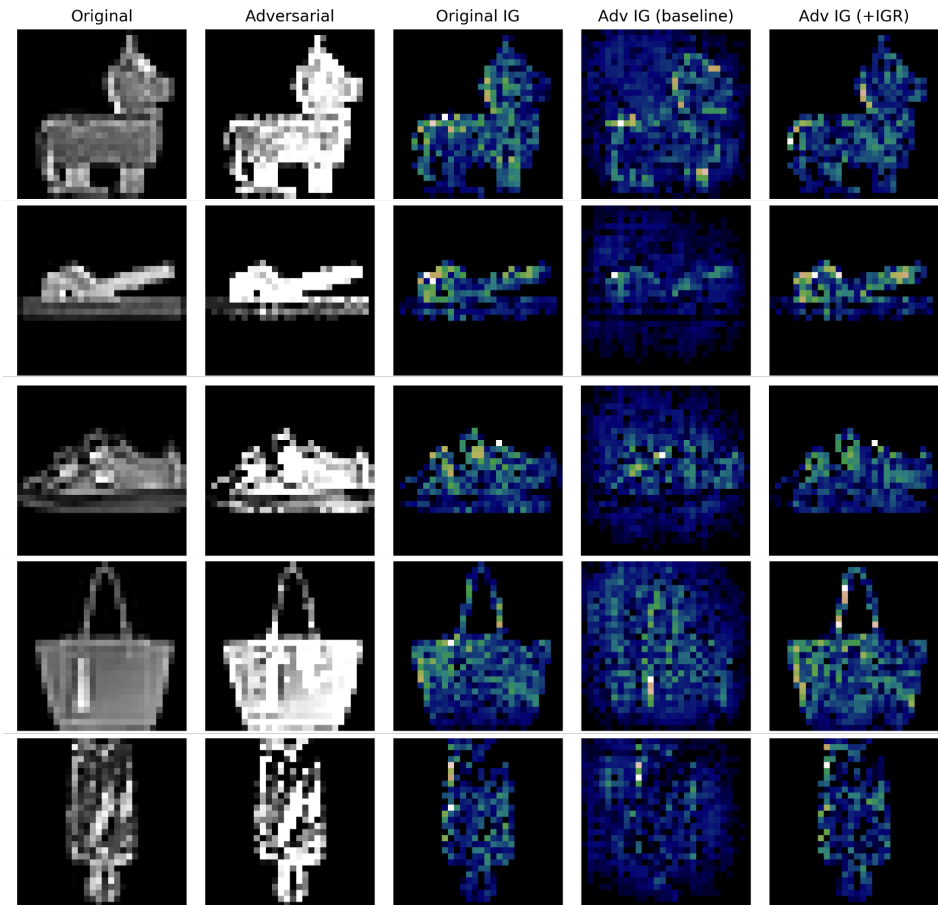


Figure 5: Additional visualization on Fashion-MNIST.

SSR [33]

$$\mathbb{E}_{\mathcal{D}} = \left[ \mathcal{L}(\mathbf{x}, y; \boldsymbol{\theta}) + \lambda s \max_i \xi_i \right]. \quad (45)$$

$\max_i \xi_i$  is the largest eigenvalue of Hessian matrix  $\tilde{\mathbf{H}}_{\mathbf{x}} = W(\text{diag}(\mathbf{p}) - \mathbf{p}^{\top} \mathbf{p})W^{\top}$ , where  $W$  is the Jacobian matrix of the logits vector and  $\mathbf{p}$  is the probits of the model.

### C.3 Additional visualization of attribution robustness

In this section, additional visualizations are provided in Fig. 5 and Fig. 6 to demonstrate that IGR improves attribution robustness. The original and adversarial images from different datasets are shown in the first two columns. The remaining three columns are IG of the original images on baseline model, IG of the adversarial images on baseline model and IG of the adversarial images on baseline+IGR model, respectively. The baseline model in the visualizations is MART.

The first two columns are the original and adversarial images from Fashion-MNIST. The third column is the IG of the original image. The last two columns are IG of adversarial examples on a baseline model and the baseline model trained with IGR. Both baseline and baseline+IGR models make the correct classifications, while only baseline+IGR protects the model interpretations.

### C.4 Visualization of Kendall’s rank correlation and Pearson’s correlation

Pearson’s correlation against Kendall’s rank correlation has been plotted in Fig. 7 under the same setting as Fig. 2a. For the same set of simulations, the corresponding Pearson’s correlations are more randomly distributed.



Figure 6: Additional visualization on CIFAR-10.

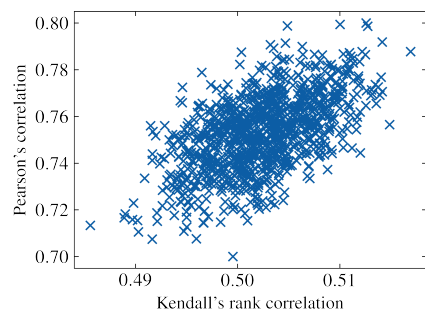


Figure 7: Visualization of Kendall's rank correlation and Pearson's correlation using simulated data.