
Supplementary Material of Learning from Crap Data via Generation

Anonymous Author(s)

Affiliation

Address

email

1 A Proofs

2 In this section, we denote $\mathbf{z} = (\mathbf{x}, y)$ and $\mathbf{z} \in \mathcal{Z}$ for convenience.

3 A.1 Proof of Theorem 1

4 **Theorem 1.** *With the optimal critical network D and the classifier C fixed, the optimization of*
 5 *generator G is equivalent to minimize $\lambda \cdot d_W(\hat{\mathbb{P}}_N, \mathbb{Q}) - D_{KL}(\mathbb{Q}||\mathbb{P}_c)$.*

6 *Proof.* Recall the objective function defined in Eq. (1),

$$\min_G \max_{D, C} U(C, G, D) = \lambda(\mathbb{E}_{\hat{\mathbb{P}}_N}[D(\mathbf{x}, y)] - \mathbb{E}_{\mathbb{Q}}[D(\mathbf{x}, y)]) - \mathbb{E}_{\mathbb{Q}}[\ell(C(\mathbf{x}), y)]. \quad (1)$$

7 Given the optimal critical network D and classifier C , the generator G is optimized by minimizing
 8 the function

$$V_{C, D}(G) = \lambda(\mathbb{E}_{\hat{\mathbb{P}}_N}[D(\mathbf{x}, y)] - \mathbb{E}_{\mathbb{Q}}[D(\mathbf{x}, y)]) - \mathbb{E}_{\mathbb{Q}}[\ell(C(\mathbf{x}), y)]. \quad (2)$$

9 As the critical network D is optimized for describe the Wasserstein, which means that

$$\mathbb{E}_{\hat{\mathbb{P}}_N}[D(\mathbf{x}, y)] - \mathbb{E}_{\mathbb{Q}}[D(\mathbf{x}, y)] = d_W(\hat{\mathbb{P}}_N, \mathbb{Q}). \quad (3)$$

10 Then we consider the last term in Eq. (2),

$$\begin{aligned} \mathbb{E}_{\mathbb{Q}}[\ell(C(\mathbf{x}), y)] &= \mathbb{E}_{\mathbb{Q}}[-\log p_c(y|\mathbf{x})] \\ &= \int p_g(\mathbf{x}, y) \log \frac{p_g(\mathbf{x}, y)}{p_c(\mathbf{x}, y)p_g(y|\mathbf{x})} d(\mathbf{x}, y) \\ &= \int p_g(\mathbf{x}, y) \log \frac{p_g(\mathbf{x}, y)}{p_c(\mathbf{x}, y)} + p_g(\mathbf{x}, y) \log \frac{1}{p_g(y|\mathbf{x})} d(\mathbf{x}, y) \\ &= D_{KL}(p_g(\mathbf{x}, y)||p_c(\mathbf{x}, y)) + H_g(y|\mathbf{x}). \end{aligned} \quad (4)$$

11 Note that the label y is provided to G during generation progress. As a result, $H_g(y|\mathbf{x})$ is irrelevant
 12 to G . By concreting Eq. (4) and Eq. (3), The proof of Theorem 1 is completed. \square

13 A.2 Proof of Theorem 2

14 **Theorem 2.** *Consider \mathbf{x} as the input samples of classifier C , and the distribution $\mathbb{Q} \in \mathbb{B}_\epsilon(\hat{\mathbb{P}}_N)$ lays*
 15 *in a Wasserstein ball centered at $\hat{\mathbb{P}}_N$ with radius ϵ . Then for any $\epsilon \geq 0$ and $\alpha \geq 1 + \beta$, we have*

$$\epsilon \|\nabla_{\mathbf{z}} \ell(\mathbf{z})\|_{\hat{\mathbb{P}}_N}^{\alpha_*} - \epsilon^{\beta+1} \|h(\mathbf{z})\|_{\hat{\mathbb{P}}_N}^{\frac{\alpha}{\alpha-\beta-1}} \leq \mathbb{E}_{\mathbb{Q}}(\ell(\mathbf{z})) - \mathbb{E}_{\hat{\mathbb{P}}_N}(\ell(\mathbf{z})) \leq \epsilon \|\nabla_{\mathbf{z}} \ell(\mathbf{z})\|_{\hat{\mathbb{P}}_N}^{\alpha_*} + \epsilon^{\beta+1} \|h(\mathbf{z})\|_{\hat{\mathbb{P}}_N}^{\alpha_*}, \quad (5)$$

16 where $\|f(\mathbf{z})\|_{\hat{\mathbb{P}}_N}^\alpha \triangleq (\frac{1}{N} \sum_{i=1, \mathbf{z} \sim \hat{\mathbb{P}}_N}^N (\|f(\mathbf{z}_i)\|^\alpha))^{1/\alpha}$, $\alpha_* = \frac{\alpha}{\alpha-1}$, and $h(\mathbf{z})$ is a function and $\beta \in$
 17 $(0, 1]$ a constant which satisfy $\|\nabla_{\mathbf{z}} \ell(\mathbf{z}_1) - \nabla_{\mathbf{z}} \ell(\mathbf{z}_2)\| \leq h(\mathbf{z}_2) \cdot \|\mathbf{z}_1 - \mathbf{z}_2\|^\beta$ for any $\mathbf{z} = (\mathbf{x}, y) \in \mathcal{Z}$.

18 In this part, we firstly proof the right part of this theorem which is an upper bound of $\mathbb{E}_{\mathbb{Q}}(\ell(\mathbf{z})) -$
 19 $\mathbb{E}_{\hat{\mathbb{P}}_N}(\ell(\mathbf{z}))$. Then we provide the proof of a lower bound of it. By combining them, the proof of
 20 theorem 2 is completed.

21 *Proof.* Considering the inner part of the proposed objective function defined as follows,

$$\sup_{\mathbb{Q} \in \mathbb{B}_\epsilon(\hat{\mathbb{P}}_N)} \mathbb{E}_{\mathbb{Q}}[\ell_\theta(\mathbf{z})] = \inf_{\lambda \geq 0} \lambda \epsilon + \sup_{\mathbb{Q} \in \mathcal{M}(\mathcal{Z})} \int_{\mathcal{Z}} \ell_\theta(\mathbf{z}) \mathbb{Q}(\mathrm{d}(\mathbf{z})) - \lambda \cdot d(\mathbb{Q}, \hat{\mathbb{P}}_N). \quad (6)$$

22 Under the situation that the network D is optimized for calculating the Wasserstein distance, we
 23 consider the network D is sufficient to describe the Wasserstein distance. Recall that we define the
 24 Wasserstein distance as

$$d_W(\mathbb{Q}_1, \mathbb{Q}_2) \triangleq \min_{\Pi \in \mathcal{M}(\Pi)} \left\{ \int_{\mathcal{Z} \times \mathcal{Z}} s(\mathbf{z}_1, \mathbf{z}_2) \Pi(\mathrm{d}(\mathbf{z}_1, \mathbf{z}_2)) \right\}. \quad (7)$$

25 Assuming the metric $s(\cdot, \cdot)$ is induced by some norm $\|\cdot\|^\alpha$, it is easy to be reformulated as follows:

$$d_W(\mathbb{Q}_1, \mathbb{Q}_2) = \min_{\Pi \in \mathcal{M}(\Pi)} \left\{ \int_{\mathcal{Z} \times \mathcal{Z}} \|\mathbf{z}_1 - \mathbf{z}_2\|^\alpha \Pi(\mathrm{d}(\mathbf{z}_1, \mathbf{z}_2)) \right\}. \quad (8)$$

26 Plugging Eq. (8) into Eq. (6) gives us

$$\mathbb{E}_{\mathbb{Q}}(\ell(\mathbf{z})) = \inf_{\lambda \geq 0} \left\{ \lambda \epsilon + \frac{1}{n} \sum_{i=1}^n \sup_{\mathbf{z} \in \mathcal{Z}} (\ell(\mathbf{z}) - \lambda \|\mathbf{z} - \mathbf{z}'_i\|^\alpha) \right\} \quad (9)$$

27 where $\mathbf{z} \sim \mathbb{Q}$ and $\mathbf{z}' \sim \hat{\mathbb{P}}_N$, then we consider a upper bound that

$$\begin{aligned} & \sup_{\mathbf{z} \in \mathcal{Z}} \{\ell(\mathbf{z}) - \ell(\mathbf{z}'_i) - \lambda \cdot \|\mathbf{z} - \mathbf{z}'_i\|^\alpha\} \\ & \leq \sup_{\mathbf{z} \in \mathcal{Z}} \{\|\nabla_{\mathbf{z}} \ell(\mathbf{z}'_i)\|_* \cdot \|\mathbf{z} - \mathbf{z}'_i\| + h(\mathbf{z}'_i) \cdot \|\mathbf{z} - \mathbf{z}'_i\|^{\beta+1} - \lambda \cdot \|\mathbf{z} - \mathbf{z}'_i\|^\alpha\} \\ & \leq \sup_{\mathbf{z} \in \mathcal{Z}} \{\|\nabla_{\mathbf{z}} \ell(\mathbf{z}'_i)\|_* \cdot \|\mathbf{z} - \mathbf{z}'_i\| + h(\mathbf{z}'_i) \cdot \|\mathbf{z} - \mathbf{z}'_i\|^{\beta+1} - \lambda \cdot \|\mathbf{z} - \mathbf{z}'_i\|^\alpha + C \cdot \|\mathbf{z} - \mathbf{z}'_i\|^{\gamma+1}\} \\ & \leq \sup_{\xi \geq 0} \{\|\nabla_{\mathbf{z}} \ell(\mathbf{z}'_i)\|_* \cdot \xi + h(\mathbf{z}'_i) \cdot \xi^{\beta+1} + C \cdot \xi^{\gamma+1} - \lambda \cdot \xi^\alpha\}, \end{aligned} \quad (10)$$

28 where $0 \leq C$, $1 < \gamma < \beta$ and $\xi := \|\mathbf{z} - \mathbf{z}'_i\|$. Following Young's inequality for products that
 29 $ab \leq \frac{a^p}{p} + \frac{b^q}{q}$, we set $p = \frac{\alpha-1}{\alpha-1-\beta}$, $q = \frac{\alpha-1}{\beta}$ satisfying $\frac{1}{p} + \frac{1}{q} = 1$ and $a = \varphi^{1/p} \xi^{1/p}$, $b = \varphi^{-1/q} \xi^{\alpha/q}$.
 30 Then for any $t > 0$ and $\varphi > 0$, it holds that

$$\xi^{\gamma+1} \leq \xi^{\beta+1} \leq \xi^{\alpha+1} \leq \frac{\alpha-1-\beta}{\alpha-1} \varphi \xi + \frac{\beta}{\alpha-1} \varphi^{-\frac{\alpha-1-\beta}{\beta}} \xi^\alpha. \quad (11)$$

31 Replacing $\xi^{\gamma+1}$ and $\xi^{\beta+1}$ with the last term of Eq. (11), it gives us

$$\begin{aligned} & \sup_{\xi \geq 0} \{\|\nabla_{\mathbf{z}} \ell(\mathbf{z}'_i)\|_* \cdot \xi + h(\mathbf{z}'_i) \cdot \xi^{\beta+1} + C \cdot \xi^{\gamma+1} - \lambda \cdot \xi^\alpha\} \\ & \leq \sup_{\xi \geq 0} \left\{ \|\nabla_{\mathbf{z}} \ell(\mathbf{z}'_i)\|_* + \frac{\alpha-\beta-1}{\alpha-1} \cdot h(\mathbf{z}'_i) \cdot \varphi_1 + \frac{\alpha-\gamma-1}{\alpha-1} \cdot C \cdot \varphi_2 \right\} \cdot \xi \\ & \quad - \left(\lambda - \frac{\beta}{\alpha-1} \cdot h(\mathbf{z}'_i) \cdot \varphi_1^{-\frac{\alpha-\beta-1}{\beta}} - \frac{\gamma}{\alpha-1} \cdot C \cdot \varphi_2^{\frac{\alpha-\gamma-1}{\gamma}} \right) \cdot \xi^\alpha \} \\ & \leq \sup_{\xi \geq 0} \{ \mathcal{G}_\varphi(\mathbf{z}'_i) \cdot \xi - (\lambda - \mathcal{N}_\varphi) \cdot \xi^\alpha \}, \end{aligned} \quad (12)$$

32 where $\mathcal{G}_\varphi(\mathbf{z}'_i) = \|\nabla_{\mathbf{z}} \ell(\mathbf{z}'_i)\|_* + \frac{\alpha-\beta-1}{\alpha-1} \cdot h(\mathbf{z}'_i) \cdot \varphi_1 + \frac{\alpha-\gamma-1}{\alpha-1} \cdot C \cdot \varphi_2$ and $\mathcal{N}_\varphi = \lambda - \frac{\beta}{\alpha-1} \cdot h(\mathbf{z}'_i) \cdot$
 33 $\varphi_1^{-\frac{\alpha-\beta-1}{\beta}} - \frac{\gamma}{\alpha-1} \cdot C \cdot \varphi_2^{\frac{\alpha-\gamma-1}{\gamma}}$. Considering the value of Eq. (12) is $+\infty$ when $\lambda \leq \mathcal{N}_\varphi$, we solve
 34 Eq. (12) over ξ and conclude that

$$\begin{aligned} & \mathbb{E}_{\mathbb{Q}}(\ell(\mathbf{z})) - \mathbb{E}_{\hat{\mathbb{P}}_N}(\ell(\mathbf{z})) \\ & \leq \inf_{\lambda \geq \mathcal{N}_\varphi} \left\{ \lambda \epsilon^\alpha + \alpha^{-1-\frac{\alpha}{\alpha-1}} (\alpha-1) (\lambda - \mathcal{N}_\varphi)^{-\frac{1}{\alpha-1}} (\|\mathcal{G}_\varphi\|_{\hat{\mathbb{P}}_N}^{\frac{\alpha}{\alpha-1}})^{\frac{\alpha}{\alpha-1}} \right\} \\ & \leq \epsilon \|\mathcal{G}_\varphi\|_{\hat{\mathbb{P}}_N}^{\frac{\alpha}{\alpha-1}} + \mathcal{N}_\varphi \epsilon^\alpha. \end{aligned} \quad (13)$$

35 Plugging \mathcal{G}_φ and \mathcal{N}_φ into Eq. (13) and solving the minimization problem on φ , we obtain the right
 36 part of Theorem 2. Next we step to the left part of Theorem 2. We firstly began with a lower bound
 37 of $\sup_{\mathbf{Q} \in \mathbb{B}_\epsilon(\hat{\mathbb{P}}_N)} \mathbb{E}_{\mathbf{Q}}(\ell(\mathbf{z})) - \mathbb{E}_{\hat{\mathbb{P}}_N}(\ell(\mathbf{z}))$ as follows:

$$\begin{aligned}
 & \sup_{\mathbf{z}_i \in \mathcal{Z}} \left\{ \frac{1}{N} \sum_{i=1}^N [\ell(\mathbf{z}_i) - \ell(\mathbf{z}'_i)] : \left(\frac{1}{N} \sum_{i=1}^N \|\mathbf{z}_i - \mathbf{z}'_i\|^\alpha \right)^{\frac{1}{\alpha}} \leq \epsilon \right\} \\
 & \geq \sup_{\mathbf{z}_i \in \mathcal{Z}} \left\{ \frac{1}{N} \sum_{i=1}^N [\nabla_{\mathbf{z}} \ell(\mathbf{z}'_i) \|\mathbf{z}_i - \mathbf{z}'_i\| - h(\mathbf{z}'_i) \|\mathbf{z}_i - \mathbf{z}'_i\|^{\beta+1}] : \left(\frac{1}{N} \sum_{i=1}^N \|\mathbf{z}_i - \mathbf{z}'_i\|^\alpha \right)^{\frac{1}{\alpha}} \leq \epsilon \right\} \\
 & \geq \sup_{\mathbf{z}_i \in \mathcal{Z}} \left\{ \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{z}} \ell(\mathbf{z}'_i) \|\mathbf{z}_i - \mathbf{z}'_i\| : \left(\frac{1}{N} \sum_{i=1}^N \|\mathbf{z}_i - \mathbf{z}'_i\|^\alpha \right)^{\frac{1}{\alpha}} \leq \epsilon \right\} \\
 & \quad - \sup_{\mathbf{z}_i \in \mathcal{Z}} \left\{ \frac{1}{N} \sum_{i=1}^N h(\mathbf{z}'_i) \|\mathbf{z}_i - \mathbf{z}'_i\|^{\beta+1} : \left(\frac{1}{N} \sum_{i=1}^N \|\mathbf{z}_i - \mathbf{z}'_i\|^\alpha \right)^{\frac{1}{\alpha}} \leq \epsilon \right\}
 \end{aligned} \tag{14}$$

38 Further we conclude that with the help of Holder's inequality,

$$\begin{aligned}
 & \sup_{\mathbf{z}_i \in \mathcal{Z}} \left\{ \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{z}} \ell(\mathbf{z}'_i) \|\mathbf{z}_i - \mathbf{z}'_i\| : \left(\frac{1}{N} \sum_{i=1}^N \|\mathbf{z}_i - \mathbf{z}'_i\|^\alpha \right)^{\frac{1}{\alpha}} \leq \epsilon \right\} \\
 & = \sup_{\xi \in \mathbb{R}} \left\{ \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{z}} \ell(\mathbf{z}'_i) \xi_i : \left(\frac{1}{N} \sum_{i=1}^N \xi_i^\alpha \right)^{\frac{1}{\alpha}} \leq \epsilon \right\} \\
 & = \epsilon \|\nabla_{\mathbf{z}} \ell(\mathbf{z})\|_{\hat{\mathbb{P}}_N}^{\alpha*}.
 \end{aligned} \tag{15}$$

39 Wee also have that

$$\begin{aligned}
 & \sup_{\mathbf{z}_i \in \mathcal{Z}} \left\{ \frac{1}{N} \sum_{i=1}^N h(\mathbf{z}'_i) \|\mathbf{z}_i - \mathbf{z}'_i\|^{\beta+1} : \left(\frac{1}{N} \sum_{i=1}^N \|\mathbf{z}_i - \mathbf{z}'_i\|^\alpha \right)^{\frac{1}{\alpha}} \leq \epsilon \right\} \\
 & = \sup_{\xi \in \mathbb{R}} \left\{ \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{z}} h(\mathbf{z}'_i) \xi_i^{\beta+1} : \left(\frac{1}{N} \sum_{i=1}^N \xi_i^\alpha \right)^{\frac{1}{\alpha}} \leq \epsilon \right\} \\
 & = \epsilon^{\beta+1} \|h(\mathbf{z})\|_{\hat{\mathbb{P}}_N}^{\frac{\alpha}{\alpha-\beta-1}}.
 \end{aligned} \tag{16}$$

40 The proof is completed. \square

41 A.3 Proof of Theorem 3

42 **Theorem 3.** For any $0 < \delta < 1$, with probability at least $1 - \delta$ with respect to the sampling,

$$\mathbb{E}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) \leq \mathbb{E}_{\mathbf{Q}}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) + \frac{12\sqrt{R}}{n} (\log \frac{n}{3\sqrt{R}} + 1) + \sqrt{\frac{8 \log(2/\delta)}{N}}, \tag{17}$$

43 and for any $\zeta > \frac{12\sqrt{R}}{n} (\log \frac{n}{3\sqrt{R}} + 1) + \sqrt{\frac{8 \log(2/\delta)}{N}}$, we have

$$P(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y) \geq \mathbb{E}_{\mathbf{Q}}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) + \zeta) \leq \frac{\mathbb{E}_{\mathbf{Q}}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) + \frac{12\sqrt{R}}{n} (\log \frac{n}{3\sqrt{R}} + 1) + \sqrt{\frac{8 \log(2/\delta)}{N}}}{\mathbb{E}_{\mathbf{Q}}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) + \zeta}. \tag{18}$$

44 where R is only related to the architecture of the neural network.

45 *Proof.* As description in the Theorem 8 in [2], for any integer N and $\delta \in (0, 1)$, the risk bounds can
 46 be written as that

$$\mathbf{EL}(Y, f(X)) \leq \hat{\mathbf{E}}_N \phi(Y, f(X)) + \mathfrak{R}_N(\tilde{\phi} \circ F) + \sqrt{\frac{8 \ln(2/\delta)}{N}}. \tag{19}$$

47 Setting the loss function \mathcal{L} and Y as $\mathcal{L}(\mathbf{x}, y) = \phi(\mathbf{x}, y) = \ell(C(\mathbf{x}, \boldsymbol{\theta}), y)$, it yields that

$$\mathbb{E}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) \leq \mathbb{E}_{\mathbb{Q}}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) + 2\Re_N(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) + \sqrt{\frac{8 \ln(2/\delta)}{N}}, \quad (20)$$

48 We use Lemma A.8 in [1] and get the Rademacher complexity estimate of neural networks C with
49 loss function ℓ as follows,

$$\Re_N(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) = \frac{12\sqrt{R}}{n}(\log \frac{n}{3\sqrt{R}} + 1), \quad (21)$$

50 where R is only related with the architecture of neural network and defined as follows and the detailed
51 notation can be find in [1],

$$R = \frac{4B^2 \ln(2W^2)}{\gamma^2 \epsilon^2} \left(\prod_{j=1}^L s_j^2 \rho_j^2 \right) \left(\sum_{i=1}^L \left(\frac{b_i}{s_i} \right)^{2/3} \right)^3. \quad (22)$$

52 Plugging result in Eq. (21) into Eq. (20) drives the first part of Theorem 3. Further, we can easily
53 obtain the second part of Theorem 3 by applying Markov's inequality here to obtain,

$$\begin{aligned} P(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y) \geq \mathbb{E}_{\mathbb{Q}}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) + \zeta) &\leq \frac{\mathbb{E}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y))}{\mathbb{E}_{\mathbb{Q}}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) + \zeta} \\ &\leq \frac{\mathbb{E}_{\mathbb{Q}}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) + \frac{12\sqrt{R}}{n}(\log \frac{n}{3\sqrt{R}} + 1) + \sqrt{\frac{8 \log(2/\delta)}{N}}}{\mathbb{E}_{\mathbb{Q}}(\ell(C(\mathbf{x}, \boldsymbol{\theta}), y)) + \zeta} \end{aligned} \quad (23)$$

54 which completes the proof. \square

55 A.4 Proof of Theorem 4

56 **Theorem 4.** For the generator G and classifier C fixed, the optimal discriminator D is

$$D_{G,C}^*(\mathbf{x}, y) = \frac{p_{data}(\mathbf{x}, y)}{p_{data}(\mathbf{x}, y) + p_g(\mathbf{x}, y)}, \quad (24)$$

57 where $p_g(\mathbf{x})$ is the distribution generated by G .

58 *Proof.* Given the generator and classifier, the loss function can be written as

$$\begin{aligned} V(D) &= \int p_{data}(\mathbf{z}) \log D(\mathbf{z}) d\mathbf{z} + \int p_g(\mathbf{z}) \log(1 - D(\mathbf{z})) d\mathbf{z} \\ &= \int p_d(\mathbf{z}) \log D(\mathbf{z}) + p_g(\mathbf{z}) \log(1 - D(\mathbf{z})) d\mathbf{z}. \end{aligned} \quad (25)$$

59 Following the proof in GAN [3], the function $V(D)$ achieves its maximum at $\frac{p_{data}(\mathbf{z})}{p_{data}(\mathbf{z}) + p_g(\mathbf{z})}$. \square

60 A.5 Proof of Theorem 5

61 **Theorem 5.** With the optimal discriminator D and the classifier C fixed, the optimization of generator
62 G is equivalent to $-\log 4 + 2JSD(\hat{\mathbb{P}}_N || \mathbb{Q}) - 1/\lambda \cdot D_{KL}(\mathbb{Q} || \mathbb{P}_c)$.

63 *Proof.* Following the conclusion obtained in Theorem 4, with the optimal $D_{G,C}^*(\mathbf{x}, y) =$
64 $\frac{p_{data}(\mathbf{x}, y)}{p_{data}(\mathbf{x}, y) + p_g(\mathbf{x}, y)}$, the minimax game for G can be reformulated as:

$$\begin{aligned} V(G, C) &= \int p_{data}(\mathbf{z}) \log \frac{p_{data}(\mathbf{z})}{p_{data}(\mathbf{z}) + p_g(\mathbf{z})} d\mathbf{z} + \int p_g(\mathbf{z}) \log \frac{p_g(\mathbf{z})}{p_{data}(\mathbf{z}) + p_g(\mathbf{z})} d\mathbf{z} + \lambda \int p_g(\mathbf{z}) \ell(\mathbf{z}) d\mathbf{z} \\ &= -\log 4 + 2JSD(p_{data}(\mathbf{z}) || p_g(\mathbf{z})) d\mathbf{z} + \lambda \int p_g(\mathbf{z}) [-\log p_c(\mathbf{z})] d\mathbf{z} \\ &= -\log 4 + 2JSD(p_{data}(\mathbf{z}) || p_g(\mathbf{z})) d\mathbf{z} + \lambda(D_{KL}((p_g(\mathbf{z}) || p_c(\mathbf{z})) + H_g(y|\mathbf{x})). \end{aligned} \quad (26)$$

65 Noting that the label y in $p_g(\mathbf{x}, y)$ is assigned during the generation, $H_g(y|\mathbf{x})$ is a constant
66 which is irrelevant with G . As a result, the generator G will be optimized by $-\log 4 +$
67 $2JSD(p_{data}(\mathbf{z}) || p_g(\mathbf{z})) d\mathbf{z} + \lambda D_{KL}((p_g(\mathbf{z}) || p_c(\mathbf{z})))$, which completes the proof. \square

68 B Algorithm

Algorithm 1 Proposed Method

Input: The batch size m , the loss balanced coefficient λ .

Initialize generator parameters θ_g for G , θ_d for D , and θ_c for C with the training set $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$.

Sample a batch of pairs $(\mathbf{x}_g, y_g) \sim p_g(\mathbf{x}, y)$ and a batch of pairs $(\mathbf{x}_d, y_d) \sim p_d(\mathbf{x}, y)$.

Update D by ascending along its gradients $\nabla_{\theta_d} \left[\frac{1}{m} \left(\sum_{(\mathbf{x}_d, y_d)} D(\mathbf{x}_d, y_d) - \sum_{(\mathbf{x}_g, y_g)} D(\mathbf{x}_g, y_g) \right) \right]$

Update G by ascending along its gradients $\nabla_{\theta_g} \left[\frac{1}{m} \left(\lambda \cdot \sum_{(\mathbf{x}_g, y_g)} D(\mathbf{x}_g, y_g) - \sum_{(\mathbf{x}_d, y_d)} C(\mathbf{x}_g, y_g) \right) \right]$

Update C by ascending along its gradients $\nabla_{\theta_c} \left[\frac{1}{m} \left(\sum_{(\mathbf{x}_d, y_d)} \ell(C(\mathbf{x}_g, \theta), y_g) \right) \right]$

Output: A distribution optimized classifier C and a worst-case distribution generator G .

69 References

- 70 [1] P. L. Bartlett, D. J. Foster, and M. J. Telgarsky. Spectrally-normalized margin bounds for neural networks.
71 In *Advances in Neural Information Processing Systems*, pages 6240–6249, 2017.
- 72 [2] P. L. Bartlett and S. Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results.
73 *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- 74 [3] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio.
75 Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.