

Appendix:

Balanced Self-Paced Learning for Generative Adversarial Clustering Network

Kamran Ghasedi Dizaji¹, Xiaoqian Wang¹, Cheng Deng², and Heng Huang^{1,3}

¹Electrical and Computer Engineering Department, University of Pittsburgh, PA, USA

²School of Electronic Engineering, Xidian University, Xi'an, Shaanxi, China

³JD Digits

kamran.ghasedi@gmail.com, xiaoqian.wang@pitt.edu, chdeng.xd@gmail.com, heng.huang@pitt.edu

The objective function of the adversarial game for *ClusterGAN* is:

$$\min_{\mathcal{G}, \mathcal{C}} \max_{\mathcal{D}} \mathbf{U}(\mathcal{D}, \mathcal{G}, \mathcal{C}) = \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} [\log \mathcal{D}(\mathcal{C}(\mathbf{x}), \mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim P(\mathbf{z})} [\log (1 - \mathcal{D}(\mathbf{z}, \mathcal{G}(\mathbf{z})))] . \quad (1)$$

Lemma 1. For any fixed \mathcal{G} and \mathcal{C} , the optimal \mathcal{D} defined by the utility function $\mathbf{U}(\mathcal{D}, \mathcal{G}, \mathcal{C})$ is:

$$\begin{aligned} \mathcal{D}^*(\mathbf{z}, \mathbf{x}) &= \frac{P(\mathbf{x})P_{\mathcal{C}}(\mathbf{z}|\mathbf{x})}{P(\mathbf{x})P_{\mathcal{C}}(\mathbf{z}|\mathbf{x}) + P(\mathbf{z})P_{\mathcal{G}}(\mathbf{x}|\mathbf{z})} \\ &= \frac{P_{\mathcal{C}}(\mathbf{z}, \mathbf{x})}{P_{\mathcal{C}}(\mathbf{z}, \mathbf{x}) + P_{\mathcal{G}}(\mathbf{z}, \mathbf{x})} \end{aligned}$$

Proof. Given the clusterer and generator, the utility function $\mathbf{U}(\mathcal{D}, \mathcal{G}, \mathcal{C})$ can be rewritten as

$$\begin{aligned} \mathbf{U}(\mathcal{D}, \mathcal{G}, \mathcal{C}) &= \iint P(\mathbf{x})P_{\mathcal{C}}(\mathbf{z}|\mathbf{x}) \log(\mathcal{D}(\mathbf{z}, \mathbf{x})) d\mathbf{x}d\mathbf{z} \quad (2) \\ &+ \iint P(\mathbf{z})P_{\mathcal{G}}(\mathbf{x}|\mathbf{z}) \log(1 - \mathcal{D}(\mathbf{z}, \mathbf{x})) d\mathbf{x}d\mathbf{z} \\ &= \iint P_{\mathcal{C}}(\mathbf{z}, \mathbf{x}) \log(\mathcal{D}(\mathbf{z}, \mathbf{x})) d\mathbf{x}d\mathbf{z} \\ &+ \iint P_{\mathcal{G}}(\mathbf{z}, \mathbf{x}) \log(1 - \mathcal{D}(\mathbf{z}, \mathbf{x})) d\mathbf{x}d\mathbf{z} \\ &= f(\mathcal{D}(\mathbf{z}, \mathbf{x})) \end{aligned}$$

For any $(P_{\mathcal{C}}(\mathbf{z}, \mathbf{x}), P_{\mathcal{G}}(\mathbf{z}, \mathbf{x})) \in \mathbb{R}^2 \setminus \{0, 0\}$, the function $f(\mathcal{D}(\mathbf{z}, \mathbf{x}))$ achieves its maximum at $\frac{P_{\mathcal{C}}(\mathbf{z}, \mathbf{x})}{P_{\mathcal{C}}(\mathbf{z}, \mathbf{x}) + P_{\mathcal{G}}(\mathbf{z}, \mathbf{x})}$. \square

Given $\mathcal{D}^*(\mathbf{x}, \mathbf{z})$, we can further replace \mathcal{D} in the utility function $\mathbf{U}(\mathcal{D}, \mathcal{G}, \mathcal{C})$ and reformulate the objective as $\mathbf{V}(\mathcal{G}, \mathcal{C}) = \max_{\mathcal{D}} \mathbf{U}(\mathcal{D}, \mathcal{G}, \mathcal{C})$.

Lemma 2. The global optimum point of $\mathbf{V}(\mathcal{G}, \mathcal{C})$ is achieved if and only if $P(\mathbf{z}, \hat{\mathbf{x}}) = P(\hat{\mathbf{z}}, \mathbf{x})$.

Proof. Given $\mathcal{D}^*(\mathbf{x}, \mathbf{z})$, the utility function $\mathbf{V}(\mathcal{G}, \mathcal{C})$ can be reformulated as:

$$\begin{aligned} \mathbf{V}(\mathcal{G}, \mathcal{C}) &= \iint P_{\mathcal{C}}(\mathbf{z}, \mathbf{x}) \log \left(\frac{P_{\mathcal{C}}(\mathbf{z}, \mathbf{x})}{P_{\mathcal{C}}(\mathbf{z}, \mathbf{x}) + P_{\mathcal{G}}(\mathbf{z}, \mathbf{x})} \right) d\mathbf{x}d\mathbf{z} \\ &+ \iint P_{\mathcal{G}}(\mathbf{z}, \mathbf{x}) \log \left(\frac{P_{\mathcal{G}}(\mathbf{z}, \mathbf{x})}{P_{\mathcal{C}}(\mathbf{z}, \mathbf{x}) + P_{\mathcal{G}}(\mathbf{z}, \mathbf{x})} \right) d\mathbf{x}d\mathbf{z} \quad (3) \end{aligned}$$

Sketching the proof in original GAN paper [1], $\mathbf{V}(\mathcal{G}, \mathcal{C})$ can be rewritten as:

$$\mathbf{V}(\mathcal{G}, \mathcal{C}) = -\log 4 + 2JSD(P_{\mathcal{C}}(\mathbf{z}, \mathbf{x}) \| P_{\mathcal{G}}(\mathbf{z}, \mathbf{x})), \quad (4)$$

where *JSD* represents the Jensen-Shannon divergence, which is always non-negative. Therefore, the unique optimum of $\mathbf{V}(\mathcal{G}, \mathcal{C})$ is achieved if and only if $P_{\mathcal{C}}(\mathbf{z}, \mathbf{x}) = P_{\mathcal{G}}(\mathbf{z}, \mathbf{x})$, or in other words

$$P(\mathbf{z}, \hat{\mathbf{x}}) = P(\hat{\mathbf{z}}, \mathbf{x})$$

\square

The optimization problem for estimating our balanced self-paced learning algorithm is:

$$\min_{\boldsymbol{\nu}} \mathbf{L}(\boldsymbol{\nu}) = \sum_{i=1}^n \nu_i l_i - \lambda_{\nu} \|\boldsymbol{\nu}\|_1 + \gamma \|\boldsymbol{\nu}\|_e \quad s.t. \quad \boldsymbol{\nu} \in [0, 1]^n. \quad (5)$$

Theorem 1. For any fixed \mathcal{C} , the optimal $\boldsymbol{\nu}$ defined by the objective function $\mathbf{L}(\boldsymbol{\nu})$ is:

$$\begin{cases} \nu_{kq}^* = 1, & \text{if } l_{kq} < \lambda_{\nu} - 2\gamma q \\ \nu_{kq}^* = \frac{\lambda_{\nu} - l_{kq}}{2\gamma} - q, & \text{if } \lambda_{\nu} - 2\gamma q \leq l_{kq} < \lambda_{\nu} - 2\gamma(q-1) \\ \nu_{kq}^* = 0, & \text{if } l_{kq} \geq \lambda_{\nu} - 2\gamma(q-1) \end{cases}$$

where $q \in \{1, \dots, n_k\}$ is the sorted index of loss values $\{l_{k1}, \dots, l_{kn_k}\}$ in the k -th group.

Proof.

$$\begin{aligned} \min_{\boldsymbol{\nu}} \mathbf{L}(\boldsymbol{\nu}) &= \sum_{k=1}^c \mathbf{L}(\boldsymbol{\nu}_k) \quad (6) \\ &= \sum_{k=1}^c \left[\sum_{i=1}^{n_k} \nu_{ki} (l_{ki} - \lambda_{\nu}) + \gamma \left(\sum_{i=1}^{n_k} |\nu_{ki}| \right)^2 \right], \quad s.t. \boldsymbol{\nu} \in [0, 1]^n, \end{aligned}$$

We can handle the c groups in Problem (6) separately. Given k , define $\mathbf{b} = \left[\frac{(l_{k1} - \lambda_{\nu})}{\gamma}, \frac{(l_{k2} - \lambda_{\nu})}{\gamma}, \dots, \frac{(l_{kn_k} - \lambda_{\nu})}{\gamma} \right]$, the optimization problem *w.r.t.* the k -th group can be formulated as follows:

$$\min_{\mathbf{u}} \mathbf{b}^T \mathbf{u} + \mathbf{u}^T \mathbf{1} \mathbf{1}^T \mathbf{u}, \quad s.t. \mathbf{0} \leq \mathbf{u} \leq \mathbf{1}, \quad (7)$$

where $\mathbf{u} = [u_{k1}, u_{k2}, \dots, u_{kn_k}]$. The Lagrangian function of Problem (7) is

$$\min_{\mathbf{u}} \mathbf{b}^T \mathbf{u} + \mathbf{u}^T \mathbf{1} \mathbf{1}^T \mathbf{u} - \eta^T \mathbf{u} - \lambda^T (\mathbf{1} - \mathbf{u}). \quad (8)$$

where $\eta \geq \mathbf{0}$ and $\lambda \geq \mathbf{0}$ are Lagrangian multipliers. Take derivate of Problem (8) *w.r.t.* \mathbf{u} and set it to zero, we get

$$\eta + \lambda - \mathbf{b} = 2m\mathbf{1}. \quad (9)$$

where $m = \mathbf{1}^T \mathbf{u}$. From the KKT condition we can derive $\eta^T \mathbf{u} = 0$ and $\lambda^T (\mathbf{1} - \mathbf{u}) = 0$. Consequently, we can derive

$$\begin{cases} u_q = 0 & \implies \eta_q > 0, \lambda_q = 0 & \implies \frac{b_q}{2} + m > 0, \\ 0 < u_q < 1 & \implies \eta_q = 0, \lambda_q = 0 & \implies \frac{b_q}{2} + m = 0, \\ u_q = 1 & \implies \eta_q = 0, \lambda_q > 0 & \implies \frac{b_q}{2} + m < 0, \end{cases} \quad (10)$$

where $q \in \{1, \dots, n_k\}$. Without loss of generality, suppose \mathbf{b} is a sorted vector such that $b_1 < b_2 < \dots < b_{n_k}$, then according to Eq. (10) we have $1 \geq u_1 \geq u_2 \geq \dots \geq u_{n_k} \geq 0$, from which we can derive

$$\begin{cases} u_q = 0 & \implies u_r = 0, \forall r \geq q & \implies m \leq q - 1, \\ 0 < u_q < 1 & \implies u_r = 0, \forall r > q, \text{ and } u_r = 1, \forall r < q & \implies q - 1 < m < q, \\ u_q = 1 & \implies u_r = 1, \forall r \leq q & \implies m \geq q. \end{cases} \quad (11)$$

Combining Eq. (10) and Eq. (11) we can derive the solution to Problem (7) as follows:

$$\begin{cases} -\frac{b_q}{2} \leq q - 1 & \implies u_q = 0, \\ q - 1 < -\frac{b_q}{2} < q & \implies u_q = -\frac{b_q}{2} - q + 1, \\ -\frac{b_q}{2} \geq q & \implies u_q = 1, \end{cases}$$

which can be rewritten based on $\boldsymbol{\nu}$ as:

$$\begin{cases} \nu_{kq}^* = 1, & \text{if } l_{kq} < \lambda_{\nu} - 2\gamma q \\ \nu_{kq}^* = \frac{\lambda_{\nu} - l_{kq}}{2\gamma} - q, & \text{if } \lambda_{\nu} - 2\gamma q \leq l_{kq} < \lambda_{\nu} - 2\gamma(q - 1) \\ \nu_{kq}^* = 0, & \text{if } l_{kq} \geq \lambda_{\nu} - 2\gamma(q - 1) \end{cases}$$

□

References

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems (NIPS)*, pages 2672–2680, 2014. 1