

# Mitigating Information Leakage in Image Representations: A Maximum Entropy Approach (supplementary Material)

Proteek Chandan Roy and Vishnu Naresh Boddeti  
Department of Computer Science and Engineering  
Michigan State University, East Lansing MI 48824  
{royprote, vishnu}@msu.edu

In this supplementary material we include proof of Theorem 1 in Section 1, Corollary 1.1 in Section 2 and finally provide the numerical values of the trade-off fronts in the CIFAR-10 and CIFAR-100 experiment in Section 3.

## 1. Proof of Theorem 1

**Theorem 1.** Given a fixed encoder  $E$ , the optimal discriminator is  $q_D(s|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) = p(s|E(\mathbf{x}; \boldsymbol{\theta}_E))$  and the optimal predictor is  $q_T(t|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_T^*) = p(t|E(\mathbf{x}; \boldsymbol{\theta}_E))$ .

*Proof.* Let,  $\mathbf{z}$  be the fixed encoder output from input  $\mathbf{x}$  i.e.  $\mathbf{z} = E(\mathbf{x}; \boldsymbol{\theta}_E)$ . Let,  $p(\mathbf{x}, t, s)$  be the true joint distribution of the variables, i.e. input  $\mathbf{x}$ , target label  $t$  and sensitive label  $s$ . The fixed encoder is a deterministic transformation of  $\mathbf{x}$  and generates an implicit distribution  $p(\mathbf{z}, t, s)$ .

*Discriminator:* The objective of the discriminator is,

$$\begin{aligned} V_1(\boldsymbol{\theta}_E, \boldsymbol{\theta}_D) &= KL(p(s|\mathbf{x}) \| q_D(s|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D)) \\ &= \mathbb{E}_{(\mathbf{z}, t, s) \sim p(\mathbf{z}, t, s)} -\log q_D(s|\mathbf{z}; \boldsymbol{\theta}_D) \\ &= -\sum_{\mathbf{x}, t, s} p(\mathbf{x}, t, s) \log q_D(s|\mathbf{z}; \boldsymbol{\theta}_D) \quad (1) \\ \text{s.t. } \sum_s q_D(s|\mathbf{z}; \boldsymbol{\theta}_D) &= 1, \quad \forall \mathbf{z} \\ q_D(s|\mathbf{z}; \boldsymbol{\theta}_D) &\geq 0, \quad \forall \mathbf{z} \end{aligned}$$

The Lagrangian dual of the problem can be written as

$$L = V_1(\boldsymbol{\theta}_E, \boldsymbol{\theta}_D) + \sum_{\mathbf{z}} \lambda(\mathbf{z}) \left( \sum_s q_D(s|\mathbf{z}; \boldsymbol{\theta}_D) - 1 \right)$$

Now we take partial derivative of  $L$  w.r.t.  $q_D(s|\mathbf{z}; \boldsymbol{\theta}_D^*)$ , the distribution of optimal discriminator. Therefore, the opti-

mal discriminator satisfies,

$$\begin{aligned} \frac{\partial L}{\partial q_D(s|\mathbf{z}; \boldsymbol{\theta}_D^*)} &= 0 \\ \Rightarrow -\frac{\sum_t p(\mathbf{z}, t, s)}{q_D(s|\mathbf{z}; \boldsymbol{\theta}_D^*)} + \lambda(\mathbf{z}) &= 0 \quad (2) \\ \Rightarrow \lambda(\mathbf{z}) q_D(s|\mathbf{z}; \boldsymbol{\theta}_D^*) &= p(\mathbf{z}, s) \end{aligned}$$

where we used the fact that,  $\sum_t p(\mathbf{z}, t, s) = p(\mathbf{z}, s)$ . Now summing w.r.t. to variable  $s$  on the both sides of last line and using the fact that  $\sum_s q_D(s|\mathbf{z}; \boldsymbol{\theta}_D^*) = 1$  we get,

$$\lambda(\mathbf{z}) = p(\mathbf{z})$$

By substituting  $\lambda(\mathbf{z})$  we obtain the solution for the optimal discriminator,

$$q_D(s|\mathbf{z}; \boldsymbol{\theta}_D^*) = \frac{p(\mathbf{z}, s)}{p(\mathbf{z})} = p(s|\mathbf{z}) \quad (3)$$

Therefore,

$$q_D(s|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) = p(s|E(\mathbf{x}; \boldsymbol{\theta}_E))$$

*Target Predictor:* The objective of the predictor is,

$$\begin{aligned} V_2(\boldsymbol{\theta}_E, \boldsymbol{\theta}_T) &= KL(p(t|\mathbf{x}) \| q_T(t|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_T)) \\ &= \mathbb{E}_{(\mathbf{z}, t, s) \sim p(\mathbf{z}, t, s)} -\log q_T(t|\mathbf{z}; \boldsymbol{\theta}_T) \\ &= -\sum_{\mathbf{x}, t, s} p(\mathbf{x}, t, s) \log q_T(t|\mathbf{z}; \boldsymbol{\theta}_T) \quad (4) \end{aligned}$$

$$\begin{aligned} \text{s.t. } \sum_t q_T(t|\mathbf{z}; \boldsymbol{\theta}_T) &= 1, \quad \forall \mathbf{z} \\ q_T(t|\mathbf{z}; \boldsymbol{\theta}_T) &\geq 0, \quad \forall \mathbf{z} \end{aligned}$$

The Lagrangian dual of the problem can be written as

$$L = V_2(\boldsymbol{\theta}_E, \boldsymbol{\theta}_T) + \sum_{\mathbf{z}} \lambda(\mathbf{z}) \left( \sum_t q_T(t|\mathbf{z}; \boldsymbol{\theta}_T) - 1 \right)$$

Now we take partial derivative of  $L$  w.r.t.  $q_T(t|\mathbf{z}; \boldsymbol{\theta}_T^*)$ , the distribution of optimal predictor. The optimal predictor satisfies the equation.

$$\begin{aligned} \frac{\partial L}{\partial q_T(t|\mathbf{z}; \boldsymbol{\theta}_T^*)} &= 0 \\ \Rightarrow -\frac{\sum_s p(\mathbf{z}, t, s)}{q_T(t|\mathbf{z}; \boldsymbol{\theta}_T^*)} + \lambda(\mathbf{z}) &= 0 \\ \Rightarrow \lambda(\mathbf{z}) q_T(t|\mathbf{z}; \boldsymbol{\theta}_T^*) &= p(\mathbf{z}, t) \end{aligned} \quad (5)$$

where we used the fact that,  $\sum_s p(\mathbf{z}, t, s) = p(\mathbf{z}, t)$ . Now summing w.r.t. to variable  $t$  on the both sides of last line and using the fact that  $\sum_t q_T(t|\mathbf{z}; \boldsymbol{\theta}_T^*) = 1$  we get,

$$\lambda(\mathbf{z}) = p(\mathbf{z})$$

By substituting  $\lambda(\mathbf{z})$  we obtain the solution of the optimal discriminator

$$q_T(t|\mathbf{z}; \boldsymbol{\theta}_T^*) = \frac{p(\mathbf{z}, t)}{p(\mathbf{z})} = p(t|\mathbf{z}) \quad (6)$$

Therefore,

$$q_T(t|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_T^*) = p(t|E(\mathbf{x}; \boldsymbol{\theta}_E))$$

□

## 2. Proof of Corollary 1.1

**Corollary 1.1.** When  $s \perp\!\!\!\perp t$ , let the optimum discriminator and predictor for an encoder  $E$  be  $q_D(s|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*)$  and  $q_T(t|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_T^*)$  respectively. The optimal encoder  $E(\cdot)$  in the MaxEnt-ARL formulation induces a uniform distribution in the discriminator  $q_D(s|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*)$  over the classes of the sensitive attribute.

*Proof.* Here we will prove that, when discriminator is fixed, then the encoder learns a representation of data  $\mathbf{x}$  such that  $q_D(s|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) = 1/m$ . First we note that although the discriminator is fixed, the discriminator probability  $q_D(s|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*)$  can change by changing the encoder parameters  $\boldsymbol{\theta}_E$ . Optimization of the encoder in MaxEnt-ARL is formulated as:

$$\begin{aligned} \min V &= \min_{\boldsymbol{\theta}_E} \mathbb{E}_{(\mathbf{x}, t, s) \sim p(\mathbf{x}, t, s)} [-\log q_T(t|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_T^*)] \\ &+ \alpha \mathbb{E}_{\mathbf{x}} \left[ \sum_{i=1}^m q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) \log q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) \right] \\ &+ \log m \end{aligned}$$

$$\text{s.t. } \sum_{i=1}^m q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) = 1$$

$$q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) \geq 0, \quad \forall i$$

(7)

The Lagrangian dual of the problem can be written as,

$$L = V - \lambda \left( \sum_{i=1}^m q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) - 1 \right)$$

Here  $\lambda$  is a Lagrangian multiplier and is assumed to be a constant in the absence of any further information. Since  $s \perp\!\!\!\perp t$ , we have  $q_T(t|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_T^*)$  is independent of  $q_D(s|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*)$  given  $E(\mathbf{x}; \boldsymbol{\theta}_E)$  from Theorem 1. Therefore, if we take derivative of  $L$  w.r.t.  $q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*)$  and set it to zero we have:

$$\begin{aligned} \frac{\partial L}{\partial q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*)} &= 0 \\ \Rightarrow 1 + \log(q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) - \lambda) &= 0 \\ \Rightarrow q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) &= \exp(\lambda - 1) \end{aligned} \quad (8)$$

Using the first (non-trivial) constraint, we have

$$\sum_{i=1}^m q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) = 1$$

$$\sum_{i=1}^m \exp(\lambda - 1) = 1$$

$$\exp(\lambda - 1) \sum_{i=1}^m 1 = 1$$

$$m(\exp(\lambda - 1)) = 1$$

$$\lambda = \log(1/m) + 1$$

Hence, the probability distribution of the discriminator after the encoder's parameters  $\boldsymbol{\theta}_E$  are optimized is  $q_D(s_i|E(\mathbf{x}; \boldsymbol{\theta}_E); \boldsymbol{\theta}_D^*) = 1/m$ . Thus, when the optimum discriminator parameters are fixed, the encoder optimizes the representation such that the discriminator does not leak any information, i.e., it induces a uniform distribution. □

## 3. CIFAR Trade-Off

We report the numerical values of the target accuracy and adversary accuracy trade-off results on the CIFAR-10 and CIFAR-100 experiments in Table 1 and Table 3, respectively. Similarly, we report the numerical values of the target accuracy and adversary entropy trade-off results on the CIFAR-10 and CIFAR-100 experiments in Table 2 and Table 4, respectively.

Target Accuracy (%)	97.75	97.73	97.68	Target Accuracy (%)	97.52	97.44	97.35	91.52	91.15	60.00
Adversary Accuracy (%)	23.44	23.09	22.68	Adversary Accuracy (%)	20.83	20.77	20.64	19.68	14.27	10.00

(a) No Privacy

Target Accuracy (%)	97.78	97.74	97.53	96.79	95.01	92.34	61.17
Adversary Accuracy (%)	23.44	22.91	21.17	21.14	19.05	12.00	10.64

(b) ML-ARL

(c) MaxEnt-ARL

Table 1: CIFAR-10: Target Accuracy (%) vs Adversary Accuracy

Target Accuracy (%)	97.75	97.73	97.71	Target Accuracy (%)	97.52	97.50	96.58	95.97	60.00
Adversary Entropy (nats)	1.65	1.65	1.67	Adversary Entropy (nats)	1.65	1.66	1.80	2.16	2.30

(a) No Privacy

Target Accuracy (%)	97.78	97.74	97.58	97.53	97.14	96.79	95.76	92.34	61.17
Adversary Entropy (nats)	1.65	1.66	1.78	2.11	2.26	2.26	2.27	2.27	2.29

(b) ML-ARL

(c) MaxEnt-ARL

Table 2: CIFAR-10: Target Accuracy (%) vs Adversary Entropy

Target Accuracy (%)	71.99	71.56	Target Accuracy (%)	71.32	70.52	70.43	69.98	69.42	24.66	22.22	5.00
Adversary Accuracy (%)	30.69	30.59	Adversary Accuracy (%)	15.43	15.09	14.84	14.60	14.41	6.81	6.72	1.00

(a) No Privacy

Target Accuracy (%)	71.17	70.80	70.50	67.63	63.81	61.98	60.03	59.11	5.37	5.00
Adversary Accuracy (%)	16.88	16.60	16.43	13.23	8.38	5.02	3.80	2.81	1.23	1.00

(b) ML-ARL

(c) MaxEnt-ARL

Table 3: CIFAR-100: Target Accuracy (%) vs Adversary Accuracy

Target Accuracy (%)	71.99	Target Accuracy (%)	71.32	64.90	56.99	54.46	24.66	22.22	5.00
Adversary Entropy (nats)	2.09	Adversary Entropy (nats)	2.50	2.51	2.68	2.88	3.77	3.88	4.60

(a) No Privacy

Target Accuracy (%)	71.17	71.05	70.80	67.63	67.38	65.71	63.81	61.98	59.11	56.32	5.37	5.00
Adversary Entropy (nats)	2.27	2.28	2.31	2.91	3.01	3.24	4.14	4.56	4.57	4.57	4.59	4.60

(b) ML-ARL

(c) MaxEnt-ARL

Table 4: CIFAR-100: Target Accuracy vs Adversary Entropy