

## A Theoretical Analysis

In this section, we provide detailed proofs of the theoretical results in the KDGAN framework. Let  $p_\alpha^e(\mathbf{y}|\mathbf{x}) = \alpha p_c(\mathbf{y}|\mathbf{x}) + (1 - \alpha)p_t^e(\mathbf{y}|\mathbf{x})$ , which is referred to as the mixture distribution. We first show that the optimal distribution of the discriminator balances between the true data distribution  $p_u(\mathbf{y}|\mathbf{x})$  and the mixture distribution  $p_\alpha^e(\mathbf{y}|\mathbf{x})$ , as stated below.

**Lemma 4.1.** *For any fixed classifier and teacher, the value function  $U(c, t, d)$  is maximized if and only if the distribution of the discriminator is given by  $p_d^e(\mathbf{x}, \mathbf{y}) = p_u(\mathbf{y}|\mathbf{x}) / (p_u(\mathbf{y}|\mathbf{x}) + p_\alpha^e(\mathbf{y}|\mathbf{x}))$ .*

*Proof.* Given the classifier  $p_c(\mathbf{y}|\mathbf{x})$  and the teacher  $p_t^e(\mathbf{y}|\mathbf{x})$ , the discriminator aims to maximize the value function  $U(c, t, d)$  of the minimax game as

$$\begin{aligned}
\max_d U(c, t, d) &= \mathbb{E}_{\mathbf{y} \sim p_u} [\log p_d^e(\mathbf{x}, \mathbf{y})] + \alpha \mathbb{E}_{\mathbf{y} \sim p_c} [\log(1 - p_d^e(\mathbf{x}, \mathbf{y}))] + (1 - \alpha) \mathbb{E}_{\mathbf{y} \sim p_t^e} [\log(1 - p_d^e(\mathbf{x}, \mathbf{y}))] \\
&= \mathbb{E}_{\mathbf{y} \sim p_u} [\log p_d^e(\mathbf{x}, \mathbf{y})] + \alpha \sum_{\mathbf{y}} p_c(\mathbf{y}|\mathbf{x}) \log(1 - p_d^e(\mathbf{x}, \mathbf{y})) + (1 - \alpha) \sum_{\mathbf{y}} p_t^e(\mathbf{y}|\mathbf{x}) \log(1 - p_d^e(\mathbf{x}, \mathbf{y})) \\
&= \mathbb{E}_{\mathbf{y} \sim p_u} [\log p_d^e(\mathbf{x}, \mathbf{y})] + \sum_{\mathbf{y}} (\alpha p_c(\mathbf{y}|\mathbf{x}) + (1 - \alpha)p_t^e(\mathbf{y}|\mathbf{x})) \log(1 - p_d^e(\mathbf{x}, \mathbf{y})) \\
&= \mathbb{E}_{\mathbf{y} \sim p_u} [\log p_d^e(\mathbf{x}, \mathbf{y})] + \sum_{\mathbf{y}} p_\alpha^e(\mathbf{y}|\mathbf{x}) \log(1 - p_d^e(\mathbf{x}, \mathbf{y})) \\
&= \sum_{\mathbf{y}} p_u(\mathbf{y}|\mathbf{x}) \log p_d^e(\mathbf{x}, \mathbf{y}) + \sum_{\mathbf{y}} p_\alpha^e(\mathbf{y}|\mathbf{x}) \log(1 - p_d^e(\mathbf{x}, \mathbf{y})) \\
&= F(p_d^e(\mathbf{x}, \mathbf{y})).
\end{aligned}$$

The function  $F(p_d^e(\mathbf{x}, \mathbf{y}))$  achieves the maximum if and only if the distribution of the discriminator is equivalent to  $p_d^e(\mathbf{x}, \mathbf{y}) = p_u(\mathbf{y}|\mathbf{x}) / (p_u(\mathbf{y}|\mathbf{x}) + p_\alpha^e(\mathbf{y}|\mathbf{x}))$ , completing the proof.  $\square$

Next, we show that the equilibrium of the minimax game is achieved if and only if both the classifier and the teacher perfectly model the true data distribution, which is summarized as follows.

**Theorem 4.2.** *The equilibrium of the minimax game  $\min_{c,t} \max_d U(c, t, d)$  is achieved if and only if  $p_c(\mathbf{y}|\mathbf{x}) = p_t^e(\mathbf{y}|\mathbf{x}) = p_u(\mathbf{y}|\mathbf{x})$ . At that point,  $U(c, t, d)$  reaches the value  $-\log(4)$ .*

*Proof.* Let  $\mathcal{L}_{\text{MD}} = \beta \mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^e(\mathbf{y}|\mathbf{x})) + \gamma \mathcal{L}_{\text{DS}}^t(p_t^e(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x}))$ . Given the optimal distribution of the discriminator in Lemma 4.1, the classifier and the teacher aim to minimize the value function  $U(c, t, d)$  of the minimax game as follows,

$$\begin{aligned}
\min_{s,t} U(c, t, d) &= \sum_{\mathbf{y}} p_u(\mathbf{y}|\mathbf{x}) \log \frac{p_u(\mathbf{y}|\mathbf{x})}{p_u(\mathbf{y}|\mathbf{x}) + p_\alpha^e(\mathbf{y}|\mathbf{x})} + \sum_{\mathbf{y}} p_\alpha^e(\mathbf{y}|\mathbf{x}) \log(1 - \frac{p_u(\mathbf{y}|\mathbf{x})}{p_u(\mathbf{y}|\mathbf{x}) + p_\alpha^e(\mathbf{y}|\mathbf{x})}) + \mathcal{L}_{\text{MD}} \\
&= \sum_{\mathbf{y}} p_u(\mathbf{y}|\mathbf{x}) \log \frac{p_u(\mathbf{y}|\mathbf{x})}{p_u(\mathbf{y}|\mathbf{x}) + p_\alpha^e(\mathbf{y}|\mathbf{x})} + \sum_{\mathbf{y}} p_\alpha^e(\mathbf{y}|\mathbf{x}) \log \frac{p_\alpha^e(\mathbf{y}|\mathbf{x})}{p_u(\mathbf{y}|\mathbf{x}) + p_\alpha^e(\mathbf{y}|\mathbf{x})} + \mathcal{L}_{\text{MD}} \\
&= -\log(4) + \mathcal{L}_{\text{KL}}(p_u(\mathbf{y}|\mathbf{x}) || \frac{p_u(\mathbf{y}|\mathbf{x}) + p_\alpha^e(\mathbf{y}|\mathbf{x})}{2}) + \mathcal{L}_{\text{KL}}(p_\alpha^e(\mathbf{y}|\mathbf{x}) || \frac{p_u(\mathbf{y}|\mathbf{x}) + p_\alpha^e(\mathbf{y}|\mathbf{x})}{2}) + \mathcal{L}_{\text{MD}} \\
&= -\log(4) + 2\mathcal{L}_{\text{JS}}(p_u(\mathbf{y}|\mathbf{x}) || p_\alpha^e(\mathbf{y}|\mathbf{x})) + \beta \mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^e(\mathbf{y}|\mathbf{x})) + \gamma \mathcal{L}_{\text{DS}}^t(p_t^e(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})).
\end{aligned}$$

Here,  $\mathcal{L}_{\text{KL}}$  is the Kullback–Leibler divergence.  $\mathcal{L}_{\text{JS}}$  is the Jensen–Shannon divergence which is non-negative and reaches zero if and only if  $p_u(\mathbf{y}|\mathbf{x}) = p_\alpha^e(\mathbf{y}|\mathbf{x})$ . The distillation losses  $\mathcal{L}_{\text{DS}}^c$  and  $\mathcal{L}_{\text{DS}}^t$  such as the L2 loss on logits [7] and the Kullback–Leibler divergence on distributions [23] achieve the minimum at zero if and only if  $p_c(\mathbf{y}|\mathbf{x}) = p_t^e(\mathbf{y}|\mathbf{x})$ . Therefore, the value function  $U(c, t, d)$  reaches the minimum at  $-\log(4)$  if and only if  $p_c(\mathbf{y}|\mathbf{x}) = p_t^e(\mathbf{y}|\mathbf{x}) = p_\alpha^e(\mathbf{y}|\mathbf{x}) = p_u(\mathbf{y}|\mathbf{x})$ , which completes the proof.  $\square$

Further, we show that the high variance of a random variance can be reduced with a low-variance random variance, which is summarized in Lemma 4.3.

**Lemma 4.3.** *Let  $X$  and  $Y$  be random variables with  $\text{Var}(X) \leq \text{Var}(Y)$ . Let  $Z = \lambda X + (1 - \lambda)Y$ , then we have  $\text{Var}(Z) \leq \text{Var}(Y)$  for all  $\lambda \in (0, 1)$ .*

*Proof.* Given  $\text{Var}(X) \leq \text{Var}(Y)$ , the covariance  $\text{Cov}(X, Y)$  is less than or equal to  $\text{Var}(Y)$  because

$$\text{Cov}(X, Y) \leq |\text{Cov}(X, Y)| \leq \sqrt{\text{Var}(X) \text{Var}(Y)} \leq \sqrt{\text{Var}(Y) \text{Var}(Y)} \leq \text{Var}(Y).$$

According to the properties of the variance, for all  $\lambda \in (0, 1)$ , we have

$$\begin{aligned} \text{Var}(Z) &= \lambda^2 \text{Var}(X) + 2\lambda(1 - \lambda) \text{Cov}(X, Y) + (1 - \lambda)^2 \text{Var}(Y) \\ &\leq \lambda^2 \text{Var}(Y) + 2\lambda(1 - \lambda) \text{Cov}(X, Y) + (1 - \lambda)^2 \text{Var}(Y) \\ &\leq \lambda^2 \text{Var}(Y) + 2\lambda(1 - \lambda) \text{Var}(Y) + (1 - \lambda)^2 \text{Var}(Y) \\ &= \text{Var}(Y), \end{aligned}$$

This completes the proof.  $\square$

## B Gradient Derivation

We provide detailed derivations of the gradient computation in the KDGAN framework. Similar to the definition of the concrete distribution  $q_c(\mathbf{y}|\mathbf{x})$  for the classifier in Equation 9, we first define a concrete distribution  $q_t^g(\mathbf{y}|\mathbf{x})$  for the teacher as follows,

$$q_t^g(\mathbf{y}|\mathbf{x}) = \text{softmax}\left(\frac{\log p_t^g(\mathbf{y}|\mathbf{x}) + \mathbf{g}}{\tau}\right), \quad \mathbf{g} \sim \text{Gumbel}(0, 1),$$

where  $\tau \in (0, +\infty)$  is a temperature parameter and  $\text{Gumbel}(0, 1)$  is the Gumbel distribution [31]. The classifier and the teacher generate continuous samples from the concrete distributions  $q_c(\mathbf{y}|\mathbf{x})$  and  $q_t^g(\mathbf{y}|\mathbf{x})$ , respectively, and then discretize the continuous samples into pseudo labels. The discriminator aims to maximize the probability of correctly identifying the true labels as positive and the pseudo labels as negative. The discriminator is trained to maximize the value function  $U(c, t, d)$  of the minimax game by ascending along its gradients

$$\begin{aligned} \nabla_d U(c, t, d) &= \nabla_d (\mathbb{E}_{\mathbf{y} \sim p_u} [\log p_d^g(\mathbf{x}, \mathbf{y})] + \alpha \mathbb{E}_{\mathbf{y} \sim p_c} [\log(1 - p_d^g(\mathbf{x}, \mathbf{y}))] + (1 - \alpha) \mathbb{E}_{\mathbf{y} \sim p_t^g} [\log(1 - p_d^g(\mathbf{x}, \mathbf{y}))]) \\ &\approx \frac{1}{k} \sum_{i=1}^k (\nabla_d \log p_d^g(\mathbf{x}, \mathbf{y}_i) + \alpha \nabla_d \log(1 - p_d^g(\mathbf{x}, \mathbf{z}_i^c)) + (1 - \alpha) \nabla_d \log(1 - p_d^g(\mathbf{x}, \mathbf{z}_i^t))). \end{aligned}$$

Here,  $k$  is the number of samples used to estimate the gradients. The true label  $\mathbf{y}_i$  is sampled from the true data distribution  $p_u(\mathbf{y}|\mathbf{x})$ .  $\mathbf{z}_i^c = \text{onehot}(\arg\max \mathbf{y}_i^c)$  and  $\mathbf{z}_i^t = \text{onehot}(\arg\max \mathbf{y}_i^t)$  are pseudo labels where  $\mathbf{y}_i^c \sim q_c(\mathbf{y}|\mathbf{x})$  and  $\mathbf{y}_i^t \sim q_t^g(\mathbf{y}|\mathbf{x})$  are continuous samples.

The classifier aims to generate the pseudo labels that resemble the true labels and predict the soft labels produced by the teacher. The classifier is trained to minimize the value function  $U(c, t, d)$  of the minimax game by descending along its gradients

$$\begin{aligned} \nabla_c U(c, t, d) &= \nabla_c (\alpha \mathbb{E}_{\mathbf{y} \sim p_c} [\log(1 - p_d^g(\mathbf{x}, \mathbf{y}))] + \beta \mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^g(\mathbf{y}|\mathbf{x}))) \\ &= \alpha \nabla_c \sum_{\mathbf{y}} p_c(\mathbf{y}|\mathbf{x}) \log(1 - p_d^g(\mathbf{x}, \mathbf{y})) + \beta \nabla_c \mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^g(\mathbf{y}|\mathbf{x})) \\ &= \alpha \sum_{\mathbf{y}} \nabla_c p_c(\mathbf{y}|\mathbf{x}) \log(1 - p_d^g(\mathbf{x}, \mathbf{y})) + \beta \nabla_c \mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^g(\mathbf{y}|\mathbf{x})) \\ &= \alpha \sum_{\mathbf{y}} p_c(\mathbf{y}|\mathbf{x}) \nabla_c \log p_c(\mathbf{y}|\mathbf{x}) \log(1 - p_d^g(\mathbf{x}, \mathbf{y})) + \beta \nabla_c \mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^g(\mathbf{y}|\mathbf{x})) \\ &= \alpha \mathbb{E}_{\mathbf{y} \sim p_c} [\nabla_c \log p_c(\mathbf{y}|\mathbf{x}) \log(1 - p_d^g(\mathbf{x}, \mathbf{y}))] + \beta \nabla_c \mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^g(\mathbf{y}|\mathbf{x})) \\ &\approx \frac{\alpha}{k} \sum_{i=1}^k \nabla_c \log q_c(\mathbf{y}_i^c|\mathbf{x}) \log(1 - p_d^g(\mathbf{x}, \mathbf{z}_i^c)) + \beta \nabla_c \mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^g(\mathbf{y}|\mathbf{x})), \end{aligned}$$

where  $\mathbf{z}_i^c = \text{onehot}(\arg\max \mathbf{y}_i^c)$  is a pseudo label and  $\mathbf{y}_i^c \sim q_c(\mathbf{y}|\mathbf{x})$  is a continuous sample. At the training of the classifier, we use a control variate [49], which is defined as

$$b_c = \mathbb{E}_{\mathbf{y} \sim p_c(\mathbf{y}|\mathbf{x})} [\log(1 - p_d^g(\mathbf{x}, \mathbf{y}))] \approx \sum_{i=1}^k \log(1 - p_d^g(\mathbf{x}, \mathbf{z}_i^c)),$$

where  $\mathbf{z}_i^c = \text{onehot}(\arg\max \mathbf{y}_i^c)$  is obtained by discretizing a continuous sample  $\mathbf{y}_i^c \sim q_c(\mathbf{y}|\mathbf{x})$ .  $\nabla_c \mathcal{L}_{\text{DS}}^c$  is the gradients of the distillation loss  $\mathcal{L}_{\text{DS}}^c$  w.r.t. the classifier, which can be easily computed

by the back-propagation algorithm. For example, if we use the L2 loss on logits [7] to define the distillation loss  $\mathcal{L}_{\text{DS}}^c$  as

$$\mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^o(\mathbf{y}|\mathbf{x})) = \frac{1}{2} \|\log p_c(\mathbf{y}|\mathbf{x}) - \log p_t^o(\mathbf{y}|\mathbf{x})\|_2^2,$$

the gradients  $\nabla_c \mathcal{L}_{\text{DS}}^c$  are computed by

$$\nabla_c \mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^o(\mathbf{y}|\mathbf{x})) = \|\log p_c(\mathbf{y}|\mathbf{x}) - \log p_t^o(\mathbf{y}|\mathbf{x})\|_2 \nabla_c \log p_c(\mathbf{y}|\mathbf{x}).$$

Similarly, the gradients to update the teacher are derived as follows,

$$\begin{aligned} \nabla_t U(c, t, d) &= \nabla_t ((1 - \alpha) \mathbb{E}_{\mathbf{y} \sim p_t^o} [\log(1 - p_d^o(\mathbf{x}, \mathbf{y}))] + \gamma \mathcal{L}_{\text{DS}}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x}))) \\ &= (1 - \alpha) \sum_{\mathbf{y}} \nabla_t p_t^o(\mathbf{y}|\mathbf{x}) \log(1 - p_d^o(\mathbf{x}, \mathbf{y})) + \gamma \nabla_t \mathcal{L}_{\text{DS}}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})) \\ &= (1 - \alpha) \sum_{\mathbf{y}} \nabla_t p_t^o(\mathbf{y}|\mathbf{x}) \log(1 - p_d^o(\mathbf{x}, \mathbf{y})) + \gamma \nabla_t \mathcal{L}_{\text{DS}}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})) \\ &= (1 - \alpha) \sum_{\mathbf{y}} p_t^o(\mathbf{y}|\mathbf{x}) \nabla_t \log p_t^o(\mathbf{y}|\mathbf{x}) \log(1 - p_d^o(\mathbf{x}, \mathbf{y})) + \gamma \nabla_t \mathcal{L}_{\text{DS}}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})) \\ &= (1 - \alpha) \mathbb{E}_{\mathbf{y} \sim p_t^o} [\nabla_t \log p_t^o(\mathbf{y}|\mathbf{x}) \log(1 - p_d^o(\mathbf{x}, \mathbf{y}))] + \gamma \nabla_t \mathcal{L}_{\text{DS}}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})) \\ &\approx \frac{1 - \alpha}{k} \sum_{i=1}^k \nabla_t \log q_i^o(\mathbf{y}_i^t|\mathbf{x}) \log(1 - p_d^o(\mathbf{x}, \mathbf{z}_i^t)) + \gamma \nabla_t \mathcal{L}_{\text{DS}}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})), \end{aligned}$$

where  $\mathbf{z}_i^t = \text{onehot}(\text{argmax } \mathbf{y}_i^t)$  is a pseudo label and  $\mathbf{y}_i^t \sim q_i^o(\mathbf{y}|\mathbf{x})$  is a continuous sample. At the training of the teacher, we also use a control variate [49], which is defined as

$$b_t = \mathbb{E}_{\mathbf{y} \sim p_t^o(\mathbf{y}|\mathbf{x})} [\log(1 - p_d^o(\mathbf{x}, \mathbf{y}))] \approx \sum_{i=1}^k \log(1 - p_d^o(\mathbf{x}, \mathbf{z}_i^t)),$$

where  $\mathbf{z}_i^t = \text{onehot}(\text{argmax } \mathbf{y}_i^t)$  is obtained by discretizing a continuous sample  $\mathbf{y}_i^t \sim q_i^o(\mathbf{y}|\mathbf{x})$ .  $\nabla_t \mathcal{L}_{\text{DS}}^t$  is the gradients of the distillation loss  $\mathcal{L}_{\text{DS}}^t$  w.r.t. the teacher. For example, the gradients  $\nabla_t \mathcal{L}_{\text{DS}}^t$  are given by

$$\nabla_t \mathcal{L}_{\text{DS}}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})) = \|\log p_t^o(\mathbf{y}|\mathbf{x}) - \log p_c(\mathbf{y}|\mathbf{x})\|_2 \nabla_t \log p_t^o(\mathbf{y}|\mathbf{x}),$$

when the distillation loss  $\mathcal{L}_{\text{DS}}^t$  is defined as the L2 loss on logits [7],

$$\mathcal{L}_{\text{DS}}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})) = \frac{1}{2} \|\log p_t^o(\mathbf{y}|\mathbf{x}) - \log p_c(\mathbf{y}|\mathbf{x})\|_2^2.$$

## C Network Architectures

We describe network architectures which we use to conduct experiments in deep model compression and image tag recommendation tasks. First, we describe the network architectures in deep model compression task on the MNIST dataset. We implement the scoring function  $h(\mathbf{x}, \mathbf{y})$  as an MLP [27]. The architecture of the MLP is given by

1. An input layer of a  $28 \times 28$  grayscale image.
2. A stack of 2 fully connected layers with 800 neurons.
3. A softmax layer with 10 classes.

We implement the scoring function  $s(\mathbf{x}, \mathbf{y})$  as a LeNet [27]. The architecture of the LeNet is given by

1. An input layer of a  $28 \times 28$  grayscale image.
2. A convolutional layer with 32 kernels of size  $5 \times 5$  and stride 1.
3. A max pooling layer with size  $2 \times 2$  and stride 2.
4. A convolutional layer with 64 kernels of size  $5 \times 5$  and stride 1.
5. A max pooling layer with size  $2 \times 2$  and stride 2.
6. A fully connected layer with 1024 neurons.
7. A softmax layer with 10 classes.

Next, we describe the network architectures in deep model compression task on the CIFAR-10 dataset. We implement  $h(\mathbf{x}, \mathbf{y})$  as a LeNet [27]. The architecture of the LeNet is given by

1. An input layer of a  $32 \times 32$  colored image.
2. A convolutional layer with 64 kernels of size  $5 \times 5$  and stride 1.
3. A max pooling layer with size  $2 \times 2$  and stride 2.
4. A convolutional layer with 128 kernels of size  $5 \times 5$  and stride 1.
5. A max pooling layer with size  $2 \times 2$  and stride 2.
6. A fully connected layer with 1024 neurons.
7. A softmax layer with 10 classes.

We implement  $s(\mathbf{x}, \mathbf{y})$  as a 101-layer ResNet [22]. The architecture of the ResNet is given by

1. An input layer of a  $32 \times 32$  colored image.
2. A convolutional layer with 16 kernels with size  $3 \times 3$  and stride 1.
3. Three stacked blocks of 3 convolutional layers which use 64 kernels of size  $1 \times 1$ , 64 kernels of size  $3 \times 3$ , and 256 kernels of size  $1 \times 1$ , respectively.
4. Four stacked blocks of 3 convolutional layers which use 128 kernels of size  $1 \times 1$ , 128 kernels of size  $3 \times 3$ , and 512 kernels of size  $1 \times 1$ , respectively.
5. Twenty three stacked blocks of 3 convolutional layers which use 256 kernels of size  $1 \times 1$ , 256 kernels of size  $3 \times 3$ , and 1024 kernels of size  $1 \times 1$ , respectively.
6. Three stacked blocks of 3 convolutional layers which use 512 kernels of size  $1 \times 1$ , 512 kernels of size  $3 \times 3$ , and 2048 kernels of size  $1 \times 1$ , respectively.
7. A global pooling layer.
8. A softmax layer with 10 classes.

Finally, we describe the network architectures in image tag recommendation task on the YFCC100M dataset. We use the same network architectures when experimenting with the two datasets of images labeled with the 200 most popular tags and 200 randomly sampled tags, respectively. We implement a VGGNet [40] to extract image features. The architecture of the VGGNet is written as

1. An input layer of a  $224 \times 224$  colored image.
2. A stack of 2 convolutional layers with 64 kernels of size  $3 \times 3$  and stride 1.
3. A max pooling layer with size  $2 \times 2$  and stride 2.
4. A stack of 2 convolutional layers with 128 kernels of size  $3 \times 3$  and stride 1.
5. A max pooling layer with size  $2 \times 2$  and stride 2.
6. A stack of 2 convolutional layers with 256 kernels of size  $3 \times 3$  and stride 1.
7. A max pooling layer with size  $2 \times 2$  and stride 2.
8. A stack of 2 convolutional layers with 512 kernels of size  $3 \times 3$  and stride 1.
9. A max pooling layer with size  $2 \times 2$  and stride 2.
10. A stack of 2 convolutional layers with 512 kernels of size  $3 \times 3$  and stride 1.
11. A max pooling layer with size  $2 \times 2$  and stride 2.
12. A fully connected layer with 4096 neurons.
13. A fully connected layer with 4096 neurons.
14. A fully connected layer with 100 neurons.

We implement a LSTM [24] to learn text features. The architecture of the LSTM is written as

$$\begin{aligned}
\mathbf{f}_t &= \text{sigmoid}(\mathbf{W}_f \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f), \\
\mathbf{i}_t &= \text{sigmoid}(\mathbf{W}_i \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i), \\
\mathbf{o}_t &= \text{sigmoid}(\mathbf{W}_o \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o), \\
\mathbf{s}_t &= \mathbf{f}_t \odot \mathbf{s}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{W}_s \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_s), \\
\mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{s}_t),
\end{aligned}$$

where  $[h, x]$  is the vector concatenation and  $\odot$  is the element-wise product. We set the hidden size of the LSTM to 100 in the experiments. Let  $v_x \in \mathbb{R}^{100}$  be an image feature vector extracted by the VGGNet and  $v_z \in \mathbb{R}^{100}$  be a text feature vector learned by the LSTM. We implement the scoring function  $h(x, y)$  as an MLP [3]. The architecture of the MLP is written as

1. An input layer of a feature vector with size 100 (i.e. the image features  $v_x$ ).
2. A stack of 2 fully connected layers with 800 neurons.
3. A softmax layer with 200 classes.

We implement the scoring function  $s(x, y)$  as an MLP [3]. The architecture of the MLP is given by

1. An input layer of a feature vector with size 100 (i.e. the element-wise product of  $v_x$  and  $v_z$ ).
2. A stack of 2 fully connected layers with 1200 neurons.
3. A softmax layer with 200 classes.

## D Additional Experiments

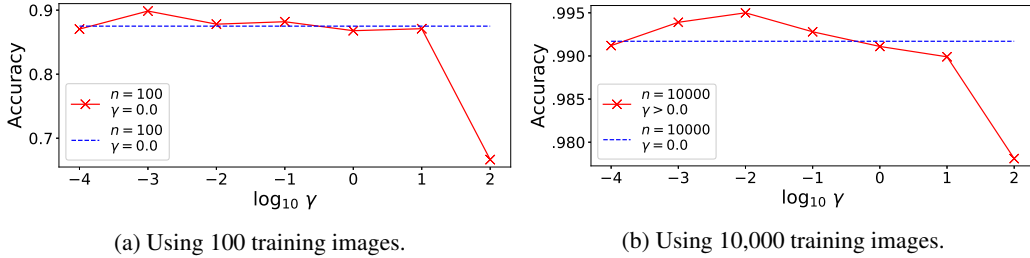


Figure 6: The accuracy of the teacher against the hyperparameter  $\gamma$  in KDGAN on MNIST. Note that  $\gamma$  controls how much the classifier distills its knowledge into the teacher.

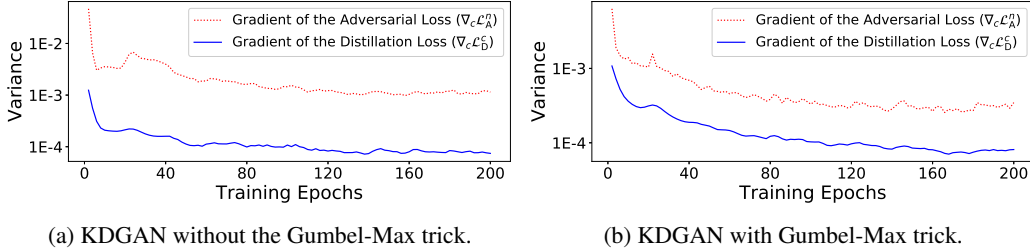


Figure 7: Variances of the gradient of the adversarial loss ( $\nabla_c \mathcal{L}_{AD}^n$ ) or the distillation loss ( $\nabla_c \mathcal{L}_{DS}^c$ ) w.r.t. the classifier. The results are obtained by training KDGAN with 100 training images on MNIST.

We study the classifier’s ratio of compression (in terms of the number of parameters) and loss of accuracy w.r.t. the teacher. The results using 5K to 50K training images on MNIST are presented in Tables 3 and 4. We can see that the loss of accuracy generally decreases as the parameter number of the classifier or the number of training examples increases. We also observe that MIMIC (1.22M parameters) achieves an accuracy of 97.93%-99.05% while KDGAN with a much smaller size (0.19M parameters) already achieves a better accuracy (98.72%-99.27%) than MIMIC.

Table 3: Model size and accuracy of the classifier and the teacher (shown in parenthesis) in KDGAN on MNIST.

#Param. (M)	$n = 5K$	$n = 10K$	$n = 50K$
0.09 (3.12)	97.96 (99.25)	98.74 (99.42)	99.03 (99.65)
0.19 (3.12)	98.72 (99.26)	98.92 (99.46)	99.27 (99.70)
1.22 (3.12)	99.01 (99.28)	99.25 (99.48)	99.54 (99.72)
2.28 (3.12)	99.04 (99.27)	99.40 (99.53)	<b>99.77 (99.78)</b>

Table 4: Average accuracy over 10 runs with varying training size ( $n$ ) on MNIST.

Method	$n = 5K$	$n = 10K$	$n = 50K$
CODIS	98.53	98.89	99.31
DISTN	98.04	98.79	99.26
MIMIC	97.93	98.65	99.05
KDGAN	<b>99.01</b>	<b>99.25</b>	<b>99.54</b>