
KDGAN: Knowledge Distillation with Generative Adversarial Networks

Xiaojie Wang

University of Melbourne
xiaojiew94@gmail.com

Rui Zhang*

University of Melbourne
rui.zhang@unimelb.edu.au

Yu Sun

Twitter Inc.
ysun@twitter.com

Jianzhong Qi

University of Melbourne
jianzhong.qi@unimelb.edu.au

김 성 철

Contents

1. Introduction
2. Methods
3. Results
4. Conclusion

Introduction

- Terminologies

- *Privileged provision*: input features or computational resources → Resources are only accessible for training!
- *Privileged information*: input features

- Example about the privileged provision between training and inference

- Image tag recommendation / face recognition

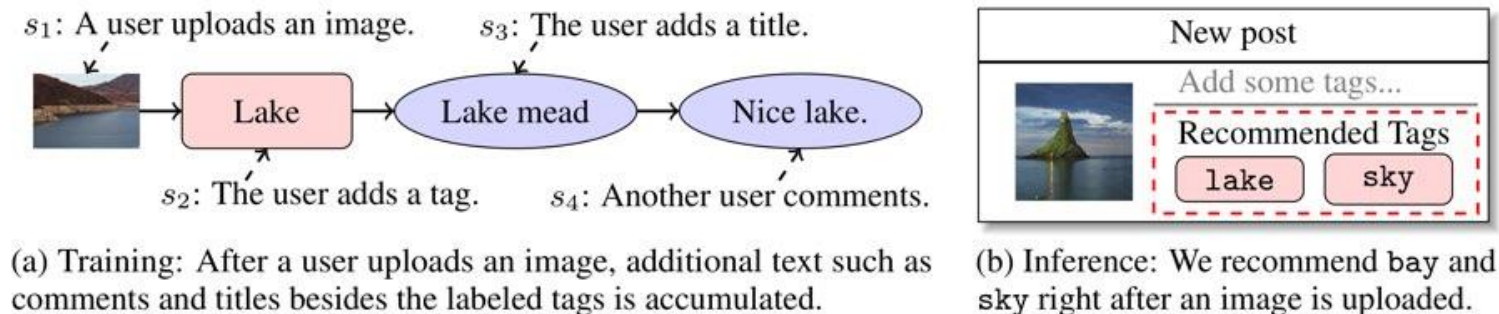
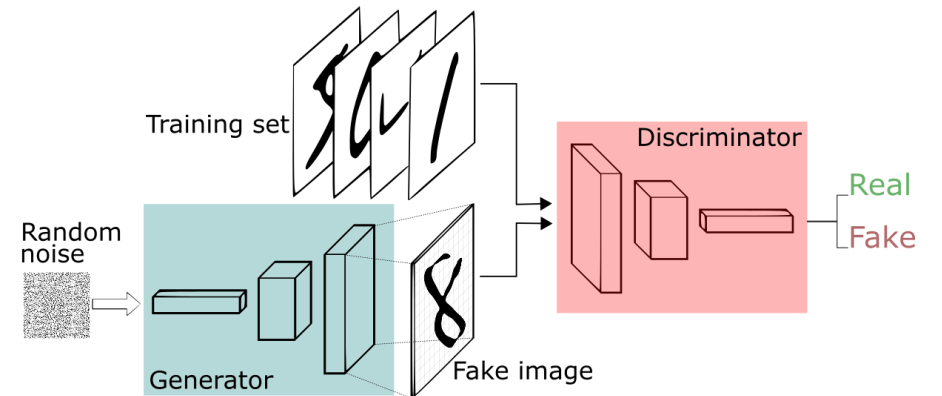
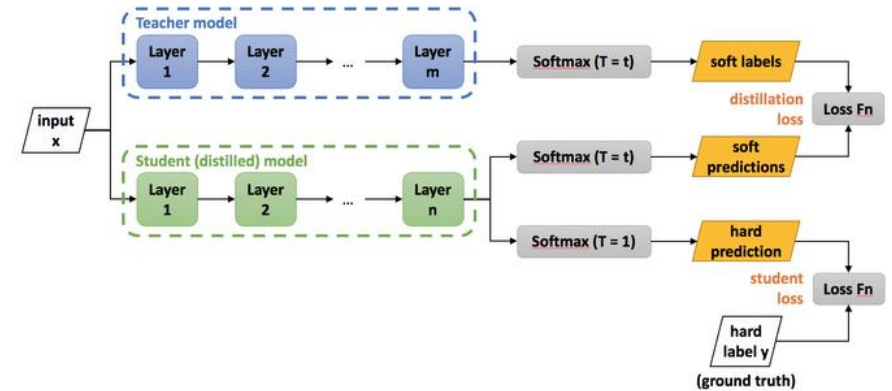


Figure 1: Image tag recommendation where the additional text is only available for training.

Introduction

- Knowledge Distillation (KD)
 - Classifier (student) / Teacher
- Generative Adversarial Networks (GAN)
 - Classifier (Generator) / Discriminator



Introduction

- Knowledge Distillation with Generative Adversarial Networks (KDGAN)
 - KD to train a lightweight classifier.
 - GAN to teach the classifier the true data distribution.

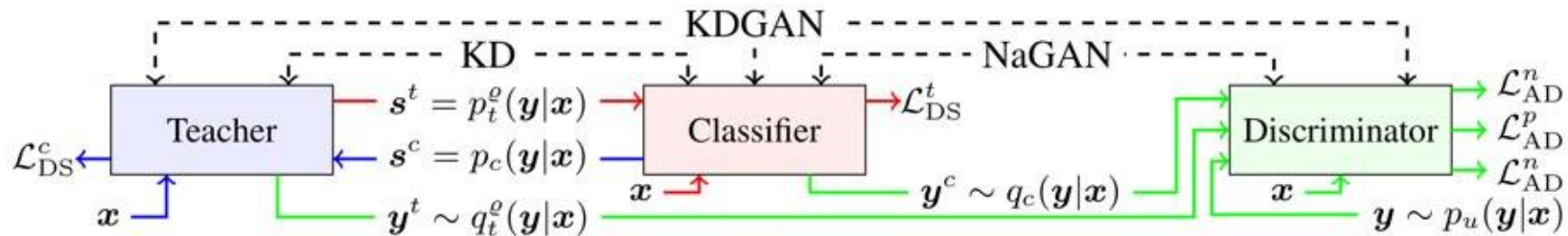


Figure 2: Comparison among KD, NaGAN, and KDGAN. The classifier (C) and the teacher (T) learn discrete categorical distributions $p_c(y|x)$ and $p_t^o(y|x)$; y is a true label generated from the true data distribution $p_u(y|x)$; y^c and y^t are continuous samples generated from concrete distributions $q_c(y|x)$ and $q_t^o(y|x)$; s^c and s^t are soft labels produced by C and T ; \mathcal{L}_{DS}^c and \mathcal{L}_{DS}^t are distillation losses for C and T ; \mathcal{L}_{AD}^p and \mathcal{L}_{AD}^n are adversarial losses for positive and negative feature-label pairs.

Methods

- Naïve GAN (NaGAN)

$$\min_c \max_d V(c, d) = \mathbb{E}_{\mathbf{y} \sim p_u} [\log p_d^g(\mathbf{x}, \mathbf{y})] + \mathbb{E}_{\mathbf{y} \sim p_c} [\log(1 - p_d^g(\mathbf{x}, \mathbf{y}))]. \quad (1)$$

Let $h(\mathbf{x}, \mathbf{y})$ and $g(\mathbf{x}, \mathbf{y})$ be the scoring functions for C and D . We define $p_c(\mathbf{y}|\mathbf{x})$ and $p_d^g(\mathbf{x}, \mathbf{y})$ as

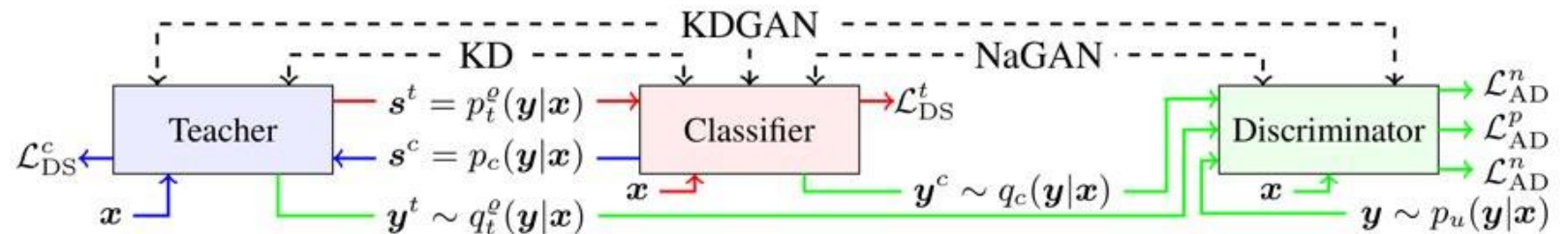
$$p_c(\mathbf{y}|\mathbf{x}) = \text{softmax}(h(\mathbf{x}, \mathbf{y})) \quad \text{and} \quad p_d^g(\mathbf{x}, \mathbf{y}) = \text{sigmoid}(g(\mathbf{x}, \mathbf{y})). \quad (2)$$

- The advantages and disadvantages of KD and NaGAN
 - KD : requires a small number of training instances and epochs, but ~.
 - NaGAN : ensures the equilibrium where $p_c(\mathbf{y}|\mathbf{x}) = p_u(\mathbf{y}|\mathbf{x})$, but ~.

Methods

- KDGAN

- a classifier C , a teacher T , a discriminator D with privileged provision ϱ
- D : maximize the probability of correctly distinguishing the true and pseudo labels.
- C, T : minimize the probability that D rejects their generated pseudo labels.
- C : learn from T by mimicking the learned distribution of T .
- T also learn from C , because a teacher's ability can also be enhanced by interacting with students. → reduce the probability of generating different pseudo labels.



Methods

- KDGAN

Algorithm 1: Minibatch stochastic gradient descent training of KDGAN.

```
1 Pretrain a classifier  $C$ , a teacher  $T$ , and a discriminator  $D$  with the training data  $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$ .
2 for the number of training epochs do
3   for the number of training steps for the discriminator do
4     Sample labels  $\{\mathbf{y}_1, \dots, \mathbf{y}_k\}$ ,  $\{\mathbf{y}_1^c, \dots, \mathbf{y}_k^c\}$ , and  $\{\mathbf{y}_1^t, \dots, \mathbf{y}_k^t\}$  from  $p_u(\mathbf{y}|\mathbf{x})$ ,  $q_c(\mathbf{y}|\mathbf{x})$ , and  $q_t^o(\mathbf{y}|\mathbf{x})$ .
5     Update  $D$  by ascending along its gradients
6      $\frac{1}{k} \sum_{i=1}^k (\nabla_d \log p_d^o(\mathbf{x}, \mathbf{y}_i) + \alpha \nabla_d \log(1 - p_d^o(\mathbf{x}, \mathbf{z}_i^c)) + (1 - \alpha) \nabla_d \log(1 - p_d^o(\mathbf{x}, \mathbf{z}_i^t)))$ .
7   for the number of training steps for the teacher do
8     Sample labels  $\{\mathbf{y}_1^t, \dots, \mathbf{y}_k^t\}$  from  $q_t^o(\mathbf{y}|\mathbf{x})$  and update the teacher by descending along its gradients
9      $\frac{1}{k} \sum_{i=1}^k (1 - \alpha) \nabla_t \log q_t^o(\mathbf{y}_i^t|\mathbf{x}) \log(1 - p_d^o(\mathbf{x}, \mathbf{z}_i^t)) + \gamma \nabla_t \mathcal{L}_{\text{DS}}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x}))$ .
10  for the number of training steps for the classifier do
11    Sample labels  $\{\mathbf{y}_1^c, \dots, \mathbf{y}_k^c\}$  from  $q_c(\mathbf{y}|\mathbf{x})$  and update  $C$  by descending along its gradients
12     $\frac{1}{k} \sum_{i=1}^k \alpha \nabla_c \log q_c(\mathbf{y}_i^c|\mathbf{x}) \log(1 - p_d^o(\mathbf{x}, \mathbf{z}_i^c)) + \beta \nabla_c \mathcal{L}_{\text{DS}}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^o(\mathbf{y}|\mathbf{x}))$ .
```

Methods

- KDGAN

- The value function $U(c, t, d)$ ($\alpha \in (0, 1), \beta \in (0, +\infty), \gamma \in (0, +\infty)$)

$$\min_{c,t} \max_d U(c, t, d) = \mathbb{E}_{\mathbf{y} \sim p_u} [\log p_d^o(\mathbf{x}, \mathbf{y})] + \alpha \mathbb{E}_{\mathbf{y} \sim p_c} [\log(1 - p_d^o(\mathbf{x}, \mathbf{y}))] + (1 - \alpha) \mathbb{E}_{\mathbf{y} \sim p_t^o} [\log(1 - p_d^o(\mathbf{x}, \mathbf{y}))] - \beta \mathcal{L}_{DS}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^o(\mathbf{y}|\mathbf{x})) + \gamma \mathcal{L}_{DS}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})), \quad (3)$$

Adversarial losses Distillation losses

- Distillation losses : L2 loss, Kullback–Leibler divergence, ...
- $\mathcal{L}_{DS}^c, \mathcal{L}_{DS}^t$ are used to train the classifier and the teacher, respectively.

Methods

- KDGAN

- The classifier perfectly learns the true data distribution.

$$\rightarrow p_{\alpha}^e(\mathbf{y}|\mathbf{x}) = \alpha p_c(\mathbf{y}|\mathbf{x}) + (1 - \alpha)p_t^e(\mathbf{y}|\mathbf{x})$$

$$\begin{aligned} & \alpha \mathbb{E}_{\mathbf{y} \sim p_c} [\log(1 - p_d^e(\mathbf{x}, \mathbf{y}))] + (1 - \alpha) \mathbb{E}_{\mathbf{y} \sim p_t^e} [\log(1 - p_d^e(\mathbf{x}, \mathbf{y}))] \\ &= \alpha \sum_{\mathbf{y}} p_c(\mathbf{y}|\mathbf{x}) \log(1 - p_d^e(\mathbf{x}, \mathbf{y})) + (1 - \alpha) \sum_{\mathbf{y}} p_t^e(\mathbf{y}|\mathbf{x}) \log(1 - p_d^e(\mathbf{x}, \mathbf{y})) \\ &= \sum_{\mathbf{y}} (\alpha p_c(\mathbf{y}|\mathbf{x}) + (1 - \alpha)p_t^e(\mathbf{y}|\mathbf{x})) \log(1 - p_d^e(\mathbf{x}, \mathbf{y})) \\ &= \mathbb{E}_{\mathbf{y} \sim p_{\alpha}^e} [\log(1 - p_d^e(\mathbf{x}, \mathbf{y}))]. \end{aligned} \tag{4}$$

$$\begin{aligned} &= \sum_{\mathbf{y}} p_u(\mathbf{y}|\mathbf{x}) \log p_d^e(\mathbf{x}, \mathbf{y}) + \sum_{\mathbf{y}} p_{\alpha}^e(\mathbf{y}|\mathbf{x}) \log(1 - p_d^e(\mathbf{x}, \mathbf{y})) \\ &= F(p_d^e(\mathbf{x}, \mathbf{y})). \end{aligned}$$

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

Methods

- KDGAN

- $\mathcal{L}_{MD} = \beta \mathcal{L}_{DS}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^o(\mathbf{y}|\mathbf{x})) + \gamma \mathcal{L}_{DS}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x}))$

- \mathcal{L}_{JS} : the Jensen-Shannon divergence

- * $\text{JSD}(p, q) = \frac{1}{2} D_{KL}(p || \frac{p+q}{2}) + \frac{1}{2} D_{KL}(q || \frac{p+q}{2})$

$$\begin{aligned} & \min_{\alpha} \max_d \mathbb{E}_{\mathbf{y} \sim p_u} [\log p_d^o(\mathbf{x}, \mathbf{y})] + \mathbb{E}_{\mathbf{y} \sim p_{\alpha}^o} [\log(1 - p_d^o(\mathbf{x}, \mathbf{y}))] + \mathcal{L}_{MD} \\ & = \min_{\alpha} 2\mathcal{L}_{JS}(p_u(\mathbf{y}|\mathbf{x}) || p_{\alpha}^o(\mathbf{y}|\mathbf{x})) + \beta \mathcal{L}_{DS}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^o(\mathbf{y}|\mathbf{x})) + \gamma \mathcal{L}_{DS}^t(p_t^o(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})) - \log(4). \end{aligned} \quad (5)$$

Methods

- KDGAN

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

Methods

- KDGAN

- \mathcal{L}_{JS} reaches the minimum if and only if $p_{\alpha}^e(\mathbf{y}|\mathbf{x}) = p_u(\mathbf{y}|\mathbf{x})$
- \mathcal{L}_{DS}^c (or \mathcal{L}_{DS}^t) reaches the minimum if and only if $p_c(\mathbf{y}|\mathbf{x}) = p_t^e(\mathbf{y}|\mathbf{x})$

→KDGAN equilibrium

$p_c(\mathbf{y}|\mathbf{x}) = p_t^e(\mathbf{y}|\mathbf{x}) = p_u(\mathbf{y}|\mathbf{x})$ (\mathcal{C} learns the true data distribution.)

$$\begin{aligned} & \min_{\alpha} \max_d \mathbb{E}_{\mathbf{y} \sim p_u} [\log p_d^e(\mathbf{x}, \mathbf{y})] + \mathbb{E}_{\mathbf{y} \sim p_{\alpha}^e} [\log(1 - p_d^e(\mathbf{x}, \mathbf{y}))] + \mathcal{L}_{MD} \quad (5) \\ & = \min_{\alpha} 2\mathcal{L}_{JS}(p_u(\mathbf{y}|\mathbf{x}) || p_{\alpha}^e(\mathbf{y}|\mathbf{x})) + \beta \mathcal{L}_{DS}^c(p_c(\mathbf{y}|\mathbf{x}), p_t^e(\mathbf{y}|\mathbf{x})) + \gamma \mathcal{L}_{DS}^t(p_t^e(\mathbf{y}|\mathbf{x}), p_c(\mathbf{y}|\mathbf{x})) - \log(4). \end{aligned}$$

Methods

- **KDGAN Training**

- The training speed is closely related to the variance of gradients.
- The KDGAN framework by design
 - can **reduce the variance of gradients!!**
 - **use a continuous space** by relaxing the discrete samples!!

Methods

- KDGAN Training

First, we show how KDGAN reduces the variance of gradients. As discussed above, C only receives gradients $\nabla_c V$ from D in NaGAN while it receives gradients $\nabla_c U$ from both D and T in KDGAN:

$$\nabla_c V = \nabla_c \mathcal{L}_{AD}^n, \quad \nabla_c U = \lambda \nabla_c \mathcal{L}_{AD}^n + (1 - \lambda) \nabla_c \mathcal{L}_{DS}^c, \quad (6)$$

where $\lambda \in (0, 1)$, $\nabla_c \mathcal{L}_{AD}^n$ and $\nabla_c \mathcal{L}_{DS}^c$ are gradients from D and T , respectively. Consistent with the findings in existing work [23, 39], we also observe that $\nabla_c \mathcal{L}_{DS}^c$ usually has a lower variance than $\nabla_c \mathcal{L}_{AD}^n$ (see Figure 7 in Appendix D for empirical evidence that the variance of $\nabla_c \mathcal{L}_{DS}^c$ is smaller than that of $\nabla_c \mathcal{L}_{AD}^n$ during the training process). Hence, it can be easily shown that the gradients w.r.t. C in KDGAN have a lower variance than that in NaGAN (refer to Lemma 4.3):

$$\text{Var}(\nabla_c \mathcal{L}_{DS}^c) \leq \text{Var}(\nabla_c \mathcal{L}_{AD}^n) \Rightarrow \text{Var}(\nabla_c U) \leq \text{Var}(\nabla_c V). \quad (7)$$

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

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

Methods

- KDGAN Training

Next, we further reduce the variance of gradients with a reparameterization trick, in particular, the Gumbel-Max trick [20, 30]. The essence of the Gumbel-Max trick is to reparameterize generating discrete samples into a differentiable function of its parameters and an additional random variable of a Gumbel distribution. To perform the Gumbel-Max trick on generating discrete samples from the categorical distribution $p_c(\mathbf{y}|\mathbf{x})$, a concrete distribution [25, 31] can be used. We use a concrete distribution $q_c(\mathbf{y}|\mathbf{x})$ to generate continuous samples and use the continuous samples to compute the gradients $\nabla_c \mathcal{L}_{\text{AD}}^n$ of the adversarial loss w.r.t. the classifier as

$$\nabla_c \mathcal{L}_{\text{AD}}^n = \nabla_c \mathbb{E}_{\mathbf{y} \sim p_c} [\log(1 - p_d^g(\mathbf{x}, \mathbf{y}))] = \mathbb{E}_{\mathbf{y} \sim q_c} [\nabla_c \log q_c(\mathbf{y}|\mathbf{x}) \log(1 - p_d^g(\mathbf{x}, \mathbf{z}))]. \quad (8)$$

Here, $\mathbf{z} = \text{onehot}(\text{argmax } \mathbf{y})$ is a discrete pseudo label where $\mathbf{y} \sim q_c(\mathbf{y}|\mathbf{x})$. We define $q_c(\mathbf{y}|\mathbf{x})$ as

$$q_c(\mathbf{y}|\mathbf{x}) = \text{softmax} \left(\frac{\log p_c(\mathbf{y}|\mathbf{x}) + \mathbf{g}}{\tau} \right), \quad \mathbf{g} \sim \text{Gumbel}(0, 1). \quad (9)$$

Here, $\tau \in (0, +\infty)$ is a temperature parameter and $\text{Gumbel}(0, 1)$ is the Gumbel distribution² [31]. We leverage the temperature parameter τ to control the variance of gradients over the training. With a high temperature, the samples from the concrete distribution are smooth, which give low-variance gradient estimates. Note that a disadvantage of the concrete distribution is that with a high temperature, it becomes a less accurate approximation to the original categorical distribution, which causes biased gradient estimates. We will discuss how to tune the temperature parameter in Section 4.

Methods

- KDGAN Training
 - Gumbel-Max trick

Gumbel-max trick(stochastic을 따로 뺌)

categorical한 variable을 reparametrization함. 요걸 쓰면 categorical에서 sample한 것과 비슷한 효과를 낸다고한다.

$x \sim \text{Cat}(\pi_\phi)$ 를 discrete categorical variable이라 해보자.

$\epsilon_k \sim \text{Gumbel}(0, 1)$ 를 가지고 Reparametrization하면

$$x = \arg \max_k (\epsilon_k + \log \pi_k) \cong g(\phi, \epsilon)$$

로 쓸 수 있다. 요것이 Gumbel-max trick

- 하지만 아직도 ϕ 에 대해 discrete하며 미분이 불가능하다.

Gumbel-softmax(continuous)

Gumbel-max의 argmax를 softmax로 바꾼 녀석. 이제 미분이 가능해진다.

$$(x_k)_{1 \leq k \leq K} = \text{softmax}((\epsilon_k + \log \pi_k)_k) \Leftrightarrow x_k = \frac{\exp((\log \pi_k + \epsilon_k)/\tau)}{\sum_j \exp((\log \pi_j + \epsilon_j)/\tau)}$$

특성들

- 여기서 τ 는 temperature parameter인데
 - 0에 가까워지면 one-hot
 - 무한대로가면 uniform distribution을 갖는다.
- $p(x_k = \max_i x_i) = \pi_k$

Results

- Experiments

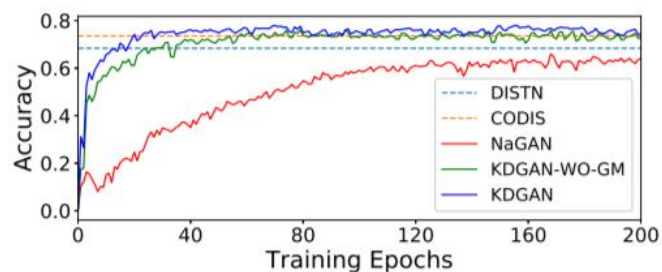
- α in $[0.0, 1.0]$
- β in $[0.001, 1000]$
- γ in $[0.0001, 100]$
- The temperature parameter τ
 - start with a large value (1.0) and exponentially decay it to a small value (0.1)

Results

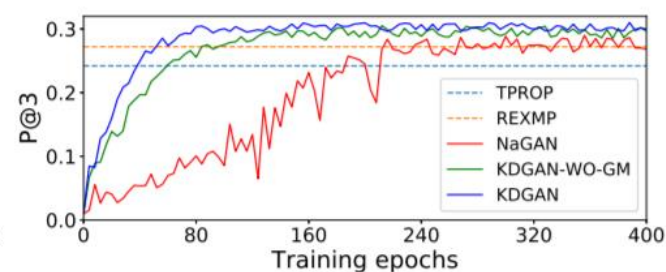
• Deep Model Compression

Table 1: Average accuracy over 10 runs in model compression (n is the number of training instances).

Method	MNIST			CIFAR-10		
	$n = 100$	$n = 1,000$	$n = 10,000$	$n = 500$	$n = 5,000$	$n = 50,000$
CODIS	74.02 ± 0.13	95.77 ± 0.10	98.89 ± 0.08	54.17 ± 0.20	77.82 ± 0.14	85.12 ± 0.11
DISTN	68.34 ± 0.06	93.97 ± 0.08	98.79 ± 0.07	50.92 ± 0.18	76.59 ± 0.15	83.32 ± 0.08
NOISY	66.53 ± 0.18	93.45 ± 0.11	98.58 ± 0.11	50.18 ± 0.28	75.42 ± 0.19	82.99 ± 0.12
MIMIC	67.35 ± 0.15	93.78 ± 0.13	98.65 ± 0.05	51.74 ± 0.23	75.66 ± 0.17	84.33 ± 0.10
NaGAN	64.90 ± 0.31	93.60 ± 0.22	98.95 ± 0.19	46.29 ± 0.32	76.11 ± 0.24	85.34 ± 0.27
KDGAN	77.95 ± 0.05	96.42 ± 0.05	99.25 ± 0.02	57.56 ± 0.13	79.36 ± 0.04	86.50 ± 0.04



(a) Deep model compression over MNIST.



(b) Image tag recommendation on YFCC100M.

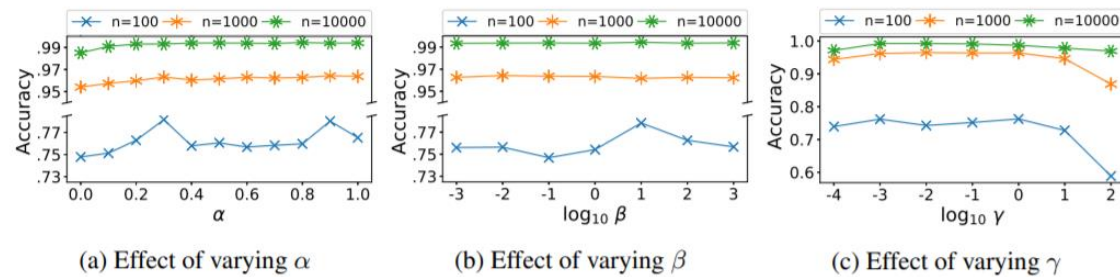


Figure 4: Effects of hyperparameters in KDGAN on MNIST for deep model compression.

Figure 3: Training curves of the classifier in the proposed NaGAN and KDGAN.

Results

- Image Tag Recommendation

Table 2: Performance of various methods on the YFCC100M dataset in tag recommendation.

Method	Most Popular Tags						Randomly Sampled Tags					
	P@3	P@5	F@3	F@5	MAP	MRR	P@3	P@5	F@3	F@5	MAP	MRR
KNN	.2320	.1680	.2339	.1633	.5755	.5852	.1623	.1198	.1575	.1088	.3970	.4092
TPROP	.2420	.1636	.2811	.1949	.6177	.6270	.1883	.1372	.1810	.1252	.4512	.4636
TFEAT	.2560	.1752	.2871	.1999	.6417	.6503	.2002	.1420	.2195	.1495	.5149	.5309
REXMP	.2720	.1800	.3324	.2295	.7015	.7122	.2228	.1378	.2427	.1669	.5205	.5331
NaGAN	.2892	.1880	.3516	.2352	.7432	.7555	.2415	.1495	.2693	.1867	.5791	.5911
KDGAN	.3047	.1968	.3678	.2526	.7787	.7905	.2572	.1666	.2946	.2009	.6302	.6452

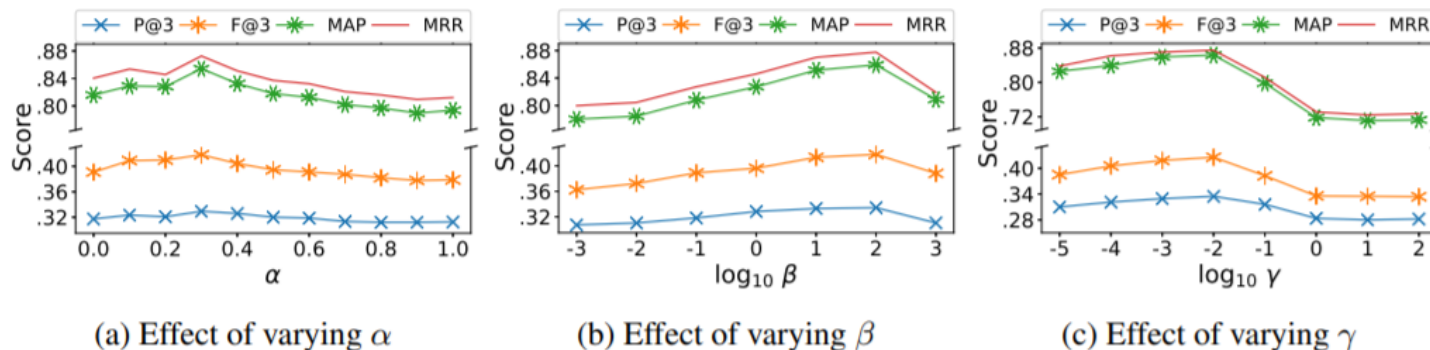


Figure 5: Effects of hyperparameters in KDGAN on YFCC100M for image tag recommendation.

감 사 합 니 다