

# TRAINING CONFIDENCE-CALIBRATED CLASSIFIERS FOR DETECTING OUT-OF-DISTRIBUTION SAMPLES

**Kimin Lee\***    **Honglak Lee<sup>§,†</sup>**    **Kibok Lee<sup>†</sup>**    **Jinwoo Shin\***

<sup>\*</sup>Korea Advanced Institute of Science and Technology, Daejeon, Korea

<sup>†</sup>University of Michigan, Ann Arbor, MI 48109

<sup>§</sup>Google Brain, Mountain View, CA 94043

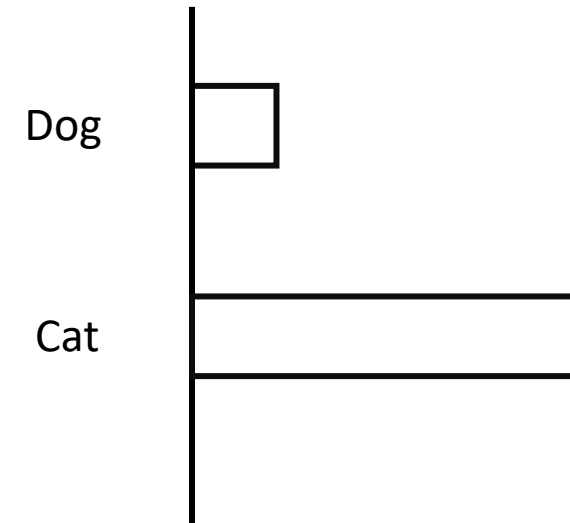
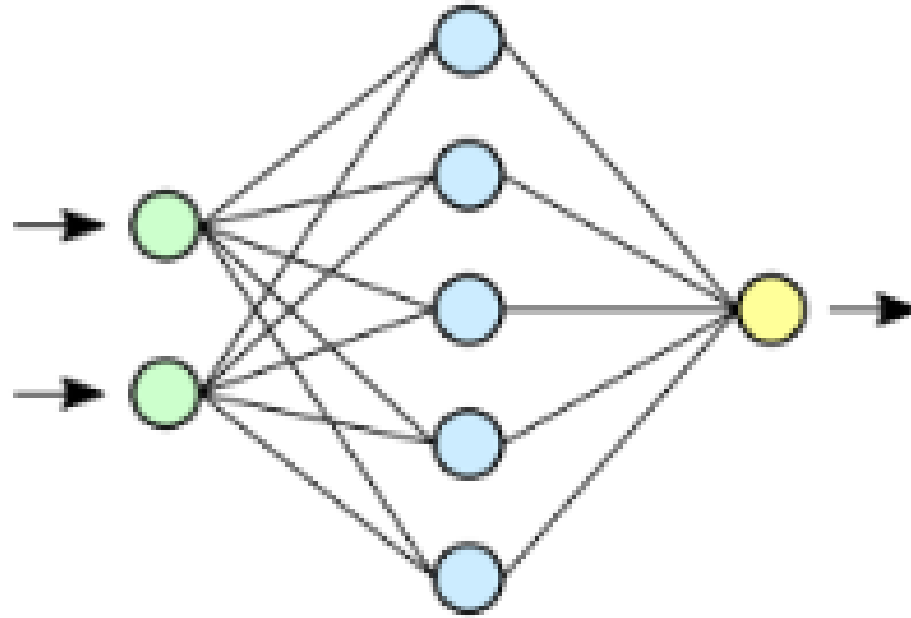
---

2023.03.24

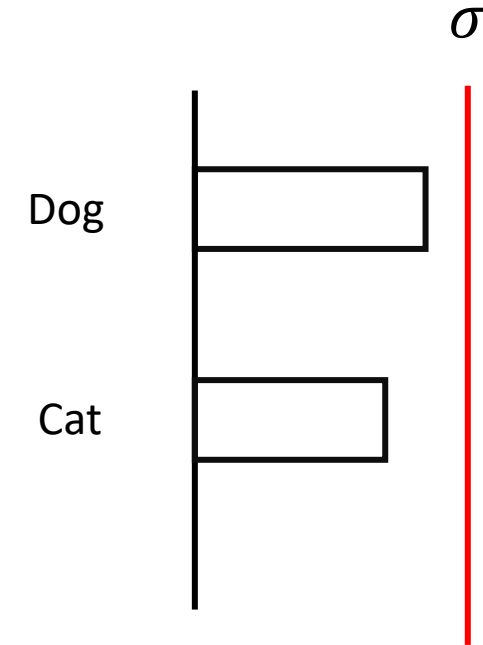
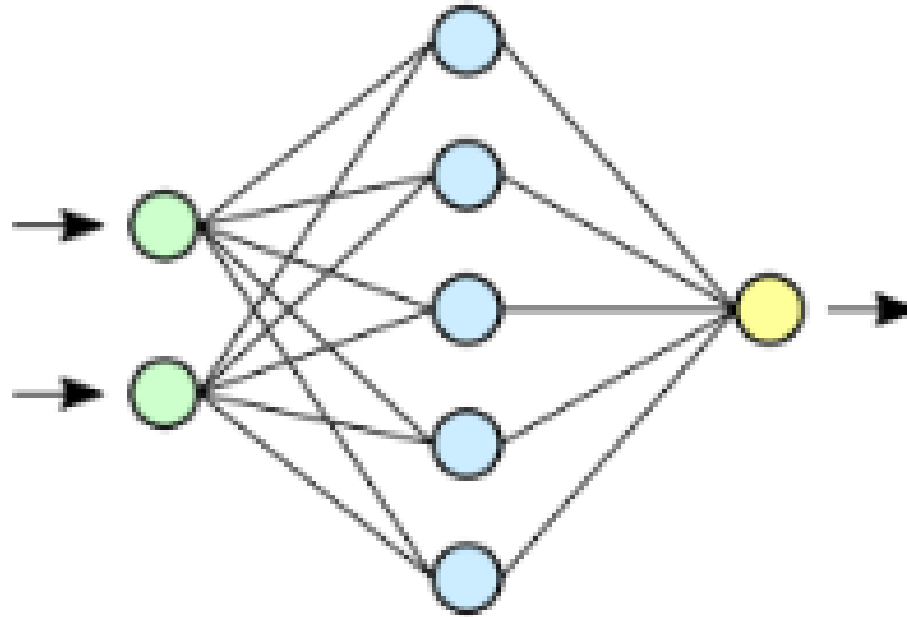
TAEWON KIM

1. DETECTING OUT-OF-DISTRIBUTION DETECTOR DATASET
2. CONTRIBUTION
3. CONFIDENCE LOSS
4. ADVERSARIAL GENERATOR FOR OUT-OF-DISTRIBUTION
5. EXPERIMENTAL

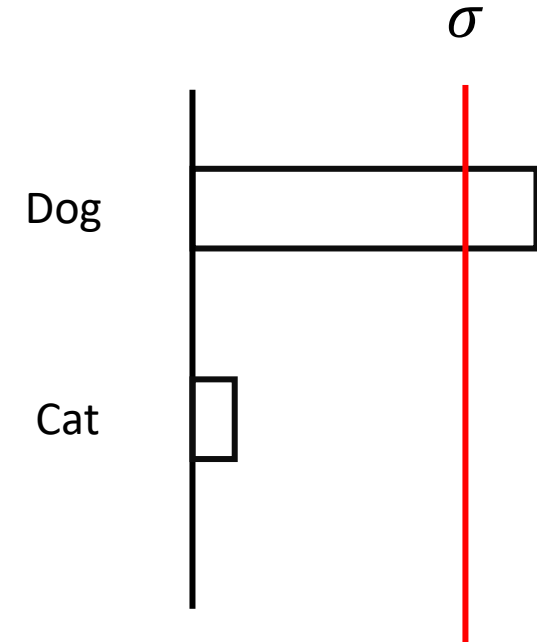
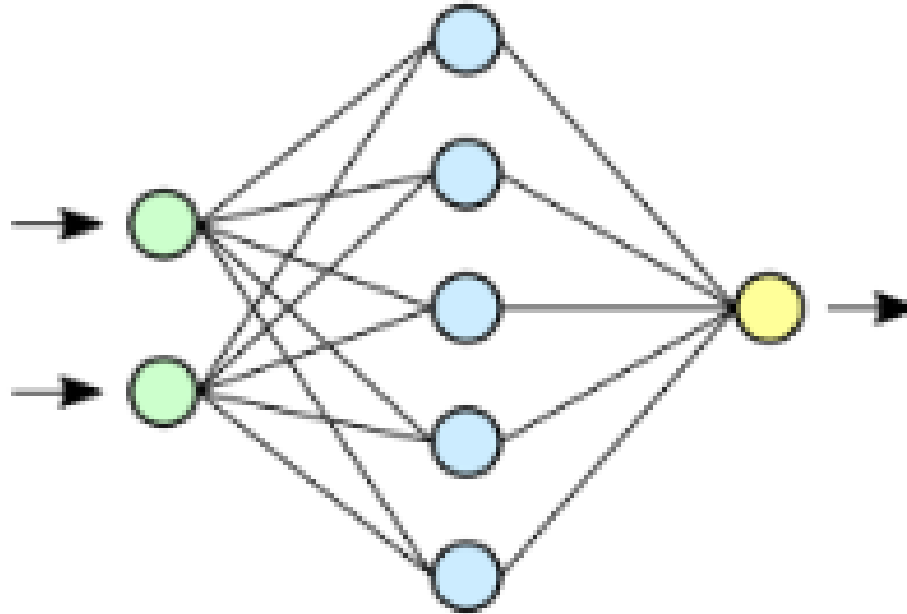
# 1. Detecting Out of distribution Detector



# 1. Detecting Out of distribution Detector



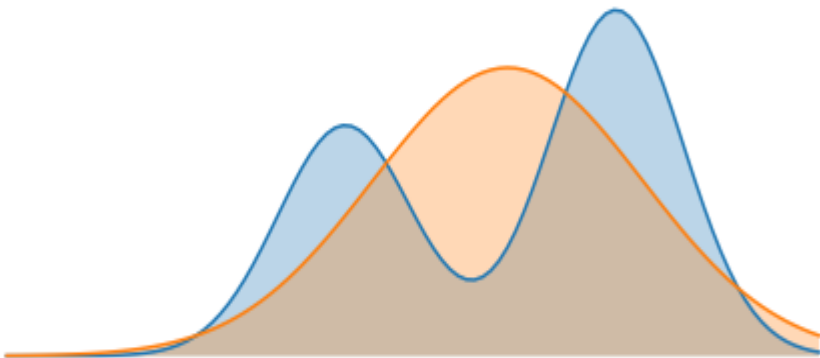
# 1. Detecting Out of distribution Detector



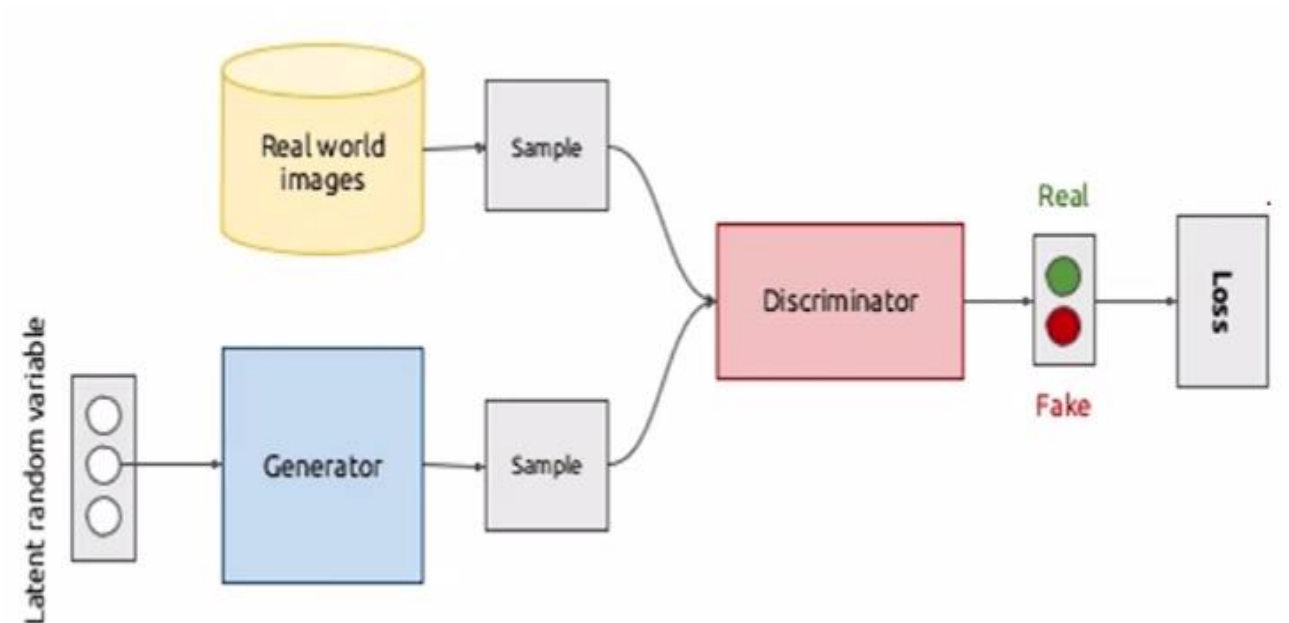
- 이전 Out of Distribution Network 는 pretrain된 모델을 기반으로 동작
- inference 단계만 사용하기 때문에 training 결과에 많이 의존

## 2. CONTRIBUTION

(KL) divergence



generative adversarial network (GAN)

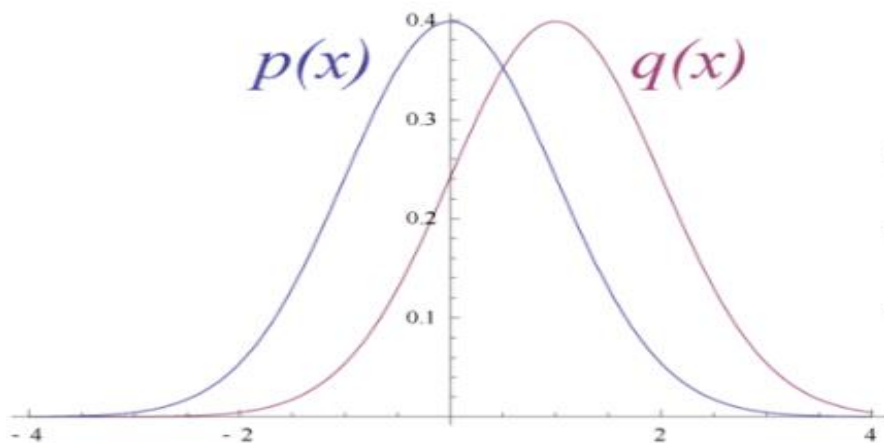


- KL-Div를 이용한 새로운 loss function을 제시
- GAN을 이용해 학습에 가장 효과적인 out-of-distribution Data를 생성

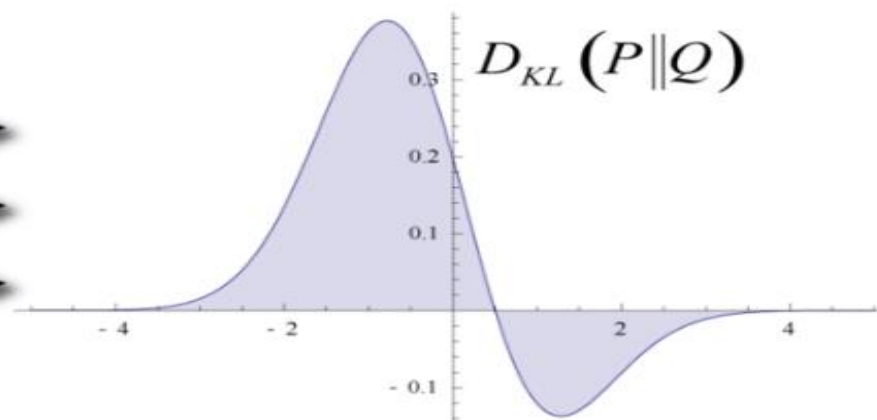
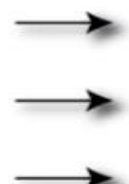
## 2. CONTRIBUTION

(KL) divergence

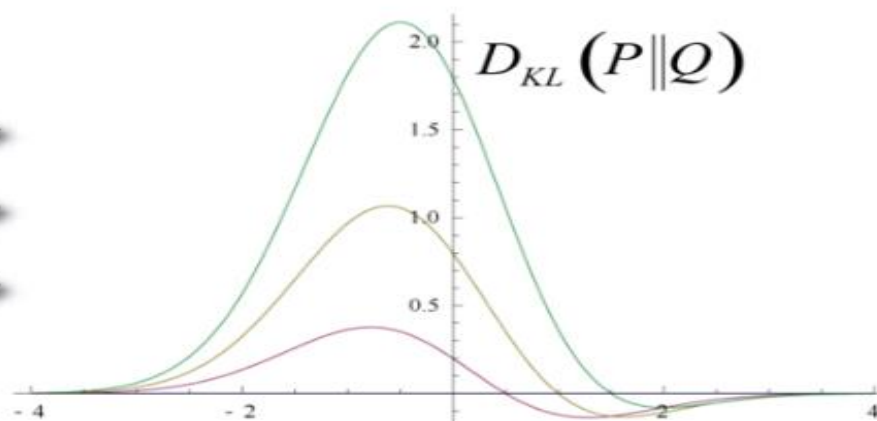
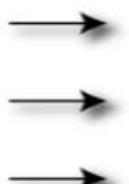
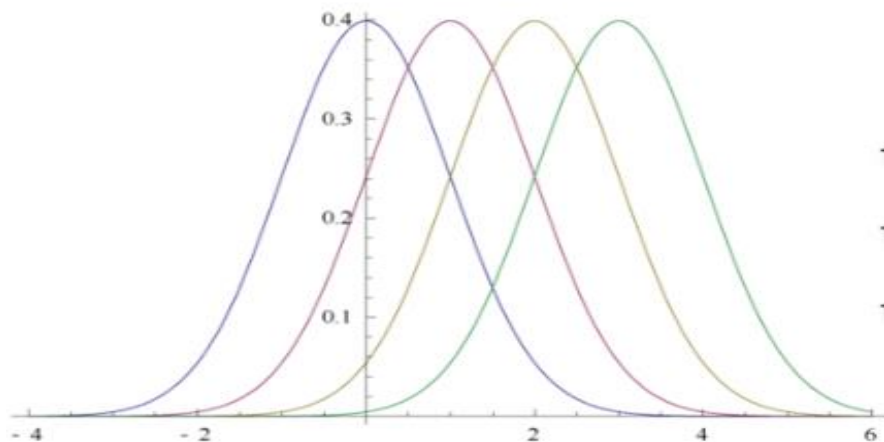
$$D_{KL}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right).$$



Original Gaussian PDF's



KL Area to be Integrated



## 2. CONFIDENCE LOSS

$$\min_{\theta} \mathbb{E}_{P_{\text{in}}(\hat{\mathbf{x}}, \hat{y})} \left[ -\log P_{\theta}(y = \hat{y} | \hat{\mathbf{x}}) \right] + \beta \mathbb{E}_{P_{\text{out}}(\mathbf{x})} \left[ KL(\mathcal{U}(y) \parallel P_{\theta}(y | \mathbf{x})) \right]$$

Cross Entropy loss

KL-Div. Term for Out of-dist. data

- Cross Entropy loss 같은 경우는 In-distribution data가 올바른 클래스를 예측하는 데 얼마나 잘 수행되는지
- KL-Div. Term 은 out-of-distribution data가 uniform한 분포를 가지도록 유도하는 과정
- 결론적으로는 In-dist 와 out-of-dist data가 모두 필요



## 2. CONFIDENCE LOSS

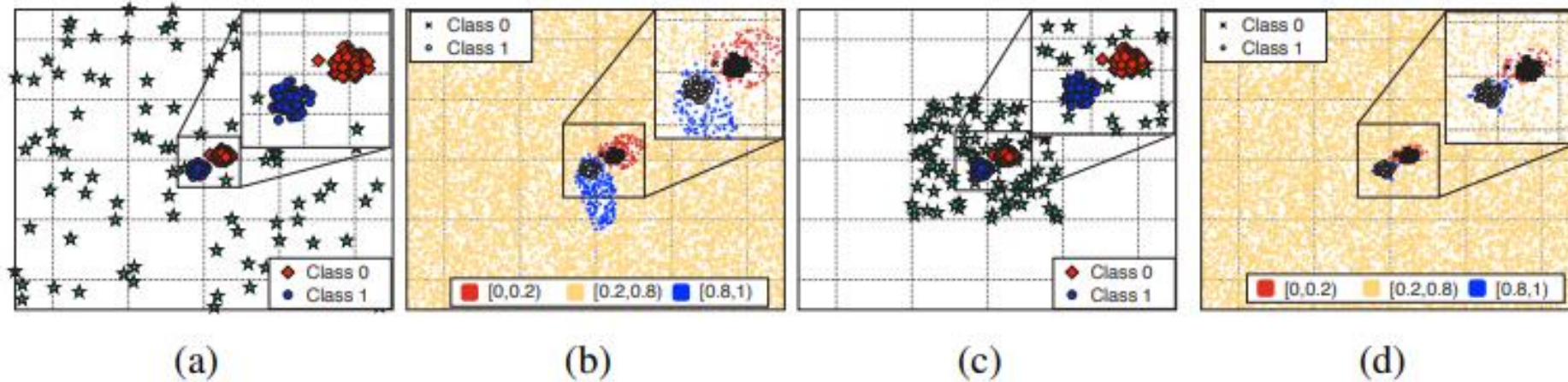
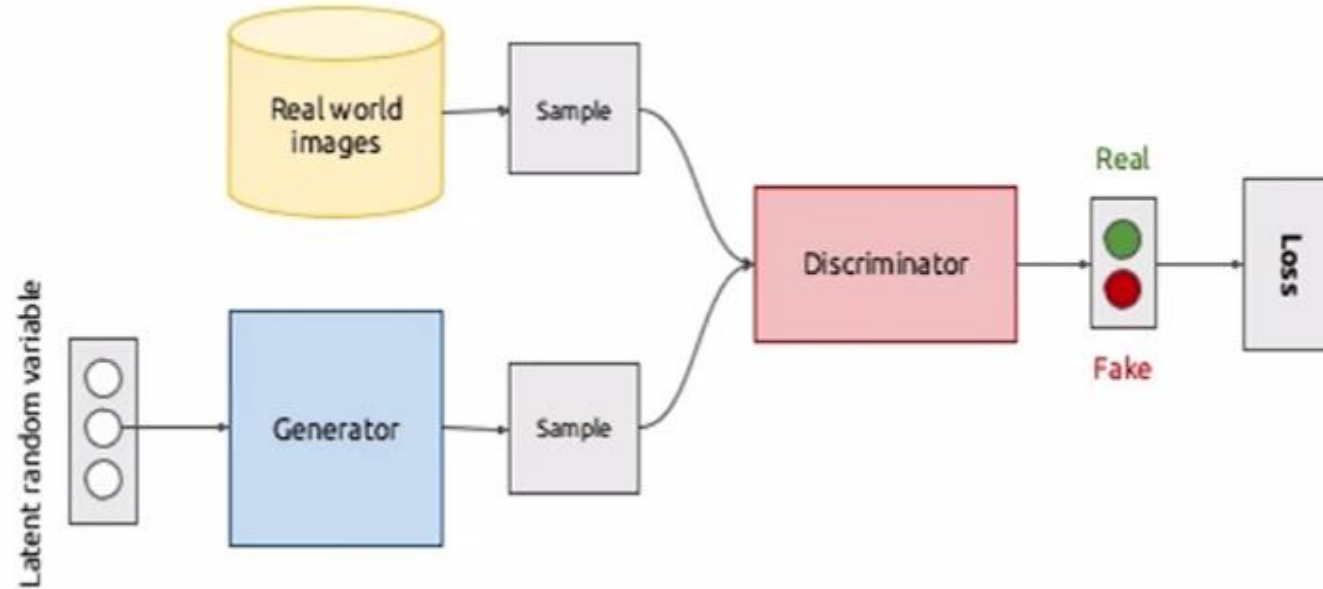


Figure 1: Illustrating the behavior of classifier under different out-of-distribution training datasets. We generate the out-of-distribution samples from (a) 2D box  $[-50, 50]^2$ , and show (b) the corresponding decision boundary of classifier. We also generate the out-of-distribution samples from (c) 2D box  $[-20, 20]^2$ , and show (d) the corresponding decision boundary of classifier.

- (a) Out-of-distribution에 대한 전체 데이터 분포를 구성
- (b) CONFIDENCE CLASSIFIERS 을 이용해 학습한 데이터 분포
- (c) In-distribution 경계에 있는 Out-of-distribution 데이터 분포를 구성

# 3. ADVERSARIAL GENERATOR

## GAN



Sample  $x$  from real data distribution

Sample latent code  $z$  from Gaussian distribution

$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)} [\log D(x)] + E_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

Maximum when  $D(x) = 1$

Maximum when  $D(G(z)) = 0$

$D$  should maximize  $V(D, G)$

# 3. ADVERSARIAL GENERATOR

## Modified GAN

$$\min_G \max_D \underbrace{\beta \mathbb{E}_{P_G(\mathbf{x})} [KL(\mathcal{U}(y) \parallel P_\theta(y|\mathbf{x}))]}_{(a)} + \underbrace{\mathbb{E}_{P_{in}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{P_G(\mathbf{x})} [\log(1 - D(\mathbf{x}))]}_{(b)},$$

- (a) data가 out of distribution 과 비슷하게 만드는 Generator loss
- (b) in-distribution data 대해 올바르게 분류하고 생성자가 생성한 가짜 데이터가 실제와 구분 안되게 학습

# 3. JOINT OBJECT FUNCTION

## JOINT TRAINING METHOD OF CONFIDENT CLASSIFIER AND ADVERSARIAL GENERATOR

$$\min_G \max_D \min_{\theta} \underbrace{\mathbb{E}_{P_{in}(\hat{\mathbf{x}}, \hat{y})} \left[ -\log P_{\theta}(y = \hat{y} | \hat{\mathbf{x}}) \right]}_{(c)} + \underbrace{\beta \mathbb{E}_{P_G(\mathbf{x})} \left[ KL(\mathcal{U}(y) \parallel P_{\theta}(y | \mathbf{x})) \right]}_{(d)} + \underbrace{\mathbb{E}_{P_{in}(\hat{\mathbf{x}})} \left[ \log D(\hat{\mathbf{x}}) \right] + \mathbb{E}_{P_G(\mathbf{x})} \left[ \log(1 - D(\mathbf{x})) \right]}_{(e)}.$$

Cross Entropy loss

KL-Div. Term for Out of-dist. data

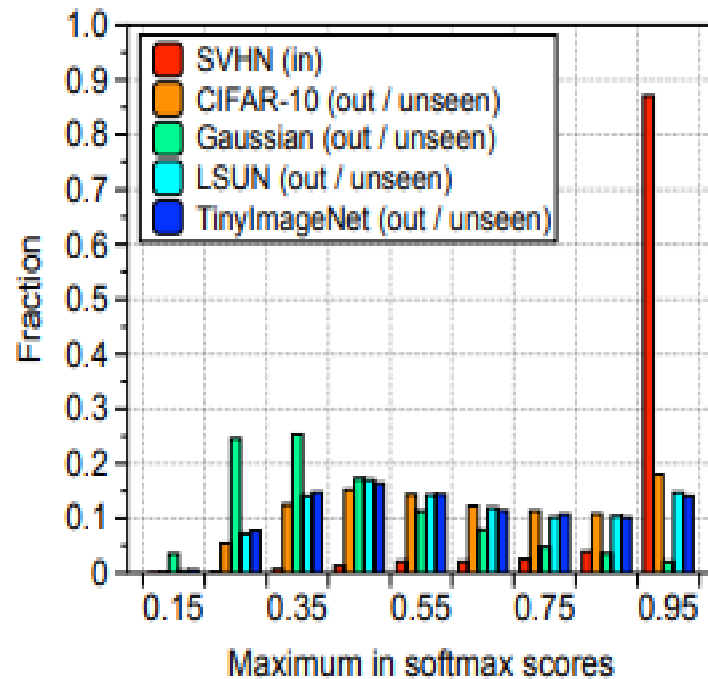
Modified GAN objective function

# 3. EXPERIMENTS

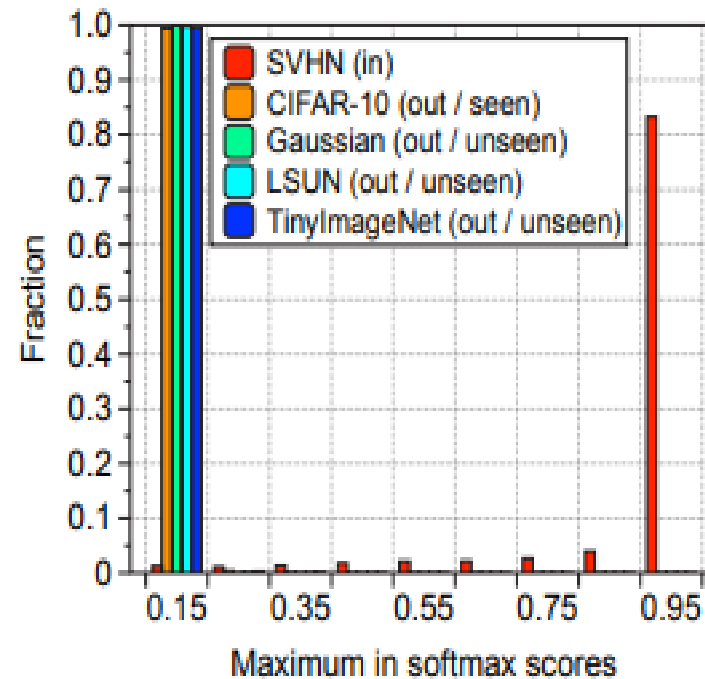
In-dist	Out-of-dist	Classification accuracy	TNR at TPR 95%	AUROC	Detection accuracy	AUPR in	AUPR out
Cross entropy loss / Confidence loss							
SVHN	CIFAR-10 (seen)	93.82 / <b>94.23</b>	47.4 / <b>99.9</b>	62.6 / <b>99.9</b>	78.6 / <b>99.9</b>	71.6 / <b>99.9</b>	91.2 / <b>99.4</b>
	TinyImageNet (unseen)		49.0 / <b>100.0</b>	64.6 / <b>100.0</b>	79.6 / <b>100.0</b>	72.7 / <b>100.0</b>	91.6 / <b>99.4</b>
	LSUN (unseen)		46.3 / <b>100.0</b>	61.8 / <b>100.0</b>	78.2 / <b>100.0</b>	71.1 / <b>100.0</b>	90.8 / <b>99.4</b>
	Gaussian (unseen)		56.1 / <b>100.0</b>	72.0 / <b>100.0</b>	83.4 / <b>100.0</b>	77.2 / <b>100.0</b>	92.8 / <b>99.4</b>
CIFAR-10	SVHN (seen)	80.14 / <b>80.56</b>	13.7 / <b>99.8</b>	46.6 / <b>99.9</b>	66.6 / <b>99.8</b>	61.4 / <b>99.9</b>	73.5 / <b>99.8</b>
	TinyImageNet (unseen)		<b>13.6</b> / 9.9	<b>39.6</b> / 31.8	<b>62.6</b> / 58.6	<b>58.3</b> / 55.3	<b>71.0</b> / 66.1
	LSUN (unseen)		<b>14.0</b> / 10.5	<b>40.7</b> / 34.8	<b>63.2</b> / 60.2	<b>58.7</b> / 56.4	<b>71.5</b> / 68.0
	Gaussian (unseen)		2.8 / <b>3.3</b>	10.2 / <b>14.1</b>	50.0 / 50.0	48.1 / <b>49.4</b>	39.9 / <b>47.0</b>

Table 1: Performance of the baseline detector (Hendrycks & Gimpel, 2016) using VGGNet. All values are percentages and boldface values indicate relative the better results. For each in-distribution, we minimize the KL divergence term in (1) using training samples from an out-of-distribution dataset denoted by “seen”, where other “unseen” out-of-distributions were only used for testing.

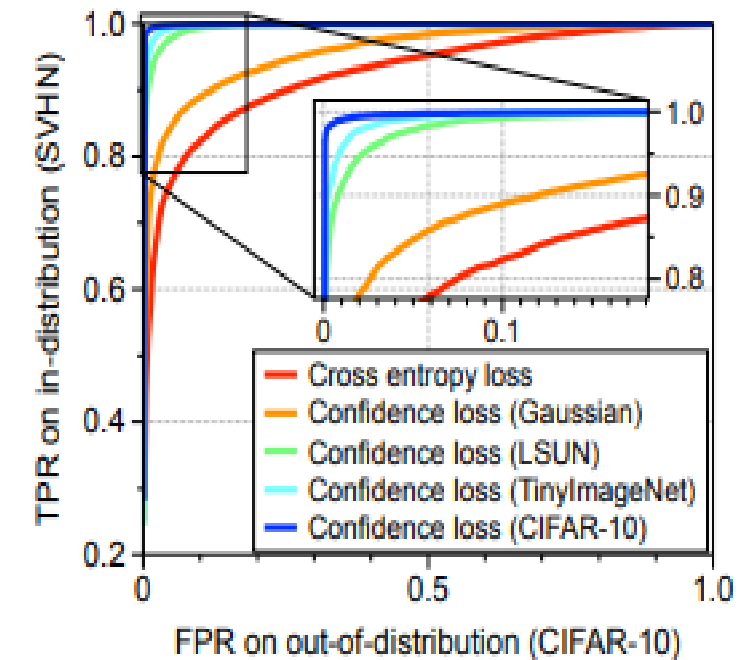
# 3. EXPERIMENTS



(a) Cross entropy loss



(b) Confidence loss in (1)



(c) ROC curve



# 3. EXPERIMENTS

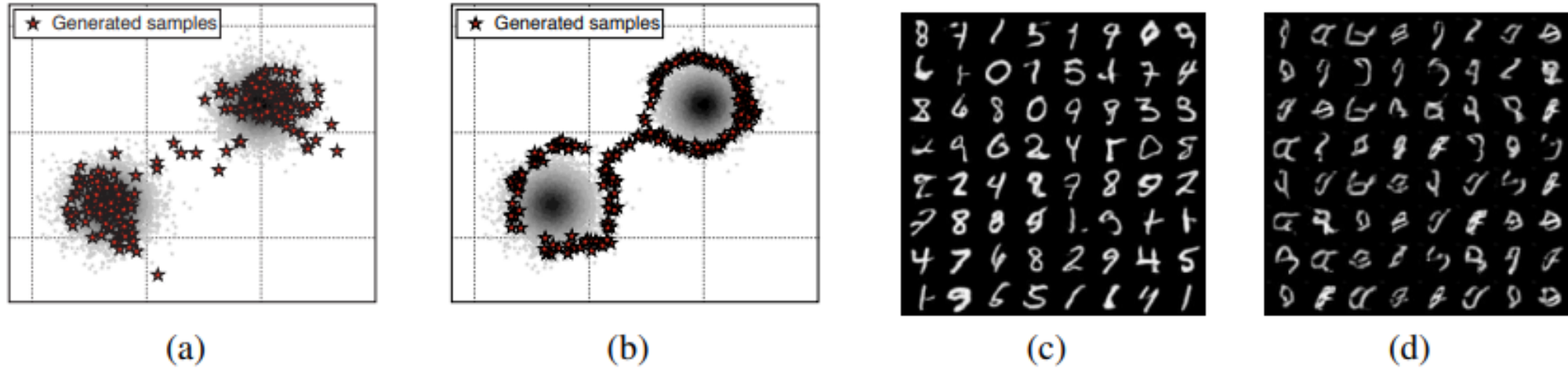
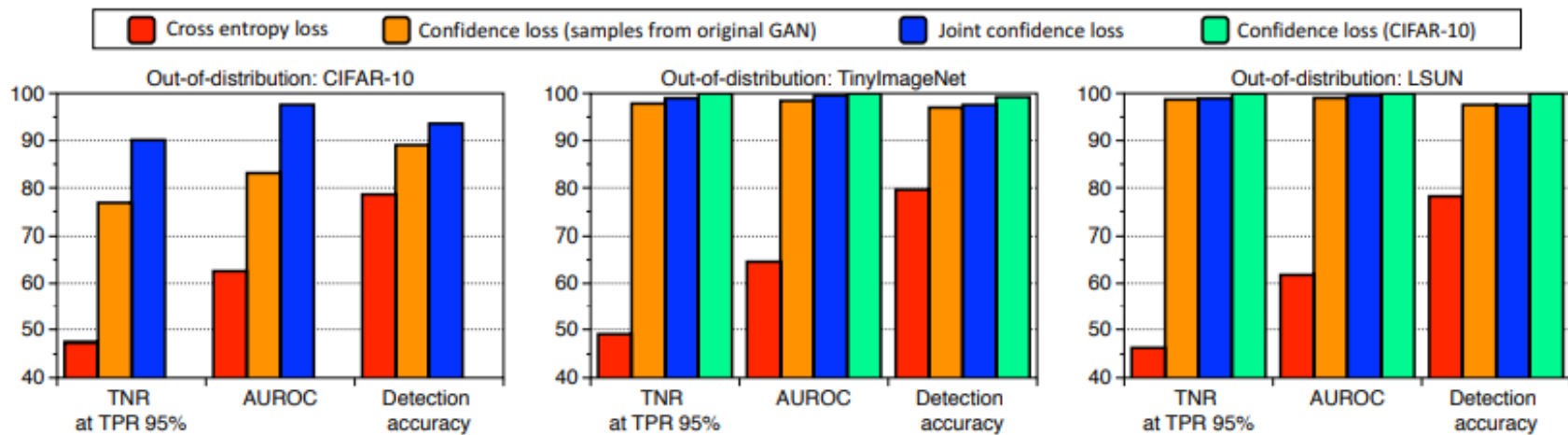
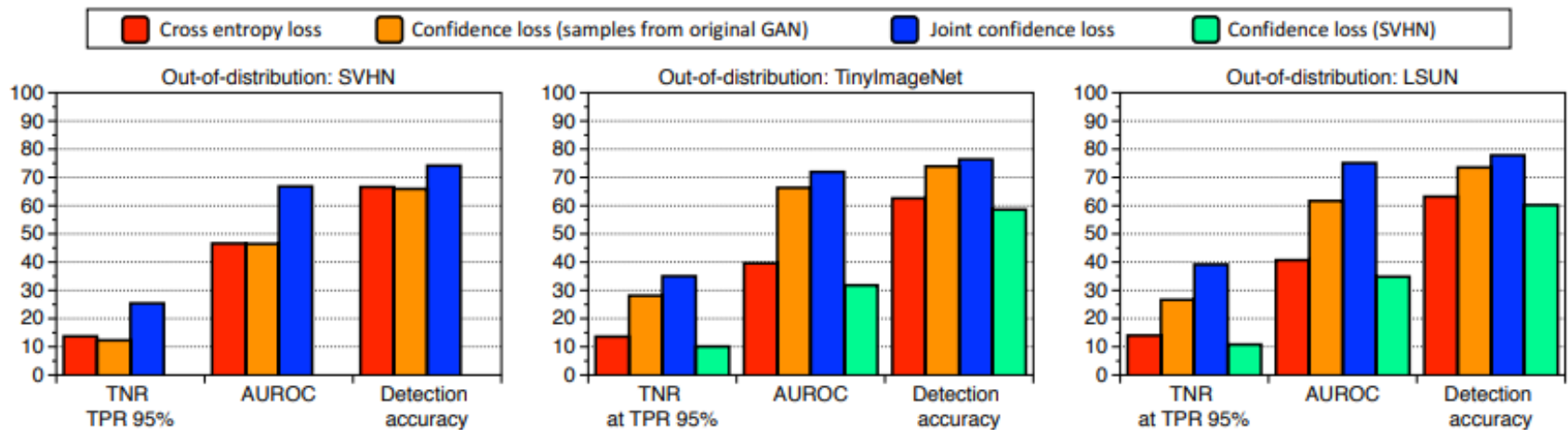


Figure 3: The generated samples from original GAN (a)/(c) and proposed GAN (b)/(d). In (a)/(b), the grey area is the 2D histogram of training in-distribution samples drawn from a mixture of two Gaussian distributions and red points indicate generated samples by GANs.

# 3. EXPERIMENTS



(a) In-distribution: SVHN



(b) In-distribution: CIFAR-10