

# GAN-BERT

accepted by ACL

Danilo Croce, Giuseppe Castelletti, Roberto Basili

채형주

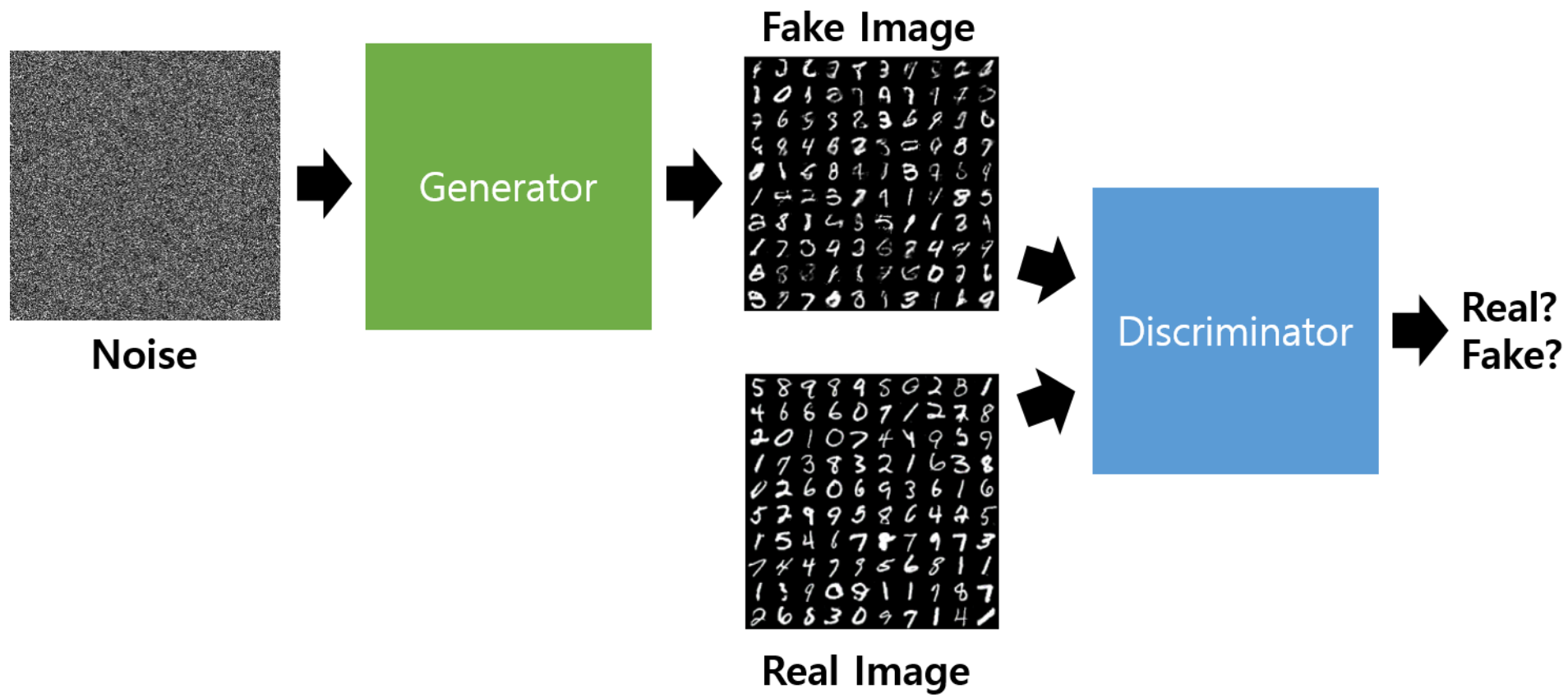
# 목차

- 제안 배경
- GAN
- SS-GAN
- GAN-BERT
- Result
- Conclusion

# 제안 배경

- 항상 annotated data는 부족하고, unlabeled data를 직접 labeling하기에는 상당한 비용이 들어간다.
- BERT는 좋은 성능을 내지만 fine-tuning을 위한 많은 양의 labeled data 필요
- BERT를 fine-tuning 할때 200개 이하의 annotated data instances로 학습시킨다면 성능의 상당한 저하가 온다.

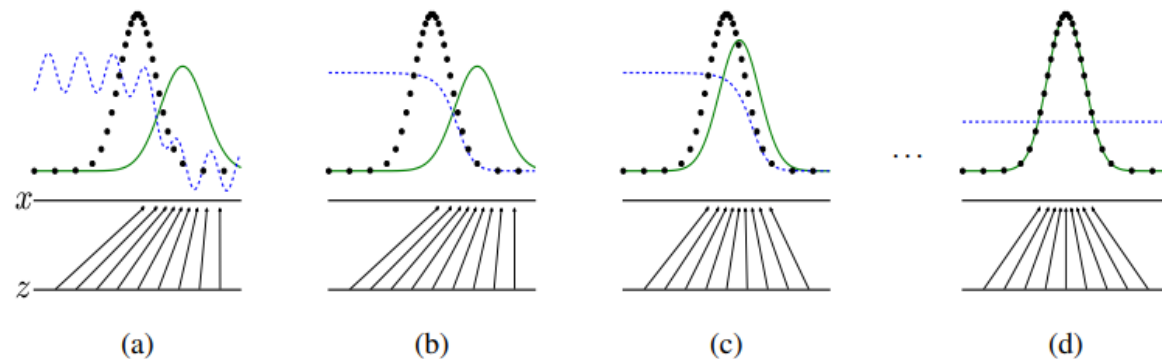
# GAN



# GAN : to find Nash Equilibrium

	Real Data	Fake Data
$D == 1$ (진짜라고 판단)		Discriminator update
$D == 0$ (가짜라고 판단)	Discriminator update	Generator update

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \quad (1)$$



## 4.1 Global Optimality of $p_g = p_{\text{data}}$

We first consider the optimal discriminator  $D$  for any given generator  $G$ .

**Proposition 1.** For  $G$  fixed, the optimal discriminator  $D$  is

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$$

# SS-GAN : semi supervised GAN

- 기존 GAN은 내쉬 균형으로 수렴이 잘 안되는 문제가 있음 (non-convex, large parameters)

- feature matching

$$||\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \mathbf{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} \mathbf{f}(G(\mathbf{z}))||_2^2$$

- semi-supervised learning

$$\begin{aligned} L &= -\mathbb{E}_{\mathbf{x}, y \sim p_{\text{data}}(\mathbf{x}, y)} [\log p_{\text{model}}(y|\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim G} [\log p_{\text{model}}(y = K + 1|\mathbf{x})] \\ &= L_{\text{supervised}} + L_{\text{unsupervised}}, \text{ where} \end{aligned}$$

$$L_{\text{supervised}} = -\mathbb{E}_{\mathbf{x}, y \sim p_{\text{data}}(\mathbf{x}, y)} \log p_{\text{model}}(y|\mathbf{x}, y < K + 1)$$

$$L_{\text{unsupervised}} = -\{\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \log[1 - p_{\text{model}}(y = K + 1|\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim G} \log[p_{\text{model}}(y = K + 1|\mathbf{x})]\},$$

# SS-GAN : semi supervised GAN



Figure 3: (*Left*) samples generated by model during semi-supervised training. Samples can be clearly distinguished from images coming from MNIST dataset. (*Right*) Samples generated with minibatch discrimination. Samples are completely indistinguishable from dataset images.

Model	Number of incorrectly predicted test examples for a given number of labeled samples			
	20	50	100	200
DGN [21]			333 $\pm$ 14	
Virtual Adversarial [22]			212	
CatGAN [14]			191 $\pm$ 10	
Skip Deep Generative Model [23]			132 $\pm$ 7	
Ladder network [24]			106 $\pm$ 37	
Auxiliary Deep Generative Model [23]			96 $\pm$ 2	
Our model	1677 $\pm$ 452	221 $\pm$ 136	93 $\pm$ 6.5	90 $\pm$ 4.2
Ensemble of 10 of our models	1134 $\pm$ 445	142 $\pm$ 96	86 $\pm$ 5.6	81 $\pm$ 4.3

# GAN-BERT

- sentence classification task
- Generator
  - input : 100-d noise vector
  - MLP
  - output : h\_fake (768-d)
- Discriminator
  - input : h\* = hCLS
  - MLP
  - output : k+1 vector of logit(trough softmax layer)
- BERT is updated when D is updated

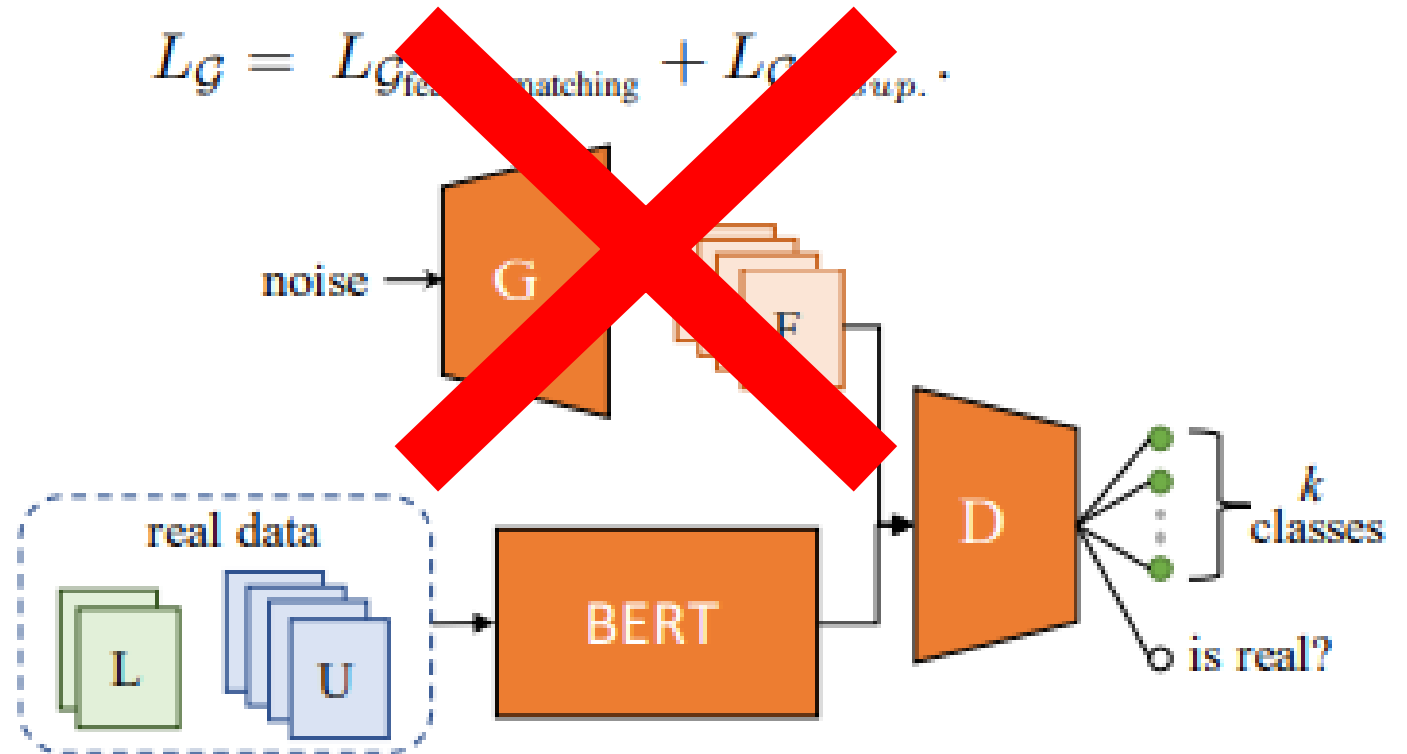
$$L_{D_{\text{supervised}}} = -\mathbb{E}_{x, y \sim p_{\text{data}}} \log [p_{\text{model}}(y = i|x, i < k + 1)]$$

$$L_{D_{\text{unsupervised}}} = -\mathbb{E}_{x \sim p_{\text{data}}} \log [1 - p_{\text{model}}(y = k + 1|x)] - \mathbb{E}_{x \sim G} \log [p_{\text{model}}(y = k + 1|x)]$$

$$L_{G_{\text{feature matching}}} = \|\mathbb{E}_{x \sim p_d} f(x) - \mathbb{E}_{x \sim G} f(x)\|_2^2$$

$$L_{G_{\text{unsup.}}} = -\mathbb{E}_{x \sim G} \log [1 - p_m(\hat{y} = y|x, y = k + 1)]$$

$$L_G = L_{G_{\text{feature matching}}} + L_{G_{\text{unsup.}}}$$





# GAN-BERT

- model specification

noise vector : from  $N(0, 1)$

$h$  : 768-d

activation function  $G, D$  : leaky relu

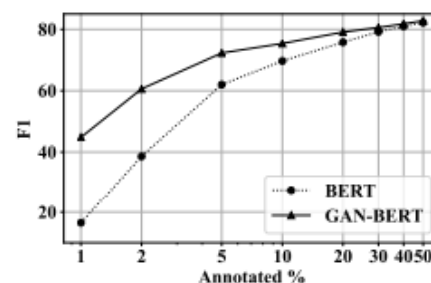
dropout : 0.1

annotated data 를 0.1% 또는 1%부터 사용하기 시작하고, 점점 늘려가면서 BERT와 차이를 비교

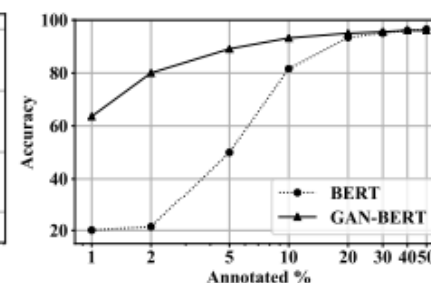
가능하면 사용한 annotated data의 100배에 해당하는 unlabeled data를 제공

# Results

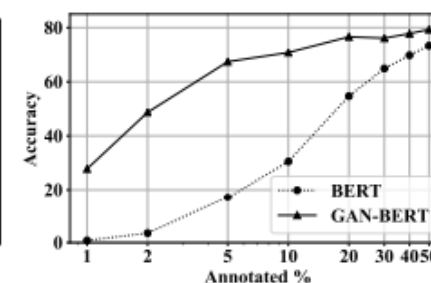
- Dataset
  - 20 NEWS Group(20N)
  - Question Classification(QC) – [fine grained / coarse grained]
  - Sentiment Analysis(SST-5)
  - MNLI



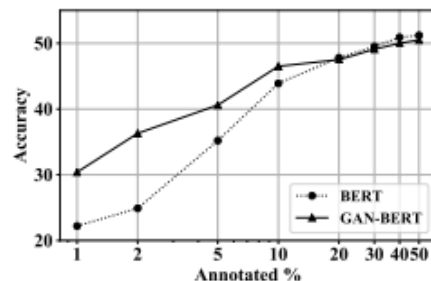
(a) 20N



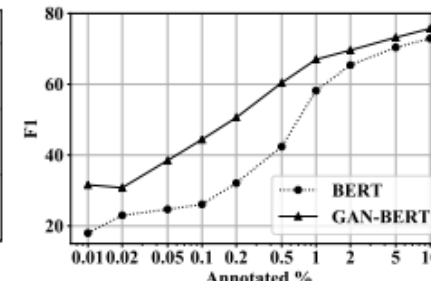
(b) QC Coarse Grained



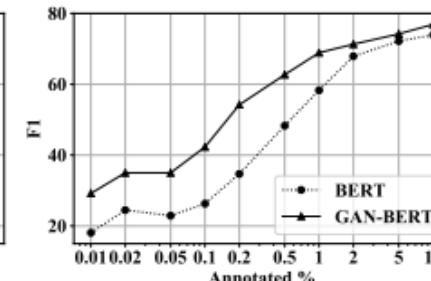
(c) QC Fine Grained



(d) SST-5



(e) MNLI Matched



(f) MNLI Mismatched

# Conclusion

- 안좋은 훈련환경에서 BERT를 능가하였다!
- labeled data개수가 적을 때에 GAN을 이용하여 성능을 비약적으로 향상시켰다.
- GPT-2, DistilBERT같은 다른 모델에도 적용가능하고, 다양한 task에도 적용할 수 있는 가능성을 열어주었다.
- BERT의 pre-training단계에도 직접적으로 적용해보고 싶다.