# Counter-Contrastive Learning for Language GANs

Yekun Chai♡♠, Haidong Zhang♠, Qiyue Yin♠, Junge Zhang♠

chaiyekun@gmail.com    haidong_zhang14@yahoo.com    {qyyin, jgzhang}@nlpr.ia.ac.cn

♠Institute of Automation, Chinese Academy of Sciences    ♡Baidu NLP

## Abstract

Generative Adversarial Networks (GANs) have achieved great success in image synthesis, but have proven to be difficult to generate natural language. Challenges arise from the uninformative learning signals passed from the discriminator. In other words, the poor learning signals limit the learning capacity for generating languages with rich structures and semantics. In this paper, we propose to adopt the counter-contrastive learning (CCL) method to support the generator's training in language GANs. In contrast to standard GANs that adopt a simple binary classifier to discriminate whether a sample is real or fake, we employ a counter-contrastive learning signal that advances the training of language synthesizers by (1) pulling the language representations of generated and real samples together and (2) pushing apart representations of real samples to compete with the discriminator and thus prevent the discriminator from being overtrained. We evaluate our method on both synthetic and real benchmarks and yield competitive performance compared to previous language GANs.

## Introduction

- Generative Adversarial Networks (GANs) hold the promise of training language models, as an alternative method to MLE. GANs learn to sample during training so as to avoid the exposure bias issue, whose aim is to train a language generator to fool the discriminator that distinguishes the fake data out of real samples.

- Previous innovations adopt various approaches to enhance the learning signals for generators, such as leaking information from the discriminator to the generator [3], directly matching the fake data distribution to that of real data [6, 1], learning to rank samples out of a collection of curated samples [4, 7], leveraging more powerful generator architectures to learning representations [5], *etc*. However, the problem of language GANs' training is far from being fully solved.

- Inspired by the recent success in contrastive learning approaches [2] in learning effective representations, we propose a counter-contrastive learning objective to aid the adversarial learning of sequence generators in language GANs. Conventional contrastive learning methods aim at pulling positive samples together and pushing away positive samples from negative ones.

## Intuition: Counter-Contrastive Learning

▷ **Contrastive Learning**

- Help $\mathcal{D}$ to discriminate positive samples from negative ones.

$$\mathcal{L}_i^{\text{CL}} = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^{N} e^{\text{sim}(\mathbf{h}_j, \mathbf{h}_j^+)/\tau}} \qquad (1)$$

- However, the generator $\mathcal{G}$ in language GANs aims to **cheat the discriminator** $\mathcal{D}$.

▷ **Counter-Contrastive Learning**

- Draw together the fake and real samples $(x_i, x_i^-)$ (to let the generator imitate the real sentences);

- Push away the real samples $(x_i, x_i^+)$ (to fool and hinder the discriminator training, thereby preventing it from fast convergence).

## Methodology

- **Positive Samples**. We construct positive pairs by applying disparate dropout masks to get positive representations for input real texts sampled from $p_{\text{data}}$. Specifically, for the same real sentence, we get positive pair representations after feed them into the discriminator twice with two different random dropout operations. Denote $\mathbf{h}_i^m = f(x_i, m)$, where $m$ is the dropout mask and $f$ is the encoder of input sentences.

- **Negative Samples**. We randomly select fake sentences generated by the generator network and feed them into the discriminator to get fake sample representations. Therefore, we choose one from positive representations and the other from the negative to construct negative pairs $(\mathbf{h}_i, \mathbf{h}_i^-)$.

- **Counter-Contrastive Learning**. Given the mini-batch of size $N$, we formulate the counter-contrastive learning objectives as:

$$\mathcal{L}_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^-)/\tau}}{\sum_{j=1}^{N} \left( e^{\text{sim}(\mathbf{h}_j, \mathbf{h}_j^-)/\tau} + e^{\text{sim}(\mathbf{h}_j, \mathbf{h}_j^+)/\tau} \right)}$$

where $\tau$ is the constant temperature.

Intuitively, this CCL objective aims to (1) force the fake representations to approach real data (the numerator), and (2) prevent the discriminator from learning effective representations of positive pairs by pushing away semantically close pairs (the right term in the denominator).

## Results

**Synthetic Data**  We evaluate the generated sequence w.r.t. both quality and diversity. It is observed that our model outperforms baseline models in terms of quality (measured by $\text{NLL}_{\text{oracle}}$) and quality-diversity trade-off (measured by $\text{NLL}_{\text{oracle}}+\text{NLL}_{\text{gen}}$), and achieves or matches the competitive results of baselines w.r.t. the diversity (indicated by $\text{NLL}_{\text{gen}}$).

| Model | $\text{NLL}_{\text{oracle}}$ (20/40) | $\text{NLL}_{\text{gen}}$ (20/40) | $\text{NLL}_{\text{oracle}} + \text{NLL}_{\text{gen}}$ (20/40) |
|---|---|---|---|
| MLE | 9.05±0.03 / 9.84±0.02 | 5.96±0.02 / 6.55±0.02 | 15.02±0.03 / 16.39±0.01 |
| SeqGAN | 8.63±0.19 / 9.63±0.04 | 6.61±0.22 / 6.98±0.08 | 15.00±0.03 / 16.35±0.02 |
| RankGAN | 8.42±0.31 / 9.52±0.11 | 7.14±0.34 / 7.05±0.12 | 15.01±0.02 / 16.37±0.02 |
| MaliGAN | 8.74±0.16 / 9.67±0.03 | 6.62±0.25 / 7.14±0.09 | 15.03±0.03 / 16.39±0.03 |
| SAL | 7.71±0.17 / 9.31±0.03 | 6.58±0.15 / 6.97±0.05 | 14.29±0.11 / 16.24±0.03 |
| Ours | **6.77**±0.34 / **6.65**±0.14 | 6.91±0.62 / 7.68±0.79 | **13.69**±0.36 / **14.33**±0.76 |

**Real Data**  Our model shows a significant improvement over previous methods, consistently achieves competitive results in terms of the sample quality (indicated by BLEU scores) while maintaining the diversity (indicated by $\text{NLL}_{\text{gen}}$).

| Model | BLEU-2 | BLEU-3 | BLEU-4 | BLEU-5 | $\text{NLL}_{\text{gen}}$ |
|---|---|---|---|---|---|
| MLE | 0.731 | 0.497 | 0.305 | 0.189 | 0.718 |
| SeqGAN | 0.745 | 0.498 | 0.294 | 0.180 | 1.082 |
| RankGAN | 0.743 | 0.467 | 0.264 | 0.156 | 1.344 |
| LeakGAN | 0.746 | 0.528 | 0.355 | 0.230 | 0.679 |
| RelGAN | 0.849±0.030 | 0.687±0.047 | 0.502±0.048 | 0.331±0.044 | 0.756±0.054 |
| SAL | 0.785±0.02 | 0.581±0.03 | 0.362±0.02 | 0.227±0.02 | 0.873±0.02 |
| Ours (CCL) | **0.871**±0.032 | **0.715**±0.050 | **0.538**±0.068 | **0.399**±0.082 | **0.630**±0.103 |

## Comparison with Language GANs without CCL

| model | Sample sentences |
|---|---|
| w/o CCL | a cat is sitting on a white plate . |
| | a cat is sitting on a bathroom sink sitting inside of a toilet . |
| | a black and white cat outside decorated in rustic kitchen . |
| | a cat is sitting on a bathroom sink sitting in a bathroom . |
| | a cat is sitting on a bathroom sink sitting on a bathroom counter . |
| | a cat sitting on a gravel ground inside of a bathroom sink . |
| | a cat is sitting on a bathroom sink sitting in a bathroom . |
| w/ CCL | a cat is sitting on top of a car . |
| | a cat is sitting on top of a car cleaning itself . |
| | a cat is sitting on top of a car roof . |
| | a cat is sitting on top of a car hood . |
| | a cat is sitting on top of a man 's head in front of a glass door . |
| | a dog sitting on top of a parked car near a cat . |
| | a cat in a white bathroom with a toilet paper beside a child . |

Fig. 3: Comparison between generated sentences from language GANs with and without counter-contrastive learning.

**Better Diversity**  For fair comparison, we select the generated sentences that contain the word "cat" from samples produced by models with and without the CCL method. It is observed that GANs with CCL tend to produce sentences with better diversity. For example, with the structure "a cat is sitting on top of a car", models w/ CCL can enrich it with different modifier words. However, after removing CCL, the model can duplicate words such as "sitting" regardless of its repetitive usage. Moreover, as shown in Fig.3, with the CCL method, language GANs tend to write semantically meaningful samples in comparison with the counterpart without CCL.

## Conclusion

In this paper, we introduce a counter-contrastive learning objective to advance the training of language GANs. It pulls the representation of generated and real samples together to promote the generator training, and pushes apart real sample pairs to depress the discriminator training as a competitor. Our work aims to integrate the prevalent contrastive learning approach in supporting the generator training, which lies in the line of methods using comparative signals or ranking classifiers, such as RankGAN and SAL. From the perspective of feature matching, the counter-contrastive learning objective can be considered as a contrastive signal to draw together the fake and real sample representations.

## References

[1] Liqun Chen et al. "Adversarial Text Generation via Feature-Mover's Distance". In: *NeurIPS*. 2018.

[2] Ting Chen et al. "A Simple Framework for Contrastive Learning of Visual Representations". In: *ArXiv* abs/2002.05709 (2020).

[3] Jiaxian Guo et al. "Long Text Generation via Adversarial Training with Leaked Information". In: *AAAI*. 2017.

[4] Kevin Lin et al. "Adversarial Ranking for Language Generation". In: *NIPS*. 2017.

[5] Weili Nie, Nina Narodytska, and Ankit B. Patel. "RelGAN: Relational Generative Adversarial Networks for Text Generation". In: *ICLR*. 2019.

[6] Yizhe Zhang et al. "Adversarial Feature Matching for Text Generation". In: *ICML*. 2017.

[7] Wangchunshu Zhou et al. "Self-Adversarial Learning with Comparative Discrimination for Text Generation". In: *ICLR*. 2020.