# Dialogue response generation with Wasserstein generative adversarial networks

Shuaib Ahmed Syed Gilani[0000-1111-2222-3333] and Edgar Jembere and  Anban W. Pillay
[1111-2222-3333-4444]

[1] Centre for Artificial Intelligence Research
[2] University of KwaZulu-Natal

**Abstract.** This research evaluates the effectiveness of a Generative Adversarial Network (GAN) for open domain dialogue response systems. The research involves developing and evaluating a Conditional Wasserstein GAN (CWGAN) for natural dialogue response generation. We begin by exploring the latest research in GANs for text generation, build a dataset for experimentation, develop our GAN architecture and model and then evaluate the results it produced.
**Keywords:** Machine Learning, Generative Adversarial Networks, Dialogue Systems, Recurrent Neural Networks.

Dialogue response systems are systems that converse and exchange information with the user. Until very recently dialogue systems were primarily built on closed-domain retrieval-based systems. Closed domain dialogue systems are restricted to a specific topic of discussion or domain of knowledge, while open domain dialogue systems do not have this restriction. In closed domain dialogue systems there are correct responses to input which can be used to measure the systems performance. However, for open domain dialogue systems there may not be a perfect response to an input. In our work we look at the open domain which has implications for the data and evaluation we use.

Machine learning approaches could not compete with retrieval-based systems until recently when sequence-to-sequence (Seq2Seq) models were introduced by Sutskever [3][17]. Seq2Seq models gained much attention in academia as well as in industrial use due to its simpler architecture while producing promising results in both closed domain and open domain dialogue systems.  However, it has been demonstrated that Seq2Seq models with the Maximum likelihood estimation (MLE) objective function suffer from exposure bias and lacks diversity of responses. Exposure bias occurs when a model is trained on the ground truth data and not exposed to its own errors, causing small errors to accumulate over time which degrades the model performance.

Generative Adversarial Networks (GANs) were invented in 2014 by Ian Goodfellow and have since produced noteworthy results in generating images [4]. GANs typically consist of 2 networks, a generator and a discriminator, that compete in a zero-sum game. The generator learns to generate text sequence responses to an input text se-

quence. The discriminator must learn to distinguish between real response sequences and generated response sequences; while the generator aims to fool the discriminator into predicting that the generated response sequence is a real response sequence. This dynamic suggests that training a GAN to generate sequences may not be prone to the same problems as MLE based models [1][2][8][16]. The generator's objective is to generate responses that fool the discriminator rather than being trained on the ground truth alone.

Training GANs on text is not trivial because discrete data is non-differentiable and an unstable training objective. To solve this problem many approaches have been taken in literature. Some of the approaches have used a policy gradient approach such as SeqGAN [7], REGS [9], RankGAN [13] and StepGAN [16]. Other approaches look to modifying the objective function such as MaliGAN [8], Gumbel-softmax [11], Soft-argmax [15] and Wasserstein distance [1][2][6].

The Wasserstein or earth mover distance with gradient penalty has shown to produce higher quality text language models [1][2][6]. The Wasserstein distance measures the work done in moving the generated response distribution towards the real response distribution. It forces the discrete data to work in a continuous field so that it is differentiable for back propagation. It makes training easier by providing a softer metric to compare the distributions [1]. We apply this objective function as well as the conditional GAN structure such that the generated output sequence is conditioned on some prior input sequence. The GAN functions as an end-to-end text-based dialogue response system.

Recurrent Neural Networks (RNN) are well suited for learning sequences because it can retain information through training time steps [1][2][3][9][14]. Thus, our system uses RNNs in the generator and discriminator models. To aid the generator in training, we employ some techniques from literature such as gradually increasing the length of generated sequences, varying the length during training and conditioning on shorter ground truth sequences [2].

In our work we evaluate the efficacy of a conditional Wasserstein GAN for dialogue response generation. A large dataset of reddit comments and replies was processed this data one million comments and replies (input and response) sequences. Smaller subsets of the dataset are taken and split into training, test and validation sets.

Evaluation of open domain dialogue systems is difficult because there are many acceptable responses to an input. We have evaluated using the Bilingual evaluation understudy (BLEU) score which compares n-grams between the generated and target response. A random example from the test set is given the comment "Droste effect, Google it." The target response is "Ok" while the generated response is "Thank you!". The BLEU score is poor because there are no matching n-grams while both responses may be deemed acceptable to human evaluators. The approach is promising and we expect good results.

# References

1. Rajeswar, Sai, et al. "Adversarial generation of natural language." *arXiv preprint arXiv:1705.10929* (2017).
2. Press, Ofir, et al. "Language generation with recurrent generative adversarial networks without pre-training." *arXiv preprint arXiv:1706.01399* (2017).
3. Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. "Sequence to sequence learning with neural networks." *Advances in neural information processing systems*. 2014.
4. Goodfellow, Ian, et al. "Generative adversarial nets." *Advances in neural information processing systems*. 2014.
5. Liu, Chia-Wei, et al. "How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation." *arXiv preprint arXiv:1603.08023* (2016).
6. Gulrajani, Ishaan, et al. "Improved training of wasserstein gans." *Advances in neural information processing systems*. 2017.
7. Yu, Lantao, et al. "Seqgan: Sequence generative adversarial nets with policy gradient." *Thirty-First AAAI Conference on Artificial Intelligence*. 2017.
8. Che, Tong, et al. "Maximum-likelihood augmented discrete generative adversarial networks." *arXiv preprint arXiv:1702.07983* (2017).
9. Li, Jiwei, et al. "Adversarial learning for neural dialogue generation." *arXiv preprint arXiv:1701.06547* (2017).
10. Fedus, William, Ian Goodfellow, and Andrew M. Dai. "MaskGAN: better text generation via filling in the_." *arXiv preprint arXiv:1801.07736* (2018).
11. Kusner, Matt J., and José Miguel Hernández-Lobato. "Gans for sequences of discrete elements with the gumbel-softmax distribution." *arXiv preprint arXiv:1611.04051* (2016).
12. Mirza, Mehdi, and Simon Osindero. "Conditional generative adversarial nets." *arXiv preprint arXiv:1411.1784* (2014).
13. Lin, Kevin, et al. "Adversarial ranking for language generation." *Advances in Neural Information Processing Systems*. 2017.
14. Mikolov, Tomáš, et al. "Recurrent neural network based language model." *Eleventh annual conference of the international speech communication association*. 2010.
15. Zhang, Yizhe, Zhe Gan, and Lawrence Carin. "Generating text via adversarial training." *NIPS workshop on Adversarial Training*. Vol. 21. 2016.
16. Tuan, Yi-Lin, and Hung-Yi Lee. "Improving conditional sequence generative adversarial networks by stepwise evaluation." *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 27.4 (2019): 788-798.
17. Vinyals, Oriol, and Quoc Le. "A neural conversational model." *arXiv preprint arXiv:1506.05869* (2015).