

Comparing the interpretability of different GAN architectures

Henry Tu

University of Toronto

Seel Patel

University of Toronto

1 Abstract

A Generative Adversarial Network (GAN) uses an adversarial process between two models which are simultaneously trained to estimate a generative model.[1] There are many variants of GAN architectures, such as COCO-GAN[2] and StyleGAN[3], which both have the ability to generate synthetic images which mimic real images.

We are exploring methods of comparing the quantitative performance of these architectures with their interpretability. Our quantitative measures include C2ST[4], image quality measures[4], Maximum Mean Discrepancy

(MMD)[4], etc.

In order to analyze these networks qualitatively, we will use methods such as SmoothGrad[5], Latent Space Exploration[6] and nearest neighbour tests[4] to decipher what the networks have learned.

By combining these two forms of analysis, we hope to gain insight into the relationship between performance metrics and the generated output of the models.

2 Introduction

bsbsbsbsbsbs

References

- [1] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. (2014) *Generative Adversarial Networks*
- [2] Chieh Hubert Lin, Chia-Che Chang, Yu-Sheng Chen, Da-Cheng Juan, Wei Wei, Hwann-Tzong Chen. (2019) *COCO-GAN: Generation by Parts via Conditional Coordinating*
- [3] Tero Karras, Samuli Laine, Timo Aila. (2018) *A Style-Based Generator Architecture for Generative Adversarial Networks*
- [4] Brownlee, Jason. (2019) *How to Evaluate Generative Adversarial Networks*.
- [5] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viegas, Martin Wattenberg. (2017) *SmoothGrad: removing noise by adding noise*
- [6] Tom White. (2016) *Sampling Generative Networks*