

Background

Generative Adversarial Networks (GAN):

- Two networks competing with each other.
- Discriminator *D* tries to distinguish between real samples and samples generated by generator G.
- *G* tries to "fool" *D*.
- G will learn to generate samples similar to real data.



Motivation

Human painters usually first draw some abstract stuff, then gradually add details.

To mimic this process, we learn a generator that first produce high-level abstract features, then gradually generate lower level features and finally the image.



Stacked Generative Adversarial Networks Xun Huang^{1,2}, Yixuan Li², Omid Poursaeed^{1,2}, John Hopcroft², Serge Belongie^{1,2} ¹Cornell Tech ²Cornell University

Architecture

<u>A stack of GANs</u>, each GAN generates <u>lower-level</u> features conditioned on higher-level features.

Each generator is trained with three loss terms:

Adversarial loss: the generated features should be indistinguishable from "real" features.

 $\mathcal{L}_{G_i}^{adv} = \mathbb{E}_{z_i \sim P_{z_i}, h_{i+1} \sim P_{data, E}} \left[-\log(D_i(G_i(h_{i+1}, z_i))) \right]$

Conditional loss: the generator should make use of the higher-level features it's conditioned on:

 $\mathcal{L}_{G_i}^{cond} = \mathbb{E}_{h_{i+1} \sim P_{data,E}, \hat{h_i} \sim P_G(\hat{h_i}|h_{i+1}))} [f(E_i(\hat{h_i}), h_{i+1})]$

Entropy loss: encourage sample diversity by maxi-mizing a variational lower bound on the entropy

$$\mathcal{L}_{G_i}^{ent} = \mathbb{E}_{z'_i \sim P_{z'_i}} [\mathbb{E}_{\hat{h_i} \sim G_i(\hat{h_i}|z'_i)} [-\log Q_i(z'_i|\hat{h_i})]]$$



Qualitative results



Generated

Real

IEEE 2017 Conference on **Computer Vision and Pattern** Recognition



Quantitative evaluations

<u>Inception score</u> on CIFAR-10:

	Score
raining [1]	4.62 ± 0.06
(as reported in [63])	5.34 ± 0.05
1] (best variant)	6.00 ± 0.19
nt-VI [4]	7.07 ± 0.10
[<mark>65</mark>]	7.17 ± 0.07
g feature matching [63]	7.72 ± 0.13
(with labels, as reported in [61])	6.58
ν [†] [<mark>61</mark>]	6.35
GAN [†] [53] (best variant)	8.09 ± 0.07
† [43]	8.25 ± 0.07
(\mathcal{L}^{adv})	6.16 ± 0.07
$(\mathcal{L}^{adv} + \mathcal{L}^{ent})$	5.40 ± 0.16
$(\mathcal{L}^{adv}+\mathcal{L}^{cond})^{\dagger}$	5.40 ± 0.08
$(\mathcal{L}^{adv} + L^{cond} + \mathcal{L}^{ent})^{\dagger}$	7.16 ± 0.10
o-joint [†]	$\textbf{8.37} \pm 0.08$
	$\textbf{8.59} \pm 0.12$
	11.24 ± 0.12

[†] Trained with labels.

Human visual Turing tests on CIFAR-10: We ask AMT workers to distinguish generated images from real images. Our samples "fool" people **24.4%** of the time, higher than our best DCGAN baseline (15.6%) and Improved GAN (21.3%).



