

Boundary-Seeking Generative Adversarial Networks (BGANs)

Hjelm, R. Devon, et al.

Presenting: Yevgeny Tkach

2019 Spring @
<https://qdata.github.io/deep2Read/>

Executive Summary

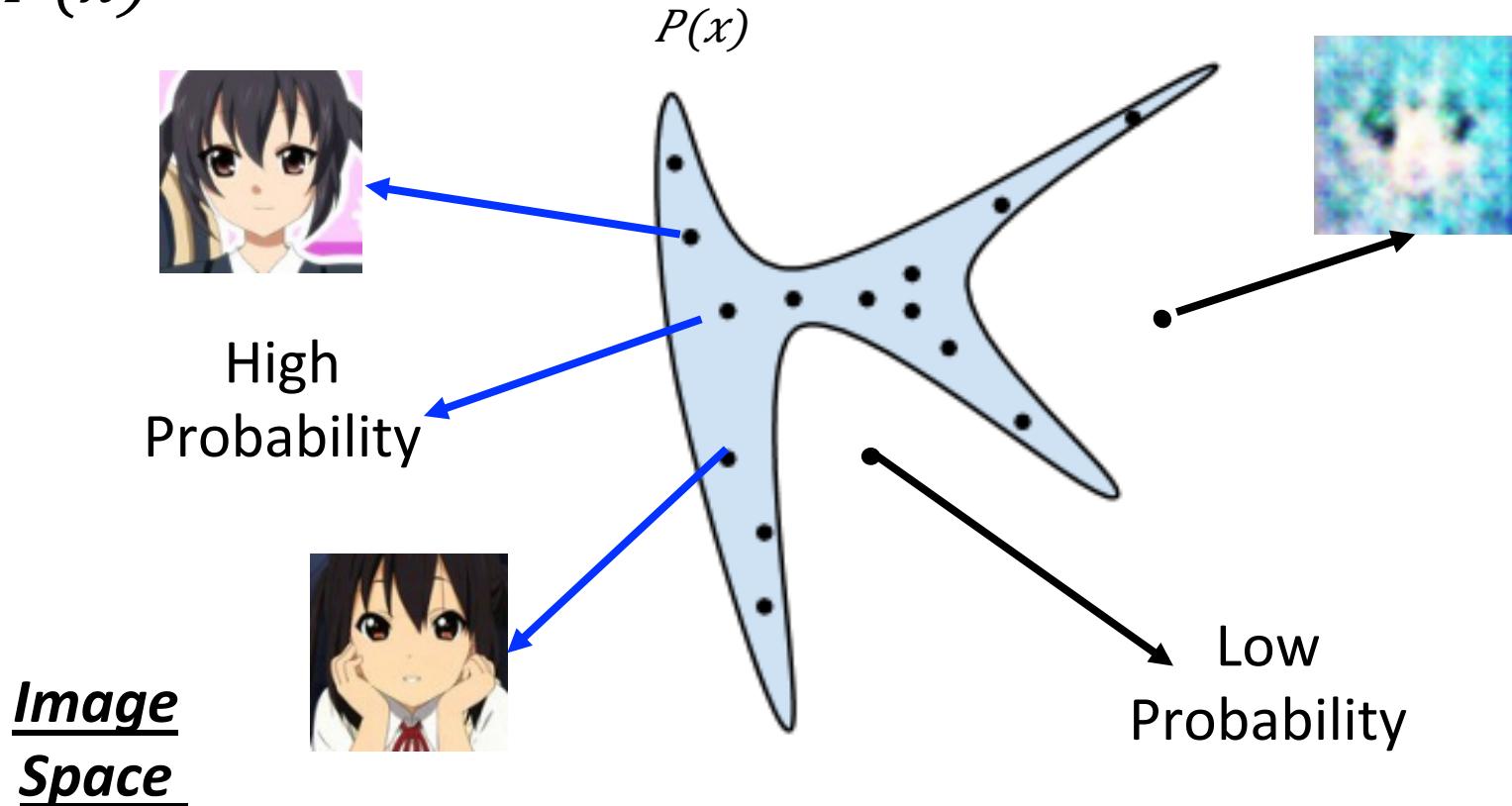
- BGAN is framework that allows GAN to generate both discrete and continuous data
- Discriminator is trained by maximizing the f-divergence between the data and generated distributions
- Generator is trained to minimize the f-divergence between the generated distribution and a self-normalized importance sampling (SIS) estimation of the data distribution
- Experiments show state of the art results in training GANs on discrete data generation and high stability in training GANs with continuous data.

Outline

- GAN – Basic Idea
- f - GAN Introduction
- Importance Sampling – Detour
- BGAN

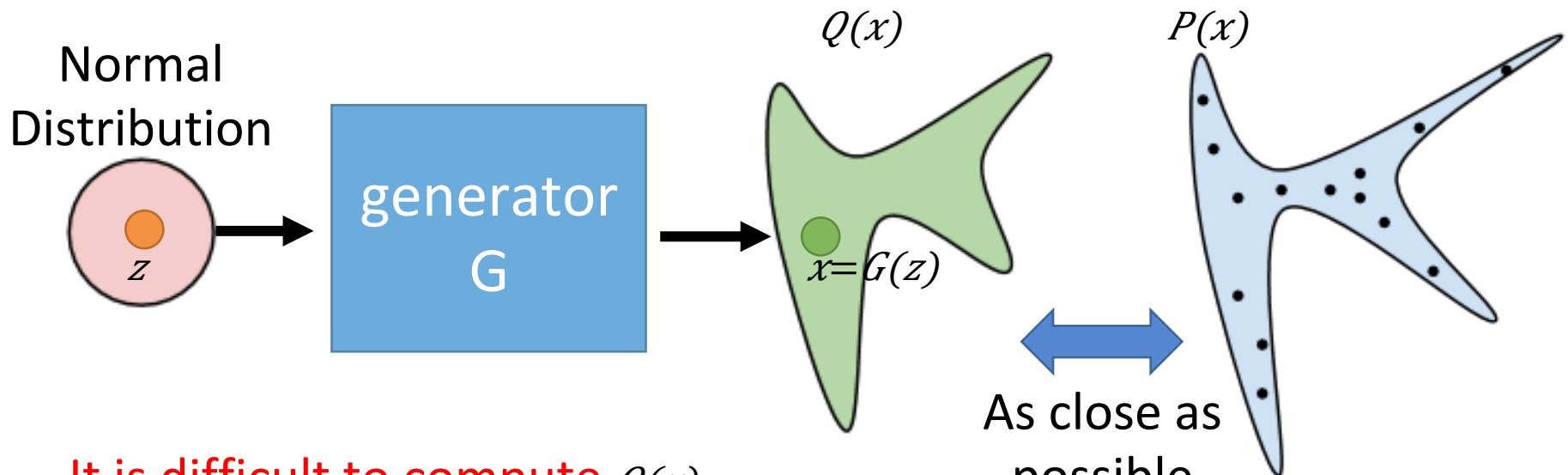
Basic Idea of GAN

- The data we want to generate has a distribution $P(x)$



Basic Idea of GAN

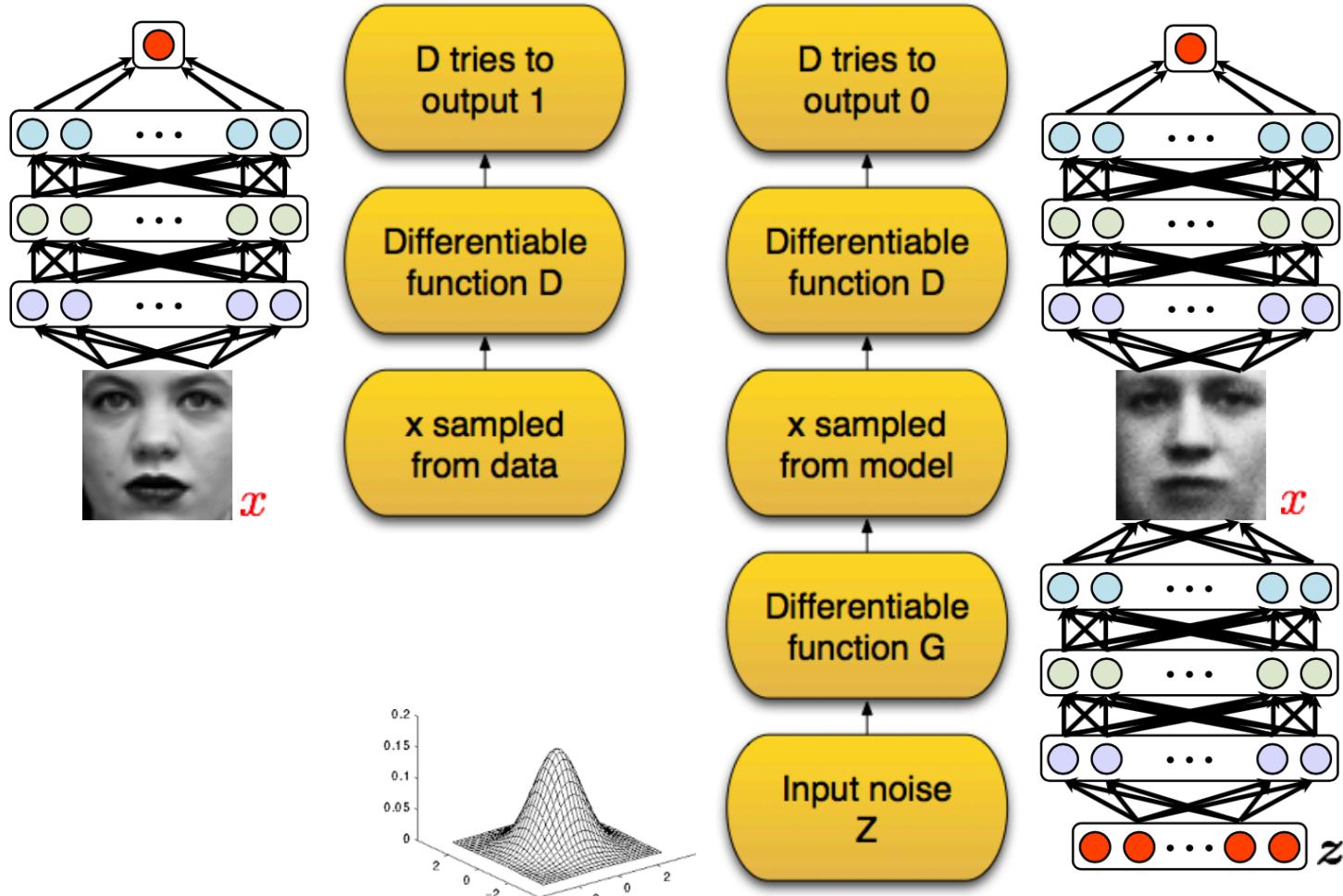
- A generator G is a network. The network defines a probability distribution.



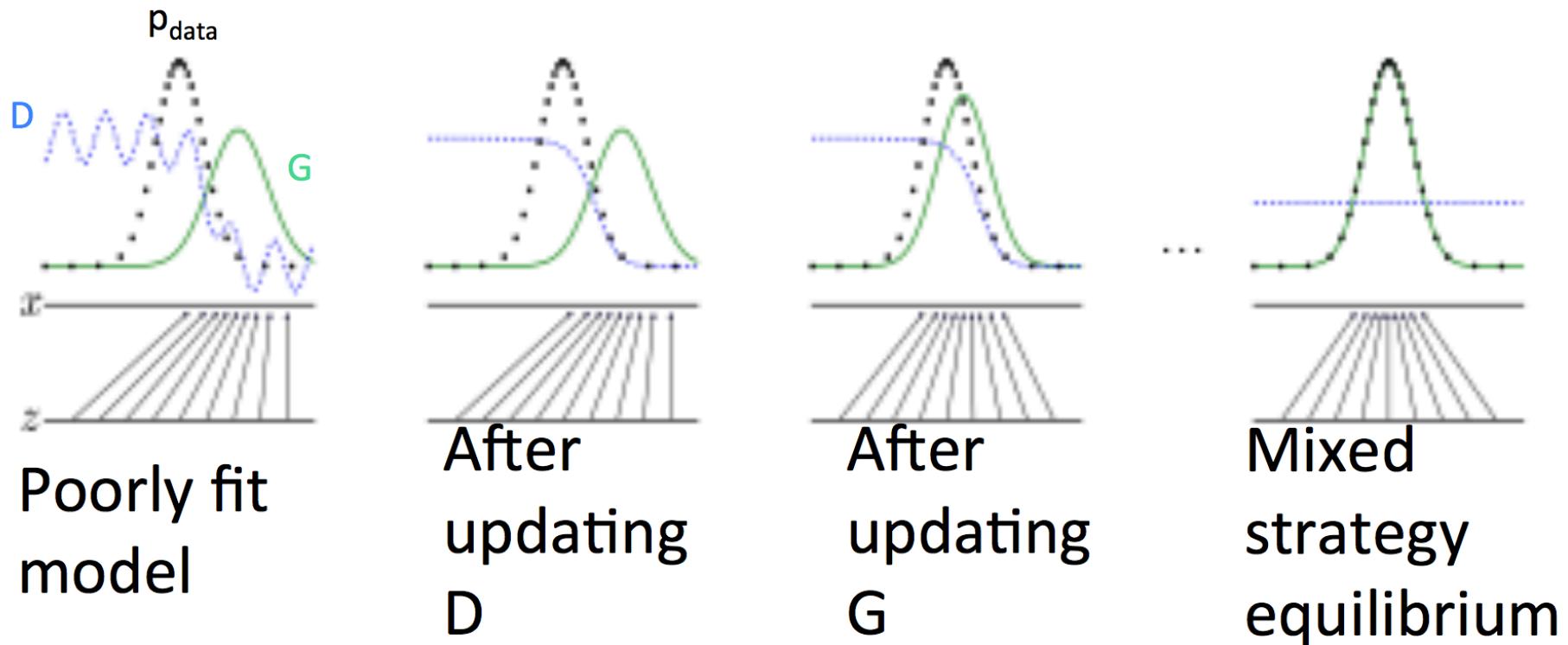
It is difficult to compute $Q(x)$

We can only sample from the distribution.

Basic Idea of GAN



GAN Intuition



GAN Formally

- Value Function:

$$\begin{aligned} V(\mathbb{P}, G \downarrow \theta, D \downarrow \phi) &= E \downarrow x \sim P [\log D(x)] + E \downarrow x \sim Q [\log(1 - D(x))] \\ &= E \downarrow x \sim P [\log D(x)] + E \downarrow z \sim h(z) [\log(1 - D(G(z)))] \end{aligned}$$

- Monte-Carlo Approximation:

$$V(\mathbb{P}, G \downarrow \theta, D \downarrow \phi) = 1/m \sum_{i=1}^m \log D(x \uparrow i) + 1/m \sum_{i=1}^m \log(1 - D(G(z \uparrow i)))$$

- Discriminator target:

$$\max_{\tau \phi} \square V(\mathbb{P}, G \downarrow \theta, D \downarrow \phi)$$

- Generator target:

$$\min_{\tau \theta} \square \max_{\tau \phi} \square V(\mathbb{P}, G \downarrow \theta, D \downarrow \phi)$$

Algorithm

Initialize $\phi \downarrow d$ for D and $\theta \downarrow g$ for G

- In each training iteration:

- Sample m examples $\{x^{\uparrow 1}, x^{\uparrow 2}, \dots, x^{\uparrow m}\}$ from data distribution $P(x)$
- Sample m noise samples $\{z^{\uparrow 1}, z^{\uparrow 2}, \dots, z^{\uparrow m}\}$ from the prior $h(z)$
- Obtaining generated data $\{x^{\uparrow 1}, x^{\uparrow 2}, \dots, x^{\uparrow m}\}, x^{\uparrow i} = G(z^{\uparrow i})$
- Update discriminator parameters $\theta \downarrow d$ to maximize
 - $V = 1/m \sum_{i=1}^m \log D(x^{\uparrow i}) + 1/m \sum_{i=1}^m \log(1 - D(x^{\uparrow i}))$
 - $\phi \downarrow d \leftarrow \phi \downarrow d + \eta \nabla V(\phi \downarrow d)$

Learning
D

Repeat
k times

Learning
G

Only
Once

- Sample another m noise samples $\{z^{\uparrow 1}, z^{\uparrow 2}, \dots, z^{\uparrow m}\}$ from the prior $P_{\text{prior}}(z)$

- Update generator parameters $\theta \downarrow g$ to minimize
 - $V = 1/m \sum_{i=1}^m \log D(x^{\uparrow i}) + 1/m \sum_{i=1}^m \log(1 - D(G(z^{\uparrow i})))$
 - $\theta \downarrow g \leftarrow \theta \downarrow g - \eta \nabla V(\theta \downarrow g)$

f - GAN Introduction

- Sebastian Nowozin, Botond Cseke, Ryota Tomioka, “*f-GAN*: Training Generative Neural Samplers using Variational Divergence Minimization”, NIPS, 2016
- One sentence: you can use any f-divergence

f-divergence

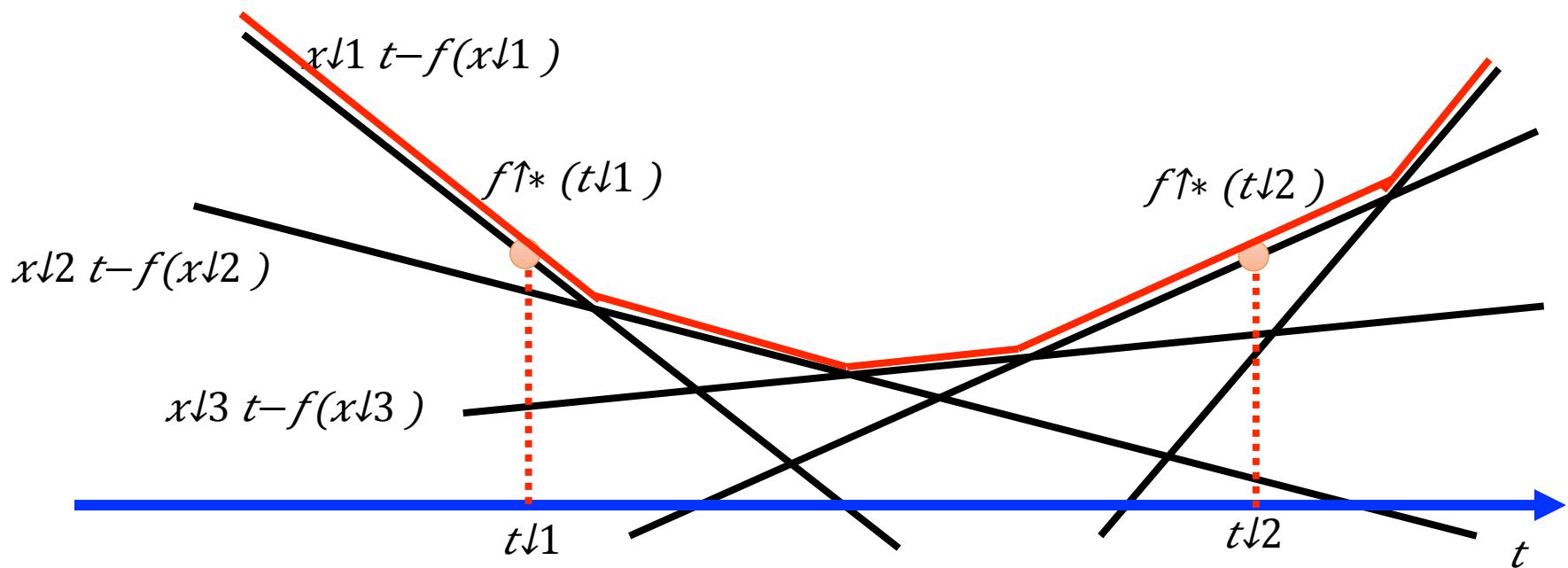
P and Q are two distributions. $p(x)$ and $q(x)$ are the density functions respectively.

$$D \downarrow f(P||Q) = \int x \uparrow q(x) f(p(x)/q(x)) dx$$

f is
convex

- Every convex function f has a conjugate function f^*

$$f^*(t) = \max_{x \in \text{dom}(f)} \{xt - f(x)\} \quad \longleftrightarrow \quad f(x) = \max_{t \in \text{dom}(f^*)} \{xt - f^*(t)\}$$



Connection with GAN

$$f^{\uparrow*}(t) = \max_{\tau} \max_{x \in \text{dom}(f)} \square\{xt - f(x)\}$$

$$f(x) = \max_{\tau} \max_{t \in \text{dom}(f^{\uparrow*})} \underline{\square}\{xt - f^{\uparrow*}(t)\}$$

$$D \downarrow f(P||Q) = \int x \uparrow \square q(x) f(p(x)/q(x)) dx$$

$$\begin{array}{c} p(x) \\ /q(x) \end{array}$$

$$\begin{array}{c} p(x) \\ /q(x) \end{array}$$

$$= \int x \uparrow \square q(x) (\max_{\tau} \max_{t \in \text{dom}(f^{\uparrow*})} \square\{p(x)/q(x) t - f^{\uparrow*}(t)\}) dx$$

— —

— —

D is a function whose input is x, and output is t

$$\geq \max_{\tau} \max_{D \in \mathcal{D}} \int x \uparrow \square q(x) (p(x)/q(x) D(x) - f^{\uparrow*}(D(x))) dx$$

— —

— —

$$= \max_{\tau} \max_{D \in \mathcal{D}} \int x \uparrow \square p(x) D(x) dx - \int x \uparrow \square q(x) f^{\uparrow*}(D(x)) dx$$

Connection with GAN

$$\geq \max_{\tau} \mathbb{D} \left\{ \int x^\top p(x) D(x) dx - \int x^\top q(x) f^{\star}(D(x)) dx \right\}$$

$$D \downarrow f(P||Q)$$

$$= \max_{\tau} \mathbb{D} \left\{ E \downarrow x \sim P [D(x)] - E \downarrow x \sim Q [f^{\star}(D(x))] \right\}$$

Samples from P Samples from Q

$$D \downarrow f(P||Q) \geq \max_{\tau} \mathbb{D} \left\{ E \downarrow x \sim P [\nu \circ D(x)] - E \downarrow x \sim Q [f^{\star}(\nu \circ D(x))] \right\}$$

$$\begin{aligned} G^{\star} &= \arg \min_{\tau} G \square D \downarrow f(P||Q) \\ &= \arg \min_{\tau} G \square \max_{\tau} D \square \left\{ E \downarrow x \sim P [\nu \circ D(x)] - E \downarrow z \sim h(z) [f^{\star}(\nu \circ D(G(z)))] \right\} \end{aligned}$$

GAN value function:

$$V(\mathbb{P}, G \downarrow \theta, D \downarrow \phi) = E \downarrow x \sim P [\log D(x)] + E \downarrow z \sim h(z) [\log(1 - D(G(z)))]$$

Importance Sampling – Detour

$$\begin{aligned} E \downarrow_{x \sim P} [f(x)] &= \int \uparrow \# f(x) p(x) dx \\ &= \int \uparrow \# f(x) p(x) / q(x) q(x) dx \\ &= \int \uparrow \# f(x) w(x) q(x) dx \\ &= E \downarrow_{x \sim Q} [f(x) w(x)] \\ &= E \downarrow_{x \sim Q} [f(x) w(x)] / E \downarrow_{x \sim Q} [w(x)] \\ w(x) &= p(x) / q(x) \end{aligned}$$

In case p or q
are scaled
density
functions

$w(x)$ - Importance
Weights

Boundary Seeking GAN - BGAN

Theorem 1: P and Q as in f-GAN, and $D \uparrow^* \in D$ satisfying:

$$D \downarrow f(P||Q) = \max_{\tau \in D} \{ E \downarrow x \sim P [D(x)] - E \downarrow x \sim Q [f \uparrow^* (D(x))] \}$$

Then: $p(x) = (\partial f \uparrow^* / \partial D)(D \uparrow^*(x))q(x)$

Proof:

$$D \downarrow f(P||Q) = E \downarrow x \sim Q [f(p(x)/q(x))] = E \downarrow x \sim Q [\sup_{t \in T} \{tp(x)/q(x) - f \uparrow^*(t)\}]$$

p re-written in terms of q and a scaling factor
 $w(x) = (\partial f \uparrow^* / \partial D)(D \uparrow^*(x))$ - Importance weights

$$\frac{p(x)}{q(x)} = \frac{\partial}{\partial t} f \uparrow^*(t)$$

Boundary Seeking GAN - BGAN

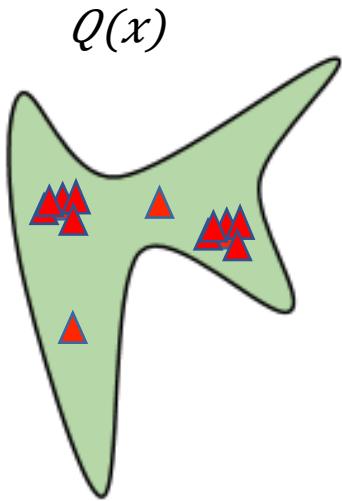
BGAN suggests to use the **divergence** between $q(x)$ and the self normalized importance sampling (IS) estimation of $p(x)$:

$$p(x) = w(x) / \beta q(x)$$

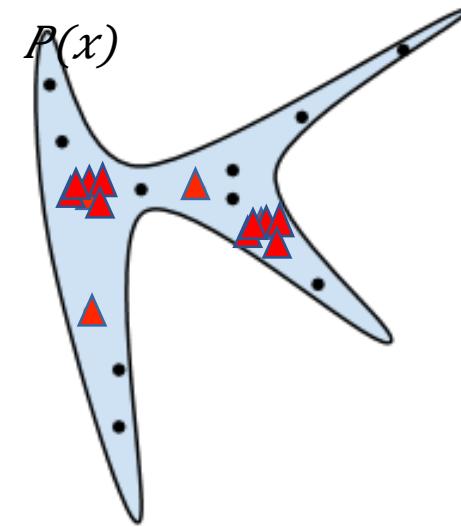
Where:

$$\beta = E_{x \sim Q} [w(x)]$$

BGAN – IS intuition



IS proxy with optimal
discriminator D



- Divergence between Δ should have lower variance than if taking arbitrary samples from $P(x)$
- Since $G(z)$ defines a distribution that x is sampled from - the variance can be further decreased by taking multiple samples from the same z

BGAN – reduced variance

We can restate everything in terms of conditional distributions:

- $q(x) = \int Z \uparrow g(x|z) h(z) dz$
- $g(x|z): Z \rightarrow [0,1]^d$ - multivariate Bernoulli distribution
- $\alpha(z) = E \downarrow x \sim g(x|z) [w(x)]$ - similar to β
- $p(x|z) = w(x) / \alpha(z)$
- $D_{KL}(p(x)||q(x)) = E \downarrow h(z) [D_{KL}(p(x|z)||q(x|z))]$
- $\nabla E \downarrow h(z) [D_{KL}(p(x|z)||q(x|z))] \text{ approximates with two MC}$

BGAN - Algorithm

Algorithm 1 . Discrete Boundary Seeking GANs

$(\theta, \phi) \leftarrow$ initialize the parameters of the generator and statistic network

repeat

$$\hat{x}^{(n)} \sim \mathbb{P}$$

▷ Draw N samples from the empirical distribution

$$z^{(n)} \sim h(z)$$

▷ Draw N samples from the prior distribution

$x^{(m|n)} \sim g_\theta(x | z^{(n)})$ ▷ Draw M samples from each conditional $g_\theta(x | z^{(m)})$ (drawn independently if \mathbb{P} and \mathbb{Q}_θ are multi-variate)

$$w(x^{(m|n)}) \leftarrow (\partial f^*/\partial T) \circ (\nu \circ F_\phi(x^{(m|n)}))$$

$\tilde{w}(x^{(m|n)}) \leftarrow w(x^{(m|n)}) / \sum_{m'} w(x^{(m'|n)})$ ▷ Compute the un-normalized and normalized importance weights (applied uniformly if \mathbb{P} and \mathbb{Q}_θ are multi-variate)

$\mathcal{V}(\mathbb{P}, \mathbb{Q}_\theta, T_\phi) \leftarrow \frac{1}{N} \sum_n F_\phi(\hat{x}^{(n)}) - \frac{1}{N} \sum_n \frac{1}{M} \sum_m w(x^{(m|n)})$ ▷ Estimate the variational lower-bound

$$\phi \leftarrow \phi + \gamma_d \nabla_\phi \mathcal{V}(\mathbb{P}, \mathbb{Q}_\theta, T_\phi)$$

▷ Optimize the discriminator parameters

$$\theta \leftarrow \theta + \gamma_g \frac{1}{N} \sum_{n,m} \tilde{w}(x^{(m|n)}) \nabla_\theta \log g_\theta(x^{(m|n)} | z)$$

▷ Optimize the generator parameters

until convergence

Boundary Seeking GAN - BGAN

$$D\downarrow f(P\downarrow \text{data} || P\downarrow G) \geq \max_{T \in \mathcal{D}} \{ E\downarrow x \sim P\downarrow \text{data} [v \circ D(x)] - E\downarrow x \sim P\downarrow G [f^{\uparrow *} (v \circ D(x))] \}$$

$$p(x) = w(x)/\beta q(x) \quad w(x) = (\partial f^{\uparrow *} / \partial D)(D^{\uparrow *} (x))$$

Table 1: Important weights and nonlinearities that ensure

Importance weights for f -divergences		
f -divergence	$v(y)$	$w(x) = (\partial f^* / \partial T)(T(x))$
GAN	$-\log(1 + e^{-y})$	$-\frac{1}{1 - e^{-T_\phi}} = e^{F_\phi(x)}$
Jensen-Shannon	$\log 2 - \log(1 + e^{-y})$	$-\frac{1}{2 - e^{-T_\phi}} = e^{F_\phi(x)}$
KL	$y + 1$	$e^{(T_\phi(x) - 1)} = e^{F_\phi(x)}$
Reverse KL	$-e^{-y}$	$-\frac{1}{T_\phi(x)} = e^{F_\phi(x)}$
Squared-Hellinger	$1 - e^{-v/2}$	$\frac{1}{(1 - T_\phi(x))^2} = e^{F_\phi(x)}$

BGAN – Experiments



Train Measure	Eval Measure (lower is better)		
	JS	reverse KL	Wasserstein
BGAN - JS	0.37 (± 0.02)	0.16 (± 0.01)	0.40 (± 0.03)
BGAN - reverse KL	0.44 (± 0.02)	0.44 (± 0.03)	0.45 (± 0.04)
WGAN-GP (samples)	0.45 (± 0.03)	1.32 (± 0.06)	0.87 (± 0.18)
WGAN-GP (softmax)	-	-	0.54 (± 0.12)

BGAN – Experiments



And it 's miant a quert could he
" We pait of condels of money wi
Lankard Avaloma was Mr. Palin ,
Thene says the sounded Sunday in
About dose and warthestrinds fro

He weirst placed produces hopesi
Sance Jory Chorotic , Sen doesin
What was like one of the July 2
The BBC nothing overton and slea
College is out in contesting rev

BGAN – Continuous case

Recall:

$$G^{\uparrow*} = \operatorname{argmin}_{\tau} G \square D \downarrow f(P \downarrow \text{data} || P \downarrow G)$$

$$D \downarrow f(P||Q) = E \downarrow x \sim Q [f(p(x)/q(x))] = E \downarrow x \sim Q [\sup_{\tau} \tau t \square \{tp(x)/q(x) - f^{\uparrow*}(t)\}]$$

⇓

$$p(x)/q(x) = (\partial f^{\uparrow*} / \partial D)(D^{\uparrow*}(x)) = w(x)$$

⇓

$$G^{\uparrow*} = \operatorname{argmin}_{\tau} G \square (\log w(G(z)))^{1/2} \quad \text{Max when } \nabla \{tp(x)/q(x) - f^{\uparrow*}(t)\} = 0$$

⇓

$$G^{\uparrow*} = \operatorname{argmin}_{\tau} G \square D(G(z))^{1/2}$$

$$p(x) = q(x) \text{ when } w(x) = 1$$

BGAN – Continuous case

f -GAN:

$$G^* = \operatorname{argmin}_G \mathbb{E} \{ E[x] - E[z] [f^*(\nu \circ D(G(z)))] \}$$

GAN (Proxy GAN):

$$G^* = \operatorname{argmin}_G \mathbb{E} \{ E[x] + E[z] [\log(1 - D(G(z)))] \}$$

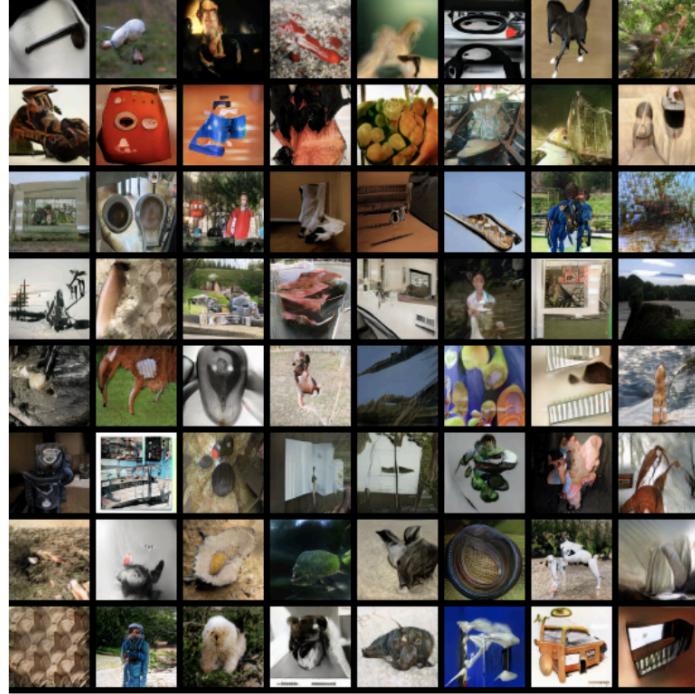
BGAN:

$$G^* = \operatorname{argmin}_G \mathbb{E} [D(G(z))] \Leftrightarrow w(x) = 1 \Leftrightarrow p(x) = q(x)$$

BGAN – Continuous Experiments



CelebA



Imagenet

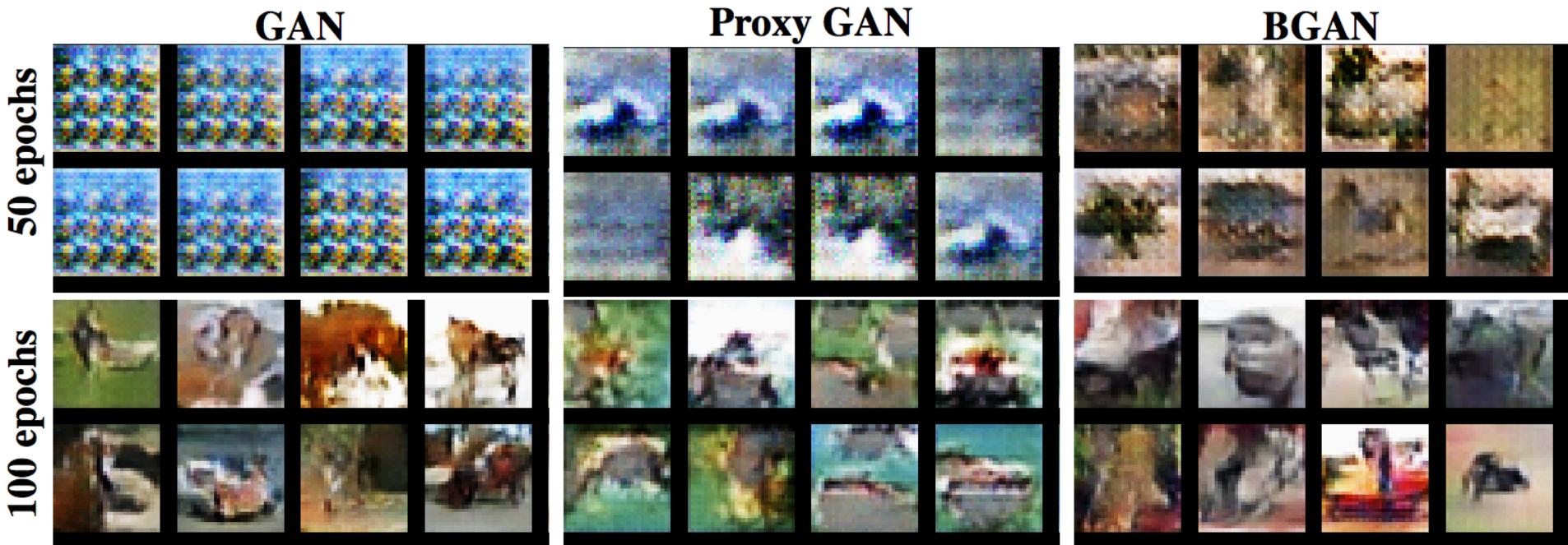


LSUN

Figure 3: Highly realistic samples from a generator trained with BGAN on the CelebA and LSUN datasets. These models were trained using a deep ResNet architecture with gradient norm regularization (Roth et al., 2017). The Imagenet model was trained on the full 1000 label dataset without conditioning.

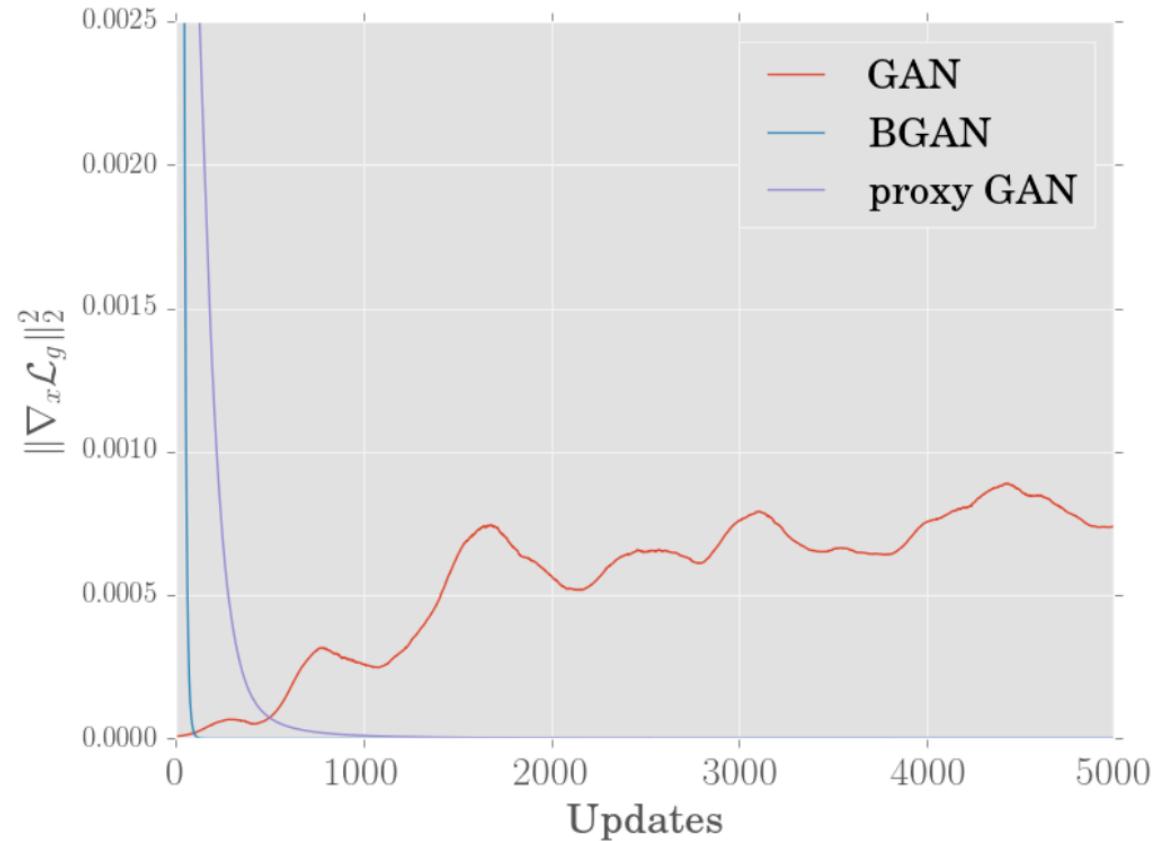
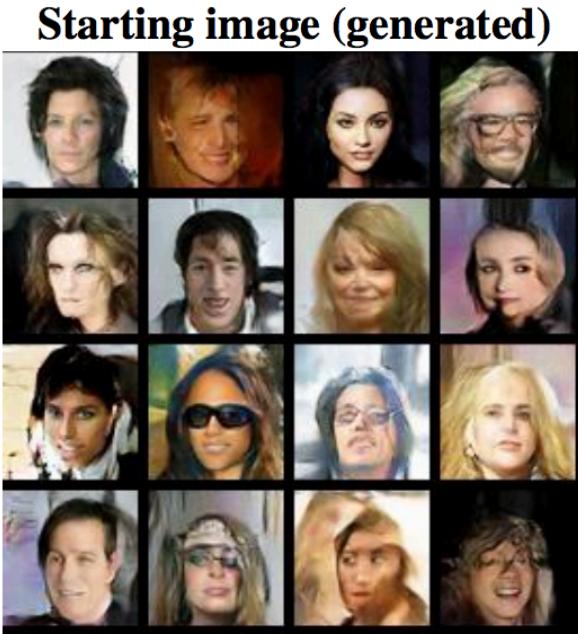
BGAN – Continuous Experiments

- Generator trained for 5 steps for every 1 step of the discriminator



BGAN – Continuous Experiments

- Train a DCGAN using the proxy loss.
- Train the discriminator for 1000 more steps
- Perform gradient descent directly on the pixels



BGAN – Continuous Experiments

10k updates

GAN



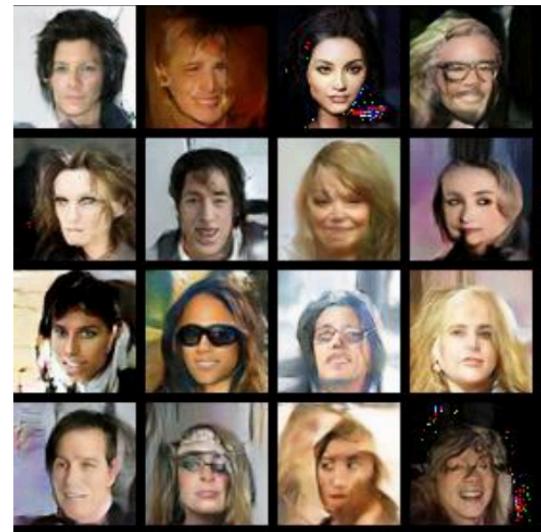
Proxy GAN



BGAN



20k updates



Discussion