Understanding Generative Adversarial Networks

Balaji Lakshminarayanan



Joint work with: Shakir Mohamed, Mihaela Rosca, Ivo Danihelka, David Warde-Farley, Liam Fedus, Ian Goodfellow, Andrew Dai & others



Problem statement

Learn a **generative model**:

$$\mathbf{x} = \mathcal{G}_{\boldsymbol{\theta}}(\mathbf{z}'); \qquad \mathbf{z}' \sim q(\mathbf{z})$$

Goal: given samples $x_1 \dots x_n$ from true distribution $p^*(x)$, find θ

p_e(**x**) is not available -> can't maximize density directly

However, we can sample from $p_{\theta}(x)$ efficiently



High level overview of GANs

Goodfellow et al. 2014

Discriminator: Train a classifier to distinguish between the two distributions *using samples*

Generator: Train to generate samples that fool the discriminator

Minimax game alternates between training discriminator and generator

- Nash equilibrium corresponds to minima of Jensen Shannon divergence
- Need a bunch of tricks to stabilize training in practice





GANs: Hope or Hype?



Ferenc Huszar

practicin' the alphabet with my son: A is for AffGAN B is for B-GAN C is for Conditional GAN D is for DCGAN E is for EBGAN F is for f-GAN



https://github.com/hindupuravinash/the-gan-zoo



https://github.com/junyanz/CycleGAN/blob/master/imgs/horse2zebra.gif



https://github.com/tkarras/progressive_growing_of_gans



Understanding GANs

How do GANs relate to other ideas in probabilistic machine learning?

Learning in implicit generative models

Shakir Mohamed* and Balaji Lakshminarayanan*



Understanding GANs

Implicit Models

Stochastic procedure that generates data

$$\mathbf{z}$$

$$\mathbf{x} = \mathcal{G}(\mathbf{z}'), \quad \mathbf{z}' \sim q(\mathbf{z})$$

$$\mathbf{x} \quad q(\mathbf{x}) = \frac{\partial}{\partial x_1} \dots \frac{\partial}{\partial x_d} \int_{\{\mathcal{G}(\mathbf{x}) \leq \mathbf{x}\}} q(\mathbf{z}) d^m \mathbf{z}$$

Examples: stochastic simulators of complex physical systems (climate, ecology, high-energy physics etc)

Prescribed Models

Provide knowledge of the probability of observations & specify a *conditional* log-likelihood function.

$$\begin{array}{c} \mathbf{z} \\ \mathbf{x}' \sim p(\mathbf{x}|\mathcal{G}(\mathbf{z}')), \quad \mathbf{z}' \sim q(\mathbf{z}) \\ \mathbf{x} \\ \end{array} \\ \mathbf{q}(\mathbf{x}) = \int p(\mathbf{x}|\mathcal{G}(\mathbf{z}))q(\mathbf{z}) \end{array}$$



Learning by Comparison





We compare the estimated distribution to the true distribution **using samples.**



Understanding GANs

Learning by Comparison

Comparison

Use a hypothesis **test or comparison** to build an auxiliary model to indicate how data simulated from the model differs from observed data.

Learning generator

Adjust model parameters to better match the data distribution using the comparison.





High-level idea

Define a joint loss function $L(\phi, \theta)$ and alternate between:

Comparison loss ("discriminator"): **arg min** $_{\omega}$ L(ϕ , θ)

Generative loss: arg min_{θ} -L(ϕ , θ)

Global optimum is $q_{\theta} = p^*$ and

- Density ratio $r_{\phi} = 1$ or
- Density difference $r_{\phi} = 0$

How do we compare distributions?



Density Ratios and Classification



		Real Data	Simulated Data	
Combine data	$\{\mathbf{x}_1,\ldots,\mathbf{x}_N\}=0$	$\{\hat{\mathbf{x}}_1,\ldots,\hat{\mathbf{x}}_{\hat{n}},$	$ ilde{\mathbf{x}}_1,\ldots, ilde{\mathbf{x}}_{ ilde{n}}\}$	
Assign labels	$\{y_1,\ldots,y_N\}=\{$	$[+1,\ldots,+1]$	$-1,\ldots,-1\}$	

Sugiyama et al, 2012



Density Ratios and Classification



Computing a density ratio is equivalent to class probability estimation.



Class Probability Estimation

$$\mathcal{D}(\mathbf{x}; \boldsymbol{\phi}) = p(\mathbf{y} = +1 | \mathbf{x}) = \frac{r}{r+1}$$

Loss	Objective Function $(\mathcal{D} := \mathcal{D}(\mathbf{x}; \boldsymbol{\phi}))$
Bernoulli loss	$\pi \mathbb{E}_{p^*(\mathbf{x})}[-\log \mathcal{D}] + (1-\pi) \mathbb{E}_{q_\theta(\mathbf{x})}[-\log(1-\mathcal{D})]$
Brier score	$\pi \mathbb{E}_{p^*(\mathbf{x})}[(1-\mathcal{D})^2] + (1-\pi)\mathbb{E}_{q_{\theta}(\mathbf{x})}[\mathcal{D}^2]$
Exponential loss	$\pi \mathbb{E}_{p^*(\mathbf{x})} \left[\left(\frac{1-\mathcal{D}}{\mathcal{D}} \right)^{\frac{1}{2}} \right] + (1-\pi) \mathbb{E}_{q_{\theta}(\mathbf{x})} \left[\left(\frac{\mathcal{D}}{1-\mathcal{D}} \right)^{\frac{1}{2}} \right]$
Misclassification	$\pi \mathbb{E}_{p^*(\mathbf{x})}[\mathbb{I}[\mathcal{D} \le 0.5]] + (1 - \pi) \mathbb{E}_{q_{\theta}(\mathbf{x})}[\mathbb{I}[\mathcal{D} > 0.\vec{5}]]$
Hinge loss	$\pi \mathbb{E}_{p^*(\mathbf{x})} \left[\max\left(0, 1 - \log \frac{\mathcal{D}}{1 - \mathcal{D}}\right) \right] + (1 - \pi) \mathbb{E}_{q_\theta(\mathbf{x})} \left[\max\left(0, 1 + \log \frac{\mathcal{D}}{1 - \mathcal{D}}\right) \right]$
Spherical	$\pi \mathbb{E}_{p^*(\mathbf{x})} \left[-\alpha \mathcal{D} \right] + (1-\pi) \mathbb{E}_{q_\theta(\mathbf{x})} \left[-\alpha (1-\mathcal{D}) \right]; \alpha = (1-2\mathcal{D}+2\mathcal{D}^2)^{-1/2}$

Table 1. Proper scoring rules that can be minimised in class probability-based learning of implicit generative models.

Other loss functions for training classifier, e.g. Brier score leads to LS-GAN

Related: Unsupervised as Supervised Learning, Classifier ABC

Divergence minimization (f-GAN)

$$D_{f}\left[p^{*}(\mathbf{x}) \| q_{\theta}(\mathbf{x})\right] = \int q_{\theta}(\mathbf{x}) f\left(\frac{p^{*}(\mathbf{x})}{q_{\theta}(\mathbf{x})}\right) d\mathbf{x}$$
$$= \mathbb{E}_{q_{\theta}(\mathbf{x})}[f(r(\mathbf{x}))]$$
$$\geq \sup_{t} \mathbb{E}_{p^{*}(\mathbf{x})}[t(\mathbf{x})] - \mathbb{E}_{q_{\theta}(\mathbf{x})}[f^{\dagger}(t(\mathbf{x}))]$$

Minimize a lower bound on f-divergence between p* and q_{ρ}

Choices of *f* recover KL(p*||q) (maximum likelihood), KL(q||p*) and JS(p*||q)

Can use different f-divergences for learning ratio vs learning generator

Density ratio estimation

$$\begin{split} B_f(r^*(\mathbf{x}) \| r_{\phi}(\mathbf{x})) \\ &= \int \left(f(r^*(\mathbf{x})) - f(r_{\phi}(\mathbf{x})) - f'(r_{\phi}(\mathbf{x})) \left[r^*(\mathbf{x}) - r_{\phi}(\mathbf{x}) \right] \right) q_{\theta}(\mathbf{x}) d\mathbf{x} \\ &= \mathbb{E}_{q_{\theta}(\mathbf{x})} \left[r_{\phi}(\mathbf{x}) f'(r_{\phi}(\mathbf{x})) - f(r_{\phi}(\mathbf{x})) \right] - \mathbb{E}_{p^*} [f'(r_{\phi}(\mathbf{x}))] + D_f[p^*(\mathbf{x}) \| q_{\theta}(\mathbf{x})] \\ &= \mathcal{L}_B(r_{\phi}(\mathbf{x})) + D_f[p^*(\mathbf{x}) \| q_{\theta}(\mathbf{x})] \end{split}$$

Optimize a Bregman divergence between r* and r_{o}

Special cases include least squares importance fitting (LSIF)

Ratio loss ends up being identical to that of f-divergence

Moment-matching

$$\begin{aligned} \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) &= (\mathbb{E}_{p^*(\mathbf{x})}[s(\mathbf{x})] - \mathbb{E}_{q_{\theta}(\mathbf{x})}[s(\mathbf{x})])^2 \\ &= (\mathbb{E}_{p^*(\mathbf{x})}[s(\mathbf{x})] - \mathbb{E}_{q(\mathbf{z})}[s(\mathcal{G}(\mathbf{z}; \boldsymbol{\theta}))])^2 \end{aligned}$$

Used by

- Generative moment matching networks
- Training generative neural networks via Maximum Mean Discrepancy optimization

Connects to optimal transport literature (e.g. Wasserstein GAN)

Summary of the approaches

Class probability estimation

- Build a classifier to distinguish real from fake samples.
- Original GAN solution.

Divergence Minimisation

- Minimise a generalised divergence between the true density p* and the product r(x)q(x).
- f-GAN approach.

Density ratio matching

- Directly minimise the expected error between the true ratio and an estimate of it.

Moment matching

- Match the moments between p* and r(x)q(x)
- MMD, optimal transport, etc.



How do we learn generator?

In GANs, the generator is **differentiable**

- Generator loss is of the following form e.g. f-divergence D_f = E_q [f(r)]
- Can apply re-parametrization trick

$$J = E_{q_{\theta}(\mathbf{x})}[\ell(\mathbf{x})] = E_{q(\mathbf{z})}[\ell(\mathcal{G}_{\theta}(\mathbf{z}))]$$
$$\nabla_{\theta}J = \nabla_{\theta}E_{q(\mathbf{z})}[\ell(\mathcal{G}_{\theta}(\mathbf{z}))] = E_{q(\mathbf{z})}[\nabla_{\theta}\ell(\mathcal{G}_{\theta}(\mathbf{z}))]$$



Choice of f-divergence



Density ratio estimation literature has investigated choices of f

However, that's only half of the puzzle. We need non-zero gradients for $D_f = E_q [f(r)]$ to learn generator

- r<<1 early on in training</p>
- Non-saturating alternative loss

We also need additional constraints on the discriminator

Figure 2. Objective functions for different choices of f.

Summary: Learning in Implicit Generative Models

Unifying view* of GANs that connects to literature on

- Density ratio estimation
 - ... but they don't focus on learning generator
- Approximate Bayesian computation (ABC) and likelihood-free inference
 - Low dimensional, better understanding of theory
 - Bayesian inference over parameters
 - Simulators are usually not differentiable (can we approximate them?)

Motivates new loss functions: can decouple generator loss from discriminator loss

GAN-like ideas can be used in other places where density ratio appears



Comparing GANs to Maximum Likelihood training using Real-NVP

Comparison of maximum likelihood and GAN-based training of Real NVPs

Ivo Danihelka, Balaji Lakshminarayanan, Benigno Uria, Daan Wierstra and Peter Dayan



Understanding GANs

Generative Models and Algorithms





Generative Models and Algorithms



DeepMind

Understanding GANs

Comparing inference algorithms for a fixed model

Generator is Real NVP (<u>Dinh et al., 2016</u>)

$$\log P(z_1) = \log P(z_0) - \log |\det \frac{dz_1}{dz_0}|$$



- 1. Train by maximum likelihood (MLE).
- 2. Train a generator by Wasserstein GAN.
- 3. Compare.

Complementary to "On the quantitative analysis of decoder-based models" by Wu et al., 2017

Wasserstein GAN

For general distributions:

$$W_{d}(P_{r}, P_{g}) = \sup_{\substack{\|f\|_{L} \leq 1}} \mathbb{E} \left[f(x_{r})\right] - \mathbb{E} \left[f(x_{g})\right]$$

$$/$$
Considering all 1-Lipschitz function
(i.e., functions with **bounded derivatives**).
$$f(x) \text{ is a "}$$

$$f(x) \text{ is a "}$$

f(x) is a *"critic"*. The critic should give high value to real samples and low value to generated samples.

Bounded by:

- a) Weight clipping (Wasserstein GAN; "WGAN").
- b) Gradient penalty (Improved Training; "WGAN-GP")

Idea: use an independent Wasserstein critic to evaluate generators

Bits/dim for NVP

Dataset: CelebA 32x32.





Wasserstein Distance for NVPs





Wasserstein Distance Minimized by WGAN





Understanding GANs

MLE vs. WGAN Training







Understanding GANs

MLE vs. WGAN Training (shallower generator)





Understanding GANs

Bits/dim for NVPs Trained by WGAN





Summary

- Wasserstein distance can compare models.
- Wasserstein distance can be approximated by training a critic.
- Training by WGAN leads to nicer samples but significantly worse log-probabilities.
- Latent codes from WGAN training are non-Gaussian





Understanding GANs

How do we combine VAEs and GANs to get the best of both worlds?

Variational approaches for auto-encoding generative adversarial networks

Mihaela Rosca*, Balaji Lakshminarayanan*, David Warde-Farley and Shakir Mohamed



Understanding GANs

Motivating problem: Mode collapse



- MoG toy example from "Unrolled GAN" paper
- VAEs have other problems, but do not suffer from mode-collapse
 - Can we add auto-encoder to GANs?

Adding auto-encoder to GANs





How does it relate to Evidence Lower Bound (ELBO) in VAEs?





Recap: Density ratio trick

Estimate the ratio of two distributions only from samples, by building a binary **classifier** to distinguish between them.





Revisiting ELBO in Variational Auto-Encoders

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \log \int p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \neq \mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] + \mathrm{KL}[q_{\eta}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$$

$$\mathsf{LIKELIHOOD TERM}$$

$$\mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] = \mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})}[\log(\frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})}{p(\mathbf{x})}p(\mathbf{x}))]$$

$$\mathbf{x} \sim p^{*}(\mathbf{x}) = \mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})}[\log\frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})}{p(\mathbf{x})}] + \mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x})]$$

$$\underbrace{\mathcal{D}} \to 0/1$$

$$\mathbf{x} \sim p_{\boldsymbol{\theta}}(\mathbf{x})$$



Revisiting ELBO in Variational Auto-Encoders

$$\log p_{\theta}(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \ge \mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathrm{KL}[q_{\eta}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$$



Encoder can be implicit!

More flexible distributions





 $\mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})} \left[\log p_{\theta}(\mathbf{x}|\mathbf{z}) \right] \approx \mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})} \left[-\lambda ||\mathbf{x} - \mathcal{G}_{\theta}(\mathbf{z})||_{1} \right]$



Combining VAEs and GANs

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \log \int p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \ge \mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \mathrm{KL}[q_{\eta}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$$

- Likelihood: Reconstruction vs "synthetic likelihood" term
- KL: Analytical vs "code discriminator"
- Can recover various hybrids of VAEs and GANs

Algorithm	Likelihood		Prior		
	Observer	Ratio estimator ("synthetic")	KL (analytic)	KL (approximate)	Ratio estimator
VAE	√		 ✓ 		
DCGAN		\checkmark			
VAE-GAN	\checkmark	*	\checkmark		
Adversarial-VB	\checkmark				\checkmark
AGE	\checkmark			\checkmark	
α -GAN (ours)	\checkmark	\checkmark			\checkmark

Table 1: Comparison of different approaches for training generative latent variable models.



Evaluating different variants



Our VAE-GAN hybrid is competitive with state-of-the-art GANs



Cifar10 - Inception score



Classifier trained on Cifar10

Improved Techniques for Training GANs T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, X. Chen



Understanding GANs

CelebA - sample diversity





Summary: VAEs and GANs





Understanding GANs

Bridging the gap between theory & practice

Many paths to equilibrium: GANs do not need to decrease a divergence at every step

William Fedus*, Mihaela Rosca*, Balaji Lakshminarayanan, Andrew Dai, Shakir Mohamed & Ian Goodfellow



Differences between GAN theory and practice

Lots of new GAN variants have been proposed (e.g. Wasserstein GAN)

- Loss functions & regularizers motivated by new theory
- Significant difference between theory and practice

How do we bridge this gap?

- Synthetic datasets where theory predicts failure
- Add new regularizers to original non-saturating GAN



Non-Saturating GAN





Understanding GANs

Gradient Penalties for Discriminators

$$\tilde{J}^{(D)}(D,G) = -\mathbb{E}_{x \sim p_{\text{data}}} \left[\log D(x) \right] - \mathbb{E}_{z \sim p_z} \left[\log(1 - D(G(z))) \right] + \lambda \mathbb{E}_{\hat{x} \sim p_{\hat{x}}} \left[(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2 \right]$$

To formalize the above, both proposed gradient penalties of the form:

$$\mathbb{E}_{\hat{x} \sim p_{\hat{x}}} \left[(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2 \right],$$

where $p_{\hat{x}}$ is defined as the distribution defined by the sampling process:

 $x \sim p_{\text{data}};$ $x_{\text{model}} \sim p_{\text{model}};$ $x_{\text{noise}} \sim p_{\text{noise}}$

DRAGAN
$$\tilde{x} = x + x_{noise}$$

WGAN-GP $\tilde{x} = x_{model}$
 $\alpha \sim U(0, 1)$
 $\hat{x} = \alpha x + (1 - \alpha)\tilde{x}.$





Comparisons on synthetic dataset where Jensen Shannon divergence fails

- Gradient penalties lead to better performance



Results on real datasets





Results on real datasets





Summary

Some surprising findings:

- Gradient penalties stabilize (non-Wasserstein) GANs as well
- Think not just about the ideal loss function but also the optimization

"In theory, there is no difference between theory and practice. In practice, there is."

- Better ablation experiments will help bridge this gap and move us closer to the holy grail



Other interesting research directions



Overloading GANs and "Adversarial training"

Originally formulated as a minimax game between a discriminator and generator

Recent insights:

- **Density ratio trick**: discriminator estimates a density ratio. Can replace density ratios and f-divergences in message passing with discriminators.

$$r_{\boldsymbol{\phi}}(\mathbf{x}) = \frac{p^{*}(\mathbf{x})}{p_{\boldsymbol{\theta}}(\mathbf{x})} = \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} = \frac{p(y=1|\mathbf{x})}{p(y=0|\mathbf{x})} = \frac{\mathcal{D}_{\boldsymbol{\phi}}(\mathbf{x})}{1 - \mathcal{D}_{\boldsymbol{\phi}}(\mathbf{x})}$$

- Implicit/Adversarial variational inference: Implicit models can be used for flexible variational inference (require only samples, no need for densities)
- Adversarial loss: Discriminator provides a mechanism to "learn" what is realistic, this is better than using a (gaussian) likelihood to train generator.

GANs for imitation learning

Use a separate network (discriminator) to "learn" what is realistic Adversarial imitation learning: RL Reward comes from a discriminator



Learning human behaviors from motion capture by adversarial imitation

Josh Merel, Yuval Tassa, Dhruva TB, Sriram Srinivasan, Jay Lemmon, Ziyu Wang, Greg Wayne, Nicolas Heess

DeepMind

Understanding GANs

Lots of other exciting research

- Research
 - Using ideas from convergence of Nash equilibria
 - Connections to RL (actor-critic methods)
 - Control theory (e.g. numerics of GANs)

• Applications

- Class-conditional generation,
- Text-to-image generation
- Image-to-image translation
- Single image super-resolution
- Domain adaptation

And many more ...

Summary

Ways to stabilize GAN training

- Combine with Auto-encoder
- Gradient penalties

Tools developed in GAN literature are intriguing even if you don't care about GANs

- Density ratio trick is useful in other areas (e.g. message passing)
- Implicit variational approximations
- Learn a realistic loss function than use a loss of convenience
- How do we handle non-differentiable simulators?
 - Search using differentiable approximations?

Thanks!

Learning in implicit generative models, Shakir Mohamed* and Balaji Lakshminarayanan*

Variational approaches for auto-encoding generative adversarial networks, Mihaela Rosca*, Balaji Lakshminarayanan*, David Warde-Farley and Shakir Mohamed

Comparison of maximum likelihood and GAN-based training of Real NVPs, Ivo Danihelka, Balaji Lakshminarayanan, Benigno Uria, Daan Wierstra and Peter Dayan

Many paths to equilibrium: GANs do not need to decrease a divergence at every step, William Fedus*, Mihaela Rosca*, Balaji Lakshminarayanan, Andrew Dai, Shakir Mohamed and Ian Goodfellow

Slide credits: Mihaela Rosca, Shakir Mohamed, Ivo Danihelka, David Warde-Farley, Danilo Rezende

Papers available on my webpage <u>http://www.gatsby.ucl.ac.uk/~balaji/</u>

