

Prof. Leal-Taixé and Prof. Niessner Figure from Ian Goodfellow, Tutorial on Generative Adversarial /networks, 2017



Prof. Leal-Taixé and Prof. Niessner Figure from Ian Goodfellow, Tutorial on Generative Adversarial /networks, 2017

2



Generative Adversarial Networks

Prof. Leal-Taixé and Prof. Niessner

Cumulative number of named GAN papers by month



Prof. Leal-Taixé and Prof. Niessner

https://github.com/hindupuravinash/the-gan-zoo 4

Convolution and Deconvolution



Convolution no padding, no stride



https://github.com/vdumoulin/conv_arithmetic

Autoencoder



Reconstruction: Autoencoder



Training Autoencoders





Reconstructed images







Interpolation between two chair models

[Dosovitsky et al. 14] Learning to Generate Chairs

Morphing between chair models









Prof. Leal-Taixé and Prof. Niessner

[Goodfellow et al. 14] GANs (slide McGuinness)



Prof. Leal-Taixé and Prof. Niessner

[Goodfellow et al. 14] GANs (slide McGuinness)



[Goodfellow et al. 14/16] GANs

GANs: Loss Functions

Discriminator loss

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \log D(\boldsymbol{x}) - \frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \log \left(1 - D\left(G(\boldsymbol{z})\right)\right)$$
Generator loss

$$J^{(G)} = -J^{(D)}$$

- Minimax Game:
 - G minimizes probability that D is correct
 - Equilibrium is saddle point of discriminator loss

-> D provides supervision (i.e., gradients) for G

[Goodfellow et al. 14/16] GANs

GANs: Loss Functions

Discriminator loss
$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \log D(\boldsymbol{x}) - \frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \log \left(1 - D\left(G(\boldsymbol{z})\right)\right)$$

Generator loss
$$J^{(G)} = -rac{1}{2} \mathbb{E}_{oldsymbol{z}} \log D\left(G(oldsymbol{z})
ight)$$

- Heuristic Method (often used in practice)
 - G maximizes the log-probability of D being mistaken
 - G can still learn even when D rejects all generator samples

Alternating Gradient Updates

• Step 1: Fix G, and perform gradient step to

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \log D(\boldsymbol{x}) - \frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \log \left(1 - D\left(G(\boldsymbol{z})\right)\right)$$

• Step 2: Fix D, and perform gradient step to

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \log D\left(G(\boldsymbol{z})\right)$$

Vanilla GAN

for number of training iterations do

for k steps do

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$abla_{ heta_d} rac{1}{m} \sum_{i=1}^m \left[\log D\left(oldsymbol{x}^{(i)}
ight) + \log \left(1 - D\left(G\left(oldsymbol{z}^{(i)}
ight)
ight)
ight)
ight].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

$$abla_{ heta_g} rac{1}{m} \sum_{i=1}^m \log\left(1 - D\left(G\left(oldsymbol{z}^{(i)}
ight)
ight)
ight).$$

end for

Prof. Leal-Taixé and Prof. Niessner

19

Training a GAN



https://medium.com/ai-society/gans-from-scratch-1-a-deep-introduction-with-code-in-pytorch-and-tensorflow-cb03cdcdba0f

Prof. Leal-Taixé and Prof. Niessner

GANs: Loss Functions



[Goodfellow et al. 14/16] GANs

DCGAN: Generator



DCGAN: https://github.com/carpedm20/DCGAN-tensorflow



Results on MNIST

Prof. Leal-Taixé and Prof. Niessner

DCGAN: <u>https://github.com/carpedm20/DCGAN-tensorflow</u>



Results on CelebA (200k relatively well aligned portrait photos)

DCGAN: https://github.com/carpedm20/DCGAN-tensorflow



Asian face dataset

Prof. Leal-Taixé and Prof. Niessner

DCGAN: https://github.com/carpedm20/DCGAN-tensorflow





Prof. Leal-Taixé and Prof. Niessner

DCGAN: https://github.com/carpedm20/DCGAN-tensorflow

"Bad" Training Curves



https://stackoverflow.com/questions/44313306/dcgans-discriminator-getting-too-strong-too-quickly-to-allow-generator-to-learn

Prof. Leal-Taixé and Prof. Niessner

"Good" Training Curves







Discriminator's Error through Time

Prof. Leal-Taixé and Prof. Niessner https://medium.com/ai-society/gans-from-scratch-1-a-deep-introduction-with-code-in-pytorch-and-tensorflow-cb03cdcdba0f 29

"Good" Training Curves



Prof. Leal-Taixé and Prof. Niessner https://stackoverflow.com/questions/42690721/how-to-interpret-the-discriminators-loss-and-the-generators-loss-in-generative

30

Training Schedules

• Adaptive schedules

• For instance:

while loss_discriminator > t_d:
 train discriminator
while loss_generator > t_g:
 train generator

Weak vs Strong Discriminator

Need balance 😊

- Discriminator too weak?
 - No good gradients (cannot get better than teacher...)
- Generator too weak?
 - Discriminator will always be right

Mode Collapse

- $\min_{G} \max_{D} V(G,D) \neq \max_{D} \min_{G} V(G,D)$
- *D* in inner loop -> convergence to correct dist.
- G in inner loop -> easy to convergence to one sample



Prof. Leal-Taixé and Prof. Niessner

[Metz et al. 16] 33

Mode Collapse

• Data dim. Fixed (512)

 Performance correlates with # of modes



-> More modes, smaller recovery rate! -> part of the reason, why we often see GAN-results on specific domains (e.g., faces) Prof. Leal-Taixé and Prof. Niessner

Slide credit Ming-Yu Liu 34

Mode Collapse



Mode recovery vs manifold dimension

 Performance correlates with dim of manifold

-> Larger latent space, more mode collapse

Prof. Leal-Taixé and Prof. Niessner

Problems with Global Structure





(Goodfellow 2016)

Prof. Leal-Taixé and Prof. Niessner
Problems with Counting













(Goodfellow 2016)

Prof. Leal-Taixé and Prof. Niessner

- Main difficulty of GANs: we don't know how good they are
- People cherry pick results in papers -> some of them will always look good, but how to quantify?
- Do we only memorize or do we generalize?
- GANs are difficult to evaluate! [This et al., ICLR 2016]

Human evaluation:

- Every n updates, show a series of predictions
- Check train curves
- What does 'look good' mean at the beginning?
 - Need variety!
 - But don' t have 'realistic' predictions yet...
- If it doesn' t look good? Go back, try different hyperparameters...

Inception Score (IS)

- Measures saliency and diversity

- Train an accurate classifier
- Train a image generation model (conditional)
- Check how accurate the classifier can recognize the generated images
- Makes some assumptions about data distributions...

Inception Score (IS)

- Saliency: check whether the generated images can be classified with high confidence (i.e., high scores only on a single class)

- Diversity: check whether we obtain samples from all classes

What if we only have one good image per class?

- Could also look at discriminator
 - If we end up with a strong discriminator, then generator must also be good
 - Use D features, for classification network
 - Only fine-tune last layer
 - If high class accuracy -> we have a good D and G

Prof. Leal-Taixé and Prof. Niessner Caveat: not sure if people do this... Couldn' t find paper 43

Next: Making GANs Work in Practice

• Training / Hyperparameters (most important)

• Choice of loss function

• Choice of architecture

GAN Hacks: Normalize Inputs

• Normalize the inputs between -1 and 1

• Tanh as the last layer of the generator output

• No-brainer 😊

45

GAN Hacks: Sampling

- Use a spherical z
- Don' t sample from a uniform distribution
- Sample from a Gaussian Distribution



• When doing interpolations, do the interpolation via a great circle, rather than a straight line from point A to point B

• Tom White's <u>Sampling Generative</u> <u>Networks</u> ref code <u>https://github.com/dribnet/plat</u> has more details

GAN Hacks: BatchNorm

• Use Batch Norm

 Construct different minibatches for real and fake, i.e. each minibatch needs to contain only all real images or all generated images.



GAN Hacks: Use ADAM

• See Adam usage [Radford et al. 15]

• SGD for discriminator

• ADAM for generator

GAN Hacks: One-sided Label Smoothing

• Prevent discriminator from giving too large gradient signal to generator:

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \log D(\boldsymbol{x}) - \frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \log \left(1 - D\left(G(\boldsymbol{z})\right)\right)$$

Some value smaller than 1; e.g.,0.9

-> reduces confidence; i.e., makes disc. 'weaker'-> encourages 'extreme samples' (prevents extrapolating)

Prof. Leal-Taixé and Prof. Niessner

GAN Hacks: Historical Generator Batches



Help stabilize discriminator training in early stage

Prof. Leal-Taixé and Prof. Niessner

Srivastava et al. 17 "Learning from Simulated and Unsupervised Images through Adversarial Training" 50

GAN Hacks: Avoid Sparse Gradients

- Stability of GAN game suffers if gradients are sparse
- LeakyReLU -> good in both G and D
- Downsample -> use average pool, conv+stride
- Upsample -> deconv+stride, PixelShuffle



Exponential Averaging of Weights

• Problem: discriminator is noisy due to SGD

- Rather than taking final result of a GAN, would be biased on last latest iterations (i.e., latest training samples),
 - -> exponential average of weights

-> keep second 'vector' of weights that are averaged

-> almost no cost, average of weights from last n iters

New Objective Functions

New Objective Functions

"heuristic is standard..."

EBGAN: "Energy-based Generative Adversarial Networks"

- BEGAN: "Boundary Equilibrium GAN"
- WGAN: "Wasserstein Generative Adversarial Networks" LSGAN: "Least Squares Generative Adversarial Networks"

Prof. Leal-Taixe and Prof. Niessner

- Discriminator is AE (Energy-based GAN)
- a good autoencoder: we want the reconstruction cost D(x) for real images to be low.
- a good critic: we want to penalize the discriminator if the reconstruction error for generated images drops below a value m. $\mathcal{L}_D(x,z) = D(x) + [m - D(G(z))]^+$

$$\mathcal{L}_D(x,z) = D(x) + [m - D(G(z))]^+$$
$$\mathcal{L}_G(z) = D(G(z))$$

55

$$D(x) = ||Dec(Enc(x)) - x||$$

where
$$[u]^+ = max(0, u)$$

Prof. Leal-Taixé and Prof. Niessner

https://medium.com/@jonathan_hui/gan-energy-based-gan-ebgan-boundary-equilibrium-gan-began-4662cceb7824

• Similar to EBGAN

 Instead of reconstruction loss, measure difference in data distribution of real and generated images



56

Prof. Leal-Taixé and Prof. Niessner

https://medium.com/@jonathan_hui/gan-energy-based-gan-ebgan-boundary-equilibrium-gan-began-4662cceb7824

• Earth Mover Distance / Wasserstein Distance



Minimum amount of work to move earth from p(x) to q(x)

Prof. Leal-Taixé and Prof. Niessner

https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490

• Formulate EMD via it's dual:

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \le 1} \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f(x)]$$

$$|f(x_1) - f(x_2)| \leq |x_1 - x_2|.$$

1-Lipschitz function: upper bound between densities

Prof. Leal-Taixé and Prof. Niessner

https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490

$$|f(x_1)-f(x_2)|\leq |x_1-x_2|.$$

f is a critic function, defined by a neural network -> f needs to be 1-Lipschitz; WGAN restricts max weight value in f; weights of the discriminator must be within a certain range controlled by hyperparameters c

$$w \leftarrow w + \alpha \cdot \operatorname{RMSProp}(w, g_w)$$

 $w \leftarrow \operatorname{clip}(w, -c, c)$



Prof. Leal-Taixé and Prof. Niessner

https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490

Discriminator/Critic

Generator

$$\begin{aligned} \mathbf{GAN} & \nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D\left(\boldsymbol{x}^{(i)} \right) + \log \left(1 - D\left(G\left(\boldsymbol{z}^{(i)} \right) \right) \right) \right] & \nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m -\log \left(D\left(G\left(\boldsymbol{z}^{(i)} \right) \right) \right) \\ \mathbf{WGAN} & \nabla_w \frac{1}{m} \sum_{i=1}^m \left[f\left(\boldsymbol{x}^{(i)} \right) - f\left(G\left(\boldsymbol{z}^{(i)} \right) \right) \right] & \nabla_\theta \frac{1}{m} \sum_{i=1}^m -f\left(G\left(\boldsymbol{z}^{(i)} \right) \right) \end{aligned}$$

Prof. Leal-Taixé and Prof. Niessner

https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490

Algorithm 1 WGAN, our proposed algorithm. All experiments in the paper used the default values $\alpha = 0.00005$, c = 0.01, m = 64, $n_{\text{critic}} = 5$.

Require: : α , the learning rate. c, the clipping parameter. m, the batch size. n_{critic} , the number of iterations of the critic per generator iteration. **Require:** : w_0 , initial critic parameters. θ_0 , initial generator's parameters.

1: while θ has not converged **do**

2: for
$$t = 0, ..., n_{\text{critic}}$$
 do

3: Sample
$$\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$$
 a batch from the real data

4: Sample
$$\{z^{(i)}\}_{i=1}^m \sim p(z)$$
 a batch of prior samples.

5:
$$g_w \leftarrow \nabla_w \left[\frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$$

$$b: \qquad w \leftarrow w + \alpha \cdot \mathrm{RMSProp}(w, g_w)$$

7:
$$w \leftarrow \operatorname{clip}(w, -c, c)$$

8: end for

9: Sample
$$\{z^{(i)}\}_{i=1}^m \sim p(z)$$
 a batch of prior samples.

10:
$$g_{\theta} \leftarrow -\nabla_{\theta} \frac{1}{m} \sum_{i=1}^{m} f_w(g_{\theta}(z^{(i)}))$$

11:
$$\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_{\theta})$$

Prof. Leal-Tai: 12: end while



Prof. Leal-Taixé and Prof. Niessner

https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490



Prof. Leal-Taixé and Prof. Niessner

https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490

- + mitigates mode collapse
- + generator still learns when critic performs well
- + actual convergence
- Enforcing Lipschitz constraint is difficult
- Weight clipping is "terrible"
 - -> too high: takes long time to reach limit limit; slow training
 - -> too small: vanishing gradients when layersi s big

Prof. Leal-Taixé and Prof. Niessner

GAN Losses

• Many more variations!!!

• High-level understanding: "loss" is a meta loss to train the actual loss (i.e., D) to provide gradients for G

 Always start simple: if things don' t converge, don' t randomly shuffle loss around; always try easy things first (AE, VAE, 'simple heuristic' GAN)

GAN Architectures

Multiscale GANs



Multiscale GANs



Denton et al, NIPS 2015

Prof. Leal-Taixé and Prof. Niessner

69 Credit: Li/Karpathy/Johnson

Progressive Growing GANs





https://github.com/tkarras/progressive growing of gans [Karras et al. 17]






































5 min 00 sec



Fixed resolution

Progressive growing

Progressive Growing GANs

CelebA-HQ 1024 × 1024

Latent space interpolations

https://github.com/tkarras/progressive growing of gans [Karras et al. 17]

Lots of GAN Variations

- Hundreds of GAN papers in the last two years
 - > Mostly with different losses
 - > Extremely hard to train and evaluate

Are GANs Created Equal? A Large-Scale Study

Mario Lucic* Karol Kurach* Marcin Michalski Sylvain Gelly Olivier Bousquet Google Brain

Abstract

Generative adversarial networks (GAN) are a powerful subclass of generative models. Despite a very rich research

GAN algorithm(s) perform objectively better than the others. That's partially due to the lack of robust and consistent metric, as well as limited comparisons which put all algorithms on equal footage, including the computational

Next lectures

- Next Monday 17th, more on Generative models
 Conditional GANs (cGANs)!
- We are still working on feedback for presentations will send around asap.

• Keep working on the projects!