

# Deep Generative Models

## 8. Generative Adversarial Networks



- 국가수리과학연구소 산업수학혁신센터 김민중

# Recap

- Model Families

- Autoregressive Models

$$p_{\theta}(\mathbf{x}) = \prod_{i=1}^d p_{\theta}(x_i | \mathbf{x}_{<i})$$

- Variational Autoencoders

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

- Normalizing Flow Models

$$p_X(\mathbf{x}; \theta) = p_Z\left(\mathbf{f}_{\theta}^{-1}(\mathbf{x})\right) \left| \det \left( \frac{\partial \mathbf{f}_{\theta}^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right) \right|$$

---

## Recap

- All the above families are trained by minimizing **KL divergence**  $D(p_{data} \parallel p_{\theta})$  or equivalently maximizing **likelihoods** (or approximations)

# Why maximum likelihood?

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{i=1}^M \log p_{\theta}(\mathbf{x}^{(i)}), \quad \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(M)} \sim p_{data}$$

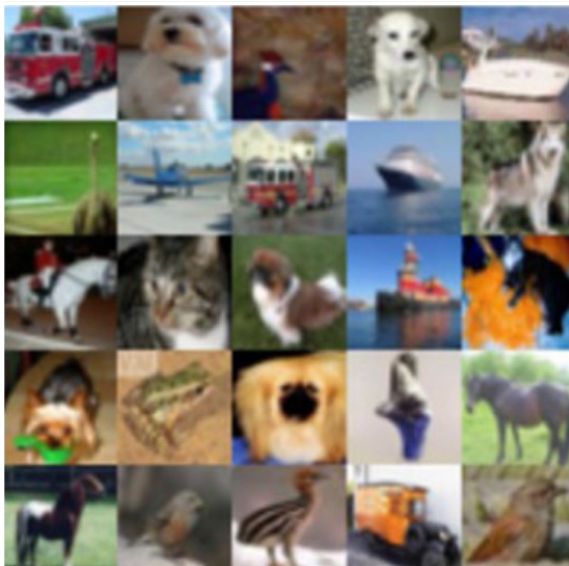
- Optimal statistical efficiency
  - Assume sufficient model capacity, such that there exists a unique  $\theta^* \in \mathcal{M}$  that satisfies  $p_{\theta^*} = p_{data}$
  - The convergence of  $\hat{\theta}$  to  $\theta^*$  when  $M \rightarrow \infty$  is the “fastest” among all statistical methods when using maximum likelihood training
- Higher likelihood = better lossless compression
- Is the likelihood a good indicator of the quality of samples generated by the model?

# Recap

- Model Families
  - Autoregressive Models:  $p_{\theta}(\mathbf{x}) = \prod_{i=1}^d p_{\theta}(x_i | \mathbf{x}_{<i})$
  - Variational Autoencoders:  $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$
  - Normalizing Flow Models:  $p_X(\mathbf{x}; \theta) = p_Z\left(\mathbf{f}_{\theta}^{-1}(\mathbf{x})\right) \left| \det\left(\frac{\partial \mathbf{f}_{\theta}^{-1}(\mathbf{x})}{\partial \mathbf{x}}\right) \right|$
- All the above families are trained by minimizing KL divergence  $D(p_{data} \parallel p_{\theta})$  or equivalently maximizing likelihoods (or approximations)
- Today: alternative for  $D(p_{data} \parallel p_{\theta})$

# Comparing distributions via samples

- Given a **finite set of samples** from two distributions  $S_1 = \{x \sim P\}$  and  $S_2 = \{x \sim Q\}$ , how can we tell if these samples are from the same distribution? (i.e.,  $P = Q$ ?)



$$S_1 = \{x \sim P\}$$

vs.



$$S_2 = \{x \sim Q\}$$

---

## Two-sample tests

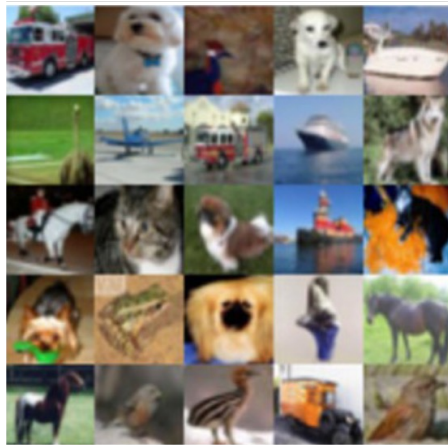
- $S_1 = \{x \sim P\}$  and  $S_2 = \{x \sim Q\}$
- Test statistic  $T$  compares  $S_1$  and  $S_2$ . Using  $T$ , determine  $P = Q$  or not
  - E.g.,

$$T(S_1, S_2) = \left| \frac{1}{|S_1|} \sum_{x \in S_1} x - \frac{1}{|S_2|} \sum_{x \in S_2} x \right|$$

- If  $T$  is large enough, then we determine  $P \neq Q$  otherwise we say  $P = Q$
- Key observation: Test statistic is likelihood-free since it does not involve the densities  $P$  or  $Q$  (only samples)

# Generative modeling and two-sample tests

- A priori, we assume direct access to  $S_1 = D = \{x \sim p_{data}\}$
- In addition, we have a model distribution  $p_\theta$
- Assume that the model distribution permits efficient sampling. Let  $S_2 = \{x \sim p_\theta\}$
- Alternative notion of distance between distributions:
  - Train the generative model to minimize a two-sample test objective between  $S_1$  and  $S_2$


$$S_1 = \{x \sim p_{data}\}$$

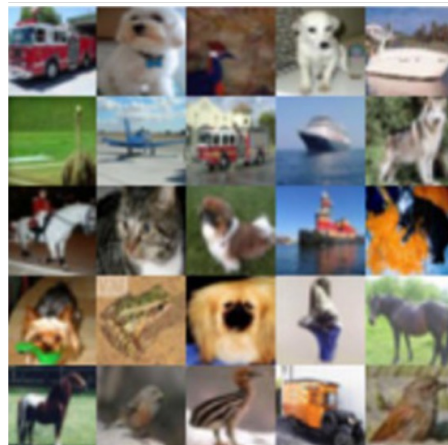
**VS.**


$$S_2 = \{x \sim p_\theta\}$$



# Two-sample test

- In the generative model setup, we know that  $S_1$  and  $S_2$  come from different distributions  $p_{data}$  and  $p_{\theta}$  respectively
- **Key idea:** Learn a statistic to automatically identify in what way the two sets of samples  $S_1$  and  $S_2$  differ from each other
- How? Train a classifier (called a discriminator)!



$$S_1 = \{x \sim p_{data}\}$$

vs.



$$S_2 = \{x \sim p_{\theta}\}$$

---

# Two-sample test via a discriminator

- Any **binary classifier**  $D_\phi$  (e.g., neural network) which tries to distinguish “real” ( $y = 1$ ) samples from the dataset and “fake” ( $y = 0$ ) samples generated from the model
- Test statistic:  $-\text{loss}$  of the classifier.
  - **Low loss**, real and fake samples are **easy to distinguish** (different)
  - **High loss**, real and fake samples are **hard to distinguish** (similar)
- **Goal**
  - Maximize the two-sample test statistic (in support of the alternative hypothesis  $p_{data} \neq p_\theta$ ), or equivalently minimize classification loss

# Two-sample test via a discriminator

- Training objective for discriminator

$$\begin{aligned}\max_{D_\phi} V(p_\theta, D_\phi) &= \max_{D_\phi} E_{x \sim p_{data}} [\log D_\phi(x)] + E_{x \sim p_\theta} [\log (1 - D_\phi(x))] \\ &\approx \max_{D_\phi} \sum_{x \in S_1} \log D_\phi(x) + \sum_{x \in S_2} \log (1 - D_\phi(x))\end{aligned}$$

- For a fixed generative model  $p_\theta$ , the discriminator is performing binary classification with the cross-entropy objective
  - Assign probability 1 to true data points  $x \sim p_{data}$  (in set  $S_1$ )
  - Assign probability 0 to fake samples  $x \sim p_\theta$  (in set  $S_2$ )

# Two-sample test via a discriminator

- Training objective for discriminator

$$\begin{aligned}\max_{D_\phi} V(p_\theta, D_\phi) &= \max_{D_\phi} E_{x \sim p_{data}} [\log D_\phi(x)] + E_{x \sim p_\theta} [\log (1 - D_\phi(x))] \\ &\approx \max_{D_\phi} \sum_{x \in S_1} \log D_\phi(x) + \sum_{x \in S_2} \log (1 - D_\phi(x))\end{aligned}$$

- For a fixed generative model  $p_\theta$ , the optimal discriminator is given by

$$D_\theta^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_\theta(x)}$$

- If  $p_\theta = p_{data}$ , classifier cannot do better than chance ( $D_\theta^*(x) = 1/2$ )

---

# Generative Adversarial Networks

- A two-player minimax game between a generator and a discriminator
- **Generator**
  - Directed latent variable model with a deterministic mapping between  $\mathbf{z}$  and  $\mathbf{x}$  given by  $G_\theta$ 
    - Sample  $\mathbf{z} \sim p_Z$ , where  $p_Z$  is a simple prior, e.g., Gaussian
    - Set  $\mathbf{x} = G_\theta(\mathbf{z})$
  - Like a flow model, but mapping  $G_\theta$  need not be invertible
  - Distribution over  $p_\theta(\mathbf{x})$  over  $\mathbf{x}$  is implicitly defined (no likelihood!)
  - Minimizes a two-sample test objective (in support of the null hypothesis  $p_{data} = p_\theta$ )

# Example of GAN objective

- Training objective for generator

$$\min_{\mathbf{G}} \max_{\mathbf{D}} V(\mathbf{G}, \mathbf{D}) = \min_{\mathbf{G}} \max_{\mathbf{D}} E_{\mathbf{x} \sim p_{data}} [\log \mathbf{D}(\mathbf{x})] + E_{\mathbf{x} \sim p_{\mathbf{G}}} [\log(1 - \mathbf{D}(\mathbf{x}))]$$

- For the optimal discriminator  $D_G^*(\cdot)$ , we have

$$\begin{aligned} V(\mathbf{G}, D_G^*) &= E_{\mathbf{x} \sim p_{data}} \left[ \log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \right] + E_{\mathbf{x} \sim p_G} \left[ \log \frac{p_G(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \right] \\ &= E_{\mathbf{x} \sim p_{data}} \left[ \log \frac{p_{data}(\mathbf{x})}{\frac{p_{data}(\mathbf{x}) + p_G(\mathbf{x})}{2}} \right] + E_{\mathbf{x} \sim p_G} \left[ \log \frac{p_G(\mathbf{x})}{\frac{p_{data}(\mathbf{x}) + p_G(\mathbf{x})}{2}} \right] - \log 4 \\ &= D \left( p_{data} \parallel \frac{p_{data} + p_G}{2} \right) + D \left( p_G \parallel \frac{p_{data} + p_G}{2} \right) - \log 4 \\ &= 2JSD(p_{data} \parallel p_G) - \log 4 \end{aligned}$$

# Jensen-Shannon Divergence

- Also called as the symmetric KL divergence

$$JSD(p \parallel q) = \frac{1}{2} D \left( p \parallel \frac{p+q}{2} \right) + \frac{1}{2} D \left( q \parallel \frac{p+q}{2} \right)$$

- Properties

- $JSD(p \parallel q) \geq 0$
- $JSD(p \parallel q) = 0$  iff  $p = q$
- $JSD(p \parallel q) = JSD(q \parallel p)$
- $\sqrt{JSD(p \parallel q)}$  satisfies triangle inequality. I.e., it is a distance

- Optimal generator for the JSD/Negative Cross Entropy GAN

$$p_G = p_{data}$$

- For the optimal discriminator  $D_{G^*}^*(\cdot)$  and generator  $G^*(\cdot)$ , we have

$$V(G^*, D_{G^*}^*) = -\log 4$$

---

# Recap of GANs

- Choose  $d(p_{data}, p_{\theta})$  to be a two-sample test statistic
  - Learn the statistic by training a classifier (discriminator)
  - Under ideal conditions, equivalent to choosing  $d(p_{data}, p_{\theta})$  to be  $JSD(p_{data} \parallel p_{\theta})$
- Generator  $G_{\theta}$  (e.g., neural network) is a mapping that generates  $\mathbf{x}$  from the latent variable  $\mathbf{z}$  and is trained to make it difficult for the classifier to distinguish
- Pros:
  - Loss only requires samples from  $p_{\theta}$ . No likelihood needed!
  - Lots of flexibility for the neural network architecture, any  $G_{\theta}$  defines a valid sampling procedure
  - Fast sampling (single forward pass)
- Cons: very difficult to train in practice



# The GAN training algorithm

- Sample minibatch of  $n$  training points  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$  from  $p_{data}$
- Sample minibatch of  $n$  noise vectors  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(n)}$  from  $p_Z$
- Update the discriminator parameters  $\phi$  by stochastic gradient ascent

$$\nabla_{\phi} V(G_{\theta}, D_{\phi}) = \frac{1}{n} \nabla_{\phi} \sum_{i=1}^n \left[ \log D_{\phi}(\mathbf{x}^{(i)}) + \log \left( 1 - D_{\phi}(G_{\theta}(\mathbf{z}^{(i)})) \right) \right]$$

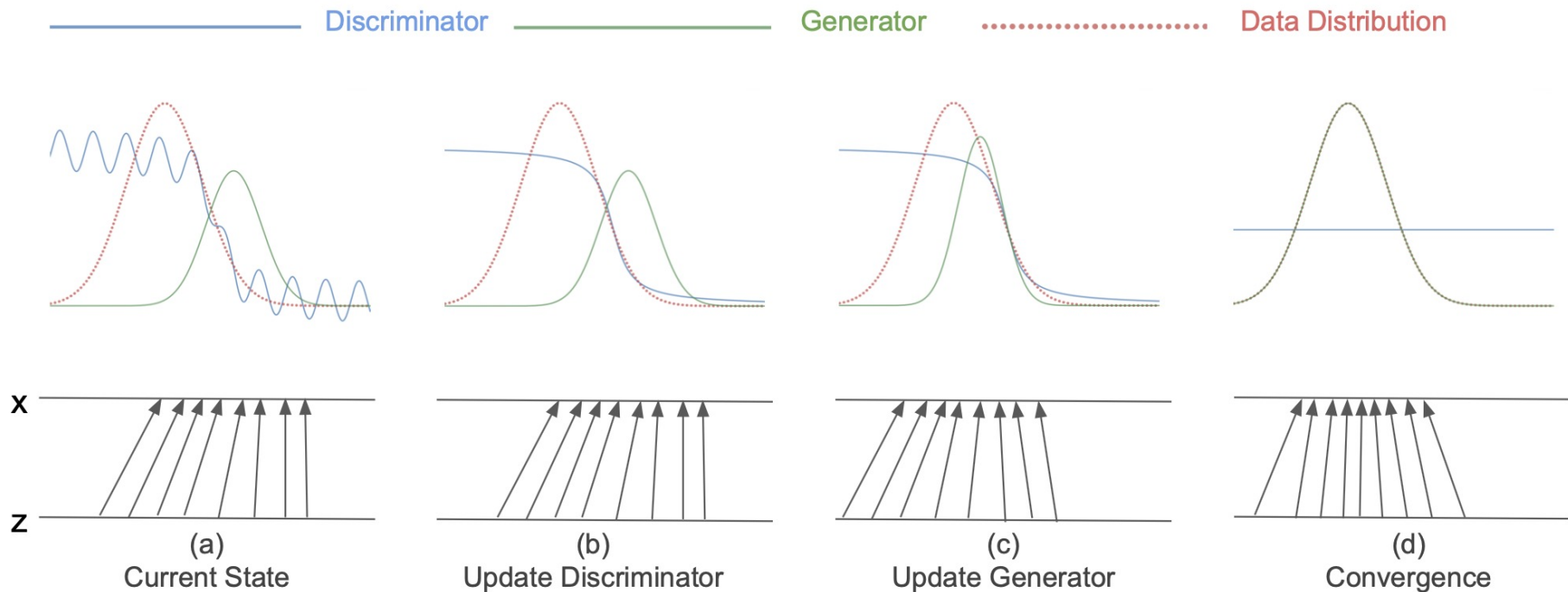
- Update the generator parameters  $\theta$  by stochastic gradient descent

$$\nabla_{\theta} V(G_{\theta}, D_{\phi}) = \frac{1}{n} \nabla_{\theta} \sum_{i=1}^n \log \left( 1 - D_{\phi}(G_{\theta}(\mathbf{z}^{(i)})) \right)$$

- Repeat for fixed number of epochs

# Alternating optimization in GANs

$$\min_{\theta} \max_{\phi} V(G_{\theta}, D_{\phi}) = E_{x \sim p_{data}} [\log D_{\phi}(x)] + E_{z \sim p_z} [\log (1 - D_{\phi}(G_{\theta}(z)))]$$



# Frontiers in GAN research

- GANs have been successfully applied to several domains and tasks
- However, working with GANs can be very challenging in practice
  - Unstable optimization
  - Mode collapse Evaluation
  - Bag of tricks needed to train GANs successfully



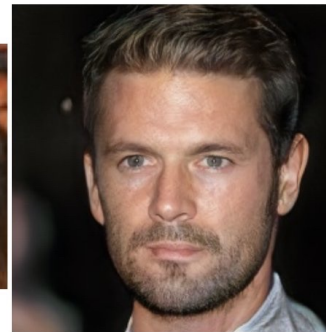
2014



2015



2016



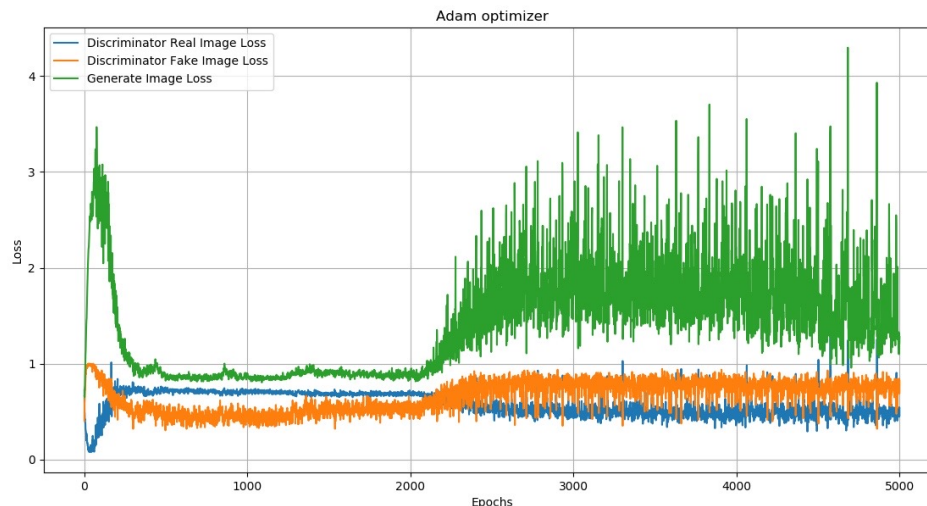
2017



2018

# Optimization challenges

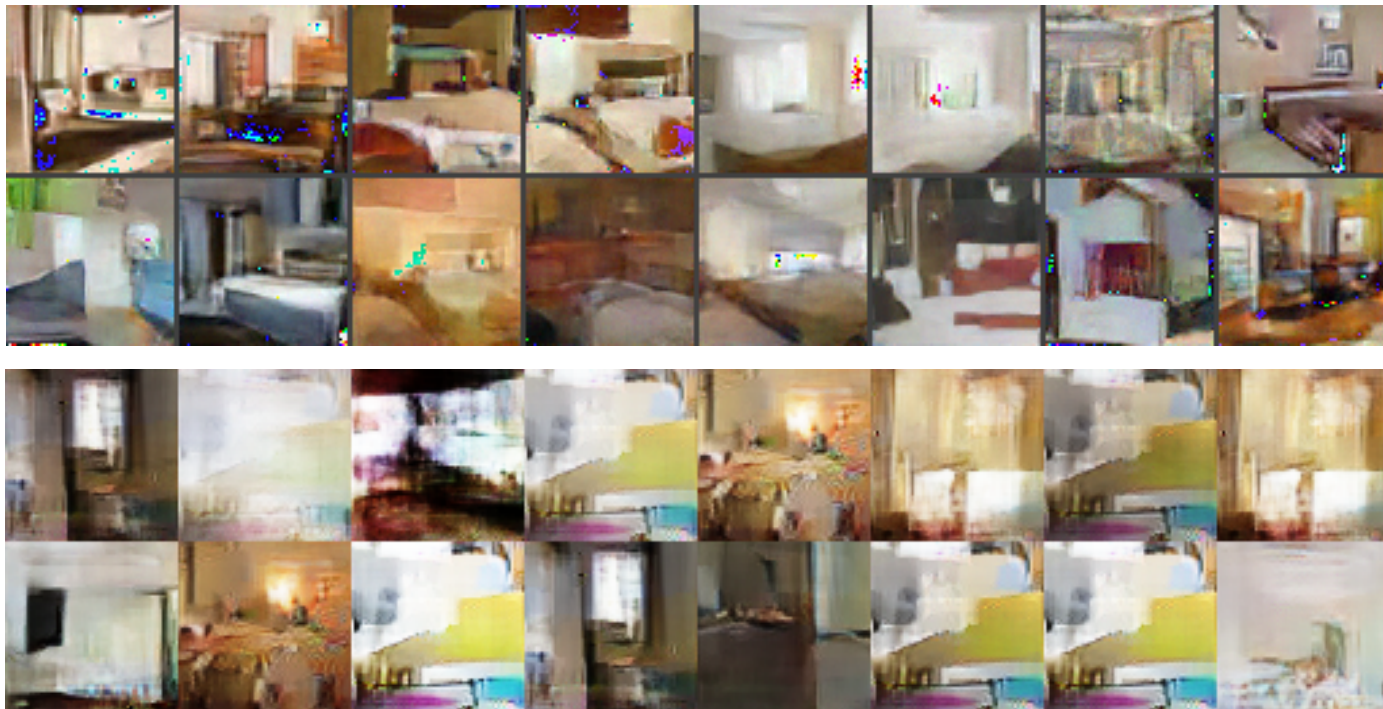
- **Theorem (informal):** If the generator updates are made in function space and discriminator is optimal at every step, then the generator is guaranteed to converge to the data distribution
- Unrealistic assumptions!
- In practice, the generator and discriminator loss keeps oscillating during GAN training
- No robust stopping criteria in practice (unlike MLE)



Source: Mirantha Jayathilaka

# Mode Collapse

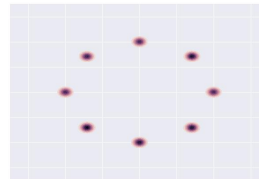
- GANs are notorious for suffering from mode collapse
- Intuitively, this refers to the phenomena where the generator of a GAN collapses to one or few samples (dubbed as “modes”)



Source: Arjovsky et al., 2017

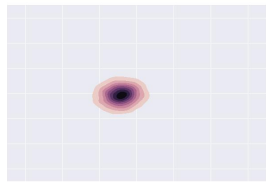
# Mode Collapse

- True distribution is a mixture of Gaussians

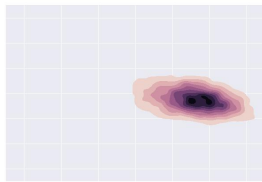


Target

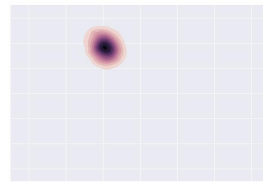
Source: Metz et al., 2017



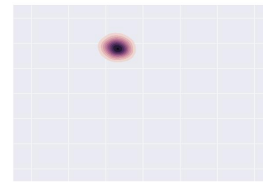
Step 0



Step 5k



Step 10k



Step 15k



Step 20k



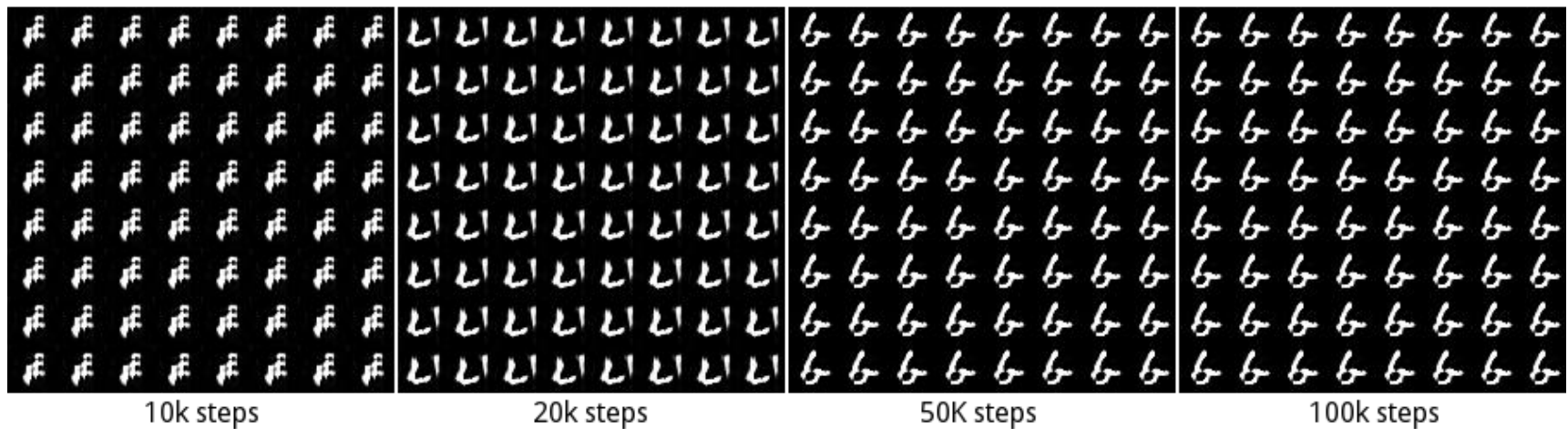
Step 25k

- The generator distribution keeps oscillating between different modes

# Mode Collapse

- Fixes to mode collapse are mostly empirically driven alternative architectures, alternative GAN loss, adding regularization terms, etc.
- How to Train a GAN? Tips and tricks to make GANs work by Soumith Chintala
  - <https://github.com/soumith/ganhacks>

Source: Metz et al., 2017





# Recap

- Likelihood-free training
- Training objective for GANs

$$\min_G \max_D V(G, D) = E_{x \sim p_{data}} [\log D(x)] + E_{x \sim p_G} [\log(1 - D(x))]$$

- With the optimal discriminator  $D_G^*$ , we see GAN minimizes a scaled and shifted Jensen-Shannon divergence

$$\min_G 2JSD(p_{data} \parallel p_G) - \log 4$$

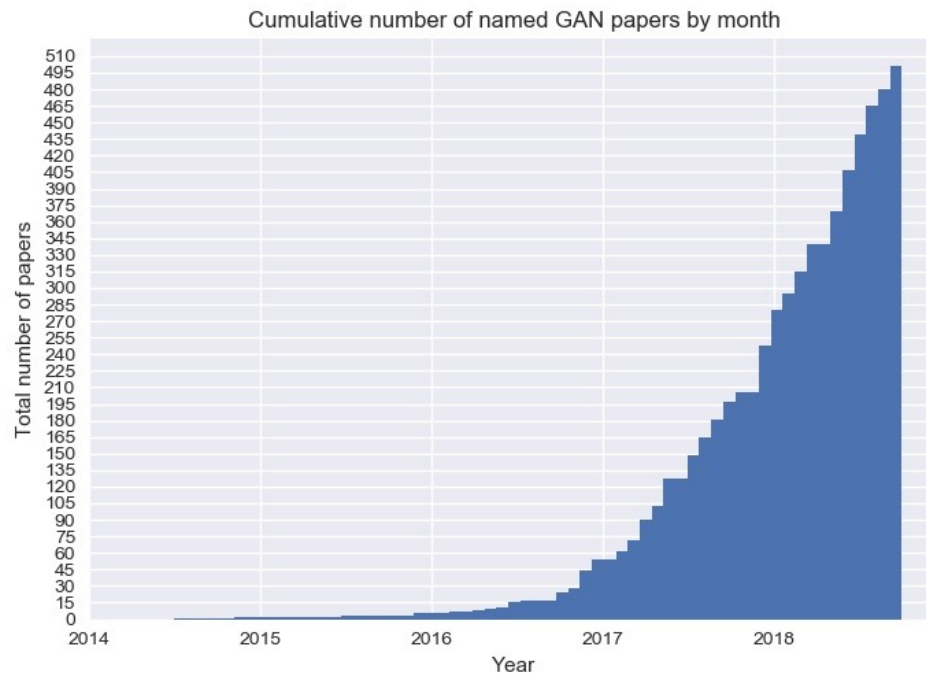
- Parameterize  $D$  by  $\phi$  and  $G$  by  $\theta$
- Prior distribution  $p_Z$

$$\min_{\theta} \max_{\phi} E_{x \sim p_{data}} [\log D_{\phi}(x)] + E_{z \sim p_Z} [\log(1 - D_{\phi}(G_{\theta}(z)))]$$



# GAN Zoo

- GAN Zoo: List of all named GANs
  - <https://github.com/hindupuravinash/the-gan-zoo>



---

# Beyond KL and Jensen-Shannon Divergence

- What choices do we have for  $d(\cdot)$ ?
  - KL divergence: Autoregressive Models, Flow models
  - (scaled and shifted) Jensen-Shannon divergence (approximately): original GAN objective

# $f$ -divergences

- What choices do we have for  $d(\cdot)$ ?
- Given two densities  $p$  and  $q$ , the  $f$ -divergence is given by

$$D_f(p, q) = E_{x \sim q} \left[ f \left( \frac{p(x)}{q(x)} \right) \right]$$

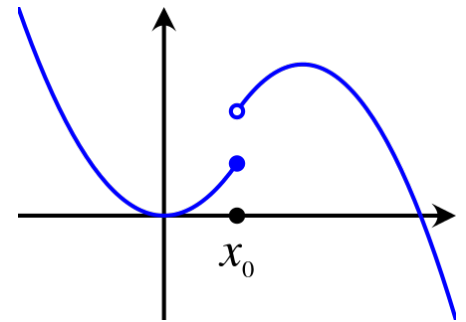
- Where  $f$  is any convex, lower-semicontinuous function with  $f(1) = 0$
- Convex: Line joining any two points lies above the function
- Lower-semicontinuous

$$\liminf_{x \rightarrow x_0} f(x) \geq f(x_0)$$

- for any point  $x_0$
- Jensen's inequality

$$E_{x \sim q} \left[ f \left( \frac{p(x)}{q(x)} \right) \right] \geq f \left( E_{x \sim q} \left[ \frac{p(x)}{q(x)} \right] \right) = f \left( \int p(x) dx \right) = f(1) = 0$$

- Example: KL divergence with  $f(u) = u \log u$



# *f*-divergences

Name	$D_f(P  Q)$	Generator $f(u)$
Total variation	$\frac{1}{2} \int  p(x) - q(x)  \, dx$	$\frac{1}{2} u - 1 $
Kullback-Leibler	$\int p(x) \log \frac{p(x)}{q(x)} \, dx$	$u \log u$
Reverse Kullback-Leibler	$\int q(x) \log \frac{q(x)}{p(x)} \, dx$	$-\log u$
Pearson $\chi^2$	$\int \frac{(q(x)-p(x))^2}{p(x)} \, dx$	$(u - 1)^2$
Neyman $\chi^2$	$\int \frac{(p(x)-q(x))^2}{q(x)} \, dx$	$\frac{(1-u)^2}{u}$
Squared Hellinger	$\int \left( \sqrt{p(x)} - \sqrt{q(x)} \right)^2 \, dx$	$(\sqrt{u} - 1)^2$
Jeffrey	$\int (p(x) - q(x)) \log \left( \frac{p(x)}{q(x)} \right) \, dx$	$(u - 1) \log u$
Jensen-Shannon	$\frac{1}{2} \int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} \, dx$	$-(u + 1) \log \frac{1+u}{2} + u \log u$
Jensen-Shannon-weighted	$\int p(x) \pi \log \frac{p(x)}{\pi p(x) + (1-\pi)q(x)} + (1 - \pi)q(x) \log \frac{q(x)}{\pi p(x) + (1-\pi)q(x)} \, dx$	$\pi u \log u - (1 - \pi + \pi u) \log(1 - \pi + \pi u)$
GAN	$\int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} \, dx - \log(4)$	$u \log u - (u + 1) \log(u + 1)$
$\alpha$ -divergence ( $\alpha \notin \{0, 1\}$ )	$\frac{1}{\alpha(\alpha-1)} \int \left( p(x) \left[ \left( \frac{q(x)}{p(x)} \right)^\alpha - 1 \right] - \alpha(q(x) - p(x)) \right) \, dx$	$\frac{1}{\alpha(\alpha-1)} (u^\alpha - 1 - \alpha(u - 1))$

Source: Nowozin et al., 2017

# Thanks

---