

Generative Adversarial Networks

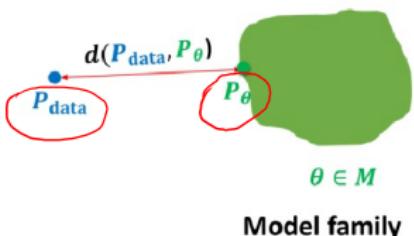
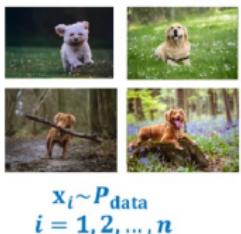
Stefano Ermon

Stanford University

Lecture 9

Recap

$$\hat{\theta}_{ML} = \arg \min_{\theta} KL(P_{\text{data}} \| P_{\theta}) = \arg \max_{\theta} E_{x \sim p_{\text{data}}} [\log P_{\theta}(x)]$$



$$\simeq \arg \max_{\theta} \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(x_i)$$

- Model families

- Autoregressive Models: $p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_{<i})$
- Variational Autoencoders: $p_{\theta}(x) = \int p_{\theta}(x, z) dz$
- Normalizing Flow Models: $p_{\theta}(x; \theta) = p_Z(f_{\theta}^{-1}(x)) \left| \det \left(\frac{\partial f_{\theta}^{-1}(x)}{\partial x} \right) \right|$
- All the above families are trained by minimizing KL divergence $D_{KL}(p_{\text{data}} \| p_{\theta})$, or equivalently maximizing likelihoods (or approximations)

Why maximum likelihood?

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{i=1}^M \log p_{\theta}(\mathbf{x}_i), \quad \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M \sim p_{\text{data}}(\mathbf{x})$$

- **Optimal statistical efficiency.**
 - Assume sufficient model capacity, such that there exists a unique $\theta^* \in \mathcal{M}$ that satisfies $p_{\theta^*} = p_{\text{data}}$.
 - The convergence of $\hat{\theta}$ to θ^* 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?

Towards likelihood-free learning

- **Case 1:** Optimal generative model will give best **sample quality** and highest test **log-likelihood**
- For imperfect models, achieving high log-likelihoods might not always imply good sample quality, and vice-versa (Theis et al., 2016)

$$\underbrace{E_{P_{\text{data}}} [\log P_{\text{data}}]} - \log 100 \leq \underbrace{\log_{\text{likelihood}}}_{\text{log-likelihood}} \leq \underbrace{E_{P_{\text{data}}} [\log P_{\text{data}}]}$$

Towards likelihood-free learning

- **Case 2:** Great test log-likelihoods, poor samples. E.g., For a discrete noise mixture model $p_\theta(\mathbf{x}) = 0.01p_{\text{data}}(\mathbf{x}) + 0.99p_{\text{noise}}(\mathbf{x})$
 - 99% of the samples are just noise (most samples are poor)
 - Taking logs, we get a lower bound

$$\begin{aligned}\log p_\theta(\mathbf{x}) &= \log[0.01p_{\text{data}}(\mathbf{x}) + 0.99p_{\text{noise}}(\mathbf{x})] \\ &\geq \log 0.01p_{\text{data}}(\mathbf{x}) = \log p_{\text{data}}(\mathbf{x}) - \log 100\end{aligned}$$

- For expected log-likelihoods, we know that

- Lower bound *likelihood*

$$\rightarrow E_{p_{\text{data}}} [\log p_\theta(\mathbf{x})] \geq E_{p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x})] - \log 100$$

- Upper bound (via non-negativity of $D_{KL}(p_{\text{data}} \| p_\theta) \geq 0$)

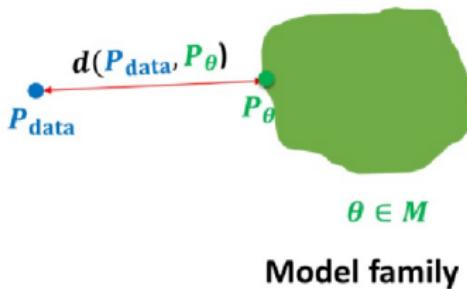
$$E_{p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x})] \geq E_{p_{\text{data}}} [\log p_\theta(\mathbf{x})] \quad \text{log-likelihood}$$

- As we increase the dimension n of $\mathbf{x} = (x_1, \dots, x_n)$, absolute value of $\log p_{\text{data}}(\mathbf{x}) = \sum_{i=1}^n \log p_\theta(x_i | \mathbf{x}_{<i})$ increases proportionally to n but $\log 100$ remains constant. Hence, likelihoods are great
 $E_{p_{\text{data}}} [\log p_\theta(\mathbf{x})] \approx E_{p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x})]$ in very high dimensions

Towards likelihood-free learning

- **Case 3:** Great samples, poor test log-likelihoods. E.g., Memorizing training set
 - • Samples look exactly like the training set (cannot do better!)
 - • Test set will have zero probability assigned (cannot do worse!)
- The above cases suggest that it might be useful to disentangle likelihoods and sample quality
- **Likelihood-free learning** consider alternative training objectives that do not depend directly on a likelihood function

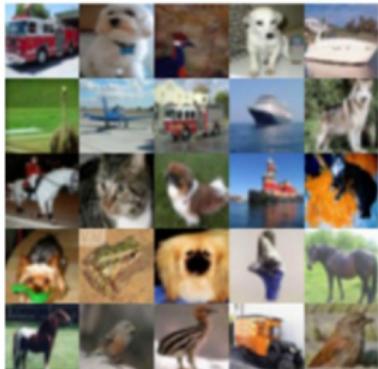
Recap



$$\mathbf{x}_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$

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

Comparing distributions via samples



vs.



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



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



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

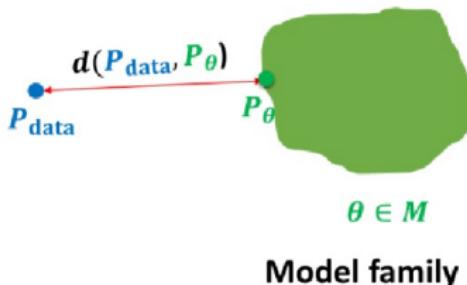
Two-sample tests

- Given $S_1 = \{x \sim P\}$ and $S_2 = \{x \sim Q\}$, a **two-sample test** considers the following hypotheses
 - Null hypothesis $H_0: P = Q$
 - Alternative hypothesis $H_1: P \neq Q$
- Test statistic T compares S_1 and S_2 . For example: difference in means, variances of the two sets of samples
 - $T(S_1, S_2) = \left| \frac{1}{|S_1|} \sum_{x \in S_1} x - \frac{1}{|S_2|} \sum_{x \in S_2} x \right|$
- If T is larger than a threshold α , then reject H_0 otherwise we say H_0 is consistent with observation.
- 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



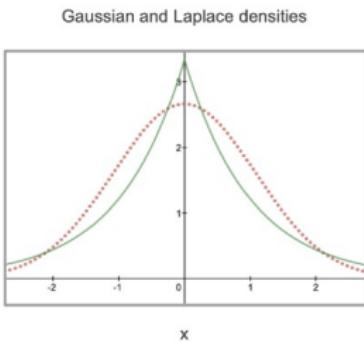
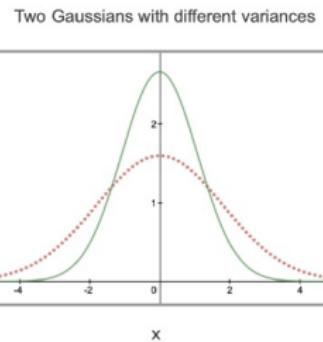
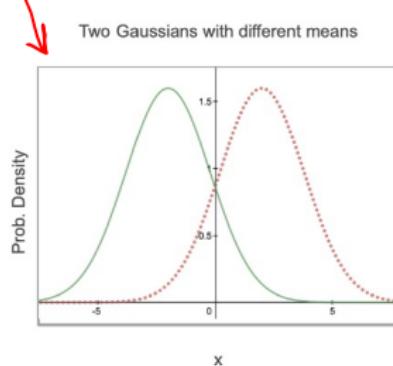
$$\mathbf{x}_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$



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

Two-Sample Test via a Discriminator

- Finding a good two-sample test objective in high dimensions is hard

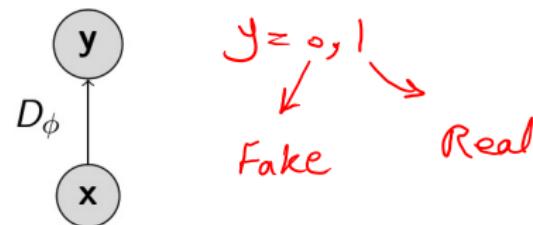


- In the generative model setup, we know that S_1 and S_2 come from different distributions p_{data} and p_{θ} respectively
- Key idea:** Learn a statistic to automatically identify in what way the two sets of samples S_1 and S_2 differ from each other
- How? Train a classifier (called a discriminator)!

Two-Sample Test via a Discriminator

$$D = \{x_1, \dots, x_n\} \quad x_i \sim P_{\text{data}}$$

$$G = \{\tilde{x}_1, \dots, \tilde{x}_n\} \quad \tilde{x}_i \sim P_{\theta}$$



- **Two-Sample Test via a Discriminator**

- Any binary classifier D_ϕ (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: -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 $p_{\text{data}} \neq p_\theta$), or equivalently minimize classification loss

Two-Sample Test via a Discriminator

- Training objective for discriminator:

$$\begin{aligned}\max_{D_\phi} V(p_\theta, D_\phi) &= E_{\mathbf{x} \sim p_{\text{data}}} [\log D_\phi(\mathbf{x})] + E_{\mathbf{x} \sim p_\theta} [\log(1 - D_\phi(\mathbf{x}))] \\ &\approx \sum_{\substack{\mathbf{x} \in S_1 \\ \text{real}}} \underbrace{\log D_\phi(\mathbf{x})}_1 + \sum_{\substack{\mathbf{x} \in S_2 \\ \text{fake}}} \underbrace{[\log(1 - D_\phi(\mathbf{x}))]}_0\end{aligned}$$

- For a fixed generative model p_θ , the discriminator is performing binary classification with the cross entropy objective
 - Assign probability 1 to true data points $\mathbf{x} \sim p_{\text{data}}$ (in set S_1)
 - Assign probability 0 to fake samples $\mathbf{x} \sim p_\theta$ (in set S_2)
 - Optimal discriminator
- $$D^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_\theta(\mathbf{x})}$$
- Sanity check: if $p_\theta = p_{\text{data}}$, classifier cannot do better than chance ($D^*(\mathbf{x}) = 1/2$)

$$\text{Loss} = E_{P_{\text{data}}} \left[\log D(x) \right] + E_{P_{\theta}} \left[\log (1 - D(x)) \right]$$

$$\text{Loss} = \int P_{\text{data}}(x) \log \underline{D(x)} dx + \int P_{\theta}(x) \log (1 - \underline{D(x)}) dx$$

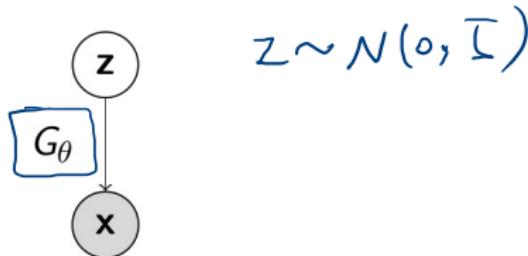
$$G[f] = \int g(f, x) dx$$

$$\frac{\delta G}{\delta f} = \frac{\delta g(f, x)}{\delta f}$$

$$\frac{P_{\text{data}}(x)}{D(x)} + \frac{-P_{\theta}(x)}{1 - D(x)} = 0$$

Generative Adversarial Networks

- A two player minimax game between a **generator** and a **discriminator**



- **Generator**

- Directed, latent variable model with a deterministic mapping between z and x given by G_θ
 - Sample $z \sim p(z)$, where $p(z)$ is a simple prior, e.g. Gaussian
 - Set $x = G_\theta(z)$
- Similar to a flow model, but mapping G_θ need not be invertible
- Distribution over $p_\theta(x)$ over x is implicitly defined (no likelihood!)
- Minimizes a two-sample test objective (in support of the null hypothesis $p_{\text{data}} = p_\theta$)

?

Example of GAN objective

- Training objective for generator:

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

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

$$D_G^* = \frac{p_{\text{data}}}{p_{\text{data}} + p_G}$$

$$\begin{aligned} & V(G, D_G^*(\mathbf{x})) \\ &= E_{\mathbf{x} \sim p_{\text{data}}} \left[\log \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})} \right] + E_{\mathbf{x} \sim p_G} \left[\log \frac{p_G(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})} \right] \\ &= E_{\mathbf{x} \sim p_{\text{data}}} \left[\log \frac{p_{\text{data}}(\mathbf{x})}{\frac{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})}{2}} \right] + E_{\mathbf{x} \sim p_G} \left[\log \frac{p_G(\mathbf{x})}{\frac{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})}{2}} \right] - \log 4 \\ &\xrightarrow{\text{red arrow}} = D_{KL} \left[p_{\text{data}}, \frac{p_{\text{data}} + p_G}{2} \right] + D_{KL} \left[p_G, \frac{p_{\text{data}} + p_G}{2} \right] - \log 4 \\ &\quad \underbrace{2 \times \text{Jensen-Shannon Divergence (JSD)}}_{= 2D_{JSD}[p_{\text{data}}, p_G] - \log 4} \end{aligned}$$

Jenson-Shannon Divergence

- Also called as the symmetric KL divergence

$$D_{JSD}[p, q] = \frac{1}{2} \left(D_{KL}\left[p, \frac{p+q}{2}\right] + D_{KL}\left[q, \frac{p+q}{2}\right] \right)$$

- Properties

- $D_{JSD}[p, q] \geq 0$

- $D_{JSD}[p, q] = 0$ iff $p = q$

- $D_{JSD}[p, q] = D_{JSD}[q, p]$

- $\sqrt{D_{JSD}[p, q]}$ satisfies triangle inequality \rightarrow Jenson-Shannon Distance

- Optimal generator for the JSD/Negative Cross Entropy GAN

$$p_G = p_{\text{data}}$$

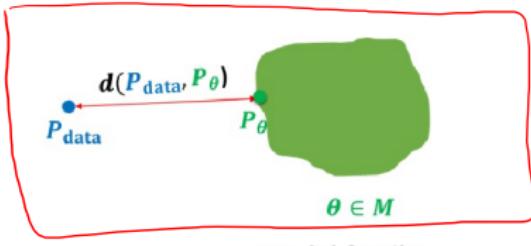
- For the optimal discriminator $D_{G^*}^*(\cdot)$ and generator $G^*(\cdot)$, we have

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

Recap of GANs



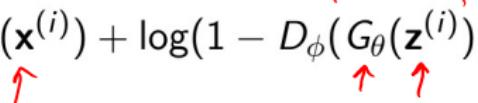
$$x_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$



- Choose $d(p_{\text{data}}, p_\theta)$ to be a two-sample test statistic
 - Learn the statistic by training a classifier (discriminator)
 - Under ideal conditions, equivalent to choosing $d(p_{\text{data}}, p_\theta)$ to be $D_{\text{JS}}[p_{\text{data}}, p_\theta]$
- Pros:
 - Loss only requires samples from p_θ . No likelihood needed!
 - Lots of flexibility for the neural network architecture, any G_θ defines a valid sampling procedure
 - Fast sampling (single forward pass)
- Cons: very difficult to train in practice

The GAN training algorithm

- Sample minibatch of m training points $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}$ from \mathcal{D}
- Sample minibatch of m noise vectors $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(m)}$ from p_z
- Update the discriminator parameters ϕ by stochastic gradient ascent

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


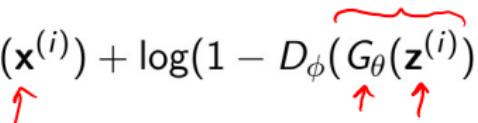
- Update the generator parameters θ by stochastic gradient descent

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

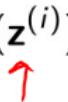

- Repeat for fixed number of epochs

The GAN training algorithm

- Sample minibatch of m training points $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}$ from \mathcal{D}
- Sample minibatch of m noise vectors $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(m)}$ from p_z
- Update the discriminator parameters ϕ by stochastic gradient ascent

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


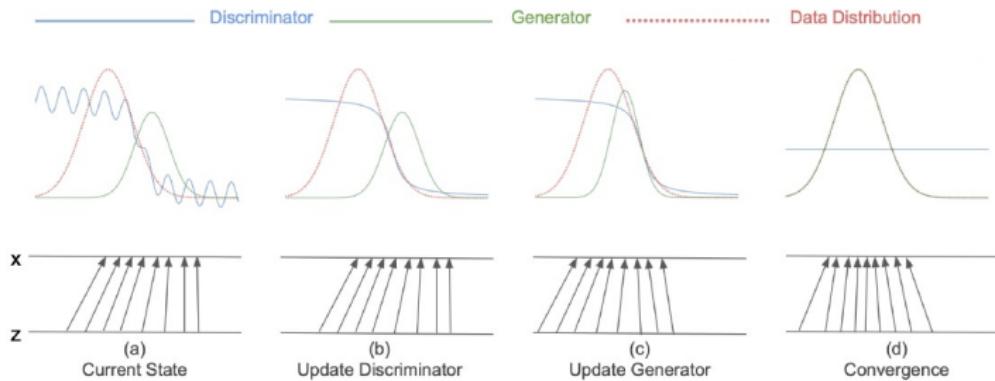
- Update the generator parameters θ by stochastic gradient descent

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


- Repeat for fixed number of epochs

Alternating optimization in GANs

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



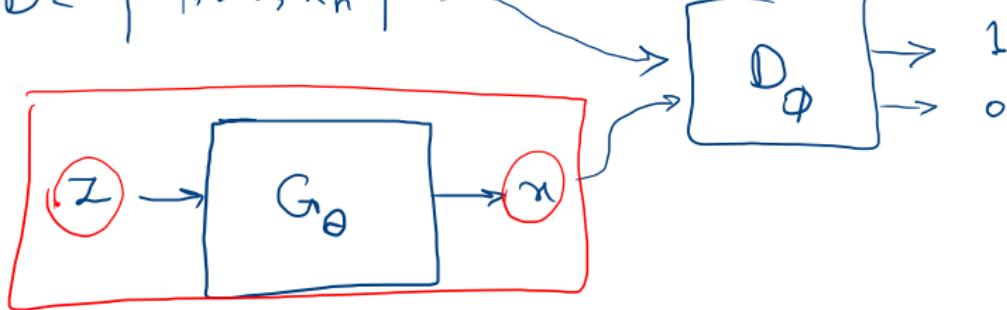
Which one is real?



Source: Karras et al., 2018; The New York Times

Both images are generated via GANs!

$$D = \{n_1, \dots, n_n\}$$



Frontiers in GAN research

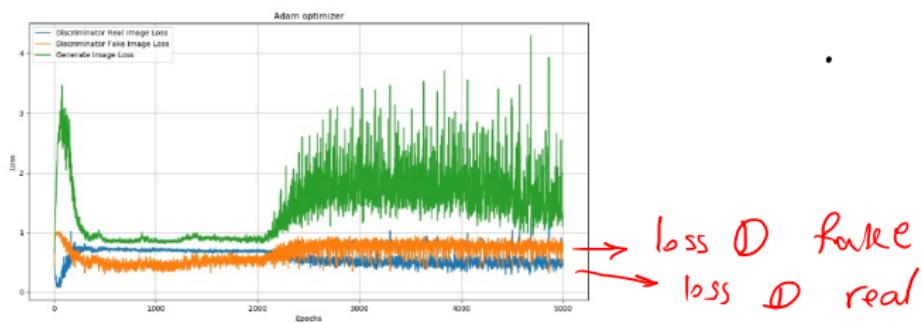


- GANs have been successfully applied to several domains and tasks
- However, working with GANs can be very challenging in practice
 - • Unstable optimization
 - • Mode collapse
 - • Evaluation
- Bag of tricks needed to train GANs successfully

Image Source: Ian Goodfellow. Samples from Goodfellow et al., 2014, Radford et al., 2015, Liu et al., 2016, Karras et al., 2017, Karras et al., 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



Source: Mirantha Jayathilaka

- No robust stopping criteria in practice (unlike MLE)

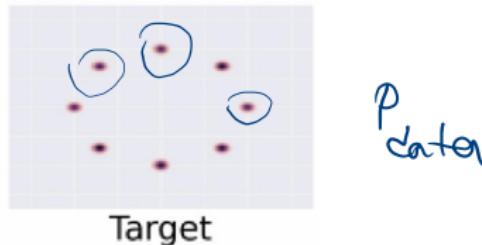
Mode Collapse

- 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”)

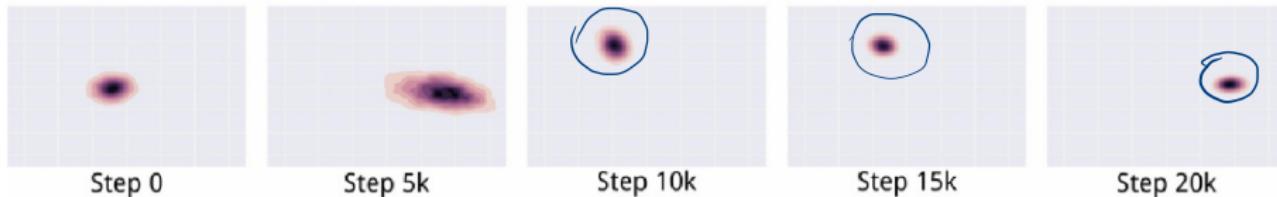


Arjovsky et al., 2017

Mode Collapse



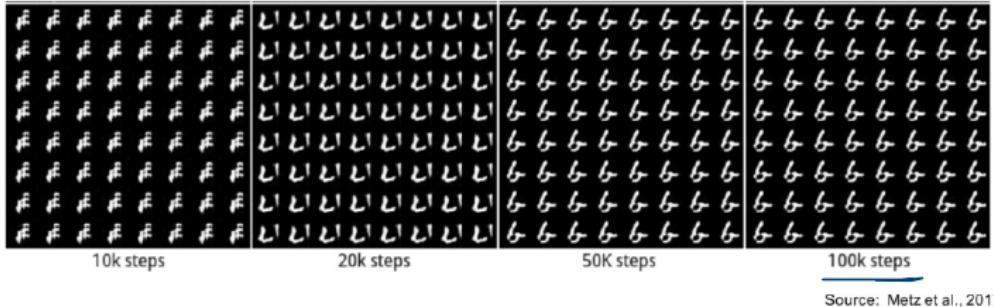
- True distribution is a mixture of Gaussians



Source: Metz et al., 2017

- The generator distribution keeps oscillating between different modes

Mode Collapse



Source: Metz et al., 2017

- Fixes to mode collapse are mostly empirically driven: alternative architectures, alternative GAN loss, adding regularization terms, etc.
- • <https://github.com/soumith/ganhacks> 2020
How to Train a GAN? Tips and tricks to make GANs work by Soumith Chintala