

# Deep Generative Models

## 9. Generative Adversarial Networks



- 국가수리과학연구소 산업수학혁신센터 김민중

# 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

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

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

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

# 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

---

## 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

---

# Recap of GANs

- 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



# Summary of GANs

- Likelihood-free training
- Training objective for GANs

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

- 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  $p_Z$

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

- I.e.,

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

---

# 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

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

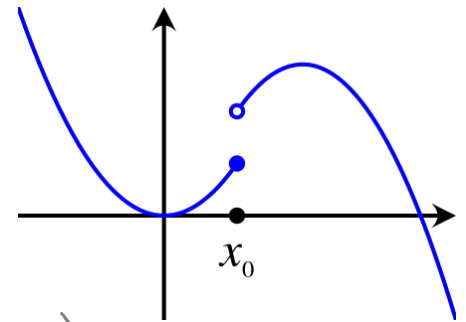
- Where  $f$  is any convex, lower-semicontinuous function with  $f(1) = 0$
- Convex: Line joining any two points lies above the function
- Lower-semicontinuous

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

- 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

# Training with $f$ -divergences

- Given  $p_{data}$  and  $p_{\theta}$ , we could minimize  $D_f(p_{data}, p_{\theta})$  or  $D_f(p_{\theta}, p_{data})$  as learning objectives. Non-negative and zero if  $p_{\theta} = p_{data}$
- However, it depends on the density ratio which is **unknown**

$$D_f(p_{\theta}, p_{data}) = E_{x \sim p_{data}} \left[ f \left( \frac{p_{\theta}(x)}{p_{data}(x)} \right) \right]$$

$$D_f(p_{data}, p_{\theta}) = E_{x \sim p_{\theta}} \left[ f \left( \frac{p_{data}(x)}{p_{\theta}(x)} \right) \right]$$

- To use  $f$ -divergences as a two-sample test objective for likelihood-free learning, we need to be able to estimate the objective using only samples (e.g., training data and samples from the model)

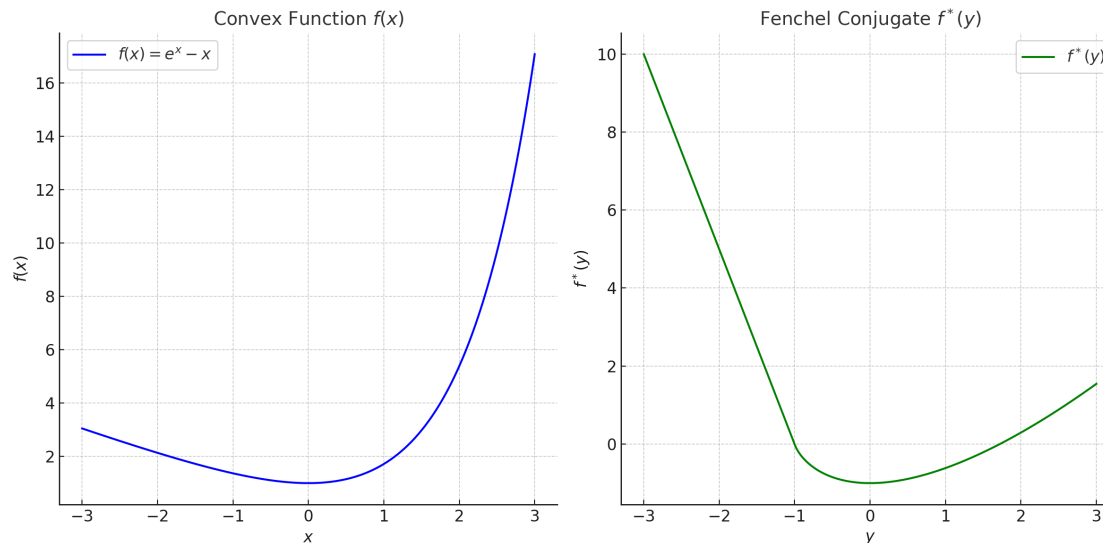
# Towards Variational Divergence Minimization

- Fenchel conjugate: for any function  $f(\cdot)$ , its convex conjugate is

$$f^*(t) := \sup_{u \in \text{dom}_f} (ut - f(u))$$

where  $\text{dom}_f$  is the domain of the function  $f$

- $f^*$  is convex (pointwise supremum of convex functions is convex) and lower semi-continuous



# Towards Variational Divergence Minimization

- Let  $f^{**}$  be the Fenchel conjugate of  $f^*$

$$f^{**}(u) := \sup_{t \in \text{dom}_{f^*}} (tu - f^*(t))$$

- $f^{**} \leq f$ . Proof: By definition, for all  $t, u$

$$f^*(t) \geq ut - f(u) \text{ or equivalently } f(u) \geq ut - f^*(t)$$

$$f(u) \geq \sup_{t \in \text{dom}_{f^*}} (ut - f^*(t)) = f^{**}(u)$$

- Strong Duality:  $f^{**} = f$  when  $f(\cdot)$  is convex and lower semicontinuous

# $f$ -GAN: Variational Divergence Minimization

- We obtain a lower bound to an  $f$ -divergence via Fenchel conjugate

$$\begin{aligned} D_f(p, q) &= E_{\mathbf{x} \sim q} \left[ f \left( \frac{p(\mathbf{x})}{q(\mathbf{x})} \right) \right] = E_{\mathbf{x} \sim q} \left[ f^{**} \left( \frac{p(\mathbf{x})}{q(\mathbf{x})} \right) \right] \\ &= E_{\mathbf{x} \sim q} \left[ \sup_{t \in \text{dom}_{f^*}} \left( t \frac{p(\mathbf{x})}{q(\mathbf{x})} - f^*(t) \right) \right] \\ &\geq \sup_{T \in \mathcal{T}} \int_{\mathcal{X}} [T(\mathbf{x})p(\mathbf{x}) - f^*(T(\mathbf{x}))q(\mathbf{x})] d\mathbf{x} \\ &= \sup_{T \in \mathcal{T}} (E_{\mathbf{x} \sim p}[T(\mathbf{x})] - E_{\mathbf{x} \sim q}[f^*(T(\mathbf{x}))]) \end{aligned}$$

- where  $\mathcal{T}: \mathcal{X} \rightarrow \mathbb{R}$  is an arbitrary class of functions
- Note:** Lower bound is likelihood-free w.r.t.  $p$  and  $q$



# $f$ -GAN: Variational Divergence Minimization

- Variational lower bound

$$D_f(p, q) \geq \sup_{T \in \mathcal{T}} (E_{x \sim p}[T(x)] - E_{x \sim q}[f^*(T(x))])$$

- Choose an  $f$ -divergence
- Let  $p = p_{data}$  and  $q = p_G$
- Parameterize  $T$  by  $\phi$  and  $G$  by  $\theta$
- Consider the following  $f$ -GAN object

$$\min_{\theta} \max_{\phi} F(\theta, \phi) = \min_{\theta} \max_{\phi} E_{x \sim p_{data}}[T_{\phi}(x)] - E_{x \sim p_{G_{\theta}}}[f^*(T_{\phi}(x))]$$

- Generator  $G_{\theta}$  tries to minimize the divergence estimate and discriminator  $T_{\phi}$  tries to tighten the lower bound
- Substitute any  $f$ -divergence and optimize the  $f$ -GAN objective
- Prior distribution  $p_Z$

$$\min_{\theta} \max_{\phi} E_{x \sim p_{data}}[T_{\phi}(x)] - E_{z \sim p_Z}[f^*(T_{\phi}(G_{\theta}(z)))]$$

# Example: Univariate Mixture of Gaussians

- $p_{data}$ : a mixture of Gaussians
- Model  $Q_{\theta}$  using linear transformation of a standard normal  $z \sim N(0,1)$  and outputs  $G_{\theta}(z) = \mu + \sigma z$ , where  $\theta = (\mu, \sigma)$

	KL	KL-rev	JS	Jeffrey	Pearson
$D_f(P  Q_{\theta^*})$	0.2831	0.2480	0.1280	0.5705	0.6457
$F(\hat{\omega}, \hat{\theta})$	0.2801	0.2415	0.1226	0.5151	0.6379
$\mu^*$	1.0100	1.5782	1.3070	1.3218	0.5737
$\hat{\mu}$	1.0335	1.5624	1.2854	1.2295	0.6157
$\sigma^*$	1.8308	1.6319	1.7542	1.7034	1.9274
$\hat{\sigma}$	1.8236	1.6403	1.7659	1.8087	1.9031

train \ test	KL	KL-rev	JS	Jeffrey	Pearson
KL	<b>0.2808</b>	0.3423	0.1314	0.5447	0.7345
KL-rev	0.3518	<b>0.2414</b>	0.1228	0.5794	1.3974
JS	0.2871	0.2760	<b>0.1210</b>	0.5260	0.92160
Jeffrey	0.2869	0.2975	0.1247	<b>0.5236</b>	0.8849
Pearson	0.2970	0.5466	0.1665	0.7085	<b>0.648</b>

**Table 3:** Gaussian approximation of a mixture of Gaussians. Left: optimal objectives, and the learned mean and the standard deviation:  $\hat{\theta} = (\hat{\mu}, \hat{\sigma})$  (learned) and  $\theta^* = (\mu^*, \sigma^*)$  (best fit). Right: objective values to the true distribution for each trained model. For each divergence, the lowest objective function value is achieved by the model that was trained for this divergence.

Source: Nowozin et al., 2017

# Beyond KL and Jensen-Shannon Divergence

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

$$\min_{\theta} \max_{\phi} F(\theta, \phi) = \min_{\theta} \max_{\phi} E_{x \sim p_{data}} [T_{\phi}(x)] - E_{x \sim p_{G_{\theta}}} [f^*(T_{\phi}(x))]$$

# Wasserstein GAN: beyond $f$ -divergence

- The  $f$ -divergence is defined as

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

- The support of  $q$  must cover the support of  $p$ . Otherwise, discontinuity arises in  $f$ -divergences

- E.g.,

- Let  $p(x) = \begin{cases} 1, & x = 0 \\ 0, & x \neq 0 \end{cases}$  and  $q_\theta(x) = \begin{cases} 1, & x = \theta \\ 0, & x \neq \theta \end{cases}$

- $D_{KL}(p, q_\theta) = \begin{cases} 0, & \theta = 0 \\ \infty, & \theta \neq 0 \end{cases}$

- $D_{JS}(p, q_\theta) = \begin{cases} 0, & \theta = 0 \\ \log 2, & \theta \neq 0 \end{cases}$

- We need a “smoother” distance  $D(p, q)$  that is defined when  $p$  and  $q$  have disjoint supports

# Wasserstein (Earth-Mover) distance

- Introduced by Leonid Vaseršteĭn(Russia)
- 1<sup>st</sup> Wasserstein distance

$$\begin{aligned} D_w(p, q) &:= \inf_{\gamma \in \Gamma(p, q)} \int_{K \times K} |\mathbf{x} - \mathbf{y}| d\gamma(\mathbf{x}, \mathbf{y}) \\ &= \inf_{\gamma \in \Gamma(p, q)} \sum_{\mathbf{x}, \mathbf{y}} |\mathbf{x} - \mathbf{y}| \gamma(\mathbf{x}, \mathbf{y}) \end{aligned}$$

- where  $\Gamma(p, q)$  contains all joint distributions of  $(\mathbf{x}, \mathbf{y})$  where the marginal of  $\mathbf{x}$  is  $p(\mathbf{x})$  and the marginal of  $\mathbf{y}$  is  $q(\mathbf{y})$
- $\gamma(\mathbf{y}|\mathbf{x})$ : a probabilistic earth moving plan that warps  $p(\mathbf{x})$  to  $q(\mathbf{y})$
- Let  $p(x) = \begin{cases} 1, & x = 0 \\ 0, & x = 1 \end{cases}$  and  $q_\theta(x) = \begin{cases} 1, & x = \theta \\ 0, & x \neq \theta \end{cases}$
- $D_w(p, q_\theta) = |\theta|$

# Wasserstein GAN (WGAN)

- Kantorovich-Rubinstein duality

$$D_w(p, q) = \sup_{\|f\|_L \leq 1} E_{x \sim p}[f(x)] - E_{x \sim q}[f(x)]$$

- $\|f\|_L \leq 1$  means the Lipschitz constant of  $f(x)$  is 1. i.e.,

$$|f(x) - f(y)| \leq \|x - y\|_1 \quad \forall x, y$$

- Intuitively,  $f$  cannot change too rapidly
- Wasserstein GAN with discriminator  $D_\phi(x)$  and generator  $G_\theta(z)$

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

- Lipschitzness of  $D_\phi(x)$  is enforced through **weight clipping** or **gradient penalty** on  $\nabla_x D_\phi(x)$
- To enforce Lipschitz constraint, clip the weights of the critic to lie within a compact space  $[-c, c]$ . The set of functions satisfying this constraint is a subset of the  $K$ -Lipschitz functions for some  $K(c)$
- If we replace  $\|f\|_L \leq 1$  for  $\|f\|_L \leq K$ , then we end up with  $K \cdot D_w(p, q)$

# Wasserstein GAN Gradient Penalty

- Wasserstein GAN-GP

$$\min_{\theta} \max_{\phi} E_{x \sim p_{data}} [D_{\phi}(x)] - E_{z \sim p_Z} [D_{\phi}(G_{\theta}(z))] \\ - \lambda E_{\hat{x} \sim p_{\hat{x}}} \left[ \left( \|\nabla_{\hat{x}} D_{\phi}(\hat{x})\|_2 - 1 \right)^2 \right]$$

**Proposition 1.** Let  $\mathbb{P}_r$  and  $\mathbb{P}_g$  be two distributions in  $\mathcal{X}$ , a compact metric space. Then, there is a 1-Lipschitz function  $f^*$  which is the optimal solution of  $\max_{\|f\|_L \leq 1} \mathbb{E}_{y \sim \mathbb{P}_r}[f(y)] - \mathbb{E}_{x \sim \mathbb{P}_g}[f(x)]$ . Let  $\pi$  be the optimal coupling between  $\mathbb{P}_r$  and  $\mathbb{P}_g$ , defined as the minimizer of:  $W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\pi \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \pi} [\|x - y\|]$  where  $\Pi(\mathbb{P}_r, \mathbb{P}_g)$  is the set of joint distributions  $\pi(x, y)$  whose marginals are  $\mathbb{P}_r$  and  $\mathbb{P}_g$ , respectively. Then, if  $f^*$  is differentiable<sup>‡</sup>,  $\pi(x = y) = 0$ <sup>§</sup>, and  $x_t = tx + (1 - t)y$  with  $0 \leq t \leq 1$ , it holds that  $\mathbb{P}_{(x,y) \sim \pi} \left[ \nabla f^*(x_t) = \frac{y - x_t}{\|y - x_t\|} \right] = 1$ .

**Corollary 1.**  $f^*$  has gradient norm 1 almost everywhere under  $\mathbb{P}_r$  and  $\mathbb{P}_g$ .

# Wasserstein GAN Gradient Penalty

- Wasserstein GAN-GP

$$\min_{\theta} \max_{\phi} E_{x \sim p_{data}} [D_{\phi}(x)] - E_{z \sim p_Z} [D_{\phi}(G_{\theta}(z))] \\ - \lambda E_{\hat{x} \sim p_{\hat{x}}} \left[ \left( \|\nabla_{\hat{x}} D_{\phi}(\hat{x})\|_2 - 1 \right)^2 \right]$$

- Sampling distribution**  $p_{\hat{x}}$ : uniformly along straight lines between pairs of points sampled from the data distribution and the generator distribution

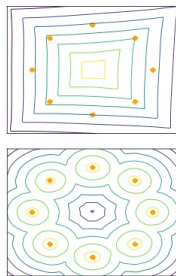


# Wasserstein GAN Gradient Penalty

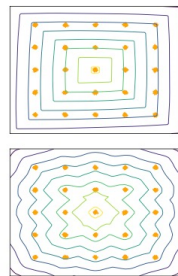
- Wasserstein GAN-GP

$$\min_{\theta} \max_{\phi} E_{x \sim p_{data}} [D_{\phi}(x)] - E_{z \sim p_Z} [D_{\phi}(G_{\theta}(z))] - \lambda E_{\hat{x} \sim p_{\hat{X}}} \left[ \left( \|\nabla_{\hat{x}} D_{\phi}(\hat{x})\|_2 - 1 \right)^2 \right]$$

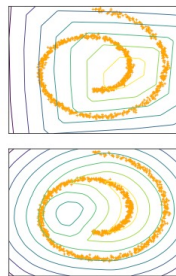
8 Gaussians



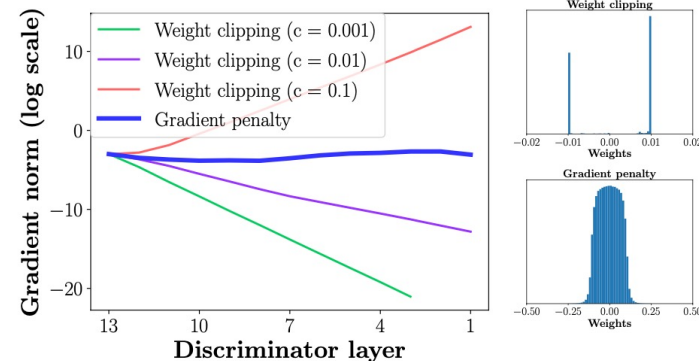
25 Gaussians



Swiss Roll



(a) Value surfaces of WGAN critics trained to optimality on toy datasets using (top) weight clipping and (bottom) gradient penalty. Critics trained with weight clipping fail to capture higher moments of the data distribution. The ‘generator’ is held fixed at the real data plus Gaussian noise.



(b) (left) Gradient norms of deep WGAN critics during training on the Swiss Roll dataset either explode or vanish when using weight clipping, but not when using a gradient penalty. (right) Weight clipping (top) pushes weights towards two values (the extremes of the clipping range), unlike gradient penalty (bottom).

Figure 1: Gradient penalty in WGANs does not exhibit undesired behavior like weight clipping.

---

# Inferring latent representations in GANs

- The generator of a GAN is typically a directed, latent variable model with latent variables  $\mathbf{z}$  and observed variables  $\mathbf{x}$
- How can we infer the **latent feature representations** in a GAN?
- Unlike a normalizing flow model, the mapping  $G: \mathbf{z} \rightarrow \mathbf{x}$  need not be invertible
- Unlike a variational autoencoder, there is no inference network  $q(\cdot)$  which can learn a variational posterior over latent variables

---

# Inferring latent representations in GANs

- **Solution 1:** For any point  $x$ , use the activations of the prefinal layer of a **discriminator** as a feature representation
- **Intuition:** Like supervised deep neural networks, the discriminator would have learned useful representations for  $x$  while distinguishing real and fake

---

# Inferring latent representations in GANs

- If we want to directly infer the latent variables  $\mathbf{z}$  of the generator, we need a **different learning algorithm**
- A regular GAN optimizes a two-sample test objective that compares samples of  $\mathbf{x}$  from the generator and the data distribution
- **Solution 2:** To infer latent representations, we will compare samples of  $\mathbf{x}$ ,  $\mathbf{z}$  from the joint distributions of observed and latent variables as per the model and the data distribution
- For any  $\mathbf{x}$  generated via the model, we have access to  $\mathbf{z}$  (sampled from a simple prior  $p(\mathbf{z})$ )
- For any  $\mathbf{x}$  from the data distribution,  $\mathbf{z}$  is however unobserved (latent). Need an encoder

---

# Summary of Generative Adversarial Networks

- **Key observation:** Samples and likelihoods are not correlated in practice
- Two-sample test objectives allow for learning generative models only via samples (likelihood-free)
- Wide range of two-sample test objectives covering  $f$ -divergences and Wasserstein distances (and more)

# Thanks

---