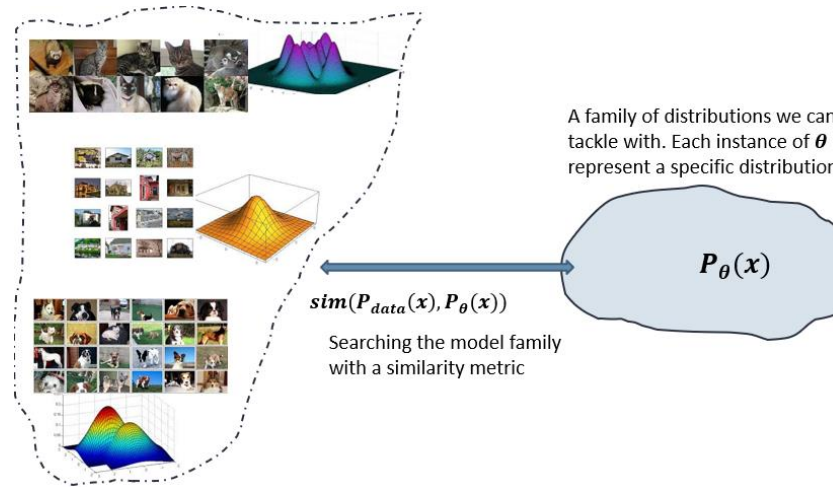# Generative adversarial networks

22-808: Generative models
Sharif University of Technology
Fall 2024

Fatemeh Seyyedsalehi

# Recap



A family of distributions we can tackle with. Each instance of $\theta$ represent a specific distribution

$$P_\theta(x)$$

$$sim(P_{data}(x), P_\theta(x))$$

Searching the model family with a similarity metric

▸ We need a framework to interact with distributions for statistical generative models.

    ▸ Probabilistic generative models

    ▸ Deep generative models

        ▸ Autoregressive models   $p_\theta(\mathbf{x}) = \prod_{i=1}^{n} p_\theta(x_i | \mathbf{x}_{<i})$

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

        ▸ **Generative adversarial networks**

    ▸ Both AR and VAE model families attempted to minimize the KL divergence between model family and data distribution, or equivalently attempt to maximize the likelihood.

    ▸ In GAN we are going to use an alternative choice for the similarity measure between model distribution and data distribution.
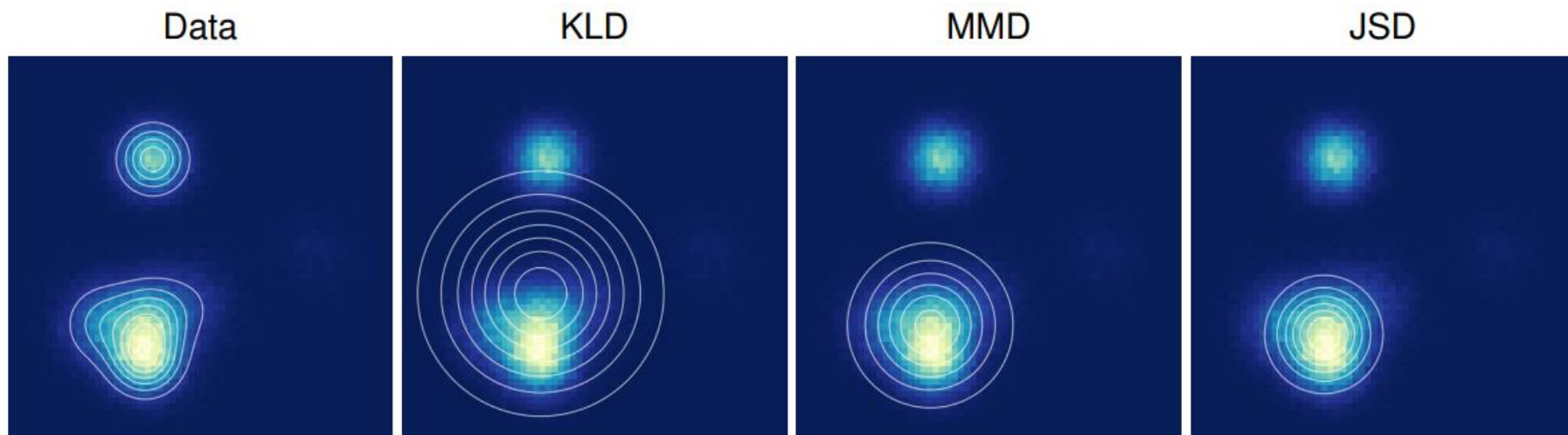
2

# Maximizing the likelihood

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

▸ Optimal statistical efficiency

  ▸ Assume sufficient model capacity, such that there exists a unique $\theta^*$ that satisfy $p_{\theta^*} = p_{data}$.

  ▸ The convergence of $\hat{\theta}$ to $\theta^*$ when $M \to \infty$, is the fastest among all statistical methods when using maximum likelihood training.
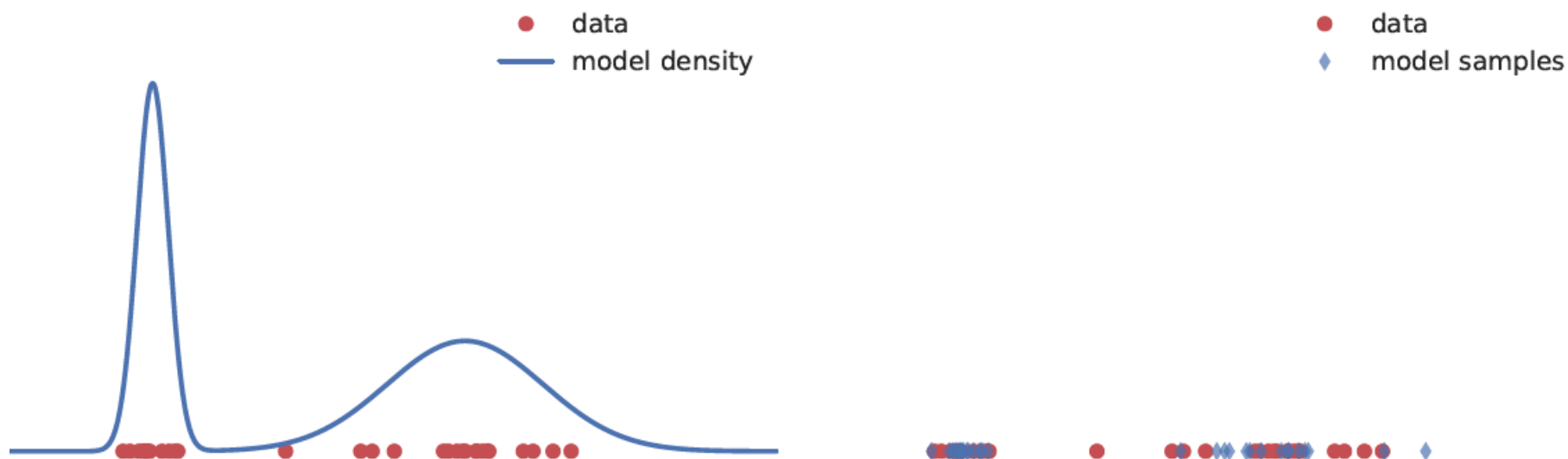
# Maximizing the likelihood

▸ For imperfect models, achieving high log-likelihoods might not always imply good sample quality.



An isotropic Gaussian distribution was fit to data drawn from a mixture of Gaussians by either minimizing KL divergence (KLD), maximum mean discrepancy (MMD), or Jensen-Shannon divergence (JSD). The different fits demonstrate different tradeoffs made by the three measures of distance between distributions.

# Implicit generative models

▸ Kind of probabilistic generative models without an explicit likelihood function

▸ We use a likelihood-free approach to train these models

  ▸ Training by comparing samples



Explicit models vs. implicit models

# Learning by comparing samples

‣ We should define a distance(similarity) measure between two distributions that:

  ‣ Provides guarantees about learning the data distribution.

$$\operatorname*{argmin}_{p_\theta} D(p_{data}, p_\theta) = p_{data}$$

  ‣ Can be evaluated only using samples from the data and model distribution.

  ‣ Are computationally cheap to evaluate.

‣ Many distributional distances and divergences fail to satisfy the later two requirements

# Learning by comparing samples

▶ The main approach to overcome these challenges is to approximate the desired quantity through optimization by introducing a comparison model, often called a **discriminator** or a **critic** $D$, such that:

$$\mathcal{D}(p^*, q) = \operatorname*{argmax}_{D} \mathcal{F}(D, p^*, q)$$

▶ where $\mathcal{F}$ is a functional that can be estimated using only samples from $p^*(p_{data})$ and $q$. One way is that it depends on distributions only in expectations.

    ▶ Therefore, it can be estimated using Monte Carlo estimation.

# Learning by comparing samples

▸ As we usually use parametric functions (ex. Neural networks) for both the model and discriminator.

▸ Therefore, by the following optimization we estimate the distance measure $\mathcal{D}(p^*, q_\theta)$

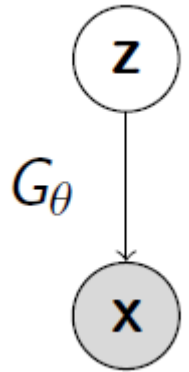$$\mathrm{argmax}_{\boldsymbol{\phi}}\, \mathcal{F}(D_{\boldsymbol{\phi}}, p^*, q_\theta)$$

▸ Then, instead of optimizing the exact objective $\mathcal{D}(p^*, q_\theta)$

we use the tractable approximation provided through the optimal $D_{\boldsymbol{\phi}}$.

# Generative adversarial networks (Goodfellow GAN)

▶ A finite number of samples from the desired real distribution is available: $x_1, x_2, \ldots, x_n$

▶ Like VAEs, we consider a latent variable model

for the model generation process and attempt to

learn $G_\theta$. However, here we learn this function by

Comparing samples.

# The Goodfellow GAN
## The probabilistic classification view

▶ Assuming $D(x)$ as a binary classifier which predicts whether a given point $x$ was sampled from the real distribution or it is a fake sample from the generator $G_\theta$.

▶ A cross entropy loss to train this classifier:

$$E_{\mathbf{x} \sim p_{\text{data}}}[\log D_\phi(\mathbf{x})] + E_{\mathbf{x} \sim p_\theta}[\log(1 - D_\phi(\mathbf{x}))]$$

▶ We can see that the optimal discriminator for a fixed generator $G_\theta$ is:

$$\frac{p(x)}{p(x) + p_\theta(x)}$$

# The Goodfellow GAN
## The objective function

▸ By substitution the optimal discriminator into the cross-entropy loss, we have:

$$V^*(q_\theta, p^*) = \frac{1}{2}\mathbb{E}_{p^*(\boldsymbol{x})}[\log \frac{p^*(\boldsymbol{x})}{p^*(\boldsymbol{x}) + q_\theta(\boldsymbol{x})}] + \frac{1}{2}\mathbb{E}_{q_\theta(\boldsymbol{x})}[\log(1 - \frac{p^*(\boldsymbol{x})}{p^*(\boldsymbol{x}) + q_\theta(\boldsymbol{x})})]$$

$$= \frac{1}{2}\mathbb{E}_{p^*(\boldsymbol{x})}[\log \frac{p^*(\boldsymbol{x})}{\frac{p^*(\boldsymbol{x})+q_\theta(\boldsymbol{x})}{2}}] + \frac{1}{2}\mathbb{E}_{q_\theta(\boldsymbol{x})}[\log(\frac{q_\theta(\boldsymbol{x})}{\frac{p^*(\boldsymbol{x})+q_\theta(\boldsymbol{x})}{2}})] - \log 2$$

$$= \frac{1}{2}D_{\mathbb{KL}}\left(p^* \| \frac{p^* + q_\theta}{2}\right) + \frac{1}{2}D_{\mathbb{KL}}\left(q_\theta \| \frac{p^* + q_\theta}{2}\right) - \log 2$$

$$= JSD(p^*, q_\theta) - \log 2$$

where JSD is the Jensen-Shannon divergence.

# The Goodfellow GAN

▸ This establishes a connection between optimal binary classification and distributional divergences.

▸ By using binary classification, we were able to compute the distributional divergence using only samples, which is the important property needed for learning implicit generative models

▸ We have turned an intractable estimation problem (how to estimate the JSD divergence) into an optimization problem (how to learn a classifier) which can be used to approximate that divergence.

# The Goodfellow GAN

▸ With optimal discriminator, we attempt to find the generative model $G_\theta$ that minimizes the JSD divergence.

$$\min_{\boldsymbol{\theta}} JSD(p^*, q_\theta) = \min_{\boldsymbol{\theta}} V^*(q_\theta, p^*) + \log 2$$

$$= \min_{\boldsymbol{\theta}} \frac{1}{2} \mathbb{E}_{p^*(\boldsymbol{x})} \log D^*(\boldsymbol{x}) + \frac{1}{2} \mathbb{E}_{q_\theta(\boldsymbol{x})} \log(1 - D^*(\boldsymbol{x})) + \log 2$$

# Training procedure of GAN

Sample minibatch of $m$ training points $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots, \mathbf{x}^{(m)}$ from $\mathcal{D}$

Sample minibatch of $m$ noise vectors $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \ldots, \mathbf{z}^{(m)}$ from $p_z$

Update the discriminator parameters $\phi$ 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 $\theta$ 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

Activa

# Training convergence

▶ If G and D have enough capacity, and at each step of training procedure, the discriminator is allowed to reach its optimum for a specific $G_\theta$ , and then $p_\theta$ is updated so as to improve
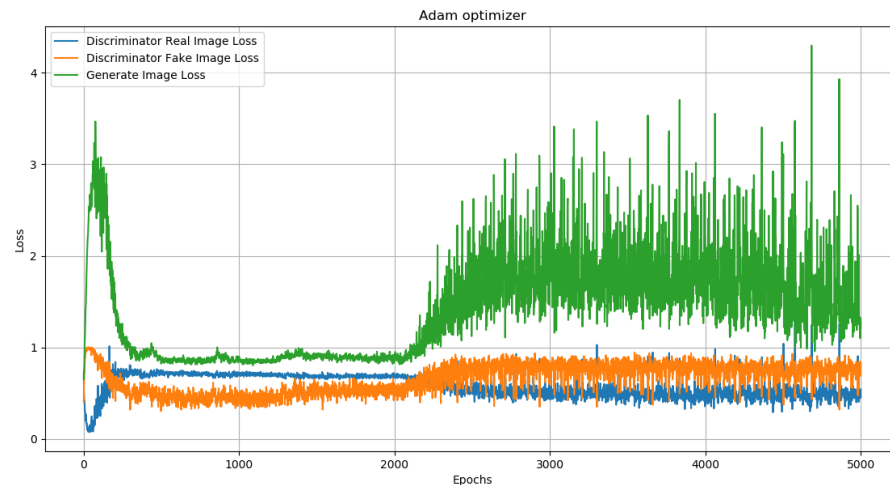
$$\mathbb{E}_{\boldsymbol{x} \sim p_{data}}[\log D_G^*(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{x} \sim p_g}[\log(1 - D_G^*(\boldsymbol{x}))]$$

then $p_\theta$ converges to $p_{data}$.

▶ Unrealistic assumptions ☹

# Training convergence

▸ However, we do not have access to the optimal discriminator and only we can approximate it with a parametrized function: neural network $D_\phi$

  ▸ No guarantee for convergence

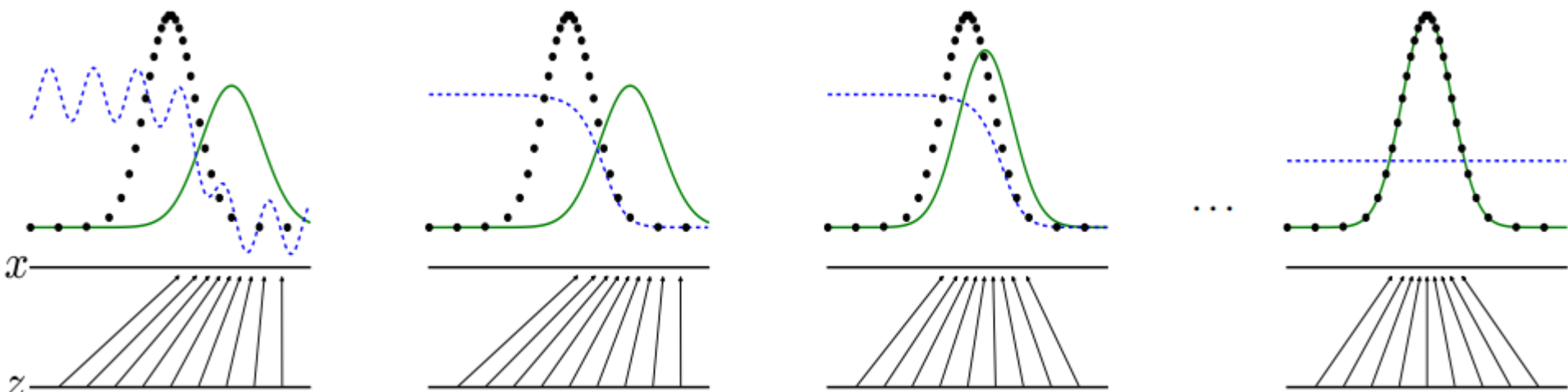  ▸ In practice, the generator and discriminator loss keeps oscillating during GAN training

# The min-max game

▶ The minmax game

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

  ▶ It is a game not an optimization problem
  ▶ It should reach to a Nash equilibria

# Example

▸ Which one is real?

# F-divergence

▸ Let $f: R \to R$ be a convex lower-semicontinuous function, such that $f(1) = 0$. We define the *f-divergence* between two distributions with densities $p$ and $q$ by:

$$D_f(p \parallel q) \equiv \int_{\mathcal{X}} q(x) f\left(\frac{p(x)}{q(x)}\right) dx.$$

▸ What's interesting about *f-divergence* is that we can construct a variational representation for it.

  ▸ Alternating the integral to an optimization

# Fenchel duality

▸ The idea is to use the convex conjugate of the function $f$, which is defined as follows:

$$f^*(t) \equiv \sup_{x}\{tx - f(x)\}.$$

▸ Fenchel duality: repeat application of the conjugate operation to convex lower-semicontinuous function $f$ yields $f^{**} = f$. Therefore, we have:

$$f(x) = \sup_{t}\{tx - f^*(t)\}.$$

# Variational representation of *F-divergence*

▸ Using Fenchel duality, we obtain the variational representation of the *f-divergence.*

$$
\begin{aligned}
\boxed{D_f(p \parallel q)} &= \int_{\mathcal{X}} q(x) \sup_t \left[ t\frac{p(x)}{q(x)} - f^*(t) \right] dx \\
&= \int_{\mathcal{X}} \sup_t \left[ tp(x) - f^*(t)q(x) \right] dx \\
&= \sup_{T:\mathcal{X}\to\mathbb{R}} \int_{\mathcal{X}} \left( T(x)p(x) - f^*(T(x))q(x) \right) dx \\
&= \boxed{\sup_{T:\mathcal{X}\to\mathbb{R}} \left[ \mathbb{E}_{x\sim p} T(x) - \mathbb{E}_{x\sim q} f^*(T(x)) \right]}.
\end{aligned}
$$

# F-GAN

▸ The dual form can be approximated using Monte Carlo estimation.

▸ Assuming a parametric family of functions $T\varphi$ (ex. a neural network) and the generator function $g_\theta$, and a valid f-divergence, the F-GAN objective is,

$$\theta_f = \arg\min_\theta \sup_\varphi \left[ \mathbb{E}_{x \sim p} T_\varphi(x) - \mathbb{E}_{x \sim p_\theta} f^*(T_\varphi(x)) \right]$$

$$= \arg\min_\theta \sup_\varphi \left[ \mathbb{E}_{x \sim p} T_\varphi(x) - \mathbb{E}_{z \sim q} f^*(T_\varphi(g_\theta(z))) \right].$$

▸ Generator $g_\theta$ tries to minimize the divergence estimate and discriminator $T\varphi$ tries to tighten the lower bound

# *F-divergence*

| distance or divergence | corresponding $g(t)$ $(t = \frac{p_i(x)}{p_j(x)})$ |
| --- | --- |
| Bhattacharyya distance [1] | $\sqrt{t}$ |
| KL-divergence | $t \log(t)$ |
| Symmetric KL-divergence | $t \log(t) - \log(t)$ |
| Hellinger distance | $(\sqrt{t} - 1)^2$ |
| Total variation | $|t - 1|$ |
| Pearson divergence | $(t - 1)^2$ |
| Jensen-Shannon divergence | $\frac{1}{2}(t \log \frac{2t}{t+1} + \log \frac{2}{t+1})$ |

# The Goodfellow GAN as F-GAN

▶ The Goodfellow GAN is an instances of the *f-GAN.*

▶ Modified version of the Jensen-Shannon

$$2\mathrm{JSD}(p, q) - \log(4) = D_{\mathrm{KL}}\left(p\left\|\frac{p + q}{2}\right.\right) + D_{\mathrm{KL}}\left(p_g\left\|\frac{p + q}{2}\right.\right) - \log(4).$$

▶ The *f-divergence:*

$$f(x) = x \log x - (x+1) \log(x+1)$$

$$f^*(t) = -\log(1 - e^t).$$

$$T_\varphi(x) = \log(d_\varphi(x))$$

▶ *We can obtain the* Goodfellow GAN :

$$\theta_f = \arg \min_\theta \sup_\varphi \left[ \mathbb{E}_{x \sim p} \log d_\varphi(x) + \mathbb{E}_{z \sim q} \log(1 - d_\varphi(g_\theta(z))) \right]$$