

Two Generator Game: Learning to Sample via Linear Goodness-of-Fit Test

Lizhong Ding, Mengyang Yu, Li Liu, Fan Zhu, Yong Liu, Yu Li, Ling Shao

Inception Institute of Artificial Intelligence (IIAI), Abu Dhabi, UAE.

Institute of Information Engineering, CAS, China.

King Abdullah University of Science and Technology (KAUST), Saudi Arabia.



Background

- Learning the probability distribution of high-dimensional data.
- GANs are a type of implicit generative models (IGMs).
- Existing GAN models are fundamentally two-sample test problems.
- We summarize existing GAN models into two categories:
 - ▷ Integral Probability Metric (IPM).

$$\delta(\mathbf{p}, \mathbf{q}) = \sup_{f \in \mathcal{F}} \left| \int_{\mathcal{X}} f(\mathbf{x}) \mathbf{p}(\mathbf{x}) d\mathbf{x} - \int_{\mathcal{X}} f(\mathbf{x}') \mathbf{q}(\mathbf{x}') d\mathbf{x}' \right|.$$

Wasserstein GANs (WGANs), MMD-GAN

- ▷ ζ -divergence

$$\delta_{\zeta}(\mathbf{p}, \mathbf{q}) = \int_{\mathcal{X}} \mathbf{q}(\mathbf{x}) \zeta \left(\frac{\mathbf{p}(\mathbf{x})}{\mathbf{q}(\mathbf{x})} \right) d\mathbf{x},$$

where ζ is a convex, lower-semicontinuous, satisfying $\zeta(\mathbf{1}) = \mathbf{0}$.
GAN, least squares GAN

- Analyzing and comparing distributions without imposing any parametric assumptions

- ▷ Two-sample test:

Whether two distributions \mathbf{p} and \mathbf{q} are different based on

$$\begin{aligned} \mathcal{D}_{\mathbf{x}} &= \{\mathbf{x}_i\}_{i=1}^n \subset \mathcal{X} \subseteq \mathbb{R}^d \sim \mathbf{p} \\ \mathcal{D}_{\mathbf{y}} &= \{\mathbf{y}_j\}_{j=1}^m \subset \mathcal{Y} \subseteq \mathbb{R}^d \sim \mathbf{q} \end{aligned}$$

- ▷ Goodness-of-fit test:

How well a given model density \mathbf{p} fits a set of given samples

$$\mathcal{D}_{\mathbf{x}} = \{\mathbf{x}_i\}_{i=1}^n \subset \mathcal{X} \subseteq \mathbb{R}^d \sim \mathbf{q}$$

Contributions

- Deep energy adversarial network (DEAN):
A new paradigm that casts the generative adversarial learning as a goodness-of-fit (GOF) test problem.
- A novel two generator game via GOF tests:
 - ▷ One explicit generator is designed to learn an energy-based distribution (EBD), which maps the real data to a scalar energy-based probability.
 - ▷ The other implicit generator is trained by minimizing the vFSSD between the EBD and the generated data.
- A two-level alternative optimization procedure to train the explicit and implicit generative networks, such that the hyper-parameters can also be automatically learned.

Energy Estimator Network

Energy-based models $\mathcal{E}_{\theta}(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}$ associate an energy value with a sample \mathbf{x} , where θ are the parameters. We can obtain a distribution based on $\mathcal{E}_{\theta}(\mathbf{x})$,

$$\mathbf{p}(\mathbf{x}; \theta) = \frac{1}{Z_{\theta}} \exp(-\mathcal{E}_{\theta}(\mathbf{x})).$$

Now we define the loss function of the explicit generative network (EGN) of DEAN as follows:

$$\min_{\theta_e} \mathcal{E}(\mathbf{x}; \theta_e) + \left[\gamma - \mathcal{E}(\mathbf{G}(\mathbf{z}; \theta_g^*); \theta_e) \right]^+, \quad (1)$$

where $\mathcal{E}(\mathbf{x}; \theta_e)$ is an *energy model* parameterized by θ_e , $[\cdot]^+ = \max(\cdot, 0)$ and γ is a given positive margin. When the network parameters θ_e^* are optimized, we can define a probability distribution

$$\mathbf{p}(\mathbf{x}; \theta_e^*) = \frac{1}{Z_{\theta_e^*}} \exp(-\mathcal{E}(\mathbf{x}; \theta_e^*)).$$

Energy Estimator Network

We consider a deep auto-encoder as a more complex energy model

$$\mathcal{E}(\mathbf{x}; \theta) = \|\mathbf{x} - \mathbf{AE}(\mathbf{x}; \theta_e)\|,$$

where $\mathbf{AE}(\mathbf{x}; \theta_e)$ denotes a deep auto-encoder parameterized by θ_e .

For the optimized parameters θ_e^* , we can define

$$\mathbf{p}(\mathbf{x}; \theta_e^*) = \frac{1}{Z_{\theta_e^*}} \exp(-\mathcal{E}(\mathbf{x}; \theta_e^*)).$$

GOF-driven Generator Network

Kernel Stein operator can be written as

$$(\mathbf{T}_{\mathbf{p}}\mathbf{f})(\mathbf{x}) = \sum_{i=1}^d \left(\frac{\partial \log \mathbf{p}(\mathbf{x})}{\partial \mathbf{x}_i} \mathbf{f}_i(\mathbf{x}) + \frac{\partial \mathbf{f}_i(\mathbf{x})}{\partial \mathbf{x}_i} \right) = \langle \mathbf{f}, \omega_{\mathbf{p}}(\mathbf{x}, \cdot) \rangle_{\mathcal{F}^d}.$$

Now we introduce the kernel Stein discrepancy (KSD), which is formulated as

$$\text{KSD}[\mathcal{F}^d, \mathbf{p}, \mathcal{D}_{\mathbf{x}}] = \sup_{\|\mathbf{f}\|_{\mathcal{F}^d} \leq 1} \langle \mathbf{f}, \mathbf{E}_{\mathbf{x} \sim \mathbf{q}} \omega_{\mathbf{p}}(\mathbf{x}, \cdot) \rangle := \|\mathbf{g}(\cdot)\|_{\mathcal{F}^d}, \quad (2)$$

where $\mathbf{g}(\cdot) = \mathbf{E}_{\mathbf{x} \sim \mathbf{q}} \omega_{\mathbf{p}}(\mathbf{x}, \cdot)$. The statistic of the finite set Stein discrepancy (FSSD) is defined as

$$\text{FSSD}[\mathcal{F}^d, \mathbf{p}, \mathcal{D}_{\mathbf{x}}] = \frac{1}{dJ} \sum_{i=1}^d \sum_{j=1}^J \mathbf{g}_i^2(\mathbf{v}_j).$$

The unbiased estimator of FSSD is defined as

$$\widehat{\text{FSSD}}^2[\mathcal{F}^d, \mathbf{p}, \mathcal{D}_{\mathbf{x}}] = \frac{2}{n(n-1)} \sum_{i < j} \Delta(\mathbf{x}_i, \mathbf{x}_j),$$

where $\Delta(\mathbf{x}, \mathbf{y}) = \tau(\mathbf{x})^T \tau(\mathbf{y})$. To alleviate the impact of dimension and stabilize the statistic, we introduce

$$\text{vFSSD}[\mathbf{p}, \mathcal{D}_{\mathbf{x}}] = \frac{1}{\hat{\sigma}_{\mathbf{H}_1}} \widehat{\text{FSSD}}^2[\mathbf{p}, \mathcal{D}_{\mathbf{x}}].$$

GOF-driven Generator Network

Now we define the loss function of the IGN as follows:

$$\min_{\theta_g} \max_{\xi} \text{vFSSD}_{\xi}[\mathbf{p}(\mathbf{x}; \theta_e^*), \mathcal{D}_{\mathbf{x}'}], \quad (3)$$

where $\xi = \{\{\mathbf{v}_i\}_{i=1}^J, \sigma_k\}$ denotes the hyper-parameters of vFSSD, including the kernel parameter σ_k and J test locations $\{\mathbf{v}_i\}_{i=1}^J$. We present the following two objectives, (4) and (5), to optimize Equation (3) and improve the test power of DEAN.

$$\max_{\xi} \text{vFSSD}_{\xi}[\mathbf{p}(\mathbf{y}; \theta_e^*), \mathcal{D}_{\mathbf{x}'}^*], \quad (4)$$

where $\mathcal{D}_{\mathbf{x}'^*} = \{\mathbf{x}_i'^* := \mathbf{G}(\mathbf{z}_i; \theta_g^*)\}_{i=1}^n$ and $\mathbf{G}(\mathbf{z}_i; \theta_g^*)$ is a deep network with the optimized parameter θ_g^* . The hyper-parameters $\xi = \{\{\mathbf{v}_i\}_{i=1}^J, \sigma_k\}$ will be optimized in Equation (4).

$$\min_{\theta_g} \text{vFSSD}_{\xi^*}[\mathbf{p}(\mathbf{y}; \theta_e^*), \mathcal{D}_{\mathbf{x}'}], \quad (5)$$

where $\xi^* = \{\{\mathbf{v}_i^*\}_{i=1}^J, \sigma_k^*\}$ denotes the optimized hyper-parameters, and the parameters θ_g for $\mathcal{D}_{\mathbf{x}'} = \{\mathbf{x}_i' := \mathbf{G}(\mathbf{z}_i; \theta_g)\}_{i=1}^n$ will be optimized.

Theorem

We assume that $\mathcal{D}_{\mathbf{x}'}$ is drawn from $\mathbf{p}_{\mathbf{x}'}$. If κ is a universal and analytic kernel;
 $\mathbf{E}_{\mathbf{a} \sim \mathbf{p}_{\mathbf{x}'}} \mathbf{E}_{\mathbf{b} \sim \mathbf{p}_{\mathbf{e}}} \left[\mathbf{s}^T(\mathbf{a}) \mathbf{s}(\mathbf{b}) \kappa(\mathbf{a}, \mathbf{b}) + \mathbf{s}^T(\mathbf{b}) \nabla_{\mathbf{a}} \kappa(\mathbf{a}, \mathbf{b}) + \mathbf{s}^T(\mathbf{a}) \nabla_{\mathbf{b}} \kappa(\mathbf{a}, \mathbf{b}) + \sum_{i=1}^d \frac{\partial^2 \kappa(\mathbf{a}, \mathbf{b})}{\partial \mathbf{a}_i \partial \mathbf{b}_i} \right] < \infty$
with $\mathbf{s}(\mathbf{a}) = \nabla_{\mathbf{a}} \log \mathbf{p}_{\mathbf{e}}(\mathbf{a})$; $\mathbf{E}_{\mathbf{a} \sim \mathbf{p}_{\mathbf{x}'}} \|\nabla_{\mathbf{a}} \log \mathbf{p}_{\mathbf{e}}(\mathbf{a}) - \nabla_{\mathbf{a}} \log \mathbf{p}_{\mathbf{x}'}(\mathbf{a})\|^2 < \infty$;
 $\lim_{\|\mathbf{a}\| \rightarrow \infty} \mathbf{p}_{\mathbf{e}}(\mathbf{a}) \mathbf{g}(\mathbf{a}) = \mathbf{0}$, where $\mathbf{g}(\cdot)$ is given in Eq. (2) in Section 4.2; for any $J \geq 1$, almost surely $\text{FSSD}[\mathbf{p}_{\mathbf{e}}, \mathcal{D}_{\mathbf{x}'}] = 0$ if and only if $\mathbf{p}_{\mathbf{x}'} = \mathbf{p}_{\mathbf{e}}$.

Theorem

Let $\Lambda(\theta_e) = \mathcal{E}(\mathbf{x}; \theta_e) + \left[\gamma - \mathcal{E}(\mathbf{G}(\mathbf{z}; \theta_g^*); \theta_e) \right]^+$. The minimum of $\Lambda(\theta_e)$ is achieved if and only if $\mathbf{p}_{\mathbf{e}} = \mathbf{p}_{\mathbf{x}}$. With the optimized θ_e^* , $\int_{\mathbf{x}, \mathbf{z}} \Lambda(\theta_e^*) \mathbf{p}_{\mathbf{x}}(\mathbf{x}) \mathbf{p}_{\mathbf{z}}(\mathbf{z}) d\mathbf{x} d\mathbf{z} = \gamma$.