

Generative Adversarial Networks

Yutao Chen

✘ May 13 2026

✘ May 13 2026

Contents

Noise Contrastive Estimation	1
Generative Adversarial Networks	2

Generative adversarial networks (GANs) (Goodfellow et al., 2014) are generative models consisting of:

- A generator G that approximates the data distribution, and
- A discriminator D that learns to distinguish real / generated data.

Noise Contrastive Estimation

Generative adversarial networks are closely related to noise contrastive estimation (NCE) (Gutmann & Hyvärinen, 2010), a parameter estimation method for (unnormalized) probabilistic models.

Assume we have samples from an unknown data distribution $p_{\text{data}}(\mathbf{x})$, which we want to approximate with a probabilistic model $p_{\theta}(\mathbf{x})$. Noise contrastive estimation learns $p_{\theta}(\mathbf{x})$ by comparison:

Given a noise distribution $q(\mathbf{x})$ ¹ with known density, we can learn $p_{\theta}(\mathbf{x})$ by estimating the *density ratio*

$$r(\mathbf{x}) = p_{\text{data}}(\mathbf{x})/q(\mathbf{x}).$$

That is, if we know $r(\mathbf{x})$ and $q(\mathbf{x})$, we have $p_{\theta}(\mathbf{x}) \approx r(\mathbf{x})q(\mathbf{x})$.

Furthermore, estimating the density ratio $r(\mathbf{x})$ can be casted as *binary classification*. Assuming we label data samples $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$ with $y = 1$ and noise $\mathbf{x} \sim q(\mathbf{x})$ with $y = 0$, we have

$$p_{\text{data}}(\mathbf{x}) = p(\mathbf{x}|y = 1), \quad q(\mathbf{x}) = p(\mathbf{x}|y = 0).$$

¹It must hold that $q(\mathbf{x})$ is non-zero everywhere $p_{\text{data}}(\mathbf{x})$ is non-zero.

By the Bayes theorem, we have $p(y|\mathbf{x}) = p(\mathbf{x}|y) \cdot p(y)/p(\mathbf{x})$. Therefore, the density ratio $r(\mathbf{x})$ can be equivalently written as

$$r(\mathbf{x}) = \frac{p_{\theta}(\mathbf{x})}{q(\mathbf{x})} = \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} = \frac{D^*(\mathbf{x})}{1 - D^*(\mathbf{x})} \cdot \frac{1 - \pi}{\pi},$$

where $D^*(\mathbf{x}) := p(y=1|\mathbf{x})$ and $\pi := p(y=1)^2$. Note that $D^*(\mathbf{x})$ is the optimal binary classifier for $p_{\text{data}}(\mathbf{x})$ and $q(\mathbf{x})$, and

$$D^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + q(\mathbf{x})} = \frac{1}{1 + \exp(-\log r(\mathbf{x}))}.$$

Finally, we can find $D^*(\mathbf{x})$ by maximizing the Bernoulli log-likelihood for some parametric classifier $D_{\theta}(\mathbf{x})$

$$\arg \max_{\theta} \mathbb{E}_{p_{\text{data}}(\mathbf{x})}[\log D_{\theta}(\mathbf{x})] + \mathbb{E}_{q(\mathbf{x})}[\log(1 - D_{\theta}(\mathbf{x}))].$$

If we find $D_{\theta}(\mathbf{x}) \approx D^*(\mathbf{x})$, we also find $r(\mathbf{x})$ and $p_{\theta}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x})$.

Generative Adversarial Networks

NCE turns density estimation of $p_{\theta}(\mathbf{x})$ into binary classification $D_{\theta}(\mathbf{x})$. The noise distribution $q(\mathbf{x})$ is often chosen to be simplistic and tractable (e.g. Gaussians). In contrast, *generative adversarial networks* learn the noise distribution $q_{\phi}(\mathbf{x})$ jointly with the binary classifier $D_{\theta}(\mathbf{x})$:

$$\min_{\phi} \max_{\theta} \mathbb{E}_{p_{\text{data}}(\mathbf{x})}[\log D_{\theta}(\mathbf{x})] + \mathbb{E}_{q_{\phi}(\mathbf{x})}[\log(1 - D_{\theta}(\mathbf{x}))]$$

With the above minimax optimization framework, a unique equilibrium exists at $q_{\phi}(\mathbf{x}) = p_{\text{data}}(\mathbf{x})$ and $D_{\theta}(\mathbf{x}) = \frac{1}{2}$ for all \mathbf{x} .

Proof. Note that for any fixed ϕ , the optimal θ corresponds to

$$D_{\theta}(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + q_{\phi}(\mathbf{x})}.$$

Furthermore, for any optimal θ , the optimal ϕ minimizes the following objective

²We assume without loss of generality $\pi = \frac{1}{2}$ hereafter.

$$\begin{aligned}
& \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[\log \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + q_{\phi}(\mathbf{x})} \right] + \mathbb{E}_{q_{\phi}(\mathbf{x})} \left[\log \frac{q_{\phi}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + q_{\phi}(\mathbf{x})} \right] \\
= & \mathbb{D}_{\text{KL}} \left(p_{\text{data}} \parallel \frac{p_{\text{data}} + q_{\phi}}{2} \right) + \mathbb{D}_{\text{KL}} \left(q_{\phi} \parallel \frac{p_{\text{data}} + q_{\phi}}{2} \right) - 2 \log 2 \\
= & 2 \text{JSD} (p_{\text{data}}(\mathbf{x}) \parallel q_{\phi}(\mathbf{x})) - \log 4,
\end{aligned}$$

where $\text{JSD}(p \parallel q)$ denotes the Jensen-Shannon divergence, minimized when $p = q$. Therefore, the optimal ϕ for any optimal θ is given by $p_{\text{data}}(\mathbf{x}) = q_{\phi}(\mathbf{x})$. ■

The noise distribution $q_{\phi}(\mathbf{x})$ is often parametrized as a neural network $G_{\phi}(z)$, where $z \sim p(z)$ and $p(z)$ is a simple prior noise distribution.

We refer the readers to [Murphy \(2023\)](#) Section 26.2 for further analysis on connecting GAN, *proper scoring rule*, and statistical divergence.

REFERENCES

- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative Adversarial Nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, & K. Weinberger (Eds.), *Advances in Neural Information Processing Systems: Vol. 27. Advances in Neural Information Processing Systems*.
- Gutmann, M., & Hyvärinen, A. (2010). Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In Y. W. Teh & M. Titterton (Eds.), *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics: Vol. 9. Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*. <https://proceedings.mlr.press/v9/gutmann10a.html>
- Murphy, K. P. (2023). *Probabilistic Machine Learning: Advanced Topics*. MIT Press. <http://probml.github.io/book2>