

Improved Techniques for Training Score-Based Generative Models

Yang Song (Stanford University)
Stefano Ermon (Stanford University)



Overview

Score-based generative models [1] can produce high-quality samples comparable to GANs without requiring adversarial training.

How they work:

1. Perturb the data distribution with multiple scales of noise.
2. Jointly estimate the score (gradient of log probability density) of each noise-perturbed data distribution by training a noise-conditional model with score matching.
3. Generate samples by running Langevin MCMC on noise-conditional score models while gradually annealing down the noise scales.

Our contributions:

1. Theoretically-guided methods for choosing noise scales and setting the hyperparameters of Langevin MCMC.
2. Improve the performance of previous models, scaling the resolution of samples to 256 x 256.

Background

Score: $\nabla_{\mathbf{x}} \log p(\mathbf{x})$

Score-Based Model: $\mathbf{s}_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x})$

Noise-Perturbed Distribution:

$$p_{\sigma}(\mathbf{x}) := \int p(\mathbf{x}') \mathcal{N}(\mathbf{x}; \mathbf{x}', \sigma^2 \mathbf{I}) d\mathbf{x}'$$

Noise-Conditional Score-Based Model: $\mathbf{s}_{\theta}(\mathbf{x}, \sigma) \approx \nabla_{\mathbf{x}} \log p_{\sigma}(\mathbf{x})$

Multiple Scales of Noise Perturbation:

$$\sigma_1 < \sigma_2 < \dots < \sigma_N$$

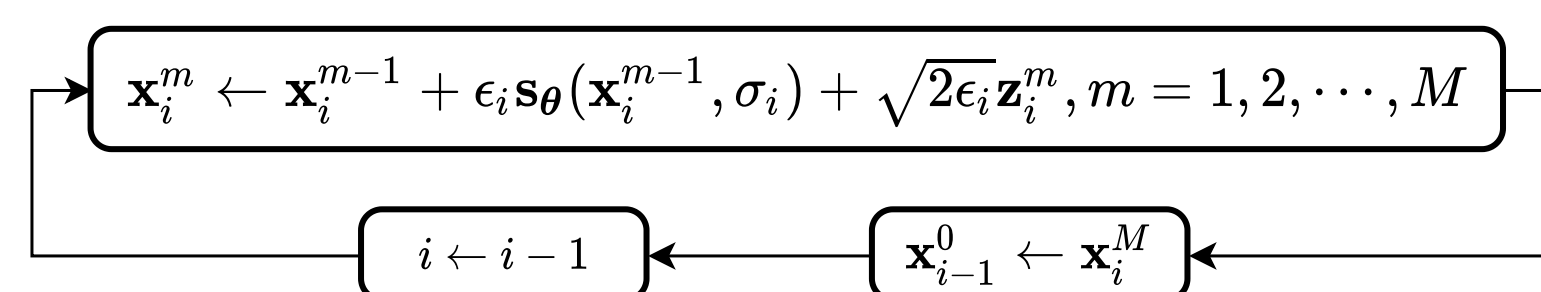
$$\mathbf{s}_{\theta}(\mathbf{x}, \sigma_i) \approx \nabla_{\mathbf{x}} \log p_{\sigma_i}(\mathbf{x}), \quad i = 1, 2, \dots, N$$

Training:

$$\frac{1}{2N} \sum_{i=1}^N \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \mathbb{E}_{\mathcal{N}(\tilde{\mathbf{x}}; \mathbf{x}, \sigma_i^2 \mathbf{I})} \left[\left\| \sigma_i \mathbf{s}_{\theta}(\tilde{\mathbf{x}}, \sigma_i) + \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma_i} \right\|_2^2 \right]$$

denoising score matching

Annealed Langevin dynamics:

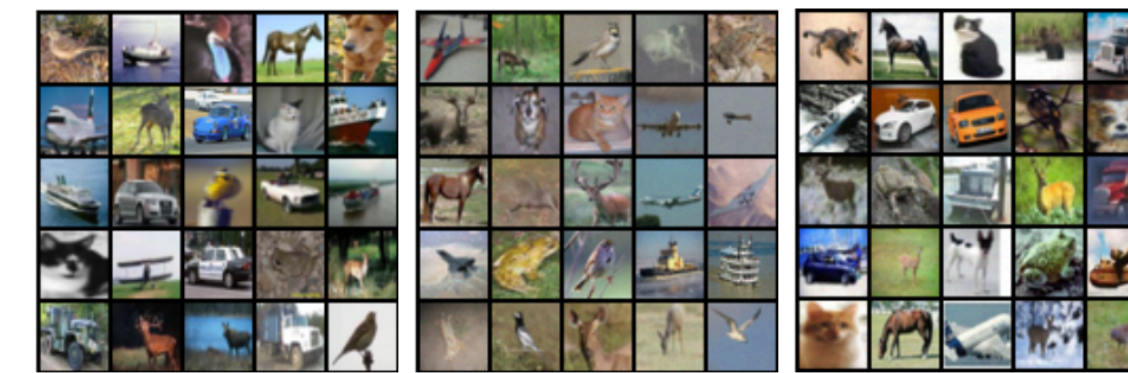


Repeat for N times

Choosing Noise Scales

Initial Noise Scale: σ_N

Theoretical analysis assuming data distribution is a mixture of Gaussian:



(a) Data (b) $\sigma_N = 1$ (c) $\sigma_N = 50$

Sampling from a mixture of Gaussian centered at test images with annealed Langevin dynamics using different initial noise scales.

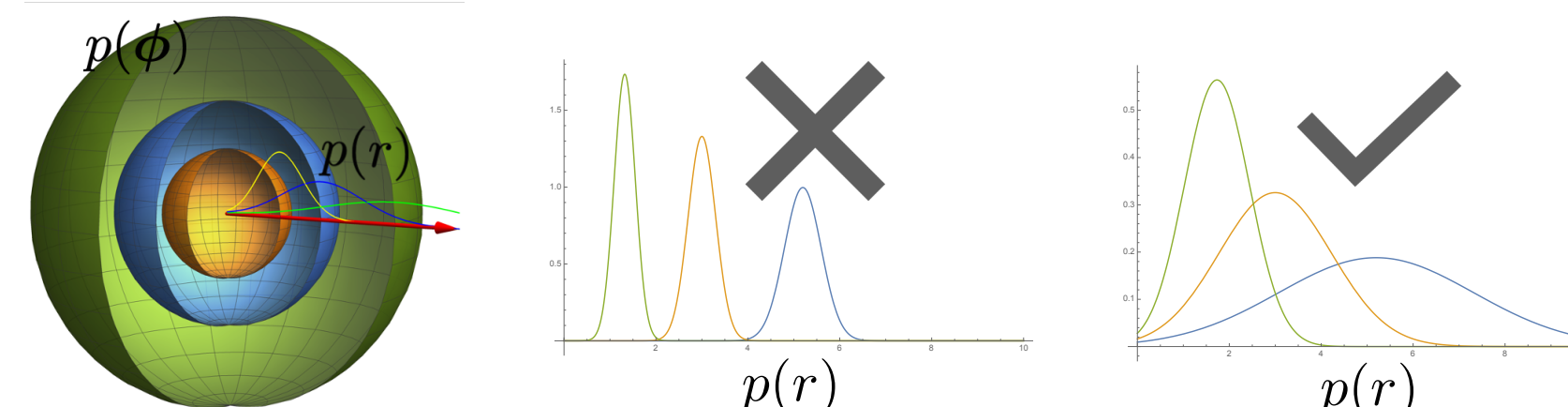
Technique 1 (Initial noise scale). Choose σ_N to be as large as the maximum Euclidean distance between all pairs of training data points.

Final Noise Scale: $\sigma_1 = 0.01$ (small enough to make the smallest noise-perturbed data distribution indiscernible from clean data)

Other Noise Scales: choose N and $\sigma_2 < \sigma_3 < \dots < \sigma_{N-1}$

Intuition: adjacent noise-perturbed distributions should have sufficient **overlap**.

Analysis: assuming data distribution is a single Gaussian.



Technique 2 (Other noise scales). Choose $\{\sigma_i\}_{i=1}^N$ as a geometric progression with common ratio γ , such that $\Phi(\sqrt{2D}(\gamma - 1) + 3\gamma) - \Phi(\sqrt{2D}(\gamma - 1) - 3\gamma) \approx 0.5$.

Configuring Annealed Langevin Dynamics

Theoretical analysis assuming data distribution is a single Gaussian.

Proposition 3. Let $\gamma = \frac{\sigma_i}{\sigma_{i-1}}$, and we choose the step size $\epsilon_i = \epsilon \cdot \frac{\sigma_i^2}{\sigma_N^2}$. After running Langevin MCMC, we have the sample $\mathbf{x}^M \sim \mathcal{N}(\mathbf{0}, s_i^2 \mathbf{I})$, where

$$\frac{s_i^2}{\sigma_i^2} = \left(1 - \frac{\epsilon}{\sigma_N^2}\right)^{2M} \left(\gamma^2 - \frac{2\epsilon}{\sigma_N^2 - \sigma_N^2 \left(1 - \frac{\epsilon}{\sigma_N^2}\right)^2}\right) + \frac{2\epsilon}{\sigma_N^2 - \sigma_N^2 \left(1 - \frac{\epsilon}{\sigma_N^2}\right)^2}.$$

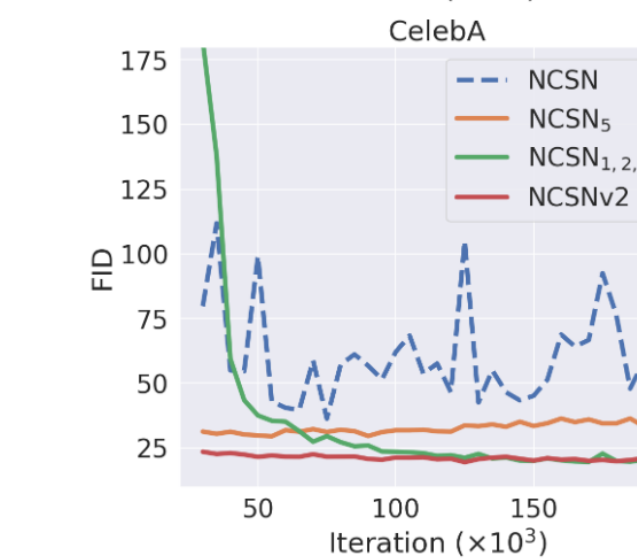
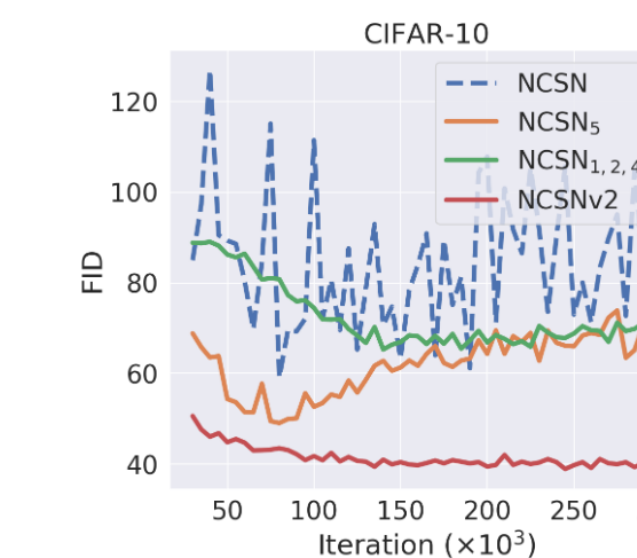
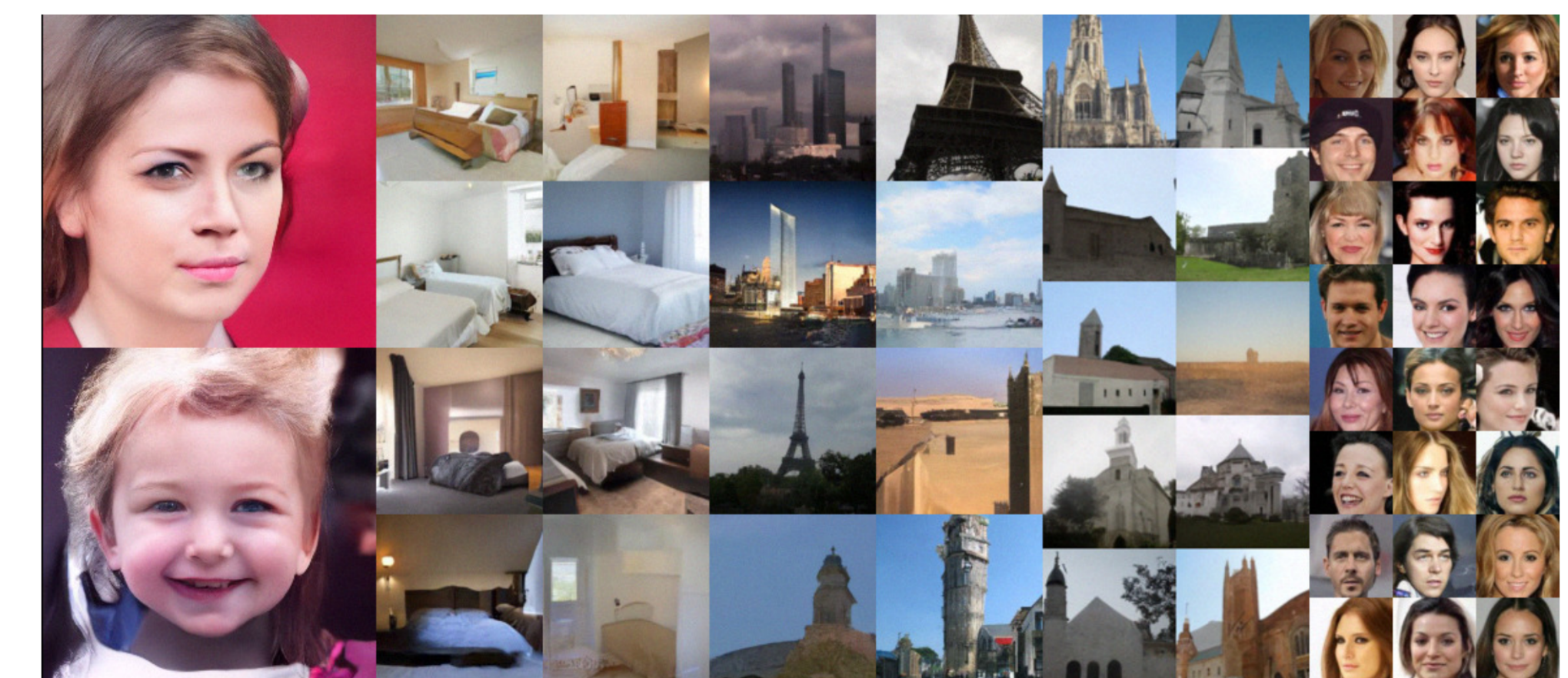
Technique 4 (selecting M and ϵ). Choose M as large as allowed by a computing budget and then select an ϵ that makes $\frac{s_i^2}{\sigma_i^2}$ maximally close to 1.

Other Improved Techniques

Technique 3 (Noise conditioning). Parameterize the NCSN with $\mathbf{s}_{\theta}(\mathbf{x}, \sigma) = \mathbf{s}_{\theta}(\mathbf{x})/\sigma$, where $\mathbf{s}_{\theta}(\mathbf{x})$ is an unconditional score network.

Technique 5 (EMA). Apply exponential moving average to parameters when sampling.

Experimental Results



| Model | Inception \uparrow | FID \downarrow |
|---|----------------------------------|------------------|
| CIFAR-10 Unconditional | | |
| PixelCNN [17] | 4.60 | 65.93 |
| IGBM [18] | 6.02 | 40.58 |
| WGAN-GP [19] | 7.86 \pm .07 | 36.4 |
| SNGAN [20] | 8.22 \pm .05 | 21.7 |
| NCSN [1] | 8.87 \pm .12 | 25.32 |
| NCSN (w/ denoising) | 7.32 \pm .12 | 29.8 |
| NCSNv2 (w/o denoising) | 8.73 \pm .13 | 31.75 |
| NCSNv2 (w/ denoising) [2] | 8.40 \pm .07 | 10.87 |
| CelebA 64 \times 64 | | |
| NCSN (w/o denoising) | - | 26.89 |
| NCSN (w/ denoising) [2] | - | 25.30 |
| NCSNv2 (w/o denoising) | - | 28.86 |
| NCSNv2 (w/ denoising) [2] | - | 10.23 |

References

[1] Song, Y. and Ermon, S., 2019. Generative modeling by estimating gradients of the data distribution. In Advances in Neural Information Processing Systems (pp. 11918-11930).

[2] Jolicœur-Martineau, A., Piché-Taillefer, R., Combes, R.T.D. and Mitliagkas, I., 2020. Adversarial score matching and improved sampling for image generation. In arXiv preprint arXiv:2009.05475.



Paper



Code