

BOOT: Data-free Distillation of Denoising Diffusion Models with Bootstrapping

Pramook Khungurn

June 27, 2023

This note was written as I read the paper “BOOT: Data-free Distillation of Denoising Diffusion Models with Bootstrapping” by Gu et al. [GZZ⁺23].

1 Introduction

- It is a known problem that diffusion models are slow compared other generative models such that GANs or VAEs.
- There are techniques that distills diffusion models into models that can generate data in few model evaluation steps. These include:
 - Straightforward knowledge distillation [LL21].
 - Progressive distillation [SH22, MRG⁺23].
 - Consistency models [SDCS23].
- However, these techniques require either (1) generating a lot of synthetic data with the teacher model or (2) using the training dataset in the process of distillation.
- The requirements above make it hard to apply for text-condition diffusion models, which currently requires a very large dataset to train.
- The paper proposes BOOT, which can distill diffusion models without needing any data.
- It is inspired by consistency models [SDCS23].
 - The process of sampling a diffusion model can be thought of as tracing a trajectory of a particle moving through a velocity field defined by the probability flow ODE [SSDK⁺21].
 - In the consistency model work, Song et al. observes that all points in a specific trajectory form an equivalent class. What we want to do is to find the terminal point from any other points the trajectory.
 - So, a consistency model wants any \mathbf{x}_t in the same trajectory to map to the same \mathbf{x}_0 .
- On the other hand, BOOT predicts all possible \mathbf{x}_t given the Gaussian noise $\boldsymbol{\xi}$ and time t .
- The name “BOOT” comes from the word “bootstrapping,” which is used in the meaning that it becomes easier to predict \mathbf{x}_t if the model has already learned to predict $\mathbf{x}_{t'}$ with $t' > t$.

2 Background

2.1 Diffusion Model

- The paper uses the continuous-time [SSDK⁺21, KAAL22] and variance-preserving [SH22] formulation.
- A data point is denoted by $\mathbf{x} \in \mathbb{R}^N$.
- The forward process generates a stochastic process $\{\mathbf{x}_t : t \in [0, T]\}$ with $\mathbf{x}_0 = \mathbf{x} \sim p_{\text{data}}$ governed by the **noise schedule** α_t and σ_t with $\alpha_t^2 + \sigma_t^2 = 1$ (because we are using the variance-preserving formulation). The stochastic process has the following properties:

$$\begin{aligned} q(\mathbf{x}_t | \mathbf{x}) &= \mathcal{N}(\mathbf{x}_t; \alpha_t \mathbf{x}, \sigma_t^2 I), \\ q(\mathbf{x}_t | \mathbf{x}_s) &= \mathcal{N}(\mathbf{x}_t; \alpha_{t|s} \mathbf{x}_s, \sigma_{t|s}^2 I) \end{aligned}$$

where

$$\begin{aligned} \alpha_{t|s} &= \frac{\alpha_t}{\alpha_s} \\ \sigma_{t|s}^2 &= \sigma_t^2 - \alpha_{t|s}^2 \sigma_s^2 \end{aligned}$$

for $s < t$.

- The quality α_t^2 / σ_t^2 is called the **signal-to-noise ratio (SNR)**. It decreases monotonically with t .
- A diffusion model is denoted by \mathbf{f}_ϕ . In this paper, it predicts the denoised data \mathbf{x} from \mathbf{x}_t and t . It is trained with the following objective:

$$\mathcal{L}_\phi^{\text{diff}} = E_{\mathbf{x} \sim p_{\text{data}}, t \sim [0, T], \mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x})} [w_t \|\mathbf{f}_\phi(\mathbf{x}_t, t) - \mathbf{x}\|^2]$$

where w_t is the weight used to “balance perceptual quality and diversity.”

- Given a diffusion model, a sample can be generated deterministically from a Gaussian noise by using the DDIM sampler [SME20], which has the following update rule:

$$\mathbf{x}_s = \frac{\sigma_s}{\sigma_t} \mathbf{x}_t + \left(\alpha_s - \alpha_t \frac{\sigma_s}{\sigma_t} \right) \mathbf{f}_\phi(\mathbf{x}_t, t)$$

where $s < t$. The process starts with sampling $\mathbf{x}_T = \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)$, and then we can obtain \mathbf{x}_t with smaller and smaller t until we reach $t \approx 0$. The catch is that the steps size $\delta = t - s$ must be small enough.

2.2 Knowledge Distillation

- The student model is denoted by \mathbf{g}_θ . This is used in contrast with the teacher diffusion model \mathbf{f}_ϕ .
- The most straightforward form of distillation is direct distillation where the student model is trained with the following loss function:

$$\mathcal{L}_\theta^{\text{direct}} = E_{\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\mathbf{g}_\theta(\boldsymbol{\xi}) - \text{ODE-Solver}(\mathbf{f}_\phi, \boldsymbol{\xi}, T \rightarrow 0)\|^2]$$

where **ODE-Solver** is any sampler like the DDIM sampler in the last section.

- The drawback of the above approach is that the **ODE-Solver** needs many steps in order to make the student model generate high quality data. This can make the training process slow.
- However, notice that the approach does not require access to training data at all.

- Other approaches such as progressive distillation [SH22], consistency models [SDCS23], and TRACT [BAL⁺23] avoid running the full **ODE-Solver** from T to 0.
- For the consistency model, the student model is conditioned on time, which means that we want $\mathbf{g}_\theta(\mathbf{x}_t, t)$ to predict \mathbf{x} . The model is trained with the following loss function:

$$\mathcal{L}_\theta^{\text{CM}} = E_{\mathbf{x} \sim p_{\text{data}}, s, t \sim [0, T], s < t, \mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x})} [\|\mathbf{g}_\theta(\mathbf{x}_t, t) - \mathbf{g}_{\theta_{\text{EMA}}}(\mathbf{x}_s, s)\|^2]$$

where $\mathbf{x}_s = \text{ODE-Solver}(\mathbf{f}_\phi, \mathbf{x}_t, t \rightarrow s)$ and θ_{EMA} is the exponential moving average of the student parameters θ .

- While training a consistency model does not require executing the **ODE-Solver** from start to finish, it requires access to the training data.

3 Method

- The goal of booth is to train a time-condition model $\mathbf{g}_\theta(\xi, t)$ that predicts $\mathbf{x}_t = \text{ODE-Solver}(\mathbf{f}_\phi, \theta, T \rightarrow t)$ when $\xi \sim \mathcal{N}(\mathbf{0}, I)$.
- We can generate a sample by evaluating $\mathbf{g}_\theta(\xi, 0)$ after sampling $\theta \sim \mathcal{N}(\mathbf{0}, I)$.
- Since the model always takes a Gaussian noise as input, there is no need for a dataset during the training process.

3.1 Signal-ODE

- Predicting \mathbf{x}_t is hard because it is a noisy image.
- It is much easier to predict $\mathbf{y}_t = (\mathbf{x}_t - \sigma_t \theta) / \alpha_t$, which is supposed to represent a predicted denoised image or the “signal component” of \mathbf{x}_t .
- Moreover, we know that

$$\begin{aligned} \mathbf{y}_s &= \frac{\mathbf{x}_s - \sigma_s \xi}{\alpha_s} \\ &= \frac{1}{\alpha_s} \mathbf{x}_s - \frac{\sigma_s}{\alpha_s} \xi \\ &= \frac{1}{\alpha_s} \left[\frac{\sigma_s}{\sigma_t} \mathbf{x}_t + \left(\alpha_s - \alpha_t \frac{\sigma_s}{\sigma_t} \right) \mathbf{f}_\phi(\mathbf{x}_t, t) \right] - \frac{\sigma_s}{\alpha_s} \xi \\ &= \left(1 - \frac{\sigma_s}{\alpha_s} \frac{\alpha_t}{\sigma_t} \right) \mathbf{f}_\phi(\mathbf{x}_t, t) + \frac{\sigma_s}{\alpha_s} \left(\frac{\mathbf{x}_t}{\sigma_t} - \xi \right) \\ &= \left(1 - \frac{\sigma_s}{\alpha_s} \frac{\alpha_t}{\sigma_t} \right) \mathbf{f}_\phi(\mathbf{x}_t, t) + \frac{\sigma_s}{\alpha_s} \left(\frac{\mathbf{x}_t - \sigma_t \xi}{\sigma_t} \right) \\ &= \left(1 - \frac{\sigma_s}{\alpha_s} \frac{\alpha_t}{\sigma_t} \right) \mathbf{f}_\phi(\mathbf{x}_t, t) + \frac{\sigma_s}{\alpha_s} \frac{\alpha_t}{\sigma_t} \left(\frac{\mathbf{x}_t - \sigma_t \xi}{\alpha_t} \right) \\ &= \left(1 - \frac{\sigma_s}{\alpha_s} \frac{\alpha_t}{\sigma_t} \right) \mathbf{f}_\phi(\mathbf{x}_t, t) + \frac{\sigma_s}{\alpha_s} \frac{\alpha_t}{\sigma_t} \mathbf{y}_t. \end{aligned}$$

Define $\lambda_t = -\log(\alpha_t / \sigma_t)$, we have that the above equation can be rewritten as:

$$\mathbf{y}_s = (1 - e^{\lambda_s - \lambda_t}) \mathbf{f}_\phi(\mathbf{x}_t, t) + e^{\lambda_s - \lambda_t} \mathbf{y}_t.$$

- The above equation can be turned into an ODE.

$$\begin{aligned}\mathbf{y}_s - \mathbf{y}_t &= (1 - e^{\lambda_s - \lambda_t}) \mathbf{f}_\phi(\mathbf{x}_t, t) - (1 - e^{\lambda_s - \lambda_t}) \mathbf{y}_t \\ \mathbf{y}_s - \mathbf{y}_t &= (1 - e^{\lambda_s - \lambda_t}) [\mathbf{f}_\phi(\mathbf{x}_t, t) - \mathbf{y}_t].\end{aligned}$$

Dividing both sides by $s - t$ and taking the limit as $s \rightarrow t^-$, we have that

$$\begin{aligned}\lim_{s \rightarrow t^-} \frac{\mathbf{y}_s - \mathbf{y}_t}{s - t} &= \left(\lim_{s \rightarrow t^-} \frac{1 - e^{\lambda_s - \lambda_t}}{s - t} \right) [\mathbf{f}_\phi(\mathbf{x}_t, t) - \mathbf{y}_t] \\ \frac{d\mathbf{y}_t}{dt} &= \left(\lim_{s \rightarrow t^-} \frac{1 - e^{\lambda_s - \lambda_t}}{s - t} \right) [\mathbf{f}_\phi(\mathbf{x}_t, t) - \mathbf{y}_t].\end{aligned}$$

The limit on the RHS is an indeterminate form $0/0$, so we will apply l'Hôpital's rule.

$$\lim_{s \rightarrow t^-} \frac{1 - e^{\lambda_s - \lambda_t}}{s - t} = \lim_{s \rightarrow t^-} \frac{\{1 - e^{\lambda_s - \lambda_t}\}'}{\{s - t\}'} = \lim_{s \rightarrow t^-} \frac{-e^{\lambda_s - \lambda_t} \{\lambda_s - \lambda_t\}'}{1} = \lim_{s \rightarrow t^-} \frac{-e^{\lambda_s - \lambda_t} \lambda_s'}{1} = -\lambda_t'.$$

As a result,

$$\frac{d\mathbf{y}_t}{dt} = -\lambda_t' [\mathbf{f}_\phi(\mathbf{x}_t, t) - \mathbf{y}_t].$$

The above ODE is called the **signal-ODE**. It is supposed to be integrated with the boundary condition $\mathbf{y}_T = \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)$ with time from $t = T$ to $t = 0$. Once we get \mathbf{y}_0 , we can output this as a sampled data because $\mathbf{y}_0 = \mathbf{x}_0$.

3.2 Learning

- We would like to train a neural network $\mathbf{y}_\theta(\boldsymbol{\xi}, t)$ that approximate \mathbf{y}_t . This network can be trained with the following loss:

$$\mathcal{L}_\theta^{\text{DE}} = E_{\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I), t \sim [0, T]} \left[\left\| \frac{d\mathbf{y}_\theta(\boldsymbol{\xi}, t)}{dt} + \lambda_t' (\mathbf{f}_\phi(\hat{\mathbf{x}}_t, t) - \mathbf{y}_\theta(\boldsymbol{\xi}, t)) \right\|^2 \right].$$

Here, we use the estimate $\hat{\mathbf{x}}_t = \alpha_t \mathbf{y}_\theta(\boldsymbol{\xi}, t) + \sigma_t \boldsymbol{\xi}$.

- While computing the time derivative $d\mathbf{y}_\theta(\boldsymbol{\xi}, t)/deet$ is possible with forward mode differentiation, computing the gradient through the derivative can be expensive. Hence, we can approximate the derivative with the difference equation:

$$\begin{aligned}\mathbf{y}_\theta(\boldsymbol{\xi}, s) &\approx \mathbf{y}_\theta(\boldsymbol{\xi}, t) - (s - t) \lambda_t' (\mathbf{f}_\phi(\hat{\mathbf{x}}_t, t) - \mathbf{y}_\theta(\boldsymbol{\xi}, t)) \\ &= \mathbf{y}_\theta(\boldsymbol{\xi}, t) + (t - s) \lambda_t' (\mathbf{f}_\phi(\hat{\mathbf{x}}_t, t) - \mathbf{y}_\theta(\boldsymbol{\xi}, t)) \\ &= \mathbf{y}_\theta(\boldsymbol{\xi}, t) + \delta \lambda_t' (\mathbf{f}_\phi(\hat{\mathbf{x}}_t, t) - \mathbf{y}_\theta(\boldsymbol{\xi}, t))\end{aligned}$$

where $\delta = t - s$ is the step size.

- So, we can train the network with the following loss instead:

$$\mathcal{L}_\theta^{\text{BS}} = E_{\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I), t \sim [\delta, T]} \left[\frac{\tilde{w}_t}{\delta^2} \left\| \mathbf{y}_\theta(\boldsymbol{\xi}, s) - \text{SG}[\mathbf{y}_\theta(\boldsymbol{\xi}, t) + \delta \lambda_t' (\mathbf{f}_\phi(\hat{\mathbf{x}}_t, t) - \mathbf{y}_\theta(\boldsymbol{\xi}, t))] \right\|^2 \right].$$

where $\text{SG}[\cdot]$ is the stop-gradient operator, and \tilde{w}_t is a time-dependent weight, which is implicit in how time t is sampled.

- A challenge in training \mathbf{y}_θ error accumulation: errors in prediction of \mathbf{y}_t might propagate to prediction of \mathbf{y}_s with $s < t$.

- The paper proposes two ways to mitigate error accumulation.
 - Sample time t uniformly, despite potential slow down in convergence.
 - Use higher-order solvers such as Heun’s method to compute the expected value of $\mathbf{y}_\theta(\boldsymbol{\xi}, s)$ instead of just using the Euler’s method like in $\mathcal{L}_\theta^{\text{BS}}$.
- Another problem that must be address is numerical issues caused by the fact that λ'_t is unbounded at $t = T$ and $t = 0$.
 - The student model must be trained in the truncated range $t \in [t_{\min}, t_{\max}]$.
 - We must also ensure that the student model behaves in the same way as the teacher model at t_{\max} . This is done by minimizing the loss:

$$\mathcal{L}_\theta^{\text{BC}} = E_{\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} \left[\|\mathbf{f}_\phi(\boldsymbol{\xi}, t_{\max}) - \mathbf{y}_\theta(\boldsymbol{\xi}, t_{\max})\|^2 \right]$$

- As a result, the overall training loss for BOOT is $\mathcal{L}_\theta = \mathcal{L}_{ves\theta}^{\text{BS}} + \beta \mathcal{L}_\theta^{\text{BC}}$ where β is a hyperparameter.

3.3 Adapting to Guided Diffusion Models

- Conditional diffusion models $\mathbf{f}_\phi(\mathbf{x}_t, t, \mathbf{c})$ are often trained to be able to perform classifier-free guidance:

$$\tilde{\mathbf{f}}_\phi(\mathbf{x}_t, t, \mathbf{c}) = \mathbf{f}_\phi(\mathbf{x}_t, t, \emptyset) + w(\mathbf{f}_\phi(\mathbf{x}_t, t, \mathbf{c}) - \mathbf{f}_\phi(\mathbf{x}_t, t, \emptyset))$$

where w is the guidance scale, and \emptyset denotes a conditioning input that makes the model unconditional.

- It is straightforward to train a (conditional) student model that follows the behavior of a teacher model with classifier free guidance. We can either
 - Fix the guidance scale w and use $\tilde{\mathbf{f}}_\phi(\mathbf{x}_t, t, \mathbf{c})$.
 - Like Meng et al. [MRG⁺23], train a student model that also receives w as input. (However, this requires architecture change.)

4 Experiments

- The authors performed experiments on the following datasets:
 - FFHQ 64×64
 - Class-conditional ImageNet 64×64 .
 - LSUN Bedroom 256×256 .
- For class-conditional ImageNet, the paper evaluates the student model on random guidance scale in the range $w \in [1, 5]$.
- The authors also distilled two text-to-image models: DeepFloyd-IF and Stable Diffusion.
 - Text prompts were obtained from DiffusionDB [WMM⁺23].
- In most experiments, the student model \mathbf{y}_θ would have the same architecture as the teacher model \mathbf{f}_θ , and the parameters of the student model would be initalized to those of the teacher models.
- The exception to the above practice is when the model must also be conditioned on the guidance scale w (i.e., class-conditonal ImageNet). When this is the case, w is taken into account via ADAIN layers.

- When looking at FID scores on the three datasets, BOOT-distilled models does not perform as well as 50-step DDIM sampler but better than 10-step DDIM sampler.
- Without the boundary condition loss, the student model’s outputs are consistently sharp across time steps, indicating mode collapse.
- We may train a BOOT model by progressively decreasing time during training. However, the paper found that this progressive-time scheme leads to more artifacts. The authors surmise that progressive-time training tends to accumulate irreversible errors.

References

- [BAL⁺23] David Berthelot, Arnaud Autef, Jierui Lin, Dian Ang Yap, Shuangfei Zhai, Siyuan Hu, Daniel Zheng, Walter Talbott, and Eric Gu. Tract: Denoising diffusion models with transitive closure time-distillation, 2023.
- [GZZ⁺23] Jiatao Gu, Shuangfei Zhai, Yizhe Zhang, Lingjie Liu, and Josh Susskind. Boot: Data-free distillation of denoising diffusion models with bootstrapping, 2023.
- [KAAL22] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models, 2022.
- [LL21] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed, 2021.
- [MRG⁺23] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik P. Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models, 2023.
- [SDCS23] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models, 2023.
- [SH22] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *CoRR*, abs/2202.00512, 2022.
- [SME20] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2020.
- [SSDK⁺21] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- [WMM⁺23] Zijie J. Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau. Diffusiondb: A large-scale prompt gallery dataset for text-to-image generative models, 2023.