

# Denoising Likelihood Score Matching for Conditional Score-based Data Generation

Chen-Hao Chao<sup>1,2</sup>, Wei-Fang Sun<sup>1</sup>, Bo-Wun Cheng<sup>1,2</sup>, Yi-Chen Lo<sup>2</sup>, Chia-Che Chang<sup>2</sup>, Yu-Lun Liu<sup>2</sup>, Yu-Lin Chang<sup>2</sup>, Chia-Ping Chen<sup>2</sup>, and Chun-Yi Lee<sup>1</sup>

<sup>1</sup> Elsa Lab, Department of Computer Science, National Tsing Hua University, Hsinchu, Taiwan

<sup>2</sup> MediaTek Inc., Hsinchu, Taiwan

## Abstract

Many existing conditional score-based data generation methods utilize Bayes' theorem to decompose the gradients of a log posterior density into a mixture of scores. These methods facilitate the training procedure of conditional score models, as a mixture of scores can be separately estimated using a score model and a classifier. However, our analysis indicates that the training objectives for the classifier in these methods may lead to a **serious score mismatch issue**, which corresponds to the situation that the estimated scores deviate from the true ones. To resolve it, we theoretically formulate a novel training objective, called **Denoising Likelihood Score Matching (DLSM)** loss, for the classifier to match the gradients of the true log likelihood density. Our experimental evidences show that the proposed method outperforms the previous methods on benchmarks noticeably in terms of several key evaluation metrics.

## Background

### Langevin Diffusion

Langevin diffusion [1,2] can be used to **generate data samples** from an unknown data distribution  $p(\tilde{\mathbf{x}})$  using only the score function  $\nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{x}})$ , where  $\tilde{\mathbf{x}} \in \mathbb{R}^d$  represents a data sample:

$$\tilde{\mathbf{x}}_t = \tilde{\mathbf{x}}_{t-1} + \frac{\epsilon^2}{2} \nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{x}}_{t-1}) + \epsilon \mathbf{z}_t, \quad (1)$$

where  $\tilde{\mathbf{x}}_0$  is sampled from an arbitrary distribution,  $\epsilon$  is a fixed positive step size, and  $\mathbf{z}_t$  is a noise vector sampled from a normal distribution  $N(\mathbf{0}, \mathbf{I}_{d \times d})$  for simulating a  $d$ -dimensional standard Brownian motion.

### Parzen Density Estimation

Previous literature [3] utilized Parzen density estimation to replace the **Dirac function**  $p(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \delta(\|\mathbf{x} - \mathbf{x}^{(i)}\|)$  with **isotropic Gaussian smoothing kernels**  $p_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \sigma^d} e^{-\frac{1}{2\sigma^2} \|\tilde{\mathbf{x}} - \mathbf{x}\|^2}$  with

variance  $\sigma^2$ , where  $\mathbf{x} \in \mathbb{R}^d$  represents a data sample, and  $\{\mathbf{x}^{(i)}\}_{i=1}^m$  are identically distributed data points. Specifically, Parzen density estimation enables the calculation of  $p(\tilde{\mathbf{x}}) = \frac{1}{m} \sum_{i=1}^m p_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x}^{(i)})$ .

### Denoising Score Matching

To train a score model  $s(\tilde{\mathbf{x}}; \phi)$ , a feasible approach is to minimize the Explicit Score-Matching (ESM) loss  $L_{\text{ESM}}$  [3], represented as:

$$L_{\text{ESM}} = \mathbb{E}_{p_{\sigma}(\tilde{\mathbf{x}})} \left[ \frac{1}{2} \|\mathbf{s}(\tilde{\mathbf{x}}; \phi) - \nabla_{\tilde{\mathbf{x}}} \log p_{\sigma}(\tilde{\mathbf{x}})\|^2 \right]. \quad (2)$$

Based on Parzen density estimation, an efficient alternative, called **Denoising Score-Matching (DSM)** loss [3], is proposed to efficiently calculate the equivalent loss  $L_{\text{DSM}}$ , expressed as:

$$L_{\text{DSM}} = \mathbb{E}_{p_{\sigma}(\tilde{\mathbf{x}}, \mathbf{x})} \left[ \frac{1}{2} \|\mathbf{s}(\tilde{\mathbf{x}}; \phi) - \nabla_{\tilde{\mathbf{x}}} \log p_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})\|^2 \right], \quad (3)$$

where  $\nabla_{\tilde{\mathbf{x}}} \log p_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})$  is simply  $\frac{1}{\sigma^2} (\mathbf{x} - \tilde{\mathbf{x}})$ .

### Conditional Score Decomposition

Score models can be extended to conditional models when conditioned on some extra information  $\tilde{\mathbf{y}}$ . A popular approach adopted by researchers utilizes **Bayes' theorem** to decompose the conditional score  $\nabla_{\tilde{\mathbf{x}}} \log p_{\sigma, \tau}(\tilde{\mathbf{x}}|\tilde{\mathbf{y}})$  as a **mixture of scores**, i.e.,

$$\nabla_{\tilde{\mathbf{x}}} \log p_{\sigma, \tau}(\tilde{\mathbf{x}}|\tilde{\mathbf{y}}; \theta, \phi) = \nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}; \theta) + \mathbf{s}(\tilde{\mathbf{x}}; \phi), \quad (4)$$

where  $p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}; \theta)$  is trained with  $L_{\text{CE}} = H(p_{\sigma, \tau}(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}), p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}; \theta))$ , and  $\mathbf{s}(\tilde{\mathbf{x}}; \phi)$  is trained with  $L_{\text{DSM}}$ .

## Motivational Example

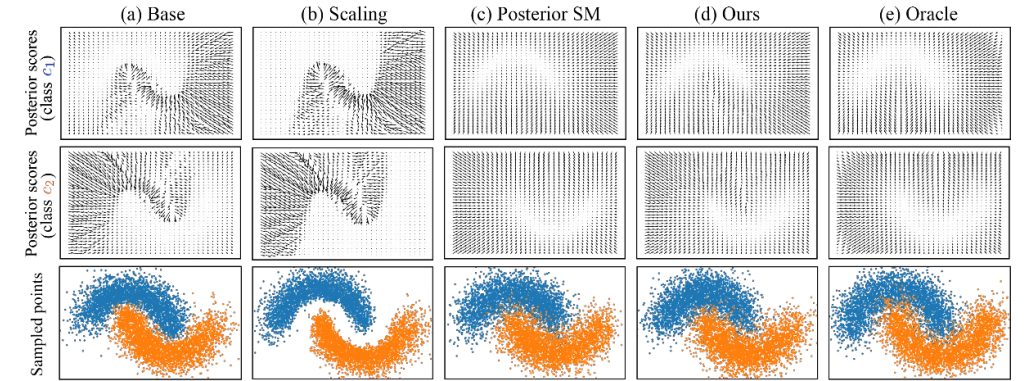


Figure 1. The visualized results on the inter-twining moon dataset.

- (a) **Base method**:  $\nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{x}}|\tilde{\mathbf{y}}; \theta, \phi) = \nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}; \theta) + \mathbf{s}(\tilde{\mathbf{x}}; \phi)$ , described in Eq. (4).  
(b) **Scaling method**:  $\nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{x}}|\tilde{\mathbf{y}}; \theta, \phi) = \alpha \nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}; \theta) + \mathbf{s}(\tilde{\mathbf{x}}; \phi)$ , where  $\alpha > 0$ .  
(c) **Posterior SM method**: The posterior scores for different class conditions, in this case,  $c_1$  and  $c_2$ , are separately estimated using different score models  $\mathbf{s}(\tilde{\mathbf{x}}; \phi_1)$  and  $\mathbf{s}(\tilde{\mathbf{x}}; \phi_2)$  trained with  $L_{\text{DSM}}$ .  
(d) **Ours**: The posterior scores are estimated in a similar fashion as method (a) except that the classifier is trained with the proposed loss function.  
(e) **Oracle**: The oracle posterior scores are directly computed.

## Methodology

### Denoising Likelihood Score Matching

A score model trained with the score-matching objective can potentially be beneficial in producing a better posterior score estimation. In light of this, a classifier may be enhanced if the score-matching process is involved during its training procedure. An intuitive way to accomplish this aim is through minimizing the explicit likelihood score-matching loss  $L_{\text{ELSM}}$ , which is defined as the following:

$$L_{\text{ELSM}} = \mathbb{E}_{p_{\sigma, \tau}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})} \left[ \frac{1}{2} \|\nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}; \theta) - \nabla_{\tilde{\mathbf{x}}} \log p_{\sigma, \tau}(\tilde{\mathbf{y}}|\tilde{\mathbf{x}})\|^2 \right]. \quad (5)$$

This loss term, however, requires the derivation of the true likelihood score  $\nabla_{\tilde{\mathbf{x}}} \log p_{\sigma, \tau}(\tilde{\mathbf{y}}|\tilde{\mathbf{x}})$ , which is computationally unaffordable during the classifier training. To reduce the computational cost, we follow the derivation of DSM as well as Bayes' theorem, and formulate an alternative objective  $L_{\text{DLSM}}$ :

$$L_{\text{DLSM}} = \mathbb{E}_{p_{\sigma, \tau}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}, \mathbf{y})} \left[ \frac{1}{2} \|\nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}; \theta) + \nabla_{\tilde{\mathbf{x}}} \log p_{\sigma}(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log p_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})\|^2 \right]. \quad (6)$$

In contrast to  $L_{\text{ELSM}}$ ,  $L_{\text{DLSM}}$  is more **computationally feasible** since  $\nabla_{\tilde{\mathbf{x}}} \log p_{\sigma}(\tilde{\mathbf{x}})$  can be estimated using a score model  $\mathbf{s}(\tilde{\mathbf{x}}; \phi)$  trained with  $L_{\text{DSM}}$ , and  $\nabla_{\tilde{\mathbf{x}}} \log p_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})$  equals to  $\frac{1}{\sigma^2} (\mathbf{x} - \tilde{\mathbf{x}})$ .

### The Introduction of the DLSM Loss in the Classifier Training Process

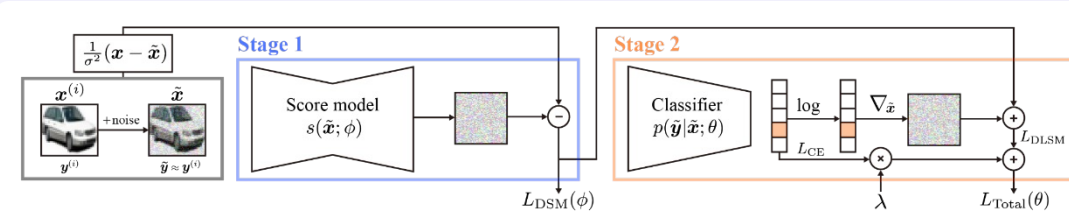


Figure 2. The training procedure of the proposed methodology.

In practice, the total training objective of the classifier is written as follows:

$$L_{\text{Total}} = L_{\text{DLSM}} + \lambda L_{\text{CE}}, \quad (7)$$

where  $\lambda > 0$  is a balancing coefficient. Fig. 2 depicts a **two-stage training procedure** adopted in this work. In stage 1, a score model  $\mathbf{s}(\tilde{\mathbf{x}}; \phi)$  is updated using  $L_{\text{DLSM}}$  to match  $\nabla_{\tilde{\mathbf{x}}} \log p_{\sigma}(\tilde{\mathbf{x}})$ . In stage 2, the weights of the trained score model are fixed, and a classifier  $\nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}; \theta)$  is updated using  $L_{\text{DLSM}}$ . After these two training stages,  $\mathbf{s}(\tilde{\mathbf{x}}; \phi)$  and  $\nabla_{\tilde{\mathbf{x}}} \log p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}; \theta)$  can then be added together based on Eq. (4) to perform conditional sampling.

## Experimental Results

		Cifar-10			Cifar-100		
		Base	Scaling	Ours	Base	Scaling	Ours
FID	↓	4.10	12.48	<b>2.25</b>	4.52	12.58	<b>3.86</b>
IS	↑	9.08	9.37	<b>9.90</b>	11.53	11.59	<b>11.62</b>
Precision	↑	0.67	<b>0.75</b>	0.65	0.62	<b>0.65</b>	0.61
Recall	↑	0.61	0.49	<b>0.62</b>	0.62	0.52	<b>0.63</b>
Density	↑	1.05	<b>1.36</b>	0.96	0.84	<b>0.93</b>	0.82
Coverage	↑	0.80	0.75	<b>0.81</b>	<b>0.71</b>	0.61	<b>0.71</b>
CAS	↑	0.38	0.46	<b>0.58</b>	0.15	0.24	<b>0.28</b>
(CW) Precision	↑	0.51	<b>0.70</b>	0.56	0.33	<b>0.59</b>	0.42
(CW) Recall	↑	0.59	0.42	<b>0.61</b>	<b>0.68</b>	0.41	<b>0.68</b>
(CW) Density	↑	0.63	<b>1.23</b>	0.76	0.30	<b>0.78</b>	0.43
(CW) Coverage	↑	0.60	0.66	<b>0.71</b>	0.38	0.51	<b>0.52</b>

Table 1. The evaluation results on the Cifar-10 and Cifar-100 datasets. The P / R / D / C metrics with '(CW)' in the last four rows represents the average class-wise P / R / D / C metrics. The arrow symbols ↑ / ↓ represent that a higher / lower evaluation result correspond to a better performance.

- Our method outperforms the other two methods with substantial margins in terms of **Frechet Inception Distance (FID)** [5] and **Inception Score (IS)** [6], indicating that the generated samples bear closer resemblance to the real data.
- The evaluation results on the **classification accuracy score (CAS)** [7] and the **class-wise (CW) Precision, Recall, Density, and Coverage (P / R / D / C)** [8] metrics suggest that our method offers a superior ability for a classifier to learn accurate class information as compared to the base method.

[1] G. O. Roberts and J. S. Rosenthal. Optimal scaling of discrete approximations to langevin diffusions. Journal of the Royal Statistical Society. 1998.  
[2] G. O. Roberts and R. L. Tweedie. Exponential convergence of Langevin distributions and their discrete approximations. Bernoulli. 1996.

[3] P. Vincent. A connection between score matching and denoising autoencoders. Neural Computation. 2011.  
[4] P. Dhariwal and A. Nichol. Diffusion models beat GANs on image synthesis. arXiv:2105.05233. 2021.  
[5] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. NeurIPS. 2017.

[6] S. Barratt and R. Sharma. A note on the inception score. ICML Workshop. 2018.  
[7] S. Ravuri and O. Vinyals. Classification accuracy score for conditional generative models. NeurIPS. 2019.  
[8] M. F. Naeem, S. J. Oh, Y. Uh, Y. Choi, and J. Yoo. Reliable fidelity and diversity metrics for generative models. ICML. 2020.