

DDPM Score Matching and Distribution Learning



Sinho Chewi

Yale University



Alkis Kalavasis

Yale University



Omar Montasser

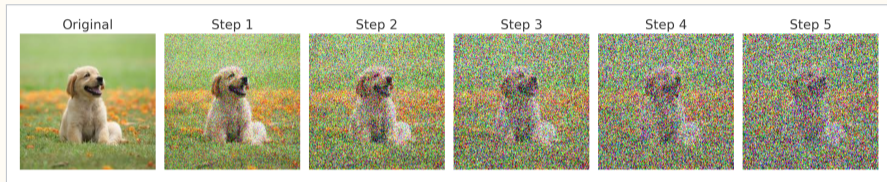
Yale University

Anay Mehrotra

Stanford University

COLT, July 2, 2026

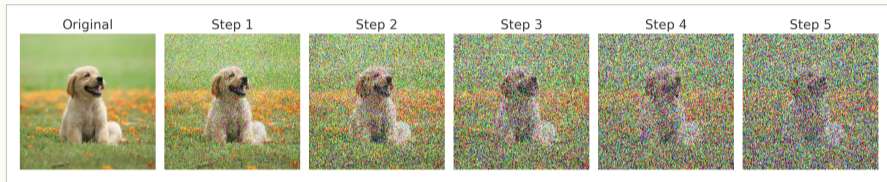
DDPM Diffusion Models



forward noising



DDPM Diffusion Models



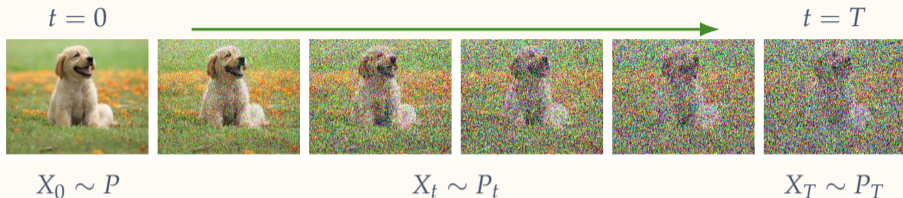
forward noising



reverse generation via scores



DDPM Diffusion Models



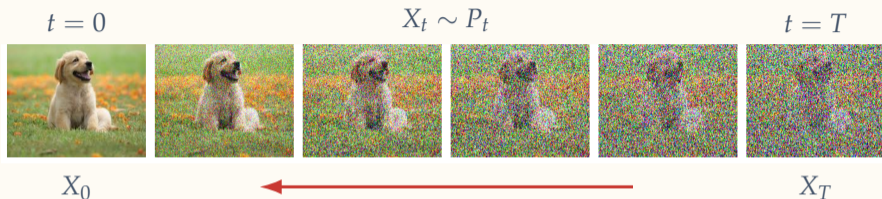
Noising process.

Start with data $X_0 \sim P$ and slowly add Gaussian noise:

$$\begin{aligned} X_t &= e^{-t} X_0 + \sqrt{1 - e^{-2t}} Z_t, & Z_t &\sim \mathcal{N}(0, I_d), \\ dX_t &= -X_t dt + \sqrt{2} dB_t, & X_0 &\sim P. \end{aligned}$$

This SDE smoothly transforms complex data-distribution P to $\mathcal{N}(0, I)$ by injecting noise.

DDPM Diffusion Models



Forward SDE can be **time-reversed** given accurate estimates of scores.

Reverse SDE with terminal time T .

$$dX_{T-t} = \{X_{T-t} + 2\nabla \log P_{T-t}(X_{T-t})\} dt + \sqrt{2} d\bar{B}_t, \quad X_T \sim P_T.$$

What does it mean to learn a distribution?

Question.

Let \mathcal{P} be a class of distributions over \mathbb{R}^d . Given i.i.d. samples from P , what does it mean to “learn” P ?

Density Estimation

Density Evaluator [Kearns–Mansour–Ron–Rubinfeld–Schapire–Sellie, STOC'94].

An *evaluator* \hat{P} for P is a function $\mathbb{R}^d \rightarrow \mathbb{R}$ such that on any input $x \in \mathbb{R}^d$, it outputs an estimate $\hat{P}(x)$ of the log-density $\log P(x)$.

Density Estimation

Density Evaluator [Kearns–Mansour–Ron–Rubinfeld–Schapire–Sellie, STOC'94].

An *evaluator* \hat{P} for P is a function $\mathbb{R}^d \rightarrow \mathbb{R}$ such that on any input $x \in \mathbb{R}^d$, it outputs an estimate $\hat{P}(x)$ of the log-density $\log P(x)$.

Let \mathcal{P} be a class of distributions over \mathbb{R}^d .

Density Estimation.

Given i.i.d. samples from $P \in \mathcal{P}$, compute an ε -evaluator \hat{P} of P

$$\mathbb{E}_{x \sim P} \left| (-\log P(x)) - \hat{P}(x) \right| \leq \varepsilon.$$

Learning to Sample

Generator [Kearns–Mansour–Ron–Rubinfeld–Schapire–Sellie, STOC'94].

A *generator* \hat{P} for P is a function $\mathbb{R}^d \rightarrow \mathbb{R}^d$ such that on random seed $z \sim \mathcal{N}(0, I)$, it outputs $x \sim \hat{P}(z)$ such that x has law close to P .

Learning to Sample

Generator [Kearns–Mansour–Ron–Rubinfeld–Schapire–Sellie, STOC'94].

A generator \hat{P} for P is a function $\mathbb{R}^d \rightarrow \mathbb{R}^d$ such that on random seed $z \sim \mathcal{N}(0, I)$, it outputs $x \sim \hat{P}(z)$ such that x has law close to P .

Let \mathcal{P} be a class of distributions over \mathbb{R}^d .

Learning to Sample.

Given i.i.d. samples S from $P \in \mathcal{P}$, compute an ε -generator \hat{P} of P

$$\text{Law}(\hat{P}(Z) \mid S) \approx_{\varepsilon} P \quad \text{on random seed } Z.$$

What does it mean to learn a distribution?

Question.

Let \mathcal{P} be a class of distributions over \mathbb{R}^d . Given i.i.d. samples from P , what does it mean to “learn” P ?

Learning to Sample.

Given a random seed z , output a sample close to P .

Density Estimation.

Given x , estimate the log-density $\log P(x)$.

Different Lenses on Distribution Learning

Statistical lens.

With sufficient time and samples, generators \Leftrightarrow evaluators.

Different Lenses on Distribution Learning

Statistical lens.

With sufficient time and samples, generators \rightleftarrows evaluators.

Computational lens.

With computational restrictions, generators $\not\leftrightarrow$ evaluators.

Score Estimation \rightarrow Sampling

Informal Theorem [Chen et al., ICLR'23; Lee-Lu-Tan, ALT'23].

Suppose score is L -Lipschitz, $m_2(P) \leq \text{poly}(d)$, and score estimates satisfy

$$\sup_{t \geq 0} \mathbb{E}_{X_t \sim P_t} \|\hat{s}_t(X_t) - \nabla \log P_t(X_t)\|^2 \leq \tilde{O}(\varepsilon^2).$$

Then, DDPM using \hat{s}_t outputs a distribution \hat{P} such that

$$d_{\text{TV}}(\hat{P}, P) \leq \varepsilon.$$

It uses $\text{poly}(d, L, 1/\varepsilon)$ score-oracle calls and additional computation.

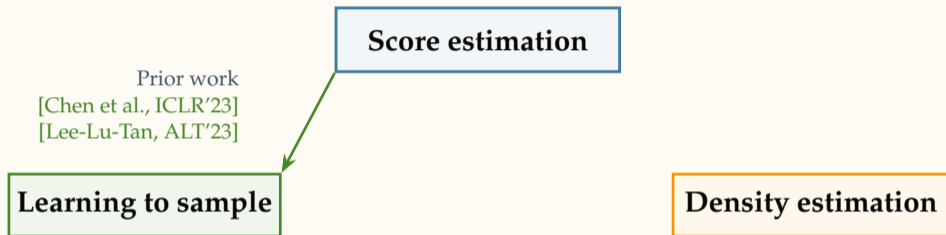


Takeaway

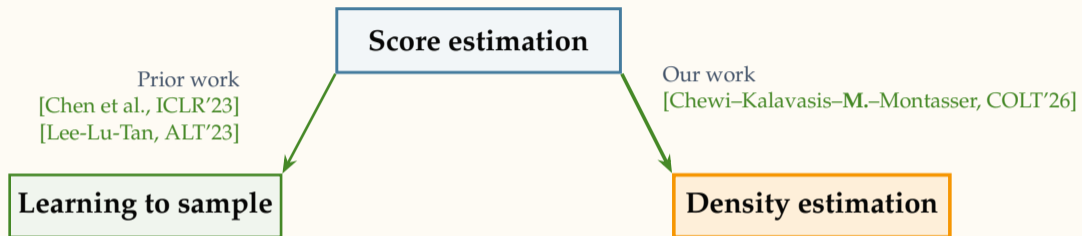
Learning to sample

Density estimation

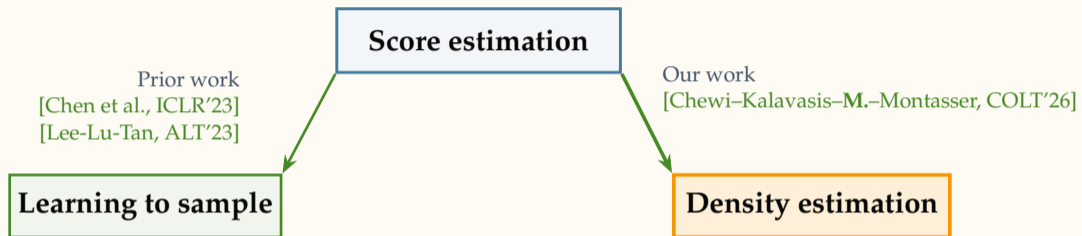
Takeaway




Takeaway



Takeaway



 **Takeaway.** Accurate DDPM score estimation \rightarrow log-density evaluator.

Our Result

Informal Theorem [Chewi–Kalavasis–M.–Montasser, COLT'26].

Let scores are L -Lipschitz, $m_2(P) \leq \text{poly}(d)$, and a score-estimation, for suitable $0 < \tau < T$,

$$\int_{\tau}^T \mathbb{E}_{X_t \sim P_t} \|\hat{s}_t(X_t) - \nabla \log P_t(X_t)\|^2 dt \leq \tilde{O}(\varepsilon^2/d).$$

Then there is an evaluator $\hat{\ell}$ for P such that

$$\mathbb{E}_{x_0 \sim P} |\hat{\ell}(x_0) - \log P(x_0)| \leq \varepsilon.$$

It uses $\text{poly}(d, L, 1/\varepsilon)$ score-oracle calls and additional computation.



Algorithmic Implications

Corollary 1.

Efficient score estimation \longrightarrow Efficient density estimation

Gaussian location mixtures.

Let \mathcal{P} be the non-parametric class of Gaussian location mixtures on \mathbb{R}^d : $P = Q \star \mathcal{N}(0, \sigma^2 I_d)$, where Q is supported on a union of k constant-radius balls, each carrying mass at least w_{\min} .

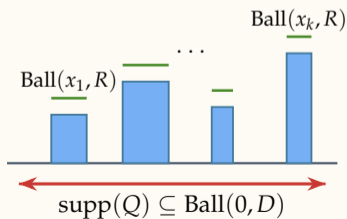
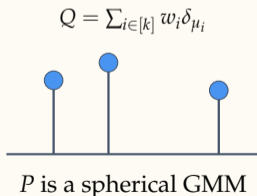
Algorithmic Implications

Corollary 1.

Efficient score estimation \rightarrow Efficient density estimation

Gaussian location mixtures.

Let \mathcal{P} be the non-parametric class of Gaussian location mixtures on \mathbb{R}^d : $P = Q \star \mathcal{N}(0, \sigma^2 I_d)$, where Q is supported on a union of k constant-radius balls, each carrying mass at least w_{\min} .



Algorithmic Implications

Let P be a Gaussian location mixture on \mathbb{R}^d with k components and $w_{\min} > \varepsilon$.

[Gatmiry–Kelner–Lee, COLT'25] gave a $d^{\text{polylog}(d,k,1/\varepsilon)}$ algorithm to learn to sample from P .

Algorithmic Implications

Let P be a Gaussian location mixture on \mathbb{R}^d with k components and $w_{\min} > \varepsilon$.

[Gatmiry–Kelner–Lee, COLT'25] gave a $d^{\text{polylog}(d,k,1/\varepsilon)}$ algorithm to learn to sample from P .

Theorem (Score Estimation) [Gatmiry–Kelner–Lee, COLT'25].

There is a $d^{\text{polylog}(d,k,1/\varepsilon)}$ time algorithm for estimating the *scores* along the DDPM path, using samples from P .

Algorithmic Implications

Let P be a Gaussian location mixture on \mathbb{R}^d with k components and $w_{\min} > \varepsilon$.

[Gatmiry–Kelner–Lee, COLT'25] gave a $d^{\text{polylog}(d,k,1/\varepsilon)}$ algorithm to learn to sample from P .

Theorem (Score Estimation) [Gatmiry–Kelner–Lee, COLT'25].

There is a $d^{\text{polylog}(d,k,1/\varepsilon)}$ time algorithm for estimating the *scores* along the DDPM path, using samples from P .

Theorem (Density Estimation) [Chewi–Kalavasis–M.–Montasser, COLT'26].

There is an algorithm that outputs an ε -evaluator for P in $d^{\text{polylog}(d,k,1/\varepsilon)}$ time, using samples from P .

Lower Bound Implications

Corollary 2.

If the family of distributions \mathcal{P} satisfies regularity conditions, then:

Lower bound for density estimation for \mathcal{P}

→ Lower bound for score estimation for \mathcal{P}

Instantiating this for GMMs conceptually recovers [Song, NeurIPS'24]'s result

Conclusion and open problems

Score estimation is not only sufficient for sampling but also for PAC log-density estimation

Conclusion and open problems

Score estimation is not only sufficient for sampling but also for PAC log-density estimation

Using score estimation, is it possible to output a *proper* density estimator?

Analogous results in discrete domains, infinite-dimensional spaces, manifolds?

Algorithms for density estimation of well-conditioned GMMs?

Beyond GMMs, when is score estimation cryptographically hard?