

# A Variational Perspective on Diffusion-Based Generative Models and Score Matching ?

Imane Elbacha, Nathan Weill

MVA 22-23

24/05/2023

# Content

## 1 Introduction

## 2 Paper analysis

- Likelihood of diffusion processes
- Connection to score matching

## 3 Experiments

- Model overview
- Datasets
- Results analysis

## 4 Conclusion

# Motivation

- Providing a theoretical grounding in MLE of denoising diffusion probabilistic models (DDPM)

# Main contributions

- Estimate the likelihood of continuous-time diffusion models
- Derive an efficient & scalable training method
- Proving the equivalence between maximizing the CT-ELBO and score matching

# Likelihood of diffusion processes

Neural ODE & continuous normalizing flows

- ① The flow is defined as the solution of  
 $dX = \mu(X, t)dt; \quad X(0) \sim p_0$
- ② Its density is tractable with the instantaneous change of variable :  
 $p_T(x) = p_0(Y(1))e^{\int_0^T -\nabla \cdot \mu(Y(s), T-s)ds}$   
 $dY = -\mu(Y, T-s)ds; \quad Y(0) = x$
- ③ Need to solve the ODE with numerical integration, not scalable

# Likelihood of diffusion processes

## Density of diffusion processes

- ① Now take a stochastic diffusion process instead defined by :  
$$dX = \mu dt + \sigma dB_t$$
- ② E.g. Ornstein-Uhlenbeck or Geometric Brownian Motion
- ③ The density is not tractable anymore but the authors derive a continuous lower bound !

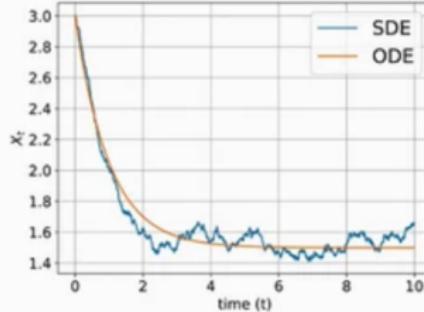


Figure – ODE vs SDE evolution

# Likelihood of diffusion processes

## CT-ELBO

- $dX = \mu(X, t)dt + \sigma(t)dB_t$
- $X_0 \sim p(\cdot, 0)$
- Fokker-Planck PDE :
$$\begin{aligned}\partial_t p(x, t) &= -\nabla \cdot (p(x, t)\mu(x, t)) + \frac{1}{2}\sigma^2(t)\Delta p(x, t) \\ &= -(\nabla \cdot \mu)p - \mu^\top \nabla p + \frac{1}{2}\sigma^2\Delta p\end{aligned}$$
- And Feynman-Kac representation as an expectation :
$$p(x, T) = \mathbb{E} \left[ p(Y_T, 0) e^{\int_0^T -\nabla \cdot \mu(Y_s, T-s) ds} \mid Y_0 = x \right]$$
$$dY = -\mu(Y, T-s)ds + \sigma(T-s)dB'_s$$
$$Y_0 = x$$

# Likelihood of diffusion processes

## CT-ELBO (1)

- still intractable so variational inference

$$\log p(x, T) = \log \mathbb{E}_{Y \sim \mathbb{P}} \left[ p(Y_T, 0) e^{\int_0^T -\nabla \cdot \mu(Y_s, T-s) ds} \mid Y_0 = x \right]$$

Marginalization

$$= \log \mathbb{E}_{Y \sim \mathbb{Q}} \left[ \frac{d\mathbb{P}}{d\mathbb{Q}} p(Y_T, 0) e^{\int_0^T -\nabla \cdot \mu(Y_s, T-s) ds} \mid Y_0 = x \right]$$

Change of measure

$$\geq \mathbb{E}_{Y \sim \mathbb{Q}} \left[ \log \frac{d\mathbb{P}}{d\mathbb{Q}} + \log p(Y_T, 0) - \int_0^T \nabla \cdot \mu ds \mid Y_0 = x \right]$$

Jensen's inequality

KL divergence

# Likelihood of diffusion processes

## CT-ELBO (2)

- Using Girsanov theorem

$$\left. \begin{array}{l} \mathbb{P}: dY = -\mu ds + \sigma dB'_s \\ \mathbb{Q}: dY = (-\mu + \sigma a)ds + \sigma d\hat{B}_s \end{array} \right\} \Rightarrow D_{\text{KL}}(\mathbb{Q} || \mathbb{P}) := \mathbb{E}_{\mathbb{Q}} \left[ \log \frac{d\mathbb{Q}}{d\mathbb{P}} \right] = \mathbb{E}_{\mathbb{Q}} \left[ \int_0^T \|a(Y_s, s)\|^2 ds \right]$$

Inference SDE

$$\begin{aligned} \log p(x, T) &\geq \mathbb{E}_{Y \sim \mathbb{Q}} \left[ \log \frac{d\mathbb{P}}{d\mathbb{Q}} + \log p(Y_T, 0) - \int_0^T \nabla \cdot \mu ds \mid Y_0 = x \right] \\ &= \mathbb{E}_{Y \sim \mathbb{Q}} \left[ \log p(Y_T, 0) - \int_0^T \nabla \cdot \mu + \|a\|^2 ds \mid Y_0 = x \right] \end{aligned}$$

25

Figure – Girsanov theorem

# Likelihood of diffusion processes

## CT-ELBO (3)

- MCMC sampling

$$\begin{aligned}\mathcal{E}^\infty &:= \mathbb{E}_{Y \sim Q} \left[ \log p(Y_T, 0) - \int_0^T \nabla \cdot \mu + \|a\|^2 dr \middle| Y_0 = x \right] \\ &= \mathbb{E}_{Y_T} [\log p(Y_T, 0)] - \int_0^T \mathbb{E}_{Y_s} [\|a\|^2 + \nabla \cdot (ga - f) | Y_0 = x] ds\end{aligned}$$

No need to integrate if we can sample

Figure – MCMC sampling

# Likelihood of diffusion processes

## CT-ELBO (4)

- Noisier but faster than numerical integration

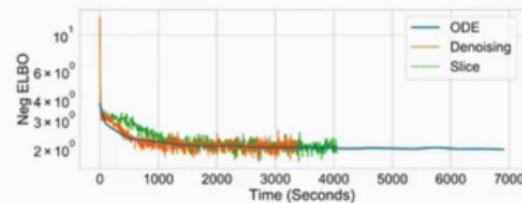
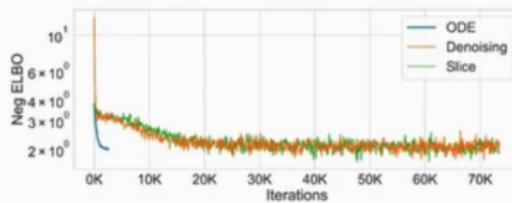


Figure – Negative ELBO evolution over training

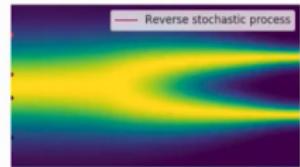
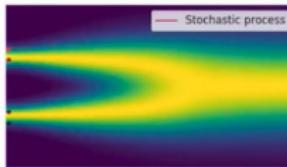
### Connection to score matching

The reverse plug-in SDE

- score matching with specific weighting and with ISM loss is found to be equivalent to maximizing the CT-ELBO of the reverse plug-in SDE

$$\begin{aligned} & (\text{setting } a = g\mathbf{s}_\theta) \\ -\mathbb{E}_{Y_0}[\mathcal{E}^\infty] & \stackrel{+\text{cst}}{=} \int_0^T g^2 \mathbb{E}_{Y_s} [\|\mathbf{s}_\theta\|^2 + \nabla \cdot \mathbf{s}_\theta] \, ds \quad (\text{Implicit score matching}) \\ & \stackrel{+\text{cst}}{=} \int_0^T g^2 \mathbb{E}_{Y_s} [\|\mathbf{s}_\theta(Y_s, s) - \nabla \log q(Y_s, s)\|^2] \, ds \quad (\text{Explicit score matching}) \quad [\text{Hyvärinen 05}^*] \end{aligned}$$

$$dY = f ds + g d\hat{B}_s \quad dX = \left( g^2 \nabla \log q(y, s) - f \right) dt + g dB_t$$



# Connection to score matching

Parametrized SDE family

$$dX = \left( \left( 1 - \frac{\lambda}{2} \right) g^2 \frac{\approx s_\theta(y, s)}{\nabla \log q} - f \right) dt + \sqrt{1 - \lambda} g dB_t \quad dY = \left( f - \frac{\lambda}{2} g^2 \nabla \log q \right) ds + \sqrt{1 - \lambda} g d\hat{B}_s$$

- $\lambda = 0$  corresponds to the previous SDE
  - $\lambda = 1$  corresponds to the ODE
  - Connection to what we saw in class (same marginal densities but different dynamics)
- As a consequence the two following dynamics have the **same** marginal densities.
- ▶  $d\mathbf{X}_t = b(t, \mathbf{X}_t)dt + c^{1/2}dB_t$
  - ▶  $d\mathbf{X}_t = \{b(t, \mathbf{X}_t) - \frac{c}{2}\nabla \log p_t(\mathbf{X}_t)\}dt$ .
  - ▶ One is **deterministic**, the other is **stochastic**.

# Connection to score matching

Link with discrete ELBO

- we saw in class variational inference in discrete formulation
- arose a connection with the DSM loss
- Similar result found in ?

# Experiments

## Overview

- ① Implement reverse SDE training on three datasets : moon dataset, Swiss roll dataset, and S-curve dataset.
- ② Use the plugging reverse SDE training to create respective models for each dataset.
- ③ Perform sampling on the generated models using Euler-Mayumar method.

# Experiments

## Implementation of plugging reverse SDE

- VariancePreservingSDE : Implements a variance-preserving stochastic differential equation (SDE) with key parameters, drift and diffusion functions, and sampling methods.
- MLP : Represents a customizable multilayer perceptron architecture with Swish activation function, used for modeling the drift function in PluginReverseSDE.
- PluginReverseSDE : Defines a reverse SDE by inverting a base SDE, utilizing drift and diffusion functions from both SDEs. Provides methods for denoising score matching loss and estimating the evidence lower bound (ELBO).

# Experiments

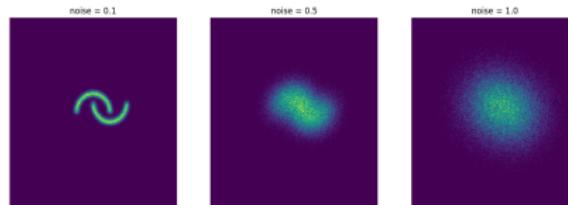
## Advantages of plugging reverse SDE

- The VariancePreservingSDE ensures accurate representation of the distribution by preserving the variance of the underlying process during sampling.
- The MLP architecture within the PluginReverseSDE effectively models complex patterns, enabling the generation of samples from the desired distribution.
- The denoising score matching loss in the PluginReverseSDE further refines the generative model by estimating the discrepancy between true and estimated distributions, improving the quality of generated samples.

# Experiments

## Datasets

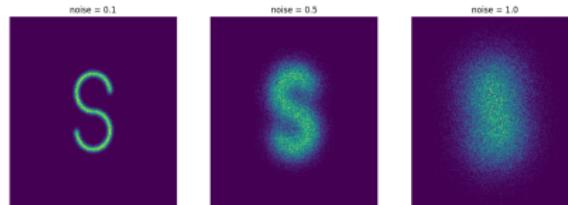
- Moon dataset : Used to evaluate the model's performance in handling non-linear and complex patterns, providing insights into its ability to capture intricate relationships and make accurate predictions in challenging scenarios.



# Experiments

## Datasets

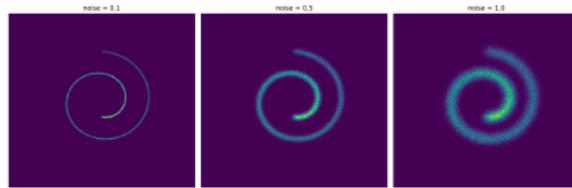
- S curve dataset : Employed to assess the model's capability to handle data with abrupt changes and sharp transitions, examining its adaptability to sudden variations and its capacity to capture the underlying dynamics accurately.



# Experiments

## Datasets

- Swiss roll dataset : Utilized to evaluate the model's ability to handle manifold learning and non-linear dimensionality reduction, examining its effectiveness in capturing the underlying structure of the data and recovering the original geometric shape.



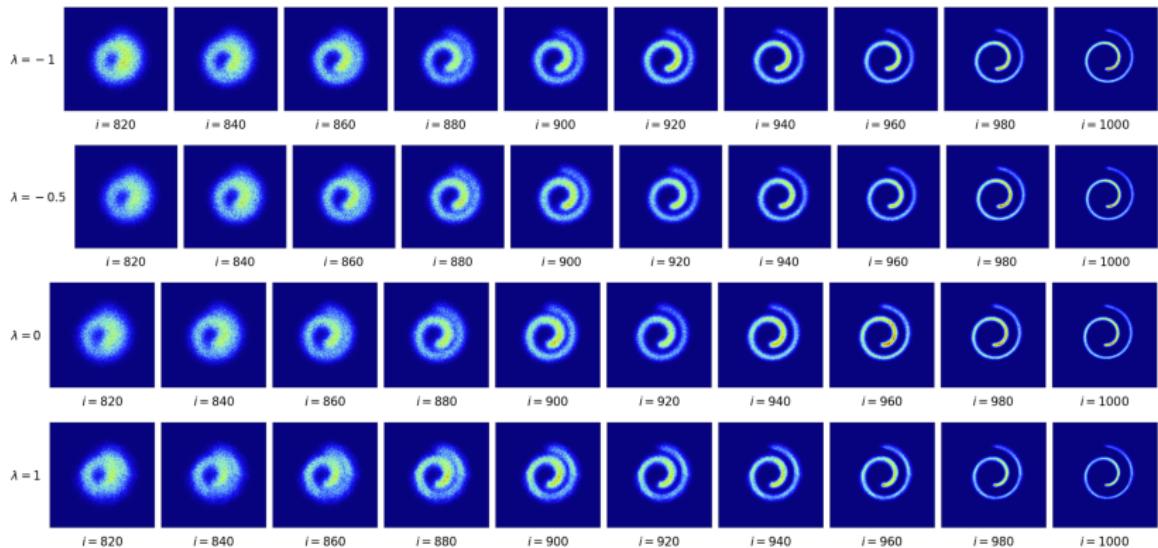
# Experiments

## Sampling results

- The proposed SDE model was applied to three distinct datasets : Moon dataset, Swiss Roll dataset, and S-curve dataset.
- The model demonstrated promising results for the Swiss Roll and S-curve datasets, accurately capturing the patterns and structures present in the data.
- However, the model encountered difficulties in generating samples that accurately represented the distribution of the Moon dataset, indicating challenges in capturing its unique shape and multi-modal nature.

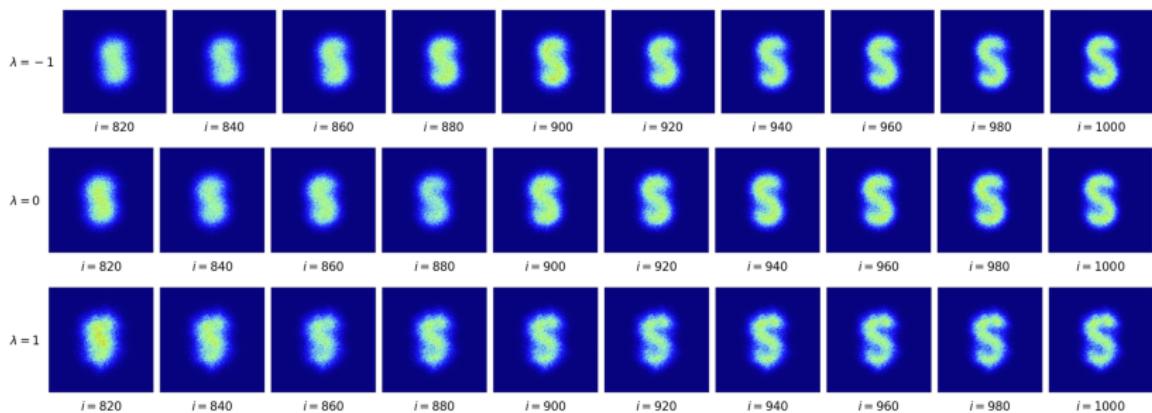
# Experiments

## Sampling results



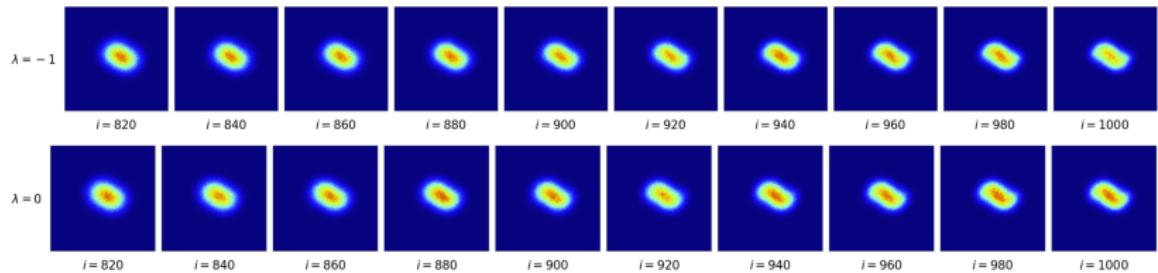
# Experiments

## Sampling results



# Experiments

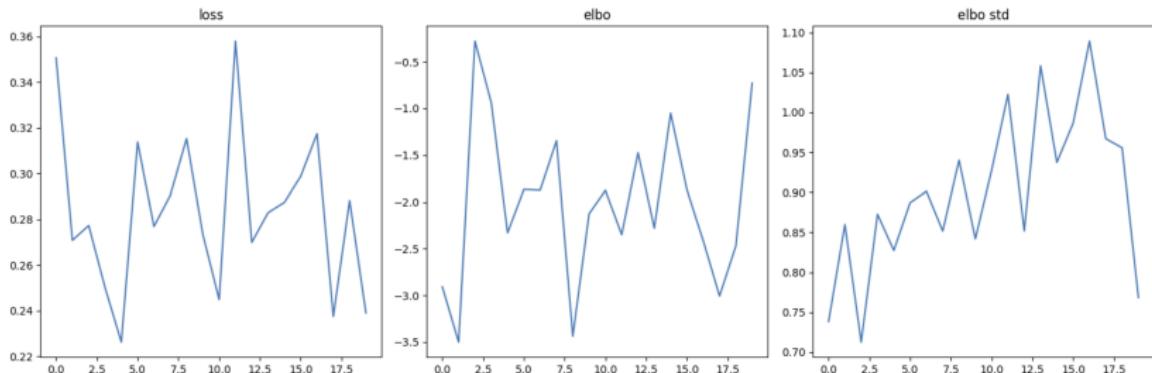
## Sampling results



# Experiments

## Loss and Elbo discrepancy

- A gap was observed between the score matching loss and the ELBO score during training, indicating a discrepancy between the quality of individual samples and the model's ability to capture the global characteristics of the data distribution.



# Conclusion

What is the importance of a variational framework ?

- The variational framework connects the score-matching loss to maximum likelihood estimation in the plug-in reverse diffusion equations.
- It provides a specific weighting in the importance sampling (ISM) loss for achieving maximum likelihood in generative modeling.
- Different generative models, such as normalizing flows and diffusion models, are interconnected.
- Beyond likelihood, alternative loss functions like quantile loss offer less restrictive parameter training, as suggested by Ostrovski et al. (2018).

# References I