

Dynamical Systems in Machine Learning

Cédric Simal

23/02/2023

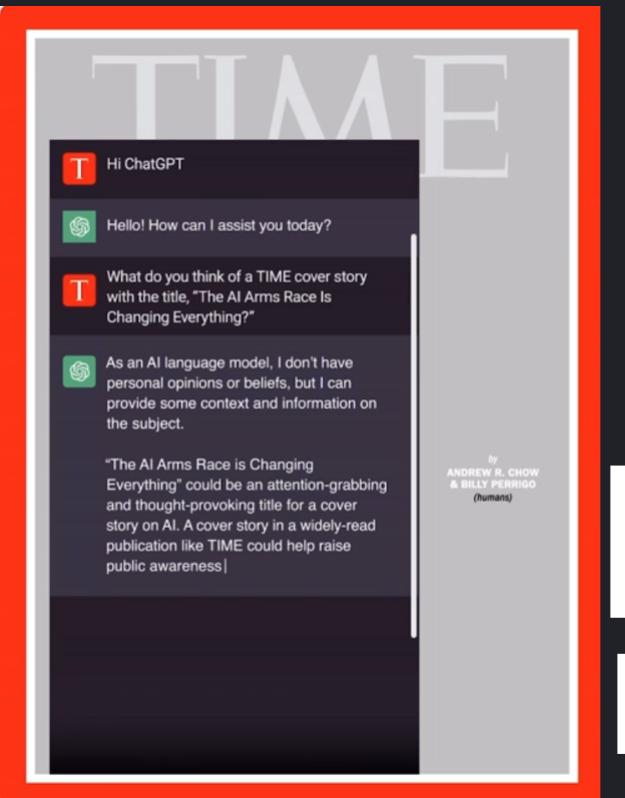
PSA



"A parody uncle sam poster with the message "I WANT YOU for SYSDYN", parody poster, digital painting, HQ, 4k", Made with Stable Diffusion

Motivation

The encroaching critical transition



ROBBIE THE WRITING ROBOT —

BuzzFeed preps AI-written content while CNET fumbles

200 percent BuzzFeed stock rise might signal start of a "pivot to AI" media trend.

BENJ EDWARDS - 1/27/2023, 7:10 PM

TECH • ARTIFICIAL INTELLIGENCE

The New AI-Powered Bing Is Threatening Users. That's No Laughing Matter

Lawsuits over Stability AI's Stable Diffusion could threaten the future of AI-generated art

<https://arstechnica.com/information-technology/2023/01/pivot-to-chatgpt-buzzfeed-preps-for-ai-written-content-while-cnet-fumbles/>

<https://www.businessinsider.com/stable-diffusion-lawsuit-getty-images-stability-ai-art-future-2023-1?op=1&r=US&IR=T>

AI at work: Coca-Cola



With ChatGPT and DALL·E, we're helping Coca-Cola to augment its world-class brands, marketing, and consumer experiences in industry-leading ways. We're also using OpenAI technology to improve business operations and capabilities.



"Man-made horrors beyond your comprehension, Zdzisław Beksiński, The Coca Cola Company, Horror, Matte painting, HQ, 4k", made with Stable Diffusion

Today's Program

0. Basics of Neural Networks
1. Optimization Algorithms
2. Generative Adversarial Networks
3. Recurrent Neural Networks
4. ResNets and Neural ODEs
5. Diffusion Models

What's a Neural Network again?

$$\Phi(x) = (\varphi_L \circ \dots \circ \varphi_1)(x)$$

↳ depth

Layers $\varphi_\ell(x_\ell) = \sigma_\ell \cdot \left(w_\ell x_\ell + b_\ell \right)$

activation function
(applied pointwise)

Weights ($\mathbb{R}^{n_{\ell+1} \times n_\ell}$) bias ($\mathbb{R}^{n_{\ell+1}}$)

Parameters θ

The diagram illustrates the structure of a neural network. At the top, the total function $\Phi(x)$ is shown as a composition of L layers: $\Phi(x) = (\varphi_L \circ \dots \circ \varphi_1)(x)$. Below this, the depth of the network is indicated by an arrow pointing to the right. The middle section shows a single layer $\varphi_\ell(x_\ell)$ being computed. It starts with input x_ℓ , which is multiplied by weight matrix w_ℓ and added to bias vector b_ℓ . The result is then passed through an activation function σ_ℓ . The bottom part of the diagram groups the weights and bias under the heading "Parameters θ ".

How to train your Neural Net

$$\min_{\theta} \underbrace{\mathcal{L}(x, y, \theta)}_{\text{Loss Function}} = \sum_{i \in I} L(y_i, \Phi_{\theta}(x_i))$$

$\hookrightarrow \text{dataset } \{(x_i, y_i)\}_{i \in I}$

trained with Gradient Descent

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta_k} \mathcal{L}$$



Learning Rate

Vignette 1: Optimization Algorithms

Pitfalls of training NNs

- For large datasets, evaluating the loss gets expensive
- The loss function is often non-convex
- No line search
- Training can take days, or months

Optimization Algorithms

Stochastic Gradient Descent

Replace the loss gradient by its expectation over the dataset

$$\theta_{k+1} = \theta_k - \eta \mathbb{E}_{i \in I} [\nabla_{\theta_k} L]$$

- **Stochastic** can get out of local minima
- Loss gets evaluated on random samples without replacement
- Still an ODE as $\eta \rightarrow 0$

Optimization Algorithms

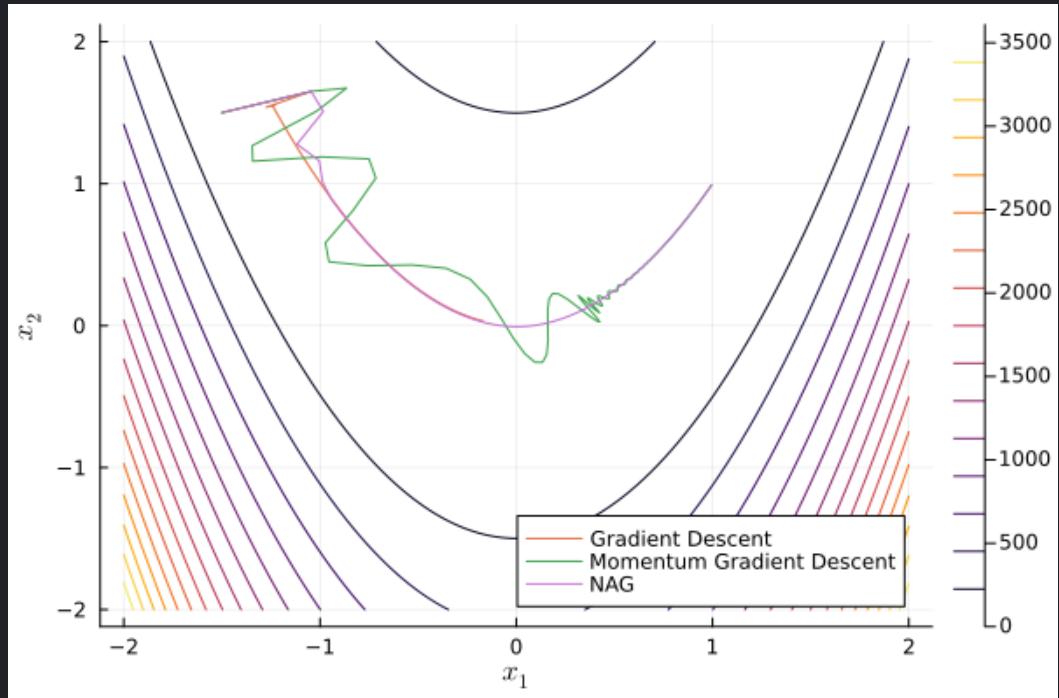
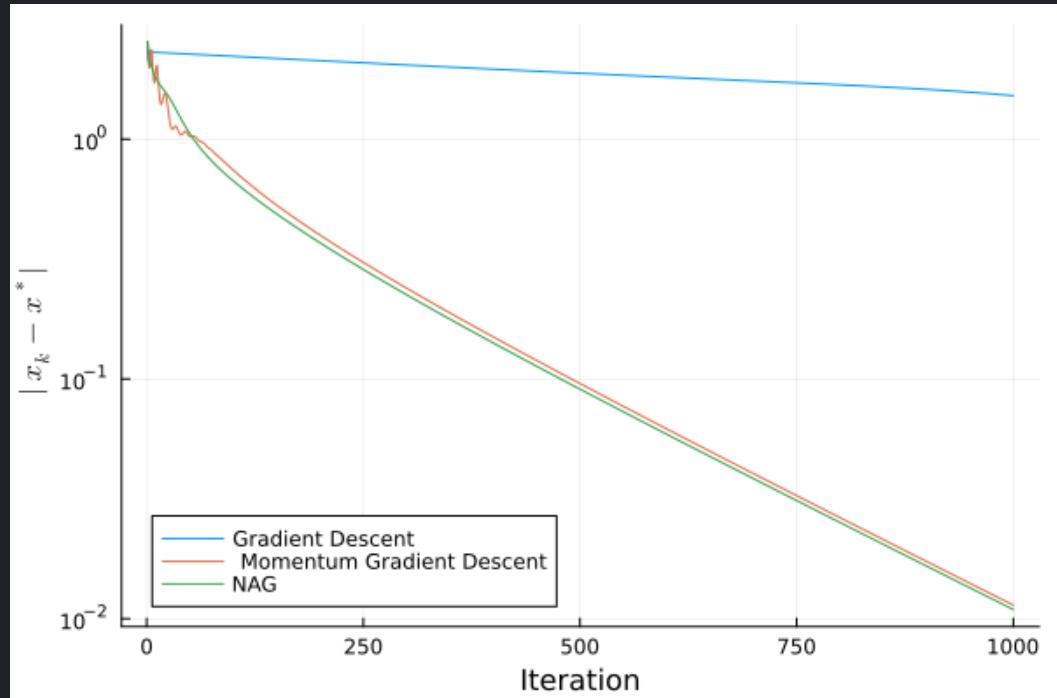
Momentum methods

Add scaled past gradients

$$\begin{aligned}v_{k+1} &= \gamma v_k + \eta \nabla_{\theta_k} \mathcal{L} \\ \theta_{k+1} &= \theta_k - v_{k+1}\end{aligned}$$

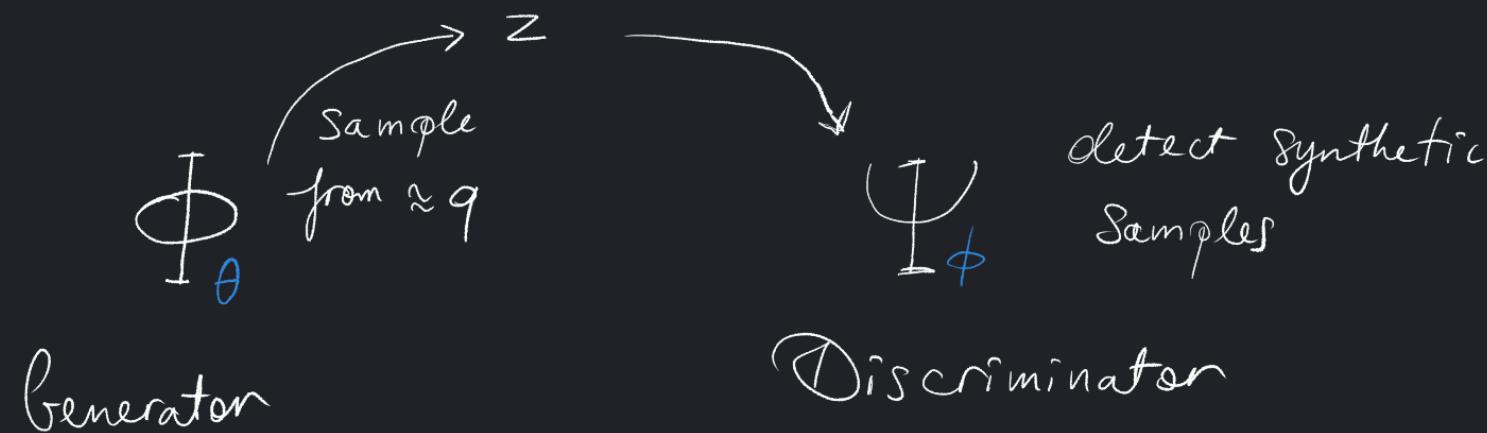
- Second order system with gradient as acceleration
- In the ODE limit, this is a Hamiltonian system [1]

[1] J. Diakonikolas, M.I. Jordan, *Generalized Momentum-Based Methods: A Hamiltonian Perspective*, SIAM J. Optim., 31 (1), (2021)



Vignette 2: Generative Adversarial Networks

Dataset $X \sim q$ data distribution



Simultaneous Gradient Descent

$$\begin{cases} \theta_{k+1} = \theta_k - \eta \nabla_{\theta_k} f(\theta_k, \phi_h) \\ \phi_{k+1} = \phi_k - \eta \nabla_{\phi_h} g(\theta_k, \phi_h) \end{cases}$$

Net a Gradient Flow !

Vignette 3: Recurrent Neural Networks

$$x_{n+1} = \phi(x_n)$$

$$= \phi^{on}(x_0)$$

Problem : gradient $\nabla_{\theta} \mathcal{L}(x_n)$ is prone to
converge to zero / diverge

$$\nabla_{\theta} \mathcal{L}(x_n) = \frac{\partial \mathcal{L}}{\partial x_n} \frac{\partial x_n}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial x_n} \left[\frac{\partial \phi}{\partial \theta} + \frac{\partial \phi}{\partial x_{n-1}} \cdot \frac{\partial x_{n-1}}{\partial \theta} \right]$$

⚠ Notation Abuse !

$$= \frac{\partial \mathcal{L}}{\partial x} \left[\frac{\partial \phi}{\partial \theta} \sum_{k=0}^{n-1} \left[\frac{\partial \phi}{\partial x} \right]^k \right] \rightarrow \text{Jacobian}$$

Vignette 4: Residual Networks

$$x_{\ell+1} = x_\ell + \textcolor{brown}{h}\varphi_\ell(x_\ell)$$

Motivation: force φ_ℓ to stay "close" to identity

Alternative: Looks like Euler's method! \rightarrow force φ_ℓ to be "regular"

Idea 1: Higher-order discretizations

Neural ODEs

Idea 2: $h \rightarrow 0$

$$\dot{x} = \Phi(x, t, \theta)$$

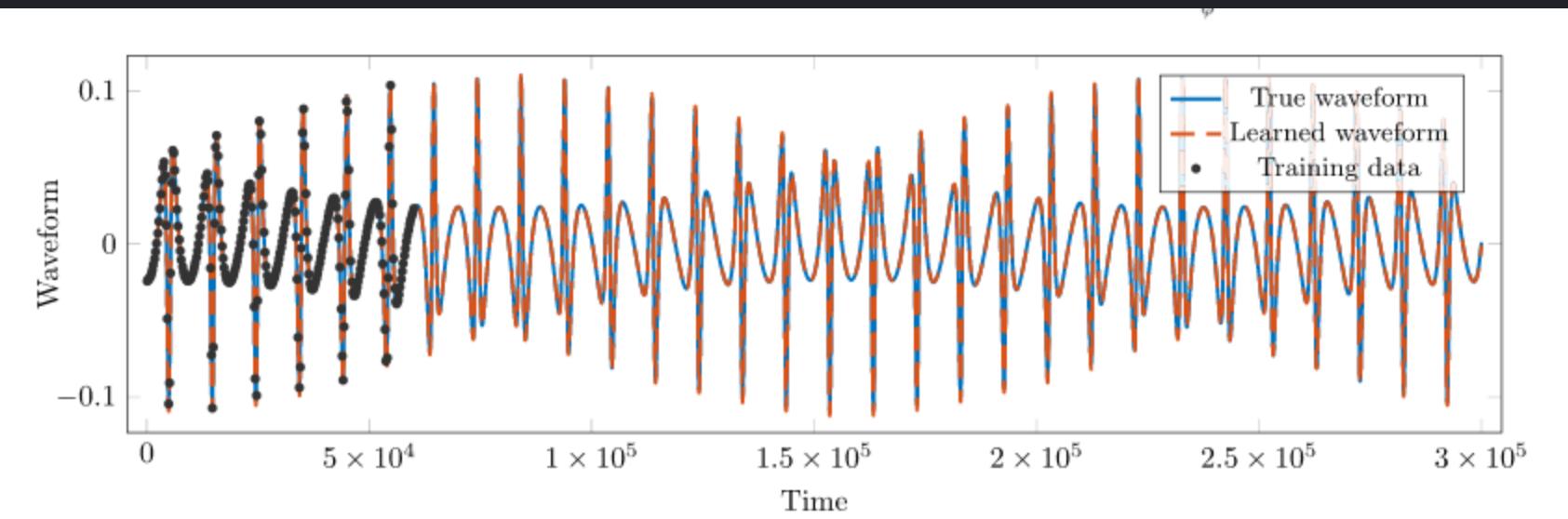
Advantages

- Can use black-box ODE solvers
- Adaptive time steps → Can handle irregularly sampled time series
- Controller design, Physical Modeling, ...
- Fit missing/unknown terms in existing scientific models (Scientific ML)

Open Access

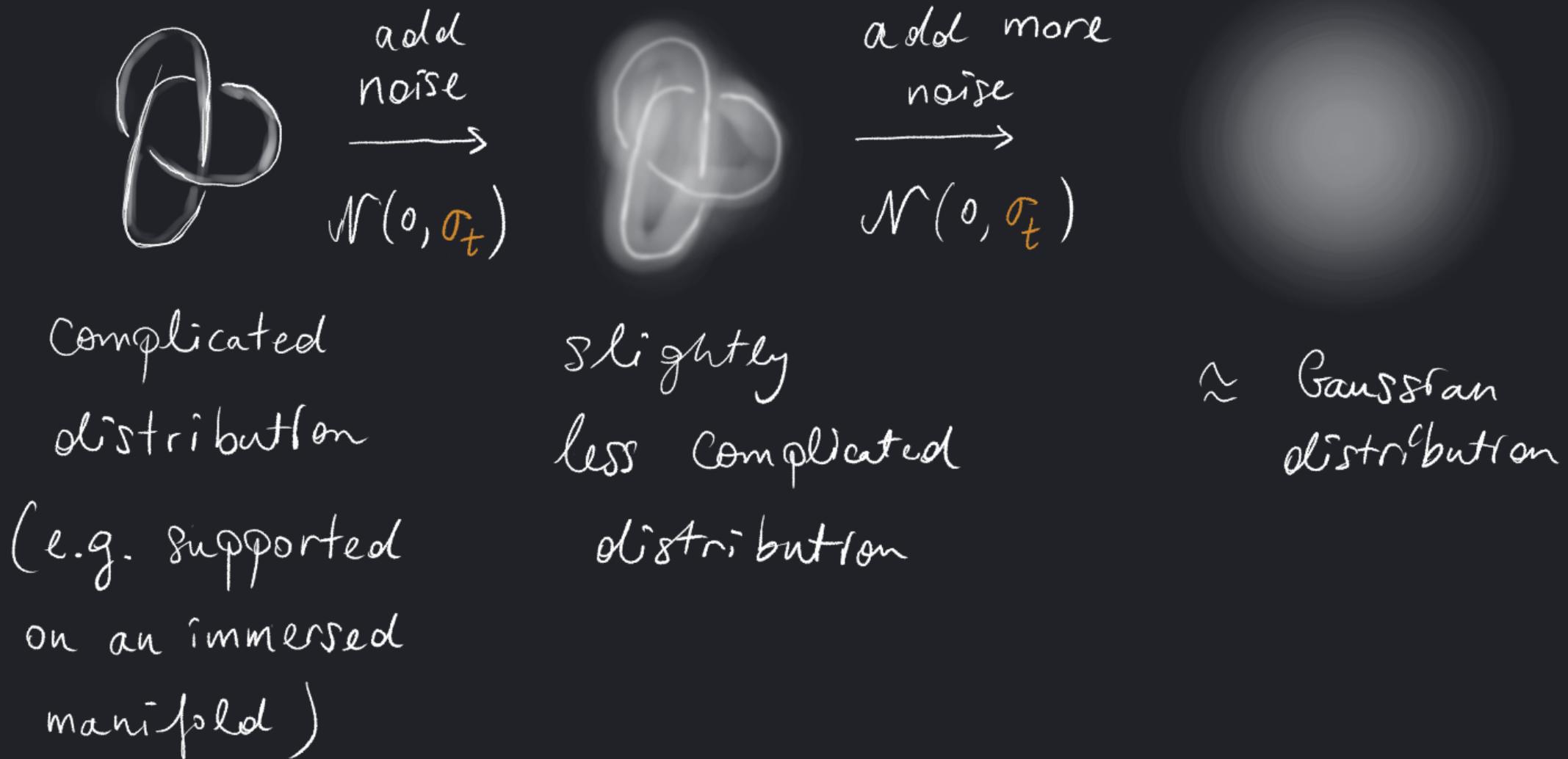
Learning orbital dynamics of binary black hole systems from gravitational wave measurements

Brendan Keith, Akshay Khadse, and Scott E. Field
Phys. Rev. Research **3**, 043101 – Published 9 November 2021



<https://phys.org/news/2021-11-machine-derive-black-hole-motion.html>

Vignette 5: Diffusion Models



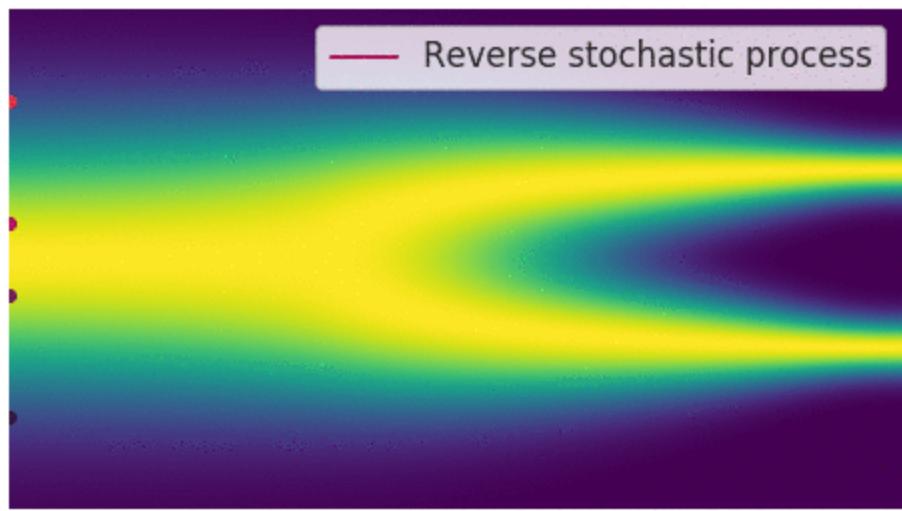
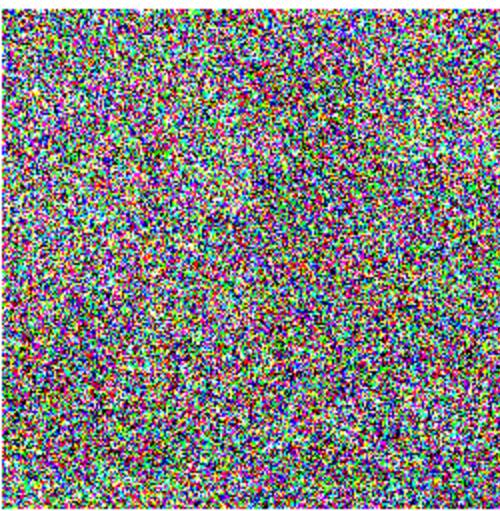
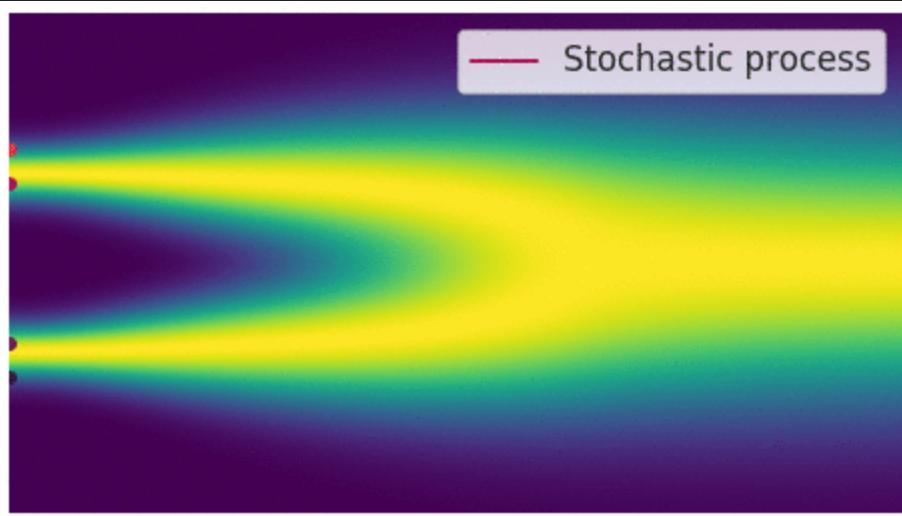
Forward process

Mathematically, the forward process can be expressed as an SDE, e.g.

$$\begin{aligned} dy(t) &= dW_t \\ y(0) &\sim q \end{aligned}$$

We can reverse this SDE to obtain another SDE, or an ODE

$$dy(t) = -\frac{1}{2} \nabla_y \log p_t(y) dt,$$



<https://yang-song.net/blog/2021/score/#score-based-generative-modeling-with-stochastic-differential-equations-sdes>

Score-based Diffusion Models

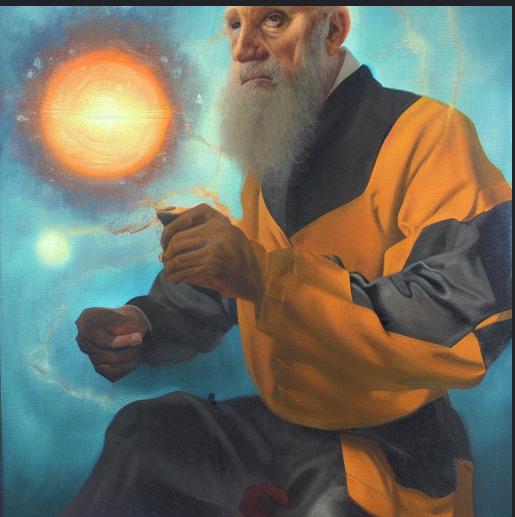
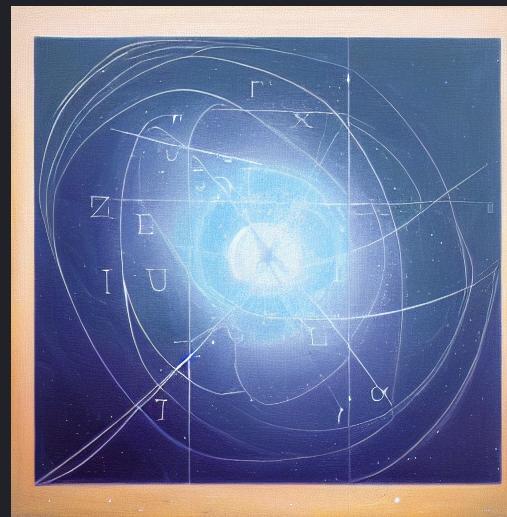
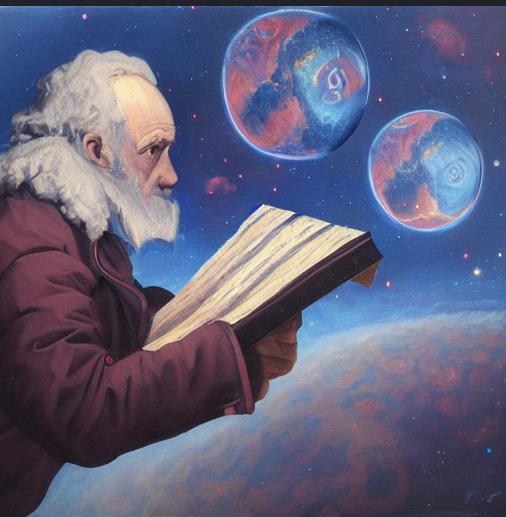
Approximate the score $\nabla_y \log p_t(y)$ by a neural network $s(y, t)$
(This is a Neural ODE!)

Controlled Generation

Use conditional score

$$\nabla_y \log p_t(y|x) = \nabla_y \log p_t(y) + \nabla_y \log p_t(x|y)$$

Closing words



'The Euler-Lagrange Equations, oil painting, trending on artstation, HQ", made with Stable Diffusion

Further Reading

- P. Kidger, *On Neural Differential Equations*, PhD Thesis (2022)
- Rackauckas et al. *Universal Differential Equations for Scientific Machine Learning*, (2020)
- P. Kidger - Score based diffusions explained in just one paragraph
- A. Ananthaswamy - The Physics Principle That inspired Modern AI Art - Quanta Magazine
- Cheatsheet: Recurrent Neural Networks