



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Enhancing Equity Data Availability using Synthetic Time Series Generation

Master Thesis

Azamat Zhaksylykov

September 4, 2025

Advisors: Dr. Florian Ofenheimer-Krach⁽³⁾, Damian Tschirky⁽¹⁾, Dr. Elizabeth Ren Rui Xing⁽²⁾, Prof. Dr. Josef Teichmann⁽³⁾,

IHH, ⁽¹⁾ Zürcher Kantonalbank

LBCH, ⁽²⁾ Zürcher Kantonalbank

Department of Mathematics, ⁽³⁾, ETH Zürich

Abstract

Recent advances in neural network techniques have opened new avenues in financial modeling, particularly in quantitative finance. A data-driven hedging models, such as Deep Hedging (Bühler et al., 2019), have gained popularity for their ability to learn and capture complex market dynamics directly from data. However, these models require large datasets for training, which are often scarce or inconsistent in historical market data. This thesis addresses the data scarcity issue by investigating the feasibility of adapting the Neural Jump Ordinary Differential Equations (NJ-ODEs) framework (F. Krach et al., 2022), originally designed for prediction, into a novel generative model. The proposed generative model was evaluated using two standard stochastic processes: the Geometric Brownian motion and the Ornstein-Uhlenbeck process.

Contents

Contents	iii
1 Introduction	1
2 Background*	3
2.1 Synthetic Time Series Generation in Finance	3
2.2 Neural Networks	4
2.2.1 Feedforward Neural Networks	4
2.2.2 Residual Neural Networks	5
2.2.3 Recurrent Neural Networks	5
2.2.4 Neural Ordinary Differential Equations	6
2.3 Mathematical setup	6
2.4 Neural Jump Ordinary Differential Equations	7
2.4.1 Introduction	7
2.4.2 The Model Framework	8
2.4.3 Objective Function	9
3 Methodology	11
3.1 Itô Diffusion Processes	11
3.2 NJ-ODE as Generative Method I	13
3.3 NJ-ODE as a Generative Method II	14
4 Experiments for Method I	19
4.1 Implementation Details	19
4.1.1 Datasets	19
4.1.2 Architecture	20
4.1.3 Training	20
4.2 Data Generation: Ornstein-Uhlenbeck 3D	21
4.3 Data Generation: Geometric Brownian Motion 1D	24
4.3.1 Separate Modeling of X and X^2	25

5 Experiments for Method II	27
5.1 Implementation Details	27
5.1.1 Datasets	27
5.1.2 Architecture and Training	28
5.2 Data Generation: Geometric Brownian Motion 1D	28
5.3 Data Generation: Ornstein-Uhlenbeck 1D	32
6 Conclusion	35
A Appendix	37
A.0.1 Stochastic Process	37
A.0.2 Conditional Expectation	39
A.0.3 Signature	40
A.0.4 Parameters Estimation for GBM	40
A.0.5 Parameters Estimation for OU	41
Bibliography	43

Chapter 1

Introduction

The last decade has seen remarkable progress in machine learning and artificial intelligence. In quantitative finance, these advancements have opened new possibilities for financial modeling and managing a derivatives portfolio. Neural networks, known for their ability to model complex and non-linear relationships, have become powerful tools for pricing and hedging financial instruments (Bühler et al., 2019). Their growing popularity is driven by their potential to learn market dynamics directly from historical data, providing an alternative to traditional models that depend on parametric assumptions.

Neural network-based financial models rely heavily on large, high-quality datasets to ensure accuracy and reliability (Ruf and Wang, 2020). However, obtaining extensive historical data in financial markets is challenging due to issues like non-stationarity, regime shifts, sparse data, irregular sampling intervals, and structural breaks. Moreover, financial data are often proprietary, costly, and restricted by privacy constraints, worsening the data scarcity problem. These barriers limit the practical use and effectiveness of neural network models in finance (Bühler et al., 2020, Potluru et al., 2023). Therefore, addressing data limitations while preserving statistical accuracy is an important research direction.

Addressing these limitations, this master's thesis attempts to adapt the NJ-ODEs framework into a generative model. NJ-ODEs, as defined by Herrera et al. (2021), are models designed for the forecast and filtering of continuous-time stochastic processes, providing theoretical optimality guarantees. This research specifically leverages the Path-Dependent Neural Jump Ordinary Differential Equations (PD-NJ-ODEs), an extension of NJ-ODEs capable of handling non-Markovian and discontinuous stochastic processes through the use of the signature transform. From this point onward, when referring to NJ-ODEs in this thesis, we mean the extended PD-NJ-ODEs framework rather than the original NJ-ODE model.

This thesis investigates the generative potential of the NJ-ODEs framework through empirical evaluations. The model's performance was tested using standard benchmark stochastic processes widely used in financial modeling, namely Geometric Brownian Motion (GBM) and the Ornstein-Uhlenbeck (OU) process. These processes are foundational in quantitative finance due to their analytical tractability and relevance to financial applications (Black and Scholes, 1973, Uhlenbeck and Ornstein, 1930, Vasicek, 1977). The experiments evaluated the NJ-ODE framework's ability to capture the key dynamics and statistical properties of these processes.

The rest of this thesis is structured as follows. Chapter 2 provides a theoretical background, covering synthetic time-series generation in finance, neural network architectures, mathematical foundations, and the NJ-ODE framework. Chapter 3 outlines the methodologies developed to adapt the NJ-ODE framework into generative modeling. Chapters 4 and 5 present the experimental setups, implementations, and results for two distinct generative approaches. Finally, Chapter 7 concludes with a summary of the key findings and challenges encountered throughout the research.

Chapter 2

Background*

This chapter opens with a concise introduction to synthetic time-series generation research in finance. It proceeds with a brief overview of feedforward, residual, and recurrent neural networks. It then introduces the key mathematical concepts relevant to this thesis, followed by a presentation of the NJ-ODE framework.

2.1 Synthetic Time Series Generation in Finance

The generation of synthetic financial time series has become increasingly important in modern quantitative finance. It facilitates various tasks such as model validation, stress testing, hedging derivatives portfolio, and the development of algorithmic trading strategies (Horvath et al., 2025). Traditionally, such data has been simulated using parametric models, such as the Black–Scholes model (Black and Scholes, 1973), the Heston model (Heston, 1993), ARIMA (Box et al., 2015), and GARCH(Bollerslev, 1986). These models depend on strict structural assumptions about the underlying stochastic processes and require calibrating a limited set of parameters to historical market data. Although these models provide analytical interpretability, they often struggle to reproduce the complex nonlinear dependencies and stylized features commonly observed in empirical financial data(Horvath et al., 2025).

Advances in machine learning have enabled the development of non-parametric, data-driven generative models that seek to learn the underlying distribution of financial time series directly from historical observations, without relying on rigid distributional assumptions(Bühler et al., 2020). They are called Market Generators. They span a variety of neural architectures, including

¹Some parts of this chapter are adapted from prior work, including Herrera et al. (2021), F. T. O. Krach (2025) and Andersson (2024).

Generative Adversarial Networks (GANs) (Goodfellow et al., 2014), Variational Autoencoders (VAEs) (Kingma and Welling, 2019), and transformer-based models (Vaswani et al., 2017). The representational power of these models has demonstrated encouraging results in capturing the complex dynamics of financial markets. This holds even in scenarios with limited data availability or non-stationarity (Horvath et al., 2025).

2.2 Neural Networks

Neural networks constitute a class of computational models inspired by the structure and functioning of biological neural systems. They have gained widespread usage across machine learning, artificial intelligence, and financial modeling due to their capacity to identify and model complex nonlinear relationships in datasets. Their strength lies in their ability to learn complex, nonlinear relationships from data through iterative training (Goodfellow et al., 2016).

2.2.1 Feedforward Neural Networks

Feedforward neural networks, also known as multi-layer perceptrons (MLPs), are foundational neural network architectures widely used in various applications, including financial modeling tasks. These networks consist of an input layer, multiple hidden layers, and an output layer, with data flowing unidirectionally from input to output without looping back (Goodfellow et al., 2016).

Mathematically, a feedforward neural network with n hidden layers can be expressed as:

$$h_{(i)} = \sigma \left(W_{(i)} h_{(i-1)} + b_{(i)} \right), \quad i = 1, \dots, n, \quad (2.1)$$

where $h_{(0)} = x$ is the input vector, $W_{(i)}$ are weight matrices, $b_{(i)}$ are bias vectors, and σ is a nonlinear activation function, such as sigmoid, tanh, or ReLU.

For example, a two-layer feedforward network is defined as:

$$y = W_{(2)} \sigma \left(W_{(1)} x + b_{(1)} \right) + b_{(2)}. \quad (2.2)$$

The training of such networks involves adjusting parameters ($W_{(i)}$ and $b_{(i)}$) to minimize a loss function. This optimization is typically done using gradient-based methods like stochastic gradient descent (SGD) or adaptive techniques such as Adam (Kingma and Ba, 2014). The backpropagation algorithm (Rumelhart et al., 1986) efficiently computes the gradients needed for parameter updates. This facilitates the widespread use of feedforward neural networks in practical applications.

2.2.2 Residual Neural Networks

Despite their ability to model complex patterns, deep feedforward networks frequently encounter challenges such as vanishing and exploding gradients, particularly in very deep architectures (Goodfellow et al., 2016). Residual neural networks (ResNets), introduced by He et al. (2016), address these challenges through skip or residual connections. These connections allow gradients to bypass specific layers directly, making it easier to train very deep networks.

A residual network updates layer representations according to:

$$h_{(i+1)} = h_{(i)} + \sigma \left(W_{(i)} h_{(i)} + b_{(i)} \right). \quad (2.3)$$

This formulation allows gradients to propagate more effectively through deeper networks, enabling the training of models with many layers without performance degradation (Goodfellow et al., 2016).

2.2.3 Recurrent Neural Networks

Recurrent neural networks (RNNs) are designed to handle sequential data, making them well-suited for time-series prediction tasks common in finance. RNNs maintain an internal state, or memory, that captures temporal dependencies in data (Goodfellow et al., 2016).

A standard recurrent network is defined as follows:

$$h_{i+1} = \sigma (Ux_{i+1} + Vh_i + b), \quad y_{i+1} = g(h_{i+1}), \quad (2.4)$$

where x_i is the input at time t_i , h_i is the hidden state, U and V are input-to-hidden and hidden-to-hidden weight matrices, respectively, and g maps hidden states to outputs (Goodfellow et al., 2016).

For this thesis, we adopt the following notation. The RNN employs a neural network ρ_θ , where θ represents the trainable parameters, to update a discrete latent variable h using discrete observations x_i on a regular grid, according to the update rule $h_{i+1} := \rho_\theta(h_i, x_{i+1})$.

Traditional RNNs often experience difficulties with capturing long-term dependencies due to gradient instability. Advanced variants, such as Long Short-Term Memory (LSTM) units (Hochreiter and Schmidhuber, 1997) and Gated Recurrent Units (GRU) (Cho et al., 2014), effectively mitigate these challenges through specialized gating mechanisms. These enhancements improve model performance on tasks involving lengthy temporal sequences.

2.2.4 Neural Ordinary Differential Equations

A key limitation of traditional RNNs is their inability to handle irregularly sampled data or model continuous-time dynamics. To address these challenges, Chen et al. (2018) introduced the Neural Ordinary Differential Equations (Neural ODEs) framework. They define the evolution of hidden states through a continuous-time differential equation:

$$\frac{dh_t}{dt} = f(h_t, t, \theta), \quad (2.5)$$

where h_t denotes the hidden state at time t , and $f(\cdot, \cdot, \theta) = f_\theta$ is a neural network parameterized by trainable weights θ .

The hidden state at any time $t \geq t_0$ can then be obtained by solving the initial value problem:

$$h_t = h_{t_0} + \int_{t_0}^t f(h_s, s, \theta) ds, \quad (2.6)$$

which enables continuous-time modeling of latent dynamics.

This solution is computed using numerical ODE solvers such as Euler or Runge–Kutta methods (Chen et al., 2018). The update rule is written as:

$$h_t := \text{ODESolve}(f, h_{t_0}, (t_0, t)), \quad (2.7)$$

where `ODESolve` represents a chosen numerical integration method applied over the time interval $[t_0, t]$.

2.3 Mathematical setup

We consider a continuous-time stochastic process $X = (X_t)_{t \in [0, T]}$ taking values in \mathbb{R}^d . The process X is assumed to be càdlàg, allowing for both continuous evolution and abrupt jumps.

We typically do not observe the full trajectory of X . Instead, we observe its values at a finite, random number of irregularly spaced time points. Let $n \in \mathbb{N}$ denote the number of observation times, and let $\{t_i\}_{i=1}^n \subset [0, T]$ denote the corresponding sequence of observation times.

To account for missing data, we introduce an observation mask $M = (M_i)_{i \in \mathbb{N}}$, where each $M \in \{0, 1\}^d$ is a random variable indicating which components of X_{t_i} are observed at time t_i . Specifically, if $M_{i,j} = 1$, then the j -th component $X_{t_i,j}$ is observed at time t_i . Otherwise, the corresponding component is considered missing. Furthermore, let $\tau(t) := \max\{t_i : 0 \leq k \leq n, t_i \leq t\}$ represent the most recent observation time (or zero if no observations have been made) before time $t \in [0, T]$.

We define the available information up to time t as the σ -algebra generated by all observations made up to and including time t . A natural filtration $(\mathcal{A}_t)_{t \in [0, T]}$ representing all available information given by

$$\mathcal{A}_t := \sigma (X_{t_i, j}, t_i, M_i | i \leq \kappa(t), j \in \{1 \leq l \leq d | M_{i, l} = 1\})$$

We denote the conditional expectation process of X by $\hat{X} := (\hat{X}_t)_{t \in [0, T]}$, where

$$\hat{X}_t := \mathbb{E} [X_t | \mathcal{A}_t].$$

As shown by Herrera et al. (2021), the NJ-ODE framework yields predictions that are L^2 optimal, as they approximate this conditional expectation.

All mathematical objects introduced in this section are formally defined in the appendix for clarity and reference.

2.4 Neural Jump Ordinary Differential Equations

In this section, I summarize the NJ-ODE framework developed by Florian Krach and collaborators. For a more detailed exposition is available at: <https://floriankrach.github.io/njode/>.

2.4.1 Introduction

The Neural Jump Ordinary Differential Equations (NJ-ODE) framework is a model for learning the latent dynamics of time-evolving systems from irregular and partially observed data. NJ-ODEs estimate the conditional expectation of future states by combining continuous latent evolution with discrete jumps at observation times (Herrera et al., 2021). The NJ-ODE architecture consists of three modular neural networks, each responsible for a distinct aspect of the system dynamics:

- **ODE Dynamics Network:** $f_{\theta_1} : \mathbb{R}^{d_H} \times \mathbb{R}^{d_X} \times [0, T] \times [0, T] \rightarrow \mathbb{R}^{d_H}$, governing the continuous evolution of the hidden state between observation times.
- **Jump Network:** $\rho_{\theta_2} : \mathbb{R}^{d_X} \rightarrow \mathbb{R}^{d_H}$, responsible for discrete updates of the hidden state at observation times, introducing a jump in the latent dynamics.
- **Readout Network:** $g_{\theta_3} : \mathbb{R}^{d_H} \rightarrow \mathbb{R}^{d_Y}$, mapping the hidden state to the observable output space.

Here, \mathbb{R}^{d_X} , \mathbb{R}^{d_H} , and \mathbb{R}^{d_Y} denote the input, latent, and output spaces, respectively, with $d_X, d_H, d_Y \in \mathbb{N}$.

2.4.2 The Model Framework

The NJ-ODE model is defined as follows:

$$\begin{aligned} H_0 &= \rho_{\theta_2}(X_0, 0), \\ dH_t &= f_{\theta_1}(H_{t-}, X_{\tau(t)}, \tau(t), t - \tau(t)) dt + (\rho_{\theta_2}(X_t, H_{t-}) - H_{t-}) du_t, \\ Y_t &= g_{\theta_3}(H_t), \end{aligned}$$

where $H_t \in \mathbb{R}^{d_H}$ denotes the hidden (latent) state process and $Y_t \in \mathbb{R}^{d_Y}$ represents the model's output. The function $u_t := \sum_{i=1}^n 1_{[t_i, \infty)}(t)$ is a counting process that increases at discrete observation times t_i , while $\tau(t)$ denotes the most recent observation before time t . The set of parameters $\theta := (\theta_1, \theta_2, \theta_3)$ encompasses all learnable weights across the three neural networks and optimized during training.

This formulation allows the model to evolve continuously between observations via the ODE network f_{θ_1} , while the jump network ρ_{θ_2} deterministically updates the hidden state at each observation time, enabling the NJ-ODE to capture complex dynamics.

Another way to write the NJ-ODE:

$$\begin{aligned} h_{t_{i+1}}^- &:= \text{ODESolve}(f_{\theta_1}, (h_t, x_{t_i}, t_i, t - t_i), (t_i, t_{i+1})), \\ h_{t_{i+1}} &:= \rho_{\theta_2}(x_{t_{i+1}}). \end{aligned}$$

Herrera et al. (2021) demonstrated that, when trained using suitable objectives, the NJ-ODE effectively approximates the conditional expectation of the underlying process, yielding optimal L^2 predictions. This ensures that the model not only fits observed data points but also generalizes for forecasting tasks.

Building on this framework, F. Krach et al. (2022) introduced the path-dependent NJ-ODE, extending the original model to capture temporal dependencies beyond the current input. This extension is motivated by applications in which the dynamics of a system depend not only on the most recent input but also on the entire history of the input trajectory. To model such dependencies, the authors incorporate the signature transform from rough path theory. This technique provides a compact, non-parametric representation of the history of a path, allowing the model to capture complex temporal dependencies (Chevyrev and Kormilitzin, 2016).

The resulting path-dependent NJ-ODE is formulated as:

$$\begin{aligned} H_0 &= \rho_{\theta_2}(0, 0, \pi_m(0), M_0 \odot X_0), \\ dH_t &= f_{\theta_1}(H_{t-}, t, \tau(t), \pi_m(\tilde{X}_{\leq \tau(t)})) dt + [\rho_{\theta_2}(H_{t-}, t, \pi_m(\tilde{X}_{\leq \tau(t)}), M_t \odot X_t) - H_{t-}] du_t, \end{aligned}$$

$$Y_t = g_{\theta_3}(H_t),$$

where $\pi_m(\tilde{X}_{\leq \tau(t)})$ denotes the truncated signature of the interpolated input path $\tilde{X}_{\leq \tau(t)}$ up to order m . The path $\tilde{X}_{\leq \tau(t)}$ is constructed via forward-fill interpolation, which preserves causality by ensuring that no future information leaks into the model at any time t . The element-wise product $M_t \odot X_t$ enables masking of certain input features, allowing the model to incorporate irregularly sampled or missing data.

2.4.3 Objective Function

The training objective of the NJ-ODE framework is to approximate the conditional expectation of the target process X , which corresponds to the optimal predictor in the L^2 -norm sense. The objective function Ψ is defined as follows:

$$\Psi(Y) := \mathbb{E}_{\mathbb{P} \times \tilde{\mathbb{P}}} \left[\frac{1}{n} \sum_{i=1}^n \left(\|M_i \odot (X_{t_i} - Y_{t_i})\|_2 + \|M_i \odot (Y_{t_i} - Y_{t_i^-})\|_2 \right)^2 \right], \quad (2.8)$$

where Y is A-adapted process, M_i is a masking operator that accounts for missing observations, and the expectation is taken over the product probability space $(\Omega \times \tilde{\Omega}, \mathcal{F} \otimes \tilde{\mathcal{F}}, \mathbb{P} \times \tilde{\mathbb{P}})$.

The training loss function \mathcal{L} is defined by evaluating this objective on the model's output for a given parameter configuration θ :

$$\mathcal{L}(\theta) := \Psi(Y^\theta(X)),$$

where $Y^\theta(X)$ represents the output of the NJ-ODE model with parameters θ , conditioned on the observed data from the process X .

This loss function incorporates two key components:

- The *jump term* $\|X_{t_i} - Y_{t_i}\|_2$, which encourages the jump network ρ_{θ_2} to produce accurate state updates upon receiving new observations.
- The *continuity term* $\|Y_{t_i} - Y_{t_i^-}\|_2$, which penalizes discontinuities in the model output. This term drives the ODE dynamics network f_{θ_1} to evolve the latent state H_t smoothly over time to follow the conditional expectation path closely.

Both terms also contribute to training the readout network g_{θ_3} , which maps the latent state to the output space. The model is thus trained to maintain accurate predictions both at and between observation times.

This formulation corresponds to the loss function referred to as "standard" in the original NJ-ODE implementation.

Chapter 3

Methodology

This chapter presents the methodology to adapt the NJ-ODE model for generative purposes. We begin by introducing the Itô diffusion process, which serves as the model for the dynamics of the underlying stochastic process. We then describe the modifications made to the original NJ-ODE framework to enable synthetic data generation. In particular, we detail how the model was adapted to estimate both the drift and diffusion components of stochastic processes, a key step in extending NJ-ODE from a forecasting tool to a generative framework.

3.1 Itô Diffusion Processes

Itô diffusion processes form a fundamental class of continuous-time stochastic processes that are widely used to model systems evolving under both deterministic trends and random perturbations. These processes are described by stochastic differential equations (SDEs) of the general form:

$$dX_t = \mu(t, X_t) dt + \sigma(t, X_t) dW_t, \quad (3.1)$$

Let $\{W_t\}_{t \in [0, T]}$ be a d_W -dimensional Brownian motion defined on a filtered probability space $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$, and let $X := (X_t)_{t \in [0, T]}$ be the solution process, where $X_t \in \mathbb{R}^{d_X}$ denotes the state at time t . The measurable functions $\mu : [0, T] \times \mathbb{R}^{d_X} \rightarrow \mathbb{R}^{d_X}$ and $\sigma : [0, T] \times \mathbb{R}^{d_X} \rightarrow \mathbb{R}^{d_X \times d_W}$ are referred to as the drift and diffusion coefficients, respectively (Bass, 2011; Björk, 2009; Shreve, 2004). Further details are provided in the Appendix.

The drift term $\mu(t, X_t)$ encodes the deterministic direction of movement, while the diffusion term $\sigma(t, X_t) dW_t$ captures the randomness introduced by the Brownian motion.

The dependence on the entire path history $X_{[0,t]}$ allows for non-Markovian dynamics in processes named as Itô processes. These enhance the modeling flexibility by incorporating memory effects. Itô processes are central to many applications in science and engineering, particularly in quantitative finance, where they are used to model asset prices, interest rates, and volatility surfaces (Björk, 2009). The ability to accommodate both temporal and historical path dependence makes Itô processes a robust framework for modeling stochastic systems.

Geometric Brownian Motion

GBM stands as a fundamental stochastic process in quantitative finance, particularly for modeling the evolution of asset prices. Its popularity arises from its analytical tractability. Additionally, it plays a central role in classical financial models, most notably the Black–Scholes option pricing framework (Black and Scholes, 1973).

The SDE describes the dynamics of a GBM

$$dX_t = \mu X_t dt + \sigma X_t dW_t,$$

where $X_t \in \mathbb{R}_{>0}$ represents the asset price at time t , μ is the drift parameter, $\sigma > 0$ is the volatility parameter, and W_t is a standard Brownian motion.

The solution to this SDE has a closed-form expression:

$$X_t = X_0 \exp \left(\left(\mu - \frac{1}{2} \sigma^2 \right) t + \sigma W_t \right),$$

which shows that the logarithm of X_t is normally distributed, implying that X_t itself follows a log-normal distribution. This property makes GBM particularly appealing for modeling financial quantities that must remain strictly positive, such as stock prices.

Despite its simplicity and assumptions like constant drift and volatility, GBM is a cornerstone of financial mathematics. It acts as a benchmark for more complex models, such as those incorporating stochastic volatility and jump-diffusion processes (Björk, 2009; Shreve, 2004).

Ornstein–Uhlenbeck Process

The OU process is a classical example of a mean-reverting stochastic process. In finance, the OU process has been used to model interest rates, stochastic volatility, and other mean-reverting quantities (Björk, 2009; Uhlenbeck and Ornstein, 1930; Vasicek, 1977).

The SDE governs the dynamics of the OU process:

$$dX_t = k(m - X_t) dt + \sigma dW_t,$$

where $X_t \in \mathbb{R}$ denotes the state of the process at time t , $k > 0$ is the rate of mean reversion, m is the long-term mean level, $\sigma > 0$ is the volatility parameter, and W_t is a standard Brownian motion.

The exact solution of the OU process is given by:

$$X_t = X_0 e^{-kt} + m(1 - e^{-kt}) + \sigma \int_0^t e^{-k(t-s)} dW_s,$$

The OU process plays a central role in various financial models. For instance, it forms the basis of the Vasicek interest rate model (Vasicek, 1977).

3.2 NJ-ODE as Generative Method I

In this section, we present a method for extracting instantaneous drift and diffusion estimates from the predictions of a neural network trained to approximate the conditional first and second moments of a stochastic process. Our approach assumes the underlying dynamics follow an Itô process of the general form:

$$dX_t = \mu(t, X_t) dt + \sigma(t, X_t) dW_t, \quad (3.2)$$

where $X_t \in \mathbb{R}_X^d$ is the state of the process at time t , μ is the drift vector, σ is the diffusion matrix, and W_t is a standard d_W -dimensional Brownian motion. The functions μ and σ may depend on the full path history X_t , allowing for non-Markovian dynamics.

To approximate this process numerically, we discretize the dynamics using the Euler–Maruyama scheme (Kloeden and Platen, 1992) :

$$X_{t+\Delta t} \approx X_t + \mu_t \Delta t + \sigma_t \Delta W_t, \quad (3.3)$$

where $\mu_t = \mu(t, X_t)$, $\sigma_t = \sigma(t, X_t)$, and $\Delta W_t \sim \mathcal{N}(0, \Delta t \cdot I_{d_W})$ represents a Brownian increment. Our goal is to estimate μ_t and σ_t based on the conditional expectations of $X_{t+\Delta t}$ and $X_{t+\Delta t} X_{t+\Delta t}^\top$, learned by our NJ-ODE model.

Drift Estimation from First Moments

The NJ-ODE model is trained to predict the conditional expectation $\mathbb{E}[X_{t+\Delta t} | \mathcal{A}_t]$, where \mathcal{A}_t represents the available information up to time t , and noting that $\mathbb{E}[\Delta W_t | \mathcal{A}_t] = 0$, we have:

$$\mathbb{E}[X_{t+\Delta t} | \mathcal{A}_t] \approx X_t + \mathbb{E}[\mu_t | \mathcal{A}_t] \Delta t. \quad (3.4)$$

Rearranging terms yields an estimate for the conditional drift:

$$\hat{\mu}_t := \mathbb{E}[\mu_t | \mathcal{A}_t] \approx \frac{\mathbb{E}[X_{t+\Delta t} | \mathcal{A}_t] - X_t}{\Delta t}. \quad (3.5)$$

Diffusion Estimation from Second Moments

To estimate the instantaneous diffusion, we consider the conditional expectation of second-order products. Applying Itô's formula to the product $X_t^i X_t^j$ yields:

$$d(X^i X^j)_t = X_t^j dX_t^i + X_t^i dX_t^j + d[X_t^i, X_t^j].$$

Substituting the SDE terms gives:

$$d(X^i X^j)_t = X_t^j \mu_t^i dt + X_t^j \sigma_t^i dW_t + X_t^i \mu_t^j dt + X_t^i \sigma_t^j dW_t + \sigma_t^i \sigma_t^j dt.$$

Discretising using the Euler–Maruyama scheme and taking conditional expectations leads to:

$$\mathbb{E}[(X^i X^j)_{t+\Delta t} \mid \mathcal{A}_t] \approx (X^i X^j)_t + X_t^j \hat{\mu}_t^i \Delta t + X_t^i \hat{\mu}_t^j \Delta t + \mathbb{E}[\sigma_t^i \sigma_t^j \mid \mathcal{A}_t] \Delta t.$$

Rearranging terms yields the estimate of the covariance matrix $\Sigma_t = \sigma_t \sigma_t^\top$ as:

$$(\hat{\Sigma}_t)_{ij} := \mathbb{E}[\sigma_t^i \sigma_t^j \mid \mathcal{A}_t] \approx \frac{\mathbb{E}[(X^i X^j)_{t+\Delta t} \mid \mathcal{A}_t] - (X^i X^j)_t - X_t^i \hat{\mu}_t^j - X_t^j \hat{\mu}_t^i}{\Delta t}. \quad (3.6)$$

To generate a new point using the Euler–Maruyama scheme, we require a matrix $\hat{\sigma}_t$ such that $\hat{\sigma}_t \hat{\sigma}_t^\top = \hat{\Sigma}_t$. If $\hat{\Sigma}_t$ is positive semi-definite, this can be computed via the Cholesky decomposition. Theoretically, this condition is satisfied since, for any non-zero vector $x \in \mathbb{R}^{d_x}$,

$$x^\top \hat{\Sigma}_t x = \mathbb{E}[x^\top \Sigma_t x \mid \mathcal{A}_t] \geq 0.$$

However, due to estimation errors and numerical approximations, the computed matrix $\hat{\Sigma}_t$ may not be positive semi-definite in practice.

3.3 NJ-ODE as a Generative Method II

In this section, we present a second generative modeling approach using the NJ-ODE framework, which enables the direct estimation of the instantaneous covariance structure of a stochastic process. Unlike Method I, which derives diffusion estimates from second-order moment predictions, Method II only uses first-order moment predictions.

Direct Estimation of Diffusion

It can be shown that Equation (3.6) can be equivalently written as

$$\hat{\Sigma}_t = \frac{1}{\Delta t} \mathbb{E} \left[(X_{t+\Delta t} - X_t)(X_{t+\Delta t} - X_t)^\top \mid \mathcal{A}_t \right],$$

which expresses the instantaneous covariance in terms of the conditional second moment of the process increments.

Based on this representation, we define a process Z_t as

$$Z_t := (X_{t+\Delta t} - X_t)(X_{t+\Delta t} - X_t)^\top$$

The above definition is valid under the assumption that the process X is fully observed. However, in practical scenarios where observations of X are incomplete or irregularly spaced, we generalize the definition of Z_t by replacing X_t with the most recent observation available before time t , denoted $\tau(t)$. The generalized process is given by

$$Z_t = (X_t - X_{\tau(t)})(X_t - X_{\tau(t)})^\top.$$

For each observation time t_i , we define:

- At the observation point:

$$Z_{t_i} = (X_{t_i} - X_{\tau(t_i)})(X_{t_i} - X_{\tau(t_i)})^\top = \mathbf{0},$$

since $\tau(t_i) = t_i$.

- Just before the observation:

$$Z_{t_i^-} = (X_{t_i^-} - X_{\tau(t_i^-)})(X_{t_i^-} - X_{\tau(t_i^-)})^\top = (X_{t_i} - X_{t_{i-1}})(X_{t_i} - X_{t_{i-1}})^\top,$$

where $\tau(t_i^-) = t_{i-1}$.

Therefore, when the dataset for the Z process was created, simultaneously two datasets were created, one for Z_t and the other for Z_t^- .

The goal is to train the NJ-ODE model to estimate the conditional expectation of the Z process, using a neural network architecture that ensures the positive semi-definiteness of the covariance matrix. It is important to note that if $X \in \mathbb{R}^{d_X}$, then the corresponding process $Z \in \mathbb{R}^{d_X \times d_X}$ is a matrix-valued quantity. To enable processing within the NJ-ODE framework, each Z matrix is first flattened into a vector of length d_X^2 , which serves as the model input. The NJ-ODE model then outputs a vector of the same dimension, which is subsequently reshaped back into a $d_X \times d_X$ matrix denoted as Y_t , representing the model's output at time t . The estimated covariance is defined as:

$$\hat{\Sigma}_t := \frac{Y_t Y_t^\top}{\Delta t}.$$

This formulation guarantees that the estimated covariance matrix is symmetric and positive semi-definite by construction. The volatility may depend on the level of the process (e.g., for the Black-Scholes volatility term is given

by $\sigma_t X_t$). Therefore, training on the Z -process alone would be insufficient and better option is to train using the dataset $[Z, Z^-, X]$. To implement this, we rely on an extension of the NJ-ODE framework introduced by Heiss et al. (2024). This new framework enables variable-sized inputs and outputs, allowing for greater flexibility. In particular, this extension enables the explicit selection of input variables used during training.

New Loss Function

To guide the NJ-ODE model during training, we begin by recalling the original loss function, which is designed to approximate the conditional expectation of the target process X in the sense of the L_2 -norm. The loss functional is given by:

$$\Psi(X, Y) := \mathbb{E}_{\mathbb{P} \times \tilde{\mathbb{P}}} \left[\frac{1}{n} \sum_{i=1}^n \left(\|M_i \odot (X_{t_i} - Y_{t_i})\|_2 + \|M_i \odot (Y_{t_i} - Y_{t_i^-})\|_2 \right)^2 \right],$$

where $M_i \in \{0, 1\}^d$ is a binary mask that encodes which components of the process are observed at time t_i . The first term penalizes the model for inaccurate predictions at observed time points, while the second term regularizes the jump size between consecutive model outputs.

In the generative setting, our goal extends beyond estimating the conditional expectation of X . We also seek to estimate the instantaneous volatility. Since the volatility structure is captured by the Z_t process, we adapt the loss function accordingly. The model takes as input a flattened representation of the Z data and outputs a flattened version of Y . The output Y is reshaped into matrix form, its product $Y_t Y_t^\top$ is computed, and this result is flattened again. The modified loss function, designed to recover the diffusion component, is defined as follows:

$$\Psi(Z, Y) := \mathbb{E}_{\mathbb{P} \times \tilde{\mathbb{P}}} \left[\frac{1}{n} \sum_{i=1}^n \left(\|M_i \odot (Y_{t_i^-} Y_{t_i^-}^\top - Z_{t_i^-})\|_2 + \|M_i \odot (Y_{t_i} Y_{t_i}^\top)\|_2 \right)^2 \right],$$

This formulation allows the NJ-ODE framework to learn second-order structure from the data. In the NJ-ODE codebase, this loss function is referred to as "vola".

Drift and Volatility Estimation

The drift and volatility components are learned using two separate NJ-ODE models. The first model is trained on the original process X to predict the

conditional mean $\mathbb{E}[X_{t+\Delta t} \mid \mathcal{A}_t]$, from which the drift is estimated as:

$$\hat{\mu}_t \approx \frac{\mathbb{E}[X_{t+\Delta t} \mid \mathcal{A}_t] - X_t}{\Delta t}.$$

The second model targets the diffusion term, learning to predict the conditional expectation of Z_t . The model is trained by minimizing the "vola" loss function, such that the $Y_t Y_t^\top$ approximates $[Z_{t+\Delta t} \mid \mathcal{A}_t]$.

Given the model output Y_t , the instantaneous covariance matrix is estimated as:

$$\hat{\Sigma}_t = \frac{Y_t Y_t^\top}{\Delta t}.$$

The corresponding volatility $\hat{\sigma}_t$ is then obtained as the Cholesky factor of $\hat{\Sigma}_t$, which is equal to $\frac{Y_t}{\sqrt{\Delta t}}$.

Generating Synthetic Data

Once the models for $\hat{\mu}_t$ and $\hat{\Sigma}_t$ are trained, the generative process proceeds via the Euler–Maruyama scheme:

$$X_{t+\Delta t} = X_t + \hat{\mu}_t \Delta t + \hat{\sigma}_t \Delta W_t,$$

where $\Delta W_t \sim \mathcal{N}(0, \Delta t \cdot I_{d_W})$, and $\hat{\sigma}_t = \frac{Y_t}{\sqrt{\Delta t}}$. This formulation overcomes the issues related to negative eigenvalues in the covariance matrix.

By iteratively applying this scheme, the NJ-ODE model generates synthetic sample paths. This approach provides a framework for generative modeling of continuous-time stochastic processes.

Chapter 4

Experiments for Method I

All implementations were done using PyTorch. The code is available at: <https://github.com/zhaksylykov/thesis>.

4.1 Implementation Details

4.1.1 Datasets

Multiple datasets were generated by simulating $N = 10,000$ sample paths of a 3-dimensional OU process and a 1-dimensional GBM. The Euler-Maruyama discretisation scheme was used for the simulations. Each simulation was carried out on an equidistant time grid with step size $\Delta t = 0.01$ over the interval $[0, T]$, with $T = 1$. All sample paths are fully observed at every time point, meaning the trajectories contain no missing data.

Each dataset was subsequently split into training and testing subsets using an 80%/20% partition. The specific parameters used for a 3-dimensional OU process and a 1-dimensional GBM are detailed below.

Ornstein-Uhlenbeck 3D

- SDE:

$$dX_t = -\Theta(X_t - \mu) dt + \Sigma dW_t,$$

where $X_t \in \mathbb{R}^3$, $\Theta \in \mathbb{R}^{3 \times 3}$ is the mean-reversion matrix, $\mu \in \mathbb{R}^3$ is the long-term mean vector, $\Sigma \in \mathbb{R}^{3 \times 3}$ is the volatility matrix, and $W_t \in \mathbb{R}^3$ is a 3-dimensional Brownian motion.

- Conditional expectation:

$$\mathbb{E}[X_{t+s} \mid \mathcal{A}_t] = e^{-\Theta s} X_t + (I - e^{-\Theta s}) \mu.$$

- **Parameters used:**

$$\mu = \begin{bmatrix} 1.2 \\ 1.0 \\ 1.5 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 0.2 & 0.1 & 0.1 \\ 0.1 & 0.25 & 0.1 \\ 0.1 & 0.1 & 0.3 \end{bmatrix}, \quad \Theta = \begin{bmatrix} 0.3 & 0 & 0 \\ 0 & 0.3 & 0 \\ 0 & 0 & 0.3 \end{bmatrix}, \quad X_0 = \begin{bmatrix} 1.0 \\ 1.5 \\ 2.0 \end{bmatrix},$$

Geometric Brownian Motion 1D

- **SDE:**

$$dX_t = \mu X_t dt + \sigma X_t dW_t,$$

where W_t is a 1-dimensional Brownian motion.

- **Conditional expectation:**

$$\mathbb{E}[X_{t+s} \mid \mathcal{A}_t] = X_t e^{\mu s},$$

- **Conditional variance:**

$$\text{Var}[X_{t+s} \mid \mathcal{A}_t] = X_t^2 e^{2\mu s} (e^{\sigma^2 s} - 1),$$

- **Parameters used:** $\mu = 2, \sigma = 0.3, X_0 = 1$.

4.1.2 Architecture

In all experiments, the dimension of the latent (hidden) state is set to $d_H = 10$. The functions f_{θ_1} , \tilde{g}_{θ_3} , and $\tilde{\rho}_{\theta_2}$ are implemented as feedforward neural networks, each comprising two hidden layers with 50 units per layer and using the tanh activation function. The readout network g_{θ_3} and the jump network ρ_{θ_2} , are then constructed as residual versions of \tilde{g}_{θ_3} and $\tilde{\rho}_{\theta_2}$.

To prevent overfitting and improve generalization, dropout with a rate of 0.1 is applied after each nonlinearity within the networks. The continuous-time dynamics of the latent state are solved using the explicit Euler method.

4.1.3 Training

The neural networks were trained using the Adam optimizer with a learning rate of 0.001. Training was conducted over 100 epochs with a batch size of 200. All network parameters were initialized randomly. No additional hyperparameter tuning was performed. We employed a modified version of the original NJ-ODE loss function, referred to as the "easy" loss in the implementation. It is defined as:

$$\Psi(Y) := \mathbb{E}_{\mathbb{P} \times \tilde{\mathbb{P}}} \left[\frac{1}{n} \sum_{i=1}^n \left(\|M_i \odot (X_{t_i} - Y_{t_i})\|_2 + \|M_i \odot (X_{t_i} - Y_{t_i^-})\|_2 \right)^2 \right], \quad (4.1)$$

This loss differs from the standard formulation by replacing the term $\|Y_{t_i} - Y_{t_i^-}\|_2$ with $\|X_{t_i} - Y_{t_i^-}\|_2$. The adjustment mitigates the risk of local minima and enhances training outcomes. This was empirically demonstrated by (F. T. O. Krach (2025)).

By directly comparing both the model's jump prediction Y_{t_i} and the pre-jump state $Y_{t_i^-}$ against the true observation X_{t_i} , this formulation penalizes discrepancies more effectively. In particular, it discourages the model from producing intermediate values $Y_{t_i} \neq X_{t_i}$ that reduce the loss artificially. As a result, the "easy" loss improves stability during training and promotes better convergence to the true conditional expectation (F. T. O. Krach, 2025).

4.2 Data Generation: Ornstein-Uhlenbeck 3D

We employ Method I to generate data based on a 3-dimensional OU process $X = (X_1, X_2, X_3)$, where $X \in \mathbb{R}$ with $d = 3$. In the initial setup, the NJ-ODE model was trained directly on the raw process values of X , such that the input and output dimensions matched the process dimensionality. While this configuration is sufficient for modeling the first-order dynamics of X_t , it is inadequate for capturing the second-order dynamics of X_t . Modeling dynamics of $(XX^\top)_t$ is central to our generative modeling approach.

We extended the input representation to include pairwise products and squared terms of the components of X , thereby capturing both linear and second-order interactions. Specifically, we define the augmented input vector as

$$\tilde{X} = (X_1, X_2, X_3, X_1^2, X_2^2, X_3^2, X_1X_2, X_1X_3, X_2X_3),$$

which consists of all unique monomials of degree at most two. This transformation results in an input and output dimensionality of $\frac{d^2+3d}{2}$; for $d = 3$, this yields 9 features.

This extended feature representation allows the model to learn the evolution of second-moment structure over time. The training performance on test data using the augmented input representation is illustrated in the figure below. The close alignment between predictions and ground truth suggests successful learning of the process dynamics.

4. EXPERIMENTS FOR METHOD I

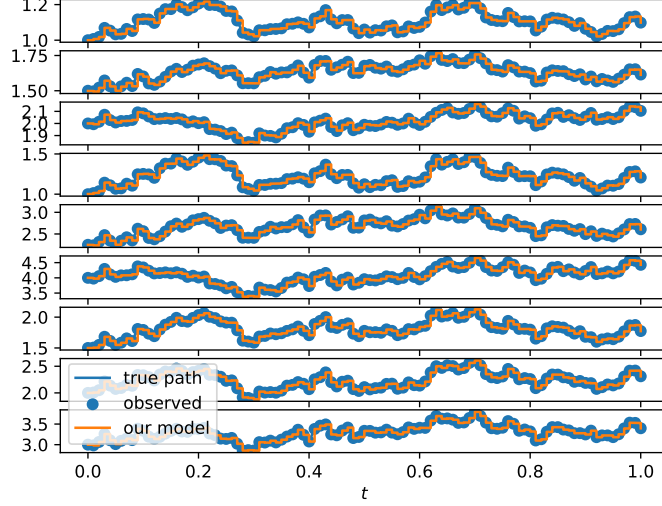


Figure 4.1: Training results for the NJ-ODE model with augmented input features.

To iteratively generate new sample paths, we begin from a given initial state X_0 and evolve the process forward in time using the conditional drift and covariance estimated by the trained NJ-ODE model. At each time step t , the model produces an estimate of the conditional drift $\hat{\mu}_t$ and the conditional covariance matrix $\hat{\Sigma}_t$. New data points are then generated using the Euler–Maruyama discretisation scheme:

$$X_{t+\Delta t} = X_t + \hat{\mu}_t \Delta t + \hat{\sigma}_t \Delta W_t,$$

where $\Delta W_t \sim \mathcal{N}(0, \Delta t \cdot I)$ denotes a Brownian increment and $\hat{\sigma}_t$ is a matrix satisfying $\hat{\sigma}_t \hat{\sigma}_t^\top = \hat{\Sigma}_t$. In practice, $\hat{\sigma}_t$ is computed via the Cholesky decomposition of $\hat{\Sigma}_t$, which requires $\hat{\Sigma}_t$ to be positive semi-definite.

Example

To illustrate the data generation procedure based on the NJ-ODE model, consider the first iteration, starting from the initial condition

$$X_0 = \begin{bmatrix} 1.0 \\ 1.5 \\ 2.0 \end{bmatrix}.$$

From this input, we construct the augmented input vector $\tilde{X}_0 \in \mathbb{R}^9$, incorporating both linear and second-order features as follows:

$$\tilde{X}_0 = (1.0, 1.5, 2.0, 1.0, 2.25, 4.0, 1.5, 2.0, 3.0).$$

This vector is then provided as input to the trained NJ-ODE model, which outputs a predicted vector $\hat{\hat{X}}_0 \in \mathbb{R}^9$:

$$\hat{\hat{X}}_0 = (1.0002, 1.4997, 1.9997, 1.0004, 2.2492, 3.9987, 1.4999, 2.0003, 2.9990).$$

From the first three components of $\hat{\hat{X}}$, we extract the model's estimate of the conditional mean:

$$\mathbb{E}[X_{0.01} \mid \mathcal{A}_0] = \begin{bmatrix} 1.0002 \\ 1.4997 \\ 1.9997 \end{bmatrix}.$$

The remaining six components of $\hat{\hat{X}}$ correspond to second-order statistics. Using them, we reconstruct the estimated second-moment matrix:

$$\mathbb{E}[X_{0.01} X_{0.01}^\top \mid \mathcal{A}_0] = \begin{bmatrix} 1.0004 & 1.4999 & 2.0003 \\ 1.4999 & 2.2492 & 2.9990 \\ 2.0003 & 2.9990 & 3.9987 \end{bmatrix}.$$

We compute the drift and covariance estimates using the following formulas:

$$\hat{\mu}_0 = \frac{\mathbb{E}[X_{0.01} \mid \mathcal{A}_0] - X_0}{0.01}, \quad (\hat{\Sigma}_0)_{i,j} = \frac{\mathbb{E}[X_i X_j \mid \mathcal{A}_0] - X_i X_j}{0.01} - X_i \hat{\mu}_j - X_j \hat{\mu}_i.$$

Evaluating this expression numerically yields:

$$\hat{\mu}_0 = \begin{bmatrix} 0.0205 \\ -0.0347 \\ -0.0345 \end{bmatrix}, \quad \hat{\Sigma}_0 = \begin{bmatrix} 0.0002 & -0.0038 & 0.0224 \\ -0.0038 & 0.0235 & 0.0171 \\ 0.0224 & 0.0171 & 0.0880 \end{bmatrix}.$$

These estimates of the conditional drift and covariance form the basis for generating the next point using the Euler-Maruyama scheme. However, a critical requirement is to ensure that $\hat{\Sigma}_0$ is positive semi-definite. In this example, the eigenvalues of the estimate are approximately equal to $(-0.023, 0.017, 0.038)$, indicating a violation of this condition. Consequently, the Cholesky decomposition is not feasible at this step. A common strategy to ensure positive definiteness is to apply a regularisation by adding a small multiple of the identity matrix (Ledoit and Wolf, 2004):

$$\hat{\Sigma}_t^{(\alpha)} = \hat{\Sigma}_t + \alpha I_d, \quad \text{for some } \alpha > 0.$$

However, this method is ineffective here as the magnitude of the negative eigenvalue was comparable to the largest positive eigenvalue.

This issue highlights a practical challenge in modeling conditional covariance structures with neural networks, while ensuring that the predicted matrices remain symmetric and positive semi-definite.

4.3 Data Generation: Geometric Brownian Motion 1D

In this section, we apply Method I to generate synthetic data based on a one-dimensional GBM process $X \in \mathbb{R}^1$. The goal is to model both the process X_t and its square X_t^2 at each time step. In this section, we consider two approaches. First is a joint model that learns both X and X^2 simultaneously. And the second one is a separate modeling approach where two independent models are trained for X and X^2 , respectively.

Joint Modeling of X and X^2

In the first approach, we generate a dataset containing the values of X_t at each time step. The NJ-ODE framework supports flexible input and output dimensionalities, which allows us to extend the model's input representation. Specifically, we utilize the function called 'func_app1_X' in the NJ-ODE framework. By setting it to "power-2", automatically augments the input by including second-order monomials such as X_t^2 . The model is trained on the augmented dataset that includes both the original process X_t and its square X_t^2 . The resulting architecture thus learns to jointly model both the first-order and the second-order moments of the process.

After training, the learned model is used to generate synthetic data. Starting from an initial value X_0 , the model is used to estimate $\hat{\mu}_t$ and $\hat{\Sigma}_t$ at each time step. These estimates are then used within the Euler–Maruyama scheme to simulate the next point:

$$X_{t+\Delta t} = X_t + \hat{\mu}_t \Delta t + \hat{\sigma}_t \Delta W_t,$$

where $\hat{\sigma}_t$ is obtained via taking square root of $\hat{\Sigma}_t$, and $\Delta W_t \sim \mathcal{N}(0, \Delta t)$ denotes a standard Brownian increment.

This process is repeated iteratively until the final time horizon $T = 1$ is reached. At each step, the model conditions its prediction on previously generated values, enabling the construction of fully synthetic trajectories.

As shown in Figure 4.2, the synthetic trajectories display higher volatility than the reference GBM paths. This suggests that the model tends to overestimate the variability of the process. Furthermore, the generated paths terminate prematurely before reaching the final time step. This early termination is attributed to numerical instabilities encountered due to estimated covariance matrices $\hat{\Sigma}_t$ that are not positive. As a result, the computation of $\hat{\sigma}_t$ fails, which in turn prevents the generation of the next point in the trajectory via the Euler–Maruyama scheme.

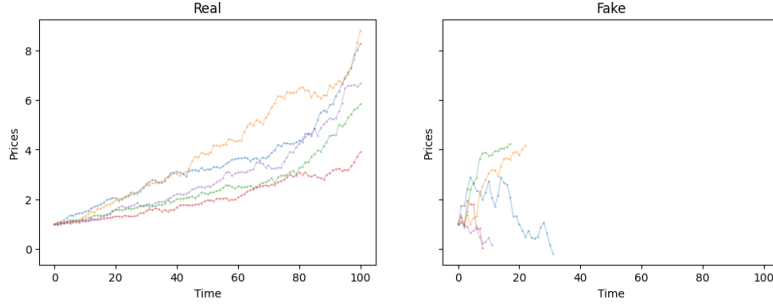


Figure 4.2: Comparison of 5 real and generated paths

4.3.1 Separate Modeling of X and X^2

In this approach, we construct two distinct datasets. One containing the values of the process X_t at each discrete time step, and another containing the squared values X_t^2 at the corresponding time steps. The goal is to develop separate models for the first moment of the corresponding process using the NJ-ODE framework.

We train two independent NJ-ODE models. One on the dataset corresponding to X_t , and the other on the dataset corresponding to X_t^2 . This decoupled modeling strategy allows each network to specialize in learning the temporal dynamics of its respective target.

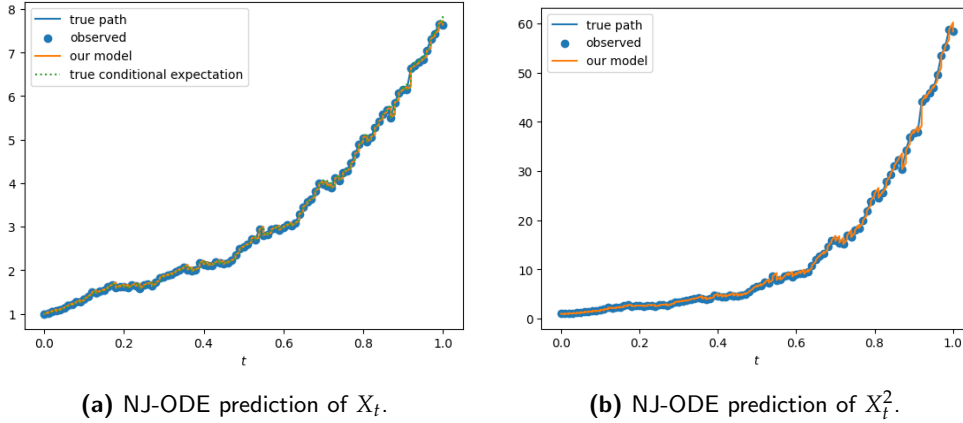

 Figure 4.3: Training results for NJ-ODE models trained separately on X_t and X_t^2 .

Figure 4.3 shows the training performance of both models on the test set. The close alignment between predicted and actual trajectories suggests that each model successfully captured the underlying structure of its respective process.

4. EXPERIMENTS FOR METHOD I

Following training, the NJ-ODE models are used to estimate the instantaneous drift $\hat{\mu}_t$ and diffusion $\hat{\Sigma}_t$ at each time step. These quantities are then employed within the Euler–Maruyama scheme to simulate synthetic trajectories.

The generation procedure begins from an initial condition X_0 . At each time step, the model estimates $\hat{\mu}_t$ and $\hat{\Sigma}_t$, which are used to compute the next value $X_{t+\Delta t}$. This process is repeated iteratively until the final time $T = 1$ is reached, resulting in a complete synthetic path. This generative procedure was repeated using NJ-ODE models with hidden state dimensions of 10, 30, 40, and 50. For each configuration, 10 generated trajectories were compared to 10 reference paths simulated from the true GBM model. The results are shown in Figure 4.4.

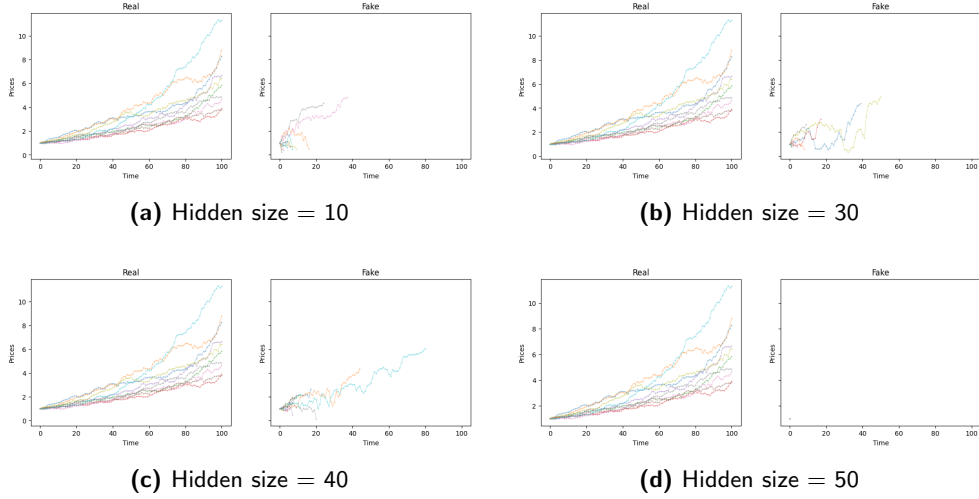


Figure 4.4: Comparison of 10 real and generated paths for varying hidden state sizes

The results indicate that the hidden state dimension of the NJ-ODE model affects the variability of the generated paths. Specifically, increasing the hidden size leads to reduced volatility. However, for the model with hidden size 50, a critical issue arises. The estimated initial covariance matrix $\hat{\Sigma}_0$ is not positive semi-definite. Consequently, its Chelosky decomposition is not possible, and this results in a failure of path generation.

This limitation motivates the development of a modified approach. The next chapter presents a modified generative approach that explicitly enforces the positive semi-definiteness of $\hat{\Sigma}_t$ during training by redefining the modeling and loss function design.

Experiments for Method II

All implementations were done using PyTorch. The code is available at: <https://github.com/zhaksylykov/thesis>.

5.1 Implementation Details

5.1.1 Datasets

In this chapter, we adopt the same dataset simulation procedures as described previously, with the key difference being that only ten percent of points are observable. Additionally, we introduce a new dataset based on a one-dimensional OU process.

For the second generative, we construct a dataset for the Z process. Where each has 100 time points compared to 101 for X . To ensure compatibility with the NJ-ODE framework, we prepend a zero matrix so $Z_0 = 0$.

Ornstein–Uhlenbeck 1D

- **SDE:**

$$dX_t = -k(X_t - m) dt + \sigma dW_t,$$

where W_t is a 1-dimensional Brownian motion.

- **Conditional expectation:**

$$\mathbb{E}[X_{t+s} \mid \mathcal{A}_t] = X_t e^{-ks} + m(1 - e^{-ks}),$$

- **Conditional variance:**

$$\text{Var}[X_{t+s} \mid \mathcal{A}_t] = \frac{\sigma^2}{2k}(1 - e^{-2ks}),$$

- **Parameters used:** $k = 2$, $m = 4$, $\sigma = 0.3$, $X_0 = 1$.

5.1.2 Architecture and Training

The model architecture and training procedure used in this method for the drift estimation model follow the setup described in the previous chapter. However, the critical differences are a hidden size of 100, 11 training epochs, a learning rate of 0.01, and the ReLU activation function is employed.

An additional critical consideration is to avoid applying input data scaling. This can be achieved by setting the parameter `'ode_input_scaling_func'` to the identity function. In prior experiments, tanh scaling was used.

The volatility estimation model employs the same architecture and training procedure as the drift estimation model. However, the difference in its use of *"vola"* loss function and is trained for 40 epochs.

5.2 Data Generation: Geometric Brownian Motion 1D

In this section, we generate synthetic data using Method II. While the overall procedure follows the same steps as described in the previous chapter, the key difference lies in the use of a two-model approach and a new loss function for the volatility estimator model. This framework ensures that the volatility estimator model outputs a scaled Chelosky decomposition of the estimated covariance matrix.

Drift Estimator Model

One NJ-ODE model is trained to learn the dynamics of X . As shown in the figure 5.1, the model successfully learns the underlying process.

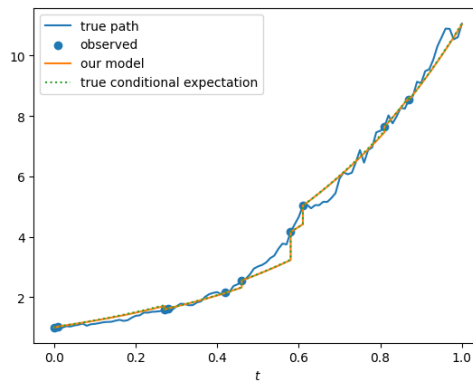


Figure 5.1: Training results.

For the one-dimensional GBM, the drift term is given by $\mu_t X_t$. The drift

term is estimated using the following equation

$$\frac{\mathbb{E}[X_{t+\Delta t}] - X_t}{\Delta t} = \frac{X_t(e^{\mu\Delta t} - 1)}{\Delta t}.$$

Figure 5.2 compares this theoretical drift with the estimates produced by the trained model. The results indicate that the model learns the drift dynamics of the underlying process.

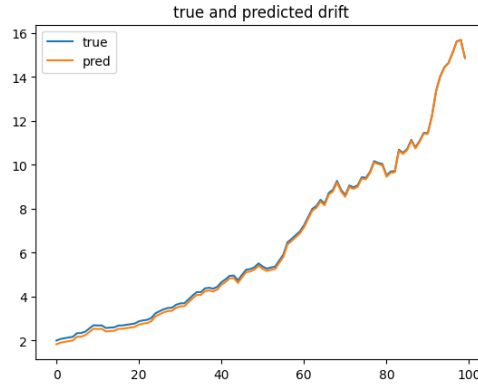


Figure 5.2: Comparison of true and estimated drift.

Volatility Estimator Model

A separate model is trained to estimate the diffusion component of the process X , using the dataset $[Z, Z^-, X]$. However, the training results appear inconclusive, as illustrated in Figure 5.3, indicating a discrepancy between the true and estimated volatility.

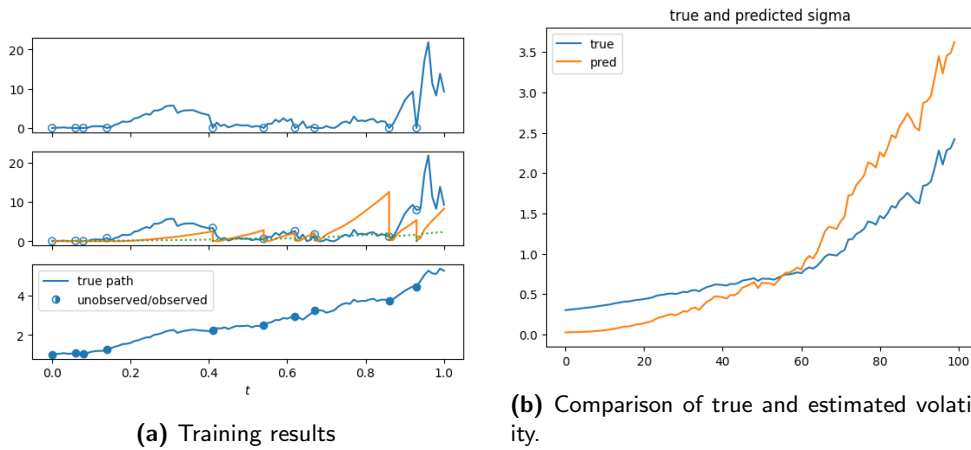


Figure 5.3: Volatility estimator training and evaluation.

5. EXPERIMENTS FOR METHOD II

To overcome this issue, the increment $X_t - X_{\tau(t)}$ is adjusted for drift, subtracting the drift term $X_{\tau(t)}(e^{\mu\Delta t} - 1)$. This leads to a modified definition of the Z-process

$$Z_t := (X_t - X_{\tau(t)} - X_{\tau(t)}(e^{\mu\Delta t} - 1))(X_t - X_{\tau(t)} - X_{\tau(t)}(e^{\mu\Delta t} - 1))^\top.$$

While the above definition relies on known parameters of the underlying process, a data-driven alternative uses the NJ-ODE framework to estimate the conditional expectation, \hat{X}_t . Subsequently, Z can be defined as

$$Z_t := (X_t - \hat{X}_t)(X_t - \hat{X}_t)^\top.$$

In this section, a model-dependent version of the Z process is used.

Figure 5.3 presents the training outcomes and resulting volatility estimates under this revised formulation. With the adjusted definition of Z, the trained model successfully captures the dynamics of the process. Also, the estimated volatility closely aligns with the theoretical values.

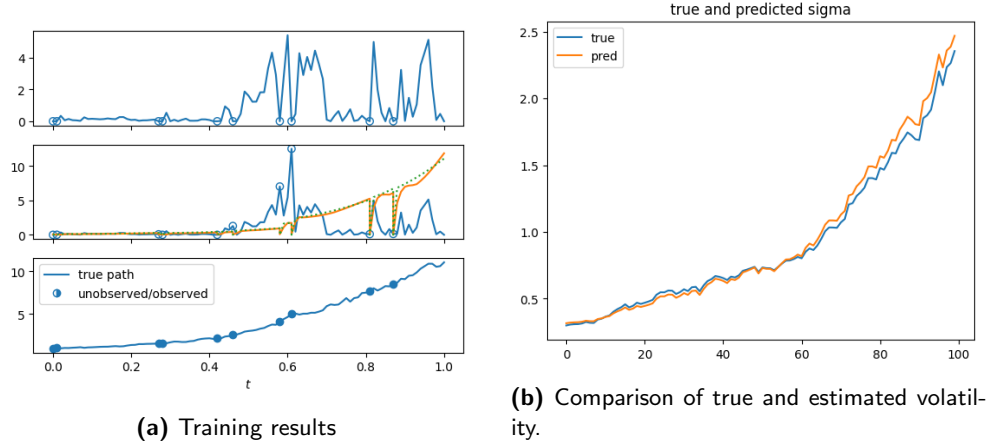


Figure 5.4: Volatility estimator training and evaluation.

Generation of Paths

Once trained, the NJ-ODE models are used to generate 1,000 synthetic paths via the Euler–Maruyama discretisation scheme. Starting from an initial value, the drift and volatility estimates are iteratively computed at each time step, allowing for the simulation of a complete trajectory. The volatility estimator model correctly learns to drop to zero and then recover. However, right after the jump, its first prediction is often too small. Using a “3 steps ahead” prediction instead fixes this and makes the results more stable.

Figure 5.5 compares 1k generated paths with 1k reference paths from the true GBM process.

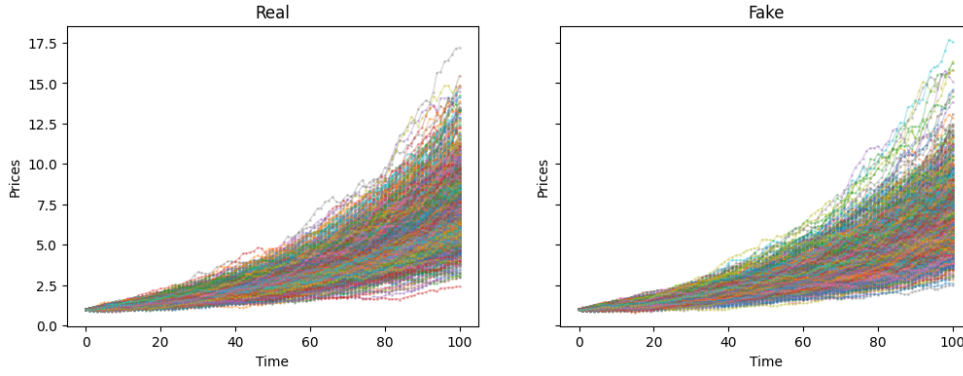


Figure 5.5: Comparison of 1k true and generated paths.

Figure 5.5 shows a histogram that compares the distribution of the true data to that of the generated data.

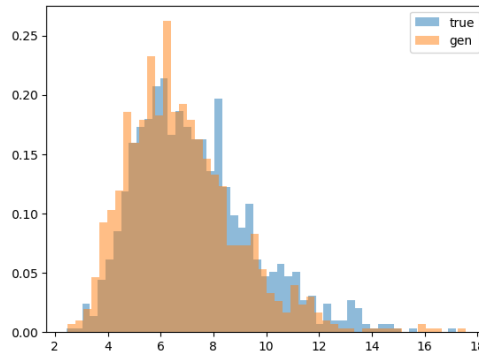


Figure 5.6: Histogram of true and generated data.

The parameters were then estimated, the details of which are provided in the Appendix. The results are presented in Table 5.1.

Parameter	True Data	Generated Data
Drift μ	1.985	1.919
Volatility σ	0.294	0.300

Table 5.1: Estimated parameters for true and generated GBM paths.

These results indicate that the proposed generative NJ-ODE (GenNJODE) framework works. It was able to produce synthetic data that closely resembles the original dataset.

5.3 Data Generation: Ornstein-Uhlenbeck 1D

In this section, we apply the NJ-ODE framework to simulate data from the one-dimensional OU process and evaluate the quality of the generated trajectories. The procedure follows the same steps as described in the previous section.

Drift Estimator Model

One NJ-ODE model is trained to learn the dynamics of X . As shown in the figure 5.7, the model successfully learns the underlying process.

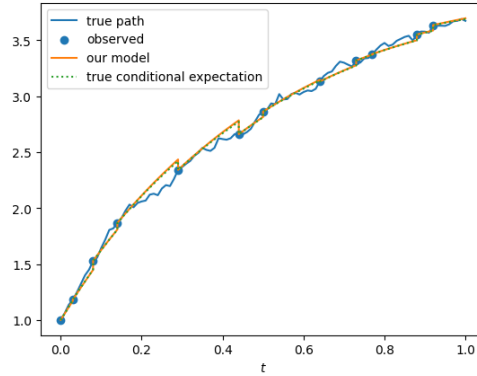


Figure 5.7: Training results.

For the one-dimensional OU, the drift term is given by $-k(X_t - m)$. The drift term is estimated using the following equation

$$\frac{\mathbb{E}[X_{t+\Delta t}] - X_t}{\Delta t} = \frac{(X_t - m)e^{-k\Delta t} + m - X_t}{\Delta t}.$$

Figure 5.8 compares this theoretical drift with the estimates produced by the trained model. The results indicate that the model learns the drift dynamics of the underlying process.

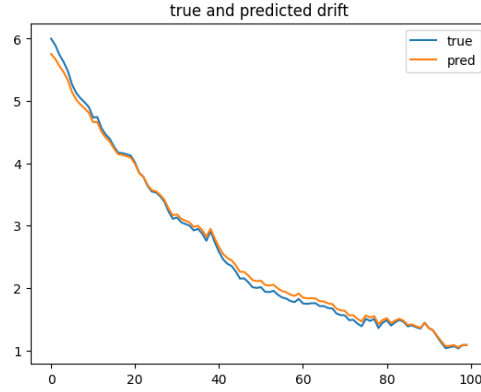


Figure 5.8: Comparison of true and estimated drift.

Volatility Estimator Model

After bias correction, the modified definition of the Z-process

$$Z_t = (X_t - (X_{\tau(t)} - m)e^{-k\Delta t} - m)(X_t - (X_{\tau(t)} - m)e^{-k\Delta t} - m)^\top.$$

A data-driven approach is to use NJ-ODE to estimate the conditional expectation \hat{X}_t . Then, Z could be defined as:

$$Z_t := (X_t - \hat{X}_t)(X_t - \hat{X}_t)^\top.$$

In this section, a model-dependent version of the Z process is used.

Figure 5.9 presents the training results and volatility estimates compared to the true values. The trained model successfully captures the dynamics of the process. However, the estimated volatility of 0.29 is not in close agreement with the theoretical value of 0.3.

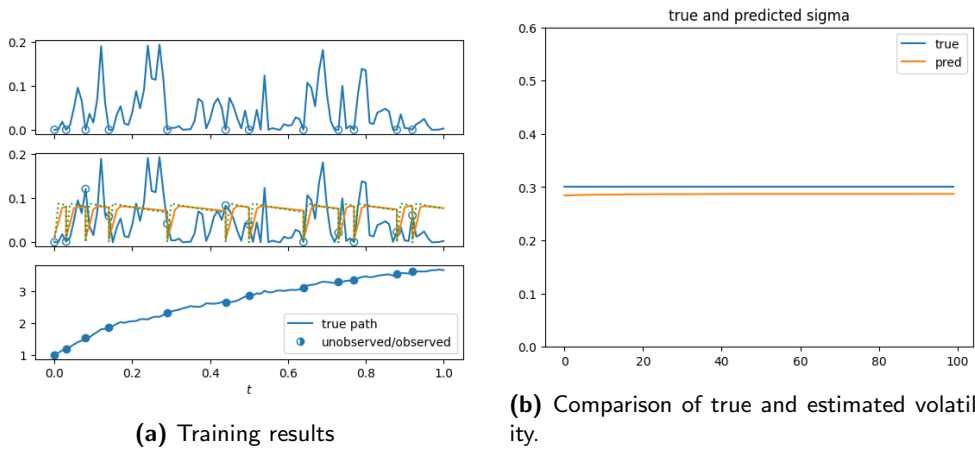


Figure 5.9: Volatility estimator training and evaluation.

Generation of Paths

Once trained, the NJ-ODE models are used to generate 10,000 synthetic paths via the Euler–Maruyama discretisation scheme. Starting from an initial value, the drift and volatility estimates are iteratively computed at each time step, allowing for the simulation of a complete trajectory. The “3-step-ahead” prediction technique is also applied in this section.

Figure 5.10 compares 10k generated paths with 10k reference paths from the true OU process.

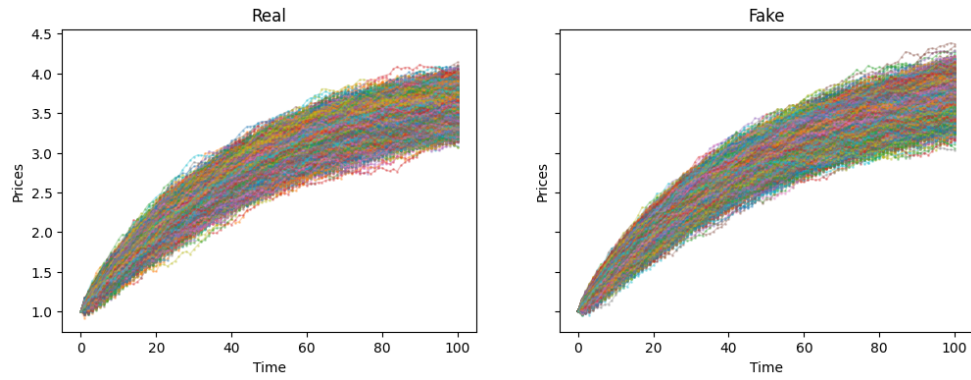


Figure 5.10: Comparison of 1k true and generated paths.

The parameters were then estimated, the details of which are provided in the Appendix. The results for both the real and synthetic datasets are presented in Table 5.2.

Parameter	True Data	Generated Data
Mean-Reversion Rate k	2.018	1.9272
Long-Term Mean m	4.0037	4.1277
Volatility σ	0.3030	0.2895

Table 5.2: Estimated OU process parameters for real and generated data.

Both the visual and quantitative results confirm that the GenNJODE framework is capable of approximating the OU process. However, there is potential for improvement in the volatility estimator model.

Conclusion

In this thesis, we explored the potential of adapting the NJ-ODEs framework from predictive to generative modeling. Two distinct methodologies for generating synthetic data were proposed. Method I aimed to estimate the instantaneous drift and instantaneous covariance directly from the predicted conditional moments. Method II independently modeled the instantaneous drift and instantaneous diffusion terms. This approach ensured that the estimated instantaneous covariance matrices are positive semi-definite.

While Method I was conceptually straightforward, it proved impractical in practice due to the lack of guaranteed positive semi-definiteness in the estimated instantaneous covariance matrices. That led to unstable and invalid synthetic data generation. Method II, on the other hand, showed some promising results. We refer to this approach as the GenNJODE framework.

Empirical assessments showed that the GenNJODE framework successfully learned both the drift and volatility of the underlying process. However, obtaining accurate volatility estimates required defining the Z process with a bias correction term.

Appendix A

Appendix

Since this thesis builds upon the NJ-ODE framework, the definitions presented in this appendix are from Florian's doctoral work (F. T. O. Krach, 2025).

A.0.1 Stochastic Process

We define the stochastic process $X := (X_t)_{t \in [0, T]} \in \mathcal{Q} \subset \mathbb{R}^{d_X}$, where \mathcal{Q} is closed, convex, and bounded. The process X is càdlàg and adapted to the filtered probability space

$$(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P}).$$

Since X may exhibit jumps, we define its left-limit at time t as

$$X_{t-} := \lim_{\varepsilon \downarrow 0} X_{t-\varepsilon}.$$

We define additional random variables on a separate probability space

$$(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{F}}, \tilde{\mathbb{P}})$$

as follows:

- $n : \tilde{\Omega} \rightarrow \mathbb{N}_{\geq 0}$ is a random variable representing the random number of observations, with $\mathbb{E}_{\tilde{\mathbb{P}}}[n] < \infty$.
- $t_i : \tilde{\Omega} \rightarrow [0, T]$ for $0 \leq i \leq n$ are sorted random variables representing the random observation times.

We define τ as the time of the last observation before a given time t , expressed as:

$$\tau : [0, T] \times \tilde{\Omega} \rightarrow [0, T], \quad (t, \tilde{\omega}) \mapsto \tau(t, \tilde{\omega}) := \max\{t_i(\tilde{\omega}) \mid 0 \leq i \leq n(\tilde{\omega}), t_i(\tilde{\omega}) \leq t\}.$$

Itô diffusion

Let $X := (X_t)_{t \in [0, T]}$ be a stochastic process satisfying the following SDE:

$$dX_t = \mu(t, X_t) dt + \sigma(t, X_t) dW_t, \quad (\text{A.1})$$

where $\{W_t\}_{t \in [0, T]}$ be a d_W -dimensional Brownian motion defined on a filtered probability space $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$, and $X_0 = x \in \mathbb{R}^{d_X}$ is starting point. The measurable functions $\mu : [0, T] \times \mathbb{R}^{d_X} \rightarrow \mathbb{R}^{d_X}$ and $\sigma : [0, T] \times \mathbb{R}^{d_X} \rightarrow \mathbb{R}^{d_X \times d_W}$ are referred to as the drift and diffusion coefficients, respectively.

The following regularity conditions are imposed:

- X is continuous and square integrable
- μ and σ are globally Lipschitz continuous their second component.
- μ is bounded and continuous in time, uniformly over space.
- The diffusion coefficient σ is càdlàg in time and square-integrable with respect to the Brownian motion W .

Information σ -algebra

We now define the information structure on the product probability space

$$(\Omega \times \tilde{\Omega}, \mathcal{F} \otimes \tilde{\mathcal{F}}, \mathbb{F} \otimes \tilde{\mathbb{F}}, \mathbb{P} \times \tilde{\mathbb{P}}),$$

where the filtration $(\mathcal{F} \otimes \tilde{\mathcal{F}})_t := \mathcal{F}_t \otimes \tilde{\mathcal{F}}_t$ for all $t \in [0, T]$ consists of the tensor-product σ -algebras.

We define the filtration of currently available information $\mathcal{A} := (\mathcal{A}_t)_{t \in [0, T]}$ by

$$\mathcal{A}_t := \sigma(X_{t_{i,j}}, t_i \mid t_i \leq t, j \in \{\ell \in \{1, \dots, d_X\} \mid M_{t_i, \ell} = 1\}),$$

where t_i are the random observation times, and $\sigma(\cdot)$ denotes the generated σ -algebra.

By the definition of τ , we have:

$$\mathcal{A}_t = \mathcal{A}_{\tau(t)} \quad \text{for all } t \in [0, T].$$

Observation mask

We define the observation mask $M = (M_k)_{0 \leq k \leq K}$ as a sequence of random variables defined on $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$, where each M_k takes values in $\{0, 1\}^{d_X}$ and is $\tilde{\mathcal{F}}_{t_k}$ -measurable. The j -th component of M_k , denoted $M_{k,j}$, indicates whether the j -th component of the stochastic process X_{t_k} was observed:

$$M_{k,j} = \begin{cases} 1, & \text{if } X_{t_k,j} \text{ is observed,} \\ 0, & \text{otherwise.} \end{cases}$$

By abuse of notation, we also write $M_{t_k} := M_k$.

A.0.2 Conditional Expectation

We denote by $\hat{X} := (\hat{X}_t)_{0 \leq t \leq T}$ the conditional expectation process of X , defined with respect to the information filtration \mathcal{A}_t as:

$$\hat{X}_t := \mathbb{E}_{\mathbb{P} \times \tilde{\mathbb{P}}} [X_t \mid \mathcal{A}_t].$$

We recall several fundamental properties of conditional expectations, which are central to the theory of stochastic processes. Let $X, X_1, X_2 \in \mathcal{F}$ be random variables such that $\mathbb{E}[|X|] < \infty$, and similarly for X_1 and X_2 . Let $\mathcal{A} \subseteq \mathcal{F}$ be a sub- σ -algebra. Then the following results hold (see Durrett, 2010):

1. The conditional expectation $\mathbb{E}[X \mid \mathcal{A}]$ exists.
2. It is unique up to sets of measure zero under \mathbb{P} .
3. $\mathbb{E}[X \mid \mathcal{A}]$ is integrable.
4. If X is \mathcal{A} -measurable, then $\mathbb{E}[X \mid \mathcal{A}] = X$ (i.e., full information).
5. If X is independent of \mathcal{A} , then $\mathbb{E}[X \mid \mathcal{A}] = \mathbb{E}[X]$ (i.e., no information).
6. Linearity: For $a, b \in \mathbb{R}$,

$$\mathbb{E}[aX_1 + bX_2 \mid \mathcal{A}] = a\mathbb{E}[X_1 \mid \mathcal{A}] + b\mathbb{E}[X_2 \mid \mathcal{A}].$$

7. Monotonicity: If $X_1 \leq X_2$ almost surely, then

$$\mathbb{E}[X_1 \mid \mathcal{A}] \leq \mathbb{E}[X_2 \mid \mathcal{A}].$$

8. Jensen's Inequality: For any convex function φ such that $\mathbb{E}[|\varphi(X)|] < \infty$,

$$\varphi(\mathbb{E}[X \mid \mathcal{A}]) \leq \mathbb{E}[\varphi(X) \mid \mathcal{A}].$$

9. Multiplicative Property: If $Y \in \mathcal{A}$ and $\mathbb{E}[|XY|] < \infty$, then

$$\mathbb{E}[XY \mid \mathcal{A}] = Y \cdot \mathbb{E}[X \mid \mathcal{A}].$$

10. Tower Property (Iterated Expectations): If $\mathcal{A}_1 \subseteq \mathcal{A}_2 \subseteq \mathcal{F}$, then

$$\mathbb{E}[\mathbb{E}[X \mid \mathcal{A}_2] \mid \mathcal{A}_1] = \mathbb{E}[X \mid \mathcal{A}_1], \quad \text{and} \quad \mathbb{E}[\mathbb{E}[X \mid \mathcal{A}_1] \mid \mathcal{A}_2] = \mathbb{E}[X \mid \mathcal{A}_1].$$

A.0.3 Signature

The signature transform provides a rich and universal representation of a path by encoding its sequential structure through an infinite collection of iterated integrals. This representation is particularly powerful because it allows any continuous functional on the path to be approximated by a linear function of the path's signature terms.

Definition A.1 (Signature) *Let $J \subset \mathbb{R}$ be a closed interval and let $X : J \rightarrow \mathbb{R}^d$ be a continuous path of finite variation. The signature of X over the interval J is defined as the infinite sequence:*

$$S(X) := \left(1, X_J^1, X_J^2, \dots\right),$$

where for each $m \geq 1$, the m -th level signature term is given by

$$X_J^m := \int_{u_1 < \dots < u_m, u_1, \dots, u_m \in J} dX_{u_1} \otimes \dots \otimes dX_{u_m} \in (\mathbb{R}^d)^{\otimes m}.$$

The map from a path X to its signature $S(X)$ is referred to as the *signature transform*. While the full signature is infinite-dimensional, in practical applications it is common to consider only a finite number of signature terms.

Definition A.2 (Truncated Signature) *Let $X : J \rightarrow \mathbb{R}^d$ be a continuous path of finite variation. The truncated signature of order $m \in \mathbb{N}$ is defined as:*

$$\pi_m(X) := \left(1, X_J^1, X_J^2, \dots, X_J^m\right),$$

i.e., the first $m + 1$ levels of the full signature of X .

For a path taking values in \mathbb{R}^d , the total dimension of the truncated signature of order m is given by:

$$\dim \pi_m(X) = \begin{cases} m + 1, & \text{if } d = 1, \\ \frac{d^{m+1} - 1}{d - 1}, & \text{if } d > 1. \end{cases}$$

A.0.4 Parameters Estimation for GBM

To estimate the drift and volatility parameters of the GBM, we followed a methodology similar to that described in Azimzadeh (2020).

Initially, the log-returns for each path were calculated at each time step. These were then aggregated across both time steps and paths.

The log-return at time t is defined as:

$$r_t := \log \left(\frac{X_{t+\Delta t}}{X_t} \right).$$

Under the GBM model, the log-returns r_t follow a normal distribution:

$$r_t \sim \mathcal{N}\left(\left(\mu - \frac{1}{2}\sigma^2\right)\Delta t, \sigma^2\Delta t\right),$$

where μ is the drift coefficient and σ is the volatility of the process.

After computing all log-returns, we flatten the matrix of returns into a single vector. Let m and s denote the sample mean and standard deviation of these aggregated log-returns:

$$m := \frac{1}{n} \sum_{i=1}^n r_i, \quad s := \sqrt{\frac{1}{n} \sum_{i=1}^n (r_i - m)^2},$$

where n is the total number of computed returns.

From these empirical moments, we estimate the GBM parameters as follows:

- The volatility σ is estimated as:

$$\hat{\sigma} = \frac{s}{\sqrt{\Delta t}}.$$

- The drift μ is estimated as:

$$\hat{\mu} = \frac{m}{\Delta t} + \frac{1}{2}\hat{\sigma}^2.$$

A.0.5 Parameters Estimation for OU

We estimate the parameters of the OU process using ordinary least squares (OLS) regression. This approach relies on least squares estimation applied to the discretised OU process, as outlined in Cantaro (n.d.).

The assumed discrete-time dynamics are given by:

$$X_{t+\Delta t} = \beta X_t + \alpha + \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, \sigma_\varepsilon^2),$$

where β and α are regression coefficients, and ε_t represents Gaussian noise.

Given simulated paths $\{X_t\}_{t=0}^T$, we define the vectors:

$$x = [X_0, X_1, \dots, X_{T-1}], \quad y = [X_1, X_2, \dots, X_T],$$

and fit the linear model $y = \beta x + \alpha + \varepsilon$ across all trajectories.

The continuous-time OU parameters are then recovered as:

- **Mean-reversion speed \hat{k} :**

$$\hat{k} = -\frac{1}{\Delta t} \log(\beta),$$

- **Long-term mean \hat{m} :**

$$\hat{m} = \frac{\alpha}{1 - \beta'},$$

- **Volatility $\hat{\sigma}$:**

$$\hat{\sigma} = \hat{\sigma}_\varepsilon \cdot \sqrt{\frac{2\hat{k}}{1 - \beta^2}}, \quad \text{where} \quad \hat{\sigma}_\varepsilon = \sqrt{\frac{1}{N-2} \sum_{i=1}^N (\varepsilon_i)^2}.$$

Bibliography

- Andersson, W. (2024). *Pd-nj-ode for predictions in convex spaces* [Master's Thesis]. ETH Zürich. <https://doi.org/10.3929/ethz-b-000690272>
- Azimzadeh, P. (2020). Maximum likelihood estimation of geometric brownian motion parameters [Available at https://parsiad.ca/blog/2020/maximum_likelihood_estimation_of_geometric_brownian_motion_parameters/].
- Bass, R. F. (2011). *Stochastic processes*. Cambridge University Press.
- Björk, T. (2009). *Arbitrage theory in continuous time* (3rd ed.). Oxford University Press.
- Black, F., & Scholes, M. (1973). The pricing of options and corporate liabilities. *Journal of Political Economy*, 81(3), 637–654. <https://doi.org/10.1086/260062>
- Bollerslev, T. (1986). Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics*, 31(3), 307–327.
- Box, G. E. P., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015). *Time series analysis: Forecasting and control* (5th). John Wiley & Sons.
- Bühler, H., Gonon, L., Teichmann, J., & Wood, B. (2019). Deep hedging. *Quantitative Finance*, 19(8), 1271–1291. <https://doi.org/10.1080/14697688.2019.1571683>
- Bühler, H., Horvath, B., Lyons, T., Perez Arribas, I., & Wood, B. (2020). A data-driven market simulator for small data environments. <https://doi.org/10.2139/ssrn.3632431>
- Cantaro, C. (n.d.). Ornstein–uhlenbeck process and applications [jupyter notebook] [Accessed: 2025-04-26]. <https://github.com/cantaro86/Financial-Models-Numerical-Methods/blob/master/6.1%20Ornstein-Uhlenbeck%20process%20and%20applications.ipynb>
- Chen, R. T. Q., Rubanova, Y., Bettencourt, J., & Duvenaud, D. (2018). Neural ordinary differential equations. *Advances in Neural Information Processing Systems (NeurIPS)*, 31, 6571–6583.
- Chevyrev, I., & Kormilitzin, A. (2016). A primer on the signature method in machine learning. <https://doi.org/10.48550/arXiv.1603.03788>

- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using rnn encoder-decoder for statistical machine translation. <https://arxiv.org/abs/1406.1078>
- Durrett, R. (2010). *Probability: Theory and examples* (4th ed.). Cambridge University Press.
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press. <https://www.deeplearningbook.org>
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. *Advances in Neural Information Processing Systems (NeurIPS)*, 27.
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778. <https://doi.org/10.1109/CVPR.2016.90>
- Heiss, J., Krach, F., Schmidt, T., & Tambe-Ndonfack, F. B. (2024). Nonparametric filtering, estimation and classification using neural jump odes. <https://arxiv.org/abs/2412.03271>
- Herrera, C., Krach, F., & Teichmann, J. (2021). Neural jump ordinary differential equations: Consistent continuous-time prediction and filtering. *International Conference on Learning Representations*. <https://openreview.net/forum?id=JFKR3WqwyXR>
- Heston, S. L. (1993). A closed-form solution for options with stochastic volatility with applications to bond and currency options. *The Review of Financial Studies*, 6(2), 327–343.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- Horvath, B., Plenk, J., Vuletić, M., & Saqur, R. (2025). Generative models in finance: Market generators, a paradigm shift in financial modeling. <https://doi.org/10.2139/ssrn.5284313>
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization.
- Kingma, D. P., & Welling, M. (2019). An introduction to variational autoencoders. *arXiv preprint arXiv:1906.02691*.
- Kloeden, P. E., & Platen, E. (1992). *Numerical solution of stochastic differential equations* (Vol. 23). Springer. <https://doi.org/10.1007/978-3-662-12616-5>
- Krach, F., Nübel, M., & Teichmann, J. (2022). Optimal estimation of generic dynamics by path-dependent neural jump odes. <https://arxiv.org/abs/2206.14284>
- Krach, F. T. O. (2025). *Neural jump ordinary differential equations* [Doctoral dissertation, ETH Zurich] [Doctoral thesis]. <https://doi.org/10.3929/ethz-b-000720717>

- Ledoit, O., & Wolf, M. (2004). A well-conditioned estimator for large-dimensional covariance matrices. *Journal of Multivariate Analysis*, 88(2), 365–411. [https://doi.org/10.1016/S0047-259X\(03\)00096-4](https://doi.org/10.1016/S0047-259X(03)00096-4)
- Potluru, V. K., Borrajo, D., Coletta, A., Dalmasso, N., El-Laham, Y., Fons, E., Ghassemi, M., Gopalakrishnan, S., Gosai, V., Kreačić, E., Mani, G., Obitalayo, S., Paramanand, D., & Veloso, T. B. (2023). Synthetic data applications in finance. <https://doi.org/10.48550/arXiv.2401.00081>
- Ruf, J., & Wang, W. (2020). Neural networks for option pricing and hedging: A literature review.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533–536. <https://doi.org/10.1038/323533a0>
- Shreve, S. E. (2004). *Stochastic calculus for finance ii: Continuous-time models*. Springer.
- Uhlenbeck, G. E., & Ornstein, L. S. (1930). On the theory of the brownian motion. *Physical Review*, 36(5), 823–841. <https://doi.org/10.1103/PhysRev.36.823>
- Vasicek, O. (1977). An equilibrium characterization of the term structure. *Journal of Financial Economics*, 5(2), 177–188. [https://doi.org/10.1016/0304-405X\(77\)90016-2](https://doi.org/10.1016/0304-405X(77)90016-2)
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 30.