

# Machine Learning for Stochastics

## A short intro to stochastic processes

Thorsten Schmidt

Abteilung für Mathematische Stochastik

[www.stochastik.uni-freiburg.de](http://www.stochastik.uni-freiburg.de)  
[thorsten.schmidt@stochastik.uni-freiburg.de](mailto:thorsten.schmidt@stochastik.uni-freiburg.de)

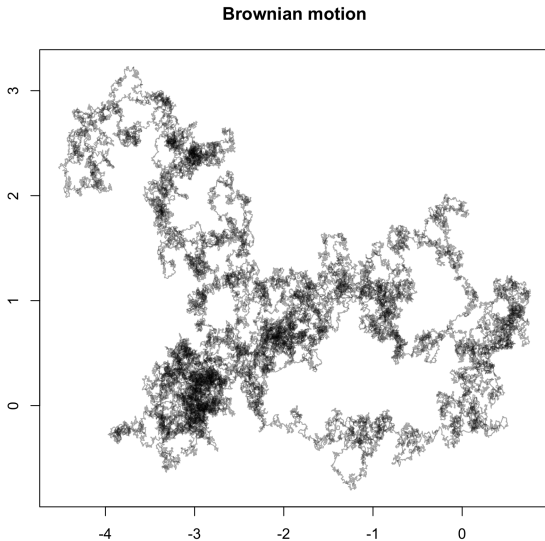
SS 2025

In Week 7 we cover the following topics:

- ▶ Stochastic processes in discrete / continuous time
- ▶ Brownian motion, Poisson process, Lévy process affine processes
- ▶ What are semimartingales?
- ▶ Stochastic differential equations
- ▶ Neural SDEs
- ▶ Signatures
- ▶ Path-dependency using signatures
- ▶ Neural SDEs

## Motivation

The botanist Robert Brown observed in 1827 the movement of a particle on water.



- ▶ Our goal is to make a precise mathematical framework for understanding and describing such phenomena.
- ▶ This is content of a full course (Stochastic Processes), so we can only scratch the surface. Many good books are out, for example: I. Karatzas and S. E. Shreve (1988). **Brownian Motion and Stochastic Calculus**. Springer Verlag. Berlin Heidelberg New York, J. Jacod and A.N. Shiryaev (2003). **Limit Theorems for Stochastic Processes**. 2nd. Berlin: Springer Verlag, Philip Protter (2004). **Stochastic Integration and Differential Equations**. 2nd. Springer Verlag. Berlin Heidelberg New York, D. Revuz and Marc Yor (2005). **Continuous Martingales and Brownian Motion**. 3rd ed. p. cm. Springer Verlag. Berlin Heidelberg New York. Also lecture notes of my course are available.
- ▶ Lets start in discrete time.

## Discrete time

- ▶ Discrete time is much simpler. A stochastic process (on a Polish space  $E$ ) is a sequence of random variables, i.e.

$$S = (S_t)_{t=0,1,\dots}.$$

- ▶ Examples include

$$S_t = \sum_{i=1}^t X_i,$$

where  $X_i$  are i.i.d., for example  $X_1 \sim \mathcal{N}(0, 1)$ . We can also have other distributions ! (essentially any ...)

- ▶ These processes have independent and stationary increments and are Markovian.
- ▶ The Polish space guarantees that we can compute conditional expectations for example. But also more general spaces are possible.

## Continuous time

- ▶ In continuous time, a stochastic process (say on  $\mathbb{R}^d$  for simplicity) is a family of random variables,

$$S = (S_t)_{t \geq 0}.$$

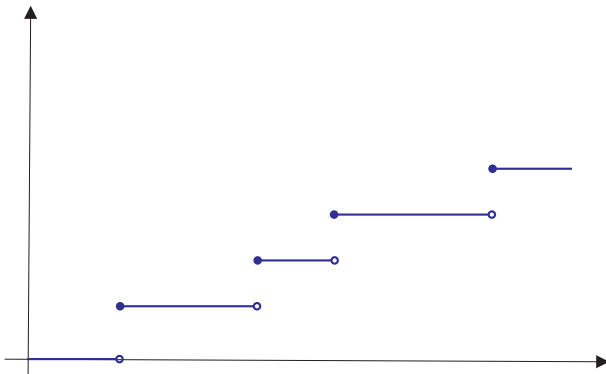
- ▶ We say that  $S$  has independent increments if  $S_t - S_{t'}$  is independent from  $S_s - S_{s'}$  whenever  $s' \leq s \leq t' \leq t$ .
- ▶ We say that the increments are stationary if  $S_{t+h} - S_t$  has the same distribution as  $S_h - S_0$  for all  $t \geq 0$  and all  $h \geq 0$ .
- ▶ A process with independent and stationary increments is called **Lévy process**.
- ▶ We have a number of examples: the Brownian motion has Gaussian increments, i.e.

$$W_t - W_s \sim \mathcal{N}(0, t - s)$$

- ▶ The **Poisson process** takes values in  $\mathbb{N}$  and

$$N_t - N_s \sim \text{Poisson}(\lambda(t - s)).$$

## Poisson process



Many interesting observations / extensions of the Poisson process are possible:  
Compound Poisson, time-inhomogeneous Poisson, doubly stochastic Poisson,  
Semi-Markov processes, Shot-Noise processes, ...

## Compound Poisson process

- ▶ A typical example of a process with jumps is the **compound Poisson process**.
- ▶ Consider a Poisson process  $N$  and independent i.i.d. random variables  $\xi_1, \xi_2, \dots$  and let

$$J_t = \sum_{i=1}^{N_t} \xi_i, \quad t \geq 0.$$

- ▶ Then  $J$  has independent increments, is of finite variation and  $W + J$  is a prototype of a semimartingale.



## The Itô theory

- ▶ We work on a filtered probability space  $(\Omega, \mathcal{F}, \mathbb{F}, P)$ , satisfying the usual conditions. A filtration  $\mathbb{F}$  is a family of increasing sub- $\sigma$ -fields.
- ▶ For a wide class of processes, i.e. **semimartingales** one can construct a stochastic integral. As usually, this is done by starting with simple processes.
- ▶ A simple process is given by

$$H = Y \mathbb{1}_{[S, T[},$$

where  $Y$  is  $\mathcal{F}_S$ -measurable and  $S \leq T$  are finite stopping times.

- ▶ A simple predictable process is given by

$$H = Y \mathbb{1}_{]S, T]}.$$

Note that if we have jumps, the integrand needs to be predictable, otherwise adapted is fine.

- ▶ For a semimartingale (a càdlàg process given as a sum of a finite variation process and a local martingale), we define

$$(H \cdot X)_t := \int_0^t H_s dX_s = \xi (X_{t \wedge T} - X_{t \wedge S}),$$

where we note that if  $X$  takes values in  $\mathbb{R}^d$ , we view  $\xi \in \mathbb{R}^d$  as an element of the dual space  $(\mathbb{R}^d)^* = \mathbb{R}^d$ . The stochastic integral in this form is then real-valued. Vector-valued stochastic integration is a bit more general.

## Theorem

*Let  $X$  be a semimartingale. The mapping  $H \mapsto H \cdot X$  has an extension from the simple processes to the space of locally bounded, predictable processes, such that*

- (i)  $H \cdot X$  is adapted and càdlàg.*
- (ii)  $H \mapsto H \cdot X$  is linear.*
- (iii) If predictable  $(H^n)$  converge pointwise to  $H$ , and  $|H^n| \leq K$  with a locally bounded, predictable process  $K$ , then*

$$(H^n \cdot X)_t \xrightarrow{P} (H \cdot X)_t \quad \forall t > 0.$$

So the limits of simple integrands build a well-defined theory for stochastic integrals, which is very powerful. For a proof we refer to Jacod and Shiryaev (2003), *op. cit.*

## The continuous case

If  $X$  is moreover continuous, we can extend the class even to locally square-integrable processes.

- We say a property of  $X$  holds locally, if there exists a sequence of stopping times  $(T_n) \rightarrow \infty$  such that the property holds for all  $X^{T_n}$ . The stopped process is defined by

$$X_t^{T_n} = X_{T_n \wedge t}, \quad t \geq 0.$$

One of the most powerful consequences is the following extension of the chain rule. We have to introduce the quadratic variation of the semimartingale  $X$ , which is given as the limit of square increments, i.e.

$$\langle X \rangle_t = \sum_{\Delta \rightarrow 0} (\Delta X_s)^2.$$

We know for example, that

$$\langle W \rangle_t = \langle N \rangle_t = t.$$

### Theorem (Itô-formula)

Let  $X = (X^1, \dots, X^d)$  be a  $d$ -dimensional semimartingale and  $f \in \mathcal{C}^2$ . Then  $f(X)$  is again a semimartingale and

$$\begin{aligned} f(X) &= f(X_0) + \sum_{i \leq d} D_i f(X_-) \cdot X^i \\ &\quad + \frac{1}{2} \sum_{i, j \leq d} D_{ij} f(X_-) \cdot \langle X^{i,c}, X^{j,c} \rangle \\ &\quad + \sum_{0 \leq s \leq \cdot} \left( f(X_s) - f(X_{s-}) - \sum_{i \leq d} D_i f(X_{s-}) \Delta X_s^i \right). \end{aligned} \tag{1}$$

Lets look at simpler special cases. For example, if  $X = W$ , then

$$\begin{aligned} f(X) &= f(X_0) + f'(X) \cdot X + \frac{1}{2} f''(X) \cdot \langle X \rangle \\ &= f(X_0) + \int_0^\cdot f'(X_s) dX_s + \frac{1}{2} \int_0^\cdot f''(X_s) ds \end{aligned} \quad (2)$$

- ▶ The Itô formula also opens the door to stochastic differential equations, our extension of ODEs.
- ▶ As an example think of

$$X_t = e^{W_t}, \quad t \geq 0$$

where  $W$  is a Brownian motion.

- ▶ Now we can apply the Itô formula and obtain ...

- ▶ The Itô formula also opens the door to stochastic differential equations, our extension of ODEs.
- ▶ As an example think of

$$X_t = e^{W_t}, \quad t \geq 0$$

where  $W$  is a Brownian motion.

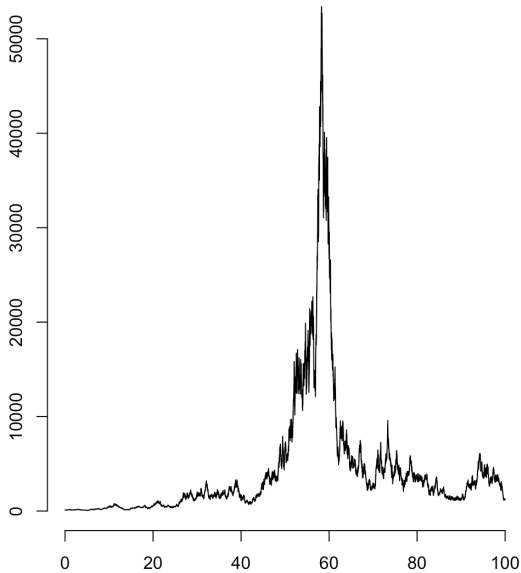
- ▶ Now we can apply the Itô formula and obtain ...
- ▶

$$X_t = X_0 + \int_0^t X_t dW_t + \frac{1}{2} \int_0^t X_t dt,$$

which could be read as

$$dX_t = X_t dW_t + \frac{1}{2} X_t dt.$$

We give an example of  $e^{\sigma W}$  with  $\sigma = 0.4$





## The martingale property

- ▶ We can already guess that this may cause problems. Note that  $W$  itself is a martingale, i.e. a process where

$$E[M_t | \mathcal{F}_s] = M_s, \quad 0 \leq s \leq t$$

- ▶ but we would expect that this fails for  $e^W$ .
- ▶ Indeed, we know that  $(H \cdot X)$  is a (local) martingale if  $X$  is a (local martingale) and as such we need to compensate  $e^W$  for the upward drift.
- ▶ Indeed, we can show with the Itô formula that

$$e^{W_t - t/2}, \quad t \geq 0$$

is a local martingale and we can calculate directly that it is indeed a martingale.

- ▶ However, it converges to 0 almost surely, since it is not closeable (i.e. there exist no  $M_\infty$  which conserves the martingale property).

We are interested in processes  $X$  which satisfy

$$\begin{aligned}dX_t &= b(t, X_t)dt + \sigma(t, X_t)dW_t, \\ X_0 &= \xi\end{aligned}\tag{3}$$

which we always understand as an abbreviation of

$$X_t = X_0 + \int_0^t b(s, X_s)ds + \int_0^t \sigma(s, X_s)dW_s, \quad t \geq 0.$$

## Definition

A **strong solution** of (3) is a continuous and adapted process  $X$ , such that

- (i)  $X$  is  $\mathbb{F}$ -adapted,
- (ii)  $P(X_0 = \xi) = 1$ ,
- (iii) for all  $0 \leq t < \infty$ ,  $1 \leq i \leq d$ ,  $1 \leq j \leq r$ , it hold that

$$\int_0^t (|b_i(s, X_s)| + \sigma_{ij}^2(s, X_s)) ds < \infty$$

a.s. and

- (iv) the first equation in (3) holds a.s.

Condition (iii) ensures that the integral is well-defined. The first property ensures that the  $ds$ -integral always exists, and the second guarantees local square integrability (which suffices for the stochastic integral when  $X$  is continuous).

# Uniqueness and Existence

## Theorem

*Are the coefficients  $b$  and  $\sigma$  locally Lipschitz-continuous, then strong uniqueness holds for the SDE (3).*

We say the global Lipschitz property holds, if

$$\|b(t, x) - b(t, y)\| + \|\sigma(t, x) - \sigma(t, y)\| \leq K\|x - y\|.$$

for all  $t \geq 0$  and  $x, y \in \mathbb{R}^d$ .

## Theorem

*Assume the global Lipschitz property holds together with*

$$\|b(t, x)\|^2 + \|\sigma(t, x)\|^2 \leq K^2(1 + \|x\|^2) \quad (4)$$

*and  $E[\|\xi\|^2] < \infty$ . Then there exists a strong solution.*

The idea of the proof is to use a modification of the Picard-Lindelöf iteration.

## Linear equations

If we consider the linear equation

$$\begin{aligned}dX_t &= \left( A(t)X_t + a(t) \right) dt + \sigma(t)dW_t, & 0 \leq t < \infty, \\X_0 &= \xi\end{aligned}\tag{5}$$

we expect that everything remains Gaussian and we can obtain a nice theory. As in the deterministic case,

$$X_t = \Phi(t) \left[ X_0 + \int_0^t \Phi^{-1}(s)a(s)ds + \int_0^t \Phi^{-1}(s)\sigma(s)dW_s \right], \quad 0 \leq t < \infty.\tag{6}$$

where

$$\dot{\Phi}(t) = A(t)\Phi(t)\tag{7}$$

is a fundamental solution of the homogeneous equation (this is a Matrix-differential equations, which has an easy solution for  $d = 1$ ).

# The Ornstein-Uhlenbeck process

Here we look for a solution of

$$dX_t = -\alpha X_t dt + \sigma dW_t,$$

with  $\alpha, \sigma > 0$ . This equation was already studied 1930 by the dutch physicists Leonard Ornstein and George Uhlenbeck.

With our solution method

$$X_t = X_0 + \int_0^t e^{-\alpha(t-s)} dW_s$$

is a solution (which is easily verified by the Itô-formula).

## The Brownian bridge

For the Brownian bridge  $B_t = W_t - t/T W_T$  one obtains an adapted representation via the SDE

$$dX_t = \frac{b - X_t}{T - t} dt + dW_t, \quad 0 \leq t < T, \quad X_0 = a,$$

with  $a, b \in \mathbb{R}$  and  $T > 0$ . As solution we obtain

$$X_t = a \left(1 - \frac{t}{T}\right) + \frac{b}{T}t + (T - t) \int_0^t \frac{dW_s}{T - s}.$$

We have the covariance function for  $a = b = 0$

$$\rho(s, t) = (s \wedge t) - \frac{st}{T},$$

which characterizes the standard Brownian bridge (and coincides with the covariance function of  $B$ ).

## The one-dimensional case

The one-dimensional case can be solved completely: consider

$$dX_t = (A(t)X_t + a(t))dt + (\Sigma(t)^\top X_t + \sigma(t)^\top) dW_t \quad (8)$$

with an  $r$ -dimensional Brownian motion  $W$ . We only assume that  $A, a$  and  $\Sigma, \sigma$  are  $\mathbb{F}$ -adapted, measurable and locally bounded.

Set

$$\begin{aligned}\zeta_t &= \int_0^t \Sigma(s)^\top dW_s - \frac{1}{2} \int_0^t \Sigma(s)^\top \Sigma(s) ds, \\ Z_t &= \exp \left( \int_0^t A(s) ds + \zeta_t \right).\end{aligned}$$

### Satz

*The unique strong solution of (8) is given by*

$$X_t = Z_t \left[ X_0 + \int_0^t \frac{1}{Z_s} (a(s) - \Sigma(s)^\top \sigma(s)) ds + \int_0^t \frac{1}{Z_s} \Sigma(s)^\top dW_s \right].$$



## Affine processes

A much more flexible class is the class of affine processes. They are well described in many textbooks, see for example [Damir Filipović \(2009\)](#). **Term Structure Models: A Graduate Course**. Springer Verlag. Berlin Heidelberg New York. Simulations are studied in [Aurélien Alfonsi \(2015\)](#). **Affine diffusions and related processes: simulation, theory and applications**. Springer. The general  $d$ -dimensional semimartingale case was studied in [Martin Keller-Ressel, Thorsten Schmidt, Robert Wardenga, et al. \(2019\)](#). „Affine processes beyond stochastic continuity“. In: **The Annals of Applied Probability** 29.6, pp. 3387–3437.

Lets keep it simple and consider the strong solution of the SDE

$$dX_t = b(X_t)dt + \sigma(X_t)dW_t,$$

with initial condition  $X_0 = x$ . We set  $a = \sigma\sigma^\top$ .

We call  $X$  affine if

$$E\left[e^{iu^\top X_T}|\mathcal{F}_t\right] = \exp\left(\phi(T-t, u) + \psi(T-t, u)^\top X_t\right)$$

for all  $u \in \mathbb{R}^d$  and all  $0 \leq t \leq T$ .

We assume that

$$\phi: \mathbb{C}^d \times \mathbb{R}_{\geq 0} \rightarrow \mathbb{C}$$

$$\psi: \mathbb{C}^d \times \mathbb{R}_{\geq 0} \rightarrow \mathbb{C}^d$$

are continuously differentiable functions.

## Theorem

Assume that  $X$  is affine. Then  $b$  and  $a$  are affine:

$$b(x) = \beta_0 + \sum_{i=1}^d \beta_i x_i$$

$$c(x) = \gamma_0 + \sum_{i=1}^d \gamma_i x_i$$

for  $\beta_0, \dots, \beta_d \in \mathbb{R}^d$  and  $\gamma_0, \dots, \gamma_d \in \mathbb{R}^{d \times d}$ .

Moreover,  $\phi$  and  $\psi$  solve the following Riccati-Equations :

$$\partial_t \phi(u, t) = \beta_0^\top \psi(u, t) + \frac{1}{2} \psi(u, t)^\top \gamma_0 \psi(u, t)$$

$$\phi(u, 0) = 0$$

$$\partial_t \psi_i(u, t) = \beta_i^\top \psi(u, t) + \frac{1}{2} \psi(u, t)^\top \gamma_i \psi(u, t)$$

$$\psi(u, 0) = u$$

Riccati equations are equations of the type  $f' = af^2$ .

We also have the converse direction

### Theorem

*Assume that  $b$  and  $a$  are affine and the Riccati equations have solutions such that*

$$\operatorname{real}(\phi(u, t) + \psi(u, t)^{\top} x) \leq 0 \quad \forall u \in i\mathbb{R}^d$$

*with  $t \geq 0$ ,  $x \in E$ , then  $X$  is affine.*

- ▶ We now come to the ML application to SDEs.
- ▶ Starting point is [Ricky TQ Chen et al. \(2018\)](#). „Neural ordinary differential equations“. In: [Advances in neural information processing systems 31](#), where neural ODEs have been introduced. Quite a number of authors generalize this to neural SDEs, for example [Belinda Tzen and Maxim Raginsky \(2019\)](#). „Neural stochastic differential equations: Deep latent gaussian models in the diffusion limit“. In: [arXiv preprint arXiv:1905.09883](#); [Junteng Jia and Austin R Benson \(2019\)](#). „Neural jump stochastic differential equations“. In: [Advances in Neural Information Processing Systems 32](#).

- ▶ We consider a driving  $r$ -dimensional semimartingale  $Z$  and study the solution to the SDE

$$dX_t = b_\theta(t, X_t)dt + \sigma_\theta(t, X_{t-})dZ_t,$$

$X_0 = x$  for some  $x \in \mathbb{R}^d$ .

- ▶ This is called a **neural SDE**, if  $b_\theta$  and  $\sigma_\theta$  are given by (deep) neural nets and trained on some given data.
- ▶ Typical examples involve that  $Z$  is a Brownian motion (classical neural SDEs) or  $Z$  is a pure-jump process (neural jump SDEs) or other variants.
- ▶ An first application to finance can be found in [Samuel N Cohen, Christoph Reisinger, and Sheng Wang \(2023\)](#). „Arbitrage-free neural-SDE market models“. In: [Applied Mathematical Finance](#) 30.1, pp. 1–46, however in the diffusion setting only.

## Path-dependence

- ▶ We note that there is in principle no difficulty to extend the setting to the path-dependent case, where



$$dX_t = b_\theta(t, X_{[0,t]})dt + \sigma_\theta(t, X_{[0,t]})dZ_t,$$

and  $X_{[0,t]} = \{X_s : s \in [0, t]\}$  denotes the path of  $X$  from 0 to  $t$ . Now the functions  $b$  and  $\sigma$  live on the path space, here the space of càdlàg functions over  $[0, t]$ .

- ▶ One question is how to find a nice representation of these functions, which can be (in an explainable way) be done by signatures.
- ▶ Existence and uniqueness is guaranteed by classical restrictions on the functions.
- ▶ The universal approximation theorem (for Banach spaces) shows that all functions on  $D([0, T])$  can be arbitrarily well approximated by a neural network.
- ▶ However - the training is more complicated. A classical question is of course when you observe only **one** path, how can you learn the parameters of  $X$  in a good way. This is typically done via maximum likelihood (which can also be done for stochastic processes through the Girsanov theorem).

- Of course, one can also consider stochastic control in the following form

$$dX_t = b_\theta(t, X_t, a_t)dt + \sigma_\theta(t, X_{t-}, a_{t-})dZ_t,$$

with a stochastic control  $a$ . One can derive neural HJB equations or solve using methods from reinforcement learning ( $X$  here is a continuous-time Markov decision process - if  $Z$  is in discrete time, we recover our setting from the first part of our lectures).