

09 Apr 2021

# REVERSE TIME STOCHASTIC DIFFERENTIAL EQUATIONS [ FOR GENERATIVE MODELLING ]

'If I Could Turn Back Time' by Cher (1989)

$$dX_{\tau} = \left( -\mu(x_{1-\tau}) + \frac{1}{p(x_{1-\tau})} \partial_{x_{1-\tau}} \left[ \sigma^2(x_{1-\tau}) p(x_{1-\tau}) \right] \right) d\tau + \sigma(x_{1-\tau}) dW_{\tau}$$

The Kolmogorov forward equation is identical to the Fokker Planck equation and states

$$\partial_t p(x_t) = -\partial_{x_t} [\mu(x_t) p(x_t)] + \frac{1}{2} \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] \quad (2)$$

It describes the evolution of a probability distribution  $p(x_t)$  forward in time. We can quite frankly think of it as, for example, a Normal distribution being slowly transformed into an arbitrary complex distribution according to the drift and diffusion parameters  $\mu(x_t)$  and  $\sigma(x_t)$ .

The Kolmogorov backward equation for  $s \geq t$  is defined as

$$-\partial_t p(x_s | x_t) = \mu(x_t) \partial_{x_t} p(x_s | x_t) + \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) \quad (3)$$

and it basically answers the question how the probability of  $x_s$  at a later point in time changes as we change  $x_t$  at an earlier point in time. The Kolmogorov backward equation is somewhat confounding with respect to time as we're taking the partial derivative with respect to the earlier time step  $t$  on which we are also conditioning. But we can think of it as asking 'How does the probability of  $x_s$  at the later point in time  $s$  change, as we slowly evolve the probability distribution backwards through time and condition on  $x_t$ '.

Taking inspiration from our crude example earlier, the backward equation offers a partial differential equation which we can solve backward in time, which would correspond to evolving the arbitrarily complex distribution backwards to our original Normal distribution. Unfortunately there is no corresponding stochastic differential equation with a drift and diffusion term that describes the evolution of a random variable backwards through time in terms of a stochastic differential equation.

This is where the remarkable result from Anderson (1982) comes into play.

The granddaddy of all probabilistic equations, Bayes theorem, tells us that a joint distribution can be factorized by conditioning:  $p(x_s, x_t) = p(x_s | x_t) p(x_t)$  with the time ordering  $t \leq s$ . Why do we invoke the joint probability  $p(x_s, x_t)$  we might ask? What we're trying to achieve is to derive a stochastic differential equation that tells us from what values of  $x_t$  we can arrive at  $x_s$ . We can ask ourselves what the partial differential equation would be that describes the evolution of the joint distribution over time. First multiplying both sides of Bayes theorem with minus one and taking the derivative with respect to time  $t$ , we obtain via the product rule

$$-\partial_t p(x_s, x_t) = -\partial_t [p(x_s | x_t) p(x_t)] \quad (4)$$

$$= \underbrace{-\partial_t p(x_s | x_t) p(x_t)}_{\text{KBE}} - \underbrace{p(x_s | x_t) \partial_t p(x_t)}_{\text{KFE}} \quad (5)$$

into which we can plug in the Kolmogorov forward (KFE) and Kolmogorov backward (KBE) equations,

$$-\partial_t p(x_s | x_t) p(x_t) - p(x_s | x_t) \partial_t p(x_t) \quad (6)$$

$$= \left( \mu(x_t) \partial_{x_t} p(x_s | x_t) + \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) \right) p(x_t) \quad (7)$$

$$+ p(x_s | x_t) \left( \partial_{x_t} [\mu(x_t) p(x_t)] - \frac{1}{2} \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] \right) \quad (8)$$

The derivative occurring in the backward Kolmogorov equation are

$$\partial_{x_t} p(x_s | x_t) = \partial_{x_t} \left[ \frac{p(x_s, x_t)}{p(x_t)} \right] \quad (9)$$

$$= \frac{\partial_{x_t} p(x_s, x_t) p(x_t) - p(x_s, x_t) \partial_{x_t} p(x_t)}{p^2(x_t)} \quad (10)$$

$$= \frac{\partial_{x_t} p(x_s, x_t)}{p(x_t)} - \frac{p(x_s, x_t) \partial_{x_t} p(x_t)}{p^2(x_t)} \quad (11)$$

The next step is to evaluate the derivative of the products in the forward Kolmogorov equation.

$$\partial_{x_t} [\mu(x_t) p(x_t)] = \partial_{x_t} \mu(x_t) p(x_t) + \mu(x_t) \partial_{x_t} p(x_t) \quad (12)$$

$$\partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] = \partial_{x_t}^2 \sigma^2(x_t) p(x_t) + 2 \partial_{x_t} \sigma^2(x_t) \partial_{x_t} p(x_t) + \sigma^2(x_t) \partial_{x_t}^2 p(x_t) \quad (13)$$

Substituting the derivatives of the probability distributions accordingly we obtain

$$-\partial_t p(x_s, x_t) = -\partial_t [p(x_s | x_t) p(x_t)] \quad (14)$$

$$= -\partial_t p(x_s | x_t) p(x_t) - p(x_s | x_t) \partial_t p(x_t) \quad (15)$$

$$= \left( \mu(x_t) \partial_{x_t} p(x_s | x_t) + \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) \right) p(x_t) \quad (16)$$

$$+ p(x_s | x_t) \left( \partial_{x_t} [\mu(x_t) p(x_t)] - \frac{1}{2} \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] \right) \quad (17)$$

$$= \mu(x_t) \partial_{x_t} p(x_s | x_t) p(x_t) + \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) p(x_t) \quad (18)$$

$$+ p(x_s | x_t) \partial_{x_t} \mu(x_t) p(x_t) + p(x_s | x_t) \mu(x_t) \partial_{x_t} p(x_t) \quad (19)$$

$$- \frac{1}{2} p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] \quad (20)$$

$$= \mu(x_t) \left( \frac{\partial_{x_t} p(x_s, x_t)}{p(x_t)} - \frac{p(x_s, x_t) \partial_{x_t} p(x_t)}{p^2(x_t)} \right) p(x_t) \quad (21)$$

$$+ p(x_s | x_t) \partial_{x_t} \mu(x_t) p(x_t) + p(x_s | x_t) \mu(x_t) \partial_{x_t} p(x_t) \quad (22)$$

$$+ \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) p(x_t) - \frac{1}{2} p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] \quad (23)$$

$$= \mu(x_t) \left( \frac{\partial_{x_t} p(x_s, x_t)}{p(x_t)} - \frac{p(x_s, x_t) \partial_{x_t} p(x_t)}{p^2(x_t)} \right) \quad (24)$$

$$+ p(x_s | x_t) \partial_{x_t} \mu(x_t) p(x_t) + p(x_s | x_t) \mu(x_t) \partial_{x_t} p(x_t) \quad (25)$$

$$+ \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) p(x_t) - \frac{1}{2} p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] \quad (26)$$

$$= \mu(x_t) \left( \frac{\partial_{x_t} p(x_s, x_t)}{p(x_t)} - \frac{p(x_s | x_t) \partial_{x_t} p(x_t)}{p^2(x_t)} \right) \quad (27)$$

$$+ p(x_s | x_t) \partial_{x_t} \mu(x_t) + \frac{p(x_s | x_t) \mu(x_t) \partial_{x_t} p(x_t)}{p^2(x_t)} \quad (28)$$

$$+ \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) p(x_t) - \frac{1}{2} p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] \quad (29)$$

$$= \underbrace{\mu(x_t) \partial_{x_t} p(x_s, x_t) + p(x_s, x_t) \partial_{x_t} \mu(x_t)}_{\text{product rule}} \quad (30)$$

$$+ \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) p(x_t) - \frac{1}{2} p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] \quad (31)$$

$$= \partial_{x_t} [\mu(x_t) p(x_s, x_t)] \quad (32)$$

$$+ \underbrace{\frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) p(x_t)}_{(1)} - \underbrace{\frac{1}{2} p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)]}_{(2)} \quad (33)$$

In order to transform the partial differential equation above into a form from which we can deduce an equivalent stochastic differential equation, we match the terms of the second order derivatives with the following identity,

$$\frac{1}{2} \partial_{x_t}^2 [p(x_s, x_t) \sigma^2(x_t)] \quad (34)$$

$$= \frac{1}{2} \partial_{x_t}^2 [p(x_s | x_t) p(x_t) \sigma^2(x_t)] \quad (35)$$

$$= \frac{1}{2} \partial_{x_t}^2 p(x_s | x_t) p(x_t) \sigma^2(x_t) + \partial_{x_t} [p(x_t) \sigma^2(x_t)] \partial_{x_t} p(x_s | x_t) + \frac{1}{2} \partial_{x_t}^2 [p(x_t) \sigma^2(x_t)] p(x_s | x_t) \quad (36)$$

$$= \underbrace{\frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) p(x_t) + \partial_{x_t} [p(x_t) \sigma^2(x_t)] \partial_{x_t} p(x_s | x_t)}_{(1)} + \underbrace{\frac{1}{2} p(x_s | x_t) \partial_{x_t}^2 [p(x_t) \sigma^2(x_t)]}_{(2)} \quad (37)$$

by observing that the terms (1) and (2) occur in both equations. We can see from the expansion of the derivative above that we can combine the terms in our derivation if we expand the "center term". Furthermore we can employ the identity  $-\frac{1}{2} X = -X + \frac{1}{2} X$  to obtain

$$-\partial_t p(x_s, x_t) = \partial_{x_t} [\mu(x_t) p(x_s, x_t)] \quad (38)$$

$$+ \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) p(x_t) - \frac{1}{2} p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] \quad (39)$$

$$= \partial_{x_t} [\mu(x_t) p(x_s, x_t)] \quad (40)$$

$$+ \frac{1}{2} \sigma^2(x_t) p(x_t) \partial_{x_t}^2 p(x_s | x_t) - \underbrace{\frac{1}{2} p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)]}_{-\frac{1}{2} X = -X + \frac{1}{2} X} \quad (41)$$

$$\underbrace{\pm \partial_{x_t} p(x_s | x_t) \partial_{x_t} [p(x_t) \sigma^2(x_t)]}_{\text{complete the square}} \quad (42)$$

$$= \partial_{x_t} [\mu(x_t) p(x_s, x_t)] + \frac{1}{2} \sigma^2(x_t) \partial_{x_t}^2 p(x_s | x_t) p(x_t) \quad (43)$$

$$- \underbrace{p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] + \frac{1}{2} p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)]}_{-\frac{1}{2} X = -X + \frac{1}{2} X} \quad (44)$$

$$\pm \partial_{x_t} p(x_s | x_t) \partial_{x_t} [p(x_t) \sigma^2(x_t)] \quad (45)$$

$$= \partial_{x_t} [\mu(x_t) p(x_s, x_t)] + \frac{1}{2} \partial_{x_t}^2 [p(x_s | x_t) p(x_t) \sigma^2(x_t)] \quad (46)$$

$$- \underbrace{p(x_s | x_t) \partial_{x_t}^2 [\sigma^2(x_t) p(x_t)] - \partial_{x_t} p(x_s | x_t) \partial_{x_t} [p(x_t) \sigma^2(x_t)]}_{-\partial_{x_t} [p(x_s | x_t) \partial_{x_t} [\sigma^2(x_t) p(x_t)]] \text{ (product rule)}} \quad (47)$$

$$= \partial_{x_t} [\mu(x_t) p(x_s, x_t)] + \frac{1}{2} \partial_{x_t}^2 [p(x_s, x_t) \sigma^2(x_t)] \quad (48)$$

$$- \partial_{x_t} [p(x_s | x_t) \partial_{x_t} [\sigma^2(x_t) p(x_t)]] \quad (49)$$

What remains to be done is to combine the joint probability and the conditional probability in the first order derivative terms to combine them,

$$-\partial_t p(x_s, x_t) = \partial_{x_t} [\mu(x_t) p(x_s, x_t) - p(x_s | x_t) \partial_{x_t} [\sigma^2(x_t) p(x_t)]] \quad (50)$$

$$+ \frac{1}{2} \partial_{x_t}^2 [p(x_s, x_t) \sigma^2(x_t)] \quad (51)$$

$$= \partial_{x_t} \left[ p(x_s, x_t) \left( \mu(x_t) - \frac{1}{p(x_t)} \partial_{x_t} [\sigma^2(x_t) p(x_t)] \right) \right] \quad (52)$$

$$+ \frac{1}{2} \partial_{x_t}^2 [p(x_s, x_t) \sigma^2(x_t)] \quad (53)$$

$$= -\partial_{x_t} \left[ p(x_s, x_t) \left( -\mu(x_t) + \frac{1}{p(x_t)} \partial_{x_t} [\sigma^2(x_t) p(x_t)] \right) \right] \quad (54)$$

$$+ \frac{1}{2} \partial_{x_t}^2 [p(x_s, x_t) \sigma^2(x_t)] \quad (55)$$

the result of which is in the form of a Kolmogorov forward equation, although using the joint probability distribution  $p(x_s, x_t)$ . For the time ordering of  $t \leq s$ , we can observe that the term  $-\partial_t p(x_s, x_t)$  describes the change of the probability distribution as we move backward in time. In accordance with Leibniz' rule we can marginalize over  $x_s$  without interfering with the partial derivative  $\partial_t$ , to obtain

$$-\partial_t p(x_t) = -\partial_{x_t} \left[ p(x_t) \left( -\mu(x_t) + \frac{1}{p(x_t)} \partial_{x_t} [\sigma^2(x_t) p(x_t)] \right) \right] \quad (56)$$

$$+ \frac{1}{2} \partial_{x_t}^2 [p(x_t) \sigma^2(x_t)] \quad (57)$$

and introduce the time reversal  $\tau \doteq 1 - t$  which, with respect to the integration with respect to the flow of time, yields

$$-\partial_t p(x_t) = \partial_{\tau} p(x_{1-\tau}) \quad (58)$$

$$= -\partial_{x_{1-\tau}} \left[ p(x_{1-\tau}) \left( -\mu(x_{1-\tau}) + \frac{1}{p(x_{1-\tau})} \partial_{x_{1-\tau}} [\sigma^2(x_{1-\tau}) p(x_{1-\tau})] \right) \right] \quad (59)$$

$$+ \frac{1}{2} \partial_{x_{1-\tau}}^2 [p(x_{1-\tau}) \sigma^2(x_{1-\tau})] \quad (60)$$

which finally gives us a stochastic differential equation analogous to the Fokker-Planck/forward Kolmogorov equation that we can solve backward in time:

$$dX_{\tau} = \left( -\mu(x_{1-\tau}) + \frac{1}{p(x_{1-\tau})} \partial_{x_{1-\tau}} [\sigma^2(x_{1-\tau}) p(x_{1-\tau})] \right) d\tau + \sigma(x_{1-\tau}) dW_{\tau} \quad (61)$$

where  $\tilde{W}_t$  is a Wiener process that flows backward in time.

By keeping the  $\sigma^2(x_t)$  constant and independent of  $x_t$  and applying the log-derivative trick, the drift simplifies to

$$dX_{\tau} = \left( -\mu(x_{1-\tau}) + \frac{1}{p(x_{1-\tau})} \partial_{x_{1-\tau}} \left[ \overbrace{\sigma^2(x_{1-\tau})}^{=\sigma^2} p(x_{1-\tau}) \right] \right) d\tau + \sigma(x_{1-\tau}) dW_{\tau} \quad (62)$$

$$= \left( -\mu(x_{1-\tau}) + \frac{\sigma^2}{p(x_{1-\tau})} \partial_{x_{1-\tau}} p(x_{1-\tau}) \right) d\tau + \sigma(x_{1-\tau}) dW_{\tau} \quad (63)$$

$$= \left( -\mu(x_{1-\tau}) + \sigma^2 \partial_{x_{1-\tau}} \log p(x_{1-\tau}) \right) d\tau + \sigma(x_{1-\tau}) d\tilde{W}_{\tau} \quad (64)$$