

# Using the time evolution of probabilities

**Concept.** To extend the formalism that we developed in previous chapters to understand the time evolution of state probabilities. With this knowledge we shall derive popular machine learning methods like Gaussian processes, as well as creating numerical procedures to directly solve for discrete state spaces. In the latter case these can be very intensive computations to perform, so in this chapter we shall also discuss the practical limitations of a fully numerical approach. For the mathematically-inclined, this chapter will take a detailed look at how our formalism can be extended to focus on the time evolution of probabilities. For the programmers, the software described in this chapter lives in this public Git repository: <https://github.com/umbralcalc/dennm-torch>.

## 2.1 Probabilistic formalism

In this section we will return to the stochadex formalism that we introduced in the first chapter of this book. As we discussed at that point; this formalism is appropriate for sampling from nearly every stochastic phenomenon that one can think of. We are going to extend this description to consider what happens to the probability that the state history matrix takes a particular set of values over time.

So, how do we begin? In the first chapter, we defined the general stochastic process with the formula  $X_{t+1}^i = F_{t+1}^i(X_{0:t}, z, t)$ . This equation also has an implicit *master equation* associated to it that fully describes the time evolution of the *probability density function*  $P_{t+1}(X|z)$  of  $X_{0:t+1} = X$  given that the parameters of the process are  $z$ . This can be written as

$$P_{t+1}(X|z) = P_t(X'|z)P_{(t+1)t}(x|X', z), \quad (2.1)$$

where for the time being we are assuming the state space is continuous in each of the matrix elements and  $P_{(t+1)t}(x|X', z)$  is the conditional probability that  $X_{t+1} = x$  given that  $X_{0:t} = X'$  at time  $t$  and the parameters of the process are  $z$ .

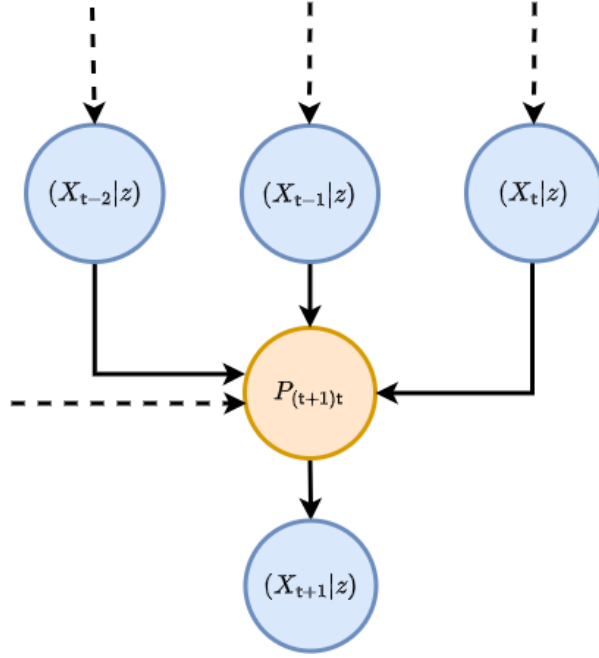


Figure 2.1: Graph representation of Eqs. (2.1) and (2.5).

Consider what happens when we extend the chain of conditional probabilities in Eq. (2.1) back in time by one step. In doing so, we retrieve a joint probability of rows  $X_{t+1} = x$  and  $X_t = x'$  on the right hand side of the expression

$$P_{t+1}(X|z) = P_{t-1}(X''|z)P_{(t+1)t(t-1)}(x, x'|X'', z). \quad (2.2)$$

Since Eqs. (2.1) and (2.2) are both valid ways to obtain  $P_{t+1}(X|z)$  we can average between them without loss of generality in the original expression, like this

$$P_{t+1}(X|z) = \frac{1}{2} [P_t(X'|z)P_{(t+1)t}(x|X', z) + P_{t-1}(X''|z)P_{(t+1)t(t-1)}(x, x'|X'', z)]. \quad (2.3)$$

Following this line of reasoning to its natural conclusion, Eq. (2.1) can hence be generalised to consider all possible joint distributions of rows at different timesteps like this

$$P_{t+1}(X|z) = \frac{1}{t} \sum_{t''=0}^t P_{t''}(X''|z)P_{(t+1)t...t''}(x, x', \dots |X'', z). \quad (2.4)$$

If we wanted to just look at the distribution over the latest row  $X_{t+1} = x$ , we could achieve this through marginalisation over all of the previous matrix rows in Eq. (2.1) like this

$$P_{t+1}(x|z) = \int_{\Omega_t} dX' P_{t+1}(X|z) = \int_{\Omega_t} dX' P_t(X'|z)P_{(t+1)t}(x|X', z). \quad (2.5)$$

But what is  $\Omega_t$ ? You can think of this as just the domain of possible matrix  $X'$  inputs into the integral which will depend on the specific stochastic process we are looking at.

The symbol  $dX'$  in Eq. (2.5) is our shorthand notation throughout the book for computing the sum of integrals over previous state history matrices which can further be reduced via Eq. (2.4) into a product of sub-domain integrals over each matrix row

$$P_{t+1}(x|z) = \frac{1}{t} \sum_{t''=0}^t \left\{ \int_{\omega_{t'}} d^n x' \dots \int_{\Omega_{t''}} dX'' \right\} P_{t''}(X''|z) P_{(t+1)t \dots t''}(x, x', \dots | X'', z) \quad (2.6)$$

$$= \frac{1}{t} \sum_{t''=0}^t \int_{\Omega_{t''}} dX'' P_{t''}(X''|z) P_{(t+1)t''}(x|X'', z), \quad (2.7)$$

where each row measure is a Cartesian product of  $n$  elements (a Lebesgue measure), i.e.,

$$d^n x = \prod_{i=0}^n dx^i, \quad (2.8)$$

and lowercase  $x, x', \dots$  values will always refer to individual rows within the state matrices. Note that  $1/t$  here is a normalisation factor — this just normalises the sum of all probabilities to 1 given that there is a sum over  $t'$ . Note also that, if the process is defined over continuous time, we would need to replace

$$\frac{1}{t} \sum_{t'=0}^t \rightarrow \frac{1}{t(t)} \sum_{t'=0}^t \delta t(t'). \quad (2.9)$$

To try and understand what Eqs. (2.1) and (2.5) are saying, we find it's helpful to think of an iterative relationship between probabilities; each of which is connected by their relative conditional probabilities. This kind of thinking is also illustrated in Fig. 2.1.

We can approximate the probability  $P_{t+1}(X|z)$  with a logarithmic expansion like this

$$\ln P_{t+1}(X|z) \simeq \ln P_{t+1}(X_*|z) + \frac{1}{2} \sum_{t'=0}^{t+1} \sum_{i=0}^n \sum_{j=0}^n (x - x_*)^i \mathcal{I}_{(t+1)t'}^{ij}(x_*, x'_*)(x' - x'_*)^j \quad (2.10)$$

$$\mathcal{I}_{(t+1)t'}^{ij}(x_*, x'_*) = \frac{\partial}{\partial x^i} \frac{\partial}{\partial (x')^j} \ln P_{t+1}(X|z) \Big|_{X=X_*}, \quad (2.11)$$

where the values for  $X = X_*$  are defined by the vanishing of the first derivative, i.e., these are chosen such that

$$\frac{\partial}{\partial x^i} \ln P_{t+1}(X|z) \Big|_{X=X_*} = 0. \quad (2.12)$$

If we keep the truncation up to second order in Eq. (2.10), note that this expression implies a pairwise correlation structure of the form

$$P_{t+1}(X|z) \rightarrow \prod_{t'=0}^t P_{(t+1)t'}(x, x'|z) = \prod_{t'=0}^t P_{t'}(x'|z) P_{(t+1)t'}(x|x', z). \quad (2.13)$$

Given this pairwise temporal correlation structure, Eq. (2.7) reduces to this simpler sum of integrals

$$P_{t+1}(x|z) = \frac{1}{t} \sum_{t'=0}^t \int_{\omega_{t'}} d^n x' P_{t'}(x'|z) P_{(t+1)t'}(x|x', z). \quad (2.14)$$

In a similar fashion, we can increase the expansion order of Eq. (2.10) to include third-order correlations such that

$$P_{t+1}(X|z) \rightarrow \prod_{t'=0}^t \prod_{t''=0}^{t'-1} P_{t't''}(x', x''|z) P_{(t+1)t't''}(x|x', x'', z), \quad (2.15)$$

and, in this instance, one can show that Eq. (2.7) reduces to

$$P_{t+1}(x|z) = \frac{1}{t} \sum_{t'=0}^t \frac{1}{t'-1} \sum_{t''=0}^{t'-1} \int_{\omega_{t'}} d^n x' \int_{\omega_{t''}} d^n x'' P_{t't''}(x', x''|z) P_{(t+1)t't''}(x|x', x'', z). \quad (2.16)$$

Using  $P_{t't''}(x', x''|z) = P_{t''}(x''|z) P_{t't''}(x'|x'', z)$  one can also show that this integral is a marginalisation of this expression

$$P_{(t+1)t''}(x|x'', z) = \frac{1}{t} \sum_{t'=0}^t \int_{\omega_{t'}} d^n x' P_{t't''}(x'|x'', z) P_{(t+1)t't''}(x|x', x'', z), \quad (2.17)$$

which describes the time evolution of the conditional probabilities.

Let's imagine that  $x$  is just a scalar (as opposed to a row vector) for simplicity in the expressions. We can then discretise the 1D space over  $x$  into separate  $i$ -labelled regions such that  $[P]_{t+1}^i - [P]_t^i = \mathcal{J}_{t+1}^i$ , where the probability current  $\mathcal{J}_{t+1}^i$  for the factorised processes above would be defined as

$$\mathcal{J}_{t+1}^i = -[P]_t^i + \frac{1}{t} \sum_{t'=0}^t \sum_{i'=0}^N \Delta x [P]_{t'}^{i'} [P]_{(t+1)t'}^{ii'} \quad (2.18)$$

$$\mathcal{J}_{t+1}^i = -[P]_t^i + \frac{1}{t} \sum_{t'=1}^t \frac{1}{t'-1} \sum_{t''=0}^{t'-1} \sum_{i'=0}^N \sum_{i''=0}^N \Delta x^2 [P]_{t''}^{i''} [P]_{t't''}^{i'i''} [P]_{(t+1)t't''}^{ii'i''}. \quad (2.19)$$

The  $[P]_{(t+1)t't''}^{ii'i''}$  tensor, in particular, will have  $N^3 t(t^2 - 1)$  elements. Note that the third-order temporal correlations can be evolved by identifying the pairwise conditional probabilities as time-dependent state variables and evolving them according to the following relation

$$[P]_{(t+1)t't''}^{ii''} = \frac{1}{t} \sum_{t'=1}^t \sum_{i'=0}^N \Delta x [P]_{t't''}^{i'i''} [P]_{(t+1)t't''}^{ii'i''}. \quad (2.20)$$

What other classes of process can be described by Eqs. (2.1) and (2.5)? For Markovian phenomena, the equations no longer depend on timesteps older than the immediately previous one, hence Eq. (2.5) reduces to just

$$P_{t+1}(x|z) = \int_{\omega_t} d^n x' P_t(x'|z) P_{(t+1)t}(x|x', z). \quad (2.21)$$

By performing a Kramers-Moyal expansion on the  $P_{(t+1)t}(x|x', z)$  distribution up to second order,<sup>1</sup> we can approximate the right hand side of Eq. (2.21) like this

$$P_{(t+1)}(x|z) \simeq P_t(x|z) - \sum_{i=0}^n \frac{\partial}{\partial x^i} \left[ f_t(x, z) P_t(x|z) \right] + \frac{1}{2} \sum_{i=0}^n \sum_{j=0}^n \frac{\partial}{\partial x^i} \frac{\partial}{\partial x^j} \left[ K_t(x, z) P_t(x|z) \right], \quad (2.22)$$

where the components of  $f_t(x, z)$  and  $K_t(x, z)$  have been defined as

$$f_t^i(x, z) = \int_{\omega_{t+1}} d^n x' (x' - x)^i P_{(t+1)t}(x'|x, z) \quad (2.23)$$

$$K_t^{ij}(x, z) = \int_{\omega_{t+1}} d^n x' (x' - x)^i (x' - x)^j P_{(t+1)t}(x'|x, z). \quad (2.24)$$

By inspection of Eq. (2.22) we can define the ‘probability current’

$$J_t(x, z) = f_t(x, z) P_t(x|z) - \frac{1}{2} \sum_{j=0}^n \frac{\partial}{\partial x^j} \left[ K_t(x, z) P_t(x|z) \right]. \quad (2.25)$$

If the probability current vanishes  $J_t(x, z) = 0$  individually (this also implies that the distribution is stationary such that  $P_{t+1}(x|z) = P_t(x|z)$ ), the implicit solution of Eq. (2.22) can be found to be

$$P_t(x|z) \propto K_t^{-1}(x, z) \exp \left[ \int d^n x K_t^{-1}(x, z) f_t(x, z) \right]. \quad (2.26)$$

An analog of Eq. (2.5) exists for discrete state spaces as well. We just need to replace the integral with a sum and the schematic would look something like this

$$P_{t+1}(x|z) = \sum_{\Omega_t} P_t(X'|z) P_{(t+1)t}(x|X', z), \quad (2.27)$$

where we note that the  $P$ ’s in the expression above all now refer to *probability mass functions*.

- Add in mean and covariance calculation of the process to generalise it continuous data and dev the Learnadex code using Gotch to support this with different data-linking distributions.
- Add some diagrams for the higher-order correlation expressions.
- Add a software design section and some examples.

---

<sup>1</sup>The Pawula theorem [1] states that this is, in fact, as far as we can truncate up to unless we include an infinite number of expansion terms.



# Bibliography

- [1] R. Pawula, *Generalizations and extensions of the fokker-planck-kolmogorov equations*, *IEEE Transactions on Information Theory* **13** (1967) 33–41.