

Probabilistic learning algorithms

Rewrite from here...

Concept. To extend the formalism that we developed in the previous chapter to help us better understand the time evolution of state probabilities. With this knowledge we shall derive popular machine learning methods like Gaussian processes, as well as discussing how some phenomena permit a convenient description via ‘mean-field’ equations or regression of distribution moments. 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, with a view to statistical inference later on. For the programmers, all of the relevant software lives in this public Git repository: <https://github.com/umbralcalc/learnadex>.

2.1 Probabilistic formalism

Let’s start by returning to the formalism that we introduced in the previous chapter. 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 $X_{t+1}^i = F_{t+1}^i(X_{0:t}, z, \mathbf{t})$. 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 . To try and understand what Eq. (2.1) is saying, we find it’s helpful to think of an iterative relationship between probabilities; each of which is connected by their relative conditional probabilities. We’ve also illustrated this kind of thinking in Fig. 2.1.

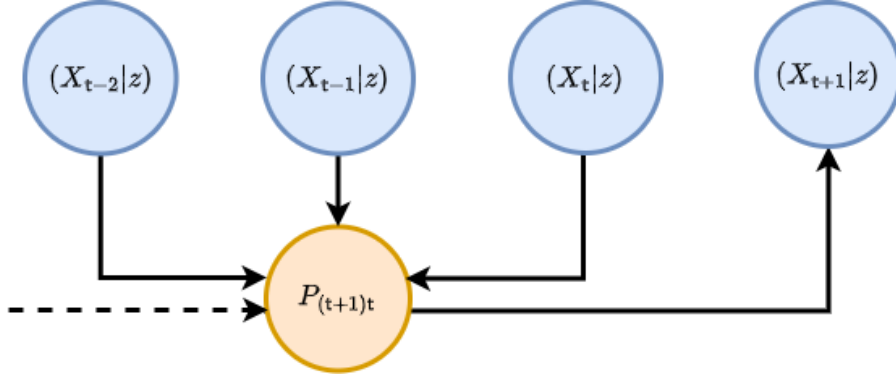


Figure 2.1: Graph representation of Eq. (2.1).

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 Ω_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

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

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

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

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

Let's go through some examples. Non-Markovian phenomena with continuous state spaces can have quite complex master equations. A relatively simple example is that of pure diffusion processes which exhibit stochastic resetting at a rate r to a remembered location from the trajectory history [1]

$$P_{t+1}(x|z) = (1-r)P_t(x|z) + \sum_{i=0}^n \sum_{j=0}^n \frac{\partial}{\partial x^i} \frac{\partial}{\partial x^j} \left[D_t(x, z) P_t(x|z) \right] + r \sum_{t'=0}^t \delta t(t') K[t(t)-t(t')] P_{t'}(x|z), \quad (2.10)$$

where here K is some memory kernel. For Markovian phenomena which have a continuous state space, Eqs. (2.1) and (2.5) no longer depend on timesteps older than the immediately previous one, hence, e.g., Eq. (2.5) reduces to just

$$P_{t+1}(x|z) = \int_{\omega_t} d^n x' P_t(x'|z) P_{(t+1)t}(x|x', z). \quad (2.11)$$

A famous example of this kind of phenomenon arises from approximating Eq. (2.11) with an expansion (Kramers-Moyal [2, 3]) up to second-order, yielding the Fokker-Planck equation

$$P_{t+1}(x|z) = P_t(x|z) - \sum_{i=0}^n \frac{\partial}{\partial x^i} \left[\mu_t(x, z) P_t(x|z) \right] + \sum_{i=0}^n \sum_{j=0}^n \frac{\partial}{\partial x^i} \frac{\partial}{\partial x^j} \left[D_t(x, z) P_t(x|z) \right], \quad (2.12)$$

which describes a process undergoing drift-diffusion.

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

$$P_{t+1}(x|z) = \sum_{\Omega_t} P_t(X'|z) P_{(t+1)t}(x|X', z), \quad (2.13)$$

where we note that the P 's in the expression above all now refer to *probability mass functions*. In what follows, discrete state space can always be considered by replacing the integrals with summations over probability masses in this manner; we only use the continuous state space formulation for our notation because one could argue it's a little more general.

Analogously to continuous state spaces, we can give some examples of master equations for phenomena with a discrete state space as well. In the Markovian case, we need look no further than a simple time-dependent Poisson process

$$P_{t+1}(x|z) = \lambda(t) \delta t(t+1) P_t(x-1|z) + [1 - \lambda(t) \delta t(t+1)] P_t(x|z). \quad (2.14)$$

For such an example of a non-Markovian system, a Hawkes process [4] master equation would look something like this

$$\begin{aligned}
P_{\mathbf{t}+1}(x|z) &= \mu \delta t(\mathbf{t}+1) P_{\mathbf{t}}(x-1|z) + [1 - \mu \delta t(\mathbf{t}+1)] P_{\mathbf{t}}(x|z) \\
&+ \sum_{x'=0}^{\infty} \sum_{x''=0}^{\infty} \sum_{\mathbf{t}'=0}^{\mathbf{t}} \phi[t(\mathbf{t}) - t(\mathbf{t}')] (x' - x'') \delta t(\mathbf{t}+1) P_{\mathbf{t}\mathbf{t}'(\mathbf{t}'-1)}(x-1, x', x''|z) \\
&+ \sum_{x'=0}^{\infty} \sum_{x''=0}^{\infty} \left\{ 1 - \sum_{\mathbf{t}'=0}^{\mathbf{t}} \phi[t(\mathbf{t}) - t(\mathbf{t}')] (x' - x'') \delta t(\mathbf{t}+1) \right\} P_{\mathbf{t}\mathbf{t}'(\mathbf{t}'-1)}(x, x', x''|z), \quad (2.15)
\end{aligned}$$

where we note the complexity in this expression arises because it has to include a coupling between the rate at which events occur and an explicit memory of when the previous ones did occur (recorded by differencing the count $x' - x''$).

2.2 Derived learning algorithms

So now that we have become familiar with the notation used by Eq. (2.5), let's use it to derive a few important statistical insights. We can begin by pointing out that by integrating over x , we compute the mean of the distribution

$$M_{\mathbf{t}+1}(z) = \int_{\omega_{\mathbf{t}+1}} d^n x x P_{\mathbf{t}+1}(x|z) = \frac{1}{\mathbf{t}} \sum_{\mathbf{t}''=0}^{\mathbf{t}} \int_{\Omega_{\mathbf{t}''}} dX'' P_{\mathbf{t}''}(X''|z) M_{(\mathbf{t}+1)\mathbf{t}''}(X'', z), \quad (2.16)$$

where you can view the $M_{(\mathbf{t}+1)\mathbf{t}''}(X'', z)$ values as either terms in some regression model, or derivable explicitly from a known master equation. The latter of these provides one approach to statistically infer the states and parameters of stochastic simulations from data: one begins by knowing what the master equation is, uses this to compute the time evolution of the mean (and potentially higher-order statistics) and then connects these \mathbf{t} and z -dependent statistics back to the likelihood of observing the data. This is what is commonly known as the ‘mean-field’ inference approach; averaging over the available degrees of freedom in the statistical moments of distributions. Knowing what precisely the master equation for an arbitrarily-defined stochastic phenomenon can be quite a tall order to fill, however.

Given that the mean-field approach isn't always going to be viable as an inference method, we should also consider other ways to describe the shape and time evolution characteristics of $P_{\mathbf{t}+1}(X|z)$. For continuous state spaces, it's possible to approximate this whole distribution with a logarithmic expansion like so

$$\ln P_{\mathbf{t}+1}(X|z) \simeq \ln P_{\mathbf{t}+1}(X_*|z) + \frac{1}{2} \sum_{\mathbf{t}'=0}^{\mathbf{t}+1} \sum_{i=0}^n \sum_{j=0}^n (x - x_*)^i \mathcal{H}_{(\mathbf{t}+1)\mathbf{t}'}^{ij}(x_*, x'_*)(x' - x'_*)^j \quad (2.17)$$

$$\mathcal{H}_{(\mathbf{t}+1)\mathbf{t}'}^{ij}(x_*, x'_*) = \frac{\partial}{\partial x^i} \frac{\partial}{\partial (x')^j} \ln P_{\mathbf{t}+1}(X|z) \Big|_{X=X_*}, \quad (2.18)$$

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

$$\frac{\partial}{\partial x^i} \ln P_{\mathbf{t}+1}(X|z) \Big|_{X=X_*} = 0. \quad (2.19)$$

If we keep the truncation up to second order in Eq. (2.17), note that this expression implies a pairwise correlation structure of the form

$$P_{t+1}(X|z) \rightarrow \prod_{t'=0}^t P_{(t+1)t'}(x, x'|z) = \prod_{t'=0}^t P_{t'}(x'|z) P_{(t+1)t'}(x|x', z). \quad (2.20)$$

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

$$P_{t+1}(x|z) = \frac{1}{t} \sum_{t'=0}^t \int_{\omega_{t'}} d^n x' P_{t'}(x'|z) P_{(t+1)t'}(x|x', z). \quad (2.21)$$

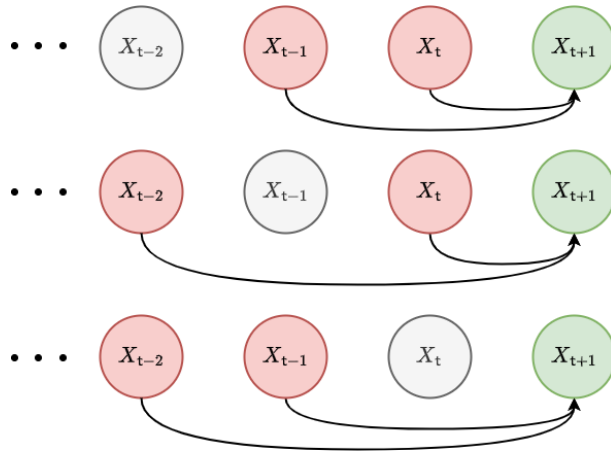


Figure 2.2: A graph representation of the correlations in Eq. (2.21).

In a similar fashion, we can increase the expansion order of Eq. (2.17) to include third-order correlations such that

$$P_{t+1}(X|z) \rightarrow \prod_{t'=0}^t \prod_{t''=0}^{t'-1} P_{t't''}(x', x''|z) P_{(t+1)t't''}(x|x', x'', z), \quad (2.22)$$

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

$$P_{t+1}(x|z) = \frac{1}{t} \sum_{t'=0}^t \frac{1}{t'-1} \sum_{t''=0}^{t'-1} \int_{\omega_{t'}} d^n x' \int_{\omega_{t''}} d^n x'' P_{t't''}(x', x''|z) P_{(t+1)t't''}(x|x', x'', z). \quad (2.23)$$

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

$$P_{(t+1)t't''}(x|x'', z) = \frac{1}{t} \sum_{t'=0}^t \int_{\omega_{t'}} d^n x' P_{t't''}(x'|x'', z) P_{(t+1)t't''}(x|x', x'', z), \quad (2.24)$$

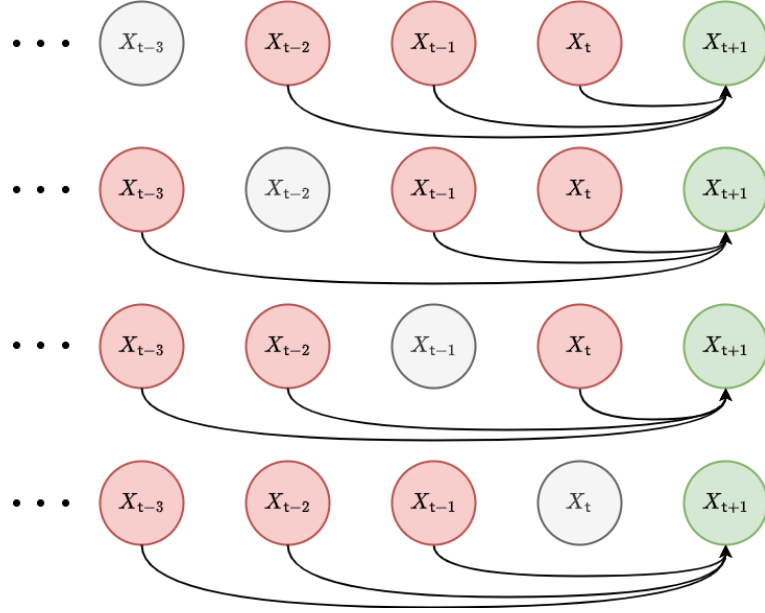


Figure 2.3: A graph representation of the correlations in Eq. (2.23).

which describes the time evolution of the conditional probabilities.

The key distinction between the methods that we will develop in this chapter and the ones in the proceeding chapters is in their utility when faced with the problem of attempting to model real-world data. In the proceeding chapter, we shall describe some powerful techniques that can be used most effectively when the researcher is aware of the family of models that generated the data. In the present chapter, we will go into the details of how a more ‘empirical’ approach can be derived for dynamical process modeling in a probabilistic framework which locally adapts the model to the data through time.

While we think that it’s worth going into some mathematical detail to give a better sense of where our formalism comes from; we want to emphasise that the framework we discuss here is not new to the technical literature at all. Our overall framework draws on influences from Empirical Dynamical Modeling (EDM) [5], some classic nonparametric local regression techniques — such as LOWESS/Savitzky-Golay filtering [6] — and also Gaussian processes [7] as well. The novelties here, instead, lie more in the specifics of how we combine some of these ideas together when referencing the stochadex formalism, and how this manifests in designing more generally-applicable software for the user. In addition, readers well-versed in machine learning will note that our software design and mathematical formalism reflects our preference for ‘online’ learning¹ in the context of time series prediction, in contrast to some of the more standard frameworks.

When trying robustly assess how far a model is from accurately describing a set of real-world data, trying to use only generated samples of the model process can be difficult. Instead, in this section, we are going to extend this formalism to look at how probability theory can help with this data comparison problem in a systematic way. In order to do this, we need to return to the

¹We’ll explain this later in the chapter.

probabilistic formalism which we discussed in the previous chapter.

Given the general master equation $P_{t+1}(X|z) = P_t(X'|z)P_{(t+1)t}(x|X', z)$, if we wanted to compute the mean $M_{t+1}(z)$ of the distribution over the matrix row corresponding to time $(t + 1)$, it would be straightforward to just multiply both sides by x and integrate over it in its ω_{t+1} sub-domain. However, there is another similar expression for the mean that we can derive under certain conditions which will be valuable to us when developing the empirical reweighting. If the probability distribution over each row of the state history matrix is *stationary* — meaning that $P_{t+1}(x|z) = P_t(x|z)$ — it's possible to derive

$$M_{t+1}(z) = \int_{\omega_{t+1}} d^n x x P_{t+1}(x|z) = \frac{1}{t} \sum_{t'=0}^t \int_{\omega_{t'}} d^n x' x' P_{t'}(x'|z) \int_{\omega_{t+1}} d^n x P_{(t+1)t'}(x|x', z). \quad (2.25)$$

To see that Eq. (2.25) is true, first note that a joint distribution over both x and x' can be derived like this $P_{(t+1)t'}(x, x'|z) = P_{(t+1)t'}(x|x', z)P_{t'}(x'|z)$. Secondly, note that this joint distribution will always allow variable swaps trivially like this $P_{(t+1)t'}(x, x'|z) = P_{t'}(x'|z)P_{(t+1)t}(x|x', z)$. Then, lastly, note that stationarity of $P_{t+1}(x|z) = P_t(x|z)$ means

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

where we've used the trivial variable swap and integration variable relabelling to arrive at the second equality in the expressions above.

The standard covariance matrix elements can also be computed in a similar fashion

$$\begin{aligned} C_{t+1}^{ij}(z) &= \int_{\omega_{t+1}} d^n x [x - M_{t+1}(z)]^i [x - M_{t+1}(z)]^j P_{t+1}(x|z) \\ &= \frac{1}{t} \sum_{t'=0}^t \int_{\omega_{t'}} d^n x' [x' - M_{t+1}(z)]^i [x' - M_{t+1}(z)]^j P_{t'}(x'|z) \int_{\omega_{t+1}} d^n x P_{(t+1)t'}(x|x', z). \end{aligned} \quad (2.26)$$

While they look quite abstract, Eqs. (2.25) and (2.26) express the core idea behind how our probabilistic reweighting will function. By assuming a stationary distribution, we gain the ability to directly estimate the statistics of the probability distribution of the next sample from the stochastic process $P_{t+1}(x|z)$ from past samples it may have in empirical data; which are represented here by $P_{t'}(x'|z)$. More on this later.

It may seem needlessly more complex to deal with probability distributions over matrices instead of marginal distributions over separate vectors which represent the rows of these matrices. However, the matrix description is more general, and it turns out to be quite neat to describe correlations across time (which would be lost by marginal distributions over row vectors). In order to study these out-of-time-order correlations, we need only consider using the statistical moments computed across

the columns of X . For example, pairwise correlations can be analysed through the out-of-time-order covariance matrix elements

$$C_{(t+1)t'}^{ij}(z) = \int_{\omega_{t+1}} d^n x \int_{\omega_{t'}} d^n x' [x - M_{t+1}(z)]^i [x' - M_{t'}(z)]^j P_{t+1}(X|z). \quad (2.27)$$

Probabilistic reweighting depends on the stationarity of $P_{t+1}(x|z) = P_{t'}(x|z)$ such that, e.g., Eq. (2.25) is applicable. The core idea behind it is to represent the past distribution of state values $P_{t'}(x'|z)$ with the samples from a real time series dataset. If the user then specifies a good model for the relationships in this data by providing a weighting function which returns the conditional probability mass

$$\mathbf{w}_{t'}(y, z) = \int_{\omega_{t+1}} d^n x P_{(t+1)t'}(x|x'=y, z), \quad (2.28)$$

we can apply this as a *reweighting* of the historical time series samples to estimate any statistics of interest. Taking Eqs. (2.25) and (2.26) as the examples; we are essentially approximating these integrals through weighted sample estimations like this

$$M_{t+1}(z) \simeq \frac{1}{t} \sum_{t'=0}^t Y_{t'} \mathbf{w}_{t'}(Y_{t'}, z) \quad (2.29)$$

$$C_{t+1}^{ij}(z) \simeq \frac{1}{t} \sum_{t'=0}^t [Y_{t'} - M_{t+1}(z)]^i [Y_{t'} - M_{t+1}(z)]^j \mathbf{w}_{t'}(Y_{t'}, z), \quad (2.30)$$

where we have defined the data matrix Y with rows Y_{t+1}, Y_t, \dots , each of which representing specific observations of the rows in X at each point in time from a real dataset.

Our goal in this section will be to learn the optimal reweighting function $\mathbf{w}_{t'}(Y_{t'}, z)$ with respect to z , i.e., the ones which most accurately represent a provided dataset. But before we think about the various kinds of conditional probability we could use, we need to think about how to connect the post-reweighting statistics to the data by defining an objective function.

If the mean is a sufficient statistic for the distribution which describes the data, a choice of, e.g., Exponential, Poisson or Binomial distribution could be used where the mean is estimated directly from the time series using Eq. (2.25), given a conditional probability $P_{(t+1)t'}(x|x', z)$. Extending this idea further to include distributions which also require a variance to be known, e.g., the Normal, Gamma or Negative Binomial distributions could be used where the variance (and/or covariance) could be estimated using Eq. (2.26). These are just a few simple examples of distributions that can link the estimated statistics from Eqs. (2.25) and (2.26) to a time series dataset. However, the algorithmic framework is very general to whatever choice of ‘data linking’ distribution that a researcher might need.

We should probably make what we’ve just said a little more mathematically concrete. We can define $P_{t+1}[y; M_{t+1}(z), C_{t+1}(z), \dots]$ as representing the likelihood of $y = Y_{t+1}$ given the estimated statistics from Eqs. (2.25) and (2.26) (and maybe higher-orders). Note that in order to do this, we need to identify the x' and t' values that are used to estimate, e.g., $M_{t+1}(z)$ with the past data values which are observed in the dataset time series itself. Now that we have this likelihood, we can immediately evaluate an objective function (a cumulative log-likelihood) that we might seek to

optimise over for a given dataset

$$\ln \mathcal{L}_{t+1}(Y|z) = \sum_{t'=0}^{t+1} \ln P_{t'}[y; M_{t'}(z), C_{t'}(z), \dots], \quad (2.31)$$

where the summation continues until all of the past measurements Y_{t+1}, Y_t, \dots which exist as rows in the data matrix Y have been taken into account. The code to compute this objective function follows the schematic we have provided in Fig. 2.4.

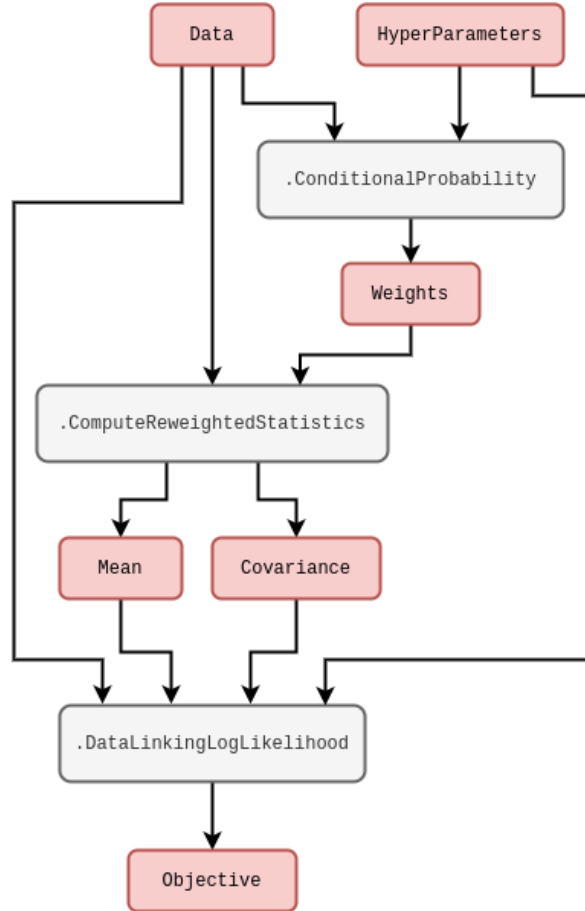


Figure 2.4: Code schematic of the probability reweighting optimisation.

In order to specify what $P_{(t+1)t'}(x|x', z)$ is, it's quite natural to define a set of hyperparameters for the elements of z . To get a sense of how the data-linking function relates to these hyperparameters, it's instructive to consider an example. One generally-applicable option for the conditional probability could be a purely time-dependent kernel

$$P_{(t+1)t'}(x|x', z) \propto \mathcal{K}(z, t+1, t'), \quad (2.32)$$

and the data-linking distribution, e.g., could be a Gaussian

$$P_{t+1}[y; M_{t+1}(z), C_{t+1}(z), \dots] = \text{MultivariateNormalPDF}[y; M_{t+1}(z), C_{t+1}(z)]. \quad (2.33)$$

It's worth pointing out that other machine learning frameworks could easily be used to model these conditional probabilities. For example, neural networks could be used to infer the optimal reweighting scheme and this would still allow us to use the data-linking distribution.² It would still be desirable to keep the data-linking distribution as it can usually be sampled from very easily — something that can be quite difficult to achieve with a purely machine learning-based representation of the distribution. Sampling itself could even be made more flexible by leveraging a Variational Autoencoder (VAE) [9]; these use neural networks not just on the compression (or ‘encode’) step to estimate the statistics but also use them as a layer between the sample from the data distribution model and the output (the ‘decode’ step).

In the case of Eqs. (2.32) and (2.33) above, the hyperparameters that would be optimised could relate to the kernel in a wide variety of ways. Optimising them would make our optimised reweighting very similar to (but not quite the same as) evaluating maximum a posteriori (MAP) of a Gaussian process regression. The main differences here are that the mean of a Gaussian process as a function of time is typically included within z , and hence must be obtained through optimisation. In contrast, our methodology relies on the fact that the mean estimator can be computed directly by weighted sample estimation and can then be fed to some data-linking distribution. By doing it this way, we enable many different kinds of data to be described by the same underlying conditional probability-based reweighting and also make incorporating future data into the current model much easier than is the case for a standard Gaussian process (which would require re-optimising with respect to the next data point every time). The time-dependent kernel we have chosen in Eq. (2.32) also only represents one particular choice, but we could consider a wide range of state-dependent conditional probability weightings for the algorithm as well.

As another form of flexibility; we could also try adapting the data-linking distributions to include an intercept term and linear coefficients for the statistics which are passed to it. These could then be treated as additional hyperparameters and optimised jointly with the others if there is sufficient constraining power in the data.

The optimisation approach that we choose to use for obtaining the best hyperparameters in the conditional probability of Eq. (2.31) will depend on a few factors. For example, if the number of hyperparameters is relatively low, but their gradients are difficult to calculate exactly; then a gradient-free optimiser (such as the Nelder-Mead [10] method or something like a particle swarm [11, 12]) would likely be the most effective choice. On the other hand, when the number of hyperparameters ends up being relatively large, it's usually quite desirable to utilise the gradients in algorithms like vanilla Stochastic Gradient Descent [13] (SGD) or Adam [14].

If the gradients of Eq. (2.31) are needed, we can always factorise each derivative with respect

²One can think of using this neural network-based reweighting scheme as similar to constructing a normalising flow model [8] with an autoregressive layer. Invertibility and further network structural constraints mean that these are not exactly equivalent, however.

to hyperparameter z^i in the following way through the chain rule

$$\begin{aligned} \frac{\partial}{\partial z^i} \ln \mathcal{L}_{t+1}(Y|z) &= \sum_{t'=0}^{t+1} \frac{\partial M_{t'}}{\partial z^i} \frac{\partial}{\partial M_{t'}} \ln P_{t'}[y; M_{t'}(z), C_{t'}(z), \dots] \\ &\quad + \sum_{t'=0}^{t+1} \frac{\partial C_{t'}}{\partial z^i} \frac{\partial}{\partial C_{t'}} \ln P_{t'}[y; M_{t'}(z), C_{t'}(z), \dots]. \end{aligned} \quad (2.34)$$

By factoring derivatives in this manner, the computation can be separated into two parts: the derivatives with respect to $M_{t'}$ and $C_{t'}$, which are typically quite straightforward; and the derivatives with respect to z elements, which typically need a more involved calculation depending on the model. Incidentally, this separation also neatly lends itself to abstracting gradient calculations as having a simpler, general purpose component that can be built directly into a library of data models and a more complex, model-specific component that the user must specify.

Before moving on to the software design aspects, we need to consider how we might structure learning by optimisation of Eq. (2.31) for a sequence of observations in time. One of the issues that can arise when learning streams of data is ‘concept drift’. In our context, this would be when the optimal value for z does not match the optimal value at some later point in time. In order to solve this issue, our learning algorithm should track an up-to-date optimal value for z as data is continually passed into it. Iteratively updating the optimal parameters as new data is ingested into the objective function is typically called ‘online learning’ [15, 16], in contrast to ‘offline learning’ which would correspond to learning an optimal z only once with the entire dataset provided upfront.³

³This book will mostly be focussed on using online learning techniques due to the inherently sequential nature of stochastic processes.

Bibliography

- [1] D. Boyer, M. R. Evans, and S. N. Majumdar, “Long time scaling behaviour for diffusion with resetting and memory,” *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2017, no. 2, p. 023208, 2017.
- [2] H. A. Kramers, “Brownian motion in a field of force and the diffusion model of chemical reactions,” *Physica*, vol. 7, no. 4, pp. 284–304, 1940.
- [3] J. Moyal, “Stochastic processes and statistical physics,” *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 11, no. 2, pp. 150–210, 1949.
- [4] A. G. Hawkes, “Spectra of some self-exciting and mutually exciting point processes,” *Biometrika*, vol. 58, no. 1, pp. 83–90, 1971.
- [5] G. Sugihara and R. M. May, “Nonlinear forecasting as a way of distinguishing chaos from measurement error in time series,” *Nature*, vol. 344, no. 6268, pp. 734–741, 1990.
- [6] A. Savitzky and M. J. Golay, “Smoothing and differentiation of data by simplified least squares procedures,” *Analytical chemistry*, vol. 36, no. 8, pp. 1627–1639, 1964.
- [7] K. P. Murphy, *Machine learning: a probabilistic perspective*. MIT press, 2012.
- [8] I. Kobyzev, S. J. Prince, and M. A. Brubaker, “Normalizing flows: An introduction and review of current methods,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 11, pp. 3964–3979, 2020.
- [9] L. Pinheiro Cinelli, M. Araújo Marins, E. A. Barros da Silva, and S. Lima Netto, “Variational autoencoder,” in *Variational Methods for Machine Learning with Applications to Deep Networks*. Springer, 2021, pp. 111–149.
- [10] J. A. Nelder and R. Mead, “A simplex method for function minimization,” *The computer journal*, vol. 7, no. 4, pp. 308–313, 1965.
- [11] J. Kennedy and R. Eberhart, “Particle swarm optimization,” in *Proceedings of ICNN’95-international conference on neural networks*, vol. 4. IEEE, 1995, pp. 1942–1948.
- [12] Y. Shi and R. Eberhart, “A modified particle swarm optimizer,” in *1998 IEEE international conference on evolutionary computation proceedings. IEEE world congress on computational intelligence (Cat. No. 98TH8360)*. IEEE, 1998, pp. 69–73.

- [13] H. Robbins and S. Monro, “A stochastic approximation method,” *The annals of mathematical statistics*, pp. 400–407, 1951.
- [14] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [15] E. Hazan *et al.*, “Introduction to online convex optimization,” *Foundations and Trends® in Optimization*, vol. 2, no. 3-4, pp. 157–325, 2016.
- [16] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.