

Variational inference for Markov jump processes

Guido Sanguinetti

Joint work with Manfred Opper (T.U. Berlin)



The
University
Of
Sheffield.

Talk plan

- Markov jump processes in the life sciences
- Basics of jump processes
- Variational mean field
- Results
- Future extensions?

Markov Jump Processes

- Describe the dynamics of populations of interacting species
- Changes in populations are discrete (birth/death) and happen at random times
- Rates of births and deaths for a species depend on population levels for all species
- Many examples: chemical kinetics, predator-prey models, telecoms, etc.

Relevance to life sciences

- Environmental systems
- Single-cell dynamics
- Low concentration pathogens (see F. Bois talk)
- Usually dynamics is simulated using (variants of) Gillespie's algorithm (SSA)
- SSA does not include observations and requires full specification of system
- Inference is hard: one MCMC approach (Boys et al; "A bloody mess", A. Golightly, PESB2007)

Mathematical notation

- Mathematically, a Markov Jump Process (MJP) is a family of discrete random variables indexed by time $\mathbf{x}(t)$.

- Markov property:

$$p(\mathbf{x}(t_N) | \mathbf{x}(t_{N-1}), \dots, \mathbf{x}(t_0)) = p(\mathbf{x}(t_N) | \mathbf{x}(t_{N-1})).$$

- Process rates $f(\mathbf{x}'|\mathbf{x})$ defined by

$$\lim_{\delta t \rightarrow 0} p(\mathbf{x}'(t + \delta t) | \mathbf{x}(t)) = \delta_{\mathbf{x}'\mathbf{x}} + f(\mathbf{x}'|\mathbf{x}).$$

- Process rates satisfy $\sum f(\mathbf{x}'|\mathbf{x})=0$ (normalisation)

Master equation

- The process rates and the time-dependent probabilities $p(\mathbf{x})$ are linked by the *Master equation*

$$\frac{dp_t(\mathbf{x})}{dt} = \sum_{\mathbf{x}' \neq \mathbf{x}} \left[-p_t(\mathbf{x}) f(\mathbf{x}'|\mathbf{x}) + p_t(\mathbf{x}') f(\mathbf{x}|\mathbf{x}') \right].$$

- Notice that the Master equation is in fact a system of S^D ODEs, with S the number of states accessible to each species, and D the number of species

Variational Inference

- Deterministic approximate inference technique to approximate a probability distribution p with a simpler one q .
- The measure of proximity between distributions is the Kullback-Leibler (KL) divergence

$$KL(q||p) = - \int dq \log \frac{p}{q}$$

- Notice that you do not need explicit knowledge of p , just of the expectation of $\log(p)$ under q

KL divergence between MJPs

- Consider two MJPs p and q , with rates f and g .
- A MJP can be viewed as a probability distribution over trajectories of the system; denote a trajectory as $\mathbf{x}_{0:K} = (\mathbf{x}(t_0), \dots, \mathbf{x}(t_0 + K\delta t))$.
- The KL divergence between these processes is

$$\begin{aligned}
 KL(q||p) &= \sum_{\mathbf{x}_{0:K}} q(\mathbf{x}_{0:K}) \ln \frac{q(\mathbf{x}_{0:K})}{p(\mathbf{x}_{0:K})} = \\
 &= \sum_{k=0}^{K-1} \sum_{\mathbf{x}_k} q(\mathbf{x}_k) \sum_{\mathbf{x}_{k+1}} q(\mathbf{x}_{k+1}|\mathbf{x}_k) \ln \frac{q(\mathbf{x}_{k+1}|\mathbf{x}_k)}{p(\mathbf{x}_{k+1}|\mathbf{x}_k)} + K_0
 \end{aligned}$$

Continuous limit

- By taking the limit $\delta t \rightarrow 0$, we can rewrite the KL divergence in terms of process rates.
- Setting the initial KL $K_0=0$, we get

$$KL(q||p) = \int_0^T dt \sum_{\mathbf{x}} q_t(\mathbf{x}) \times \sum_{\mathbf{x}': \mathbf{x}' \neq \mathbf{x}} \left\{ g(\mathbf{x}'|\mathbf{x}) \ln \frac{g(\mathbf{x}'|\mathbf{x})}{f(\mathbf{x}'|\mathbf{x})} + f(\mathbf{x}'|\mathbf{x}) - g(\mathbf{x}'|\mathbf{x}) \right\}$$

Posterior processes

- We assume now to have discrete-time, noise-corrupted observations \mathbf{y}_l of the process.
- The noise model is given by $\hat{p}(\mathbf{y}_l | \mathbf{x}(t_l))$.
- The posterior process is still Markovian and is

$$p_{post}(\mathbf{x}_{0:K}) = \frac{1}{Z} p_{prior}(\mathbf{x}_{0:K}) \times \prod_{l=1}^N \hat{p}(\mathbf{y}_l | \mathbf{x}(t_l)).$$

- The KL divergence with a process q is

$$KL(q || p_{post}) = \ln Z + KL(q || p_{prior}) - \sum_{l=1}^N E_q [\ln p(\mathbf{y}_l | \mathbf{x}(t_l))].$$

Mean-field approximation

- The main assumption we will make is that the approximating process q is factorised

$$q_t(\mathbf{x}) = \prod_{i=1}^D q_{it}(x_i) \quad g_t(\mathbf{x}'|\mathbf{x}) = \sum_{i=1}^D \prod_{j \neq i} \delta_{x'_j, x_j} g_{it}(x'_i|x_i)$$

- The KL-divergence then becomes

$$KL(q||p_{post}) = \ln Z - \sum_{l=1}^N E_q [\ln \hat{p}(y_l|\mathbf{x}(t_l))] + \int_0^T dt \sum_i \sum_x q_{it}(x) \sum_{x':x' \neq x} \left\{ g_{it}(x'|x) \ln \frac{g_{it}(x'|x)}{\hat{f}_i(x'|x)} + \tilde{f}_i(x'|x) - g_{it}(x'|x) \right\}$$

so that it decomposes as a sum over species i .

Constraints

- Each of the factors in the approximating distribution satisfy a 1-dimensional Master equation.
- The complexity is decreased from S^D to DS.
- Using Lagrange multiplier functions $\lambda_i(x,t)$, we obtain the Lagrangian

$$L = KL(q||p_{post}) +$$

$$- \sum_i \int_0^T dt \sum_x \lambda_i(x,t) \left(\partial_t q_{it}(x) - \sum_{x' \neq x} \{g_{it}(x|x') q_{it}(x') - g_{it}(x'|x) q_{it}(x)\} \right).$$

with the standard boundary condition $\lambda_i(x,T)=0$.

Functional derivatives

- To optimise the Lagrangian, we compute its functional derivatives

$$\frac{\delta L}{\delta q_{it}(x)} = \sum_{x' \neq x} \left[g_{it}(x'|x) \ln \frac{g_{it}(x'|x)}{\hat{f}_i(x'|x)} - g_{it}(x'|x) + \tilde{f}_i(x'|x) \right] + \partial_t \lambda_i(x, t) +$$

$$\sum_{x'} g_{it}(x'|x) \left\{ \lambda_i(x', t) - \lambda_i(x, t) \right\} - \sum_l \ln \hat{p}(\mathbf{y}_l | \mathbf{x}(t)) \delta(t - t_l) = 0$$

(1)

$$\frac{\delta L}{\delta g_{it}(x'|x)} = q_{it}(x) \left(\ln \frac{g_{it}(x'|x)}{\hat{f}_i(x'|x)} + \lambda_i(x', t) - \lambda_i(x, t) \right) = 0$$

(2)

Backward equation

- Inserting equation (2) into (1) and defining $r_i(x,t) = e^{-\lambda_i(x,t)}$ we obtain

$$\frac{dr_i(x,t)}{dt} = \sum_{x' \neq x} \left(\tilde{f}_i(x'|x) r_i(x,t) - \hat{f}_i(x'|x) r_i(x',t) \right)$$

- For each species, this is a system of S linear differential equations valid at all times except the observation times.

Including observations

- We assume for simplicity that observations for different species are independent

$$\hat{p}(\mathbf{y}_l | \mathbf{x}(t)) = \prod_i \hat{p}_i(y_{il} | x_i(t_l)) \quad \forall l$$

- At observations, the Lagrange multiplier has a discontinuity

$$\lim_{t \rightarrow t_l^-} r_i(x, t) = \hat{p}_i(y_{il} | x_i(t_l)) \lim_{t \rightarrow t_l^+} r_i(x, t).$$

- Numerically, this can become difficult to solve.
- Easier to solve the ODE for the ratios of r .

Algorithm: approximating the posterior (E-step)

- We can find an approximate posterior as follows:
 1. Choose a species i ;
 2. From an initial guess of $q_i(x)$, compute the averaged rates;
 3. Solve the backward equation backward in time from $r_i(x, T)=1$;
 4. Solve the Master equation forward in time to update $q_i(x)$;
 5. Iterate 2-4 until an optimum is found;
 6. Choose another species and follow 2-5;
 7. Iterate until convergence.
- This procedure guarantees a decrease in KL

Parameter estimation (M-step)

- Prior parameters can also be estimated by minimising the KL divergence, once an approximating distribution is obtained.
- E-step and M-step can be iterated until convergence is reached.
- Local minima can be a problem, thought must be given to initialisation.

Application: Lotka-Volterra

- Transition rates are given by

$$f_{prey}(p + 1|p) = \alpha p \quad f_{prey}(p - 1|p) = \beta P p$$

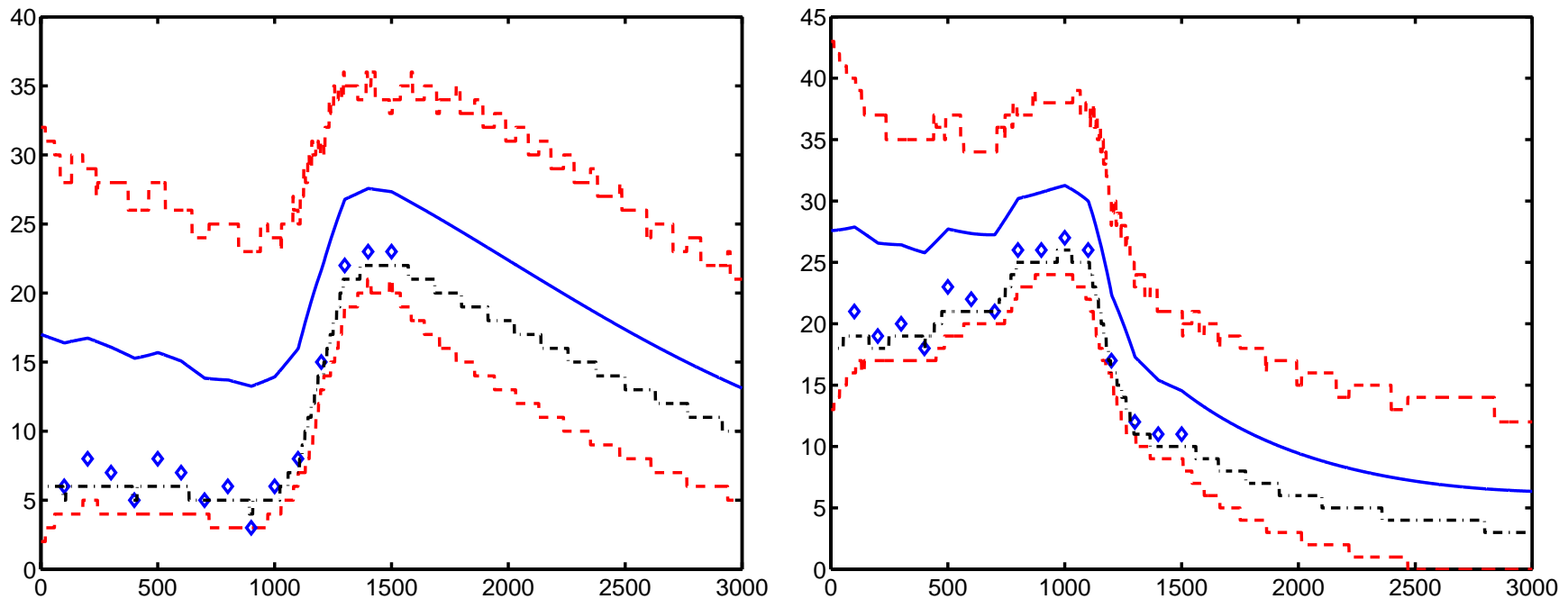
$$f_{predator}(P + 1|P) = \delta P p \quad f_{predator}(P - 1|P) = \gamma P$$

- M-step equations are given by

$$\alpha = \frac{\int_0^T \langle g_{preyt}(x + 1|x) \rangle_{preyt} dt}{\int_0^T dt \langle x \rangle_{preyt}}, \quad \beta = \frac{\int_0^T \langle g_{preyt}(x - 1|x) \rangle_{preyt} dt}{\int_0^T dt \langle x \rangle_{preyt} \langle y \rangle_{predatort}},$$

$$\gamma = \frac{\int_0^T \langle g_{predatort}(y - 1|y) \rangle_{predatort} dt}{\int_0^T dt \langle y \rangle_{predatort}}, \quad \delta = \frac{\int_0^T \langle g_{predatort}(y + 1|y) \rangle_{predatort} dt}{\int_0^T dt \langle y \rangle_{predatort} \langle x \rangle_{preyt}}.$$

Results: Lotka-Volterra



Posterior predator (left) and prey (right) distribution. Diamonds are data, dashed-dotted posterior mode, solid posterior mean. The noise model is asymmetric and given by

$$\hat{p}_i (y_{il} | x_i (t_l)) \propto \left[\frac{1}{5|y_{il} - x_i(t_l)|} + 0.001 \right]$$

Parameter estimates

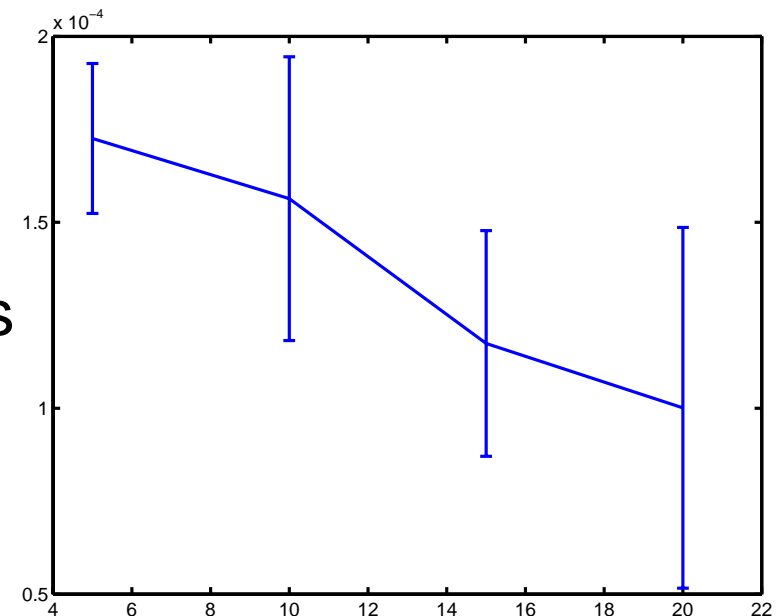
- Estimates of the parameters were reasonable

$$\alpha = 5.14 \times 10^{-4} \quad (5 \times 10^{-4}) \quad \beta = 6.95 \times 10^{-5} \quad (1 \times 10^{-4})$$

$$\gamma = 7.26 \times 10^{-4} \quad (5 \times 10^{-4}) \quad \delta = 5.77 \times 10^{-5} \quad (1 \times 10^{-4})$$

Estimates appear to converge to the true value when more observations are available.

Results show average of five runs per data-set size.



Application: gene auto-regulation

- Autoregulation is one of the fundamental blocks in gene regulatory networks.
- Protein represses transcription of its own coding gene.
- We use a logical approximation. Process rates are

$$f_{RNA}(x+1|x, y) = \alpha (1 - 0.99 \times \Theta(y - y_c)) \quad f_{RNA}(x-1|x, y) = \beta x$$
$$f_p(y+1|x, y) = \gamma x \quad f_p(y-1|x, y) = \delta y$$

The critical parameter y_c is the integer threshold for protein count above which repression begins. Θ is the Heaviside step function.

Parameter estimation

- Fixed point equations for the parameters are obtained

$$\alpha = \frac{\int_0^T dt \langle g_{RNA}(x+1|x) \rangle_{q_{RNA}}}{T \left(1 - 0.99 \frac{1}{T} \int h(y_c) dt \right)} \quad \beta = \frac{\int_0^T dt \langle g_{RNA}(x-1|x) \rangle_{q_{RNA}}}{\int_0^T dt \langle x \rangle_{q_{RNA}}}$$

$$\gamma = \frac{\int_0^T dt \langle g_p(x+1|x) \rangle_{q_p}}{\int_0^T dt \langle x \rangle_{q_{RNA}}} \quad \delta = \frac{\int_0^T dt \langle g_p(x-1|x) \rangle_{q_p}}{\int_0^T dt \langle y \rangle_{q_p}}$$

where $h(y_c) = \sum_{y \geq y_c} q_p(y)$ is the posterior probability that

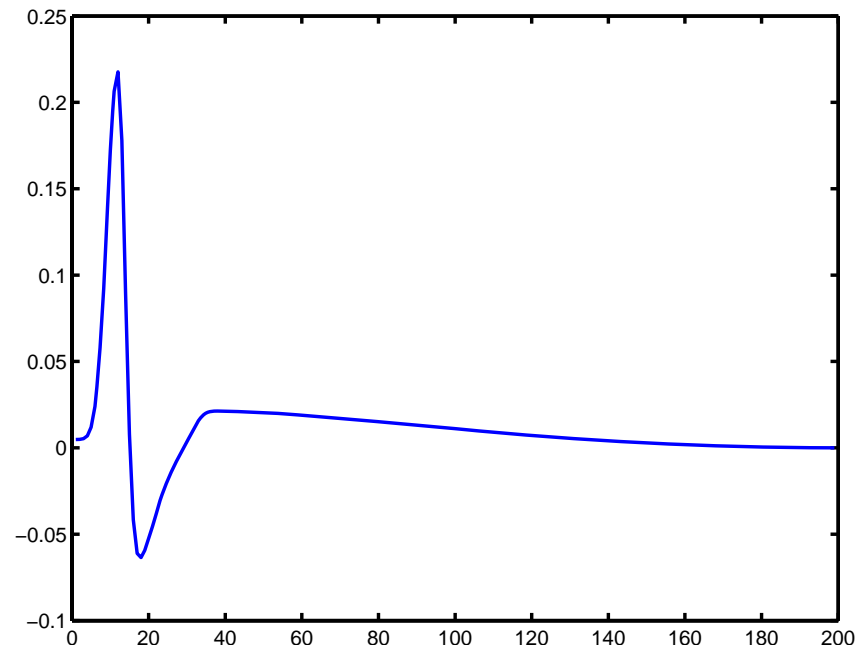
the protein levels will be above the critical threshold (function of time).

Identifiability of critical parameter

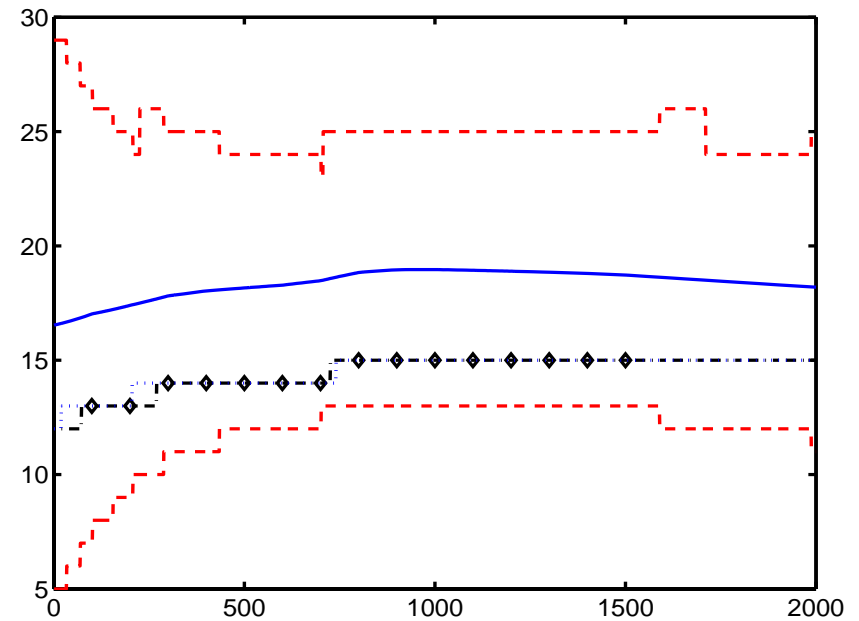
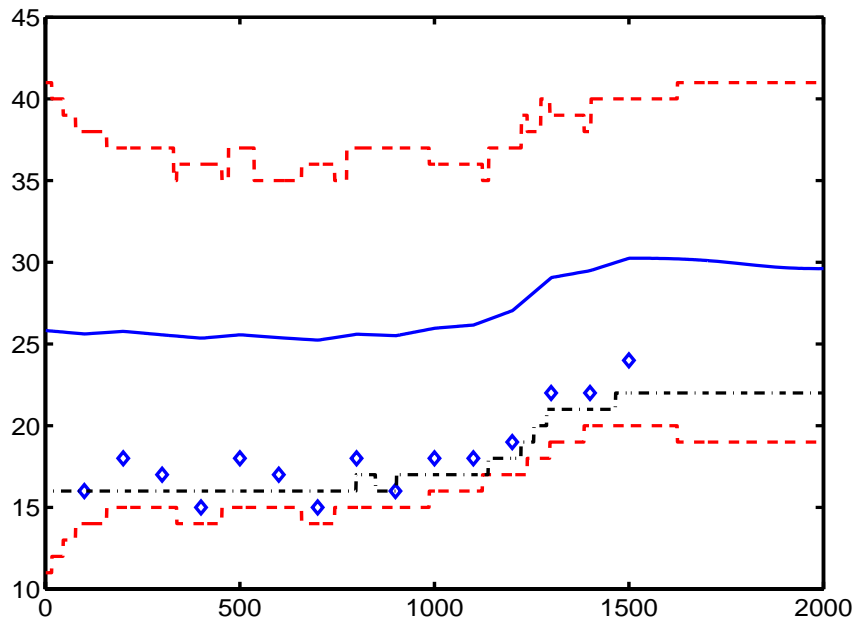
- The critical parameter y_c is found by optimising (by search) the free energy

$$\mathcal{L}_{y_c} = const + \left\{ 2 \int_0^T dt \bar{g} h(y_c) + \log \left[1 - 0.99 \frac{1}{T} \int_0^T h(y_c) dt \right] \int_0^T dt \bar{g} \right\}$$

Clearly, if the protein levels in the data never exceed the threshold, y_c is not identifiable. If it does, we get free energies with a well defined minimum.



Autoregulatory network: results



$$\alpha = 3.4 \times 10^{-3} \quad (2 \times 10^{-3})$$

$$\beta = 4.4 \times 10^{-5} \quad (6 \times 10^{-5})$$

$$\gamma = 1.9 \times 10^{-4} \quad (5 \times 10^{-4})$$

$$\delta = 6.7 \times 10^{-5} \quad (7 \times 10^{-5})$$

Conclusions

- Efficient framework for posterior inference and parameter estimation in MJPs.
- Oustrips MCMC (Boys *et al.*) by orders of magnitude.
- Numerics are tricky and could be improved.
- Issues swept under the carpet: regularisation, cut-off on number of states.
- Readily extends to missing data.

Future work

- More complex, realistic networks.
- Missing data.
- Hybrid systems: one species has large numbers, so approximate as deterministic or diffusion.
- Hybrid systems: particles also have a spatial dimension, and diffuse in the environment.