# Variational Gaussian Process Timeseries Inference

Carl Edward Rasmussen

May 7th 2017

**Abstract**

Variational inference in nonlinear dynamical models.

The model is specified by

$$f_e(\tilde{\mathbf{x}}) \sim \mathcal{GP}\big(0, k_e(\cdot, \cdot)\big), \qquad \text{where } \tilde{\mathbf{x}} = (\mathbf{x}, \mathbf{u}) \text{ and } e = 1, \ldots, E, \tag{1}$$

$$\mathbf{x}_t | \mathbf{f}_t \sim \mathcal{N}(\mathbf{x}_t | \mathbf{f}_t, Q), \qquad Q \text{ diagonal}, \tag{2}$$

$$\mathbf{y}_t | \mathbf{x}_t \sim \mathcal{N}(\mathbf{y}_t | C\mathbf{x}_t, R), \tag{3}$$

and $\mathbf{u}_t, \mathbf{y}_t, t = 1, \ldots, T$ are the control inputs and measurements (both observed), and $\mathbf{f}_t, t = 2, \ldots, T$ and $\mathbf{x}_t, t = 1, \ldots, T$ are unobserved, latent variables. The GPs implement the non-linear transition from one time point to the next conditioned on the state $\mathbf{x}_{t-1}$ and all the previous transition pairs $\mathbf{f}_{2:t-1}, \mathbf{x}_{1:t-2}$

$$\mathbf{f}_t(\mathbf{x}_{t-1}) = p(\mathbf{f}_t | \mathbf{f}_{2:t-1}, \mathbf{x}_{1:t-1}), \text{ where } t = 2, \ldots, T.$$

The joint probability of all the variables is given by the product of $T$ observation probabilities and $T - 1$ transition probabilities

$$p(\mathbf{y}, \mathbf{x}, \mathbf{f}) = \prod_{t=1}^{T} p(\mathbf{y}_t | \mathbf{x}_t) \prod_{t=2}^{T} p(\mathbf{x}_t | \mathbf{f}_t) p(\mathbf{f}_t | \mathbf{f}_{1:t-1}, \mathbf{x}_{1:t-1}).$$

Each GP is augmented with a set of $M$ inducing inputs $\mathbf{z}$ and corresponding targets $\mathbf{v}$ such that $\mathbf{v}_e = f_e(\mathbf{z}_e)$. The augmented joint is

$$p(\mathbf{y}, \mathbf{x}, \mathbf{f}, \mathbf{v}) = p(\mathbf{y} | \mathbf{x}) p(\mathbf{x}, \mathbf{f} | \mathbf{v}) p(\mathbf{v}).$$

Exact inference in the model is intractable, instead we fit the model by optimizing a variational lower bound based on an approximating distribution q, which we chose to have the following form

$$q(\mathbf{x}, \mathbf{f}, \mathbf{v}) = q(\mathbf{v}) q(\mathbf{x}) \prod_{t=2}^{T} p(\mathbf{f}_t | \mathbf{f}_{1:t-1}, \mathbf{x}_{1:t-1}, \mathbf{v}), \text{ where } q(\mathbf{x}) = \mathcal{N}(\mu_x, \Sigma_x),$$

the assumptions being that 1) the joint on $\mathbf{v}$ and $\mathbf{x}$ factorizes, 2) that $q(\mathbf{x})$ is Gaussian and 3) that the conditional $q(\mathbf{f} | \mathbf{x}, \mathbf{v})$ is chosen to be equal to the conditional *prior*. Generally, we would expect the variational bound to be tight if the approximating distribution is close to the *posterior*, but for tractability we are forced to set the conditional $q(\mathbf{f} | \mathbf{x}, \mathbf{v})$ to be equal to the conditional prior. This may still be a good approximation, since we are conditioning on the inducing targets $\mathbf{v}$. If the inducing targets are able to capture the properties of the posterior, then the bound may still be good.

The variational log marginal likelihood lower bound is a single time series (contributions for multiple time series are simply added together)

$$
\begin{aligned}
\mathcal{L}(\mathbf{y} | q(\mathbf{v}), q(\mathbf{x}), \theta) = {}& -\mathrm{KL}(q(\mathbf{v}) \| p(\mathbf{v})) + H(q(\mathbf{x})) + \sum_{t=1}^{T} \langle \log p(\mathbf{y}_t | \mathbf{x}_t) \rangle_{q(\mathbf{x}_t)} \\
& + \sum_{t=2}^{T} -\tfrac{1}{2} \mathrm{tr}(Q^{-1} \langle B_{t-1} \rangle_{q(\mathbf{x}_{t-1})}) + \langle \log \mathcal{N}(\mathbf{x}_t | A_{t-1} \mathbf{v}, Q) \rangle_{q(\mathbf{v}), q(\mathbf{x}_{t-1:t})},
\end{aligned}
\tag{4}
$$

with the following definitions

$$A_{t-1} = k(\mathbf{x}_{t-1}, \mathbf{z}) K^{-1}, \text{ and } B_{t-1} = k(\mathbf{x}_{t-1}, \mathbf{x}_{t-1}) - k(\mathbf{x}_{t-1}, \mathbf{z}) K^{-1} k(\mathbf{z}, \mathbf{x}_{t-1}).$$

A free form optimization of this bound wrt $q(\mathbf{v})$ yields independent Gaussians for each GP

$$q^*(\mathbf{v}_e) = \mathcal{N}\big(\boldsymbol{\mu}_e = K_e(K_e + \Psi_{2e})^{-1}\Psi_{1e},\ \Sigma_e = K_e(K_e + \Psi_{2e})^{-1}K_e\big), \tag{5}$$

where $K_e = k_e(\mathbf{z}_e, \mathbf{z}_e)$ and we have defined the expectations

$$\Psi_1 = \sum_{t=2}^{T}\langle k(\mathbf{z}, \mathbf{x}_{t-1})Q^{-1}\mathbf{x}_t\rangle_{q(\mathbf{x}_{t-1:t})},\ \text{and}\ \Psi_2 = \sum_{t=2}^{T}\langle k(\mathbf{z}, \mathbf{x}_{t-1})Q^{-1}k(\mathbf{x}_{t-1}, \mathbf{z})\rangle_{q(\mathbf{x}_{t-1})}, \tag{6}$$

of size $M \times E$ and $M \times M \times E$ respectively.

Plugging the optimal $q^*(\mathbf{v})$ back into the bound eq. (4), we get

$$\mathcal{L}(\mathbf{y}|q(\mathbf{x}), \theta) = -\mathrm{KL}(q^*(\mathbf{v})\|p(\mathbf{v})) + H(q(\mathbf{x})) + \sum_{t=1}^{T}\langle \log p(\mathbf{y}_t|\mathbf{x}_t)\rangle_{q(\mathbf{x}_t)}$$

$$+ \sum_{t=2}^{T} -\tfrac{1}{2}\langle \mathrm{tr}\big(Q^{-1}(B_{t-1} + A_{t-1}\Sigma A_{t-1})\big)\rangle_{q(\mathbf{x}_{t-1})} + \langle \log \mathcal{N}(\mathbf{x}_t|A_{t-1}\boldsymbol{\mu}, Q)\rangle_{q(\mathbf{x}_{t-1:t})}. \tag{7}$$

Note that except for the entropy $H(q(\mathbf{x}))$, the bound only depends on $q(\mathbf{x})$ through its pair-wise marginals. This means that the model will be identical for all $q(\mathbf{x})$ which have the same pair-wise marginals, except for an offset in the bound which depends on the entropy. We will chose $q(\mathbf{x})$ to be Markovian, ie the precision $\Sigma_x^{-1}$ is block tri-diagonal.

**Transition model**

Writing out each term from the transition model from eq. (7) in detail

$$-\mathrm{KL}(q^*(\mathbf{v})\|p(\mathbf{v})) = -\tfrac{1}{2}\sum_{e=1}^{E}\mathrm{tr}(K_e + \Psi_{2e})^{-1}K_e + \boldsymbol{\mu}_e^\top K_e^{-1}\boldsymbol{\mu}_e - M - \log|(K_e + \Psi_{2e})^{-1}K_e|, \tag{8}$$

and

$$-\tfrac{1}{2}\sum_{t=2}^{T}\mathrm{tr}Q^{-1}\langle B_{t-1}\rangle_{q(\mathbf{x}_{t-1})} = -\tfrac{T-1}{2}\mathrm{tr}Q^{-1} + \tfrac{1}{2}\sum_{t=2}^{T}\mathrm{tr}K^{-1}\langle k(\mathbf{z}, \mathbf{x}_{t-1})Q^{-1}k(\mathbf{x}_{t-1}, \mathbf{z})\rangle_{q(\mathbf{x}_{t-1})}$$

$$= \tfrac{1}{2}\sum_{e=1}^{E}\mathrm{tr}K_e^{-1}\Psi_{2e} - \tfrac{T-1}{2}\mathrm{tr}Q^{-1}, \tag{9}$$

and

$$-\tfrac{1}{2}\sum_{t=2}^{T}\mathrm{tr}Q^{-1}\langle A_{t-1}\Sigma A_{t-1}\rangle_{q(\mathbf{x}_{t-1})} = -\tfrac{1}{2}\sum_{t=2}^{T}\mathrm{tr}\big(\Sigma K^{-1}\langle k(\mathbf{x}_{t-1}, \mathbf{z})Q^{-1}k(\mathbf{z}, \mathbf{x}_{t-1})\rangle_{q(\mathbf{x}_{t-1})}K^{-1}\big)$$

$$= -\tfrac{1}{2}\sum_{e=1}^{E}\mathrm{tr}(K_e + \Psi_{2e})^{-1}\Psi_{2e}, \tag{10}$$

and

$$\sum_{t=2}^{T}\langle \log \mathcal{N}(\mathbf{x}_t|A_{t-1}\boldsymbol{\mu}, Q)\rangle_{q(\mathbf{x}_{t-1:t})}$$

$$= -\tfrac{(T-1)E}{2}\log(2\pi) - \tfrac{T-1}{2}\log|Q| - \tfrac{1}{2}\sum_{t=2}^{T}\langle(\mathbf{x}_t - A_{t-1}\boldsymbol{\mu})^\top Q^{-1}(\mathbf{x}_t - A_{t-1}\boldsymbol{\mu})\rangle_{q(\mathbf{x}_{t-1:t})}$$

$$= -\tfrac{(T-1)E}{2}\log(2\pi) - \tfrac{T-1}{2}\log|Q| - \tfrac{1}{2}\mathrm{tr}Q^{-1}\sum_{t=2}^{T}\langle \mathbf{x}_t^\top \mathbf{x}_t\rangle_{q(\mathbf{x}_t)} + \boldsymbol{\mu}^\top\langle \mathbf{x}_t Q^{-1}A_{t-1}\rangle_{q(\mathbf{x}_{t-1:t})} \tag{11}$$

$$- \tfrac{1}{2}\boldsymbol{\mu}^\top K^{-1}\langle k(\mathbf{x}_{t-1}, \mathbf{z})Q^{-1}k(\mathbf{z}, \mathbf{x}_{t-1})\rangle_{q(\mathbf{x}_{t-1})}K^{-1}\boldsymbol{\mu}$$

$$= -\tfrac{(T-1)E}{2}\log(2\pi) - \tfrac{T-1}{2}\log|Q| - \tfrac{1}{2}\mathrm{tr}Q^{-1}\sum_{t=2}^{T}\langle \mathbf{x}_t^\top \mathbf{x}_t\rangle_{q(\mathbf{x}_t)} - \tfrac{1}{2}\boldsymbol{\mu}^\top K^{-1}\Psi_2 K^{-1}\boldsymbol{\mu} + \boldsymbol{\mu}^\top K^{-1}\Psi_1.$$

2

Collecting terms form eq. (8-11) which depend on $\Psi_1$ and $\Psi_2$, two possibilies arise, either training or testing, in both cases we introduce $\mu$ and $\Sigma$ from eq. (5) for the *training* cases, giving rise to

$$
\begin{aligned}
\Psi &= \tfrac{1}{2}\sum_{e=1}^{E} \log|K_e| - \log|K_e + \Psi_{2e}| + \mathrm{tr}K_e^{-1}\Psi_{2e} + \Psi_{1e}^{\top}(K_e + \Psi_{2e})^{-1}\Psi_{1e}, \\
\Psi^* &= \tfrac{1}{2}\sum_{e=1}^{E} \mathrm{tr}K_e^{-1}\Psi_{2e}(K_e + \Psi_{2e})^{-1}\Psi_{2e}^* - \mathrm{tr}w_e w_e^{\top}\Psi_{2e}^* + 2\Psi_{1e}^{\top}(K_e + \Psi_{2e})^{-1}\Psi_{1e}^*,
\end{aligned}
\tag{12}
$$

where $w_e = (K_e + \Psi_{2e})^{-1}\Psi_{1e}$, for training and test respectively. Pulling together all terms from eq. (8-11) and eq. (12) we get the following contribution to the log likelihood

$$
\Psi - \tfrac{1}{2}\mathrm{tr}Q^{-1}\sum_{t=2}^{T}(I + \mu_t^{\top}\mu_t + \Sigma_{t,t}) - \tfrac{T-1}{2}\log|Q| - \tfrac{(T-1)E}{2}\log(2\pi).
\tag{13}
$$

## Entropy

The entropy of Markovian Gaussian with specified $E$ dimensional marginals and $2E$ dimensional consequtive pair-wise marginals and marginals is given by

$$
\mathcal{H}(q(\mathbf{x})) = \tfrac{TE}{2}(1 + \log(2\pi)) + \tfrac{1}{2}\sum_{t=2}^{T}\log|\Sigma_{t-1:t,t-1:t}| - \tfrac{1}{2}\sum_{t=2}^{T-1}\log|\Sigma_t|
\tag{14}
$$

3    ⟨*entropy* 3⟩≡            (5b)

```
1  function [L dLd dLo] = gaussMarkovEntropy(d, o);
2  [E, E, T] = size(d); dd = zeros(T,1); dp = zeros(T-1,1);
3  for t = 1:T, dd(t) = det(d(:,:,t)); end                    % det of diagonals
4  for t = 1:T-1, dp(t) = dd(t)*det(d(:,:,t+1)-o(:,:,t)'/d(:,:,t)*o(:,:,t)); end;
5  L = E*T*(1+log(2*pi))/2 + sum(log(dp))/2 - sum(log(dd(2:T-1)))/2;    % entropy
6  if nargout > 1                                              % want derivatives?
7    dLd = zeros(E,E,T); dLo = zeros(E,E,T-1);
8    for t = 1:T-1
9      dLd(:,:,t) = dLd(:,:,t) + inv(d(:,:,t)-o(:,:,t)/d(:,:,t+1)*o(:,:,t)')/2;
10     dLd(:,:,t+1) = inv(d(:,:,t+1)-o(:,:,t)'/d(:,:,t)*o(:,:,t))/2;
11     dLo(:,:,t) = -d(:,:,t)\o(:,:,t)/(d(:,:,t+1)-o(:,:,t)'/d(:,:,t)*o(:,:,t));
12   end
13   for t = 2:T-1, dLd(:,:,t) = dLd(:,:,t) - inv(d(:,:,t))/2; end
14 end
```

## Likelihood

The linear Gaussian log likelihood is

$$
\sum_{t=1}^{T}\langle\log p(\mathbf{y}_t|\mathbf{x}_t)\rangle_{q(\mathbf{x}_t)} = -\tfrac{DT}{2}\log(2\pi) - \tfrac{T}{2}\log|R| - \tfrac{1}{2}\mathrm{tr}R^{-1}\sum_{t=1}^{T}\left((\mathbf{y} - C\mu_t)(\mathbf{y} - C\mu_t)^{\top} + C\Sigma_t C^{\top}\right).
\tag{15}
$$

Maximizing the log likelihood wrt observation noise covariance $R$ and the parameters $C$ yields:

$$
R^* = \tfrac{1}{T}\left[\sum_{t=1}^{T}\mathbf{y}_t\mathbf{y}_t^{\top} - C^*\sum_{t=1}^{T}\mu_t\mathbf{y}^{\top}\right], \quad\text{and}\quad C^* = \sum_{t=1}^{T}\mathbf{y}_t\mu_t^{\top}\left[\sum_{t=1}^{T}\mu_t\mu_t^{\top} + \Sigma_{t,t}\right]^{-1},
\tag{16}
$$

and the maximum attained is

$$
\mathcal{L}^* = -\tfrac{DT}{2}(1 + \log(2\pi)) - \tfrac{T}{2}\log|R^*|,
\tag{17}
$$

with derivatives

$$
\frac{\mathcal{L}^*}{\partial\mu_t} = C^{\top}R^{-1}(\mathbf{y}_t - C\mu_t), \quad\text{and}\quad \frac{\mathcal{L}^*}{\partial\Sigma_{t,t}} = -\tfrac{1}{2}C^{\top}R^{-1}C, \quad\text{evaluated at } C = C^*, \text{ and } R = R^*.
\tag{18}
$$

In the test case, when we are inferring the latent space representation of a given short sequence, we use the $R^*$ and $C^*$ parameters derived from the training sequences, so that the contribution to the log likelihood does not take on the simple form $\mathcal{L}^*$, but must be calculated in full from eq. (16).

4a ⟨*likelihood* 4a⟩≡ (5b)

```
 1 function [L, R, C, dm, dS] = likelihood(y, qx, lat)
 2 D = size(y{1},2); E = size(qx(1).m,1); T = sum(arrayfun(@(x)size(x.m,2),qx));
 3 N = size(y,2); yy = zeros(D); ym = zeros(D, E+1); mm = zeros(E+1);
 4 for n = 1:N
 5   m = [qx(n).m' ones(size(qx(n).m,2),1)]; mm = mm + m'*m; ym = ym + y{n}'*m;
 6   yy = yy + y{n}'*y{n}; mm(1:E,1:E) = mm(1:E,1:E) + sum(qx(n).Sd,3);
 7 end
 8 if nargin < 3
 9   C = ym/mm; R = (yy - C*ym')/T;
10   L = -D*T*(1+log(2*pi))/2 - T*sum(log(diag(chol(R))));        % log likelihood
11 else %We are in the testing case, using a provided latent mapping
12   C = lat.C; R = lat.R; d = zeros(1,N);
13   for n=1:N
14     CsC = 0; for t=1:T, CsC = CsC + C(:,E)*qx(n).Sd(:,:,t)*C(:,E)'; end
15     d(1,n) = sum(sum((y{n}' - C*[qx(n).m' ones(size(qx(n).m,2),1)]).^2)) + CsC;
16   end
17   L = -D*T*log(2*pi)/2 - T*sum(log(diag(chol(R))))/2 - tr(inv(R))*(sum(d))/2;
18 end
19 if nargout > 3                                    % do we want derivatives?
20   dm = cell(N,1);
21   for n = 1:N,
22     dm{n} = C(:,1:E)'/R*(y{n}'-C*[qx(n).m; ones(1,size(qx(n).m,2))]);
23   end
24   dS = -C(:,1:E)'/R*C(:,1:E)/2;             % all dS identical, return once only
25 end
```

## The lower bound

Pulling together all terms

$$\mathcal{L}(\mathbf{y}|q(\mathbf{x}), \theta) = \frac{1}{2}\sum_{e=1}^{E} \log|(K_e + \Psi_{2e})^{-1}K_e| + \text{tr}K_e^{-1}\Psi_{2e} + \Psi_{1e}^{\top}(K_e + \Psi_{2e})^{-1}\Psi_{1e} + \frac{1}{2}\sum_{t=2}^{T}\log|\Sigma_{t-1:t,t-1:t}|$$

$$- \frac{1}{2}\sum_{t=2}^{T-1}\log|\Sigma_t| - \frac{1}{2}\text{tr}Q^{-1}\sum_{t=2}^{T}(I + \boldsymbol{\mu}_t^{\top}\boldsymbol{\mu}_t + \Sigma_{t,t}) - \frac{T-1}{2}\log|Q| - \frac{T}{2}\log|R^*| - \frac{(D-E)T}{2} - \frac{TD-E}{2}\log(2\pi).$$

(19)

4b ⟨*lower bound* 4b⟩≡ (5b)

```
 1 [L1, dnlml] = Psi(hyp, qx, z, u);
 2 T = sum(arrayfun(@(x)size(x.m,2),qx)); L2 = 0; L3 = 0;
 3 dLd = cell(1,N); dLo = cell(1,N);
 4 for n = 1:N
 5   L2 = L2 + sum(qx(n).m(:,2:end).^2,2) + diag(sum(qx(n).Sd(:,:,2:end),3));
 6   [L dLd{n} dLo{n}] = gaussMarkovEntropy(qx(n).Sd, qx(n).So); L3 = L3 + L;
 7 end
 8 L5 = -exp(-2*[hyp(:).pn]) * (L2+T-N) / 2;
 9 L4 = -(T-N)*sum([hyp(:).pn])-(T-N)*E*log(2*pi)/2;
10 [L6, R, C, dm, dS] = likelihood(y, qx);
11 nlml = -L1-L5-L3-L4-L6;
12 %keyboard
```

4c ⟨*bound derivatives* 4c⟩≡ (5b)

```
 1 for e = 1:E
 2   dnlml.hyp(e).pn = dnlml.hyp(e).pn - exp(-2*hyp(e).pn)*(L2(e)+T-N)+T-N;
 3 end
 4 iQ = diag(exp(-2*[hyp(:).pn]));
 5 for n = 1:N
 6   dnlml.qx(n).m(:,2:end) = dnlml.qx(n).m(:,2:end) + iQ*qx(n).m(:,2:end);
 7   dnlml.qx(n).m = dnlml.qx(n).m - dm{n};
```

```
 8    dnlml.qx(n).Sd(:,:,2:end) = bsxfun(@plus, dnlml.qx(n).Sd(:,:,2:end), iQ/2);
 9    dnlml.qx(n).Sd = dnlml.qx(n).Sd - bsxfun(@plus, dS, dLd{n});
10    dnlml.qx(n).So = dnlml.qx(n).So - dLo{n};
11 end
12
13 out1 = nlml; out2 = dnlml; out3 = struct('C', C, 'R', R);    % rename outputs
```

5a  ⟨*predictions* 5a⟩≡                                                                  (5b)
```
 1 [Psi1, Psi2] = Psi(hyp, qx, z, u);
```

5b  ⟨*vgpt.m* 5b⟩≡
```
 1 function [out1, out2, out3] = vgpt(p, data, x);
 2 ⟨usage 5c⟩
 3
 4 N = length(p.qx); z = p.z; [M, F, E] = size(z); D = size(data(1).y,2); hyp = p.hyp;
 5 u = arrayfun(@(x)(x.u),data,'UniformOutput',false); [qx(1:N).m] = deal(p.qx(:).m);
 6 y = arrayfun(@(x)(x.y),data,'UniformOutput',false);
 7 for n = 1:N, [qx(n).Sd qx(n).So] = convert(p.qx(n).s); end % convert covariance
 8
 9 if nargin == 2
10   ⟨lower bound 4b⟩
11   ⟨bound derivatives 4c⟩
12   ⟨revert covariances 9b⟩
13 else
14   ⟨predictions 5a⟩
15 end
16
17 ⟨Psi 5d⟩
18 ⟨entropy 3⟩
19 ⟨likelihood 4a⟩
20 ⟨convert 9a⟩
21 ⟨revert 9c⟩
22 ⟨maha 10a⟩
23 ⟨prediction 10b⟩
```

5c  ⟨*usage* 5c⟩≡                                                                  (5b)
```
 1 % Variational GP Timeseries inference. Compute the nlml lower bound and its
 2 % derivative wrt hyp hyperparameters, qx distribution and z inducing inputs.
 3 %
 4 % p                    parameter struct
 5 %   hyp      1 x E     GP hyperparameter struct
 6 %     l      F x 1     log length scale
 7 %     pn     1 x 1     log process noise std dev
 8 %   qx       1 x N     struct array for Gaussian q(x) distribution
 9 %     m      E x T_n   mean
10 %     s    2ExE x T_n  representation of covariance
11 %   z      M x F x E   inducing inputs
12 % data       1 x N     data struct
13 %   y      T_n x D     cell array of observations
14 %   u      T_n x U     cell array of control inputs
15 %
16 % Copyright (C) 2016 by Carl Edward Rasmussen, 20160530.
```

## The Ψ function

In the implementation, a function `Psi` handles the part of the (negative) log marginal likelihood which depends on the quantities $\Psi_1$ and $\Psi_2$:

$$\psi = \frac{1}{2}\sum_{e=1}^{E}\log|K_e| - \log|K_e + \Psi_{2e}| - \text{tr}K_e^{-1}\Psi_{2e} - \Psi_{1e}^{\top}(K_e + \Psi_{2e})^{-1}\Psi_{1e}. \tag{20}$$

5d  ⟨*Psi* 5d⟩≡                                                                  (5b)

```
 1 function [lml, dnlml] = Psi(hyp, qx, z, u, test);
 2 % hyp        1 x E    GP hyperparameter struct
 3 %   l        F x 1    log length scale
 4 %    pn      1 x 1    log process noise std dev
 5 % qx         1 x N    Gaussian q(x) distribution
 6 %    m       E x T_n    mean
 7 %   Sd  ExE x T_n    diagonal elements of covariance matrix
 8 %   So  ExE x T_n-1  immediately off-diagonal elements of covariance matrix
 9 % z        M x F x E  inducing inputs
10 % u        T_n x U    cell array of control inputs
11 % lml        1 x 1    contribution to the log marginal likelihood
12 % dnlml               derivatives
13 % test               flag indicating test mode
14
15 persistent K Psi1 Psi2;                     % keep these around if necessary
16 [M, F, E] = size(z);                                        % get sizes
17 if ~test %If we test, we use the stored values of Psi1, Psi2 from the training set.
18   K = zeros(M,M,E); Psi1 = zeros(M,E); Psi2 = zeros(M,M,E); Sd = zeros(F,F);
19 end
20 if nargin < 5 %Use to return stored Psi values
21   lml = Psi1; dnlml = Psi2;
22 else
23   ⟨expectations 6⟩
24   if nargout > 0
25     ⟨expectation derivatives 7a⟩
26   end
27 end
```

**The $\Psi_1$ and $\Psi_2$ expectations**

The expectations from eq. (6) and derivatives wrt hyperparameters, the parameters of the $q(\mathbf{x})$ distribution and the pseudo-inputs $\mathbf{z}$ are calculated by the Psi function. To compute these expectations, the pairwise joint

$$q(\mathbf{x}_{t-1:t}) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_{t-1} \\ \boldsymbol{\mu}_t \end{bmatrix}, \begin{bmatrix} \Sigma_{t-1,t-1} & \Sigma_{t-1,t} \\ \Sigma_{t,t-1} & \Sigma_{t,t} \end{bmatrix}\right),$$

is multiplied with the covariance function, which can be written as an un-normalized joint Gaussian

$$k_e(\mathbf{x}_{t-1}, \mathbf{z}_{ie}) = \exp\left(-\tfrac{1}{2}\begin{bmatrix} \mathbf{x}_{t-1} - \mathbf{z}_{ie} \\ \mathbf{x}_t \end{bmatrix}^\top \begin{bmatrix} \Lambda_e^{-1} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x}_{t-1} - \mathbf{z}_{ie} \\ \mathbf{x}_t \end{bmatrix}\right),$$

yielding

$$\int \mathbf{x}_t k_e(\mathbf{x}_{t-1}, \mathbf{z}_{ie}) q(\mathbf{x}_{t-1:t}) d\mathbf{x}_{t-1} d\mathbf{x}_t = \left(\boldsymbol{\mu}_t + \Sigma_{t,t-1}[\Lambda_e + \Sigma_{t-1,t-1}]^{-1}(\mathbf{z}_{ie} - \boldsymbol{\mu}_{t-1})\right)$$
$$\times |I + \Lambda_e^{-1}\Sigma_{t-1,t-1}|^{-1/2} \exp\left(-\tfrac{1}{2}(\boldsymbol{\mu}_{t-1} - \mathbf{z}_{ie})[\Lambda_e + \Sigma_{t-1,t-1}]^{-1}(\boldsymbol{\mu}_{t-1} - \mathbf{z}_{ie})\right). \tag{21}$$

For $\Psi_2$ we have from eq. (6)

$$\int k_e(\mathbf{z}_{ie}, \mathbf{x}_{t-1}) k_e(\mathbf{x}_{t-1}, \mathbf{z}_{je}) q(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1} = \exp(-(\mathbf{z}_{ie} - \mathbf{z}_{je})\Lambda_e^{-1}(\mathbf{z}_{ie} - \mathbf{z}_{je})/4)$$
$$\times |I + 2\Lambda_e^{-1}\Sigma_{t-1,t-1}|^{-1/2} \exp(-(\tfrac{\mathbf{z}_{ie}+\mathbf{z}_{je}}{2} - \boldsymbol{\mu}_{t-1})[\Lambda_e/2 + \Sigma_{t-1,t-1}]^{-1}(\tfrac{\mathbf{z}_{ie}+\mathbf{z}_{je}}{2} - \boldsymbol{\mu}_{t-1})/2). \tag{22}$$

Both $\Psi_1$ and $\Psi_2$ are computed for each GP $e = 1, \ldots, E$, each inducing input $\mathbf{z}_{ie}, i = 1, \ldots, M$, and added over (N time series and) $T_n - 1$ time points:

6    ⟨expectations 6⟩≡                      (5d)

```
1 lml = 0;
2 for e = 1:E                                         % for each GP
3   K(:,:,e) = exp(-maha(z(:,:,e),[],diag(exp(-2*hyp(e).l)))/2) + 1e-6*eye(M);
4   iL = diag(exp(-hyp(e).l)); L2 = diag(exp(2*hyp(e).l));
5   b1 = zeros(M,1); b2 = zeros(M,M);
6   for n = 1:length(qx)                              % for each time series
```

```
7       for t = 2:size(qx(n).m, 2)                              % for each time step
8          Sd(1:E,1:E) = qx(n).Sd(:,:,t-1);          % covariance in top left corner
9          r1 = prod(diag(chol(eye(F)+iL*Sd*iL)));                      % sqrt det
10         r2 = prod(diag(chol(eye(F)+2*iL*Sd*iL)));                    % sqrt det
11         s = bsxfun(@minus, z(:,:,e), [qx(n).m(:,t-1)' u{n}(t-1,:)]);
12         a = s/(L2+Sd);
13         b1 = b1 + (qx(n).m(e,t) + a(:,1:E)*qx(n).So(:,e,t-1)) ...
14                                              .*exp(-sum(a.*s,2)/2)/r1;
15         b2 = b2 + exp(-maha(s,-s,inv(L2+2*Sd))/4) / r2;
16       end
17     end
18     if test
19       w = (K(:,:,e)+Psi2(:,:,e))\Psi1(:,e);
20       W = -K(:,:,e)\Psi2(:,:,e)/(Psi2(:,:,e)+K(:,:,e)) + w*w';
21       lml = lml + exp(-2*hyp(e).pn)*(b1'*w(:,e) ...
22                       - sum(sum(b2.*exp(-maha(z(:,:,e),[],inv(L2))/4).*W))/2);
23     else
24       Psi1(:,e) = b1 * exp(-2*hyp(e).pn);
25       Psi2(:,:,e) = b2 * exp(-2*hyp(e).pn) .* exp(-maha(z(:,:,e),[],inv(L2))/4);
26       lml = lml - sum(log(diag(chol(K(:,:,e)+Psi2(:,:,e))))) + ...
27             sum(log(diag(chol(K(:,:,e))))) + trace(K(:,:,e)\Psi2(:,:,e))/2 + ...
28                            Psi1(:,e)'/(K(:,:,e)+Psi2(:,:,e))*Psi1(:,e)/2;
29     end
30 end
```

Note that in the implementation the state distribution $q(\mathbf{x})$ is concatenated with the (deterministic) control inputs $\mathbf{u}$.

### Psi **derivatives**

We need to compute the derivatives of $\Psi$ wrt the parameters of the $q(\mathbf{x})$ distribution, wrt the $\mathbf{u}$ inducing inputs and wrt the hyperparameters. First, from eq. (20) we note

$$\frac{\partial \psi}{\partial \Psi_{1e}} = -(K_e + \Psi_{2e})^{-1}\Psi_{1e} = -\mathbf{w}_e,$$

$$\frac{\partial \psi}{\partial \Psi_{2e}} = -\tfrac{1}{2}K_e^{-1}\Psi_{2e}(K_e + \Psi_{2e})^{-1} + \tfrac{1}{2}\mathbf{w}_e\mathbf{w}_e^\top = -\tfrac{1}{2}R_e(K_e + \Psi_{2e})^{-1} + \tfrac{1}{2}\mathbf{w}_e\mathbf{w}_e^\top = W_e,$$

$$\frac{\partial \psi}{\partial K_e} = -\tfrac{1}{2}R_e(K_e + \Psi_{2e})^{-1}R_e^\top + \tfrac{1}{2}\mathbf{w}_e\mathbf{w}_e^\top = V_e,$$

where we have defined $\mathbf{w}_e = (K_e + \Psi_{2e})^{-1}\Psi_{1e}$ and $R_e = K_e^{-1}\Psi_{2e}$. These can be used together with the derivatives of $\Psi_{1e}$, $\Psi_{2e}$ and $K_e$ and the chain rule to get the desired derivatives.

7a    ⟨*expectation derivatives* 7a⟩≡                                              (5d)
```
1 dnlml.z = zeros(M,F,E);
2 for n = 1:size(qx,2), dnlml.qx(n).m = 0*qx(n).m; dnlml.qx(n).So = 0*qx(n).So;
3 dnlml.qx(n).Sd = -0*qx(n).Sd; end;
4 for e = 1:E
5    w = (K(:,:,e)+Psi2(:,:,e))\Psi1(:,e);
6    R = K(:,:,e)\Psi2(:,:,e);
7    W = -R/(Psi2(:,:,e)+K(:,:,e)) + w*w';
8    ⟨hyp derivatives 7b⟩
9    ⟨Psi derivatives 7c⟩
10 end
```

7b    ⟨*hyp derivatives* 7b⟩≡                                                      (7a)
```
1 dnlml.hyp(e).pn = 2*sum(w.*Psi1(:,e)) - sum(sum(W.*Psi2(:,:,e)));
```

7c    ⟨*Psi derivatives* 7c⟩≡                                                      (7a)
```
1 iL = diag(exp(-hyp(e).l)); L2 = diag(exp(2*hyp(e).l));
2 W1 = W .* exp(-maha(z(:,:,e),[],inv(L2))/4);
3 D = zeros(M,F); H = zeros(F,1);
4 for n = 1:length(qx)
5    T = size(qx(n).m,2);
```

```matlab
 6    A = zeros(E,T); B = zeros(E,E,T); C = zeros(E,E,T-1);
 7    for t = 2:T
 8      Sd(1:E,1:E) = qx(n).Sd(:,:,t-1);            % covariance in top left corner
 9      r2 = prod(diag(chol(eye(F)+2*iL*Sd*iL)));                    % sqrt det
10      s = bsxfun(@minus, z(:,:,e), [qx(n).m(:,t-1)' u{n}(t-1,:)]);
11      a = s/(L2+Sd);
12      a2 = s/(2*L2+4*Sd);
13      SiS = (L2(1:E,1:E)+Sd(1:E,1:E))\qx(n).So(:,e,t-1);
14      r = exp(-sum(a.*s,2)/2) / prod(diag(chol(eye(F)+iL*Sd*iL)));
15      g = (qx(n).m(e,t) + a(:,1:E)*qx(n).So(:,e,t-1)).*w.*r;
16      W2 = W1.*exp(-maha(s,-s,inv(L2+2*Sd))/4);
17      X = bsxfun(@plus,permute(a2,[1 3 2]),permute(a2,[3 1 2]));
18      A(:,t-1) = A(:,t-1) + SiS*(w'*r) - a(:,1:E)'*g + ...
19                      squeeze(sum(sum(bsxfun(@times,W2,X(:,:,1:E)),2),1))/r2;
20      A(e,t) = -w'*r;
21      B(:,:,t-1) = squeeze(sum(sum(bsxfun(@times, ...
22        bsxfun(@times,W2,X(:,:,1:E)),permute(X(:,:,1:E),[1 2 4 3])),2),1))/r2;
23      B(:,:,t-1) = B(:,:,t-1) + SiS*((w.*r)'*a(:,1:E)) ...
24                             + inv(L2(1:E,1:E)+Sd(1:E,1:E))*sum(g)/2 ...
25                             - a(:,1:E)'*bsxfun(@times,g,a(:,1:E))/2 ...
26                    - inv(L2(1:E,1:E)+2*Sd(1:E,1:E))*sum(sum(W2))/r2/2;
27      C(:,e,t-1) = -bsxfun(@times,a(:,1:E),r)'*w;
28      if ~test
29        D(:,1:E) = D(:,1:E) - bsxfun(@times,w,r)*SiS';
30        D = D + bsxfun(@times,g,a) - W2*a2/r2 - bsxfun(@times,sum(W2,2),a2)/r2;
31        H = H + diag(Sd/(L2+2*Sd))*sum(sum(W2))/r2 ...
32          + exp(2*hyp(e).l).*squeeze(sum(sum(bsxfun(@times,W2,X.^2),2),1))/r2;
33        H(1:E) = H(1:E) ...
34              + 2*exp(2*hyp(e).l(1:E)).*diag(SiS*bsxfun(@times,w,r)'*a(:,1:E));
35        H = H - diag(Sd/(L2+Sd))*sum(g);
36        H = H - exp(2*hyp(e).l).*diag(a'*bsxfun(@times,g,a));
37      end
38    end
39    dnlml.qx(n).m = dnlml.qx(n).m + A * exp(-2*hyp(e).pn);
40    dnlml.qx(n).Sd = dnlml.qx(n).Sd + B * exp(-2*hyp(e).pn);
41    dnlml.qx(n).So = dnlml.qx(n).So + C * exp(-2*hyp(e).pn);
42  end
43  if ~test
44  G = W.*Psi2(:,:,e);
45  a = z(:,:,e)*diag(exp(-2*hyp(e).l)/2);
46  dnlml.z(:,:,e) = D*exp(-2*hyp(e).pn) + G*a - bsxfun(@times,sum(G,2),a);
47  B = bsxfun(@minus,permute(a,[1 3 2]),permute(a,[3 1 2]));
48  dnlml.hyp(e).l = H * exp(-2*hyp(e).pn)  ...
49          + exp(2*hyp(e).l).*squeeze(sum(sum(bsxfun(@times,G,B.^2),1),2));
50  G = (R/(K(:,:,e)+Psi2(:,:,e))*R' + w*w').*K(:,:,e);
51  a = z(:,:,e)*diag(exp(-2*hyp(e).l));
52  B = bsxfun(@minus,permute(z(:,:,e),[1 3 2]),permute(z(:,:,e),[3 1 2]));
53  dnlml.hyp(e).l = dnlml.hyp(e).l ...
54        + exp(-2*hyp(e).l).*squeeze(sum(sum(bsxfun(@times,B.^2,G),1),2))/2;
55  dnlml.z(:,:,e) = dnlml.z(:,:,e) + G*a - bsxfun(@times,sum(G,2),a);
56  end
```

## Test set calculation

The distinguishing factor between training and test set calculations is whether the inducing target distribution is updated (training set) or kept fixed (test set). For the test set the contribution from the transition model to the log probability is

$$
\begin{aligned}
\sum_{e=1}^{E} & \tfrac{1}{2}\mathrm{tr}K_e^{-1}\Psi_{2e}(K_e+\Psi_{2e})^{-1}\Psi_{2e}^* - \tfrac{1}{2}\mathrm{tr}(K_e+\Psi_{2e})^{-1}\Psi_{1e}\Psi_{1e}^\top(K_e+\Psi_{2e})^{-1}\Psi_{2e}^* + \Psi_{1e}^\top(K_e+\Psi_{2e})^{-1}\Psi_{1e}^* \\
& + \tfrac{1}{2}\sum_{t=2}^{T}\log|\Sigma_{t-1:t,t-1:t}| - \tfrac{1}{2}\sum_{t=2}^{T-1}\log|\Sigma_t| - \tfrac{1}{2}\mathrm{tr}Q^{-1}\sum_{t=2}^{T}(I+\mu_t^\top\mu_t+\Sigma_{t,t}) - \tfrac{T-1}{2}\log|Q| - \tfrac{(T-1)E}{2}\log(2\pi).
\end{aligned}
\tag{23}
$$

The testing consists of two steps. From a supplied initial set of observations, we find the most likely latent states, maximising a q distribution while keeping the hyperparameters and inducing inputs constant, using eq. (23), along with the likelihood term.

After the most likely latent state has been found, we predict forward by adding a further point to the timeseries, and maximising the likelihood. This can be solved analytically..

## Representation of the q(x) distribution

The $q(\mathbf{x})$ distribution is parameterised through its mean qx.m and the marginal and pairwise covariances. Conceptually, we wish to parameterize the E by E covariance matrices (for the marginal distibutions) which we call qx.Sd (for diagonal) and the E by E covariances between consequtive time points (for the pairwise marginals) which we call qx.So (for off-diagonal). However it is inconvenient to parametrise these matrices directly, as it would be difficult to ensure positive definiteness of the marginal and pairwise marginal covariance matrices. Instead, we use 2E by E representation qx.s such that

$$
S_{d,t} = s_t^\top s_t, \text{ and } S_{o,t-1} = s_{t-1}^\top s_t.
\tag{24}
$$

Using this representation, we can use call the optimizer with the unconstrained representation, which is the converted to the more convenient diagonal and off-diagonal representation at the beginning and the derivatives are reverted back at the end.

9a  ⟨*convert* 9a⟩≡ (5b)

```
1 function [Sd, So] = convert(s)
2 [t, E, T]= size(s); Sd = zeros(E,E,T); So = zeros(E,E,T-1);
3 for t = 1:T, Sd(:,:,t) = s(:,:,t)'*s(:,:,t); end          % diagonal terms
4 for t = 2:T, So(:,:,t-1) = s(:,:,t-1)'*s(:,:,t); end      % off-diagonal
```

9b  ⟨*revert covariances* 9b⟩≡ (5b)

```
1 out2.qx = rmfield(out2.qx,{'Sd','So'});      % change to qx.s representation
2 [out2.qx.s] = deal([]);                      % create the "s" field
3 %for n = 1:N
4 %  Sd = sum(dnlml.qx(n).Sd,3) / size(dnlml.qx(n).Sd,3);
5 %  for t = 1:size(dnlml.qx(n).Sd,3), dnlml.qx(n).Sd(:,:,t) = Sd; end
6 %end
7 for n = 1:N
8   out2.qx(n).s = revert(p.qx(n).s, dnlml.qx(n).Sd, dnlml.qx(n).So);
9 end
```

The derivatives are

$$
\begin{aligned}
\frac{\partial\mathcal{L}}{\partial s_t} &= \frac{\partial\mathcal{L}}{\partial S_{d,t}}\frac{\partial S_{d,t}}{\partial s_t} + \frac{\partial\mathcal{L}}{\partial S_{o,t-1}}\frac{\partial S_{o,t-1}}{\partial s_t} + \frac{\partial\mathcal{L}}{\partial S_{o,t}}\frac{\partial S_{o,t}}{\partial s_t} = \frac{\partial}{\partial s_t}\mathrm{tr}\Big(\frac{\partial\mathcal{L}}{\partial S_{d,t}}s_t^\top s_t\Big) \\
&+ s_{t-1}\frac{\partial\mathcal{L}}{\partial S_{o,t-1}} + s_t\Big[\frac{\partial\mathcal{L}}{\partial S_{o,t}}\Big]^\top = s_t\Big(\frac{\partial\mathcal{L}}{\partial S_{d,t}} + \Big[\frac{\partial\mathcal{L}}{\partial S_{d,t}}\Big]^\top\Big) + s_{t-1}\frac{\partial\mathcal{L}}{\partial S_{o,t-1}} + s_t\Big[\frac{\partial\mathcal{L}}{\partial S_{o,t}}\Big]^\top.
\end{aligned}
\tag{25}
$$

9c  ⟨*revert* 9c⟩≡ (5b)

```
1 function r = revert(s, dSd, dSo)
2 for t = 1:size(s,3), r(:,:,t) = s(:,:,t)*(dSd(:,:,t)+dSd(:,:,t)'); end
3 for t = 2:size(s,3)
4   r(:,:,t-1) = r(:,:,t-1) + s(:,:,t)*dSo(:,:,t-1)';
5   r(:,:,t) = r(:,:,t) + s(:,:,t-1)*dSo(:,:,t-1);
6 end
```

⟨*maha* 10a⟩≡ (5b)

```
1 % Squared Mahalanobis distance (a-b)*Q*(a-b)'; vectors are row-vectors
2 % a, b   d x n   matrices containing n length d row vectors
3 % Q       d x d   weight matrix
4 % K       n x n   squared distances
5 function K = maha(a, b, Q)
6 if isempty(b), b = a; end
7 aQ = a*Q; K = bsxfun(@plus,sum(aQ.*a,2),sum(b*Q.*b,2)')-2*aQ*b';
```

## Analytic Prediction

Predicting the next step is found by taking an existing latent representation of a timeseries, and maximising the likelihood of the timeseries found by increasing its length by one. Given a timeseries, if increase the length of the timeseries by one, we can analytically compute the parameters $\mu_t$, $\Sigma_{t-1,t}$, $\Sigma_{t,t}$ which maximise the likelihood. From equation **??**, $\frac{\partial \mathcal{L}}{\partial \mu_t} = \Psi_{1e}^\top (K_e + \Psi_{2e}^*)^{-1} \frac{\partial \Psi_{1e}^*}{\partial \mu_t} - \text{tr}(Q^{-1})\mu_t$, and using eq. 21,

$$\mu_t^* = \Psi_{1e}^\top (K_e + \Psi_{2e})^{-1/2} |I + \Lambda_e^{-1}\Sigma_{t-1,t-1}|^{-1/2} \exp\left(-\frac{1}{2}(\mu_{t-1} - z_{ie})[\Lambda_e + \Sigma_{t-1,t-1}]^{-1}(\mu_{t-1} - z_{ie})\right), \quad (26)$$

which corresponds to the mean prediction of the GP, conditioned on the inducing inputs. For the marginal $\Sigma_{t,t}$, we find that $\frac{\partial \mathcal{L}}{\partial \Sigma_{t,t}} = -\frac{1}{2}\text{tr}Q^{-1} + \frac{1}{2}\frac{\partial}{\partial \Sigma_{t,t}}|\log \Sigma_{t-1:t,t-1:t}|$. Optimising this gives us

$$\Sigma_{t,t}^* = Q + \Sigma_{t-1,t}^* \Sigma_{t-1,t-1}^{-1} \Sigma_{t,t-1}^* \quad (27)$$

For the pair-wise marginal, $\frac{\partial \mathcal{L}}{\partial \Sigma_{t-1,t}} = \Psi_{1e}^\top (K_e + \Psi_{2e})^{-1} \frac{\partial \Psi_{1e}^*}{\partial \Sigma_{t-1,t}} + \frac{1}{2}\frac{\partial}{\partial \Sigma_{t-1,t}}|\log \Sigma_{t-1:t,t-1:t}|$. This gives us

$$0 = \Psi_{1e}^\top (K_e + \Psi_{2e})^{-1} \frac{\partial \Psi_{1e}^*}{\partial \Sigma_{t-1,t}} - \Sigma_{t-1,t-1}^{-1}\Sigma_{t-1,t}^* \left(\Sigma_{t,t}^* - \Sigma_{t-1,t}^*\Sigma_{t-1,t-1}\Sigma_{t-1,t}^*\right)^{-1} \quad (28)$$

.

Substituting, we find that

$$\begin{aligned}\Sigma_{t-1,t}^* =&\Sigma_{t-1,t-1}\Psi_{1e}^\top (K_e + \Psi_{2e})^{-1} (\Lambda_e + \Sigma_{t-1,t-1})^{-1} (z - \mu_{t-1}) \\ &\times |I + \Lambda_e^{-1}\Sigma_{t-1,t-1}|^{-1/2} \exp\left(-\frac{1}{2}(\mu_{t-1} - z_{ie})[\Lambda_e + \Sigma_{t-1,t-1}]^{-1}(\mu_{t-1} - z_{ie})\right),\end{aligned} \quad (29)$$

which we can then substitute in to equation 27 to find $\Sigma_{t,t}^*$.

We can motivate our choice of $\Sigma_{t,t}^*$ by considering the prediction of the value of $x_t$ conditional on the value of $x_{t-1}$. We have a conditional distribution $q(x_t|x_{t-1}) = \mathcal{N}(\mu', \Sigma')$, with $\Sigma' = \Sigma_{t,t}^* - \Sigma_{t-1,t}^{\top *}\Sigma_{t-1,t-1}^{-1}\Sigma_{t-1,t-1}^* = Q$. We can clearly see how the derived $\Sigma_{t,t}^*$ is exactly the variance for which the extra uncertainty added to the prediction of the next timestep is Q.

The intuition for equation 29 is that we are directly calculating the covariance $\mathbb{E}((z_i - \mu_{t-1})(z_j - \mu_t)) = \mathbb{E}(z_i(z_j - \mu_{t-1}))$, from the inducing points $z$. We rescale the individual contributions by $(\Lambda_e + \Sigma_{t-1,t-1})^{-1}$, and then scale up by a factor of $\Sigma_{t-1,t-1}$ after collecting the contributions from all $z$.

⟨*prediction* 10b⟩≡ (5b)

```
1 function qx_opt = predict(hyp, qx, z, u)
2 %Given a timeseries corresponding to latent states of an observed
3 %timeseries, analytically calculate the optimal next step params
4 %Fetch the Psi1,2 values by calling Psi with 4 arguments
5 [Psi1, Psi2] = Psi(hyp, qx, z, u);
6 E = size(qx, 2); [M, F, E] = size(z);
7
8 for e = 1:E
9   K(:,:,e) = exp(-maha(z(:,:,e),[],diag(exp(-2*hyp(e).l)))/2) + ...
10       1e-6*eye(M);
11   w = (K(:,:,e)+Psi2(:,:,e))\Psi1(:,e);
12   iL = diag(exp(-hyp(e).l)); L2 = diag(exp(2*hyp(e).l));
```

```matlab
13    b1 = zeros(M,1);
14    for n = 1:size(qx,2)
15      t = size(qx.m, 2)
16      Sd(1:E,1:E) = qx(n).Sd(:,:,t);              % covariance in top left corner
17      r1 = prod(diag(chol(eye(F)+iL*Sd*iL)));                    % sqrt det
18      %  s = bsxfun(@minus, z(:,:,e), [qx(n).m(:,t-1)'
19      %                               u{n}(t-1,:)]);
20      s = bsxfun(@minus, z(:,:,e), [qx(n).m(:,t-1)' u{n}(t,:)]);
21      a = s/(L2+Sd);
22      b3 = exp(-sum(a.*s,2)/2)/r1;
23      qx_opt(n).m = b3'*w;
24      qx_opt(n).So = (qx(n).Sd(:,e,t)*w)'*(a(:,1:E).*b3);
25      qx_opt(n).Sd = exp(2*hyp(e).pn) + qx_opt(n).So * inv(qx(n).Sd(:, e, end)) ...
26      * qx_opt(n).So';
27    end
28 end
```