

Probabilistic & Unsupervised Learning

Expectation Propagation

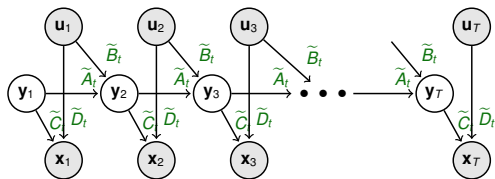
Maneesh Sahani

maneesh@gatsby.ucl.ac.uk

Gatsby Computational Neuroscience Unit, and
MSc ML/CSML, Dept Computer Science
University College London

Term 1, Autumn 2013

Nonlinear state-space model (NLSSM)



$$\mathbf{y}_{t+1} = f(\mathbf{y}_t, \mathbf{u}_t) + \mathbf{w}_t$$

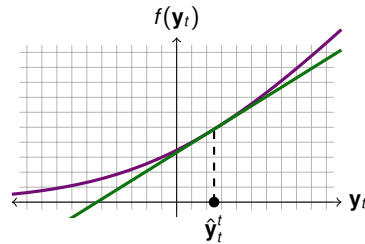
$$\mathbf{x}_t = g(\mathbf{y}_t, \mathbf{u}_t) + \mathbf{v}_t$$

$\mathbf{w}_t, \mathbf{v}_t$ usually still Gaussian.

Extended Kalman Filter (EKF): linearise nonlinear functions about current estimate, $\hat{\mathbf{y}}_t^i$:

$$\mathbf{y}_{t+1} \approx \underbrace{f(\hat{\mathbf{y}}_t^i, \mathbf{u}_t)}_{\tilde{B}_t \mathbf{u}_t} + \underbrace{\frac{\partial f}{\partial \mathbf{y}_t}}_{\tilde{A}_t} (\mathbf{y}_t - \hat{\mathbf{y}}_t^i) + \mathbf{w}_t$$

$$\mathbf{x}_t \approx \underbrace{g(\hat{\mathbf{y}}_t^i, \mathbf{u}_t)}_{\tilde{D}_t \mathbf{u}_t} + \underbrace{\frac{\partial g}{\partial \mathbf{y}_t}}_{\tilde{C}_t} (\mathbf{y}_t - \hat{\mathbf{y}}_t^i) + \mathbf{v}_t$$



Run the Kalman filter (smoother) on non-stationary linearised system $(\tilde{A}_t, \tilde{B}_t, \tilde{C}_t, \tilde{D}_t)$:

- ▶ Adaptively approximates non-Gaussian messages by Gaussians.
- ▶ Local linearisation depends on central point of distribution \Rightarrow approximation degrades with increased state uncertainty. May work acceptably for close-to-linear systems.

Can base EM-like algorithm on EKF/EKS (or alternatives).

Intractabilities and approximations

- ▶ Inference – computational intractability
 - ▶ Factored variational approx
 - ▶ Loopy BP/EP/Power EP
 - ▶ Gibbs sampling, other MCMC
- ▶ Inference – analytic intractability
 - ▶ Laplace approximation (global)
 - ▶ Parametric variational approx (for special cases).
 - ▶ Message approximations (linearised, sigma-point, Laplace)
 - ▶ Assumed-density methods and Expectation-Propagation
 - ▶ (Sequential) Monte-Carlo methods
- ▶ Learning – intractable partition function
 - ▶ Contrastive divergence
 - ▶ Sampling parameters
 - ▶ Score-matching
- ▶ Model selection
 - ▶ Laplace approximation / BIC
 - ▶ Variational Bayes
 - ▶ (Annealed) importance sampling
 - ▶ Reversible jump MCMC

Not a complete list!

Online learning with NLSSM message passing

Nonlinear message passing can also be used to implement online parameter learning in (non)linear latent state-space systems:

Eg: for linear model, augment state vector to include the model parameters and introduce **nonlinear** transition \bar{f} :

$$\bar{\mathbf{y}}_t = \begin{bmatrix} \mathbf{y}_t \\ A \\ C \end{bmatrix}$$

$$\bar{\mathbf{y}}_{t+1} = \bar{f}(\bar{\mathbf{y}}_t) + \begin{bmatrix} \mathbf{w}_t \\ 0 \\ 0 \end{bmatrix} \quad \text{with } \bar{f}(\bar{\mathbf{y}}) = \bar{f} \left(\begin{bmatrix} \mathbf{y} \\ A \\ C \end{bmatrix} \right) = \begin{bmatrix} A\mathbf{y} \\ A \\ C \end{bmatrix}$$

Use **EKF** or alternative to compute online estimates of $E[\bar{\mathbf{y}}_t | \mathbf{x}_1, \dots, \mathbf{x}_t]$ and $\text{Cov}[\bar{\mathbf{y}}_t | \mathbf{x}_1, \dots, \mathbf{x}_t]$. These now include mean and posterior variance of parameter estimates.

- ▶ Pseudo-Bayesian approach: gives Gaussian distributions over parameters.
- ▶ Can model nonstationarity by assuming non-zero innovations noise in A, C .
- ▶ Not simple to implement for Q and R (e.g. covariance constraints?).
- ▶ May be faster than EM/gradients approaches.

Sometimes called the **joint-EKF** approach.

Other message approximations

Consider the forward messages on a latent chain:

$$P(y_t | x_{1:t}) = \frac{1}{Z} P(x_t | y_t) \int dy_{t-1} P(y_t | y_{t-1}) P(y_{t-1} | x_{1:t-1})$$

We want to approximate the messages to retain a tractable form (i.e. Gaussian).

$$\tilde{P}(y_t | x_{1:t}) \approx \frac{1}{Z} P(x_t | y_t) \int dy_{t-1} \underbrace{P(y_t | y_{t-1})}_{\mathcal{N}(f(\hat{y}_{t-1}), Q)} \underbrace{\tilde{P}(y_{t-1} | x_{1:t-1})}_{\mathcal{N}(\hat{y}_{t-1}, V_{t-1})}$$

- ▶ Linearisation at the peak (EKF) is only one approach.
- ▶ Laplace filter: use mode and curvature of integrand.
- ▶ Sigma-point (“unscented”) filter:
 - ▶ Evaluate $f(\hat{y}_{t-1}), f(\hat{y}_{t-1} \pm \lambda(V)_{t-1})$ for eigenvectors λ .
 - ▶ “Fit” Gaussian to these $2K + 1$ points.
 - ▶ Equivalent to numerical evaluation of mean and covariance by Gaussian quadrature.
 - ▶ One form of “Assumed Density Filtering” and EP.
- ▶ Parametric variational: $\operatorname{argmin} \mathbf{KL}[\mathcal{N}(\hat{y}_t, \hat{V}_t) \parallel \int dy_{t-1} \dots]$. Requires Gaussian expectations of $\log \int \Rightarrow$ may be challenging.
- ▶ The other KL: $\operatorname{argmin} \mathbf{KL}[\int dy_{t-1} \parallel \mathcal{N}(\hat{y}_t, \hat{V}_t)]$ needs only first and second moments of nonlinear message \Rightarrow EP.

Approximating the posterior

Linearisation (or local Laplace, sigma-point and other such approaches) seem *ad hoc*. A more principled approach might look for an approximate q that is **closest** to P in some sense.

$$q = \operatorname{argmin}_{q \in \mathcal{Q}} D(P \leftrightarrow q)$$

Open choices:

- ▶ form of the metric D
- ▶ nature of the constraint space \mathcal{Q}
- ▶ Variational methods: $D = \mathbf{KL}[q \parallel P]$.
 - ▶ Choosing $\mathcal{Q} = \{\text{tree-factored distributions}\}$ leads to efficient message passing.
- ▶ Can we use other divergences?

Variational learning

Free energy:

$$\mathcal{F}(q, \theta) = \langle \log P(\mathcal{X}, \mathcal{Y} | \theta) \rangle_{q(\mathcal{Y} | \mathcal{X})} + \mathbf{H}[q] = \log P(\mathcal{X} | \theta) - \mathbf{KL}[q(\mathcal{Y}) \parallel P(\mathcal{Y} | \mathcal{X}, \theta)] \leq \ell(\theta)$$

E-steps:

- ▶ Exact EM: $q(\mathcal{Y}) = \operatorname{argmax}_q \mathcal{F} = P(\mathcal{Y} | \mathcal{X}, \theta)$
 - ▶ Saturates bound: converges to local maximum of likelihood.
- ▶ (Factored) variational approximation:

$$q(\mathcal{Y}) = \operatorname{argmax}_{q_1(\mathcal{Y}_1)q_2(\mathcal{Y}_2)} \mathcal{F} = \operatorname{argmin}_{q_1(\mathcal{Y}_1)q_2(\mathcal{Y}_2)} \mathbf{KL}[q_1(\mathcal{Y}_1)q_2(\mathcal{Y}_2) \parallel P(\mathcal{Y} | \mathcal{X}, \theta)]$$
 - ▶ Increases bound: converges, but not necessarily to ML.
- ▶ Other approximations: $q(\mathcal{Y}) \approx P(\mathcal{Y} | \mathcal{X}, \theta)$
 - ▶ Usually no guarantees, but if learning converges it is frequently more accurate than the factored approximation

The other KL

What about the ‘other’ KL ($q = \operatorname{argmin} \mathbf{KL}[P \parallel q]$)?

For a factored approximation the (clique) marginals obtained by minimising this KL are correct:

$$\begin{aligned} \operatorname{argmin}_{q_i} \mathbf{KL} \left[P(\mathcal{Y} | \mathcal{X}) \parallel \prod q_i(\mathcal{Y}_i | \mathcal{X}) \right] &= \operatorname{argmin}_{q_i} - \int d\mathcal{Y} P(\mathcal{Y} | \mathcal{X}) \log \prod q_i(\mathcal{Y}_i | \mathcal{X}) \\ &= \operatorname{argmin}_{q_i} - \sum_j \int d\mathcal{Y} P(\mathcal{Y} | \mathcal{X}) \log q_j(\mathcal{Y}_j | \mathcal{X}) \\ &= \operatorname{argmin}_{q_i} - \int d\mathcal{Y}_i P(\mathcal{Y}_i | \mathcal{X}) \log q_i(\mathcal{Y}_i | \mathcal{X}) \\ &= P(\mathcal{Y}_i | \mathcal{X}) \end{aligned}$$

and the marginals are what we need for learning (although if factored over disjoint sets as in the variational approximation some cliques will be missing).

Perversely, this means finding the best q for this KL is intractable!

But it raises the hope that **approximate** minimisation might still yield useful results.

Approximate optimisation

The posterior distribution in a graphical model is a (normalised) product of factors:

$$P(\mathcal{Y}|\mathcal{X}) = \frac{P(\mathcal{Y}, \mathcal{X})}{P(\mathcal{X})} = \frac{1}{Z} \prod_i P(Y_i | \text{pa}(Y_i)) \propto \prod_{i=1}^N f_i(\mathcal{Y}_i)$$

where the \mathcal{Y}_i are not necessarily disjoint. In the language of EP the f_i are called **sites**.

Consider q with the **same** factorisation, but potentially approximated sites: $q(\mathcal{Y}) \stackrel{\text{def}}{=} \prod_{i=1}^N \tilde{f}_i(\mathcal{Y}_i)$.

We would like to minimise (at least in some sense) $\text{KL}[P||q]$.

Possible optimisations:

$$\begin{aligned} \min_{q(\mathcal{Y})} \text{KL} \left[\prod_{i=1}^N f_i(\mathcal{Y}_i) \middle\| \prod_{i=1}^N \tilde{f}_i(\mathcal{Y}_i) \right] & \quad (\text{global: intractable}) \\ \min_{\tilde{f}_i} \text{KL} \left[f_i(\mathcal{Y}_i) \middle\| \tilde{f}_i(\mathcal{Y}_i) \right] & \quad (\text{local, fixed: simple, inaccurate}) \\ \min_{\tilde{f}_i} \text{KL} \left[f_i(\mathcal{Y}_i) \prod_{j \neq i} \tilde{f}_j(\mathcal{Y}_j) \middle\| \tilde{f}_i(\mathcal{Y}_i) \prod_{j \neq i} \tilde{f}_j(\mathcal{Y}_j) \right] & \quad (\text{local, contextual: iterative, accurate}) \leftarrow \text{EP} \end{aligned}$$

Local updates

Each EP update involves a KL minimisation:

$$\tilde{f}_i^{\text{new}}(\mathcal{Y}) \leftarrow \underset{f \in \{\tilde{f}\}}{\text{argmin}} \text{KL}[f(\mathcal{Y}_i) q_{-i}(\mathcal{Y}) || f(\mathcal{Y}_i) q_{-i}(\mathcal{Y})]$$

Write $q_{-i}(\mathcal{Y}) = q_{-i}(\mathcal{Y}_i) q_{-i}(\mathcal{Y}_{-i} | \mathcal{Y}_i)$. Then:

$$\begin{aligned} \min_{\tilde{f}_i} \text{KL}[f_i(\mathcal{Y}_i) q_{-i}(\mathcal{Y}) || f(\mathcal{Y}_i) q_{-i}(\mathcal{Y})] \\ &= \max_{\tilde{f}_i} \int d\mathcal{Y}_i d\mathcal{Y}_{-i} f_i(\mathcal{Y}_i) q_{-i}(\mathcal{Y}) \log f(\mathcal{Y}_i) q_{-i}(\mathcal{Y}) \\ &= \max_{\tilde{f}_i} \int d\mathcal{Y}_i d\mathcal{Y}_{-i} f_i(\mathcal{Y}_i) q_{-i}(\mathcal{Y}_i) q_{-i}(\mathcal{Y}_{-i} | \mathcal{Y}_i) (\log f(\mathcal{Y}_i) q_{-i}(\mathcal{Y}_i) + \log q_{-i}(\mathcal{Y}_{-i} | \mathcal{Y}_i)) \\ &= \max_{\tilde{f}_i} \int d\mathcal{Y}_i f_i(\mathcal{Y}_i) q_{-i}(\mathcal{Y}_i) (\log f(\mathcal{Y}_i) q_{-i}(\mathcal{Y}_i)) \int d\mathcal{Y}_{-i} q_{-i}(\mathcal{Y}_{-i} | \mathcal{Y}_i) \\ &= \min_{\tilde{f}_i} \text{KL}[f_i(\mathcal{Y}_i) q_{-i}(\mathcal{Y}_i) || f(\mathcal{Y}_i) q_{-i}(\mathcal{Y}_i)] \end{aligned}$$

$q_{-i}(\mathcal{Y}_i)$ is sometimes called the **cavity distribution**.

Expectation? Propagation?

EP is really two ideas:

- ▶ **Approximation** of factors.
 - ▶ Usually by “projection” to exponential families.
 - ▶ This involves finding expected sufficient statistics, hence **expectation**.
- ▶ **Local** divergence minimization in the context of other factors.
 - ▶ This leads to a message passing approach, hence **propagation**.

Expectation Propagation (EP)

Input $f_1(\mathcal{Y}_1) \dots f_N(\mathcal{Y}_N)$

Initialize $\tilde{f}_i(\mathcal{Y}_i) = \underset{f \in \{\tilde{f}\}}{\text{argmin}} \text{KL}[f_1(\mathcal{Y}_1) || f_1(\mathcal{Y}_1)]$, $\tilde{f}_i(\mathcal{Y}_i) = 1$ for $i > 1$, $q(\mathcal{Y}) \propto \prod_i \tilde{f}_i(\mathcal{Y}_i)$

repeat

for $i = 1 \dots N$ **do**

Delete: $q_{-i}(\mathcal{Y}) \leftarrow \frac{q(\mathcal{Y})}{\tilde{f}_i(\mathcal{Y}_i)} = \prod_{j \neq i} \tilde{f}_j(\mathcal{Y}_j)$

Project: $\tilde{f}_i^{\text{new}}(\mathcal{Y}) \leftarrow \underset{f \in \{\tilde{f}\}}{\text{argmin}} \text{KL}[f_1(\mathcal{Y}_1) q_{-i}(\mathcal{Y}) || f(\mathcal{Y}_i) q_{-i}(\mathcal{Y})]$

Include: $q(\mathcal{Y}) \leftarrow \tilde{f}_i^{\text{new}}(\mathcal{Y}_i) q_{-i}(\mathcal{Y})$

end for

until convergence

Message Passing

- ▶ The cavity distribution (in a tree) can be further broken down into a product of terms from each neighbouring clique:

$$q_{-i}(\mathcal{Y}_i) = \prod_{j \in \text{ne}(i)} M_{j \rightarrow i}(\mathcal{Y}_j \cap \mathcal{Y}_i)$$

- ▶ Once the i th site has been approximated, the messages can be passed on to neighbouring cliques by marginalising to the shared variables (SSM example follows).
⇒ belief propagation.

- ▶ In loopy graphs, we can use **loopy belief propagation**. In that case

$$q_{-i}(\mathcal{Y}_i) = \prod_{j \in \text{ne}(i)} M_{j \rightarrow i}(\mathcal{Y}_j \cap \mathcal{Y}_i)$$

becomes an approximation to the **true** cavity distribution (or we can recast the approximation directly in terms of messages ⇒ later lecture).

- ▶ For some approximations (e.g. Gaussian) may be able to compute true loopy cavity using approximate sites, even if computing exact message would have been intractable.
- ▶ In either case, message updates can be scheduled in any order.
- ▶ No guarantee of convergence (but see “power-EP” methods).

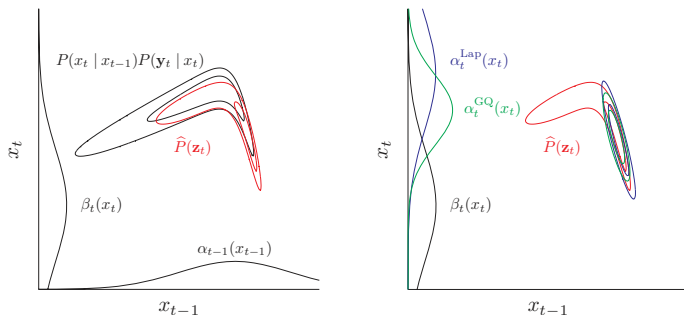
NLSSM EP message updates

$$\tilde{f}_i(\mathbf{y}_i, \mathbf{y}_{i-1}) = \underset{f \in \mathcal{N}}{\text{argmin KL}} \left[\underbrace{\phi_i(\mathbf{y}_i, \mathbf{y}_{i-1}) \psi_i(\mathbf{y}_i) \alpha_{i-1}(\mathbf{y}_{i-1}) \beta_i(\mathbf{y}_i)}_{\hat{P}(\mathbf{y}_{i-1}, \mathbf{y}_i)} \middle\| \underbrace{f(\mathbf{y}_i, \mathbf{y}_{i-1}) \alpha_{i-1}(\mathbf{y}_{i-1}) \beta_i(\mathbf{y}_i)}_{\tilde{P}(\mathbf{y}_{i-1}, \mathbf{y}_i)} \right]$$

$$\alpha_i(\mathbf{y}_i) = \tilde{P}(\mathbf{y}_{i-1}, \mathbf{y}_i) / \beta_i(\mathbf{y}_i)$$

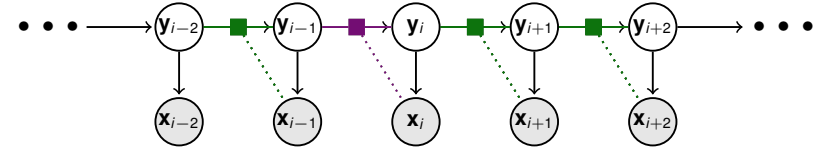
$$\beta_{i-1}(\mathbf{y}_{i-1}) = \tilde{P}(\mathbf{y}_{i-1}, \mathbf{y}_i) / \alpha_{i-1}(\mathbf{y}_{i-1})$$

For exponential family approximations, division ⇒ subtraction of natural parameters.



[Note: For this figure (from a paper): x is latent, \mathbf{y} observed and $\mathbf{z}_t = [x_{t-1}, x_t]$.]

EP for a NLSSM



$$P(\mathbf{y}_i | \mathbf{y}_{i-1}) = \phi_i(\mathbf{y}_i, \mathbf{y}_{i-1}) \quad \text{e.g. } \exp(-\|\mathbf{y}_i - h_s(\mathbf{y}_{i-1})\|^2 / 2\sigma^2)$$

$$P(\mathbf{x}_i | \mathbf{y}_i) = \psi_i(\mathbf{y}_i) \quad \text{e.g. } \exp(-\|\mathbf{x}_i - h_o(\mathbf{y}_i)\|^2 / 2\sigma^2)$$

Then $f_i(\mathbf{y}_i, \mathbf{y}_{i-1}) = \phi_i(\mathbf{y}_i, \mathbf{y}_{i-1}) \psi_i(\mathbf{y}_i)$. As ϕ_i and ψ_i are non-linear, inference is not generally tractable.

Assume $\tilde{f}_i(\mathbf{y}_i, \mathbf{y}_{i-1})$ is Gaussian. Then,

$$q_{-i}(\mathbf{y}_i, \mathbf{y}_{i-1}) = \sum_{\substack{\mathbf{y}_1 \dots \mathbf{y}_{i-2} \\ \mathbf{y}_{i+1} \dots \mathbf{y}_i}} \prod_{i' \neq i} \tilde{f}_{i'}(\mathbf{y}_{i'}, \mathbf{y}_{i'-1}) = \underbrace{\sum_{\mathbf{y}_1 \dots \mathbf{y}_{i-2}} \prod_{i' < i} \tilde{f}_{i'}(\mathbf{y}_{i'}, \mathbf{y}_{i'-1})}_{\alpha_{i-1}(\mathbf{y}_{i-1})} \underbrace{\sum_{\mathbf{y}_{i+1} \dots \mathbf{y}_i} \prod_{i' > i} \tilde{f}_{i'}(\mathbf{y}_{i'}, \mathbf{y}_{i'-1})}_{\beta_i(\mathbf{y}_i)}$$

with both α and β Gaussian.

$$\tilde{f}_i(\mathbf{y}_i, \mathbf{y}_{i-1}) = \underset{f \in \mathcal{N}}{\text{argmin KL}} \left[\phi_i(\mathbf{y}_i, \mathbf{y}_{i-1}) \psi_i(\mathbf{y}_i) \alpha_{i-1}(\mathbf{y}_{i-1}) \beta_i(\mathbf{y}_i) \middle\| f(\mathbf{y}_i, \mathbf{y}_{i-1}) \alpha_{i-1}(\mathbf{y}_{i-1}) \beta_i(\mathbf{y}_i) \right]$$

Moment Matching

Each EP update involves a KL minimisation:

$$\tilde{f}_i^{\text{new}}(\mathcal{Y}) \leftarrow \underset{f \in \{\}}{\text{argmin KL}} \left[f_i(\mathcal{Y}_i) q_{-i}(\mathcal{Y}) \middle\| f(\mathcal{Y}_i) q_{-i}(\mathcal{Y}) \right]$$

Usually, both $q_{-i}(\mathcal{Y}_i)$ and \tilde{f} are in the same exponential family. Let $q(x) = \frac{1}{Z(\theta)} e^{\mathbf{S}(x) \cdot \theta}$. Then

$$\begin{aligned} \underset{q}{\text{argmin KL}} [p(x) \middle\| q(x)] &= \underset{\theta}{\text{argmin KL}} \left[p(x) \middle\| \frac{1}{Z(\theta)} e^{\mathbf{T}(x) \cdot \theta} \right] \\ &= \underset{\theta}{\text{argmin}} - \int dx p(x) \log \frac{1}{Z(\theta)} e^{\mathbf{T}(x) \cdot \theta} \\ &= \underset{\theta}{\text{argmin}} - \int dx p(x) \mathbf{T}(x) \cdot \theta + \log Z(\theta) \\ \frac{\partial}{\partial \theta} &= - \int dx p(x) \mathbf{T}(x) + \frac{1}{Z(\theta)} \frac{\partial}{\partial \theta} \int dx e^{\mathbf{T}(x) \cdot \theta} \\ &= -\langle \mathbf{T}(x) \rangle_p + \frac{1}{Z(\theta)} \int dx e^{\mathbf{T}(x) \cdot \theta} \mathbf{T}(x) \\ &= -\langle \mathbf{T}(x) \rangle_p + \langle \mathbf{T}(x) \rangle_q \end{aligned}$$

So minimum is found by **matching sufficient stats**. This is usually **moment matching**.

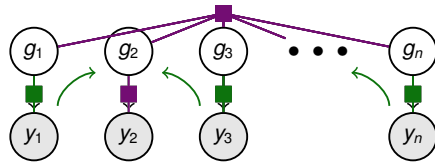
Numerical issues

How do we calculate $\langle T(x) \rangle_\rho$?

Often analytically tractable, but even if not requires a (relatively) low-dimensional integral:

- ▶ **Quadrature methods.**
 - ▶ Classical Gaussian quadrature (same Gauss, but nothing to do with the distribution) gives an iterative version of Sigma-point methods.
 - ▶ Positive definite joints, but not guaranteed to give positive definite messages.
 - ▶ Heuristics include skipping non-positive-definite steps, or damping messages by interpolation or exponentiating to power < 1 .
 - ▶ Other quadrature approaches (e.g. GP quadrature) may be more accurate, and may allow formal constraint to PD cone.
- ▶ **Laplace approximation.**
 - ▶ Equivalent to Laplace propagation.
 - ▶ As long as messages remain positive definite will converge to global Laplace approximation.

GP EP updates



- ▶ We can write the GP joint on g_i and y_i as a factor graph:

$$P(g_1 \dots g_n, y_1, \dots, y_n) = \underbrace{\mathcal{N}(g_1 \dots g_n | \mathbf{0}, K)}_{f_0(g)} \prod_i \underbrace{\mathcal{N}(y_i | g_i, \sigma_i^2)}_{f_i(g_i)}$$

- ▶ The same factorisation applies to non-Gaussian $P(y_i | g_i)$ (e.g. $P(y_i=1) = 1/(1 + e^{-g_i})$).
- ▶ EP: approximate non-Gaussian $f_i(g_i)$ by Gaussian $\tilde{f}_i(g_i) = \mathcal{N}(\tilde{\mu}_i, \tilde{\psi}_i^2)$.
- ▶ $q_{-i}(g_i)$ can be constructed by the usual GP marginalisation. If $\Sigma = K + \text{diag}[\tilde{\psi}_1^2 \dots \tilde{\psi}_n^2]$

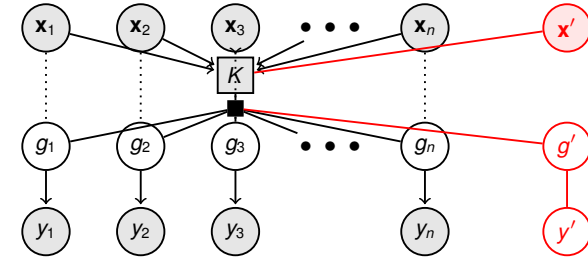
$$q_{-i}(g_i) = \mathcal{N}(\Sigma_{-i,-i} \Sigma_{-i,-i}^{-1} \tilde{\mu}_{-i}, K_{i,i} - \Sigma_{-i,-i} \Sigma_{-i,-i}^{-1} \Sigma_{-i,i})$$

- ▶ The EP updates thus require calculating Gaussian expectations of $f_i(g)g^{1,2}$:

$$\tilde{f}_i^{\text{new}}(g_i) = \mathcal{N}\left(\int dg q_{-i}(g) f_i(g) g, \int dg q_{-i}(g) f_i(g) g^2 - (\tilde{\mu}_i^{\text{new}})^2\right)$$

EP for Gaussian process classification

EP provides the most successful framework for Gaussian-process modelling of non-Gaussian observations (e.g. for classification).

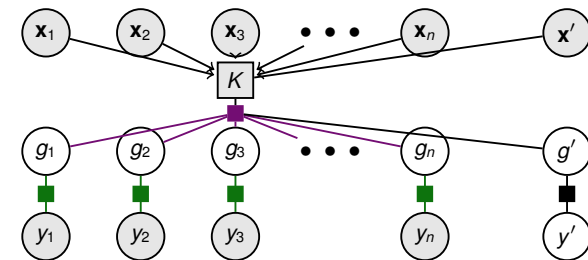


Recall:

- ▶ A GP defines a **multivariate Gaussian** distribution on any finite subset of random vars $\{g_1 \dots g_n\}$ drawn from a (usually uncountable) potential set indexed by "inputs" \mathbf{x}_i .
- ▶ The Gaussian parameters depend on the inputs: $(\mu = [\mu(\mathbf{x}_i)], \Sigma = [K(\mathbf{x}_i, \mathbf{x}_j)])$.
- ▶ If we think of the g_s as function values, a GP provides a prior over functions.
- ▶ In a GP regression model, noisy observations y_i are conditionally independent given g_i .
- ▶ No parameters to learn (though often hyperparameters); instead, we make predictions on test data directly: [assuming $\mu = 0$, and matrix Σ incorporates diagonal noise]

$$P(y' | \mathbf{x}', D) = \mathcal{N}(\Sigma_{\mathbf{x}', X} \Sigma_{X, X}^{-1} \mathbf{y}, \Sigma_{\mathbf{x}', \mathbf{x}'} - \Sigma_{\mathbf{x}', X} \Sigma_{X, X}^{-1} \Sigma_{X, \mathbf{x}'})$$

EP GP prediction



- ▶ Once approximate site potentials have stabilised, they can be used to make predictions.
- ▶ Introducing a test point changes K , but does not affect the *marginal* $P(g_1 \dots g_n)$ (by consistency of the GP).
- ▶ The unobserved output factor provides no information about g' (\Rightarrow constant factor on g').
- ▶ Thus no change is needed to the approximating potentials \tilde{f}_i .
- ▶ Predictions are obtained by marginalising the approximation: [let $\tilde{\Psi} = \text{diag}[\tilde{\psi}_1^2 \dots \tilde{\psi}_n^2]$]

$$P(y' | \mathbf{x}', D) = \int dg' P(y' | g') \mathcal{N}(g' | K_{\mathbf{x}', X} (K_{X, X} + \tilde{\Psi})^{-1} \tilde{\mu}, K_{\mathbf{x}', \mathbf{x}'} - K_{\mathbf{x}', X} (K_{X, X} + \tilde{\Psi})^{-1} K_{X, \mathbf{x}'})$$

Normalisers

- ▶ Approximate sites determined by moment matching are naturally normalised.
- ▶ For posteriors, sufficient to normalise product after convergence.
 - ▶ Often straightforward for exponential family approximations.
- ▶ To compute likelihood need to keep track of site integrals:
 - ▶ minimising “unnormalised KL”:

$$\mathbf{KL}[p||q] = \int dx p(x) \log \frac{p(x)}{q(x)} + \int dx (q(x) - p(x))$$

incorporates normaliser into each \tilde{f} (match zeroth moment, along with suff stats).

as well as the overall normaliser of $\prod_i \tilde{f}_i(\mathcal{Y}_i)$.

Alpha divergences and Power EP

- ▶ Alpha divergences $D_\alpha[p||q] = \frac{1}{\alpha(1-\alpha)} \int dx \alpha p(x) + (1-\alpha)q(x) - p(x)^\alpha q(x)^{1-\alpha}$

$$D_{-1}[p||q] = \frac{1}{2} \int dx \frac{(p(x) - q(x))^2}{p(x)}$$

$$\lim_{\alpha \rightarrow 0} D_\alpha[p||q] = \mathbf{KL}[q||p] \quad \text{Note: } \lim_{\alpha \rightarrow 0} \frac{(p(x)/q(x))^\alpha}{\alpha} = \log \frac{p(x)}{q(x)}$$

$$D_{\frac{1}{2}}[p||q] = 2 \int dx (p(x)^{\frac{1}{2}} - q(x)^{\frac{1}{2}})^2$$

$$\lim_{\alpha \rightarrow 1} D_\alpha[p||q] = \mathbf{KL}[p||q]$$

$$D_2[p||q] = \frac{1}{2} \int dx \frac{(p(x) - q(x))^2}{q(x)}$$

- ▶ Local (EP) minimisation gives fixed-point updates that blend messages (to power α) with previous site approximations.

$$\tilde{f}_i^{\text{new}} = \operatorname{argmin}_{f \in \{f\}} \mathbf{KL} [f_i(\mathcal{Y}_i)^\alpha \tilde{f}_i(\mathcal{Y}_i)^{1-\alpha} q_{-i}(\mathcal{Y}) || f(\mathcal{Y}_i) q_{-i}(\mathcal{Y})]$$

- ▶ Small changes (for $\alpha < 1$) lead to more stable updates, and more reliable convergence.