Probabilistic & Unsupervised Learning Approximate Inference

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 2023

Intractabilities and approximations

- Inference computational intractability
 - Gibbs sampling, other MCMC
 - Factored variational approx
 - Loopy BP/EP/Power EP
 - Recognition models
- Inference analytic intractability
 - Laplace approximation (global)
 - (Sequential) Monte-Carlo
 - Message approximations (linearised, sigma-point, Laplace)
 - Assumed-density methods and Expectation-Propagation
 - Parametric variational approx
 - Recognition models
- Learning intractable partition function
 - Sampling parameters
 - Constrastive divergence
 - Score-matching
- Posterior estimation and model selection
 - Laplace approximation / BIC
 - Monte-Carlo
 - (Annealed) importance sampling
 - Reversible jump MCMC
 - Variational Bayes

Not a complete list!



 $\mathbf{z}_{t+1} = f(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{w}_t$ $\mathbf{x}_t = g(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{v}_t$

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





 $\mathbf{z}_{t+1} = f(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{w}_t$ $\mathbf{x}_t = g(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{v}_t$

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

Extended Kalman Filter (EKF): linearise nonlinear functions about current estimate, \hat{z}_{i}^{t} :

$$\mathbf{z}_{t+1} \approx f(\hat{\mathbf{z}}_t^t, \mathbf{u}_t) + \left. \frac{\partial f}{\partial \mathbf{z}_t} \right|_{\mathbf{z}_t^t} (\mathbf{z}_t - \hat{\mathbf{z}}_t^t) + \mathbf{w}_t$$

$$\mathbf{x}_t pprox g(\mathbf{\hat{z}}_t^{t-1}, \mathbf{u}_t) + \left. rac{\partial g}{\partial \mathbf{z}_t} \right|_{\mathbf{\hat{z}}_t^{t-1}} (\mathbf{z}_t - \mathbf{\hat{z}}_t^{t-1}) + \mathbf{v}_t$$





 $\mathbf{z}_{t+1} = f(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{w}_t$ $\mathbf{x}_t = g(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{v}_t$

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

Extended Kalman Filter (EKF): linearise nonlinear functions about current estimate, \hat{z}_{t}^{t} :



Run the Kalman filter (smoother) on non-stationary linearised system (A_t, B_t, C_t, D_t) :



 $\mathbf{z}_{t+1} = f(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{w}_t$ $\mathbf{x}_t = g(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{v}_t$

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

Extended Kalman Filter (EKF): linearise nonlinear functions about current estimate, \hat{z}_{t}^{t} :



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

Adaptively approximates non-Gaussian messages by Gaussians.



 $\mathbf{z}_{t+1} = f(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{w}_t$ $\mathbf{x}_t = g(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{v}_t$

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

Extended Kalman Filter (EKF): linearise nonlinear functions about current estimate, \hat{z}_{t}^{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 ⇒ approximation degrades with increased state uncertainty.



 $\mathbf{z}_{t+1} = f(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{w}_t$ $\mathbf{x}_t = g(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{v}_t$

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

Extended Kalman Filter (EKF): linearise nonlinear functions about current estimate, \hat{z}_{t}^{t} :



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

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



 $\mathbf{z}_{t+1} = f(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{w}_t$ $\mathbf{x}_t = g(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{v}_t$

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

Extended Kalman Filter (EKF): linearise nonlinear functions about current estimate, \hat{z}_{i}^{t} :



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

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

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

Consider the forward messages on a latent chain:

$$P(\mathbf{z}_{t}|\mathbf{x}_{1:t}) = \frac{1}{Z}P(\mathbf{x}_{t}|\mathbf{z}_{t})\int d\mathbf{z}_{t-1} P(\mathbf{z}_{t}|\mathbf{z}_{t-1})P(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})$$

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

$$\tilde{P}(\mathbf{z}_t|\mathbf{x}_{1:t}) \approx \frac{1}{Z} P(\mathbf{x}_t|\mathbf{z}_t) \int d\mathbf{z}_{t-1} \underbrace{P(\mathbf{z}_t|\mathbf{z}_{t-1})}_{\mathcal{N}(f(\mathbf{z}_{t-1}), Q)} \underbrace{\tilde{P}(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})}_{\mathcal{N}(\hat{\mathbf{z}}_{t-1}, V_{t-1})}$$

Consider the forward messages on a latent chain:

$$P(\mathbf{z}_{t}|\mathbf{x}_{1:t}) = \frac{1}{Z}P(\mathbf{x}_{t}|\mathbf{z}_{t})\int d\mathbf{z}_{t-1} P(\mathbf{z}_{t}|\mathbf{z}_{t-1})P(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})$$

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

$$\tilde{P}(\mathbf{z}_t|\mathbf{x}_{1:t}) \approx \frac{1}{Z} P(\mathbf{x}_t|\mathbf{z}_t) \int d\mathbf{z}_{t-1} \underbrace{P(\mathbf{z}_t|\mathbf{z}_{t-1})}_{\mathcal{N}(f(\mathbf{z}_{t-1}), Q)} \underbrace{\tilde{P}(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})}_{\mathcal{N}(\hat{\mathbf{z}}_{t-1}, V_{t-1})}$$

Linearisation at the peak (EKF) is only one approach.

Consider the forward messages on a latent chain:

$$P(\mathbf{z}_{t}|\mathbf{x}_{1:t}) = \frac{1}{Z}P(\mathbf{x}_{t}|\mathbf{z}_{t})\int d\mathbf{z}_{t-1} P(\mathbf{z}_{t}|\mathbf{z}_{t-1})P(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})$$

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

$$\tilde{P}(\mathbf{z}_t|\mathbf{x}_{1:t}) \approx \frac{1}{Z} P(\mathbf{x}_t|\mathbf{z}_t) \int d\mathbf{z}_{t-1} \underbrace{P(\mathbf{z}_t|\mathbf{z}_{t-1})}_{\mathcal{N}(f(\mathbf{z}_{t-1}), Q)} \underbrace{\tilde{P}(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})}_{\mathcal{N}(\hat{\mathbf{z}}_{t-1}, V_{t-1})}$$

- Linearisation at the peak (EKF) is only one approach.
- Laplace filter: use mode and curvature of integrand.

Consider the forward messages on a latent chain:

$$P(\mathbf{z}_{t}|\mathbf{x}_{1:t}) = \frac{1}{Z}P(\mathbf{x}_{t}|\mathbf{z}_{t})\int d\mathbf{z}_{t-1} P(\mathbf{z}_{t}|\mathbf{z}_{t-1})P(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})$$

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

$$\widetilde{P}(\mathbf{z}_t|\mathbf{x}_{1:t}) \approx \frac{1}{Z} P(\mathbf{x}_t|\mathbf{z}_t) \int d\mathbf{z}_{t-1} \underbrace{P(\mathbf{z}_t|\mathbf{z}_{t-1})}_{\mathcal{N}(f(\mathbf{z}_{t-1}), Q)} \underbrace{\widetilde{P}(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})}_{\mathcal{N}(\hat{\mathbf{z}}_{t-1}, V_{t-1})}$$

- Linearisation at the peak (EKF) is only one approach.
- Laplace filter: use mode and curvature of integrand.
- Sigma-point ("unscented") filter: next slide.

Consider the forward messages on a latent chain:

$$P(\mathbf{z}_t|\mathbf{x}_{1:t}) = \frac{1}{Z} P(\mathbf{x}_t|\mathbf{z}_t) \int d\mathbf{z}_{t-1} P(\mathbf{z}_t|\mathbf{z}_{t-1}) P(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})$$

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

$$\tilde{P}(\mathbf{z}_{t}|\mathbf{x}_{1:t}) \approx \frac{1}{Z} P(\mathbf{x}_{t}|\mathbf{z}_{t}) \int d\mathbf{z}_{t-1} \underbrace{P(\mathbf{z}_{t}|\mathbf{z}_{t-1})}_{\mathcal{N}(f(\mathbf{z}_{t-1}), Q)} \underbrace{\tilde{P}(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})}_{\mathcal{N}(\hat{\mathbf{z}}_{t-1}, V_{t-1})}$$

- Linearisation at the peak (EKF) is only one approach.
- Laplace filter: use mode and curvature of integrand.
- Sigma-point ("unscented") filter: next slide.
- Parametric variational:

$$\operatorname{argmin} \mathbf{KL} \left[\mathcal{N}(\hat{\mathbf{z}}_t, \hat{V}_t) \middle\| \int d\mathbf{z}_{t-1} \ldots \right].$$

Needs Gaussian expectations of log $\int \Rightarrow$ Monte-Carlo integration (later lecture).

Consider the forward messages on a latent chain:

$$P(\mathbf{z}_t|\mathbf{x}_{1:t}) = \frac{1}{Z} P(\mathbf{x}_t|\mathbf{z}_t) \int d\mathbf{z}_{t-1} P(\mathbf{z}_t|\mathbf{z}_{t-1}) P(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})$$

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

$$\tilde{P}(\mathbf{z}_{t}|\mathbf{x}_{1:t}) \approx \frac{1}{Z} P(\mathbf{x}_{t}|\mathbf{z}_{t}) \int d\mathbf{z}_{t-1} \underbrace{P(\mathbf{z}_{t}|\mathbf{z}_{t-1})}_{\mathcal{N}(f(\mathbf{z}_{t-1}), Q)} \underbrace{\tilde{P}(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})}_{\mathcal{N}(\hat{\mathbf{z}}_{t-1}, V_{t-1})}$$

- Linearisation at the peak (EKF) is only one approach.
- Laplace filter: use mode and curvature of integrand.
- Sigma-point ("unscented") filter: next slide.
- Parametric variational:

$$\operatorname{argmin} \mathbf{KL} \left[\mathcal{N}(\hat{\mathbf{z}}_t, \hat{V}_t) \middle\| \int d\mathbf{z}_{t-1} \ldots \right].$$

Needs Gaussian expectations of log $\int \Rightarrow$ Monte-Carlo integration (later lecture). The other KL:

$$\operatorname{argmin} \mathbf{KL} \left[\int d\mathbf{z}_{t-1} \left\| \mathcal{N}(\hat{\mathbf{z}}_t, \hat{V}_t) \right]
ight]$$

needs only first and second moments of nonlinear message \Rightarrow EP.

Historical interest, but also a useful intuition for what comes next.

Historical interest, but also a useful intuition for what comes next.



Approximates pushed-forward belief from time t-1.

Historical interest, but also a useful intuition for what comes next.



- Approximates pushed-forward belief from time t-1.
- Evaluate $f(\hat{z}_{t-1}), f(\hat{z}_{t-1} \pm \sqrt{\lambda}v)$ for eigenvalues, eigenvectors $\hat{V}_{t-1}v = \lambda v$.

Historical interest, but also a useful intuition for what comes next.



- Approximates pushed-forward belief from time t-1.
- Evaluate $f(\hat{z}_{t-1}), f(\hat{z}_{t-1} \pm \sqrt{\lambda}v)$ for eigenvalues, eigenvectors $\hat{V}_{t-1}v = \lambda v$.
- Fit" Gaussian to these $2K + 1 \sigma$ -points:

$$\mathcal{N}\Big(\underbrace{\frac{1}{2K+1}\sum_{i=1}^{2K+1}\mathbf{f}(\boldsymbol{\sigma}_i)}_{\boldsymbol{\mu}}, \frac{1}{2K+1}\sum_{i=1}^{2K+1}\mathbf{f}(\boldsymbol{\sigma}_i)\mathbf{f}(\boldsymbol{\sigma}_i)^{\mathsf{T}} - \boldsymbol{\mu}\boldsymbol{\mu}^{\mathsf{T}}\Big)$$

Historical interest, but also a useful intuition for what comes next.



- Approximates pushed-forward belief from time t-1.
- Evaluate $f(\hat{z}_{t-1}), f(\hat{z}_{t-1} \pm \sqrt{\lambda}v)$ for eigenvalues, eigenvectors $\hat{V}_{t-1}v = \lambda v$.
- Fit" Gaussian to these $2K + 1 \sigma$ -points:

$$\mathcal{N}\left(\underbrace{\frac{1}{2K+1}\sum_{i=1}^{2K+1}\mathbf{f}(\boldsymbol{\sigma}_i)}_{\boldsymbol{\mu}}, \frac{1}{2K+1}\sum_{i=1}^{2K+1}\mathbf{f}(\boldsymbol{\sigma}_i)\mathbf{f}(\boldsymbol{\sigma}_i)^{\mathsf{T}} - \boldsymbol{\mu}\boldsymbol{\mu}^{\mathsf{T}} + \mathsf{Q}\right)$$

Incorporate noise process.

Historical interest, but also a useful intuition for what comes next.



- Approximates pushed-forward belief from time t-1.
- Evaluate $f(\hat{\mathbf{z}}_{t-1}), f(\hat{\mathbf{z}}_{t-1} \pm \sqrt{\lambda}\mathbf{v})$ for eigenvalues, eigenvectors $\hat{V}_{t-1}\mathbf{v} = \lambda\mathbf{v}$.
- Fit" Gaussian to these $2K + 1 \sigma$ -points:

$$\mathcal{N}\left(\underbrace{\frac{1}{2K+1}\sum_{i=1}^{2K+1}\mathbf{f}(\boldsymbol{\sigma}_i)}_{\boldsymbol{\mu}}, \frac{1}{2K+1}\sum_{i=1}^{2K+1}\mathbf{f}(\boldsymbol{\sigma}_i)\mathbf{f}(\boldsymbol{\sigma}_i)^{\mathsf{T}} - \boldsymbol{\mu}\boldsymbol{\mu}^{\mathsf{T}} + \mathsf{Q}\right)$$

- Incorporate noise process.
- Equivalent to evaluation of mean and covariance of \tilde{P}_{t-1} #f by Gaussian quadrature.

Historical interest, but also a useful intuition for what comes next.



- Approximates pushed-forward belief from time t-1.
- Evaluate $f(\hat{z}_{t-1}), f(\hat{z}_{t-1} \pm \sqrt{\lambda}v)$ for eigenvalues, eigenvectors $\hat{V}_{t-1}v = \lambda v$.
- Fit" Gaussian to these $2K + 1 \sigma$ -points:

$$\mathcal{N}\left(\underbrace{\frac{1}{2K+1}\sum_{i=1}^{2K+1}\mathbf{f}(\boldsymbol{\sigma}_i)}_{\boldsymbol{\mu}}, \frac{1}{2K+1}\sum_{i=1}^{2K+1}\mathbf{f}(\boldsymbol{\sigma}_i)\mathbf{f}(\boldsymbol{\sigma}_i)^{\mathsf{T}} - \boldsymbol{\mu}\boldsymbol{\mu}^{\mathsf{T}} + \mathsf{Q}\right)$$

- Incorporate noise process.
- Equivalent to evaluation of mean and covariance of $\tilde{P}_{t-1} \# f$ by Gaussian quadrature.
- One form of "Assumed Density Filtering" (and of calculations for EP).

Free energy:

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

Free energy:

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

E-steps:

Exact EM:
$$q(\mathcal{Z}) = \underset{q}{\operatorname{argmax}} \mathcal{F} = P(\mathcal{Z}|\mathcal{X}, \theta)$$

Free energy:

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

E-steps:

• Exact EM:
$$q(\mathcal{Z}) = \underset{q}{\operatorname{argmax}} \mathcal{F} = P(\mathcal{Z}|\mathcal{X}, \theta)$$

Saturates bound: converges to local maximum of likelihood.

Free energy:

 $\mathcal{F}(q,\theta) = \left\langle \log \mathsf{P}(\mathcal{X},\mathcal{Z}|\theta) \right\rangle_{q(\mathcal{Z}|\mathcal{X})} + \mathsf{H}[q] = \log \mathsf{P}(\mathcal{X}|\theta) - \mathsf{KL}[q(\mathcal{Z}) \| \mathsf{P}(\mathcal{Z}|\mathcal{X},\theta)] \le \ell(\theta)$

E-steps:

• Exact EM:
$$q(\mathcal{Z}) = \operatorname{argmax} \mathcal{F} = P(\mathcal{Z}|\mathcal{X}, \theta)$$

Saturates bound: converges to local maximum of likelihood.

(Factored) variational approximation:

 $q(\mathcal{Z}) = \underset{q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2)}{\operatorname{argmax}} \mathcal{F} = \underset{q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2)}{\operatorname{argmax}} \mathsf{KL}[q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2) || P(\mathcal{Z}|\mathcal{X},\theta)]$

Free energy:

 $\mathcal{F}(q,\theta) = \left\langle \log \mathsf{P}(\mathcal{X},\mathcal{Z}|\theta) \right\rangle_{q(\mathcal{Z}|\mathcal{X})} + \mathsf{H}[q] = \log \mathsf{P}(\mathcal{X}|\theta) - \mathsf{KL}[q(\mathcal{Z}) \| \mathsf{P}(\mathcal{Z}|\mathcal{X},\theta)] \le \ell(\theta)$

E-steps:

• Exact EM:
$$q(\mathcal{Z}) = \operatorname{argmax}_{a} \mathcal{F} = P(\mathcal{Z}|\mathcal{X}, \theta)$$

Saturates bound: converges to local maximum of likelihood.

(Factored) variational approximation:

$$q(\mathcal{Z}) = \operatorname*{argmax}_{q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2)} \mathcal{F} = \operatorname*{argmin}_{q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2)} \mathsf{KL}[q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2) || P(\mathcal{Z}|\mathcal{X},\theta)]$$

Increases bound: converges, but not necessarily to ML.

Free energy:

 $\mathcal{F}(q,\theta) = \left\langle \log \mathsf{P}(\mathcal{X},\mathcal{Z}|\theta) \right\rangle_{q(\mathcal{Z}|\mathcal{X})} + \mathsf{H}[q] = \log \mathsf{P}(\mathcal{X}|\theta) - \mathsf{KL}[q(\mathcal{Z})\|\mathsf{P}(\mathcal{Z}|\mathcal{X},\theta)] \leq \ell(\theta)$

E-steps:

• Exact EM:
$$q(\mathcal{Z}) = \operatorname{argmax}_{a} \mathcal{F} = P(\mathcal{Z}|\mathcal{X}, \theta)$$

Saturates bound: converges to local maximum of likelihood.

(Factored) variational approximation:

 $q(\mathcal{Z}) = \underset{q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2)}{\operatorname{argmax}} \mathcal{F} = \underset{q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2)}{\operatorname{argmax}} \mathsf{KL}[q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2) || P(\mathcal{Z}|\mathcal{X}, \theta)]$

Increases bound: converges, but not necessarily to ML.

• Other approximations: $q(\mathcal{Z}) \approx P(\mathcal{Z}|\mathcal{X}, \theta)$

Free energy:

 $\mathcal{F}(q,\theta) = \left\langle \log \mathsf{P}(\mathcal{X},\mathcal{Z}|\theta) \right\rangle_{q(\mathcal{Z}|\mathcal{X})} + \mathsf{H}[q] = \log \mathsf{P}(\mathcal{X}|\theta) - \mathsf{KL}[q(\mathcal{Z})\|\mathsf{P}(\mathcal{Z}|\mathcal{X},\theta)] \leq \ell(\theta)$

E-steps:

Exact EM:
$$q(\mathcal{Z}) = \operatorname{argmax}_{a} \mathcal{F} = P(\mathcal{Z}|\mathcal{X}, \theta)$$

Saturates bound: converges to local maximum of likelihood.

(Factored) variational approximation:

 $q(\mathcal{Z}) = \underset{q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2)}{\operatorname{argmax}} \mathcal{F} = \underset{q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2)}{\operatorname{argmax}} \mathsf{KL}[q_1(\mathcal{Z}_1)q_2(\mathcal{Z}_2) || P(\mathcal{Z}|\mathcal{X}, \theta)]$

Increases bound: converges, but not necessarily to ML.

• Other approximations: $q(\mathcal{Z}) \approx P(\mathcal{Z}|\mathcal{X}, \theta)$

Usually no guarantees, but if learning converges it may be more accurate than the factored approximation

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.

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

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.

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

Open choices:

- form of the metric D
- nature of the constraint space Q

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.

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

Open choices:

- form of the metric D
- nature of the constraint space Q

Variational methods: $D = \mathbf{KL}[q || P]$.

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 Q

- Variational methods: $D = \mathbf{KL}[q || P]$.
 - Choosing $Q = \{$ tree-factored distributions $\}$ leads to efficient message passing.

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 Q

- Variational methods: $D = \mathbf{KL}[q || P]$.
 - Choosing $Q = \{$ tree-factored distributions $\}$ leads to efficient message passing.
- Can we use other divergences?

The other KL

What about the 'other' KL ($q = \operatorname{argmin} \operatorname{KL}[P || q]$)?

The other KL

What about the 'other' KL ($q = \operatorname{argmin} \operatorname{KL}[P || q]$)?

For a factored approximation the (clique) marginals obtained by minimising this KL are correct:
What about the 'other' KL ($q = \operatorname{argmin} \operatorname{KL}[P || q]$)?

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

$$\operatorname*{argmin}_{q_i} \mathsf{KL}\Big[\mathsf{P}(\mathcal{Z}|\mathcal{X}) \Big\| \prod q_i(\mathcal{Z}_j|\mathcal{X}) \Big] = \operatorname*{argmin}_{q_i} - \int \mathsf{d}\mathcal{Z} \; \mathsf{P}(\mathcal{Z}|\mathcal{X}) \log \prod_j q_j(\mathcal{Z}_j|\mathcal{X})$$

What about the 'other' KL ($q = \operatorname{argmin} \operatorname{KL}[P || q]$)?

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

$$\begin{aligned} \operatorname*{argmin}_{q_{i}} \mathsf{KL}\Big[\mathsf{P}(\mathcal{Z}|\mathcal{X}) \Big\| \prod q_{i}(\mathcal{Z}_{j}|\mathcal{X}) \Big] &= \operatorname*{argmin}_{q_{i}} - \int d\mathcal{Z} \ \mathsf{P}(\mathcal{Z}|\mathcal{X}) \log \prod_{j} q_{j}(\mathcal{Z}_{j}|\mathcal{X}) \\ &= \operatorname*{argmin}_{q_{i}} - \sum_{j} \int d\mathcal{Z} \ \mathsf{P}(\mathcal{Z}|\mathcal{X}) \log q_{j}(\mathcal{Z}_{j}|\mathcal{X}) \end{aligned}$$

What about the 'other' KL ($q = \operatorname{argmin} \operatorname{KL}[P || q]$)?

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

$$\begin{aligned} \underset{q_i}{\operatorname{argmin}} \operatorname{\mathsf{KL}}\left[P(\mathcal{Z}|\mathcal{X}) \middle\| \prod q_i(\mathcal{Z}_j|\mathcal{X})\right] &= \underset{q_i}{\operatorname{argmin}} - \int d\mathcal{Z} \ P(\mathcal{Z}|\mathcal{X}) \log \prod_j q_j(\mathcal{Z}_j|\mathcal{X}) \\ &= \underset{q_i}{\operatorname{argmin}} - \sum_j \int d\mathcal{Z} \ P(\mathcal{Z}|\mathcal{X}) \log q_j(\mathcal{Z}_j|\mathcal{X}) \\ &= \underset{q_i}{\operatorname{argmin}} - \int d\mathcal{Z}_i \ P(\mathcal{Z}_i|\mathcal{X}) \log q_i(\mathcal{Z}_i|\mathcal{X}) \end{aligned}$$

What about the 'other' KL ($q = \operatorname{argmin} \operatorname{KL}[P || q]$)?

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

$$\begin{aligned} \underset{q_i}{\operatorname{argmin}} \operatorname{\mathsf{KL}}\left[P(\mathcal{Z}|\mathcal{X}) \middle\| \prod q_i(\mathcal{Z}_j|\mathcal{X})\right] &= \underset{q_i}{\operatorname{argmin}} - \int d\mathcal{Z} \ P(\mathcal{Z}|\mathcal{X}) \log \prod_j q_j(\mathcal{Z}_j|\mathcal{X}) \\ &= \underset{q_i}{\operatorname{argmin}} - \sum_j \int d\mathcal{Z} \ P(\mathcal{Z}|\mathcal{X}) \log q_i(\mathcal{Z}_j|\mathcal{X}) \\ &= \underset{q_i}{\operatorname{argmin}} - \int d\mathcal{Z}_i \ P(\mathcal{Z}_i|\mathcal{X}) \log q_i(\mathcal{Z}_i|\mathcal{X}) \\ &= P(\mathcal{Z}_i|\mathcal{X}) \end{aligned}$$

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).

What about the 'other' KL ($q = \operatorname{argmin} \operatorname{KL}[P || q]$)?

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

$$\begin{aligned} \underset{q_i}{\operatorname{argmin}} \operatorname{\mathsf{KL}}\left[P(\mathcal{Z}|\mathcal{X}) \middle\| \prod q_i(\mathcal{Z}_j|\mathcal{X})\right] &= \underset{q_i}{\operatorname{argmin}} - \int d\mathcal{Z} \ P(\mathcal{Z}|\mathcal{X}) \log \prod_j q_j(\mathcal{Z}_j|\mathcal{X}) \\ &= \underset{q_i}{\operatorname{argmin}} - \sum_j \int d\mathcal{Z} \ P(\mathcal{Z}|\mathcal{X}) \log q_i(\mathcal{Z}_j|\mathcal{X}) \\ &= \underset{q_i}{\operatorname{argmin}} - \int d\mathcal{Z}_i \ P(\mathcal{Z}_i|\mathcal{X}) \log q_i(\mathcal{Z}_i|\mathcal{X}) \\ &= P(\mathcal{Z}_i|\mathcal{X}) \end{aligned}$$

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!

What about the 'other' KL ($q = \operatorname{argmin} \operatorname{KL}[P || q]$)?

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

$$\begin{aligned} \underset{q_i}{\operatorname{argmin}} \operatorname{\mathsf{KL}}\left[P(\mathcal{Z}|\mathcal{X}) \middle\| \prod q_i(\mathcal{Z}_j|\mathcal{X})\right] &= \underset{q_i}{\operatorname{argmin}} - \int d\mathcal{Z} \ P(\mathcal{Z}|\mathcal{X}) \log \prod_j q_j(\mathcal{Z}_j|\mathcal{X}) \\ &= \underset{q_i}{\operatorname{argmin}} - \sum_j \int d\mathcal{Z} \ P(\mathcal{Z}|\mathcal{X}) \log q_i(\mathcal{Z}_j|\mathcal{X}) \\ &= \underset{q_i}{\operatorname{argmin}} - \int d\mathcal{Z}_i \ P(\mathcal{Z}_i|\mathcal{X}) \log q_i(\mathcal{Z}_i|\mathcal{X}) \\ &= P(\mathcal{Z}_i|\mathcal{X}) \end{aligned}$$

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.

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

$$P(\mathcal{Z}|\mathcal{X}) = rac{P(\mathcal{Z},\mathcal{X})}{P(\mathcal{X})} = rac{1}{Z}\prod_i P(Z_i|\operatorname{pa}(Z_i)) \propto \prod_{i=1}^N f_i(\mathcal{Z}_i)$$

where the Z_i are not necessarily disjoint. In the language of EP the f_i are called sites.

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

$$P(\mathcal{Z}|\mathcal{X}) = rac{P(\mathcal{Z},\mathcal{X})}{P(\mathcal{X})} = rac{1}{Z}\prod_i P(Z_i|\operatorname{pa}(Z_i)) \propto \prod_{i=1}^N f_i(\mathcal{Z}_i)$$

where the Z_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: $q(Z) \propto \prod_{i=1}^{N} \tilde{f}_i(Z_i)$. We would like to minimise (at least in some sense) $\mathbf{KL}[P||q]$.

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

$$P(\mathcal{Z}|\mathcal{X}) = rac{P(\mathcal{Z},\mathcal{X})}{P(\mathcal{X})} = rac{1}{Z}\prod_i P(Z_i|\operatorname{pa}(Z_i)) \propto \prod_{i=1}^N f_i(\mathcal{Z}_i)$$

where the \mathcal{Z}_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: $q(\mathcal{Z}) \propto \prod_{i=1}^{\infty} \tilde{f}_i(\mathcal{Z}_i)$. We would like to minimise (at least in some sense) **KL**[*P*||*q*].

Possible optimisations:

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

$$P(\mathcal{Z}|\mathcal{X}) = \frac{P(\mathcal{Z}, \mathcal{X})}{P(\mathcal{X})} = \frac{1}{Z} \prod_{i} P(Z_i| \operatorname{pa}(Z_i)) \propto \prod_{i=1}^{N} f_i(\mathcal{Z}_i)$$

where the \mathcal{Z}_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: $q(\mathcal{Z}) \propto \prod_{i=1}^{n} \tilde{f}_i(\mathcal{Z}_i)$. We would like to minimise (at least in some sense) $\mathsf{KL}[P||q]$.

Possible optimisations:

$$\min_{\{\tilde{f}_i\}} \mathsf{KL}\Big[\frac{1}{Z} \prod_{i=1}^N f_i(\mathcal{Z}_i) \Big\| \frac{1}{Z} \prod_{i=1}^N \tilde{f}_i(\mathcal{Z}_i) \Big]$$

(global: intractable)

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

$$P(\mathcal{Z}|\mathcal{X}) = \frac{P(\mathcal{Z}, \mathcal{X})}{P(\mathcal{X})} = \frac{1}{Z} \prod_{i} P(Z_i| \operatorname{pa}(Z_i)) \propto \prod_{i=1}^{N} f_i(\mathcal{Z}_i)$$

where the Z_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: $q(\mathcal{Z}) \propto \prod_{i=1}^{n} \tilde{f}_i(\mathcal{Z}_i)$. We would like to minimise (at least in some sense) $\mathsf{KL}[P||q]$.

Possible optimisations:

$$\min_{\{\tilde{l}_i\}} \mathbf{KL} \left[\frac{1}{Z} \prod_{i=1}^{N} f_i(\mathcal{Z}_i) \right\| \frac{1}{Z} \prod_{i=1}^{N} \tilde{f}_i(\mathcal{Z}_i) \right]$$
$$\min_{\tilde{l}_i} \mathbf{KL} \left[f_i(\mathcal{Z}_i) \right\| \tilde{f}_i(\mathcal{Z}_i) \right]$$

(global: intractable)

(local, fixed: simple, inaccurate)



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

$$P(\mathcal{Z}|\mathcal{X}) = \frac{P(\mathcal{Z}, \mathcal{X})}{P(\mathcal{X})} = \frac{1}{Z} \prod_{i} P(Z_i| \operatorname{pa}(Z_i)) \propto \prod_{i=1}^{N} f_i(\mathcal{Z}_i)$$

where the \mathcal{Z}_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: $q(\mathcal{Z}) \propto \prod_{i=1}^{N} \tilde{f}_i(\mathcal{Z}_i)$. We would like to minimise (at least in some sense) $\mathsf{KL}[P||q]$.

Possible optimisations:

$$\begin{split} \min_{\{\tilde{t}_i\}} \mathsf{KL} \Big[\frac{1}{\mathbb{Z}} \prod_{i=1}^{N} f_i(\mathcal{Z}_i) \Big\| \frac{1}{\mathbb{Z}} \prod_{i=1}^{N} \tilde{f}_i(\mathcal{Z}_i) \Big] \\ \min_{\tilde{t}_i} \mathsf{KL} \Big[f_i(\mathcal{Z}_i) \Big\| \tilde{f}_i(\mathcal{Z}_i) \Big] \\ \min_{\tilde{t}_i} \mathsf{KL} \Big[f_i(\mathcal{Z}_i) \prod_{j \neq i} \tilde{f}_j(\mathcal{Z}_j) \Big\| \tilde{f}_i(\mathcal{Z}_i) \prod_{j \neq i} \tilde{f}_j(\mathcal{Z}_j) \Big] \end{split}$$

(global: intractable)

(local, fixed: simple, inaccurate)

(local, contextual: iterative, accurate)



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

$$P(\mathcal{Z}|\mathcal{X}) = \frac{P(\mathcal{Z}, \mathcal{X})}{P(\mathcal{X})} = \frac{1}{Z} \prod_{i} P(Z_i| \operatorname{pa}(Z_i)) \propto \prod_{i=1}^{N} f_i(\mathcal{Z}_i)$$

where the \mathcal{Z}_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: $q(\mathcal{Z}) \propto \prod_{i=1}^{N} \tilde{f}_i(\mathcal{Z}_i)$. We would like to minimise (at least in some sense) $\mathsf{KL}[P||q]$.

Possible optimisations:

$$\begin{split} \min_{\{\tilde{t}_i\}} \mathsf{KL} \Big[\frac{1}{Z} \prod_{i=1}^{N} f_i(\mathcal{Z}_i) \Big\| \frac{1}{Z} \prod_{i=1}^{N} \tilde{f}_i(\mathcal{Z}_i) \Big] \\ \min_{\tilde{t}_i} \mathsf{KL} \Big[f_i(\mathcal{Z}_i) \Big\| \tilde{f}_i(\mathcal{Z}_i) \Big] \\ \min_{\tilde{t}_j} \mathsf{KL} \Big[f_i(\mathcal{Z}_j) \prod_{j \neq i} \tilde{f}_j(\mathcal{Z}_j) \Big\| \tilde{f}_i(\mathcal{Z}_i) \prod_{j \neq i} \tilde{f}_j(\mathcal{Z}_j) \Big] \end{split}$$

(global: intractable)

(local, fixed: simple, inaccurate)

(local, contextual: iterative, accurate) $\leftarrow EP$



EP is really two ideas:

Approximation of factors.

EP is really two ideas:

- Approximation of factors.
 - Usually by "projection" to exponential families.
 - This involves finding expected sufficient statistics, hence expectation.

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.

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.

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.

Note: we will ignore normalisation for now, but return to this later.

Each EP update involves a KL minimisation:

 $\tilde{t}_{i}^{\text{new}}(\mathcal{Z}) \leftarrow \underset{f \in \{\tilde{l}\}}{\operatorname{argmin}} \operatorname{\mathsf{KL}}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \qquad \left[q_{\neg i}(\mathcal{Z}) \stackrel{\text{def}}{=} \prod_{j \neq i} \tilde{f}_{j}(\mathcal{Z}_{j})\right]$

Each EP update involves a KL minimisation:

$$\tilde{f}_{i}^{\text{new}}(\mathcal{Z}) \leftarrow \underset{t \in \{\tilde{l}\}}{\operatorname{argmin}} \operatorname{\mathsf{KL}}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \qquad \left[q_{\neg i}(\mathcal{Z}) \stackrel{\text{def}}{=} \prod_{j \neq i} \tilde{f}_{j}(\mathcal{Z}_{j}) \right]$$

Separate the contextual factor: $q_{\neg i}(\mathcal{Z}) = q_{\neg i}(\mathcal{Z}_i)q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_i)$ $[\mathcal{Z}_{\neg i} \stackrel{\text{def}}{=} \mathcal{Z} \setminus \mathcal{Z}_i]$

Each EP update involves a KL minimisation:

$$\tilde{f}_{i}^{\text{new}}(\mathcal{Z}) \leftarrow \underset{t \in \{\tilde{l}\}}{\operatorname{argmin}} \operatorname{\mathsf{KL}}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \qquad \left[q_{\neg i}(\mathcal{Z}) \stackrel{\text{def}}{=} \prod_{j \neq i} \tilde{f}_{j}(\mathcal{Z}_{j}) \right]$$

Separate the contextual factor: $q_{\neg i}(\mathcal{Z}) = q_{\neg i}(\mathcal{Z}_i)q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_i)$ $[\mathcal{Z}_{\neg i} \stackrel{\text{def}}{=} \mathcal{Z} \setminus \mathcal{Z}_i]$

Then:

$$\min_{t} \mathsf{KL}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) || f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})]$$
$$= \max_{t} \int d\mathcal{Z} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})$$

Each EP update involves a KL minimisation:

$$\tilde{f}_{i}^{\text{new}}(\mathcal{Z}) \leftarrow \underset{t \in \{\tilde{l}\}}{\operatorname{argmin}} \operatorname{\mathsf{KL}}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \qquad \left[q_{\neg i}(\mathcal{Z}) \stackrel{\text{def}}{=} \prod_{j \neq i} \tilde{f}_{j}(\mathcal{Z}_{j}) \right]$$

Separate the contextual factor: $q_{\neg i}(\mathcal{Z}) = q_{\neg i}(\mathcal{Z}_i)q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_i)$ $[\mathcal{Z}_{\neg i} \stackrel{\text{def}}{=} \mathcal{Z} \setminus \mathcal{Z}_i]$

Then:

$$\begin{split} \min_{f} \mathbf{KL}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \\ &= \max_{f} \int d\mathcal{Z} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \\ &= \max_{f} \int d\mathcal{Z}_{i} d\mathcal{Z}_{\neg i} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i}) \big(\log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i}) + \log q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i})\big) \end{split}$$

Each EP update involves a KL minimisation:

$$\tilde{f}_{i}^{\text{new}}(\mathcal{Z}) \leftarrow \underset{t \in \{\tilde{f}\}}{\operatorname{argmin}} \operatorname{\mathsf{KL}}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \qquad \left[q_{\neg i}(\mathcal{Z}) \stackrel{\text{def}}{=} \prod_{j \neq i} \tilde{f}_{j}(\mathcal{Z}_{j}) \right]$$

Separate the contextual factor: $q_{\neg i}(\mathcal{Z}) = q_{\neg i}(\mathcal{Z}_i)q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_i)$ $[\mathcal{Z}_{\neg i} \stackrel{\text{def}}{=} \mathcal{Z} \setminus \mathcal{Z}_i]$ Then:

$$\begin{split} \min_{f} \mathsf{KL}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \\ &= \max_{f} \int d\mathcal{Z} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \\ &= \max_{f} \int d\mathcal{Z}_{i} d\mathcal{Z}_{\neg i} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i}) \big(\log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i}) + \log q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i})\big) \\ &= \max_{f} \int d\mathcal{Z}_{i} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i}) \big(\log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})\big) \int d\mathcal{Z}_{\neg i} q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i}) \end{split}$$

r

Each EP update involves a KL minimisation:

$$\tilde{f}_{i}^{\text{new}}(\mathcal{Z}) \leftarrow \underset{i \in \{\tilde{l}\}}{\operatorname{argmin}} \operatorname{KL}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \qquad \left[q_{\neg i}(\mathcal{Z}) \stackrel{\text{def}}{=} \prod_{j \neq i} \tilde{f}_{j}(\mathcal{Z}_{j}) \right]$$

Separate the contextual factor: $q_{\neg i}(\mathcal{Z}) = q_{\neg i}(\mathcal{Z}_i)q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_i)$ $[\mathcal{Z}_{\neg i} \stackrel{\text{def}}{=} \mathcal{Z} \setminus \mathcal{Z}_i]$ Then:

$$\begin{split} \min_{f} \mathsf{KL}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \\ &= \max_{f} \int d\mathcal{Z} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \\ &= \max_{f} \int d\mathcal{Z}_{i} d\mathcal{Z}_{\neg i} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i}) \big(\log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i}) + \log q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i})\big) \\ &= \max_{f} \int d\mathcal{Z}_{i} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i}) \big(\log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})\big) \int d\mathcal{Z}_{\neg i} q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i}) \\ &= \min_{f} \mathsf{KL}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})] \end{split}$$

Each EP update involves a KL minimisation:

$$\tilde{f}_{i}^{\text{new}}(\mathcal{Z}) \leftarrow \underset{f \in \{\tilde{l}\}}{\operatorname{argmin}} \operatorname{KL}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \qquad \left[q_{\neg i}(\mathcal{Z}) \stackrel{\text{def}}{=} \prod_{j \neq i} \tilde{f}_{j}(\mathcal{Z}_{j}) \right]$$

Separate the contextual factor: $q_{\neg i}(\mathcal{Z}) = q_{\neg i}(\mathcal{Z}_i)q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_i)$ $[\mathcal{Z}_{\neg i} \stackrel{\text{def}}{=} \mathcal{Z} \setminus \mathcal{Z}_i]$ Then:

$$\begin{split} \min_{f} \mathsf{KL}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})\|f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})] \\ &= \max_{f} \int d\mathcal{Z} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})\log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \\ &= \max_{f} \int d\mathcal{Z}_{i}d\mathcal{Z}_{\neg i} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i}) \big(\log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i}) + \log q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i})\big) \\ &= \max_{f} \int d\mathcal{Z}_{i} f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i}) \big(\log f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})\big) \int d\mathcal{Z}_{\neg i} q_{\neg i}(\mathcal{Z}_{\neg i}|\mathcal{Z}_{i}) \\ &= \min_{f} \mathsf{KL}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})\|f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})] \end{split}$$

 $q_{\neg i}(\mathcal{Z}_i)$ is sometimes called the cavity distribution.

Expectation Propagation (EP)

Input $f_1(\mathcal{Z}_1) \dots f_N(\mathcal{Z}_N)$ Initialize $\tilde{f}_1(\mathcal{Z}_1) = \operatorname{argmin} \operatorname{KL}[f_1(\mathcal{Z}_1) || f_1(\mathcal{Z}_1)], \ \tilde{f}_i(\mathcal{Z}_i) = 1 \text{ for } i > 1, \ q(\mathcal{Z}) \propto \prod_i \tilde{f}_i(\mathcal{Z}_i)$ $f \in \{\tilde{f}\}$ repeat for i = 1 ... N do Delete: $q_{\neg i}(\mathcal{Z}) \leftarrow \frac{q(\mathcal{Z})}{\tilde{f}_i(\mathcal{Z}_i)} = \prod_{i \neq i} \tilde{f}_i(\mathcal{Z}_i)$ Project: $\tilde{t}_{i}^{\text{new}}(\mathcal{Z}) \leftarrow \operatorname{argmin} \operatorname{KL}[t_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})||f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}_{i})]$ $f \in \{\hat{f}\}$ Include: $q(\mathcal{Z}) \leftarrow \tilde{f}_i^{\text{new}}(\mathcal{Z}_i) q_{\neg i}(\mathcal{Z})$ end for until convergence

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

$$q_{\neg i}(\mathcal{Z}_i) = \prod_{j \in \mathsf{ne}(i)} M_{j \to i}(\mathcal{Z}_j \cap \mathcal{Z}_i)$$

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

$$q_{\neg i}(\mathcal{Z}_i) = \prod_{j \in \mathsf{ne}(i)} M_{j \to i}(\mathcal{Z}_j \cap \mathcal{Z}_i)$$

 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.

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

$$q_{\neg i}(\mathcal{Z}_i) = \prod_{j \in \mathsf{ne}(i)} M_{j \to i}(\mathcal{Z}_j \cap \mathcal{Z}_i)$$

- 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

$$q_{\neg i}(\mathcal{Z}_i) = \prod_{j \in \mathsf{ne}(i)} M_{j \rightarrow i}(\mathcal{Z}_j \cap \mathcal{Z}_i)$$

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

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

$$q_{\neg i}(\mathcal{Z}_i) = \prod_{j \in \mathsf{ne}(i)} M_{j \to i}(\mathcal{Z}_j \cap \mathcal{Z}_i)$$

- 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

$$q_{\neg i}(\mathcal{Z}_i) = \prod_{j \in \mathsf{ne}(i)} M_{j \to i}(\mathcal{Z}_j \cap \mathcal{Z}_i)$$

becomes an approximation to the **true** cavity distribution (or we can recast the approximation directly in terms of messages \Rightarrow 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.

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

$$q_{\neg i}(\mathcal{Z}_i) = \prod_{j \in \mathsf{ne}(i)} M_{j \to i}(\mathcal{Z}_j \cap \mathcal{Z}_i)$$

- 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

$$q_{\neg i}(\mathcal{Z}_i) = \prod_{j \in \mathsf{ne}(i)} M_{j \to i}(\mathcal{Z}_j \cap \mathcal{Z}_i)$$

becomes an approximation to the **true** cavity distribution (or we can recast the approximation directly in terms of messages \Rightarrow 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.

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

$$q_{\neg i}(\mathcal{Z}_i) = \prod_{j \in \mathsf{ne}(i)} M_{j \to i}(\mathcal{Z}_j \cap \mathcal{Z}_i)$$

- 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

$$q_{\neg i}(\mathcal{Z}_i) = \prod_{j \in \mathsf{ne}(i)} M_{j \to i}(\mathcal{Z}_j \cap \mathcal{Z}_i)$$

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



$$P(\mathbf{z}_i | \mathbf{z}_{i-1}) = \phi_i(\mathbf{z}_i, \mathbf{z}_{i-1})$$
$$P(\mathbf{x}_i | \mathbf{z}_i) = \psi_i(\mathbf{z}_i)$$

e.g.
$$\exp(-\|\mathbf{z}_i - h_s(\mathbf{z}_{i-1})\|^2/2\sigma^2)$$

e.g. $\exp(-\|\mathbf{x}_i - h_o(\mathbf{z}_i)\|^2/2\sigma^2)$



$$P(\mathbf{z}_{i}|\mathbf{z}_{i-1}) = \phi_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) \qquad e.g. \exp(-\|\mathbf{z}_{i} - h_{s}(\mathbf{z}_{i-1})\|^{2}/2\sigma^{2}) P(\mathbf{x}_{i}|\mathbf{z}_{i}) = \psi_{i}(\mathbf{z}_{i}) \qquad e.g. \exp(-\|\mathbf{x}_{i} - h_{o}(\mathbf{z}_{i})\|^{2}/2\sigma^{2})$$

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



$$P(\mathbf{z}_{i}|\mathbf{z}_{i-1}) = \phi_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) \qquad e.g. \exp(-\|\mathbf{z}_{i} - h_{s}(\mathbf{z}_{i-1})\|^{2}/2\sigma^{2}) \\ P(\mathbf{x}_{i}|\mathbf{z}_{i}) = \psi_{i}(\mathbf{z}_{i}) \qquad e.g. \exp(-\|\mathbf{x}_{i} - h_{o}(\mathbf{z}_{i})\|^{2}/2\sigma^{2})$$

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

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

with both α and β Gaussian.



$$P(\mathbf{z}_{i}|\mathbf{z}_{i-1}) = \phi_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) \qquad e.g. \exp(-\|\mathbf{z}_{i} - h_{s}(\mathbf{z}_{i-1})\|^{2}/2\sigma^{2}) \\ P(\mathbf{x}_{i}|\mathbf{z}_{i}) = \psi_{i}(\mathbf{z}_{i}) \qquad e.g. \exp(-\|\mathbf{x}_{i} - h_{o}(\mathbf{z}_{i})\|^{2}/2\sigma^{2})$$

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

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

with both α and β Gaussian.

$$\tilde{f}_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) = \operatorname*{argmin}_{f \in \mathcal{N}} \mathsf{KL}\big[\phi_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1})\psi_{i}(\mathbf{z}_{i})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})\big\|f(\mathbf{z}_{i}, \mathbf{z}_{i-1})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})\big]$$
$$\tilde{f}_i(\mathbf{z}_i, \mathbf{z}_{i-1}) = \operatorname*{argmin}_{f \in \mathcal{N}} \mathsf{KL} \Big[f(\mathbf{z}_i, \mathbf{z}_{i-1}) q_{\neg i}(\mathbf{z}_i, \mathbf{z}_{i-1}) \big\| f(\mathbf{z}_i, \mathbf{z}_{i-1}) q_{\neg i}(\mathbf{z}_i, \mathbf{z}_{i-1}) \Big]$$

$$\tilde{f}_i(\mathbf{z}_i, \mathbf{z}_{i-1}) = \operatorname*{argmin}_{f \in \mathcal{N}} \mathsf{KL}\big[\phi_i(\mathbf{z}_i, \mathbf{z}_{i-1})\psi_i(\mathbf{z}_i)\alpha_{i-1}(\mathbf{z}_{i-1})\beta_i(\mathbf{z}_i)\big\|f(\mathbf{z}_i, \mathbf{z}_{i-1})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_i(\mathbf{z}_i)\big]$$



$$\tilde{f}_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) = \operatorname*{argmin}_{f \in \mathcal{N}} \mathsf{KL}\left[\underbrace{\phi_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1})\psi_{i}(\mathbf{z}_{i})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})}_{\widehat{P}(\mathbf{z}_{i-1}, \mathbf{z}_{i})} \middle\| \underbrace{f(\mathbf{z}_{i}, \mathbf{z}_{i-1})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})}_{P(\mathbf{z}_{i-1}, \mathbf{z}_{i})}\right]$$



$$\tilde{f}_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) = \operatorname*{argmin}_{t \in \mathcal{N}} \mathsf{KL}\left[\underbrace{\phi_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1})\psi_{i}(\mathbf{z}_{i})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})}_{\widehat{P}(\mathbf{z}_{i-1}, \mathbf{z}_{i})} \middle\| \underbrace{f(\mathbf{z}_{i}, \mathbf{z}_{i-1})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})}_{P(\mathbf{z}_{i-1}, \mathbf{z}_{i})}\right]$$

 $\tilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_i) = \underset{P \in \mathcal{N}}{\operatorname{argmin}} \operatorname{KL}[\widehat{P}(\mathbf{z}_{i-1}, \mathbf{z}_i) \| P(\mathbf{z}_{i-1}, \mathbf{z}_i)]$



$$\widetilde{f}_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) = \underset{f \in \mathcal{N}}{\operatorname{argmin}} \operatorname{KL}\left[\underbrace{\phi_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1})\psi_{i}(\mathbf{z}_{i})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})}_{\widehat{P}(\mathbf{z}_{i-1}, \mathbf{z}_{i})} \middle\| \underbrace{f(\mathbf{z}_{i}, \mathbf{z}_{i-1})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})}_{P(\mathbf{z}_{i-1}, \mathbf{z}_{i})}\right]$$

$$\widetilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_{i}) = \underset{P \in \mathcal{N}}{\operatorname{argmin}} \operatorname{KL}\left[\widehat{P}(\mathbf{z}_{i-1}, \mathbf{z}_{i}) \middle\| P(\mathbf{z}_{i-1}, \mathbf{z}_{i})\right] \qquad \widetilde{f}_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) = \frac{\widetilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_{i})}{\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})}$$



$$\tilde{f}_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) = \underset{i \in \mathcal{N}}{\operatorname{argmin}} \operatorname{KL}\left[\underbrace{\phi_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1})\psi_{i}(\mathbf{z}_{i})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})}_{\widehat{P}(\mathbf{z}_{i-1}, \mathbf{z}_{i})} \middle\| \underbrace{f(\mathbf{z}_{i}, \mathbf{z}_{i-1})\alpha_{i-1}(\mathbf{z}_{i-1})\beta_{i}(\mathbf{z}_{i})}_{P(\mathbf{z}_{i-1}, \mathbf{z}_{i})}\right]$$

$$\tilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_i) = \operatorname*{argmin}_{P \in \mathcal{N}} \mathsf{KL}\big[\widehat{P}(\mathbf{z}_{i-1}, \mathbf{z}_i)\big\| P(\mathbf{z}_{i-1}, \mathbf{z}_i)\big] \qquad \tilde{f}_i(\mathbf{z}_i, \mathbf{z}_{i-1}) = \frac{\tilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_i)}{\alpha_{i-1}(\mathbf{z}_{i-1})\beta_i(\mathbf{z}_i)}$$

$$\alpha_{i}(\mathbf{z}_{i}) = \int_{\mathbf{z}_{1}...\mathbf{z}_{i-1}} \prod_{i' < i+1} \tilde{f}_{i'}(\mathbf{z}_{i'}, \mathbf{z}_{i'-1}) = \int_{\mathbf{z}_{i-1}} \alpha_{i-1}(\mathbf{z}_{i-1}) \tilde{f}_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) = \frac{1}{\beta_{i}(\mathbf{z}_{i})} \int_{\mathbf{z}_{i-1}} \tilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_{i})$$

$$\beta_{i-1}(\mathbf{z}_{i-1}) = \int_{\mathbf{z}_{i+1}...\mathbf{z}_{i}} \prod_{i' > i} \tilde{f}_{i'}(\mathbf{z}_{i'}, \mathbf{z}_{i'-1}) = \int_{\mathbf{z}_{i}} \beta_{i}(\mathbf{z}_{i}) \tilde{f}_{i}(\mathbf{z}_{i}, \mathbf{z}_{i-1}) = \frac{1}{\alpha_{i-1}(\mathbf{z}_{i-1})} \int_{\mathbf{z}_{i}} \tilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_{i})$$





Moment Matching

Each EP update involves a KL minimisation:

$$\tilde{t}_{i}^{\text{new}}(\mathcal{Z}) \leftarrow \underset{t \in \{\tilde{l}\}}{\operatorname{argmin}} \operatorname{\mathsf{KL}}[f_{i}(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z}) \| f(\mathcal{Z}_{i})q_{\neg i}(\mathcal{Z})]$$

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

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

So minimum is found by matching sufficient stats or moment matching.

How do we calculate $\langle T(x) \rangle_{\widehat{P}}$?

How do we calculate $\langle T(x) \rangle_{\widehat{P}}$?

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

Quadrature methods.

How do we calculate $\langle T(x) \rangle_{\widehat{P}}$?

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.

How do we calculate $\langle T(x) \rangle_{\widehat{P}}$?

- 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.

How do we calculate $\langle T(x) \rangle_{\widehat{P}}$?

- 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.</p>

How do we calculate $\langle T(x) \rangle_{\widehat{p}}$?

- 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.</p>
 - Other quadrature approaches (e.g. GP quadrature) may be more accurate, and may allow formal constraint to pos-def cone.

How do we calculate $\langle T(x) \rangle_{\widehat{p}}$?

- 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.</p>
 - Other quadrature approaches (e.g. GP quadrature) may be more accurate, and may allow formal constraint to pos-def cone.
- Laplace approximation.

How do we calculate $\langle T(x) \rangle_{\widehat{p}}$?

- 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.</p>
 - Other quadrature approaches (e.g. GP quadrature) may be more accurate, and may allow formal constraint to pos-def cone.
- Laplace approximation.
 - Equivalent to Laplace propagation.

How do we calculate $\langle T(x) \rangle_{\widehat{p}}$?

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.</p>
 - Other quadrature approaches (e.g. GP quadrature) may be more accurate, and may allow formal constraint to pos-def cone.

Laplace approximation.

- Equivalent to Laplace propagation.
- As long as messages remain positive definite will converge to global Laplace approximation.

EP provides a succesful framework for Gaussian-process modelling of non-Gaussian observations (*e.g.* for classification).

EP provides a succesful 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 .

EP provides a succesful framework for Gaussian-process modelling of non-Gaussian observations (*e.g.* for classification).



- 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: $(\mu = [\mu(\mathbf{x}_i)], \Sigma = [\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)]).$

EP provides a succesful framework for Gaussian-process modelling of non-Gaussian observations (*e.g.* for classification).



- 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: $(\boldsymbol{\mu} = [\boldsymbol{\mu}(\mathbf{x}_i)], \boldsymbol{\Sigma} = [\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)]).$
- If we think of the gs as function values, a GP provides a prior over functions.

EP provides a succesful framework for Gaussian-process modelling of non-Gaussian observations (*e.g.* for classification).



- 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: $(\boldsymbol{\mu} = [\boldsymbol{\mu}(\mathbf{x}_i)], \boldsymbol{\Sigma} = [\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)]).$
- If we think of the gs 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 .

EP provides a succesful framework for Gaussian-process modelling of non-Gaussian observations (*e.g.* for classification).



- 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: $(\boldsymbol{\mu} = [\boldsymbol{\mu}(\mathbf{x}_i)], \boldsymbol{\Sigma} = [\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)]).$
- ▶ 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 μ = 0, and matrix Σ incorporates diagonal noise]

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



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

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



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

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



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

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

▶ The same factorisation applies to non-Gaussian $P(y_i|g_i)$ (e.g. $P(y_i=1) = 1/(1 + e^{-g_i})$).



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

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

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 f̃_i(g_i) = N(µ̃_i, ψ̃²_i).



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

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

- ► 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 $\tilde{f}_i(g_i) = \mathcal{N}(\tilde{\mu}_i, \tilde{\psi}_i^2)$.

• $q_{\neg i}(g_i)$ can be constructed by the usual GP marginalisation. If $\Sigma = K + \text{diag} \left[\tilde{\psi}_1^2 \dots \tilde{\psi}_n^2 \right]$

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



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

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

- 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 $\tilde{f}_i(g_i) = \mathcal{N}(\tilde{\mu}_i, \tilde{\psi}_i^2)$.
- $q_{\neg i}(g_i)$ can be constructed by the usual GP marginalisation. If $\Sigma = K + \text{diag} \left[\tilde{\psi}_1^2 \dots \tilde{\psi}_n^2 \right]$

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

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





Once appoximate site potentials have stabilised, they can be used to make predictions.

Introducing a test point changes K, but does not affect the marginal P(g₁...g_n) (by consistency of the GP).



- Introducing a test point changes K, but does not affect the marginal P(g₁...g_n) (by consistency of the GP).
- ▶ The unobserved output factor provides no information about g' (⇒ constant factor on g')



- Introducing a test point changes K, but does not affect the marginal P(g₁...g_n) (by consistency of the GP).
- The unobserved output factor provides no information about $g' (\Rightarrow \text{ constant factor on } g')$
- Thus no change is needed to the approximating potentials f_i.



- Introducing a test point changes K, but does not affect the marginal P(g₁...g_n) (by consistency of the GP).
- ▶ The unobserved output factor provides no information about g' (⇒ constant factor on g')
- Thus no change is needed to the approximating potentials f_i.
- Predictions are obtained by marginalising the approximation: [let $\tilde{\Psi} = \text{diag} \left| \tilde{\psi}_1^2 \dots \tilde{\psi}_n^2 \right|$]

$$\begin{split} \mathcal{P}(y'|\mathbf{x}',\mathcal{D}) &= \int dg' \, \mathcal{P}(y'|g') \mathcal{N}\Big(g' \mid \mathcal{K}_{x',X} (\mathcal{K}_{X,X} + \tilde{\Psi})^{-1} \tilde{\mu}, \\ & \mathcal{K}_{x',x'} - \mathcal{K}_{x',X} (\mathcal{K}_{X,X} + \tilde{\Psi})^{-1} \mathcal{K}_{X,x'}\Big) \end{split}$$

Normalisers

As long as our approximating class is a tractable exponential family, normalisers can be computed as needed.

Normalisers

- As long as our approximating class is a tractable exponential family, normalisers can be computed as needed.
- Consider an approximating class written

$$\widetilde{f}_i(\mathcal{Z}_i) \propto e^{\mathbf{T}(\mathcal{Z}) \cdot \boldsymbol{\theta}_i - \Phi(\boldsymbol{\theta}_i)}$$

i.e., define a single sufficient statistic vector on all latents, setting entries in θ_i to 0 for suff stat functions that take cliques other than Z_i .

Normalisers

- As long as our approximating class is a tractable exponential family, normalisers can be computed as needed.
- Consider an approximating class written

$$\widetilde{f}_i(\mathcal{Z}_i) \propto e^{\mathbf{T}(\mathcal{Z}) \cdot \boldsymbol{ heta}_i - \Phi(\boldsymbol{ heta}_i)}$$

i.e., define a single sufficient statistic vector on all latents, setting entries in θ_i to 0 for suff stat functions that take cliques other than Z_i .

Then

$$q(\mathcal{Z}) \propto \prod_{i} \widetilde{f}_{i} \propto e^{T(\mathcal{Z}) \cdot \sum \theta_{i} - \sum \Phi(\theta_{i})}$$

and so we can simply renormalise at the end as usual:

$$q(\mathcal{Z}) = e^{T(\mathcal{Z}) \cdot \sum \theta_i - \Phi(\sum \theta_i)}.$$
Normalisers

- As long as our approximating class is a tractable exponential family, normalisers can be computed as needed.
- Consider an approximating class written

$$\widetilde{f}_i(\mathcal{Z}_i) \propto e^{\mathbf{T}(\mathcal{Z}) \cdot \boldsymbol{ heta}_i - \Phi(\boldsymbol{ heta}_i)}$$

i.e., define a single sufficient statistic vector on all latents, setting entries in θ_i to 0 for suff stat functions that take cliques other than Z_i .

Then

$$q(\mathcal{Z}) \propto \prod_{i} \widetilde{f}_{i} \propto e^{T(\mathcal{Z}) \cdot \sum \theta_{i} - \sum \Phi(\theta_{i})}$$

and so we can simply renormalise at the end as usual:

$$q(\mathcal{Z}) = e^{T(\mathcal{Z}) \cdot \sum \theta_i - \Phi(\sum \theta_i)}.$$

► However, to compute an approximation to the likelihood $\int dZ \prod_i f_i(Z_i)$ we need to keep track of the site integrals.

Computing likelihoods – keeping track of normalisers

• Define unnormalised ExpFam approximating sites $\tilde{f}_i = \tilde{C}_i e^{\mathsf{T}(\mathcal{Z}) \cdot \theta_i}$.

Write $\theta = \sum \theta_j$ for the natural parameters of $q(\mathcal{Z})$ and $\theta_{\neg i} = \sum_{j \neq i} \theta_j$ for the natural parameters of $q_{\neg i}(\mathcal{Z})$.

Let $\Phi(\theta) = \log \int e^{T(\mathcal{Z}) \cdot \theta}$ be the (tractable) ExpFam log normaliser.

Computing likelihoods – keeping track of normalisers

• Define unnormalised ExpFam approximating sites $\tilde{f}_i = \tilde{C}_i e^{T(\mathcal{Z}) \cdot \theta_i}$.

Write $\theta = \sum \theta_j$ for the natural parameters of $q(\mathcal{Z})$ and $\theta_{\neg i} = \sum_{j \neq i} \theta_j$ for the natural parameters of $q_{\neg i}(\mathcal{Z})$.

Let $\Phi(\theta) = \log \int e^{T(\mathcal{Z}) \cdot \theta}$ be the (tractable) ExpFam log normaliser.

Now, at each EP step minimise the "unnormalised KL":

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

This matches the zeroth moment of $f_i(\mathcal{Z}_i)q_{\neg i}(\mathcal{Z})$ as well as the expected sufficient statistics as before. That is:

$$\int \tilde{C}_i \boldsymbol{e}^{\mathsf{T}(\mathcal{Z}) \cdot \theta_i} \prod_{\neg i} \tilde{C}_j \boldsymbol{e}^{\mathsf{T}(\mathcal{Z}) \cdot \theta_j} = \int f_i(\mathcal{Z}_i) \prod_{\neg i} \tilde{C}_j \boldsymbol{e}^{\mathsf{T}(\mathcal{Z}) \cdot \theta_j} \quad \Rightarrow \quad \tilde{C}_i = \boldsymbol{e}^{\Phi_i(\theta_{\neg i}) - \Phi(\theta)}$$

where Φ_i is the log-normaliser of the "tilted" ExpFam $\widehat{P}_i(\mathcal{Z}) \propto f(\mathcal{Z}_i) e^{T(\mathcal{Z}) \cdot \theta}$.

Computing likelihoods – keeping track of normalisers

• Define unnormalised ExpFam approximating sites $\tilde{f}_i = \tilde{C}_i e^{T(\mathcal{Z}) \cdot \theta_i}$.

Write $\theta = \sum \theta_j$ for the natural parameters of $q(\mathcal{Z})$ and $\theta_{\neg i} = \sum_{j \neq i} \theta_j$ for the natural parameters of $q_{\neg i}(\mathcal{Z})$.

Let $\Phi(\theta) = \log \int e^{T(\mathcal{Z}) \cdot \theta}$ be the (tractable) ExpFam log normaliser.

Now, at each EP step minimise the "unnormalised KL":

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

This matches the zeroth moment of $f_i(\mathcal{Z}_i)q_{\neg i}(\mathcal{Z})$ as well as the expected sufficient statistics as before. That is:

$$\int \tilde{C}_i \boldsymbol{e}^{\mathsf{T}(\boldsymbol{\mathcal{Z}}) \cdot \boldsymbol{\theta}_i} \prod_{\neg i} \tilde{C}_j \boldsymbol{e}^{\mathsf{T}(\boldsymbol{\mathcal{Z}}) \cdot \boldsymbol{\theta}_j} = \int f_i(\boldsymbol{\mathcal{Z}}_i) \prod_{\neg i} \tilde{C}_j \boldsymbol{e}^{\mathsf{T}(\boldsymbol{\mathcal{Z}}) \cdot \boldsymbol{\theta}_j} \quad \Rightarrow \quad \tilde{C}_i = \boldsymbol{e}^{\Phi_i(\boldsymbol{\theta}_{\neg i}) - \Phi(\boldsymbol{\theta})}$$

where Φ_i is the log-normaliser of the "tilted" ExpFam $\widehat{P}_i(\mathcal{Z}) \propto f(\mathcal{Z}_i)e^{\mathbf{T}(\mathcal{Z})\cdot\theta}$. The likelihood approximation is then:

$$\log \int \prod_{i=1}^{N} f_i(\mathcal{Z}_i) \approx \log \int \prod_{i=1}^{N} \tilde{f}_i(\mathcal{Z}_i) = \Phi(\theta) + \sum \log \tilde{C}_i \stackrel{\text{def}}{=} \tilde{\ell}$$

EP yields approximate *inferential* posteriors. To learn (hyper)parameters we can use:

Approximate Bayesian inference (analagous to VB)

- Approximate Bayesian inference (analagous to VB)
 - may be difficult to construct a coherent normalisable exponential family approximation on both latents and parameters.

EP yields approximate *inferential* posteriors. To learn (hyper)parameters we can use:

- Approximate Bayesian inference (analagous to VB)
 - may be difficult to construct a coherent normalisable exponential family approximation on both latents and parameters.

• Approximate EM – maximize $\langle \log P(\mathcal{X}, \mathcal{Z}) \rangle_{q_{FP}(\mathcal{Z})}$.

- Approximate Bayesian inference (analagous to VB)
 - may be difficult to construct a coherent normalisable exponential family approximation on both latents and parameters.
- Approximate EM maximize $\langle \log P(\mathcal{X}, \mathcal{Z}) \rangle_{q_{FP}(\mathcal{Z})}$.
 - Practical, but no coherent cost function (unlike variational inference), so no guarantee of convergence even if EP itself converges.

- Approximate Bayesian inference (analagous to VB)
 - may be difficult to construct a coherent normalisable exponential family approximation on both latents and parameters.
- Approximate EM maximize $\langle \log P(\mathcal{X}, \mathcal{Z}) \rangle_{q_{EP}(\mathcal{Z})}$.
 - Practical, but no coherent cost function (unlike variational inference), so no guarantee of convergence even if EP itself converges.
- Direct maximisation of EP log-likelihood estimate.

- Approximate Bayesian inference (analagous to VB)
 - may be difficult to construct a coherent normalisable exponential family approximation on both latents and parameters.
- Approximate EM maximize $\langle \log P(\mathcal{X}, \mathcal{Z}) \rangle_{q_{FP}(\mathcal{Z})}$.
 - Practical, but no coherent cost function (unlike variational inference), so no guarantee of convergence even if EP itself converges.
- Direct maximisation of EP log-likelihood estimate.
 - Consistent, although convergence guarantees still difficult.

- Approximate Bayesian inference (analagous to VB)
 - may be difficult to construct a coherent normalisable exponential family approximation on both latents and parameters.
- Approximate EM maximize $\langle \log P(\mathcal{X}, \mathcal{Z}) \rangle_{q_{FP}(\mathcal{Z})}$.
 - Practical, but no coherent cost function (unlike variational inference), so no guarantee of convergence even if EP itself converges.
- Direct maximisation of EP log-likelihood estimate.
 - Consistent, although convergence guarantees still difficult.
 - Seems challenging as we need to differentiate through (iteration-based) dependence of approximate q(Z) and C
 _{is}.

- Approximate Bayesian inference (analagous to VB)
 - may be difficult to construct a coherent normalisable exponential family approximation on both latents and parameters.
- Approximate EM maximize $\langle \log P(\mathcal{X}, \mathcal{Z}) \rangle_{q_{FP}(\mathcal{Z})}$.
 - Practical, but no coherent cost function (unlike variational inference), so no guarantee of convergence even if EP itself converges.
- Direct maximisation of EP log-likelihood estimate.
 - Consistent, although convergence guarantees still difficult.
 - Seems challenging as we need to differentiate through (iteration-based) dependence of approximate q(Z) and C
 _is.
 - However, proves to be simpler than it sounds.

Let true potentials f_i depend on model (hyper)parameters η .

Let true potentials f_i depend on model (hyper)parameters η . We have

$$abla_\eta ilde{\ell} =
abla_\eta \Phi(heta) + \sum_{i=1}^N
abla_\eta \log ilde{C}_i$$

Let true potentials f_i depend on model (hyper)parameters η . We have

$$\nabla_{\eta}\tilde{\ell} = \nabla_{\eta}\Phi(\theta) + \sum_{i=1}^{N} \nabla_{\eta}\log\tilde{C}_{i} = \mu \cdot \nabla_{\eta}\theta + \sum_{i=1}^{N} \nabla_{\eta}\log\tilde{C}_{i} \tag{*}$$

using the standard ExpFam moment-generating result with mean parameters $\mu = \langle T(Z) \rangle_{q(Z)}$.

Let true potentials f_i depend on model (hyper)parameters η . We have

$$abla_\eta ilde{\ell} =
abla_\eta \Phi(heta) + \sum_{i=1}^N
abla_\eta \log ilde{C}_i = \mu \cdot
abla_\eta heta + \sum_{i=1}^N
abla_\eta \log ilde{C}_i \tag{*}$$

using the standard ExpFam moment-generating result with mean parameters $\mu = \langle T(\mathcal{Z}) \rangle_{q(\mathcal{Z})}.$

Now, zeroth-moment matching implies that at EP convergence:

$$\log \tilde{C}_i = \Phi_i(\theta_{\neg i}) - \Phi(\theta)$$

Let true potentials f_i depend on model (hyper)parameters η . We have

$$\nabla_{\eta}\tilde{\ell} = \nabla_{\eta}\Phi(\theta) + \sum_{i=1}^{N} \nabla_{\eta}\log\tilde{C}_{i} = \mu \cdot \nabla_{\eta}\theta + \sum_{i=1}^{N} \nabla_{\eta}\log\tilde{C}_{i}$$
(*)

using the standard ExpFam moment-generating result with mean parameters $\mu = \langle T(\mathcal{Z}) \rangle_{q(\mathcal{Z})}$.

Now, zeroth-moment matching implies that at EP convergence:

$$\log \tilde{C}_i = \Phi_i(\theta_{\neg i}) - \Phi(\theta) \Rightarrow \nabla_\eta \log \tilde{C}_i = \nabla_\eta \Phi_i(\theta_{\neg i}) - \mu \cdot \nabla_\eta \theta \qquad (^{**})$$

Let true potentials f_i depend on model (hyper)parameters η . We have

$$\nabla_{\eta} \tilde{\ell} = \nabla_{\eta} \Phi(\theta) + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i} = \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i}$$
(*)

using the standard ExpFam moment-generating result with mean parameters $\mu = \langle T(\mathcal{Z}) \rangle_{q(\mathcal{Z})}$.

Now, zeroth-moment matching implies that at EP convergence:

$$\log \tilde{C}_i = \Phi_i(\theta_{\neg i}) - \Phi(\theta) \Rightarrow \nabla_\eta \log \tilde{C}_i = \nabla_\eta \Phi_i(\theta_{\neg i}) - \mu \cdot \nabla_\eta \theta \qquad (^{**})$$

but $\Phi_i(\theta_{\neg i}) = \log \int f_i(\mathcal{Z}_i) e^{\mathbf{T}(\mathcal{Z}) \cdot \theta_{\neg i}}$ depends on η in two ways: *directly* through f_i and *indirectly* through the converged $\theta_{\neg i}$.

Let true potentials f_i depend on model (hyper)parameters η . We have

$$\nabla_{\eta} \tilde{\ell} = \nabla_{\eta} \Phi(\theta) + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i} = \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i}$$
(*)

using the standard ExpFam moment-generating result with mean parameters $\mu = \langle T(\mathcal{Z}) \rangle_{q(\mathcal{Z})}$.

Now, zeroth-moment matching implies that at EP convergence:

$$\log \tilde{C}_i = \Phi_i(\theta_{\neg i}) - \Phi(\theta) \Rightarrow \nabla_\eta \log \tilde{C}_i = \nabla_\eta \Phi_i(\theta_{\neg i}) - \mu \cdot \nabla_\eta \theta \qquad (**)$$

but $\Phi_i(\theta_{\neg i}) = \log \int f_i(\mathcal{Z}_i) e^{\mathbf{T}(\mathcal{Z}) \cdot \theta_{\neg i}}$ depends on η in two ways: *directly* through f_i and *indirectly* through the converged $\theta_{\neg i}$.

$$\nabla_{\eta} \Phi_{i}(\theta_{\neg i}) = \partial_{\theta_{\neg i}} \Phi_{i}(\theta_{\neg i}) \cdot \nabla_{\eta} \theta_{\neg i} + e^{-\Phi_{i}(\theta_{\neg i})} \int \nabla_{\eta} f_{i}(\mathcal{Z}_{i}) \ e^{\mathbf{T}(\mathcal{Z}) \cdot \theta_{\neg i}} \ d\mathcal{Z}$$

Let true potentials f_i depend on model (hyper)parameters η . We have

$$\nabla_{\eta} \tilde{\ell} = \nabla_{\eta} \Phi(\theta) + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i} = \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i}$$
(*)

using the standard ExpFam moment-generating result with mean parameters $\mu = \langle T(Z) \rangle_{q(Z)}$. Now, zeroth-moment matching implies that at EP convergence:

 $\log \tilde{C}_i = \Phi_i(\theta_{\neg i}) - \Phi(\theta) \Rightarrow \nabla_\eta \log \tilde{C}_i = \nabla_\eta \Phi_i(\theta_{\neg i}) - \mu \cdot \nabla_\eta \theta \tag{**}$

but $\Phi_i(\theta_{\neg i}) = \log \int f_i(\mathcal{Z}_i) e^{\mathbf{T}(\mathcal{Z}) \cdot \theta_{\neg i}}$ depends on η in two ways: *directly* through f_i and *indirectly* through the converged $\theta_{\neg i}$.

$$\begin{split} \nabla_{\eta} \Phi_{i}(\theta_{\neg i}) &= \partial_{\theta_{\neg i}} \Phi_{i}(\theta_{\neg i}) \cdot \nabla_{\eta} \theta_{\neg i} + e^{-\Phi_{i}(\theta_{\neg i})} \int \nabla_{\eta} f_{i}(\mathcal{Z}_{i}) \ e^{\mathbf{T}(\mathcal{Z}) \cdot \theta_{\neg i}} \, d\mathcal{Z} \\ &= \langle T(\mathcal{Z}) \rangle_{\widehat{\mathbf{P}}_{i}} \cdot \nabla_{\eta} \theta_{\neg i} + \int \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \ f_{i}(\mathcal{Z}_{i}) e^{\mathbf{T}(\mathcal{Z}) \cdot \theta_{\neg i} - \Phi_{i}(\theta_{\neg i})} \, d\mathcal{Z} \end{split}$$

Let true potentials f_i depend on model (hyper)parameters η . We have

$$\nabla_{\eta} \tilde{\ell} = \nabla_{\eta} \Phi(\theta) + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i} = \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i}$$
(*)

using the standard ExpFam moment-generating result with mean parameters $\mu = \langle T(Z) \rangle_{q(Z)}$. Now, zeroth-moment matching implies that at EP convergence:

 $\log \tilde{C}_i = \Phi_i(\theta_{\neg i}) - \Phi(\theta) \Rightarrow \nabla_\eta \log \tilde{C}_i = \nabla_\eta \Phi_i(\theta_{\neg i}) - \mu \cdot \nabla_\eta \theta \tag{**}$

but $\Phi_i(\theta_{\neg i}) = \log \int f_i(\mathcal{Z}_i) e^{\mathbf{T}(\mathcal{Z}) \cdot \theta_{\neg i}}$ depends on η in two ways: *directly* through f_i and *indirectly* through the converged $\theta_{\neg i}$.

by EP moment matching at convergence!

(*)

$$abla_\eta ilde{\ell} = oldsymbol{\mu} \cdot
abla_\eta oldsymbol{ heta} + \sum_{i=1}^N
abla_\eta \log ilde{C}_i$$

$$\nabla_{\eta} \tilde{\ell} = \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i}$$

$$= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\nabla_{\eta} \Phi_{i}(\theta_{\neg i}) - \mu \cdot \nabla_{\eta} \theta \right)$$
(**)

$$\begin{aligned} \nabla_{\eta} \tilde{\ell} &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i} \end{aligned} \tag{*} \\ &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\nabla_{\eta} \Phi_{i}(\theta_{\neg i}) - \mu \cdot \nabla_{\eta} \theta \right) \end{aligned} \tag{*} \\ &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\mu \cdot \nabla_{\eta} \theta_{\neg i} - \mu \cdot \nabla_{\eta} \theta + \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{P}_{i}} \right) \end{aligned} \tag{*}$$

$$\begin{aligned} \nabla_{\eta} \tilde{\ell} &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i} \end{aligned} \tag{*} \\ &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\nabla_{\eta} \Phi_{i}(\theta_{\neg i}) - \mu \cdot \nabla_{\eta} \theta \right) \end{aligned} \tag{*} \\ &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\mu \cdot \nabla_{\eta} \theta_{\neg i} - \mu \cdot \nabla_{\eta} \theta + \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{\rho}_{i}} \right) \end{aligned} \tag{*} \\ &= \mu \cdot \nabla_{\eta} \left(\theta + \sum_{i=1}^{N} (\theta_{\neg i} - \theta) \right) + \sum_{i=1}^{N} \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{\rho}_{i}} \end{aligned}$$

i=1

$$\nabla_{\eta} \tilde{\ell} = \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i}$$

$$= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\nabla_{\eta} \Phi_{i}(\theta_{\neg i}) - \mu \cdot \nabla_{\eta} \theta \right)$$
(**)

$$= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\mu \cdot \nabla_{\eta} \theta_{\neg i} - \mu \cdot \nabla_{\eta} \theta + \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{\mathbf{P}}_{i}} \right)$$
(***)

$$= \mu \cdot \nabla_{\eta} \Big(\theta + \sum_{i=1}^{N} (\theta_{\neg i} - \theta) \Big) + \sum_{i=1}^{N} \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{P}_{i}}$$
$$= \mu \cdot \nabla_{\eta} \Big(\sum_{i=1}^{N} \theta_{i} + \sum_{i=1}^{N} (\theta_{\neg i} - \theta) \Big) + \sum_{i=1}^{N} \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{P}_{i}}$$

$$\begin{split} \nabla_{\eta} \tilde{\ell} &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i} \qquad (*) \\ &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\nabla_{\eta} \Phi_{i}(\theta_{\neg i}) - \mu \cdot \nabla_{\eta} \theta \right) \qquad (**) \\ &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\mu \cdot \nabla_{\eta} \theta_{\neg i} - \mu \cdot \nabla_{\eta} \theta + \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{P}_{i}} \right) \qquad (***) \\ &= \mu \cdot \nabla_{\eta} \left(\theta + \sum_{i=1}^{N} (\theta_{\neg i} - \theta) \right) + \sum_{i=1}^{N} \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{P}_{i}} \\ &= \mu \cdot \nabla_{\eta} \left(\sum_{i=1}^{N} \theta_{i} + \sum_{i=1}^{N} (\theta_{\neg i} - \theta) \right) + \sum_{i=1}^{N} \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{P}_{i}} \end{split}$$

$$= \mu \cdot \nabla_\eta \sum_{i=1}^N (\theta - \theta) + \sum_{i=1}^N \langle \nabla_\eta \log f_i(\mathcal{Z}_i) \rangle_{\widehat{\boldsymbol{P}}_i}$$

So putting it all together:

$$\begin{split} \nabla_{\eta} \tilde{\ell} &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \nabla_{\eta} \log \tilde{C}_{i} \qquad (*) \\ &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\nabla_{\eta} \Phi_{i}(\theta_{\neg i}) - \mu \cdot \nabla_{\eta} \theta \right) \qquad (*) \\ &= \mu \cdot \nabla_{\eta} \theta + \sum_{i=1}^{N} \left(\mu \cdot \nabla_{\eta} \theta_{\neg i} - \mu \cdot \nabla_{\eta} \theta + \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{P}_{i}} \right) \qquad (*) \\ &= \mu \cdot \nabla_{\eta} \left(\theta + \sum_{i=1}^{N} (\theta_{\neg i} - \theta) \right) + \sum_{i=1}^{N} \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{P}_{i}} \\ &= \mu \cdot \nabla_{\eta} \left(\sum_{i=1}^{N} \theta_{i} + \sum_{i=1}^{N} (\theta_{\neg i} - \theta) \right) + \sum_{i=1}^{N} \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{P}_{i}} \\ &= \mu \cdot \nabla_{\eta} \sum_{i=1}^{N} (\theta - \theta) + \sum_{i=1}^{N} \langle \nabla_{\eta} \log f_{i}(\mathcal{Z}_{i}) \rangle_{\widehat{P}_{i}} \end{split}$$

and the gradient can be computed provided EP converges.

Alpha divergences

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

....

Alpha divergences

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

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

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

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

$$\lim_{\alpha \to 1} D_{\alpha}[p||q] = \mathsf{KL}[p||q]$$

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

. .

Alpha divergences

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

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

$$\lim_{\alpha \to 0} D_{\alpha}[p||q] = \mathsf{KL}[q||p]$$

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

$$\lim_{\alpha \to 1} D_{\alpha}[p||q] = \mathsf{KL}[p||q]$$

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

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

$$\widetilde{f}_{i}^{\text{new}} = \operatorname*{argmin}_{f \in \{\overline{l}\}} \mathsf{KL} \Big[f_{i}(\mathcal{Z}_{i})^{\alpha} \widetilde{f}_{i}(\mathcal{Z}_{i})^{1-\alpha} q_{\neg i}(\mathcal{Z}) \Big\| f(\mathcal{Z}_{i}) q_{\neg i}(\mathcal{Z}) \Big]$$

Alpha divergences

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

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

$$\lim_{\alpha \to 0} D_{\alpha}[p||q] = \mathsf{KL}[q||p]$$

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

$$\lim_{\alpha \to 1} D_{\alpha}[p||q] = \mathsf{KL}[p||q]$$

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

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

$$\widetilde{f}_{i}^{\text{new}} = \operatorname*{argmin}_{f \in \{\overline{l}\}} \mathsf{KL}\big[f_{i}(\mathcal{Z}_{i})^{\alpha} \widetilde{f}_{i}(\mathcal{Z}_{i})^{1-\alpha} q_{\neg i}(\mathcal{Z}) \big\| f(\mathcal{Z}_{i}) q_{\neg i}(\mathcal{Z})\big]$$

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