

Variational Methods ~~Approximate~~ Approximate

↑  
not intrinsically

Probabilistic & Unsupervised Learning

But naturally leads to Approximate ; find the best class of function to fit posterior

Factored Variational Approximations  
and Variational Bayes

---

KEY to variational approximation :

Remove as few arcs as possible  
from the moral graph s.t.  
inference becomes tractable .

Maneesh Sahani

maneesh@gatsby.ucl.ac.uk

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

Term 1, Autumn 2017

## Expectations in Statistical Modelling

- ▶ Parameter estimation

$$\hat{\theta} = \operatorname{argmax}_{\theta} \int d\mathcal{Y} P(\mathcal{Y}|\theta) P(\mathcal{X}|\mathcal{Y}, \theta)$$

(or, using EM)

$$\theta^{\text{new}} = \operatorname{argmax}_{\theta} \int d\mathcal{Y} P(\mathcal{Y}|\mathcal{X}, \theta^{\text{old}}) \log P(\mathcal{X}, \mathcal{Y}|\theta)$$

- ▶ Prediction

$$p(x|\mathcal{D}, m) = \int d\theta p(\theta|\mathcal{D}, m) p(x|\theta, \mathcal{D}, m)$$

- ▶ Model selection or weighting (by marginal likelihood)

$$p(\mathcal{D}|m) = \int d\theta p(\theta|m) p(\mathcal{D}|\theta, m)$$

These integrals are often intractable:

- ▶ **Analytic intractability:** integrals may not have closed form in non-linear, non-Gaussian models  $\Rightarrow$  numerical integration.
- ▶ **Computational intractability:** Numerical integral (or sum if  $\mathcal{Y}$  or  $\theta$  are discrete) may be exponential in data or model size.

## Examples of Intractability

Bishop.

- ▶ Marginal likelihood/model evidence for Mixture of Gaussians: exact computations are exponential in number of data points

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N) &= \int d\theta \, p(\theta) \prod_{i=1}^N \sum_{s_i} p(\mathbf{x}_i | s_i, \theta) p(s_i | \theta) \\ &= \sum_{s_1} \sum_{s_2} \dots \sum_{s_N} \int d\theta \, p(\theta) \prod_{i=1}^N p(\mathbf{x}_i | s_i, \theta) p(s_i | \theta) \end{aligned}$$

- ▶ Computing the conditional probabilities in a very large multiply-connected DAG:

$$p(x_i | X_j = a) = \sum_{\text{all settings of } \mathbf{y} \setminus \{i,j\}} p(x_i, \mathbf{y}, X_j = a) / p(X_j = a)$$

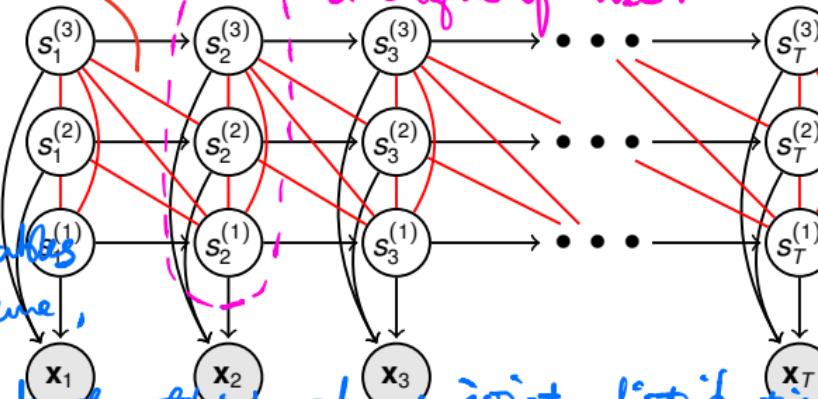
- ▶ Computing the hidden state distribution in a general nonlinear dynamical system

$$p(\mathbf{y}_t | \mathbf{x}_1, \dots, \mathbf{x}_t) \propto \int d\mathbf{y}_{t-1} p(\mathbf{y}_t | f(\mathbf{y}_{t-1})) p(\mathbf{x}_t | g(\mathbf{y}_t)) p(\mathbf{y}_{t-1} | \mathbf{x}_1, \dots, \mathbf{x}_{t-1})$$

## Distributed models

Red moralisation lines

## Factorial Hidden Markov Model



Consider a binary FHMM, with  $K=2$ . Then it takes  $2^M$  different states that the variables can take. When doing inference, we moralise, then we go back to think about joint distribution over  $K$  binary variables, which takes  $(2!)^M$

Consider an FHMM with  $M$  state variables taking on  $K$  values each.

- ▶ Moralisation puts simultaneous states  $(s_t^{(1)}, s_t^{(2)}, \dots, s_t^{(M)})$  into a single clique
- ▶ Triangulation extends cliques to size  $M + 1$
- ▶ Each state takes  $K$  values  $\Rightarrow$  sums over  $K^{M+1}$  terms.
- ▶ Factorial prior  $\neq$  Factorial posterior (explaining away).

Variational methods **approximate** the posterior, often in a factored form.

To see how they work, we need to review the free-energy interpretation of EM.

## The Free Energy for a Latent Variable Model

Observed data  $\mathcal{X} = \{\mathbf{x}_i\}$ ; Latent variables  $\mathcal{Y} = \{\mathbf{y}_i\}$ ; Parameters  $\theta$ .

**Goal:** Maximize the log likelihood wrt  $\theta$  (i.e. ML learning):

$$\ell(\theta) = \log P(\mathcal{X}|\theta) = \log \int P(\mathcal{Y}, \mathcal{X}|\theta) d\mathcal{Y}$$

Any distribution,  $q(\mathcal{Y})$ , over the hidden variables can be used to obtain a lower bound on the log likelihood using Jensen's inequality:

$$\ell(\theta) = \log \int q(\mathcal{Y}) \frac{P(\mathcal{Y}, \mathcal{X}|\theta)}{q(\mathcal{Y})} d\mathcal{Y} \geq \int q(\mathcal{Y}) \log \frac{P(\mathcal{Y}, \mathcal{X}|\theta)}{q(\mathcal{Y})} d\mathcal{Y} \stackrel{\text{def}}{=} \mathcal{F}(q, \theta)$$

$$\begin{aligned} \int q(\mathcal{Y}) \log \frac{P(\mathcal{Y}, \mathcal{X}|\theta)}{q(\mathcal{Y})} d\mathcal{Y} &= \int q(\mathcal{Y}) \log P(\mathcal{Y}, \mathcal{X}|\theta) d\mathcal{Y} - \int q(\mathcal{Y}) \log q(\mathcal{Y}) d\mathcal{Y} \\ &= \int q(\mathcal{Y}) \log P(\mathcal{Y}, \mathcal{X}|\theta) d\mathcal{Y} + \mathbf{H}[q], \end{aligned}$$

where  $\mathbf{H}[q]$  is the entropy of  $q(\mathcal{Y})$ .

So:  $\mathcal{F}(q, \theta) = \langle \log P(\mathcal{Y}, \mathcal{X}|\theta) \rangle_{q(\mathcal{Y})} + \mathbf{H}[q]$

## The E and M steps of EM

Recall:

maximising the lower

The log likelihood is bounded below by:

bound of free energy

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

is equivalent to minimizing KL divergence

EM alternates between:

**E step:** optimise  $\mathcal{F}(q, \theta)$  wrt distribution over hidden variables holding parameters fixed:

$$q^{(k)}(\mathcal{Y}) := \underset{q(\mathcal{Y})}{\operatorname{argmax}} \mathcal{F}(q(\mathcal{Y}), \theta^{(k-1)}) = P(\mathcal{Y} | \mathcal{X}, \theta^{(k-1)})$$

**M step:** maximise  $\mathcal{F}(q, \theta)$  wrt parameters holding hidden distribution fixed:

$$\theta^{(k)} := \underset{\theta}{\operatorname{argmax}} \mathcal{F}(q^{(k)}(\mathcal{Y}), \theta) = \underset{\theta}{\operatorname{argmax}} \langle \log P(\mathcal{Y}, \mathcal{X} | \theta) \rangle_{q^{(k)}(\mathcal{Y})}$$

## EM Never Decreases the Likelihood

The E and M steps together never decrease the log likelihood:

$$\ell(\theta^{(k-1)}) \underset{\text{E step}}{=} \mathcal{F}(q^{(k)}, \theta^{(k-1)}) \underset{\text{M step}}{\leq} \mathcal{F}(q^{(k)}, \theta^{(k)}) \underset{\text{Jensen}}{\leq} \ell(\theta^{(k)}),$$

- ▶ The E step brings the free energy to the likelihood.
- ▶ The M-step maximises the free energy wrt  $\theta$ .
- ▶  $\mathcal{F} \leq \ell$  by Jensen – or, equivalently, from the non-negativity of KL

If the M-step is executed so that  $\theta^{(k)} \neq \theta^{(k-1)}$  iff  $\mathcal{F}$  increases, then the overall EM iteration will step to a new value of  $\theta$  iff the likelihood increases.

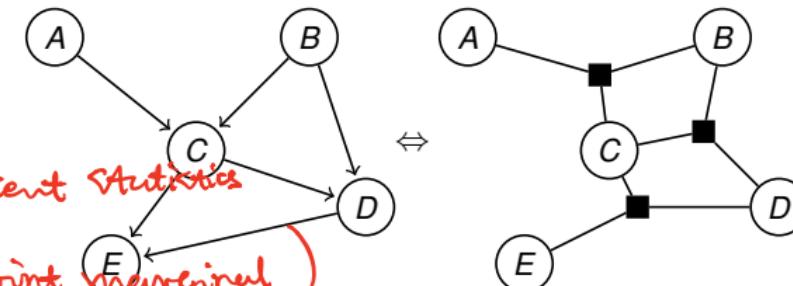
## Intractability

The M-step for a graphical model is usually (relatively) easy.

Hard thing :

Too many number of steps  
in inference

(Finding joint marginals / sufficient statistics  
of the joint marginal)



Easy thing:

Learning is quite straightforward.

$$P(A, B, C, D, E) = \underbrace{P(A)P(B)}_{f_1(A, B, C)} \underbrace{P(C|A, B)}_{f_2(B, C, D)} \underbrace{P(D|B, C)P(E|C, D)}_{f_3(C, D, E)}$$

- ▶ Need expected sufficient stats from marginal posteriors on each factor group.
- ▶ Then (at least for a DAG) can optimise each factor parameter vector separately.
- ▶ Intractability in EM comes from the difficulty of computing marginal posteriors in graphs with large tree-width or non-linear/non-conjugate conditionals.
- ▶ [For non-DAG models, partition function (normalising constant) may also be intractable.]

## Free-energy-based variational approximation

What if finding expected sufficient stats under  $P(\mathcal{Y}|\mathcal{X}, \theta)$  is computationally **intractable**?

For the **generalised EM** algorithm, we argued that intractable maximisations could be replaced by gradient M-steps.

- ▶ Each step increases the likelihood.
- ▶ A fixed point of the gradient M-step must be at a mode of the expected log-joint.

For the E-step we could:

- ▶ **Parameterise**  $q = q_\rho(\mathcal{Y})$  and take a gradient step in  $\rho$ .
- ▶ **Assume** some simplified form for  $q$ , usually **factored**:  $q = \prod_i q_i(\mathcal{Y}_i)$  where  $\mathcal{Y}_i$  partition  $\mathcal{Y}$ , and maximise within this form.

In either case, we choose  $q$  from within a limited set  $\mathcal{Q}$ :

**VE step:** maximise  $\mathcal{F}(q, \theta)$  wrt constrained latent distribution given parameters:

$$q^{(k)}(\mathcal{Y}) := \underset{q(\mathcal{Y}) \in \mathcal{Q} \leftarrow \text{Constraint}}{\operatorname{argmax}} \mathcal{F}(q(\mathcal{Y}), \theta^{(k-1)}).$$

**M step:** unchanged

$$\theta^{(k)} := \underset{\theta}{\operatorname{argmax}} \mathcal{F}(q^{(k)}(\mathcal{Y}), \theta) = \underset{\theta}{\operatorname{argmax}} \int q^{(k)}(\mathcal{Y}) \log p(\mathcal{Y}, \mathcal{X} | \theta) d\mathcal{Y},$$

Unlike in GEM, the fixed point may not be at an unconstrained optimum of  $\mathcal{F}$ .

## What do we lose?

What does restricting  $q$  to  $\mathcal{Q}$  cost us?

- Recall that the free-energy is bounded above by Jensen:

$$\mathcal{F}(q, \theta) \leq \ell(\theta^{\text{ML}})$$

Thus, as long as every step increases  $\mathcal{F}$ , convergence is still guaranteed.

- But, since  $P(\mathcal{Y}|\mathcal{X}, \theta^{(k)})$  may not lie in  $\mathcal{Q}$ , we no longer saturate the bound after the E-step. Thus, the likelihood may not increase on each full EM step.

$$\ell(\theta^{(k-1)}) \underset{\text{E step}}{\not\geq} \mathcal{F}(q^{(k)}, \theta^{(k-1)}) \underset{\text{M step}}{\leq} \mathcal{F}(q^{(k)}, \theta^{(k)}) \underset{\text{Jensen}}{\leq} \ell(\theta^{(k)}),$$

- This means we may not converge to a maximum of  $\ell$ .

The hope is that by increasing a lower bound on  $\ell$  we will find a decent solution.

[Note that if  $P(\mathcal{Y}|\mathcal{X}, \theta^{\text{ML}}) \in \mathcal{Q}$ , then  $\theta^{\text{ML}}$  is a fixed point of the variational algorithm.]

## KL divergence

Recall that

$$\begin{aligned}\mathcal{F}(q, \theta) &= \langle \log P(\mathcal{X}, \mathcal{Y}|\theta) \rangle_{q(\mathcal{Y})} + \mathbf{H}[q] \\ &= \langle \log P(\mathcal{X}|\theta) + \log P(\mathcal{Y}|\mathcal{X}, \theta) \rangle_{q(\mathcal{Y})} - \langle \log q(\mathcal{Y}) \rangle_{q(\mathcal{Y})} \\ &= \langle \log P(\mathcal{X}|\theta) \rangle_{q(\mathcal{Y})} - \mathbf{KL}[q||P(\mathcal{Y}|\mathcal{X}, \theta)].\end{aligned}$$

Thus,

E step maximise  $\mathcal{F}(q, \theta)$  wrt the distribution over latents, given parameters:

$$q^{(k)}(\mathcal{Y}) := \underset{q(\mathcal{Y}) \in \mathcal{Q}}{\operatorname{argmax}} \mathcal{F}(q(\mathcal{Y}), \theta^{(k-1)}).$$

is equivalent to:

E step minimise  $\mathbf{KL}[q||p(\mathcal{Y}|\mathcal{X}, \theta)]$  wrt distribution over latents, given parameters:

$$q^{(k)}(\mathcal{Y}) := \underset{q(\mathcal{Y}) \in \mathcal{Q}}{\operatorname{argmin}} \int q(\mathcal{Y}) \log \frac{q(\mathcal{Y})}{p(\mathcal{Y}|\mathcal{X}, \theta^{(k-1)})} d\mathcal{Y}$$

So, in each E step, the algorithm is trying to find the best approximation to  $P(\mathcal{Y}|\mathcal{X})$  in  $\mathcal{Q}$  in a KL sense. This is related to ideas in *information geometry*. It also suggests generalisations to other distance measures.

## Factored Variational E-step

The most common form of variational approximation partitions  $\mathcal{Y}$  into disjoint sets  $\mathcal{Y}_i$  with

$\mathcal{Q} = \{q \mid q(\mathcal{Y}) = \prod_i q_i(\mathcal{Y}_i)\}$ . Rather than jointly optimise the free energy w.r.t. all of the distributions, we think as iteratively update w.r.t just one of the distributions at a time. E-step itself

In this case the E-step is itself iterative:

(Factored VE step) $_i$ : maximise  $\mathcal{F}(q, \theta)$  wrt  $q_i(\mathcal{Y}_i)$  given other  $q_j$  and parameters: becomes iterative

$$q_i^{(k)}(\mathcal{Y}_i) := \underset{q_i(\mathcal{Y}_i)}{\operatorname{argmax}} \mathcal{F}(q_i(\mathcal{Y}_i) \prod_{j \neq i} q_j(\mathcal{Y}_j), \theta^{(k-1)})$$

Iteratively update E step to convergence,

then do M-step

Iterate one distribution each time.

Closed form updates for  
each one of the factors

- $q_i$  updates iterated to convergence to "complete" VE-step.
- In fact, every  $(VE)_i$ -step separately increases  $\mathcal{F}$ , so any schedule of  $(VE)_i$ - and M-steps will converge. Choice can be dictated by practical issues (rarely efficient to fully converge E-step before updating parameters).

See Ass5.  
Q1(a).

As long as we visit all factors

SEE

Matthew Beal  
2003  
Chapter 2.2

Target:  
Find functional form of  $q_i$   
that maximise the  
free energy.

Functional  
Derivative

## Factored Variational E-step

The Factored Variational E-step has a general form.

The free energy is:

$$\mathcal{F}\left(\prod_j q_j(\mathcal{Y}_j), \theta^{(k-1)}\right) = \left\langle \log P(\mathcal{X}, \mathcal{Y} | \theta^{(k-1)}) \right\rangle_{\prod_j q_j(\mathcal{Y}_j)} + \mathbf{H}\left[\prod_j q_j(\mathcal{Y}_j)\right]$$
$$= \int d\mathcal{Y}_i q_i(\mathcal{Y}_i) \left\langle \log P(\mathcal{X}, \mathcal{Y} | \theta^{(k-1)}) \right\rangle_{\prod_{j \neq i} q_j(\mathcal{Y}_j)} + \mathbf{H}[q_i] + \sum_{j \neq i} \mathbf{H}[q_j]$$

Now, taking the variational derivative of the Lagrangian (enforcing normalisation of  $q_i$ ):

$$\frac{\delta}{\delta q_i} \left( \mathcal{F} + \lambda \left( \int q_i - 1 \right) \right) = \left\langle \log P(\mathcal{X}, \mathcal{Y} | \theta^{(k-1)}) \right\rangle_{\prod_{j \neq i} q_j(\mathcal{Y}_j)} - \log q_i(\mathcal{Y}_i) - \frac{q_i(\mathcal{Y}_i)}{q_i(\mathcal{Y}_i)} + \lambda$$

Expectation w.r.t all other distribution

Constraint of being valid distribution. ↗ closed form solution, unique optimum,  
also global optimum

Function of  $\mathcal{Y}_i$  itself

In general, this depends only on the expected sufficient statistics under  $q_i$ . Thus, again, we don't actually need the entire distributions, just the relevant expectations (now for approximate inference as well as learning).

$$\mathbf{H}[q_i] = - \int q_i(\mathcal{Y}_i) \log q_i(\mathcal{Y}_i) d\mathcal{Y}_i$$

$$\mathbf{H}[q_i] + \sum_{j \neq i} \mathbf{H}[q_j]$$

All other factors apart from  $\mathcal{Y}_i$  are integrated out  
only thing left is  $\mathcal{Y}_i$

\* Functional Example : Entropy

$$H[p] = - \int p(x) \ln p(x) dx \quad H: \mathcal{P} \rightarrow \mathbb{R}^+$$

\* Derivation of  $\mathcal{F}$

$$\begin{aligned} \mathcal{F}(q(y), \theta^{(k-1)}) &= \left\langle \log P(X, Y | \theta^{(k-1)}) \right\rangle_{\prod_j q_j(y_j)} + H[q_j(y)] \\ \Leftrightarrow \mathcal{F}\left(\prod_j q_j(y_j), \theta^{(k-1)}\right) &= \left\langle \log P(X, Y | \theta^{(k-1)}) \right\rangle_{\prod_j q_j(y_j)} + H\left[\prod_j q_j(y_j)\right] \\ &= \int q_i(y_i) \cdot \left\langle \log P(X, Y | \theta^{(k-1)}) \right\rangle_{\prod_{j \neq i} q_j(y_j)} dy_i \\ &\quad + H[q_i] + \underline{\sum_{j \neq i} H[q_j]} \end{aligned}$$

Notice Lagrangian constraint :  $\int q_i(y_i) dy_i = 1$ .

Also functional derivative :  $\frac{\partial}{\partial q_i} \int q_i(y_i) dy_i = \mathcal{F}$ . And  $\frac{\partial}{\partial q_i} \int q_i(y_i) dy_i = 1$

$$\begin{aligned} &\frac{\partial}{\partial q_i} \left( \mathcal{F} + \lambda \left( \int q_i(y_i) dy_i - 1 \right) \right) \quad \text{not related} \\ &= \frac{\partial}{\partial q_i} \left( \int q_i(y_i) \left\langle \log P(X, Y | \theta^{(k-1)}) \right\rangle_{\prod_{j \neq i} q_j(y_j)} dy_i - \int q_i(y_i) \log q_i(y_i) dy_i - \sum_{j \neq i} H[q_j] \right. \\ &\quad \left. + \lambda \left( \int q_i(y_i) dy_i - 1 \right) \right) = 0 \\ &= \left\langle \log P(X, Y | \theta^{(k-1)}) \right\rangle_{\prod_{j \neq i} q_j(y_j)} - \log q_i(y_i) + \lambda = 0. \end{aligned}$$

$$\Rightarrow q_i(y_i) \propto \exp \left\{ \left\langle \log P(X, Y | \theta^{(k-1)}) \right\rangle_{\prod_{j \neq i} q_j(y_j)} \right\}$$

Functional Derivative:

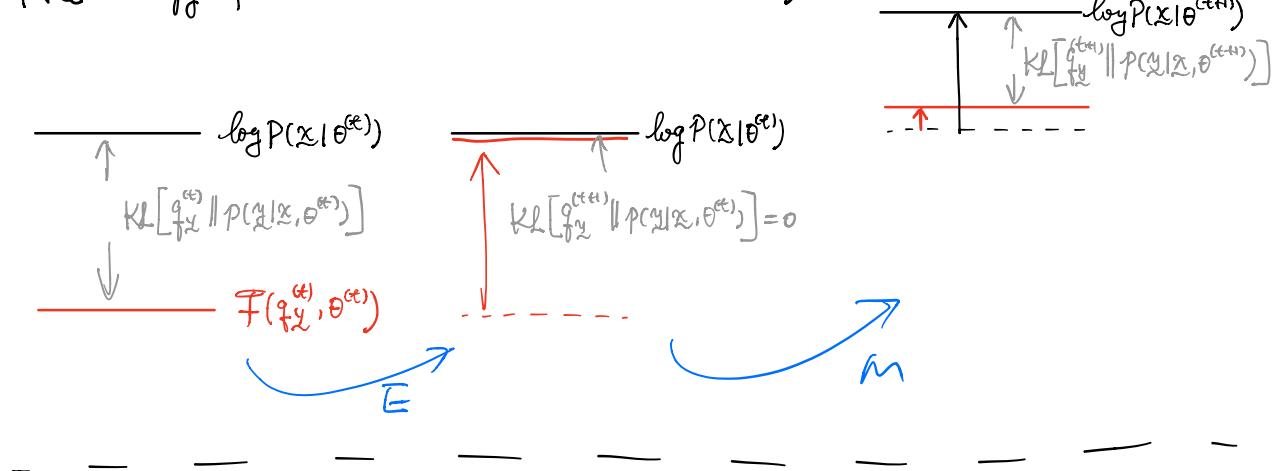
$\mathcal{F}: M \rightarrow \mathbb{R} / \mathbb{C}$  as a functional.

the functional derivative of  $\mathcal{F}(f)$ , denoted by  $\frac{\delta \mathcal{F}}{\delta f}$  is defined by:

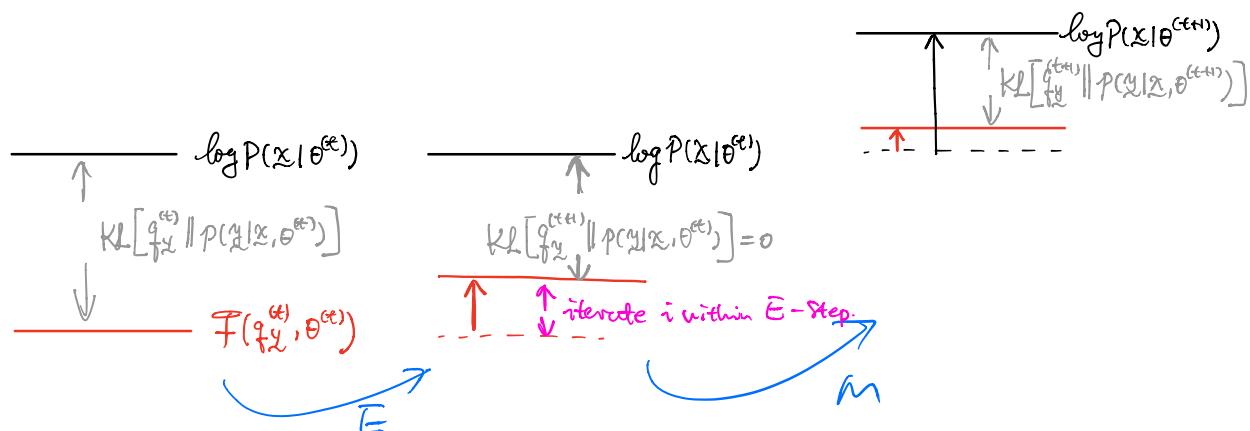
$$\frac{\delta \mathcal{F}}{\delta f} = \int \frac{\delta \mathcal{F}}{\delta f}(x) \cdot \phi(x) dx = \lim_{\epsilon \rightarrow 0} \frac{\mathcal{F}(f + \epsilon \phi) - \mathcal{F}(f)}{\epsilon} = \left[ \frac{d}{d\epsilon} \mathcal{F}(f + \epsilon \phi) \right]_{\epsilon=0}$$

where  $\phi$  is an arbitrary function, and  $\epsilon \phi$  is called the variation of  $f$ .

Free Energy for Normal EM (Unconstrained)



Free Energy for constrained EM  $Q_{\mathbf{y}}(\mathbf{y}) = \prod_i Q_{y_i}(y_i)$



**Mean-field approximations** Factored to an extreme level that each set has just one latent variable

If  $\mathcal{Y}_i = y_i$  (i.e.,  $q$  is factored over all variables) then the variational technique is often called a "mean field" approximation.

## In Boltzmann Machine:

- Suppose  $P(\mathcal{X}, \mathcal{Y})$  has sufficient statistics that are **separable** in the latent variables:  
e.g. the Boltzmann machine

Sufficient statistics

$$P(\mathcal{X}, \mathcal{Y}) = \frac{1}{Z} \exp \left( \sum_{ij} W_{ij} s_i s_j + \sum_i b_i s_i \right)$$

(Iteration of the E-step)

with some  $s_i \in \mathcal{Y}$  and others observed.

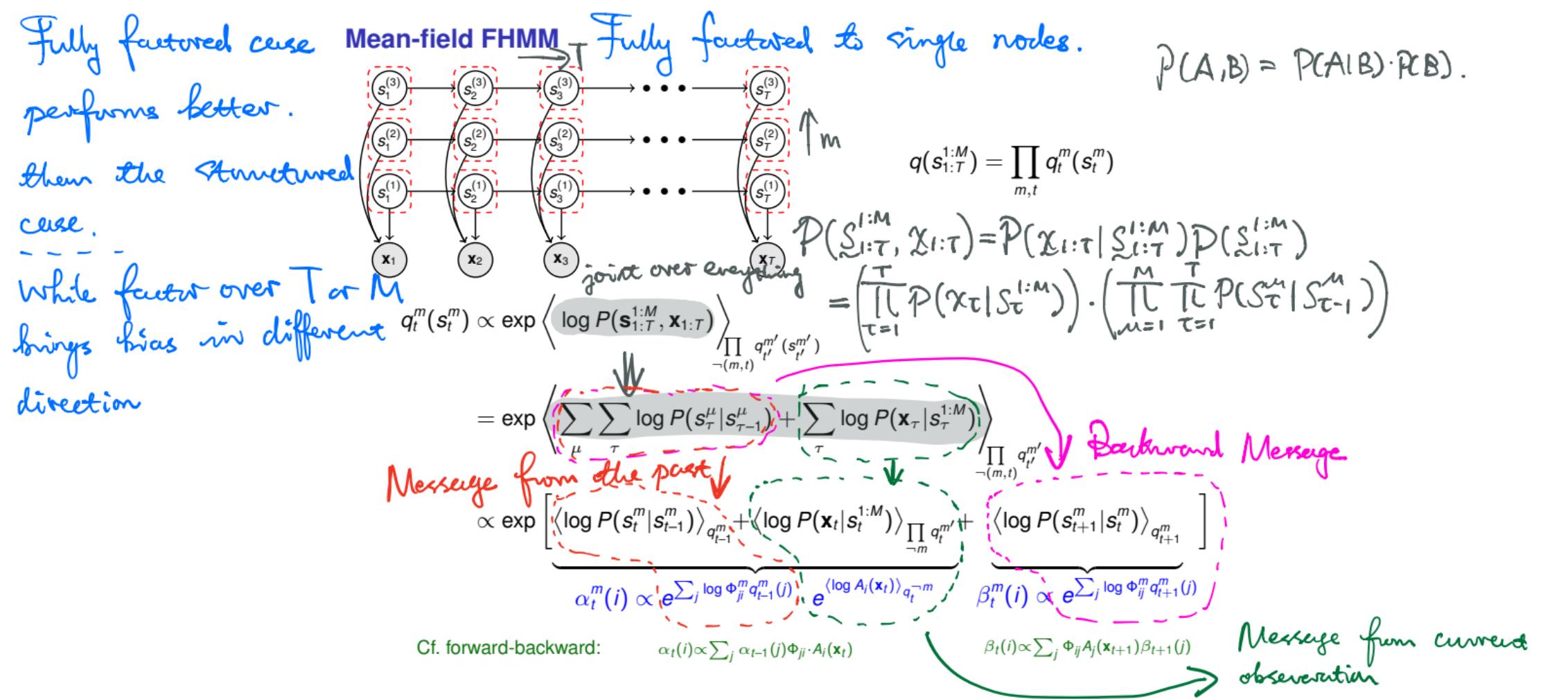
- Expectations wrt a fully-factored  $q$  distribute over all  $s_i \in \mathcal{Y}$

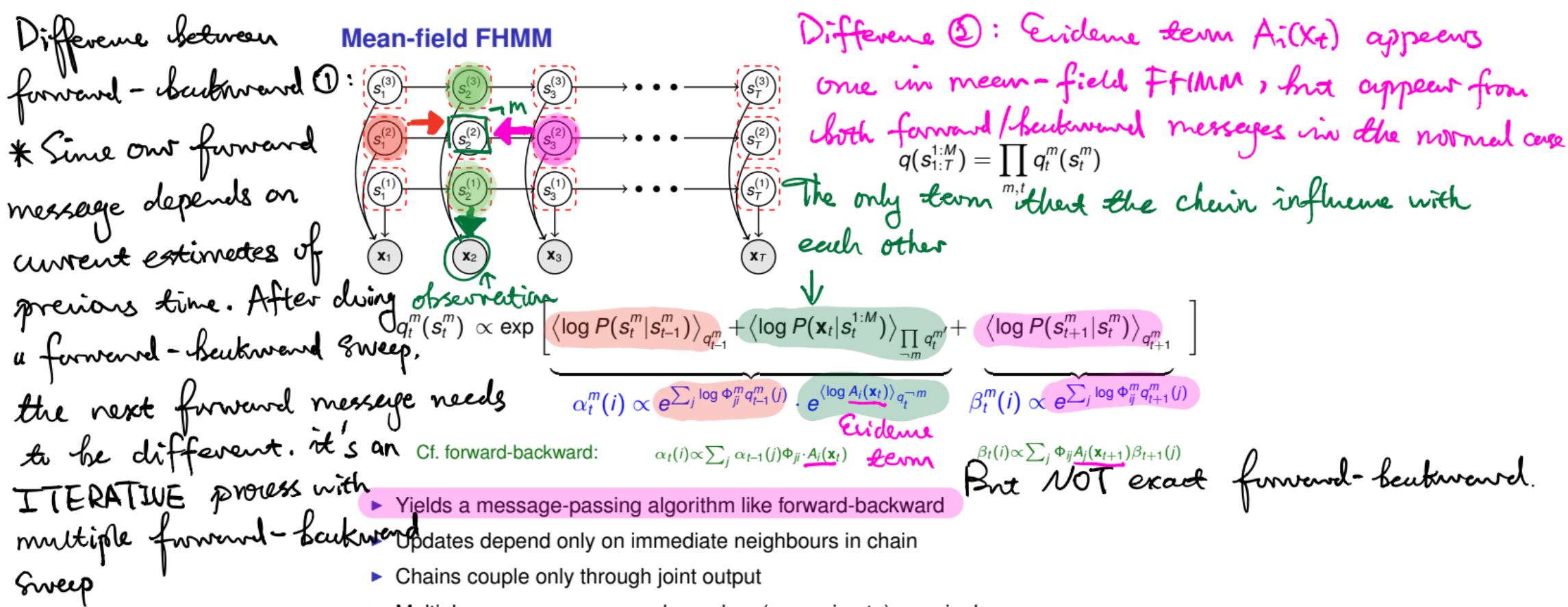
Iteration through this process eventually will find a self-consistent point.

$$\langle \log P(\mathcal{X}, \mathcal{Y}) \rangle_{\prod q_i} = \sum_{ij} W_{ij} \langle s_i \rangle_{q_i} \langle s_j \rangle_{q_j} + \sum_i b_i \langle s_i \rangle_{q_i}$$

(where  $q_i$  for  $s_i \in \mathcal{X}$  is a delta function on the observed value).

- Thus, we can update each  $q_i$  in turn given the **means** (or, in general, mean sufficient statistics) of the others.
- Each variable sees the **mean field** imposed by its neighbours, and we update these fields until they all agree.

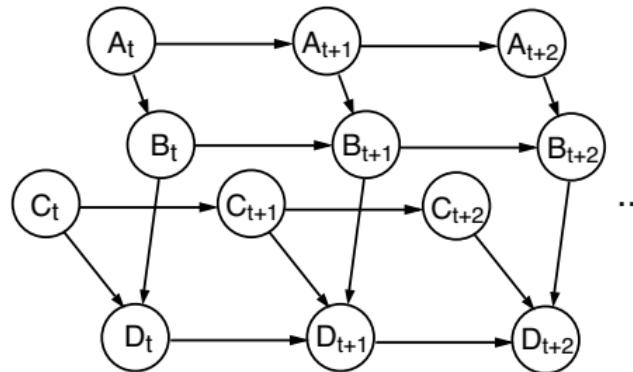


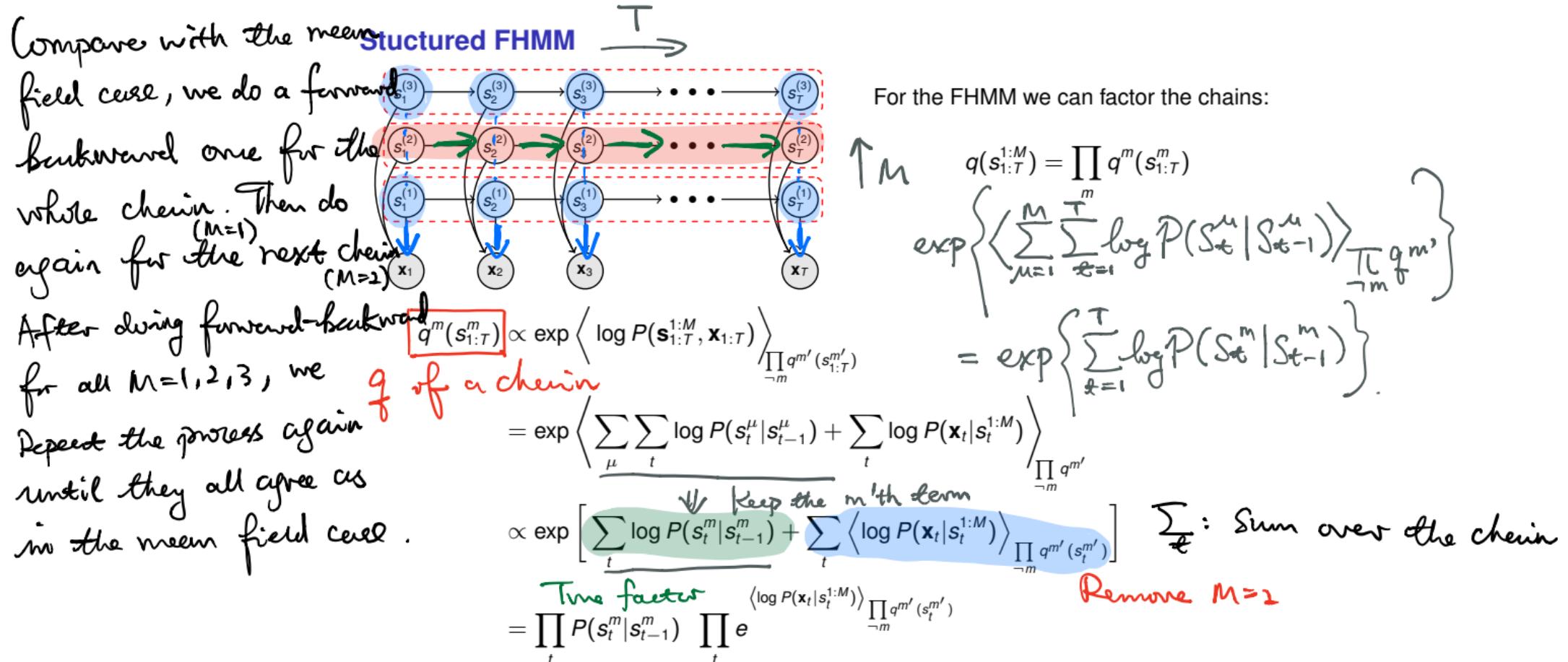


## Structured variational approximation

Theorem :

- ▶  $q(\mathcal{Y})$  need not be completely factorized.
- ▶ For example, suppose  $\mathcal{Y}$  can be partitioned into sets  $\mathcal{Y}_1$  and  $\mathcal{Y}_2$  such that computing the expected sufficient statistics under  $P(\mathcal{Y}_1|\mathcal{Y}_2, \mathcal{X})$  and  $P(\mathcal{Y}_2|\mathcal{Y}_1, \mathcal{X})$  would be tractable.
- ⇒ Then the factored approximation  $q(\mathcal{Y}) = q(\mathcal{Y}_1)q(\mathcal{Y}_2)$  is tractable.
- ▶ In particular, any factorisation of  $q(\mathcal{Y})$  into a product of distributions on **trees**, yields a tractable approximation.

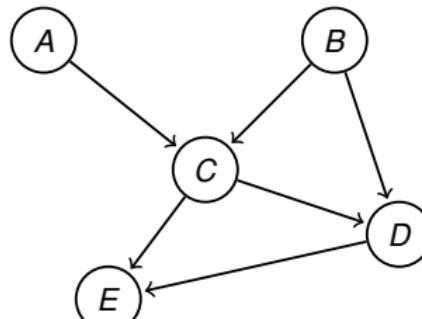




This looks like a standard HMM joint, with a modified likelihood term  $\Rightarrow$  cycle through multiple forward-backward passes, updating likelihood terms each time.

## Messages on an arbitrary graph

Consider a DAG:



$$P(\mathcal{X}, \mathcal{Y}) = \prod_k P(Z_k | \text{pa}(Z_k))$$

Further, if the conditional and let  $q(\mathcal{Y}) = \prod_i q_i(\mathcal{Y}_i)$  for disjoint sets  $\{\mathcal{Y}_i\}$ .

We have that the VE update for  $q_i$  is given by  $q_i^*(\mathcal{Y}_i) \propto \exp \langle \log p(\mathcal{Y}, \mathcal{X}) \rangle_{q_{-i}(\mathcal{Y})}$  where distribution have conjugate  $\langle \cdot \rangle_{q_{-i}(\mathcal{Y})}$  denotes averaging with respect to  $q_j(\mathcal{Y}_j)$  for all  $j \neq i$ . Then: the variational passing scheme

$$\begin{aligned} \log q_i^*(\mathcal{Y}_i) &= \left\langle \sum_k \log P(Z_k | \text{pa}(Z_k)) \right\rangle_{q_{-i}(\mathcal{Y})} + \text{const} \\ &= \sum_{j \in \mathcal{Y}_i} \langle \log P(Y_j | \underline{\text{pa}}(Y_j)) \rangle_{q_{-i}(\mathcal{Y})} + \sum_{j \in \text{ch}(\mathcal{Y}_i)} \langle \log P(Z_j | \underline{\text{pa}}(Z_j)) \rangle_{q_{-i}(\mathcal{Y})} + \text{const} \end{aligned}$$

Parents                                      Children

*co-parent of node j: the other parents of  
of the child nodes apart from j itself  
↳ co-parents.*

Theorem:

This defines messages that are passed between nodes in the graph. Each node receives messages from its Markov boundary: parents, children and parents of children (all neighbours in the corresponding factor graph).

The update of factors in the variational posterior distribution represents a local calculation on the graph.

## Non-factored variational methods

The term **variational approximation** is used whenever a bound on the likelihood (or on another estimation cost function) is optimised, but does not necessarily become tight.

Many further variational approximations have been developed, including:

- ▶ parametric forms (e.g. Gaussian) for non-linear models
- ▶ non-free-energy-based bounds (both upper and lower) on the likelihood.

*Rather than learning datasets / iid observations, it learns a separate machine takes data  
We can also see MAP- or zero-temperature EM and recognition models as parametric forms produces parameters of  
of variational inference. constrained class are just delta functions variational distribution  
as output*

*Guaranteed to converge as it always pushes up free energy*



Variational methods can also be used to find an approximate posterior on the parameters.

Previous: use factorisation to find distributions over latent variables (Approximate posterior on latent variables)  
Then: use it to find distribution over own parameters. (Approximate posterior on parameters)  
*Enlarge the graph, then factorisation makes it easy ... ??*

Previous : Variational Bayes Approximate Posterior on PARAMETERS

Integrate Free energy (log-joint) over latent.

So far, we have applied Jensen's bound and factorisations to help with integrals over latent variables.

Now : Integrate evidence / log-likelihood over latent & parameters.

We can do the same for integrals over parameters in order to bound the log marginal Bound Evidence

Evidence over a model class

$$\log P(\mathcal{X}|\mathcal{M}) = \log \iint d\mathcal{Y} d\theta P(\mathcal{X}, \mathcal{Y}, \theta | \mathcal{M}) P(\theta | \mathcal{M})$$

$$= \max_Q \iint d\mathcal{Y} d\theta Q(\mathcal{Y}, \theta) \log \frac{P(\mathcal{X}, \mathcal{Y}, \theta | \mathcal{M})}{Q(\mathcal{Y}, \theta)}$$

$$\geq \max_{Q_Y, Q_\theta} \iint d\mathcal{Y} d\theta Q_Y(\mathcal{Y}) Q_\theta(\theta) \log \frac{P(\mathcal{X}, \mathcal{Y}, \theta | \mathcal{M})}{Q_Y(\mathcal{Y}) Q_\theta(\theta)}$$

$\leftarrow$  Integrate out both the hidden variables and parameters

★ The constraint that the distribution  $Q$  must factor into the product  $Q_Y(\mathcal{Y})Q_\theta(\theta)$  leads to the variational Bayesian EM algorithm or just "Variational Bayes".

$$Q(\mathcal{Y}, \theta) = Q_Y(\mathcal{Y})Q_\theta(\theta)$$

Some call this the "Evidence Lower Bound" (ELBO). I'm not fond of that term

果然知道了。

$$\begin{aligned}
 \log P(x|M) &= \log \iint P(x, y, \theta | M) dy d\theta \\
 &= \log \iint Q(y, \theta) \cdot \frac{P(x, y, \theta | M)}{Q(y, \theta)} dy d\theta \quad \xrightarrow{\text{Jensen's Inequality}} \\
 &\geq \iint Q(y, \theta) \cdot \log \left( \frac{P(x, y, \theta | M)}{Q(y, \theta)} \right) dy d\theta \\
 &\geq \iint Q_y(y) \cdot Q_\theta(\theta) \cdot \log \left( \frac{P(x, y, \theta | M)}{Q_y(y) \cdot Q_\theta(\theta)} \right) dy d\theta \\
 &\stackrel{\text{def}}{=} \underline{F(Q_y(y), Q_\theta(\theta))} \quad \xrightarrow{\text{Factorise } Q(y, \theta)} \\
 &= Q_y(y) \cdot Q_\theta(\theta). \\
 &\quad "My distribution over latent variables are independent of parameters, hence gives a lower bound. Seldom happens... So equality almost never holds"
 \end{aligned}$$

Model-class Marginal likelihood / Evidence bounded by VB-Free Energy.

## Variational Bayesian EM ...

Coordinate maximization of the VB free-energy lower bound

$$\mathcal{F}(Q_Y, Q_\theta) = \iint dY d\theta Q_Y(Y) Q_\theta(\theta) \log \frac{P(X, Y, \theta | M)}{Q_Y(Y) Q_\theta(\theta)}$$

leads to **EM-like** updates: (Instead of updating single value, we update whole distribution)

$$\left. \begin{array}{l} Q_Y^*(Y) \propto \exp \langle \log P(Y, X | \theta) \rangle_{Q_\theta(\theta)} \\ Q_\theta^*(\theta) \propto \underline{P(\theta)} \exp \langle \log P(Y, X | \theta) \rangle_{Q_Y(Y)} \end{array} \right\} \begin{array}{l} E\text{-like step} \\ M\text{-like step} \end{array}$$

Prior on  $\theta$

Instead of optimizing single  $\theta$ , we optimise a distribution

Maximizing  $\mathcal{F}$  is equivalent to minimizing KL-divergence between the approximate posterior,  $Q(\theta)Q(Y)$  and the true posterior,  $P(\theta, Y | X)$ .

$$\log P(X) - \mathcal{F}(Q_Y, Q_\theta) = \log P(X) - \iint dY d\theta Q_Y(Y) Q_\theta(\theta) \log \frac{P(X, Y, \theta)}{Q_Y(Y) Q_\theta(\theta)}$$

$$= \iint Q_Y Q_\theta \log \left( \frac{Q_Y Q_\theta P(X)}{P(X, Y, \theta)} \right) dY d\theta = \iint dY d\theta Q_Y(Y) Q_\theta(\theta) \log \frac{Q_Y(Y) Q_\theta(\theta)}{P(Y, \theta | X)} = \text{KL}[Q_Y Q_\theta || P(Y, \theta | X)]$$

$$\log P(X) = \text{KL}(Q || P) + \mathcal{F}(Q_Y, Q_\theta)$$

Detailed learning rule

Bear 2003. Theorem 2.1 : Learning rule of Variational Bayes EM:

Let  $M$  be model with parameters  $\Theta$  giving rise to iid observed data  $X = \{x_1, \dots, x_n\}$  with latent variables  $Y = \{y_1, y_2, \dots, y_n\}$ .

We have the log marginal lower bound / Evidence lower bound / VB free energy:

$$\mathcal{F}(Q_Y(Y), Q_\Theta(\Theta)) = \iint Q_Y(Y) Q_\Theta(\Theta) \cdot \log \left( \frac{P(X, Y | \Theta, M)}{Q_Y(Y) Q_\Theta(\Theta)} \right) dY d\Theta$$

Which can be optimised by the following step, with  $t$  be number of iteration:

$$\text{VBE - Step: } q_{Y_i}^{(t+1)}(y_i) = \frac{1}{Z_{Y_i}} \exp \left( \int q_\Theta^{(t)}(\theta) \cdot \log P(x_i, y_i | \theta, M) d\theta \right) \quad \forall i \in \text{Iteration in VBE - Step.}$$

$$\text{where } q_Y^{(t+1)}(y) = \prod_{i=1}^n q_{Y_i}^{(t+1)}(y_i).$$

$$\text{VBM - Step: } q_\Theta^{(t+1)}(\theta) = \frac{1}{Z(\theta)} \underbrace{\left[ P(\theta | M) \cdot \exp \left( \int q_Y(Y) \log P(X, Y | \theta, M) dY \right) \right]}_{\substack{\uparrow \\ \text{Prior on } \theta}} \leftarrow \text{one step in VBM}$$

Moreover, the learning rule converges to a LOCAL maximum of  $\mathcal{F}(Q_Y(Y), Q_\Theta(\Theta))$

Proof: \*VBE:

$$\begin{aligned} \frac{\delta}{\delta Q_Y(Y)} \mathcal{F}(Q_Y(Y), Q_\Theta(\Theta)) &= \int Q_\Theta(\Theta) d\Theta \cdot \left[ \frac{\delta}{\delta Q_Y(Y)} \int Q_Y(Y) \cdot \log \left( \frac{P(X, Y | \Theta, M)}{Q_Y(Y) Q_\Theta(\Theta)} \right) dY \right] \\ &= \int Q_\Theta(\Theta) d\Theta \cdot \left[ \frac{\delta}{\delta Q_Y(Y)} \int Q_Y(Y) \cdot \log \left( \frac{P(X, Y | \Theta, M)}{Q_Y(Y)} \right) dY \right] \\ &= \int Q_\Theta(\Theta) d\Theta \cdot \left[ \log \left( \frac{P(X, Y | \Theta, M)}{Q_Y(Y)} \right) + Q_Y(Y) \cdot \left( -\frac{1}{Q_Y(Y)} \right) \right] \\ &= \int Q_\Theta(\Theta) d\Theta \cdot \left[ \log P(X, Y | \Theta, M) - \log Q_Y(Y) - 1 \right] = 0. \end{aligned}$$

$$\Leftrightarrow \log Q_Y^{(t+1)}(Y) = \int \log P(X, Y | \Theta, M) \cdot Q_\Theta^{(t)}(\Theta) d\Theta + \text{const.}$$

$t$ : outer loop

into factored form:

$$\Leftrightarrow \log Q_{Y_i}^{(t+1)}(y_i) = \int \log P(x_i, y_i | \theta, M) \cdot Q_\Theta^{(t)}(\theta) d\theta + \text{const.} \quad \forall i. \quad i: \text{iteration within VBE-step}$$

\* VBM:

$$\frac{\delta}{\delta Q_\theta(\theta)} \mathcal{F}(Q_x(x), Q_\theta(\theta)) = \frac{\delta}{\delta Q_\theta(\theta)} \int d\theta Q_\theta(\theta) \left[ \int Q_y(y) \log \left( \frac{P(x, y | \theta, M)}{Q_\theta(\theta) Q_x(x)} \right) dy \right]$$

$$= \frac{\delta}{\delta Q_\theta(\theta)} \int d\theta Q_\theta(\theta) \left[ \int Q_y(y) \log \left( \frac{P(x, y | \theta) \cdot P(\theta | M)}{Q_\theta(\theta) \cdot Q_x(x)} \right) dy \right] \quad \text{移掉 } x$$

$$= \frac{\delta}{\delta Q_\theta(\theta)} \int d\theta Q_\theta(\theta) \left[ \int Q_y(y) \log P(x, y | \theta) dy + \int Q_y(y) \log \left( \frac{P(\theta | M)}{Q_\theta(\theta) Q_x(x)} \right) dy \right]$$

$$= \frac{\delta}{\delta Q_\theta(\theta)} \int d\theta Q_\theta(\theta) \left[ \int Q_y(y) \log P(x, y | \theta) dy + \log \left( \frac{P(\theta | M)}{Q_\theta(\theta)} \right) \right]$$

$$= \int Q_y(y) \log P(x, y | \theta) dy + \log P(\theta | M) - \log Q_\theta(\theta) + \text{const.} = 0$$

$$\Leftrightarrow \log Q_\theta^{(t+1)}(\theta) = \underbrace{\log P(\theta | M)}_{\text{Prior of } \theta} + \underbrace{\int Q_y^{(t+1)}(y) \log P(x, y | \theta) dy}_{\text{order loop.}} + \text{const.}$$

VBM-step is just optimise over whole  $Q_\theta(\theta)$  distribution.

So just do one step.

Therefore, the VBEM algorithm is:

while VB-free energy not converge

for  $i = 1:n$

$$\log Q_{y_i}^{(t+1)}(y_i) = \int \log P(x_i, y_i | \theta, M) \cdot Q_\theta^{(t)}(\theta) d\theta + \text{const.}$$

end

$$\log Q_\theta^{(t+1)}(\theta) = \log P(\theta | M) + \int Q_y^{(t+1)}(y) \log P(x, y | \theta) dy + \text{const.}$$

end

$\overline{\log P(x | m)}$ .

$\overline{\log P(x | m)}$ .

$\overline{\log P(x | m)}$ .

$\uparrow$   
 $KL[Q_y^{(t+1)} || P(y, \theta | x)]$

$\downarrow$   
 $KL[Q_y^{(t+1)} || P(y, \theta | x)]$

$\uparrow$   
 $KL[Q_y^{(t+1)} \cdot Q_\theta^{(t+1)} || P(y, \theta | x)]$

$\downarrow$   
 $\mathcal{F}(Q_y^{(t+1)}, Q_\theta^{(t+1)})$

$\downarrow$   
 $\mathcal{F}(Q_x^{(t)}(x), Q_\theta^{(t)})$

$\uparrow$   
 $\mathcal{F}(Q_x^{(t+1)}, Q_\theta^{(t)})$

$\rightarrow$   
 $\mathcal{F}(Q_x^{(t+1)}, Q_\theta^{(t+1)})$

$\overbrace{\qquad\qquad\qquad}^{\text{VBE}}$

$\overbrace{\qquad\qquad\qquad}^{\text{iterate through } i}$

$\overbrace{\qquad\qquad\qquad}^{\text{VBM}}$

} gives fixed point iteration.

## Conjugate-Exponential models

Let's focus on *conjugate-exponential* (**CE**) latent-variable models:

- ▶ **Condition (1).** The *joint probability* over *variables* is in the *exponential family*:

$$P(\mathcal{Y}, \mathcal{X}|\theta) = f(\mathcal{Y}, \mathcal{X}) g(\theta) \exp \left\{ \phi(\theta)^T T(\mathcal{Y}, \mathcal{X}) \right\}$$

where  $\phi(\theta)$  is the vector of *natural parameters*,  $T$  are *sufficient statistics*

- ▶ **Condition (2).** The *prior* over *parameters* is *conjugate* to this joint probability:

$$P(\theta|\nu, \tau) = h(\nu, \tau) g(\theta)^\nu \exp \left\{ \phi(\theta)^T \tau \right\}$$

where  $\nu$  and  $\tau$  are hyperparameters of the prior.

Conjugate priors are computationally convenient and have an intuitive interpretation:

- ▶  $\nu$ : number of pseudo-observations
- ▶  $\tau$ : values of pseudo-observations

CE :

Exact Inference

or  
Variational Inference.

## Conjugate-Exponential examples

In the **CE** family:

- ▶ Gaussian mixtures
- ▶ factor analysis, probabilistic PCA
- ▶ hidden Markov models and factorial HMMs
- ▶ linear dynamical systems and switching models
- ▶ discrete-variable belief networks

} closed form for variational Bayes  
► LDA (can do variational LDA).

Other as yet undreamt-of models combinations of Gaussian, Gamma, Poisson, Dirichlet, Wishart, Multinomial and others.

Non-CE :

Variational Inference

Not in the **CE** family: *Cannot find normaliser*

- ▶ Boltzmann machines, MRFs (no simple conjugacy)
- ▶ logistic regression (no simple conjugacy)
- ▶ sigmoid belief networks (not exponential)
- ▶ independent components analysis (not exponential)

} Non-closed form  
Numerical Issues

Note: one can often approximate such models with a suitable choice from the **CE** family.

## Conjugate-exponential VB

Given an iid data set  $\mathcal{D} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ , if the model is **CE** then:

- $Q_\theta(\theta)$  is also conjugate, i.e.

*Conjugate Prior*

$$\begin{aligned} Q_\theta(\theta) &\propto P(\theta) \cdot \exp \left\langle \sum_i \log P(\mathbf{y}_i, \mathbf{x}_i | \theta) \right\rangle_{Q_Y} \\ &= h(\nu, \tau) g(\theta)^\nu e^{\phi(\theta)^\top \tau} \quad g(\theta)^n e^{\left\langle \log f(Y, X) \right\rangle_{Q_Y}} e^{\phi(\theta)^\top \left\langle \sum_i T(\mathbf{y}_i, \mathbf{x}_i) \right\rangle_{Q_Y}} \\ &\propto h(\tilde{\nu}, \tilde{\tau}) g(\theta)^{\tilde{\nu}} e^{\phi(\theta)^\top \tilde{\tau}} \end{aligned}$$

proportion, has no  $\theta$  dependences .

with  $\tilde{\nu} = \nu + n$  and  $\tilde{\tau} = \tau + \sum_i \langle T(\mathbf{y}_i, \mathbf{x}_i) \rangle_{Q_Y}$  ⇒ only need to track  $\tilde{\nu}, \tilde{\tau}$ .

- $Q_Y(Y) = \prod_{i=1}^n Q_{Y_i}(\mathbf{y}_i)$  takes the same form as in the E-step of regular EM

$$\begin{aligned} Q_{Y_i}(\mathbf{y}_i) &\propto \exp \langle \log P(\mathbf{y}_i, \mathbf{x}_i | \theta) \rangle_{Q_\theta} \quad \text{redefine .} \\ &\propto f(\mathbf{y}_i, \mathbf{x}_i) e^{\langle \phi(\theta) \rangle_{Q_\theta}^\top T(\mathbf{y}_i, \mathbf{x}_i)} = P(\mathbf{y}_i | \mathbf{x}_i, \bar{\phi}(\theta)) \end{aligned}$$

with natural parameters  $\bar{\phi}(\theta) = \langle \phi(\theta) \rangle_{Q_\theta}$  ⇒ inference unchanged from regular EM.

\* Complete-data likelihood:

$$P(x_i, y_i | \theta) = g(\theta) \cdot f(x_i, y_i) \cdot \exp(\phi(\theta)^T T(x_i, y_i)) \text{ where } g(\theta) = \int \int f(x_i, y_i) \exp(\phi(\theta)^T T(x_i, y_i)) dx_i dy_i$$

\* Parameter prior:

$$P(\theta | \nu, \tau) = h(\nu, \tau) g(\theta)^\nu \exp(\phi(\theta)^T \tau). \text{ where } h(\nu, \tau) = \int g(\theta)^\nu \exp(\phi(\theta)^T \tau) d\theta$$

VBE Step:

$$\begin{aligned} Q_{y_i}^{(t)}(y_i) &\propto \exp\left(\langle \log P(x_i, y_i | \theta^{(t)}) \rangle_{Q_\theta}\right) \quad \forall i \\ &\propto f(x_i, y_i) \exp(\bar{\phi}^T T(x_i, y_i)), \text{ where } \bar{\phi} = \int Q_\theta(\theta) \phi(\theta) d\theta \quad \forall i \\ &= P(y_i | x_i, \bar{\phi}^{(t)}) \quad \forall i \quad = \langle \phi(\theta) \rangle_{Q_\theta(\theta)}. \end{aligned}$$

rewrite as

VBEM Step:

$$\begin{aligned} Q_\theta(\theta) &\propto P(\theta | M) \cdot \exp\left(\langle \log P(x, y | \theta, M) \rangle_{Q_{y_i}^{(t)}(y_i)}\right) \\ &= h(\nu, \tau) \cdot g(\theta)^\nu \exp(\phi(\theta)^T \tau) \cdot \exp\left(\langle \sum_{i=1}^n \log P(x_i, y_i | \theta, M) \rangle_{Q_{y_i}^{(t)}(y_i)}\right) \\ &= h(\nu, \tau) \cdot g(\theta)^\nu \exp(\phi(\theta)^T \tau) \cdot g(\theta)^n \cdot \boxed{\exp\left(\langle \sum_{i=1}^n f(x_i, y_i) \rangle_{Q_{y_i}^{(t)}}\right)} \cdot \exp\left(\phi(\theta)^T \sum_{i=1}^n \langle T(x_i, y_i) \rangle_{Q_{y_i}^{(t)}}\right) \\ &\quad \text{no } \theta \text{ dependence.} \\ &\propto h(\tilde{\nu}, \tilde{\tau}) \cdot g(\theta)^{\nu+n} \cdot \exp\left(\phi(\theta)^T [\tau + \sum_{i=1}^n \langle T(x_i, y_i) \rangle_{Q_{y_i}^{(t)}}]\right) \\ \tilde{\nu} &= \nu + n \quad \text{and} \quad \tilde{\tau} = \tau + \sum_{i=1}^n \langle T(x_i, y_i) \rangle_{Q_{y_i}^{(t)}}. \end{aligned}$$

### Summary of VBEM for CE Models

while VB free energy not converge (loop over  $t$ ).

**VBE** [ for  $i = 1:n$   
compute sufficient statistics  $\langle T(x_i, y_i) \rangle_{Q_{y_i}^{(t)}}$   
end  
 $\nu += n$ ,  $\tau += \sum_i \langle T(x_i, y_i) \rangle_{Q_{y_i}^{(t)}}$

**VBEM** [ compute  $\bar{\phi} = \langle \phi(\theta) \rangle$  under  $\nu$  and  $\tau$   
end

## The Variational Bayesian EM algorithm

*(the  $Q_{\theta}(\mathcal{Y})$  in EM and VB-EM are exactly the same)*

EM for MAP estimation (for CE)	
Goal:	maximize $P(\theta   \mathcal{X}, m)$ wrt $\theta$
E Step:	compute $Q_{\theta}(\mathcal{Y}) \leftarrow p(\mathcal{Y}   \mathcal{X}, \theta)$
M Step:	$\theta \leftarrow \operatorname{argmax}_{\theta} \int d\mathcal{Y} Q_{\theta}(\mathcal{Y}) \log P(\mathcal{Y}, \mathcal{X}, \theta)$
	<i>one point value</i>

Variational Bayesian EM (for CE)	
Goal:	maximise bound on $P(\mathcal{X}   m)$ wrt $Q_{\theta}$
VB-E Step:	compute $Q_{\theta}(\mathcal{Y}) = p(\mathcal{Y}   \mathcal{X}, \bar{\phi})$
VB-M Step:	$Q_{\theta}(\theta) \leftarrow \exp \int d\mathcal{Y} Q_{\theta}(\mathcal{Y}) \log P(\mathcal{Y}, \mathcal{X}, \theta)$
	<i>Mean of natural parameters</i>
	<i>distribution</i>

$$\bar{\phi}(\theta) = \langle \phi(\theta) \rangle_{Q(\theta)}$$

Properties:

Pushing up  $\mathcal{F}_m$  incorporates model complexity penalty.

Dirac Delta

- ▶ Reduces to the EM algorithm if  $Q_{\theta}(\theta) = \delta(\theta - \theta^*)$ . ↴ (Because we're working on evidence not likelihood)
  - ▶  $\mathcal{F}_m$  increases monotonically, and incorporates the model complexity penalty.
  - ▶ Analytical parameter distributions (but not constrained to be Gaussian).
  - ▶ VB-E step has same complexity as corresponding E step.
  - ▶ We can use the junction tree, belief propagation, Kalman filter, etc, algorithms in the VB-E step of VB-EM, but using **expected natural parameters**,  $\bar{\phi}$ .
- ↳ Integrating over parameters incorporating a Occam's Razor effect.

## VB and model selection

- ▶ Variational Bayesian EM yields an approximate posterior  $Q_\theta$  over model parameters.
- ▶ It also yields an optimised lower bound on the model evidence

Model Selection

Free energy depends on  
 $Q_Y$ ,  $Q_\theta$  and  
Hyperparameter  $\eta$

$$\max \mathcal{F}_M(Q_Y, Q_\theta) \leq P(D|M)$$

- ▶ These lower bounds can be compared amongst models to learn the right (structure, connectivity ... of the) model

- ▶ If a continuous domain of models is specified by a hyperparameter  $\eta$ , then the VB free energy depends on that parameter:

\* Optimise w.r.t.  $Q_Y$  and  $Q_\theta$   
to obtain a lower bound

$$\mathcal{F}(Q_Y, Q_\theta, \eta) = \iint dY d\theta Q_Y(Y) Q_\theta(\theta) \log \frac{P(X, Y, \theta | \eta)}{Q_Y(Y) Q_\theta(\theta)} \leq P(X | \eta)$$

Just like mean field,  
as long as we visit (update)

$$\eta \leftarrow \operatorname{argmax}_{\eta} \iint dY d\theta Q_Y(Y) Q_\theta(\theta) \log P(X, Y, \theta | \eta)$$

each  $q(y_i)$  and  $q(\theta)$ ,  
pushes up lower bound

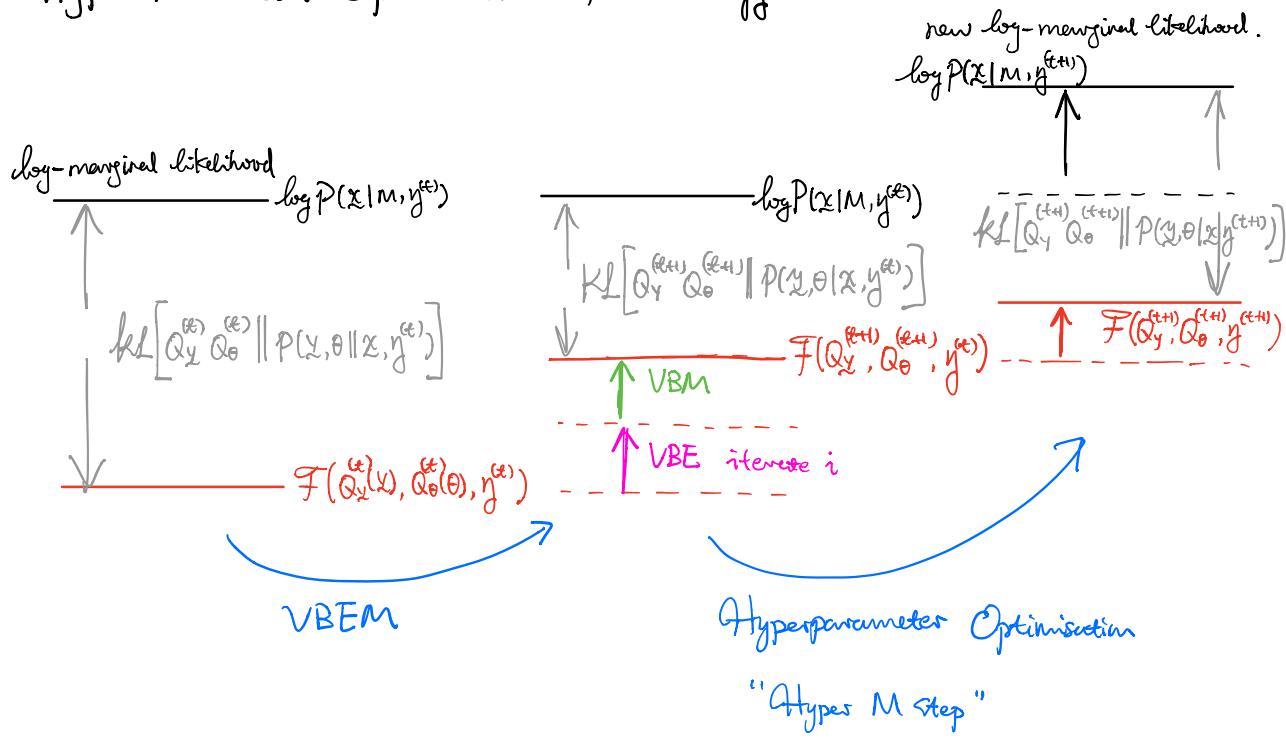
Update Hyper parameters to increase Free Energy.  
which is approximation to Model Evidence.

A hyper-M step maximises the current bound wrt  $\eta$ :

↳ Analogous to standard M step.

Hyper  
M-step.

## Hyper Parameter Optimisation Free Energy



Model Selection :

Compare the log ratio of posterior probabilities for  $M$  and  $M'$

$$\begin{aligned} \log \frac{P(M|X)}{P(M'|X)} &= \log P(M) + \log P(X|M) - \log P(M') - \log P(X|M') \\ &= \log P(M) + \mathcal{F}(Q_y, \theta) + KL[Q(y, \theta) \| P(y, \theta|x, M)] \\ &\quad - \log P(M') - \mathcal{F}'(Q_y, \theta) - KL[Q'(y, \theta) \| P(y, \theta|x, M')] \end{aligned}$$

Motivation for ARD :

(variance w.r.t. itself.  
ignore correlation  
between dimensions)

For a regression problem with too many inputs, some of which are irrelevant to the prediction. So we need a **PRIOR** ( $\Lambda$ ) over the regression parameters embodies the concept of relevance. In this case,  $\Lambda : i \sim N(\mathbf{0}, \alpha_i^{-1} \mathbf{I})$ . We hence  $\alpha_i$  associated to each input, controls weights from input to output.

Under the variational Bayes / Evidence framework, we do a hyperparameter M step, to learn  $\alpha$

Maximise Evidence w.r.t.  $\alpha$

## ARD for unsupervised learning      Automatic Way to Learn dimension

Recall that ARD (automatic relevance determination) was a hyperparameter method to select relevant or useful inputs in regression.

- ▶ A similar idea used with variational Bayesian methods can learn a **latent dimensionality**.

- ▶ Consider factor analysis: *See Ass5, Q1, (f) Select Latent Dimension*

$$\mathbf{x} \sim \mathcal{N}(\Lambda\mathbf{y}, \Psi) \quad \mathbf{y} \sim \mathcal{N}(0, I) \quad \text{with a column-wise prior } \Lambda_{:i} \sim \mathcal{N}(0, \alpha_i^{-1}I)$$

*linear*

- ▶ The VB free energy is

$$\mathcal{F}(Q_{\mathcal{Y}}(\mathcal{Y}), Q_{\Lambda}(\Lambda), \Psi, \alpha) = \langle \log P(\mathcal{X}, \mathcal{Y} | \Lambda, \Psi) + \log P(\Lambda | \alpha) + \log P(\Psi) \rangle_{Q_{\mathcal{Y}} Q_{\Lambda}} + \dots$$

and so hyperparameter optimisation requires  
*Prior depending on  
hyperparameter  $\alpha$*

$$\alpha \leftarrow \operatorname{argmax} \langle \log P(\Lambda | \alpha) \rangle_{Q_{\Lambda}}$$

- ▶ Now  $Q_{\Lambda}$  is Gaussian, with the same form as in linear regression, but with **expected moments of  $\mathbf{y}$**  appearing in place of the inputs.
- ▶ Optimisation wrt the distributions,  $\Psi$  and  $\alpha$  in turn causes some  $\alpha_i$  to diverge as in regression ARD.
- ▶ In this case, these parameters select “relevant” **latent dimensions**, effectively learning the dimensionality of  $\mathbf{y}$ .

*Optimise jointly with  $\Psi$  and  $\Lambda$*

\* Working in a model **Augmented Variational Methods**  
with observed Gaussian,  $p(x|s) = N(\dots)$   $\leftarrow$  Computational Intractability.  
Then computational complexity of inference  
scales cubically with the number of data we have (we have to inverse the matrix)

\* Analytical Intractability  
with non-conjugate models

Sometimes it may be useful to introduce additional latent variables, solely to achieve computational tractability.



Make the problem bigger  
then decompose

EP will focus on analytic intractability

Two examples are GP regression and the GPLVM.

## Sparse GP approximations

Recall.

GP predictions:

$$y' | X, Y, \mathbf{x}' \sim \mathcal{N} \left( K_{\mathbf{x}'X} (K_{XX} + \sigma^2 I)^{-1} Y, K_{\mathbf{x}'\mathbf{x}'} - K_{\mathbf{x}'X} (K_{XX}^{-1}) K_{X\mathbf{x}'} + \sigma^2 \right)$$

$\mathcal{O}(n^3)$ , where problem comes from (Computational Intractability)

Evidence (for learning kernel hyperparameters):

$$\log P(Y|X) = -\frac{1}{2} \log |2\pi(K_{XX} + \sigma^2 I)| - \frac{1}{2} Y(K_{XX} + \sigma^2 I)^{-1} Y^T$$

Computing either form requires inverting the  $N \times N$  matrix  $K_{XX}$ , in  $\mathcal{O}(N^3)$  time.

One proposal to make this more efficient is to find (or select) a smaller set of possibly fictitious measurements  $U$  at inputs  $Z$  such that

$$P(y'|Z, U, \mathbf{x}') \approx P(y'|X, Y, \mathbf{x}').$$

What values should  $U$  and  $Z$  take?

Called Sparse Approximation to GP (A reduced matrix for inversion)

$Z$ : Sendo Input  
 $U$ : Sendo Output

$X$ : Actual Input  
 $Z$ : Sudo Input  
 $Y$ : Actual Output  
 $U$ : Sudo Output.  
 $F$ : Value of  $Y$

**Variational Sparse GP approximations**

$F$  is still a LATENT variable. We observe  $Y$ , but still don't know the exact value of  $F$ .

Write  $F$  for the (smooth) GP function values that underlie  $Y$  (so  $Y \sim \mathcal{N}(F, \sigma^2 I)$ ). Introduce latent measurements  $U$  at inputs  $Z$  (and integrate over  $U$ ).

The likelihood can be written

Sudo - Observation Inputs  
 $\downarrow \downarrow$

$$P(Y|X) = \iint dF dU P(Y, F, U|X, Z) = \iint dF dU P(Y|F)P(F|U, X, Z)P(U|Z)$$

Now, both  $U$  and  $F$  are latent, so we introduce a variational distribution  $q(F, U)$  to form a free-energy.

Any set of parameters for kernel covariance.

$$\mathcal{F}(q(F, U), \theta) = \left\langle \log \frac{P(Y|F)P(F|U, X, Z)P(U|Z)}{q(F, U)} \right\rangle_{q(F, U)}$$

Now, choose the variational form  $q(F, U) = P(F|U, X, Z)q(U)$ . That is, fix  $F|U$  without reference to  $Y$  – so information about  $Y$  will need to be “compressed” into  $q(U)$ .

Then

$$\begin{aligned} \mathcal{F}(q(F, U), \theta, Z) &= \left\langle \log \frac{P(Y|F) P(F|U, X, Z) P(U|Z)}{P(F|U, X, Z) q(U)} \right\rangle_{P(F|U)q(U)} \\ &= \left\langle \langle \log P(Y|F) \rangle_{P(F|U)} + \log P(U|Z) - \log q(U) \right\rangle_{q(U)} \end{aligned}$$

① We don't need to consider joint. We have just Gaussian left out inducing points, not over number of data. So we can have more data, break data into small sets out inducing points.

## Variational Sparse GP approximations

$$\mathcal{F}(q(U), \theta, Z) = \left\langle \langle \log P(Y|F) \rangle_{P(F|U)} + \log P(U|Z) - \log q(U) \right\rangle_{q(U)}$$

Now  $P(F|U)$  is fixed by the generative model (rather than being subject to free optimisation).

So we can evaluate that expectation:

$$\begin{aligned} & \langle \log P(Y|F) \rangle_{P(F|U)} \quad \text{↳ Gaussian on } Y. \\ &= \left\langle -\frac{1}{2} \log |2\pi\sigma^2 I| - \frac{1}{2\sigma^2} \text{Tr}[(Y - F)(Y - F)^\top] \right\rangle_{P(F|U)} \\ &= -\frac{1}{2} \log |2\pi\sigma^2 I| - \frac{1}{2\sigma^2} \text{Tr}[(Y - \langle F \rangle_{P(F|U)})(Y - \langle F \rangle_{P(F|U)})^\top] - \frac{1}{2\sigma^2} \text{Tr}[\Sigma_{F|U}] \\ &= \log \mathcal{N}(Y|K_{xz}K_{zz}^{-1}U, \sigma^2 I) - \frac{1}{2\sigma^2} \text{Tr}[K_{xx} - K_{xz}K_{zz}^{-1}K_{zx}] \end{aligned}$$

So,

$$\begin{aligned} \mathcal{F}(q(U), \theta, Z) &= \left\langle \log \mathcal{N}(Y|K_{xz}K_{zz}^{-1}U, \sigma^2 I) + \log P(U|Z) - \log q(U) \right\rangle_{q(U)} \\ &\quad - \frac{1}{2\sigma^2} \text{Tr}[K_{xx} - K_{xz}K_{zz}^{-1}K_{zx}] . \end{aligned}$$

\* To maximise the Variational Sparse GP approximations free energy, we just set  $q(U)$  to the posterior of the PPCA-like model.

A PPCA free energy



$$\mathcal{F}(q(U), \theta, Z) = \left\langle \log \frac{\mathcal{N}(Y|K_{xz}K_{zz}^{-1}U, \sigma^2 I) P(U|Z)}{q(U)} \right\rangle_{q(U)} - \frac{1}{2\sigma^2} \text{Tr}[K_{xx} - K_{xz}K_{zz}^{-1}K_{zx}] .$$

The expectation is the free energy of a PPCA-like model with normal prior  $U \sim \mathcal{N}(0, K_{zz})$  and loading matrix  $K_{xz}K_{zz}^{-1}$ . The maximum of this free energy is the log-likelihood (achieved with  $q$  equal to the posterior under the PPCA-like model).

This gives

Gaussian  $Y$  with zero mean and  $K_{xz}K_{zz}^{-1}K_{zx} + \sigma^2 I$  covariance

$$\mathcal{F}(q^*(U), \theta, Z) = \log \mathcal{N}(Y|0, K_{xz}K_{zz}^{-1}K_{zx}K_{zz}^{-1}K_{zx} + \sigma^2 I) - \frac{1}{2\sigma^2} \text{Tr}[K_{xx} - K_{xz}K_{zz}^{-1}K_{zx}] .$$

Note that we have eliminated all terms in  $K_{xx}^{-1}$ .

The only inverse, comes from variational.

We can optimise the free energy numerically with respect to  $Z$  and  $\theta$  to adjust the GP prior and quality of variational approximation.

A similar approach can be used to learn  $X$  if they are unobserved (i.e. in the GPLVM).

Assume  $q(X, F, U) = q(X)P(F|X, U)q(U)$ . Then  $\mathcal{F} = \langle \log P(Y, F, U|X) \log P(X) \rangle_{q(U)q(X)}$  which simplifies into tractable components in much the same way as above.

difference  
between PPCA

B 3#

## A few references

- ▶ Jordan, Ghahramani, Jaakkola, Saul, 1999. [An introduction to variational methods for graphical models.](#) *Machine Learning* **37**:183–233.
- ▶ Attias, 2000. [A variational Bayesian framework for graphical models.](#) *NIPS 12*.  
<http://www.gatsby.ucl.ac.uk/publications/papers/03-2000.ps>
- ▶ Beal, 2003. [Variational algorithms for approximate Bayesian inference.](#) *PhD thesis*, Gatsby Unit, UCL. <http://www.cse.buffalo.edu/faculty/mbeal/thesis/>
- ▶ Winn, 2003. [Variational message passing and its applications.](#) *PhD thesis*, Cambridge.  
<http://johnwinn.org/Publications/Thesis.html>; also **VIBES** software for conjugate-exponential graphs.

Some complexities:

- ▶ MacKay, 2001. [A problem with variational free energy minimization.](#)  
<http://www.inference.phy.cam.ac.uk/mackay/minima.pdf>
- ▶ Turner, MS, 2011. [Two problems with variational expectation maximisation for time-series models.](#)  
In Barber, Cemgil, Chiappa, eds., *Bayesian Time Series Models*.  
<http://www.gatsby.ucl.ac.uk/~maneesh/papers/turner-sahani-2010-iltn.pdf>
- ▶ Berkes, Turner, MS, 2008. [On sparsity and overcompleteness in image models.](#) *NIPS 20*.  
<http://www.gatsby.ucl.ac.uk/~maneesh/papers/berkes-etal-2008-nips.pdf>
- ▶ Giordano, R, Broderick, T, and Jordan, MI, 2015. [Linear response methods for accurate covariance estimates from mean field variational Bayes.](#) *NIPS*