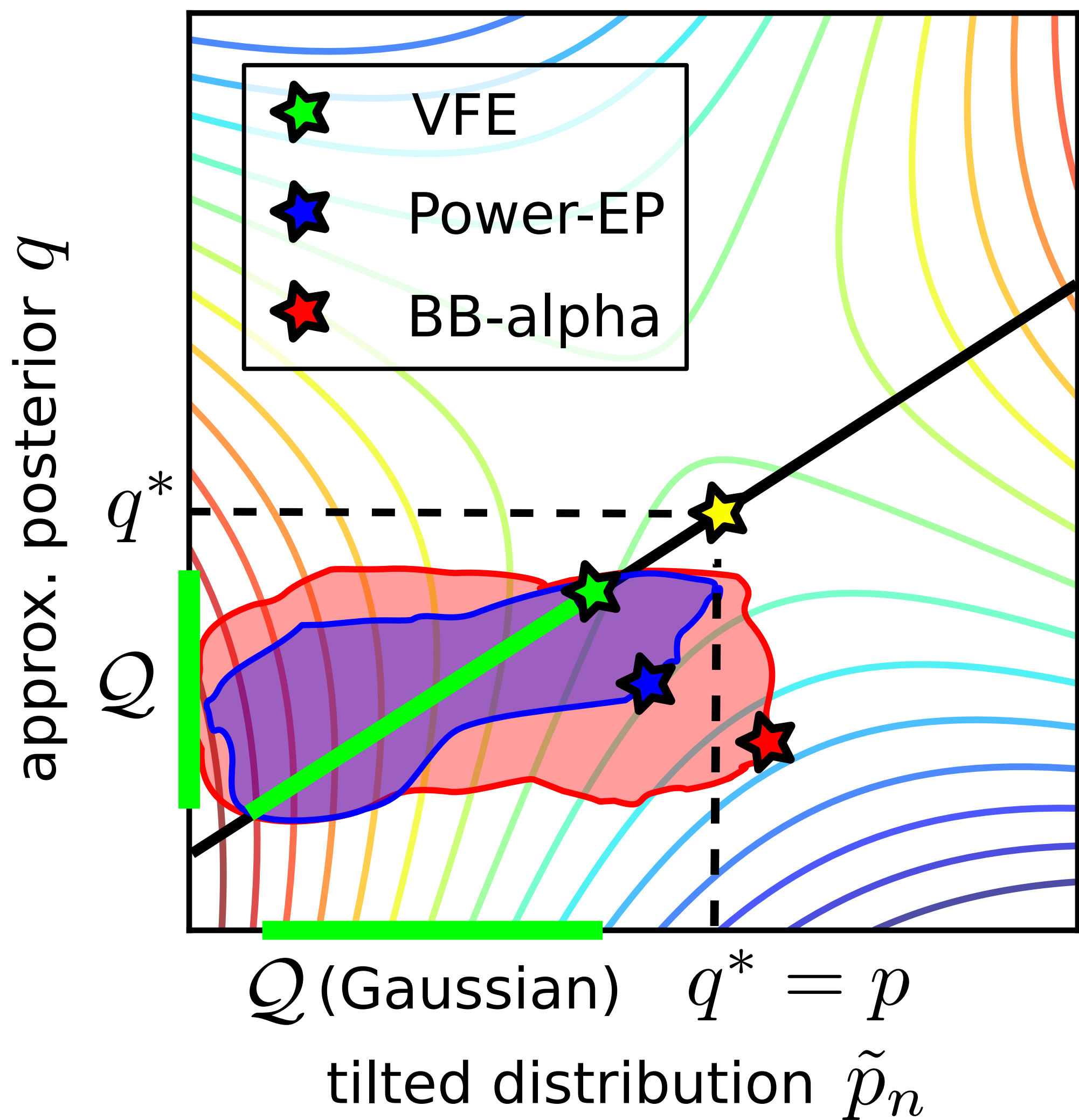


A Unifying Approximate Inference Framework from Variational Free Energy Relaxation

Yingzhen Li and Richard E. Turner, University of Cambridge

BIG PICTURE: APPROXIMATE INFERENCE AS CONSTRAINT RELAXATION



- Approximate posterior is obtained by solving a constrained minimisation problem with the following energy (with $\sum_n \frac{1}{\alpha_n} \neq 1$):

$$\min_{q, \{\tilde{p}_n\}} \mathcal{F}(q, \{\tilde{p}_n\}) = \left(1 - \sum_n \frac{1}{\alpha_n}\right) \text{KL}[q||p_0] - \sum_n \frac{1}{\alpha_n} \mathbb{E}_{\tilde{p}_n} \left[\log \frac{p_0(\theta) f_n(\theta)^{\alpha_n}}{\tilde{p}_n(\theta)} \right].$$

- VFE constraints:** $\tilde{p}_n = q, \forall n$; ($\mathcal{F}(q, \{\tilde{p}_n\})$ simplified to $\mathcal{F}_{\text{VFE}}(q)$)
- Power-EP constraints:** $\mathbb{E}_q[\phi(\theta)] = \mathbb{E}_{\tilde{p}_n}[\phi(\theta)], \forall n$; (moment matching)
- (*NEW*) BB- α constraints:** $N\mathbb{E}_q[\phi(\theta)] = \sum_n \mathbb{E}_{\tilde{p}_n}[\phi(\theta)]$; (moment averaging)
- (*NEW*)** Mixing distributed EP and BB- α ;
- (*NEW*)** Extensions to latent variable models.

FROM VFE TO POWER EP

- Target distribution $p(\theta) \propto p_0(\theta) \prod_n f_n(\theta)$, e.g., $f_n(\theta) = p(x_n|\theta)$;
- VI minimises the **variational free energy** (VFE):

$$\min_q \mathcal{F}_{\text{VFE}}(q) = \mathbb{E}_q \left[\log \frac{q(\theta)}{p_0(\theta)} - \sum_{n=1}^N \log f_n(\theta) \right] = \text{KL}[q||p] - \text{const.}$$

- An equivalent optimisation problem, **subject to** $\tilde{p}_n = q, \forall n$:

$$\min_{q, \{\tilde{p}_n\}} \left(1 - \sum_n \frac{1}{\alpha_n}\right) \text{KL}[q||p_0] - \sum_n \frac{1}{\alpha_n} \mathbb{E}_{\tilde{p}_n} \left[\log \frac{p_0(\theta) f_n(\theta)^{\alpha_n}}{\tilde{p}_n(\theta)} \right].$$

- Constraint relaxation:** from $q = \tilde{p}_n, \forall n$ to **moment matching**:

$$\mathbb{E}_q[\phi(\theta)] = \mathbb{E}_{\tilde{p}_n}[\phi(\theta)], \forall n.$$

- Introduce a new variable by using the KL duality:

$$-\text{KL}[q||p_0] = \min_{\lambda_q(\theta)} -\mathbb{E}_q[\lambda_q(\theta)] + \log \mathbb{E}_{p_0} [\exp[\lambda_q(\theta)]] .$$

- Write λ_{-n} as the Lagrange multiplier for moment matching constraints, and solve the Lagrangian:

$$\tilde{p}_n(\theta) = \frac{1}{Z_n} p_0(\theta) f_n(\theta)^{\alpha_n} \exp[\lambda_{-n}^T \phi(\theta)] ,$$

$$\left(\sum_n \frac{1}{\alpha_n} - 1 \right) \lambda_q(\theta) = \sum_n \frac{1}{\alpha_n} \lambda_{-n}^T \phi(\theta) + \text{const.}$$

- Defining $\lambda_q(\theta) = \lambda_{-n}^T \phi(\theta) + \text{const}$ and substituting in the fixed point solutions, we arrive the **power-EP (dual) energy**:

$$\min_{\lambda_q} \max_{\{\lambda_{-n}\}} \left(\sum_n \frac{1}{\alpha_n} - 1 \right) \log Z_q - \sum_n \frac{1}{\alpha_n} \log Z_n, \\ \text{subject to } \left(\sum_n \frac{1}{\alpha_n} - 1 \right) \lambda_q = \sum_n \frac{1}{\alpha_n} \lambda_{-n} .$$

- Approximation: $q(\theta) = \frac{1}{Z_q} p_0(\theta) \exp[\lambda_q^T \phi(\theta)]$.

- Local factor parameterisation returns power EP:**

define $\lambda_n = (\lambda_q - \lambda_{-n})/\alpha_n$,
then rewrite $\lambda_q = \sum_n \lambda_n$ and $\lambda_{-n} = \lambda_q - \alpha_n \lambda_n$.
Now $f_n(\theta) \approx \exp[\lambda_n^T \phi(\theta)]$.

BB- α & DISTRIBUTED ALGORITHMS

Deriving black-box alpha:

- Power-EP proposes N sets of constraints (main reason for memory overhead);
- Idea: reduce to **weighted moment averaging**:
 $\mathbb{E}_q[\phi] = \sum_n w_n \mathbb{E}_{\tilde{p}_n}[\phi], \quad \sum_n w_n = 1$;
- Choose $\alpha_n = \alpha, w_n = 1/N$ and solve the Lagrangian again, we arrive at the **BB- α (dual) energy**:

$$\min_{\lambda_q} \left(\frac{N}{\alpha} - 1 \right) \log Z_q - \frac{1}{\alpha} \sum_n \log \int p_0(\theta) f_n(\theta)^\alpha \exp[\lambda_q \phi(\theta)] d\theta .$$

Distributed Power-EP algorithms:

- In this case factor indices are divided into subsets N_1, N_2, \dots, N_K
- Rewrite $p(\theta) \propto p_0(\theta) \prod_k F_k(\theta), F_k(\theta) = \prod_{n_k \in N_k} f_{n_k}(\theta)$,
- ...and repeat the same procedure!
- Alternatively, add extra constraints $\tilde{p}_i = \tilde{p}_j \forall i, j \in N_k$.

Mixing BB- α and distributed methods:

- Distributed BB- α : let $\mathbb{E}_q[\phi(\theta)] = \frac{1}{N} \sum_k \mathbb{E}_{\tilde{p}_k}[\phi(\theta)]$.
- Nesting BB- α in distributed EP: let
 $\mathbb{E}_q[\phi(\theta)] = \frac{1}{|N_k|} \sum_{n_k \in N_k} \mathbb{E}_{\tilde{p}_{n_k}}[\phi(\theta)]$.

EXTENSION: LATENT VARIABLE MODELS

- Assume factorised approximation $q(\theta, z_n) = q(\theta) \prod_n q(z_n)$:

$$\mathcal{F}_{\text{VFE}}(q) = \mathbb{E}_q \left[\log \frac{q(\theta)}{p_0(\theta)} + \sum_n \log \frac{q(z_n)}{p_0(z_n)} - \sum_{n=1}^N \log f_n(\theta, z_n) \right] .$$

- Decouple q to \tilde{p}_n similarly, and
- ...VI considers constraints $\tilde{p}_n(\theta, z_n) = q(\theta)q(z_n), \forall n$;
- Different constraint relaxations return different algorithms!
 - Full EP: $\mathbb{E}_{\tilde{p}_n}[\phi(\theta), \psi(z_n)] = \mathbb{E}_q[\phi(\theta), \psi(z_n)]$;
 - Nesting EP in VI: $\mathbb{E}_{\tilde{p}_n}[\psi(z_n)] = \mathbb{E}_q[\psi(z_n)], q(\theta) = \tilde{p}_n(\theta)$;
 - Nesting VI in EP: $\mathbb{E}_{\tilde{p}_n}[\phi(\theta)] = \mathbb{E}_q[\phi(\theta)], \tilde{p}_n(z_n) = q(z_n), \tilde{p}_n(\theta, z_n) = \tilde{p}_n(\theta)\tilde{p}_n(z_n)$;
 - "Tilted" VMP: $\mathbb{E}_{\tilde{p}_n}[\phi(\theta), \psi(z_n)] = \mathbb{E}_q[\phi(\theta), \psi(z_n)], \tilde{p}_n(\theta, z_n) = \tilde{p}_n(\theta)\tilde{p}_n(z_n)$.