



# Approximate Inference in Deep Gaussian Processes by Minimizing Alpha Divergences

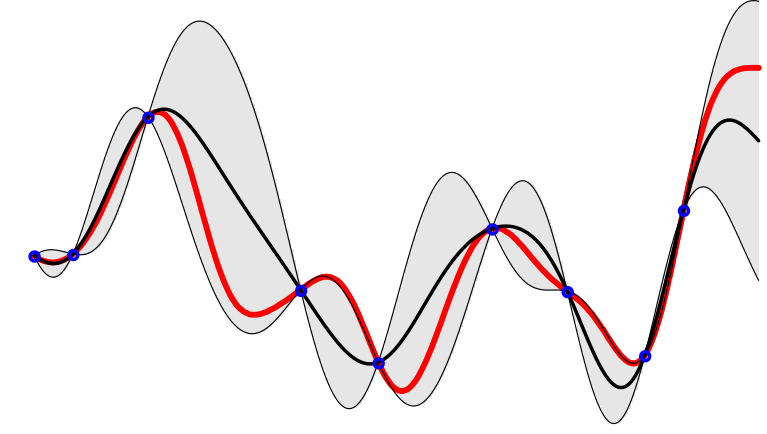
Carlos Villacampa-Calvo<sup>1</sup>, Gonzalo Hernández-Muñoz, Daniel Hernández-Lobato<sup>1</sup>  
 (1) - Universidad Autónoma de Madrid



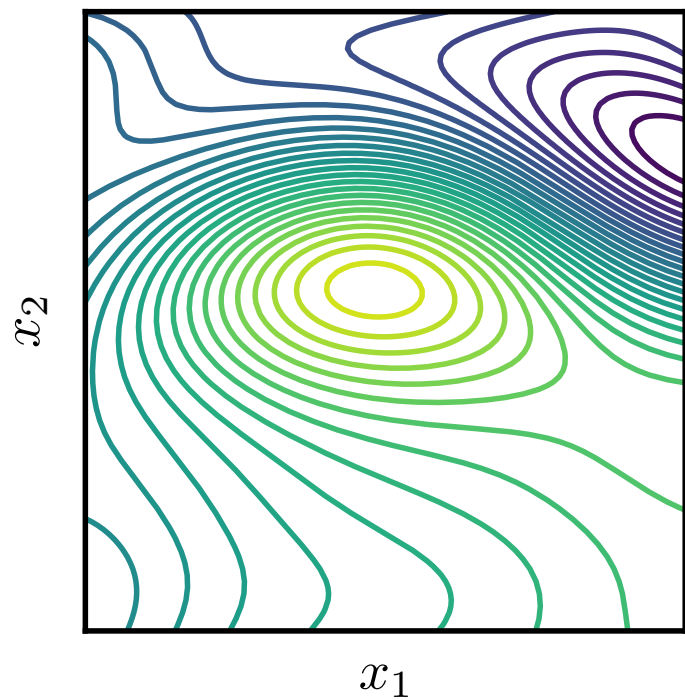
## 1 - Introduction to Gaussian Processes

GPs are flexible models that can output closed form predictive distributions:

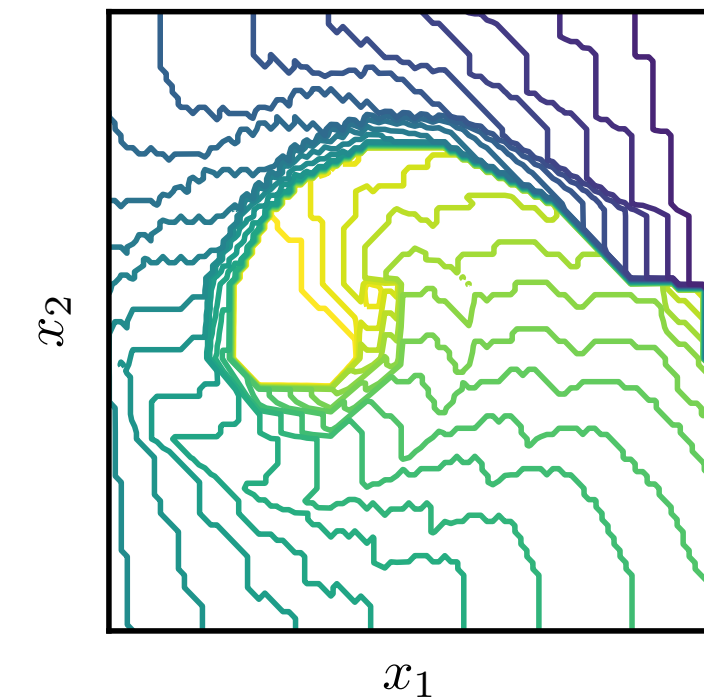
GP Predictive Distribution



GP Fit



Target Function

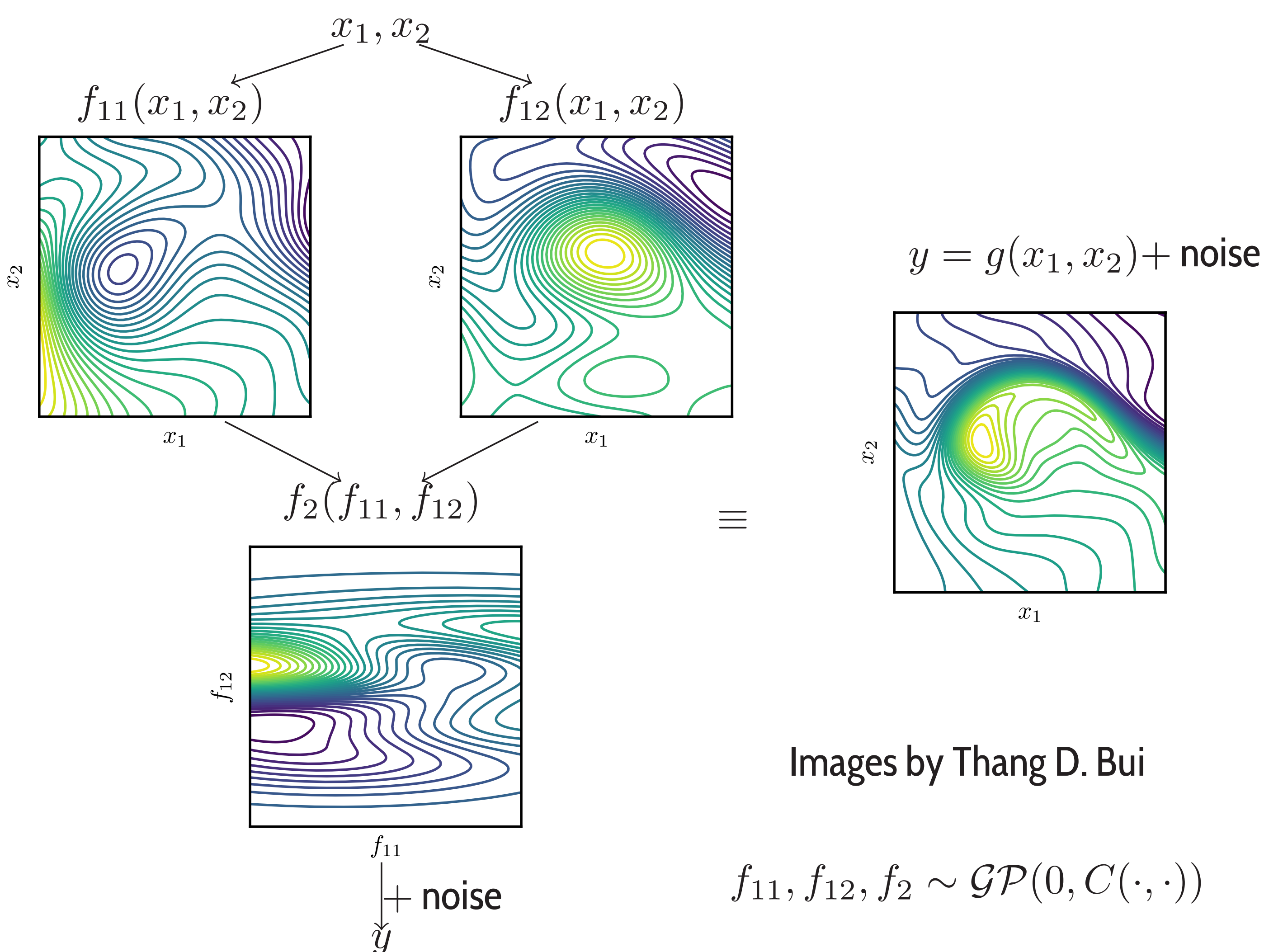
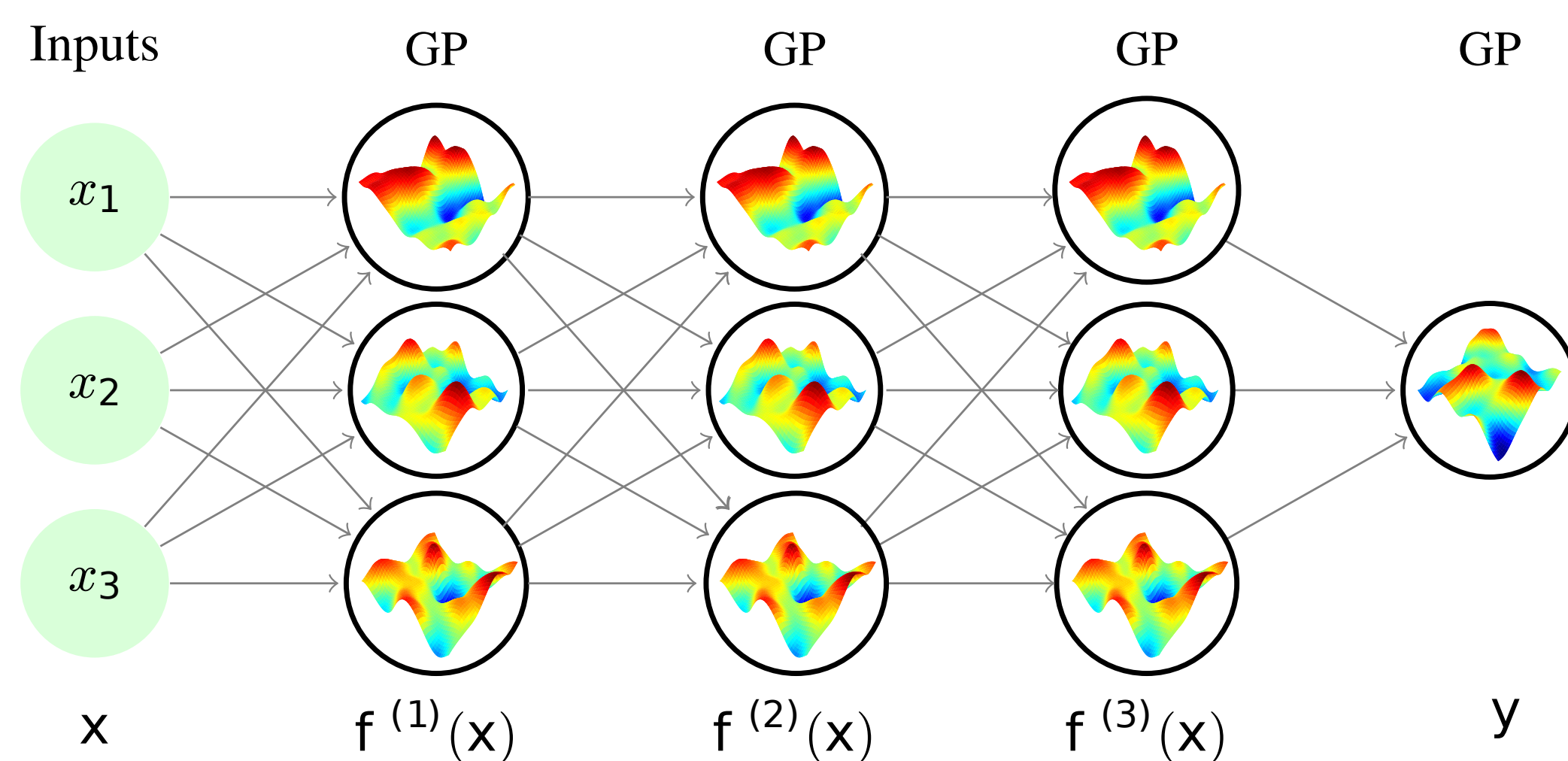


Limitations (Bui et al., 2016):

- Cubic training cost!
- Sparse approximations may limit model flexibility!
- Often the target is assumed to be stationary!

## 2 - Deep Gaussian Processes

Deep neural network in which the activations functions are GPs:



## 3 - Bayesian Inference

Posterior over latent functions (typically at the observed data  $\mathbf{X}$ ):

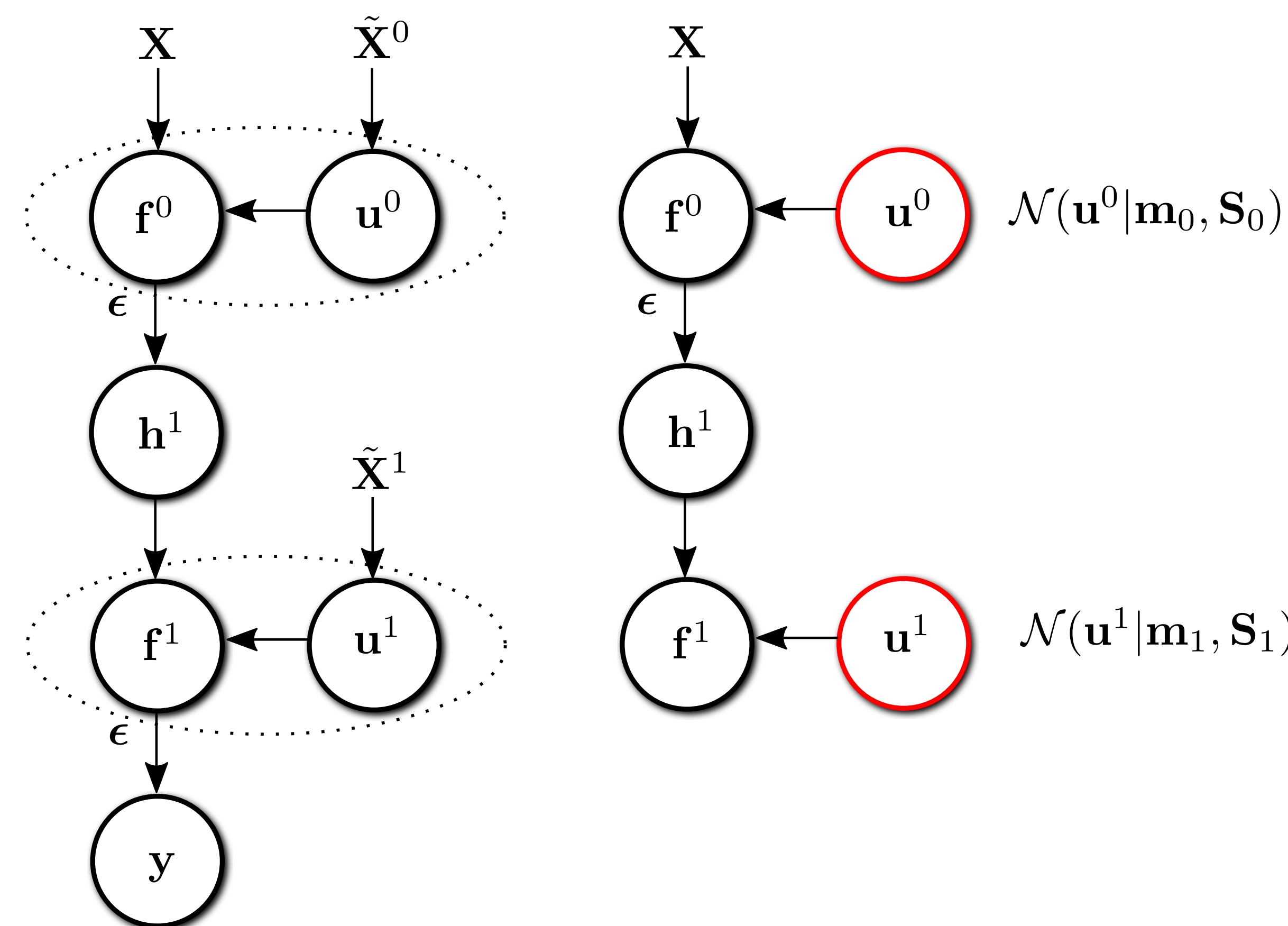
$$p(\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3 | \mathbf{Y}) = \frac{p(\mathbf{f}_1)p(\mathbf{f}_2)p(\mathbf{f}_3)p(\mathbf{Y}|\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3, \mathbf{X})}{p(\mathbf{Y})}$$

- GP priors
- Likelihood function
- Marginal likelihood

But the posterior  $p(\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3 | \mathbf{Y})$  is **intractable**.

## 4 - Sparse GPs and Posterior Approximation

Inducing points  $\tilde{\mathbf{X}}^i$  associated to each GP:  $p(\mathbf{f}^i, \mathbf{u}^i | \mathbf{h}^i) = p(\mathbf{f}^i | \mathbf{u}^i, \mathbf{h}^i)p(\mathbf{u}^i)$



$$q(\mathbf{f}^0, \mathbf{u}^0, \mathbf{h}^1, \dots, \mathbf{f}^L, \mathbf{u}^L) = q(\mathbf{u}^0)p(\mathbf{f}^0 | \mathbf{u}^0)p(\mathbf{h}^1 | \mathbf{f}^1) \dots p(\mathbf{f}^L | \mathbf{u}^L, \mathbf{h}^L)q(\mathbf{u}^L)$$

## 5 - Alpha Divergence Minimization

$$D_\alpha(p||q) = \frac{\int_\theta (\alpha p(\theta) + (1-\alpha)q(\theta) - p(\theta)^\alpha q(\theta)^{1-\alpha}) d\theta}{\alpha(1-\alpha)}$$

[Amari, 1985].

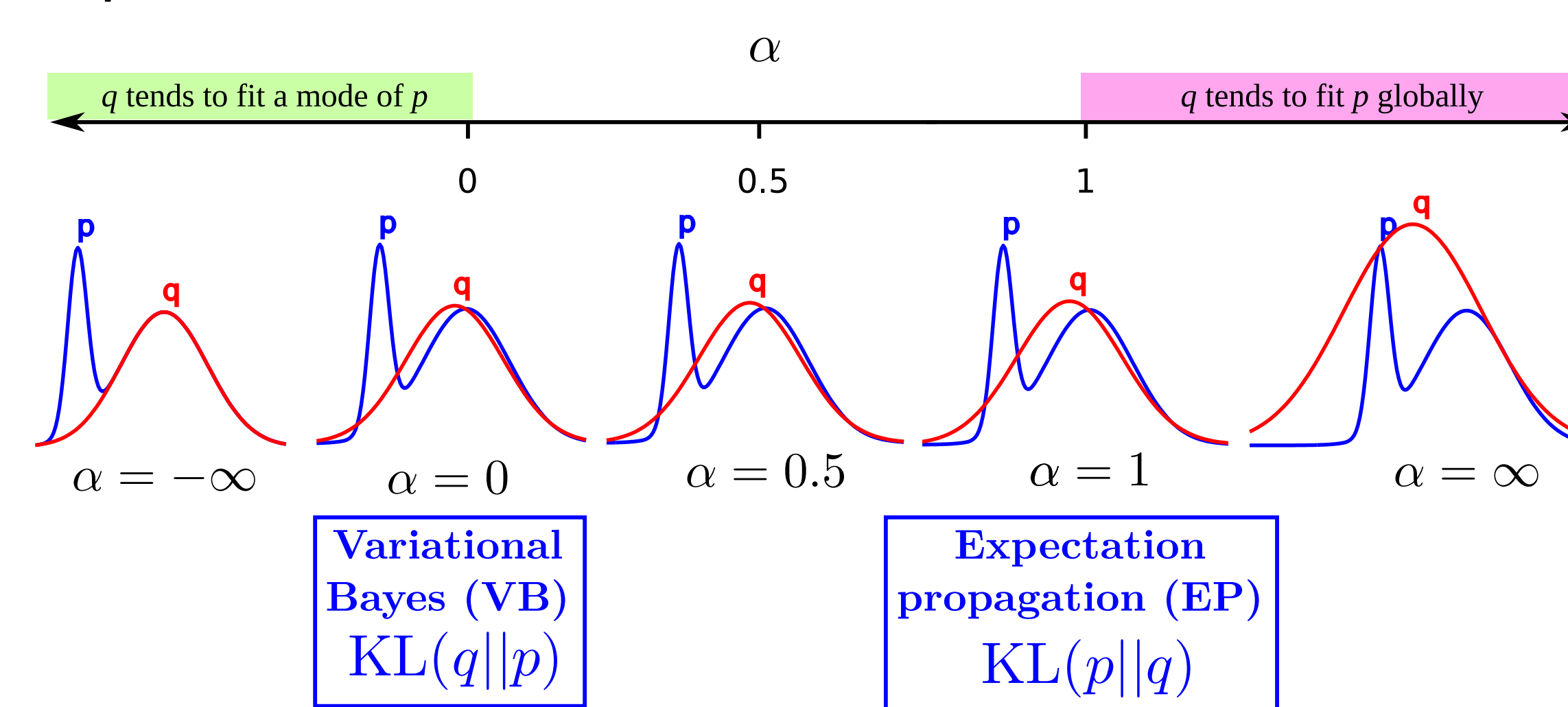


Figure source: [Minka, 2005].

## 5 - Local $\alpha$ -divergence Minimization via Power EP

Approximates  $p(\theta) \propto p_0(\theta) \prod_{n=1}^N f_n(\theta)$  with  $q(\theta) \propto p_0(\theta) \prod_{n=1}^N \tilde{f}_n(\theta)$  (Minka, 2004)

$$p(\theta) \propto p_0(\theta) f_1(\theta) f_2(\theta) f_3(\theta) \approx q(\theta) \propto p_0(\theta) \tilde{f}_1(\theta) \tilde{f}_2(\theta) \tilde{f}_3(\theta)$$

The  $\tilde{f}_n$  are tuned by minimizing the local  $\alpha$ -divergences

$$D_\alpha[p_n||q] \text{ for } n = 1, \dots, N, \text{ where } p_n(\theta) \propto f_n(\theta) \prod_{j \neq n} \tilde{f}_j(\theta) \\ q(\theta) \propto \tilde{f}_n(\theta) \prod_{j \neq n} \tilde{f}_j(\theta)$$

Power EP steps to refine  $f_n$ :

- 1 Compute cavity:  $q^{\setminus \alpha n} \propto q / \tilde{f}_n^\alpha$ .
- 2 Minimize  $\text{KL}(Z_n^{-1} f_n^\alpha q^{\setminus \alpha n} || q)$  to find  $q^{\text{new}}$ .
- 3 Update factor:  $f_n^{\text{new}} = (Z_n q^{\text{new}} / q^{\setminus \alpha n})^{\frac{1}{\alpha}}$ .

At convergence the moments of  $p_n = Z_n^{-1} f_n^\alpha q^{\setminus \alpha n}$  and  $q$  match!

$$\nabla_{\eta q} D_\alpha[p_n||q] = \frac{Z_{\tilde{p}}}{\alpha} (\mathbb{E}_q[s(\theta)] - \mathbb{E}_{\tilde{p}}[s(\theta)]) \propto \nabla_{\eta q} \text{KL}[\tilde{p}||q]$$

where  $\tilde{p} \propto (f_n q^{\setminus \alpha n})^\alpha q^{1-\alpha} = f_n^\alpha q^{\setminus \alpha n}$ .

The likelihood factors are simply  $p(y_i | \mathbf{f}_i^L)$ !

KL minimization by adding approx. factors that depend only on  $\mathbf{u}^0, \dots, \mathbf{u}^L$ !

## 6 - Approximate Power Expectation Propagation

The Power-EP approximation to the **evidence** is given by

$$\log Z_{\text{PEP}} = \log Z_q - \log Z_{\text{prior}} + \sum_{n=1}^N \frac{1}{\alpha} \log \mathbb{E}_q \left[ \left( \frac{f_n(\theta)}{\tilde{f}_n(\theta)} \right)^\alpha \right]$$

The power-EP solution for  $q$  can be obtained by solving

$$\max_q \min_{\tilde{f}_1, \dots, \tilde{f}_N} \log Z_{\text{PEP}} \text{ subject to } q(\theta) = p_0(\theta) \prod_{n=1}^N \tilde{f}_n(\theta)$$

Solved with **double-loop** algorithm (Heskes, 2002). **Too slow in practice!**

By following (Li et al., 2015):

$$p(\theta) \propto p_0(\theta) f_1(\theta) f_2(\theta) f_3(\theta) \approx q(\theta) \propto p_0(\theta) \tilde{f}_1(\theta) \tilde{f}_2(\theta) \tilde{f}_3(\theta)$$

We tie the factor approximations

$$p(\theta) \propto p_0(\theta) f_1(\theta) f_2(\theta) f_3(\theta) \approx q(\theta) \propto p_0(\theta) \tilde{f}(\theta)^N$$

- $\max_q \min_{\tilde{f}_1, \dots, \tilde{f}_N}$  problem  $\rightarrow$   $\max_q$  problem, **no double-loop needed!**

The expectations can be approximated via Monte Carlo by sampling!