



# Perturbation Theory for Variational Inference

Manfred Opper<sup>1</sup>, Marco Fraccaro<sup>2</sup>, Ulrich Paquet<sup>3</sup>, Alex Susemihl<sup>1</sup> and Ole Winther<sup>2</sup>  
<sup>1</sup>Technical University of Berlin    <sup>2</sup>Technical University of Denmark    <sup>3</sup>Apple



## Overview

- ▶ We view the variational free energy and its accompanying evidence lower bound as a first-order term from a perturbation of the true log partition function and derive a power series of corrections.
- ▶ This allows for example to get better estimates of normalizing constants or to correct predictions in latent variable models whose permutation-invariant parameters are self-pruned and ignored by the variational approximation (e.g. matrix factorization).

## Variational Inference

Given a random variable  $x$  with probability distribution  $p(x) = \frac{\mu(x)e^{-H(x)}}{Z}$ , we often need to compute

$$-\log Z = -\log \int \mu(x) e^{-H(x)} dx \quad \text{or} \quad E[f(x)] = \int f(x) \mu(x) e^{-H(x)} dx.$$

If these are not analytically tractable, we introduce a *variational approximation*  $q(x) = \frac{1}{Z_q} \mu(x) e^{-H_q(x)}$  whose parameters are found minimizing the free energy

$$F[q] = KL[q||p] - \log Z = E_q \left[ \log \frac{q(x)}{p(x)} \right] - \log Z = E_q[V(x)] - \log Z_q,$$

where we have defined  $V(x) = H(x) - H_q(x)$ .

## Perturbative corrections

*Perturbation theory* aims at finding approximate solutions to a problem given exact solutions of a simpler related sub-problem (the VB solution in our case).

**Normalizing constant.** By defining  $\hat{H}_\lambda = H_q + \lambda V = (1 - \lambda)H_q + \lambda H$  (notice that this gives  $\hat{H}_1 = H$  and  $\hat{H}_0 = H_q$ ), we can write

$$\begin{aligned} -\log \int \mu(x) e^{-\hat{H}_\lambda(x)} dx &= -\log Z_q - \log E_q[e^{-\lambda V(x)}] \\ &= \underbrace{-\log Z_q + \lambda E_q[V]}_{F[q] \text{ for } \lambda=1} - \frac{\lambda^2}{2} E_q[(V - E_q[V])^2] + \frac{\lambda^3}{3!} E_q[(V - E_q[V])^3] + \dots \end{aligned}$$

Knowing the VB solution, we can therefore correct our estimate of  $\log Z$  using higher order terms.

Note that this may not be a convergent series, but lead to an asymptotic expansion only.

**Expectations.** Defining  $E_\lambda$  as the expectation with respect to  $p_\lambda(x) = \mu(x) e^{-H_\lambda(x)}$ , so that  $E = E_1$  and  $E_q = E_0$ , we get:

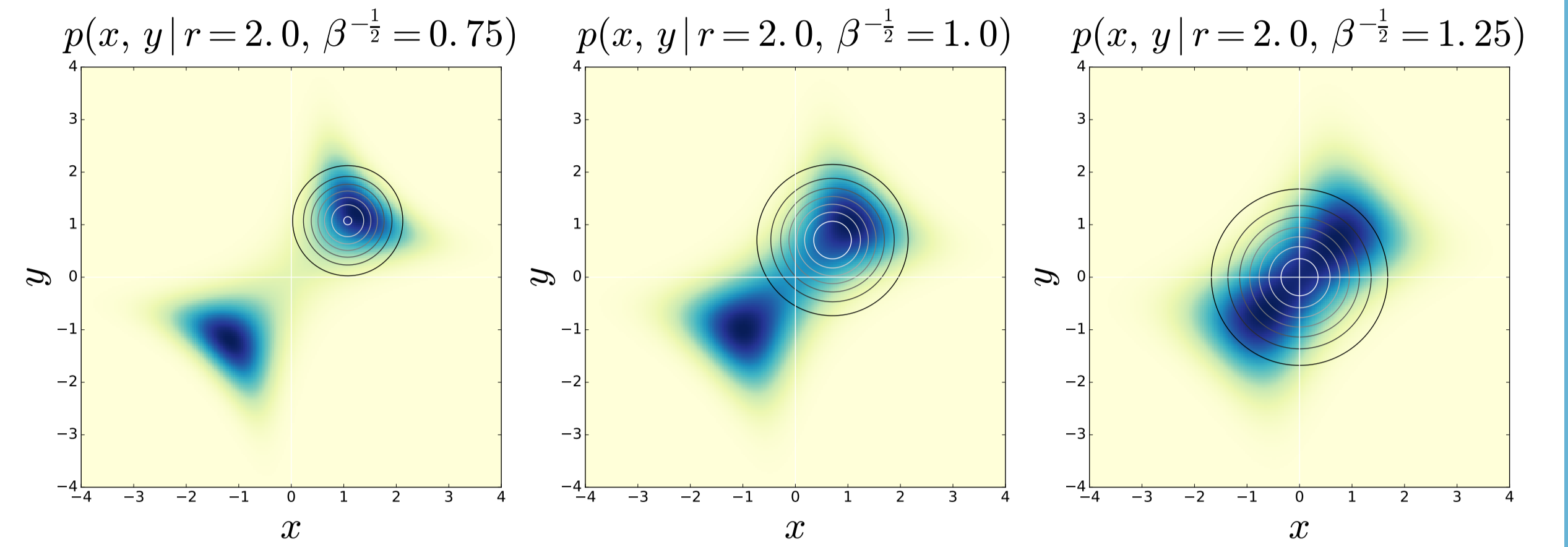
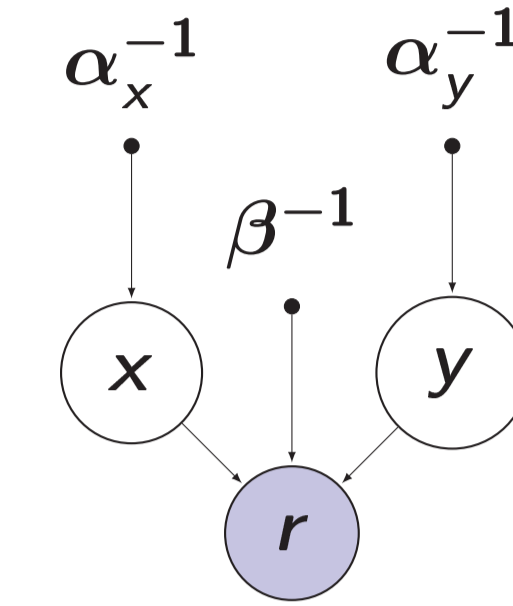
$$\begin{aligned} E_\lambda[f(x)] &= \frac{\int f(x) \mu(x) e^{-\hat{H}_0(x) - \lambda V(x)} dx}{\int \mu(x) e^{-\hat{H}_0(x) - \lambda V(x)} dx} = \frac{E_0[f(x) e^{-\lambda V(x)}]}{E_0[e^{-\lambda V(x)}]} \\ &= \frac{E_0 \left[ f(x) \left( 1 - \lambda V + \frac{\lambda^2}{2} V^2 - \frac{\lambda^3}{3!} V^3 \pm \dots \right) \right]}{E_0 \left[ \left( 1 - \lambda V + \frac{\lambda^2}{2} V^2 - \frac{\lambda^3}{3!} V^3 \pm \dots \right) \right]} \quad \frac{1}{1-z} = 1+z+z^2+\dots \\ &= \underbrace{E_0[f(x)]}_{\log Z} - \lambda \text{Cov}_0[f(x), V] - \lambda^2 E_0[V] \text{Cov}_0[f(x), V] + \frac{\lambda^2}{2} \text{Cov}_0[f(x), V^2] \pm \dots \end{aligned}$$

We can correct the expectation w.r.t the VB solution using higher order terms.

## Example: Variational Matrix Factorization

**Model:**

$$\begin{aligned} r &= xy + \epsilon \\ \epsilon &\sim \mathcal{N}(0, \beta^{-1}) \\ x &\sim \mathcal{N}(0, \alpha_x^{-1}) \\ y &\sim \mathcal{N}(0, \alpha_y^{-1}) \end{aligned}$$

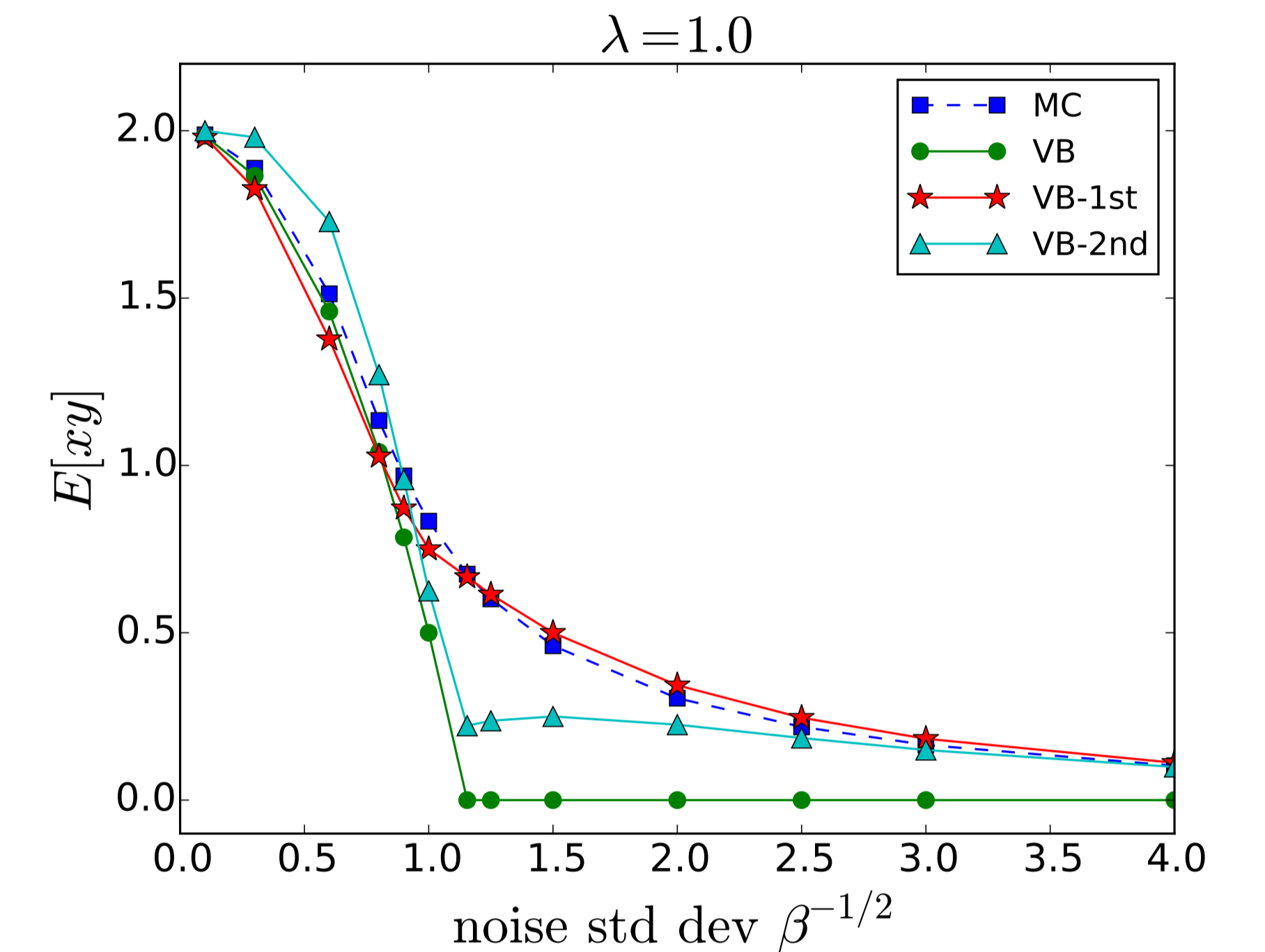


**VB solution:**

$$q(x) = \mathcal{N}(\mu_x, \gamma_x^{-1}), \quad q(y) = \mathcal{N}(\mu_y, \gamma_y^{-1})$$

**Predictions + corrections:**

$$\begin{aligned} H(x, y) &= \frac{1}{2} [\beta(r - xy)^2 + \alpha_x x^2 + \alpha_y y^2] \\ H_q(x, y) &= \frac{1}{2} [\gamma_x(x - \mu_x)^2 + \gamma_y(y - \mu_y)^2] \\ V(x, y) &= H(x, y) - H_q(x, y) \\ E_\lambda[xy] &= E_0[xy] - \lambda \text{Cov}_0[xy, V(x, y)] \pm \dots \end{aligned}$$



## Example: Variational Inference for Sparse GP Regression

- ▶ We can compute the predictive mean and covariance using a sparse approximation for the GP and then include the perturbative corrections.

$$y_i(t_i) = \underbrace{3 \text{sinc}(t_i)}_{x(t_i)} + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$x \sim GP(0, K) \text{ with RBF kernel}$$

