

Integrated Non-Factorized Variational Inference

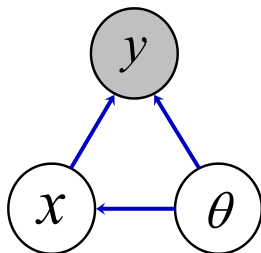
Shaobo Han, Xuejun Liao and Lawrence Carin

Duke University

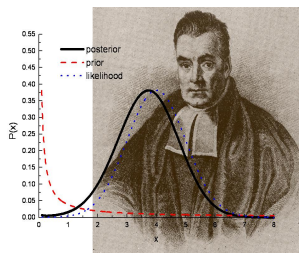
February 27, 2014



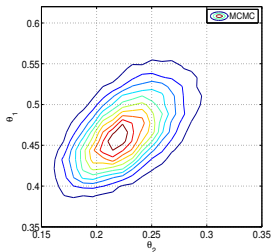
World



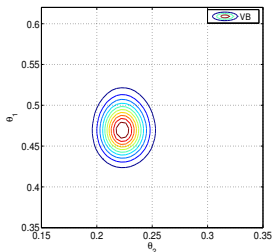
Graphical Models



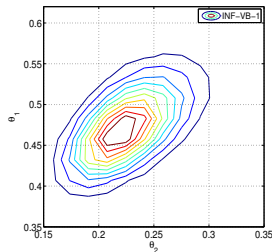
Posterior Inference



MCMC



VB



Our method

For full posterior inference, our method is

- ▶ A fast deterministic alternative to MCMC
- ▶ More accurate than mean-field variational Bayes (VB)

- ▶ Introduction
- ▶ Integrated Nested Laplace Approximation (INLA)
- ▶ **Integrated Non-Factorized Variational Bayes (INF-VB)**
- ▶ Applications
- ▶ Summary

Consider a general Bayesian hierarchical model

- ▶ Observation model: $\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$
- ▶ Latent variables: $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$
- ▶ Hyperparameters: $\boldsymbol{\theta} \sim p(\boldsymbol{\theta})$

Problem of Interest

Consider a general Bayesian hierarchical model

- ▶ Observation model: $\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$
- ▶ Latent variables: $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$
- ▶ Hyperparameters: $\boldsymbol{\theta} \sim p(\boldsymbol{\theta})$

Posterior inference:

$$p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$$

Problem of Interest

Consider a general Bayesian hierarchical model

- ▶ Observation model: $\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$
- ▶ Latent variables: $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$
- ▶ Hyperparameters: $\boldsymbol{\theta} \sim p(\boldsymbol{\theta})$

Posterior inference:

$$\begin{array}{ccc} p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) & \longrightarrow & p(\mathbf{x}|\mathbf{y}) \\ \downarrow & \searrow & \\ p(\boldsymbol{\theta}|\mathbf{y}) & & p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}) \end{array}$$

Problem of Interest

Consider a general Bayesian hierarchical model

- ▶ Observation model: $\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$
- ▶ Latent variables: $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$
- ▶ Hyperparameters: $\boldsymbol{\theta} \sim p(\boldsymbol{\theta})$

Posterior inference:

$$\begin{array}{ccc} p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) & \xrightarrow{\quad} & p(\mathbf{x}|\mathbf{y}) \\ \downarrow \uparrow & \swarrow \nearrow & \\ p(\boldsymbol{\theta}|\mathbf{y}) & + & p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}) \end{array}$$

Problem of Interest

Consider a general Bayesian hierarchical model

- ▶ Observation model: $\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$
- ▶ Latent variables: $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$
- ▶ Hyperparameters: $\boldsymbol{\theta} \sim p(\boldsymbol{\theta})$

Posterior inference:

$$\begin{array}{ccc} p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) & \xrightarrow{\quad} & p(\mathbf{x}|\mathbf{y}) \\ \downarrow \uparrow & \swarrow \nwarrow & \\ p(\boldsymbol{\theta}|\mathbf{y}) & + & p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}) \end{array}$$

The exact joint posterior

$$p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta})}{p(\mathbf{y})} = \frac{p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\mathbf{x}d\boldsymbol{\theta}}$$

can be difficult to evaluate.

Sampling based methods:

- ▶ Markov chain Monte Carlo (MCMC)

Deterministic alternatives:

- ▶ Laplace approximation (LA)
- ▶ Variational inference
- ▶ Expectation propagation (EP)
- ▶ **Integrated nested Laplace approximation (INLA)¹**

¹Rue et al., 2009

- ▶ Introduction
- ▶ **Integrated Nested Laplace Approximation (INLA)**
- ▶ Integrated Non-Factorized Variational Bayes (INF-VB)
- ▶ Applications
- ▶ Summary

INLA in a Nutshell (1/3)

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

Demo:

$$\theta_k \in \mathcal{G}$$

²Kass & Steffey, 1989

INLA in a Nutshell (1/3)

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

Demo:

$$q_G(\mathbf{x}|\mathbf{y}, \theta_k) \quad \longleftarrow \quad \theta_k \in \mathcal{G}$$

²Kass & Steffey, 1989

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

Demo:

$$q_G(\mathbf{x}|\mathbf{y}, \theta_k) \longleftarrow \theta_k \in \mathcal{G}$$

1. Laplace approximation² :

$$q_G(\mathbf{x}|\mathbf{y}, \theta_k) = \mathcal{N}(\mathbf{x}; \mathbf{x}^*(\theta_k), \mathbf{H}(\mathbf{x}^*(\theta_k))^{-1}), \quad \forall \theta_k \in \mathcal{G}$$

where $\mathbf{x}^*(\theta_k) = \operatorname{argmax}_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}, \theta_k)$ is the posterior mode, and $\mathbf{H}(\mathbf{x}^*(\theta_k))$ is the Hessian matrix of the log posterior evaluated at the mode.

²Kass & Steffey, 1989

INLA in a Nutshell (2/3)

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

Demo:

$$q_{LA}(\theta|\mathbf{y}) \longleftarrow q_G(\mathbf{x}|\mathbf{y}, \theta_k) \longleftarrow \theta_k \in \mathcal{G}$$

³Tierney & Kadane, 1986

INLA in a Nutshell (2/3)

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

Demo:

$$q_{LA}(\theta|\mathbf{y}) \longleftarrow q_G(\mathbf{x}|\mathbf{y}, \theta_k) \longleftarrow \theta_k \in \mathcal{G}$$

According to the Bayes rule,

$$p(\theta|\mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y}, \theta)}{p(\mathbf{y})p(\mathbf{x}|\mathbf{y}, \theta)}, \quad \forall \mathbf{x} \quad (1)$$

³Tierney & Kadane, 1986

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

Demo:

$$q_{LA}(\theta|\mathbf{y}) \longleftarrow q_G(\mathbf{x}|\mathbf{y}, \theta_k) \longleftarrow \theta_k \in \mathcal{G}$$

According to the Bayes rule,

$$p(\theta|\mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y}, \theta)}{p(\mathbf{y})p(\mathbf{x}|\mathbf{y}, \theta)}, \quad \forall \mathbf{x} \quad (1)$$

2. Laplace's method of integration³ :

$$q_{LA}(\theta|\mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y}, \theta)}{p(\mathbf{y})q_G(\mathbf{x}|\mathbf{y}, \theta)} \Big|_{\mathbf{x}=\mathbf{x}^*(\theta)} \quad (2)$$

³Tierney & Kadane, 1986

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

Demo:

$$q_{LA}(\theta|\mathbf{y}) \quad + \quad q_G(\mathbf{x}|\mathbf{y}, \theta_k) \quad \leftarrow \quad \theta_k \in \mathcal{G}$$

$q(\mathbf{x}|\mathbf{y})$

The diagram illustrates the relationship between the latent variable distribution, the grid-based approximation, and the observed data distribution. It shows the equation $q_{LA}(\theta|\mathbf{y}) + q_G(\mathbf{x}|\mathbf{y}, \theta_k) \leftarrow \theta_k \in \mathcal{G}$ with $q(\mathbf{x}|\mathbf{y})$ above it. Blue arrows indicate that $q_G(\mathbf{x}|\mathbf{y}, \theta_k)$ is derived from $q_{LA}(\theta|\mathbf{y})$ and θ_k , and that $q(\mathbf{x}|\mathbf{y})$ is derived from $q_G(\mathbf{x}|\mathbf{y}, \theta_k)$.

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

Demo:

$$q_{LA}(\theta|\mathbf{y}) \quad + \quad q_G(\mathbf{x}|\mathbf{y}, \theta_k) \quad \leftarrow \quad \theta_k \in \mathcal{G}$$

$q(\mathbf{x}|\mathbf{y})$

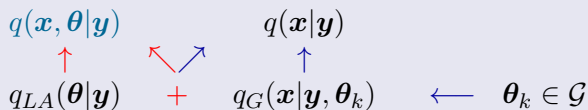
3. Numerical integration:

$$q(\mathbf{x}|\mathbf{y}) = \sum_k q_G(\mathbf{x}|\mathbf{y}, \theta_k) q_{LA}(\theta_k|\mathbf{y}) \Delta_k$$

with area weights Δ_k .

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

Demo:

$$q(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) \quad q(\mathbf{x} | \mathbf{y})$$


$q_{LA}(\boldsymbol{\theta} | \mathbf{y}) \quad + \quad q_G(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta}_k) \quad \leftarrow \quad \boldsymbol{\theta}_k \in \mathcal{G}$

3. Numerical integration:

$$q(\mathbf{x} | \mathbf{y}) = \sum_k q_G(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta}_k) q_{LA}(\boldsymbol{\theta}_k | \mathbf{y}) \Delta_k$$

with area weights Δ_k .

Benefits:

1. Preserves full posterior dependencies (i.e. joint density $q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$)
2. Computationally efficient (MCMC: hours or days, INLA: seconds or minutes)

Limitations:

1. Applies only to latent Gaussian models (LGMs)
2. No quantization for the accuracy of approximation $q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$
3. The dimension of $\boldsymbol{\theta}$ has to be no more than 5 or 6

Benefits:

1. Preserves full posterior dependencies (i.e. joint density $q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$)
2. Computationally efficient (MCMC: hours or days, INLA: seconds or minutes)

Limitations:

1. Applies only to latent Gaussian models (LGMs)
2. No quantization for the accuracy of approximation $q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$
3. The dimension of $\boldsymbol{\theta}$ has to be no more than 5 or 6

Our method addresses the first two limitations with INLA.

- ▶ Introduction
- ▶ Integrated Nested Laplace Approximation (INLA)
- ▶ **Integrated Non-Factorized Variational Bayes (INF-VB)**
- ▶ Applications
- ▶ Future Research

Variational inference turns Bayesian inference into optimization.

$$\min_{q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})} \text{KL}[q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})||p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})] \quad \text{s.t.} \quad q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) \in \mathcal{Q} \quad (3)$$

Evidence lower bound (ELBO):

Applying Jensen's inequality,

$$\begin{aligned} \ln p(\mathbf{y}) &= \ln \int \int q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) \frac{p(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})} d\mathbf{x}d\boldsymbol{\theta} \\ &\geq \int \int q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) \ln \frac{p(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})} d\mathbf{x}d\boldsymbol{\theta} := \mathcal{L} \end{aligned} \quad (4)$$

- ▶ The Jensen's gap: $\ln p(\mathbf{y}) - \mathcal{L} = \text{KL}(q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})||p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}))$
- ▶ The variational distribution $q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ is commonly restricted to tractable families \mathcal{Q}

Mean-Field Variational Bayes (VB)

Assumes **factorized** form: $q(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) = q(\mathbf{x})q(\boldsymbol{\theta})$, then

$$\begin{aligned} q^*(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) &= \operatorname{argmin}_{q(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})} \operatorname{KL}(q(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) || p(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})) \\ &= \operatorname{argmin}_{q(\mathbf{x}), q(\boldsymbol{\theta})} \int \int q(\mathbf{x})q(\boldsymbol{\theta}) \ln \frac{q(\mathbf{x})q(\boldsymbol{\theta})}{p(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})} d\mathbf{x}d\boldsymbol{\theta} \end{aligned}$$

Remarks:

- ▶ Easily derived and in close form for conjugate models
- ▶ Challenging for non-conjugate models
- ▶ Ignores posterior dependencies and impairs the accuracy
- ▶ A poor approximation for a multi-modal distribution

Mean-Field Variational Bayes (VB)

Assumes **factorized** form: $q(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) = q(\mathbf{x})q(\boldsymbol{\theta})$, then

$$\begin{aligned} q^*(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) &= \operatorname{argmin}_{q(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})} \operatorname{KL}(q(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) || p(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})) \\ &= \operatorname{argmin}_{q(\mathbf{x}), q(\boldsymbol{\theta})} \int \int q(\mathbf{x})q(\boldsymbol{\theta}) \ln \frac{q(\mathbf{x})q(\boldsymbol{\theta})}{p(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})} d\mathbf{x}d\boldsymbol{\theta} \end{aligned}$$

Remarks:

- ▶ Easily derived and in close form for conjugate models
- ▶ Challenging for non-conjugate models
- ▶ Ignores posterior dependencies and impairs the accuracy
- ▶ A poor approximation for a multi-modal distribution

Our non-factorized variational method addresses these issues with mean-field VB.

Consider **non-factorized** form:

$$q(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) = q(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta}) q_d(\boldsymbol{\theta} | \mathbf{y}) \quad (5)$$

- ▶ \mathbf{x} and $\boldsymbol{\theta}$ are still coupling
1. The **continuous** approximation $q(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta})$ is very flexible
 - ▶ Gaussian
 - ▶ Mean-Field
 2. The **discretized** approximation $q_d(\boldsymbol{\theta} | \mathbf{y})$ is a finite mixture of Dirac-delta distributions,

$$q_d(\boldsymbol{\theta} | \mathbf{y}) = \sum_k \omega_k \delta_{\boldsymbol{\theta}_k}(\boldsymbol{\theta}), \quad \omega_k = q_d(\boldsymbol{\theta}_k | \mathbf{y}), \quad \sum_k \omega_k = 1 \quad (6)$$

Within the proposed hybrid family, the optimal variational distribution is

$$\begin{aligned} q^*(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) &= \operatorname{argmin}_{q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})} \operatorname{KL}(q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) || p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})) \\ &= \operatorname{argmin}_{q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}), q_d(\boldsymbol{\theta}|\mathbf{y})} \int \int q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) q_d(\boldsymbol{\theta}|\mathbf{y}) \ln \frac{q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) q_d(\boldsymbol{\theta}|\mathbf{y})}{p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})} d\mathbf{x} d\boldsymbol{\theta} \\ &= \operatorname{argmin}_{q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k), q_d(\boldsymbol{\theta}_k|\mathbf{y})} \sum_k \int q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k) q_d(\boldsymbol{\theta}_k|\mathbf{y}) \ln \frac{q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k) q_d(\boldsymbol{\theta}_k|\mathbf{y})}{p(\mathbf{x}, \boldsymbol{\theta}_k|\mathbf{y})} d\mathbf{x} \end{aligned}$$

We give the name **integrated non-factorized variational Bayes (INF-VB)** to this method.

Variational Optimization Algorithm

- ▶ **Step 1 (Local):** For each $\theta_k \in \mathcal{G}$, independently solving,

$$q^*(\mathbf{x}|\mathbf{y}, \theta_k) = \underset{q(\mathbf{x}|\mathbf{y}, \theta_k)}{\operatorname{argmin}} \operatorname{KL}(q(\mathbf{x}|\mathbf{y}, \theta_k) || p(\mathbf{x}|\mathbf{y}, \theta_k)) \quad (7)$$

- ▶ **Step 2 (Global):** Given $\{q^*(\mathbf{x}|\mathbf{y}, \theta_k) : \theta_k \in \mathcal{G}\}$, one have

$$q_d^*(\theta_k|\mathbf{y}) \propto \exp \left(\int q^*(\mathbf{x}|\mathbf{y}, \theta_k) \ln \frac{p(\mathbf{x}, \theta_k|\mathbf{y})}{q^*(\mathbf{x}|\mathbf{y}, \theta_k)} d\mathbf{x} \right) \quad (8)$$

- ▶ INF-VB is **parallelizable**, with dominant computational load distributed on each grid point
- ▶ INF-VB requires no iteration between Step 1 and Step 2

Our approach unifies INLA under the variational framework.

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

1. Gaussian approximation

$$q_G(\mathbf{x}|\mathbf{y}, \theta_k) = \mathcal{N}(\mathbf{x}; \mathbf{x}^*(\theta_k), \mathbf{H}(\mathbf{x}^*(\theta_k))^{-1}), \quad \forall \theta_k \in \mathcal{G}$$

2. Hyperparameter learning

$$q_{LA}(\theta|\mathbf{y}) \propto \frac{p(\mathbf{x}, \mathbf{y}, \theta)}{q_G(\mathbf{x}|\mathbf{y}, \theta)} \Big|_{\mathbf{x}=\mathbf{x}^*(\theta)}$$

3. Marginal posterior of \mathbf{x}

$$q(\mathbf{x}|\mathbf{y}) = \sum_k q_G(\mathbf{x}|\mathbf{y}, \theta_k) q_{LA}(\theta_k|\mathbf{y}) \Delta_k$$

with area weights Δ_k .

Main idea: Discretizing the low-dimensional space $\boldsymbol{\theta}$ using a grid \mathcal{G}

1. Gaussian approximation

$$q_G(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k) = \mathcal{N}(\mathbf{x}; \mathbf{x}^*(\boldsymbol{\theta}_k), \mathbf{H}(\mathbf{x}^*(\boldsymbol{\theta}_k))^{-1}), \quad \forall \boldsymbol{\theta}_k \in \mathcal{G}$$

(INF-VB) Step 1: Variational Gaussian approximation

$$q_{VG}^*(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k) = \operatorname{argmin} \operatorname{KL}(q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k) || p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k)), \quad \forall \boldsymbol{\theta}_k \in \mathcal{G}$$

2. Hyperparameter learning

$$q_{LA}(\boldsymbol{\theta}|\mathbf{y}) \propto \frac{p(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta})}{q_G(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})} \Big|_{\mathbf{x}=\mathbf{x}^*(\boldsymbol{\theta})}$$

3. Marginal posterior of \mathbf{x}

$$q(\mathbf{x}|\mathbf{y}) = \sum_k q_G(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k) q_{LA}(\boldsymbol{\theta}_k|\mathbf{y}) \Delta_k$$

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

(INF-VB) Step 1: Variational Gaussian approximation

$$q_{VG}^*(\mathbf{x}|\mathbf{y}, \theta_k) = \operatorname{argmin} \operatorname{KL}(q(\mathbf{x}|\mathbf{y}, \theta_k) || p(\mathbf{x}|\mathbf{y}, \theta_k)), \quad \forall \theta_k \in \mathcal{G}$$

2. Hyperparameter learning

$$q_{LA}(\theta|\mathbf{y}) \propto \frac{p(\mathbf{x}, \mathbf{y}, \theta)}{q_G(\mathbf{x}|\mathbf{y}, \theta)} \Big|_{\mathbf{x}=\mathbf{x}^*(\theta)}$$

(INF-VB) Step 2:

$$q_d^*(\theta_k|\mathbf{y}) \propto \exp \left(\int q_{VG}^*(\mathbf{x}|\mathbf{y}, \theta_k) \ln \frac{p(\mathbf{x}, \theta_k|\mathbf{y})}{q_{VG}^*(\mathbf{x}|\mathbf{y}, \theta_k)} d\mathbf{x} \right)$$

3. Marginal posterior of \mathbf{x}

$$q(\mathbf{x}|\mathbf{y}) = \sum_k q_G(\mathbf{x}|\mathbf{y}, \theta_k) q_{LA}(\theta_k|\mathbf{y}) \Delta_k$$

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

(INF-VB) Step 1: Variational Gaussian approximation

$$q_{VG}^*(\mathbf{x}|\mathbf{y}, \theta_k) = \operatorname{argmin} \operatorname{KL}(q(\mathbf{x}|\mathbf{y}, \theta_k) || p(\mathbf{x}|\mathbf{y}, \theta_k)), \quad \forall \theta_k \in \mathcal{G}$$

(INF-VB) Step 2: Hyperparameter learning

$$q_d^*(\theta_k|\mathbf{y}) \propto \exp \left(\int q_{VG}^*(\mathbf{x}|\mathbf{y}, \theta_k) \ln \frac{p(\mathbf{x}, \theta_k|\mathbf{y})}{q_{VG}^*(\mathbf{x}|\mathbf{y}, \theta_k)} d\mathbf{x} \right)$$

3. Marginal posterior of \mathbf{x}

$$q(\mathbf{x}|\mathbf{y}) = \sum_k q_G(\mathbf{x}|\mathbf{y}, \theta_k) q_{LA}(\theta_k|\mathbf{y}) \Delta_k$$

(INF-VB) Step 3:

$$q(\mathbf{x}|\mathbf{y}) = \int q(\mathbf{x}|\mathbf{y}, \theta) q_d(\theta|\mathbf{y}) d\theta = \sum_k q_{VG}^*(\mathbf{x}|\mathbf{y}, \theta_k) q_d^*(\theta_k|\mathbf{y})$$

Main idea: Discretizing the low-dimensional space θ using a grid \mathcal{G}

(INF-VB) Step 1: Variational Gaussian approximation

$$q_{VG}^*(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k) = \operatorname{argmin} \operatorname{KL}(q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k) || p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k)), \quad \forall \boldsymbol{\theta}_k \in \mathcal{G}$$

(INF-VB) Step 2: Hyperparameter learning

$$q_d^*(\boldsymbol{\theta}_k|\mathbf{y}) \propto \exp \left(\int q_{VG}^*(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k) \ln \frac{p(\mathbf{x}, \boldsymbol{\theta}_k|\mathbf{y})}{q_{VG}^*(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k)} d\mathbf{x} \right)$$

(INF-VB) Step 3: Marginal posterior of \mathbf{x}

$$q(\mathbf{x}|\mathbf{y}) = \int q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) q_d(\boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta} = \sum_k q_{VG}^*(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k) q_d^*(\boldsymbol{\theta}_k|\mathbf{y})$$

Benefits:

- ▶ Applicable to more general scenarios
- ▶ Optimal variational distributions $q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k)$ and $q_d(\boldsymbol{\theta}|\mathbf{y})$
- ▶ Negative ELBO provides quantization of the accuracy

Limitations:

- ▶ The dimension of $\boldsymbol{\theta}$ has to be no more than 5 or 6

Application to Bayesian Lasso

1. **Non-differentiability** of the ℓ_1 norm
2. The Laplace approximation of INLA cannot be applied

Model:

$$\mathbf{y} = \Phi \mathbf{x} + \mathbf{e}, \quad \mathbf{e} \sim \mathcal{N}(\mathbf{e}; \mathbf{0}, \sigma^2 \mathbf{I}_n)$$

where $\mathbf{y} \in \mathbb{R}^n$, $\Phi \in \mathbb{R}^{n \times p}$, and $\mathbf{e} \in \mathbb{R}^n$. We assume

$$\begin{aligned} x_j | \sigma^2, \lambda^2 &\sim \frac{\lambda}{2\sqrt{\sigma^2}} \exp\left(-\frac{\lambda}{\sqrt{\sigma^2}} \|x_j\|_1\right) \\ \sigma^2 &\sim \text{InvGa}(\sigma^2; a, b) \\ \lambda^2 &\sim \text{Ga}(\lambda^2; r, s). \end{aligned}$$

Problem: Given \mathbf{y} and Φ , find posterior distributions for \mathbf{x} and $\theta = \{\lambda^2, \sigma^2\}$

⁴Park & Casella, 2008

Inference:

1. Data augmentation Gibbs sampler
2. Mean-Field VB
3. **INF-VB**

INF-VB for Bayesian Lasso

(1) $q^*(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_k)$: constrain $q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \mathbf{C}\mathbf{C}^T)$, then

$$\text{KL}(q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})||p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})) := g(\boldsymbol{\mu}, \mathbf{C}) \quad (9)$$

is concave in $(\boldsymbol{\mu}, \mathbf{C})^a$, $\mathbf{D} = \mathbf{C}\mathbf{C}^T$.

- (2) $q^*(\boldsymbol{\theta}|\mathbf{y})$: can be evaluated analytically
- (3) $q^*(\mathbf{x}|\mathbf{y})$: finite mixture of Gaussians

^aChallis & Barber, 2011

Denote $(\boldsymbol{\mu}^*, \mathbf{D}^*) = \operatorname{argmin}_{\boldsymbol{\mu}, \mathbf{D}} g(\boldsymbol{\mu}, \mathbf{D})$, the **variational Bayesian Lasso**,

$$\boldsymbol{\mu}^* = \operatorname{argmin}_{\boldsymbol{\mu}} g(\boldsymbol{\mu}), \quad g(\boldsymbol{\mu}) := \mathbb{E}_{\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \mathbf{D}^*)} (\|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + 2\lambda\sigma \|\mathbf{x}\|_1) \quad (10)$$

is a counterpart of **Lasso**⁵,

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}), \quad f(\mathbf{x}) = \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + 2\lambda\sigma \|\mathbf{x}\|_1 \quad (11)$$

Remarks:

- ▶ The conditions of Lasso **hold on average**
- ▶ Smoothing around origin and thus **differentiable**
- ▶ Optimize a non-differential function by operating on a sequence of differentiable functions

⁵Tibshirani, 1996

Results (1/4): Diabetes Dataset⁶

This benchmark dataset contains

- ▶ Measurements on $n = 442$ diabetes patients
- ▶ $p = 10$ clinical covariates (age, sex, body mass index, average blood pressure, and six blood serum measurements)
- ▶ Response variable, a quantitative measure of disease progression

Goal:

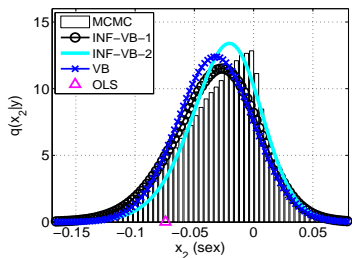
- ▶ Identify which covariates are important factors

Methods:

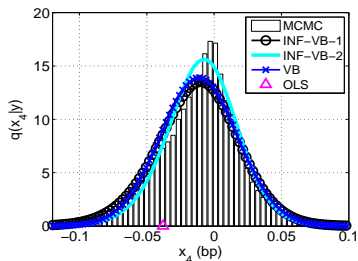
- ▶ **Intensive MCMC runs (ground truth)**
- ▶ Mean-Field VB
- ▶ **INF-VB-1**
- ▶ INF-VB-2 (INLA, replace LA with VG)
- ▶ Ordinary least square (OLS)

⁶Efron et al., 2004

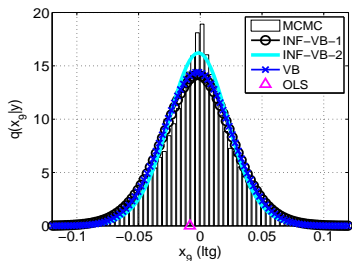
Results (2/4): Marginal Posteriors $q(x_j|\mathbf{y})$



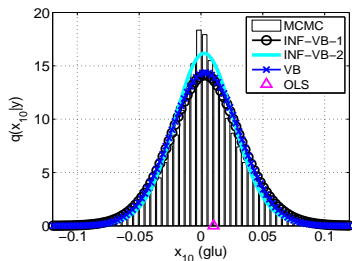
(a)



(b)

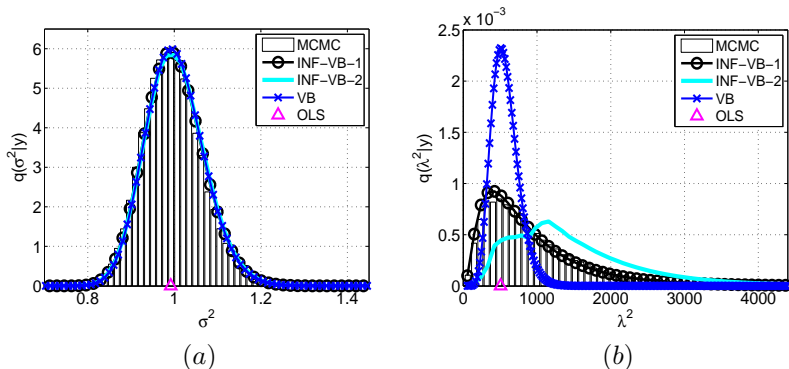


(c)



(d)

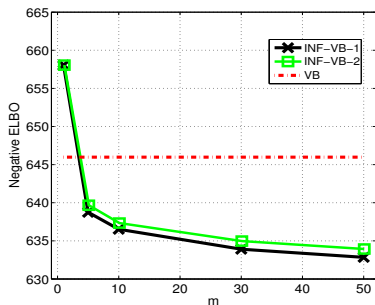
Results (3/4): Marginal Posteriors $q(\sigma^2|\mathbf{y})$ and $q(\lambda^2|\mathbf{y})$



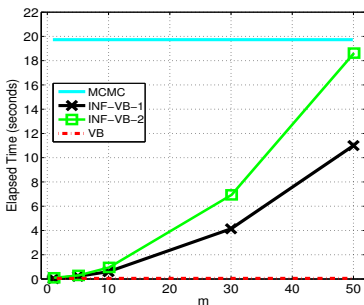
Posterior marginals of hyperparameters: (a) $q(\sigma^2|\mathbf{y})$ and (b) $q(\lambda^2|\mathbf{y})$

- ▶ Mean-Field VB could severely underestimate the posterior variance
- ▶ INF-VB-2 offers suboptimal solution

Results (4/4): Accuracy and Speed



(a) Accuracy



(b) Time

Grid size $m \times m$ and $m = 1, 5, 10, 30, 50$.

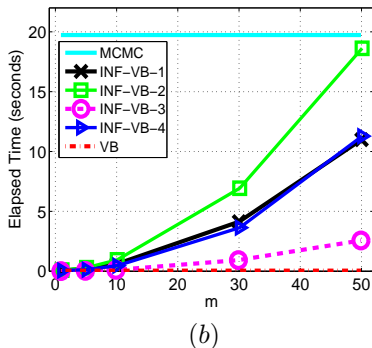
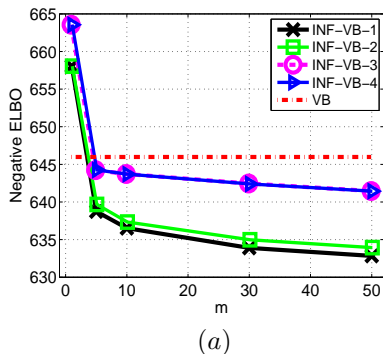
- ▶ INF-VB with a 1×1 grid: partial Bayesian learning of $q(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})$ with a fixed $\boldsymbol{\theta}$

Our method:

1. Tractable family Q : **non-factorized**
2. Conditional conjugacy: **not required**
3. Multimodal posterior: **could handle**
4. Parallelizable: **yes**

More could be done...

Q&A: Accuracy and Speed



In INF-VB-3 and INF-VB-4 (INLA, replace LA with VG), we obtain a fast VG solution by minimizing a KL divergence upper bound

Grid size $m \times m$ and $m = 1, 5, 10, 30, 50$.