

Variational Bayesian Methods Beyond Parametric and Continuous Assumptions

Mingyuan Zhou^{*§}
Joint work with Mingzhang Yin[§]

^{*}IROM Department, McCombs School of Business

[§]Department of Statistics and Data Sciences
The University of Texas at Austin

Workshop on Bayesian Nonparametrics for Signal and Image Processing
Bordeaux, July 3, 2018

Joint work with



Mingzhang Yin

PhD student (since Fall 2015)
in Statistics and Data Sciences

Bayesian Inference

- Bayes' rule:

$$P(\mathbf{z} | X) = \frac{P(X | \mathbf{z})P(\mathbf{z})}{P(X)} = \frac{P(X | \mathbf{z})P(\mathbf{z})}{\int P(X | \mathbf{z})P(\mathbf{z})d\mathbf{z}}$$

$$\text{Posterior of } \mathbf{z} \text{ given } X = \frac{\text{Conditional Likelihood} \times \text{Prior}}{\text{Marginal Likelihood}}$$

- Two main ways for approximate Bayesian inference:
 - Draw $\mathbf{z} \sim P(\mathbf{z} | X)$ using Markov chain Monte Carlo (MCMC) based methods such as **Gibbs sampling**: iteratively sample $P(z_k | X, \mathbf{z} \setminus z_k)$
 - Approximate the posterior $P(\mathbf{z} | X)$ with $Q(\mathbf{z})$, which is straightforward to sample from, using an optimization method such as Laplace approximation and **variational inference**

Inference via Gibbs sampling

- Gibbs sampling:
 - One of the simplest MCMC algorithm for multivariate distributions
 - Widely used for statistical inference
- For a multivariate distribution $P(z_1, \dots, z_K)$ that is difficult to sample from, if it is simpler to sample each of its variables conditioning on all the others, then we may use Gibbs sampling to obtain samples from this distribution as
 - Initialize $(z_1, \dots, z_K) = (z_1^0, \dots, z_K^0)$ at some values
 - For $s = 1 : S$
 - For $k = 1 : K$
 - Sample z_k^s conditioning on the others from
$$P(z_k^s | z_1^s, \dots, z_{k-1}^s, z_{k+1}^{s-1}, \dots, z_K^{s-1})$$
 - End
 - End
- Restriction of Gibbs sampling: conjugacy is often required

Variational inference

- With variational distribution $Q(\mathbf{z})$, we have

$$\begin{aligned}\ln P(X) &= \int Q(\mathbf{z}) \ln \frac{P(X, \mathbf{z})}{Q(\mathbf{z})} d\mathbf{z} + \int Q(\mathbf{z}) \ln \frac{Q(\mathbf{z})}{P(\mathbf{z} | X)} d\mathbf{z} \\ &= \mathcal{L}(Q) + \text{KL}(Q(\mathbf{z}) || P(\mathbf{z} | X)).\end{aligned}$$

- Since $\text{KL}(Q(\mathbf{z}) || P(\mathbf{z} | X)) \geq 0$, minimizing the Kullback-Leibler (KL) divergence from $P(\mathbf{z} | X)$ to $Q(\mathbf{z})$ is the same as maximizing the evidence lower bound:

$$\begin{aligned}\min_Q \text{KL}(Q(\mathbf{z}) || P(\mathbf{z} | X)) &\Leftrightarrow \max_Q \text{ELBO} \\ \text{ELBO} &= \mathcal{L}(Q) = \mathbb{E}_Q[\ln P(X, \mathbf{z})] - \mathbb{E}_Q[\ln Q(\mathbf{z})] \\ &= \mathbb{E}_Q[\ln P(X | \mathbf{z})] - \text{KL}(Q(\mathbf{z}) || P(\mathbf{z}))\end{aligned}$$

- Variational inference converts the problem of posterior inference into an optimization problem

Mean-field variational inference

- Mean-field variational inference (VI) factorizes the Q distribution of $\mathbf{z} = (z_1, \dots, z_K)^T$ as

$$Q(\mathbf{z}) = \prod_{i=1}^K q_{\phi_i}(z_i)$$

- The factorized assumption allows for closed-form coordinate ascent updates:

$$q^*(z_k) = \frac{\exp \left\{ \mathbb{E}_{q(\mathbf{z}_{-k})} [\log p(X, z_k, \mathbf{z}_{-k})] \right\}}{\int \exp \left\{ \mathbb{E}_{q(\mathbf{z}_{-k})} [\log p(X, z_k, \mathbf{z}_{-k})] \right\} dz_k}, \quad k = 1, \dots, K$$

where $\mathbf{z}_{-k} = \{z_1, \dots, z_{k-1}, z_{k+1}, \dots, z_K\}$.

- However, mean-field VI often clearly underestimates the variance of the posterior, due to the use of KL divergence and two restrictive constraints:
 - $q(z_k)$ are often restricted to the exponential family
 - The dependencies between z_k cannot be captured

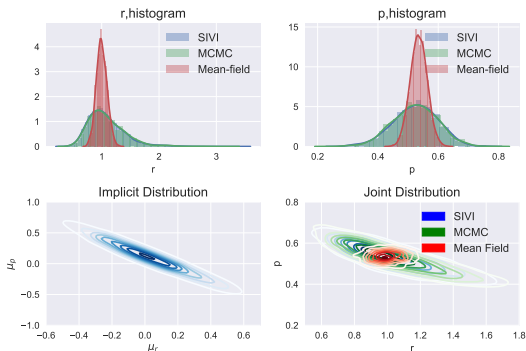
Model:

$$x_i \stackrel{i.i.d.}{\sim} \text{NB}(r, p), \quad r \sim \text{Gamma}(a, 1/b), \quad p \sim \text{Beta}(\alpha, \beta),$$

Mean-filed VI:

$$Q(r, p) = q(r)q(p) = \text{Gamma}(r; \tilde{a}, \tilde{b})\text{Beta}(p; \tilde{\alpha}, \tilde{\beta}),$$

Mean-filed VI underestimates variance (mainly due to the factorized assumption):

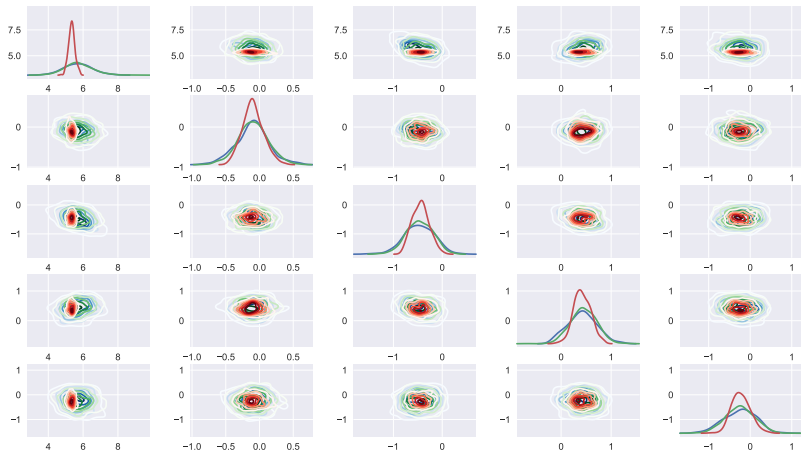


Bayesian logistic regression:

$$y_i \sim \text{Bernoulli}[(1 + e^{-x_i'\beta})^{-1}], \quad \beta \sim \mathcal{N}(\mathbf{0}, \alpha^{-1}\mathbf{I}_{V+1})$$

VI: $Q(\beta) = \mathcal{N}(\mu, \Sigma)$, which underestimates variance (mainly due to the mismatch between $\mathcal{N}(\mu, \Sigma)$ and the true posterior)

Blue: MCMC, Red: VI:



“Modern” variational inference

- Choose a more flexible $Q_\phi(\mathbf{z})$ and infer the variational parameter ϕ to maximize the ELBO via (stochastic) gradient ascent

$$\nabla_\phi \mathcal{L}(Q_\phi(\mathbf{z})) = \nabla_\phi \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z})} \left[\ln \frac{P(X, \mathbf{z})}{Q_\phi(\mathbf{z})} \right]$$

- Compute the gradient of the ELBO with respect to ϕ :
 - Score function gradient (a.k.a. REINFORCE, often suffering from high Monte Carlo estimation variance):

$$\nabla_\phi \mathcal{L}(Q_\phi(\mathbf{z})) = \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z})} \left[\ln \frac{P(X, \mathbf{z})}{Q_\phi(\mathbf{z})} \nabla_\phi \ln Q_\phi(\mathbf{z}) \right]$$

- If $\mathbf{z} \sim Q_\phi(\mathbf{z})$ is reparameterizable such that $\mathbf{z} = T_\phi(\epsilon)$, $\epsilon \sim q(\epsilon)$, then one can use the reparameterization trick:

$$\nabla_\phi \mathcal{L}(Q_\phi(\mathbf{z})) = \mathbb{E}_{\epsilon \sim q(\epsilon)} \left[\nabla_\phi \ln \frac{P(X, T_\phi(\epsilon))}{Q_\phi(T_\phi(\epsilon))} \right]$$

Challenges remain for “modern” variational inference

Parametric and continuous assumptions are commonly made, since

- There is a conflict between the ease of evaluating the log density ratio $\ln \frac{P(X, \mathbf{z})}{Q_\phi(\mathbf{z})}$ and the richness of $Q_\phi(\mathbf{z})$:
 - $\ln \frac{P(X, \mathbf{z})}{Q_\phi(\mathbf{z})}$ is straightforward to compute if $Q_\phi(\mathbf{z})$ is restricted to be analytic and point-wise evaluable
 - $\ln \frac{P(X, \mathbf{z})}{Q_\phi(\mathbf{z})}$ becomes difficult to compute if expanding the richness of the variational distribution family, e.g., allowing $Q_\phi(\mathbf{z})$ to be implicit
- One needs to control the Monte Carlo estimation variance:
 - The REINFORCE estimator often has large variance and needs to introduce appropriate control variates for variance reduction
 - The reparameterization trick often leads to low variance, but in general, it is not applicable to
 - Discrete distributions such as Bernoulli, categorical, and Poisson
 - Some commonly used continuous distributions such as gamma, beta, and Dirichlet

Relax parametric assumption with implicit distribution

Implicit distribution consists of a source of randomness $q(\epsilon)$ and a deterministic transform $T_\phi : \mathbb{R}^p \rightarrow \mathbb{R}^d$

$$\mathbf{z} = T_\phi(\epsilon), \epsilon \sim q(\epsilon)$$

- When T_ϕ is invertible and the dimension is low, the density

$$q_\phi(\mathbf{z}) = \frac{\partial}{\partial z_1} \cdots \frac{\partial}{\partial z_d} \int_{T_\phi(\epsilon) \leq \mathbf{z}} q(\epsilon) d\epsilon$$

can be calculated using change of variables

- But in general $\{T_\phi(\epsilon) \leq \mathbf{z}\}$ cannot be calculated and hence the high dimension integral is intractable, making $q_\phi(\mathbf{z})$ become implicit
- Even if $q_\phi(\mathbf{z})$ is implicit and hence difficult to evaluate, sampling $\mathbf{z} \sim q_\phi(\mathbf{z})$ is straightforward
- T_ϕ often corresponds to a **Deep Neural Network**

Why Deep Learning?

- TensorFlow, PyTorch, CNTK, Theano, . . . :
 - Automatic differentiation
 - Neural network libraries
 - Off-the-shelf optimization packages
- Exciting opportunities to be combined with Bayesian methods
 - Incorporating deep neural networks into hierarchical Bayesian models to define deep generative models
 - Empower variational inference with deep neural networks
 - Move beyond parametric and continuous assumptions
 - Define implicit distributions with MCMC
 - Build deep neural networks with many stochastic hidden layers
 - Design MCMC transition kernels with deep neural networks
 - Gradient backpropagation for discrete latent variables
 - . . .

Hierarchical variational family

The key to accurately capture the uncertainty is to capture the latent variable dependencies and expand the richness of the variational family

- One way is to add a hierarchical structure that assumes z_k to be conditional independent but marginally dependent, using

$$q(\mathbf{z} | \boldsymbol{\psi}) = \prod_{k=1}^K q(z_k | \psi_k), \quad \boldsymbol{\psi} \sim q_\phi(\boldsymbol{\psi})$$

- Marginalizing $\boldsymbol{\psi}$ out, we can view \mathbf{z} as a variable drawn from the distribution family \mathcal{H}

$$\mathcal{H} = \left\{ h_\phi(\mathbf{z}) : h_\phi(\mathbf{z}) = \mathbb{E}_{\boldsymbol{\psi} \sim q_\phi(\boldsymbol{\psi})} [q(\mathbf{z} | \boldsymbol{\psi})] = \int_{\boldsymbol{\psi}} \left[\prod_{k=1}^K q(z_k | \psi_k) \right] q_\phi(\boldsymbol{\psi}) d\boldsymbol{\psi} \right\}$$

- It is evident that $q(\mathbf{z} | \boldsymbol{\psi}) \in \mathcal{Q} \subseteq \mathcal{H}$, i.e., \mathcal{H} expands the original variational distribution family

Semi-implicit variational inference (SIVI)

- SIVI chooses $h_\phi(\mathbf{z}) = \mathbb{E}_{q_\phi(\psi)} q(\mathbf{z} | \psi)$ as its variational distribution
- Optimize ELBO = $\mathbb{E}_{h_\phi(\mathbf{z})} [\ln p(\mathbf{x}, \mathbf{z}) - \ln h_\phi(\mathbf{z})]$ for SIVI is generally intractable if $h_\phi(\mathbf{z}) = \mathbb{E}_{q_\phi(\psi)} q(\mathbf{z} | \psi)$ is not analytic
- KL convexity and Jensen's inequality lead to an ELBO lower bound:

$$\begin{aligned}\underline{\mathcal{L}}(q(\mathbf{z} | \psi), q_\phi(\psi)) &= \mathbb{E}_{\psi \sim q_\phi(\psi)} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi)} \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z} | \psi)} \\ &= - \mathbb{E}_{\psi \sim q_\phi(\psi)} \text{KL}(q(\mathbf{z} | \psi) || p(\mathbf{z} | \mathbf{x})) + \log p(\mathbf{x}) \\ &\leq - \text{KL}(\mathbb{E}_{\psi \sim q_\phi(\psi)} q(\mathbf{z} | \psi) || p(\mathbf{z} | \mathbf{x})) + \log p(\mathbf{x}) \\ &= \mathcal{L} = \mathbb{E}_{\mathbf{z} \sim h_\phi(\mathbf{z})} \log \frac{p(\mathbf{x}, \mathbf{z})}{h_\phi(\mathbf{z})}\end{aligned}$$

- Using the concavity of the logarithmic function, we have $\log h_\phi(\mathbf{z}) \geq \mathbb{E}_{\psi \sim q_\phi(\psi)} \log q(\mathbf{z} | \psi)$ and hence an ELBO upper bound:

$$\bar{\mathcal{L}}(q(\mathbf{z} | \psi), q_\phi(\psi)) = \mathbb{E}_{\psi \sim q_\phi(\psi)} \mathbb{E}_{\mathbf{z} \sim h_\phi(\mathbf{z})} \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z} | \psi)} \geq \mathcal{L}$$

- Note there is a subtle but critical difference between $\underline{\mathcal{L}}$ and $\bar{\mathcal{L}}$

To compute $\nabla_{\phi} \underline{\mathcal{L}}$, we require

- $q(\mathbf{z} | \psi)$ is reparameterizable and has an explicit probability density function
- $q_{\phi}(\psi)$ is reparameterizable and easy to sample from (does not have to be explicit), e.g., $q_{\phi}(\psi)$ can be constructed by transforming random noise ϵ via a deep neural network parameterized by ϕ

Maximizing the surrogate lower bound $\underline{\mathcal{L}}$ may lead to degeneracy that $q_{\phi}(\psi)$ converges to a point mass density:

Proposition (Degeneracy)

Let us denote $\psi^* = \arg \max_{\psi} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi)} \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z} | \psi)}$, then

$$\underline{\mathcal{L}}(q(\mathbf{z} | \psi), q_{\phi}(\psi)) \leq \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi^*)} \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z} | \psi^*)},$$

where the equality is true if and only if $q_{\phi}(\psi) = \delta_{\psi^*}(\psi)$.

Asymptotically exact ELBO (regularizing the lower bound)

- Avoid degeneracy by adding regularization $\underline{\mathcal{L}}_K = \underline{\mathcal{L}} + B_K$

$$B_K = \mathbb{E}_{\psi, \psi^{(1)}, \dots, \psi^{(K)} \sim q_\phi(\psi)} \text{KL}(q(\mathbf{z} | \psi) || \tilde{h}_K(\mathbf{z})), \quad (1)$$

where $\tilde{h}_K(\mathbf{z}) = \frac{1}{K+1} [q(\mathbf{z} | \psi) + \sum_{k=1}^K q(\mathbf{z} | \psi^{(k)})]$, $B_K \geq 0$, with $B_K = 0$ if and only if $K = 0$ or $q_\phi(\psi)$ degenerates to a point mass density

- The regularized surrogate ELBO can also be expressed as

$$\underline{\mathcal{L}}_K = \mathbb{E}_{\psi \sim q_\phi(\psi)} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi)} \mathbb{E}_{\psi^{(1)}, \dots, \psi^{(K)} \sim q_\phi(\psi)} \log \frac{p(\mathbf{x}, \mathbf{z})}{\frac{1}{K+1} [q(\mathbf{z} | \psi) + \sum_{k=1}^K q(\mathbf{z} | \psi^{(k)})]}$$

Proposition

The regularized lower bound $\underline{\mathcal{L}}_K = \underline{\mathcal{L}} + B_K$ is an asymptotically exact ELBO that satisfies $\underline{\mathcal{L}}_0 = \underline{\mathcal{L}}$ and $\lim_{K \rightarrow \infty} \underline{\mathcal{L}}_K = \underline{\mathcal{L}}$

Asymptotically exact ELBO (correcting the upper bound)

- Avoid divergence by adding correction $\bar{\mathcal{L}}_K = \bar{\mathcal{L}} - A_K$

$$A_K = \mathbb{E}_{\psi \sim q_\phi(\psi)} \mathbb{E}_{\mathbf{z} \sim h_\phi(\mathbf{z})} \mathbb{E}_{\psi^{(1)}, \dots, \psi^{(K)} \sim q_\phi(\psi)} \left[\log \left(\frac{1}{K} \sum_{k=1}^K q(\mathbf{z} | \psi^{(k)}) \right) - \log q(\mathbf{z} | \psi) \right].$$

- The corrected upper bound can be expressed as

$$\bar{\mathcal{L}}_K = \mathbb{E}_{\psi \sim q_\phi(\psi)} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi)} \mathbb{E}_{\psi^{(1)}, \dots, \psi^{(K)} \sim q_\phi(\psi)} \log \frac{p(\mathbf{x}, \mathbf{z})}{\frac{1}{K} \sum_{k=1}^K q(\mathbf{z} | \psi^{(k)})}$$

Proposition

The corrected upper bound $\bar{\mathcal{L}}_K = \bar{\mathcal{L}} - A_K$ monotonically converges from the above towards the ELBO, satisfying $\bar{\mathcal{L}}_1 = \bar{\mathcal{L}}$, $\bar{\mathcal{L}}_{K+1} \leq \bar{\mathcal{L}}_K$, and $\lim_{K \rightarrow \infty} \bar{\mathcal{L}}_K = \mathcal{L}$.

Algorithm for SIVI

Algorithm 1 Semi-Implicit Variational Inference (SIVI)

input : Data $\{x_i\}_{1:N}$, joint likelihood $p(\mathbf{x}, \mathbf{z})$, explicit variational distribution $q_\xi(\mathbf{z} | \psi)$ with reparameterization $\mathbf{z} = f(\epsilon, \xi, \psi)$, $\epsilon \sim p(\epsilon)$, implicit layer neural network $T_\phi(\epsilon)$ and source of randomness $q(\epsilon)$

output: Variational parameter ξ for the conditional distribution $q_\xi(\mathbf{z} | \psi)$, variational parameter ϕ for the mixing distribution $q_\phi(\psi)$

Initialize ξ and ϕ randomly

while *not converged* **do**

 Set $\underline{L}_{K_t} = 0$, ρ_t and η_t as step sizes, and $K_t \geq 0$ as a non-decreasing integer; Sample $\psi^{(k)} = T_\phi(\epsilon^{(k)})$, $\epsilon^{(k)} \sim q(\epsilon)$ for $k = 1, \dots, K_t$; take subsample $\mathbf{x} = \{x_i\}_{i_1:i_M}$

for $j = 1$ **to** J **do**

 Sample $\psi_j = T_\phi(\epsilon_j)$, $\epsilon_j \sim q(\epsilon)$

 Sample $\mathbf{z}_j = f(\bar{\epsilon}_j, \xi, \psi_j)$, $\bar{\epsilon}_j \sim p(\epsilon)$

$\underline{L}_{K_t} = \underline{L}_{K_t} + \frac{1}{J} \{ -\log \frac{1}{K_t+1} [\sum_{k=1}^{K_t} q_\xi(\mathbf{z}_j | \psi^{(k)}) + q_\xi(\mathbf{z}_j | \psi_j)] + \frac{N}{M} \log p(\mathbf{x} | \mathbf{z}_j) + \log p(\mathbf{z}_j) \}$

end

$t = t + 1$

$\xi = \xi + \rho_t \nabla_\xi \underline{L}_{K_t}(\{\psi^{(k)}\}_{1,K_t}, \{\psi_j\}_{1,J}, \{\mathbf{z}_j\}_{1,J})$

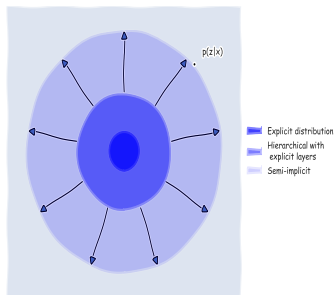
$\phi = \phi + \eta_t \nabla_\phi \underline{L}_{K_t}(\{\psi^{(k)}\}_{1,K_t}, \{\psi_j\}_{1,J}, \{\mathbf{z}_j\}_{1,J})$

end

Methods to expand variational distribution family

Expand variational family via stochastic hierarchy and/or deterministic nonlinear transform

- Hierarchy with explicit layers:
Negative Binomial \Leftrightarrow Poisson-Gamma hierarchy
- Normalizing Flow: transfer simple distribution with a chain of simple invertible mapping $\mathbf{z}_t = f_t \circ \dots \circ f_0(\mathbf{z}_0)$
- Modeling the dependencies between univariate marginals with copula
- Implicit distribution $\mathbf{z} = f(\epsilon)$, where f is not invertible
- Our approach: hierarchy with explicit conditional layer, implicit mixing layers (semi-implicit)



Hierarchical models in VI:

- Hierarchical variational models (Ranganath et al., 2016)
- Auxiliary deep generative models (Maaløe et al., 2016)
- Hierarchical implicit models and likelihood-free variational inference (Tran et al., 2017)

Implicit models in VI:

- Learning in implicit generative models (Mohamed and Lakshminarayanan, 2016)
- Variational inference using implicit distributions (Huszár, 2017)
- Implicit variational inference with kernel density ratio fitting (Shi et al., 2017)

Auxiliary deep generative models:

- Generative model (decoder): $\mathbf{x} | \mathbf{z} \sim p_{\theta}(\mathbf{x} | \mathbf{z}), \mathbf{z} \sim p(\mathbf{z})$
- Auxiliary model in the decoder: $\psi | \mathbf{z}, \mathbf{x} \sim p_{\theta}(\psi | \mathbf{z}, \mathbf{x})$
- Inference model (encoder): $q(\mathbf{z}, \psi | \mathbf{x}) = q_{\phi}(\mathbf{z} | \psi, \mathbf{x})q_{\phi}(\psi | \mathbf{x})$
- $\text{ELBO}_{\text{Auxiliary}} = \mathbb{E}_{\psi \sim q_{\phi}(\psi | \mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \psi, \mathbf{x})} \log \frac{p_{\theta}(\psi | \mathbf{z}, \mathbf{x})p_{\theta}(\mathbf{x} | \mathbf{z})p(\mathbf{z})}{q_{\phi}(\mathbf{z} | \psi, \mathbf{x})q_{\phi}(\psi | \mathbf{x})}$
- This lower bound is also used by the hierarchical variational models of Ranganath et al. (2016)

Semi-implicit variational inference (SIVI):

- Generative model (decoder): $\mathbf{x} | \mathbf{z} \sim p_{\theta}(\mathbf{x} | \mathbf{z}), \mathbf{z} \sim p(\mathbf{z})$
- Auxiliary model in the decoder: N/A
- Inference model (encoder):
$$h_{\phi}(\mathbf{z} | \mathbf{x}) = \int q_{\phi}(\mathbf{z} | \psi, \mathbf{x})q_{\phi}(\psi | \mathbf{x})d\psi = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\psi | \mathbf{x})}[q_{\phi}(\mathbf{z} | \psi, \mathbf{x})]$$
- $\text{ELBO}_{\text{SIVI}} = \mathbb{E}_{\psi \sim q_{\phi}(\psi | \mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \psi, \mathbf{x})} \log \frac{p_{\theta}(\mathbf{x} | \mathbf{z})p(\mathbf{z})}{\int q_{\phi}(\mathbf{z} | \psi, \mathbf{x})q_{\phi}(\psi | \mathbf{x})d\psi}$

- $\text{ELBO}_{\text{Auxiliary}} = \mathbb{E}_{\psi \sim q_{\phi}(\psi | \mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \psi, \mathbf{x})} \log \frac{p_{\theta}(\psi | \mathbf{z}, \mathbf{x}) p_{\theta}(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_{\phi}(\mathbf{z} | \psi, \mathbf{x}) q_{\phi}(\psi | \mathbf{x})}$
- $\text{ELBO}_{\text{SIVI}} = \mathbb{E}_{\psi \sim q_{\phi}(\psi | \mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \psi, \mathbf{x})} \log \frac{p_{\theta}(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{\int q_{\phi}(\mathbf{z} | \psi, \mathbf{x}) q_{\phi}(\psi | \mathbf{x}) d\psi}$
- Key differences:
 - SIVI has a tighter ELBO: $\log p(\mathbf{x}) \geq \text{ELBO}_{\text{SIVI}} \geq \text{ELBO}_{\text{Auxiliary}}$

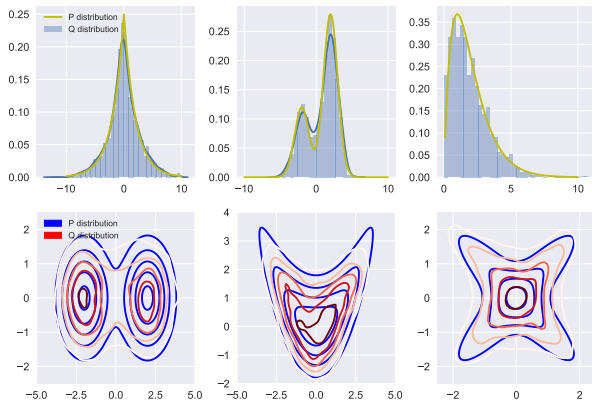
$$\text{ELBO}_{\text{SIVI}} - \text{ELBO}_{\text{Auxiliary}}$$

$$\begin{aligned}
 &= \mathbb{E}_{\psi \sim q_{\phi}(\psi | \mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \psi, \mathbf{x})} \log \frac{q_{\phi}(\mathbf{z} | \psi, \mathbf{x}) q_{\phi}(\psi | \mathbf{x})}{\int q_{\phi}(\mathbf{z} | \psi, \mathbf{x}) q_{\phi}(\psi | \mathbf{x}) d\psi} \\
 &\quad \log \frac{p_{\theta}(\psi | \mathbf{z}, \mathbf{x})}{p_{\theta}(\psi | \mathbf{z}, \mathbf{x})} \\
 &= \mathbb{E}_{\mathbf{z} \sim h_{\phi}(\mathbf{z} | \mathbf{x})} \mathbb{E}_{\psi \sim q_{\phi}(\psi | \mathbf{z}, \mathbf{x})} \log \frac{q_{\phi}(\psi | \mathbf{z}, \mathbf{x})}{p_{\theta}(\psi | \mathbf{z}, \mathbf{x})} \\
 &= \mathbb{E}_{\mathbf{z} \sim h_{\phi}(\mathbf{z} | \mathbf{x})} \text{KL}(q_{\phi}(\psi | \mathbf{z}, \mathbf{x}) || p_{\theta}(\psi | \mathbf{z}, \mathbf{x})) \\
 &\geq 0
 \end{aligned}$$

- SIVI allows $q_{\phi}(\psi | \mathbf{x})$ to be implicit
- SIVI sandwiches $\text{ELBO}_{\text{SIVI}}$ between a lower bound and an upper bound, and uses an asymptotically exact surrogate ELBO for optimization

Expressiveness of SIVI

$h(\mathbf{z}) = \mathbb{E}_{\psi \sim q(\psi)} q(\mathbf{z} \psi)$	$p(\mathbf{z})$
$z \sim \mathcal{N}(\psi, 0.1),$ $\psi \sim q(\psi)$	Laplace($z; \mu = 0, b = 2$) $0.3\mathcal{N}(z; -2, 1) + 0.7\mathcal{N}(z; 2, 1)$
$z \sim \text{Log-Normal}(\psi, 0.1),$ $\psi \sim q(\psi)$	Gamma($z; 2, 1$)
$z \sim \mathcal{N}\left(\psi, \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix}\right),$ $\psi \sim q(\psi)$	$0.5\mathcal{N}(z; -2, I) + 0.5\mathcal{N}(z; 2, I)$ $\mathcal{N}(z_1; z_1^2/4, 1)\mathcal{N}(z_2; 0, 4)$ $0.5\mathcal{N}\left(z; 0, \begin{bmatrix} 2 & 1.8 \\ 1.8 & 2 \end{bmatrix}\right) + 0.5\mathcal{N}\left(z; 0, \begin{bmatrix} 2 & -1.8 \\ -1.8 & 2 \end{bmatrix}\right)$



Expressiveness of SIVI

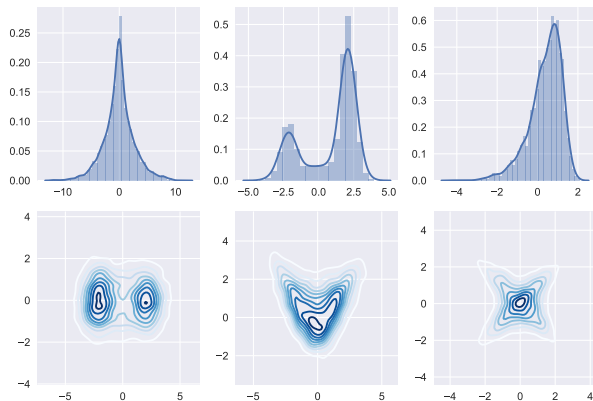


Figure: Visualization of the MLP based implicit distributions $\psi \sim q(\psi)$, which are mixed with isotropic Gaussian (or Log-Normal) distributions to approximate the target distributions.

Model:

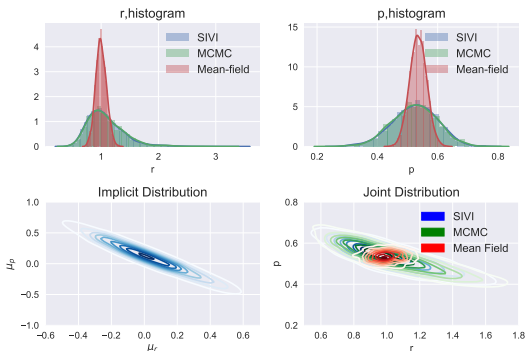
$$x_i \stackrel{i.i.d.}{\sim} \text{NB}(r, p), \quad r \sim \text{Gamma}(a, 1/b), \quad p \sim \text{Beta}(\alpha, \beta),$$

Mean-field VI:

$$Q(r, p) = q(r)q(p) = \text{Gamma}(r; \tilde{a}, \tilde{b})\text{Beta}(p; \tilde{\alpha}, \tilde{\beta}),$$

SIVI (both the conditional and mixing q distributions are reparameterizable) :

$$q(r, p | \psi) = \text{Log-Normal}(r; \mu_r, \sigma_0^2)\text{Logit-Normal}(p; \mu_p, \sigma_0^2),$$
$$\psi = (\mu_r, \mu_p) \sim q(\psi),$$



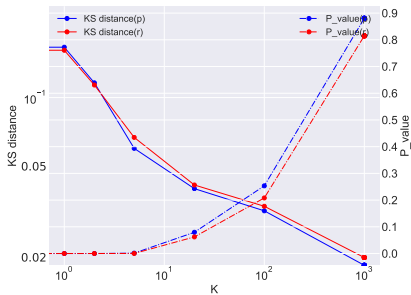
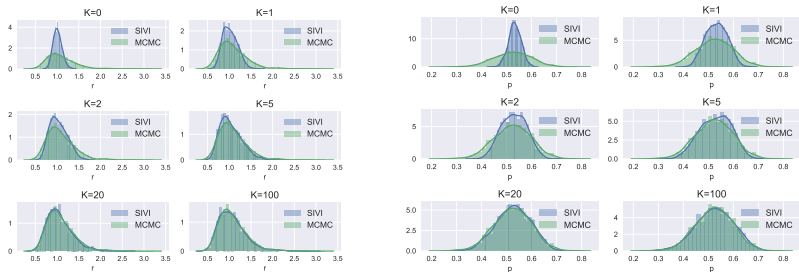


Figure: Kolmogorov-Smirnov (KS) distance and its corresponding p -value between the marginal posteriors of r and p inferred by SIVI and MCMC. SIVI rapidly improves as K increases.

Score function gradient for conjugate model

If $q(\mathbf{z} | \psi)$ is not reparameterizable, then we introduce a density ratio as

$$r_{\xi, \phi}(\mathbf{z}, \epsilon, \epsilon^{(1:K)}) = \frac{q_{\xi}(\mathbf{z} | T_{\phi}(\epsilon))}{\frac{1}{K+1}[q_{\xi}(\mathbf{z} | T_{\phi}(\epsilon)) + \sum_{k=1}^K q_{\xi}(\mathbf{z} | T_{\phi}(\epsilon^{(k)}))]}$$

and approximate the gradient of $\underline{\mathcal{L}}_K$ with respect to ϕ as

$$\begin{aligned} \nabla_{\phi} \underline{\mathcal{L}}_K \approx & \frac{1}{J} \sum_{j=1}^J \left\{ -\nabla_{\phi} \mathbb{E}_{\mathbf{z} \sim q_{\xi}(\mathbf{z} | T_{\phi}(\epsilon_j))} \log \frac{q_{\xi}(\mathbf{z} | T_{\phi}(\epsilon_j))}{p(\mathbf{x}, \mathbf{z})} \right. \\ & + \nabla_{\phi} \log r_{\xi, \phi}(\mathbf{z}_j, \epsilon_j, \epsilon^{(1:K)}) \\ & \left. + [\nabla_{\phi} \log q_{\xi}(\mathbf{z}_j | T_{\phi}(\epsilon_j))] \log r_{\xi, \phi}(\mathbf{z}_j, \epsilon_j, \epsilon^{(1:K)}) \right\}, \end{aligned}$$

- The first summation term is equivalent to the gradient of MFVI's ELBO
- Both the second and third terms correct the restrictions of $q_{\xi}(\mathbf{z} | T_{\phi}(\epsilon_j))$
- $\log r_{\xi, \phi}(\mathbf{z}, \epsilon, \epsilon^{(1:K)})$ in the third term is expected to be small regardless of convergence, effectively mitigating the variance of score function gradient estimation that is usually high in basic black-box VI

Model:

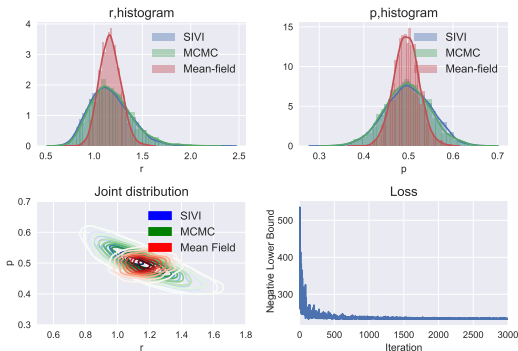
$$p(n_i, l_i | r, p) = r^{l_i} p^{n_i} (1-p)^r / Z_i, \quad r \sim \text{Gamma}(a, 1/b), \quad p \sim \text{Beta}(\alpha, \beta)$$

Mean-field VI:

$$Q(r, p) = q(r)q(p) = \text{Gamma}(r; \tilde{a}, \tilde{b})\text{Beta}(p; \tilde{\alpha}, \tilde{\beta}),$$

SIVI (non-reparameterizable conditional q distribution but conjugate model):

$$q(r, p | \psi) = \text{Gamma}(r; \psi_1, \psi_2)\text{Beta}(p; \psi_3, \psi_4), \quad \psi = (\psi_1, \psi_2, \psi_3, \psi_4) \sim q(\psi)$$

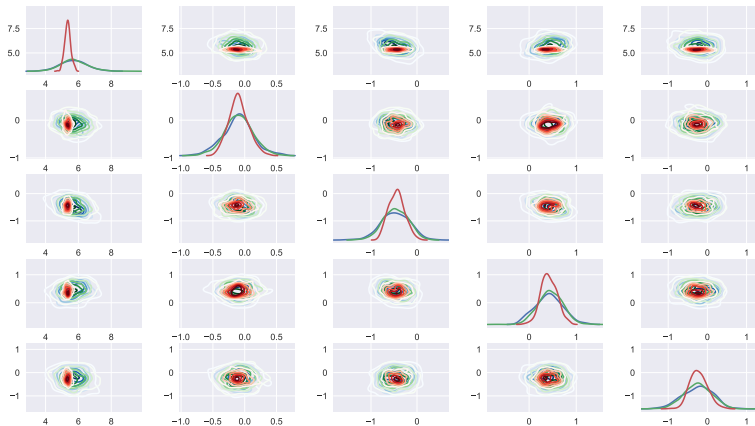


Bayesian logistic regression (pairwise joint distributions)

$$y_i \sim \text{Bernoulli}[(1 + e^{-x_i' \beta})^{-1}], \quad \beta \sim \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I}_{V+1})$$

SIVI: $q(\beta | \psi) = \mathcal{N}(\psi, \Sigma)$, $\psi \sim q_\phi(\psi)$

(Blue: MCMC, Red: VI, Green: SIVI):



Bayesian logistic regression (univariate marginals)

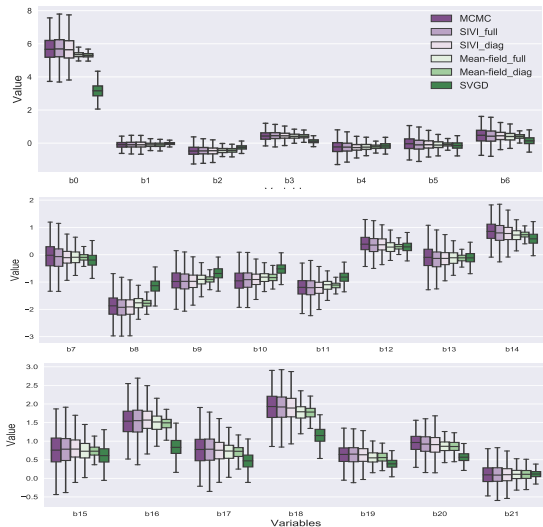


Figure: Comparison of all marginal posteriors of β_v inferred by various methods for Bayesian logistic regression on *waveform*.

Bayesian logistic regression (correlation coefficients)

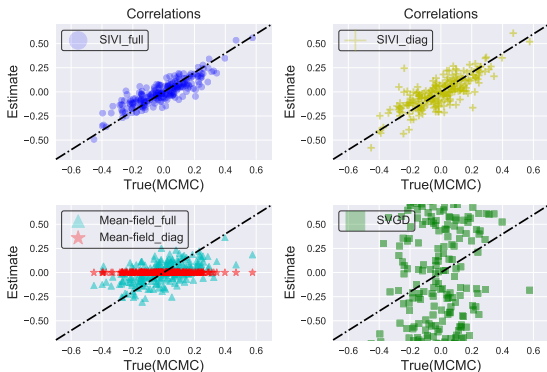


Figure: Against the correlation coefficients of β estimated from the posterior samples $\{\beta_i\}_{i=1:1000}$ of MCMC on *waveform*, top left/right plots the correlation coefficients of SIVI with a full/diagonal covariance matrix, bottom left plots these of MFVI with a full/diagonal covariance matrix, and bottom right plots these of SVGD. The closer to the dashed line the better.

Bayesian logistic regression (predictive uncertainty)

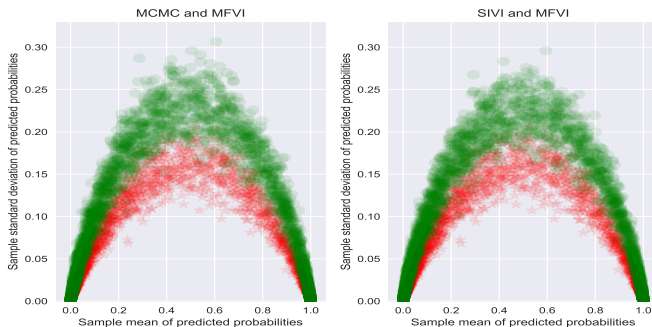


Figure: Comparison of MFVI (red) with a full covariance matrix, MCMC (green on left), and SIVI (green on right) with a full covariance matrix on quantifying predictive uncertainty for Bayesian logistic regression on *waveform*

Variational autoencoder

Variational Autoencoder (VAE) (Kingma and Welling, 2013; Rezende et al., 2014) is a popular generative model based approach for unsupervised feature learning and amortized inference.

- VAE iteratively infers the encoder parameter ϕ and decoder parameter θ to maximize the ELBO as

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \mathbf{x})} [\log(p_{\theta}(\mathbf{x} | \mathbf{z}))] - \text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) || p(\mathbf{z})).$$

- The encoder distribution $q_{\phi}(\mathbf{z} | \mathbf{x})$ is required to be reparameterizable and simple to compute its PDF, which usually restricts it to a small family of exponential distributions. A canonical form of the encoder is

$$q_{\phi}(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}(\mathbf{x}; \phi), \boldsymbol{\Sigma}(\mathbf{x}; \phi)),$$

where the Gaussian distribution parameters are deterministically transformed from the observed data \mathbf{x} , via non-probabilistic deep neural networks parameterized by ϕ .

- Thus, given observation \mathbf{x}_i , its corresponding code \mathbf{z}_i is forced to follow a Gaussian distribution, no matter how powerful the deep neural networks are. The Gaussian assumption, however, is often too restrictive to model skewed, heavy-tailed, and/or multi-modal distributions.

Semi-implicit variational autoencoder

We construct semi-implicit VAE (SIVAE) by using a hierarchical encoder that injects random noise at M different stochastic layers as

$$\begin{aligned} \ell_t &= T_t(\ell_{t-1}, \epsilon_t, \mathbf{x}; \phi), \quad \epsilon_t \sim q_t(\epsilon), \quad t = 1, \dots, M, \\ \boldsymbol{\mu}(\mathbf{x}, \phi) &= f(\ell_M, \mathbf{x}; \phi), \quad \boldsymbol{\Sigma}(\mathbf{x}, \phi) = g(\ell_M, \mathbf{x}; \phi), \\ q_\phi(\mathbf{z} | \mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}, \phi), \boldsymbol{\Sigma}(\mathbf{x}, \phi)), \end{aligned}$$

where $\ell_0 = \emptyset$ and T_t , f , and g are all deterministic neural networks. Note given data \mathbf{x}_i , $\boldsymbol{\mu}(\mathbf{x}_i, \phi)$, $\boldsymbol{\Sigma}(\mathbf{x}_i, \phi)$ are now random variables rather than following vanilla VAE to assume deterministic values. This clearly moves the encoder variational distribution beyond a simple Gaussian form.

Semi-implicit variational autoencoder

Methods	$-\log p(\mathbf{x})$
<i>Results below form Burda et al. (2015)</i>	
VAE + IWAE	= 86.76
IWAE + IWAE	= 84.78
<i>Results below form Salimans et al. (2015)</i>	
DLGM + HVI (1 leapfrog step)	= 88.08
DLGM + HVI (4 leapfrog step)	= 86.40
DLGM + HVI (8 leapfrog steps)	= 85.51
<i>Results below form Rezende & Mohamed (2015)</i>	
DLGM+NICE (Dinh et al., 2014) (k = 80)	≤ 87.2
DLGM+NF (k = 40)	≤ 85.7
DLGM+NF (k = 80)	≤ 85.1
<i>Results below form Gregor et al. (2015)</i>	
DLGM	≈ 86.60
NADE	= 88.33
DBM 2hl	≈ 84.62
DBN 2hl	≈ 84.55
EoNADE-5 2hl (128 orderings)	= 84.68
DARN 1hl	≈ 84.13
<i>Results below form Maaløe et al. (2016)</i>	
Auxiliary VAE (L=1, IW=1)	≤ 84.59
<i>Results below form Mescheder et al. (2017)</i>	
VAE + IAF (Kingma et al., 2016)	$\approx 84.9 \pm 0.3$
Auxiliary VAE (Maaløe et al., 2016)	$\approx 83.8 \pm 0.3$
AVB + AC	$\approx 83.7 \pm 0.3$
SIVI (3 stochastic layers)	= 84.07
SIVI (3 stochastic layers)+ IW($\tilde{K} = 10$)	= 83.25

Summary for SIVI

- Uncertainty estimation is difficult but important in Variational Inference
- A key to achieve an accurate uncertainty estimation is to construct a flexible variational distribution that can capture the dependencies between latent variables
- Combining the advantages of having analytic point-wise evaluable density ratios and tractable computation via Monte Carlo estimation, semi-implicit variational inference (SIVI) can approach the accuracy of MCMC in quantifying posterior uncertainty, but often pays a lower computational cost and can generate independent posterior samples on the fly via the inferred stochastic variational inference network

- Variational sampling
 - MCMC is able to accurately capture posterior uncertainty, and still hard to be replaced by variational inference in many challenging statistical inference problems
 - For a latent variable whose prior is not conjugate to the likelihood, the corresponding MCMC transition kernel may be difficult to design
 - Our idea is to use variational inference + deep neural network to learn MCMC transition kernels
- Augment-REINFORCE-merge gradient for discrete latent variable models
 - Backpropagate unbiased and low-variance gradient for discrete latent variables