

Semi-Implicit Variational Inference

Mingzhang Yin[§]

Joint work with Mingyuan Zhou^{*§}

The University of Texas at Austin

[§]Department of Statistics and Data Sciences

^{*}IROM Department, McCombs School of Business

International Conference on Machine Learning
Stockholm, Sweden, July 11, 2018

- 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 $P(\mathbf{z}_k | X, \mathbf{z} \setminus \mathbf{z}_k)$
 - 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**

Variational inference

- Evidence and ELBO:

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

- 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

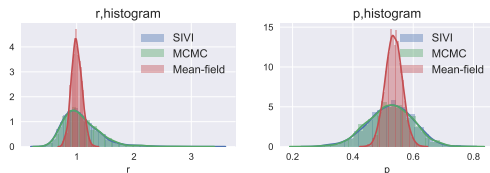
Model:

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

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):



“Modern” variational inference

Choose a more flexible $Q_\phi(\mathbf{z})$ and infer the variational parameter ϕ via (stochastic) gradient descent (by reparameterization or score method)

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

- There are two major flexibilities we want $Q_\phi(\mathbf{z})$ have:
 - We wish $Q_\phi(\mathbf{z})$ is not restricted to have an analytic density (but should be easy to sample)
 - We wish $Q_\phi(\mathbf{z})$ to incorporate dependencies of latent variables
- We also want to maintain computational tractability for a flexible inference distribution

To achieve the computation and accuracy balance, we use [the neural network implicit distribution](#) in [a hierarchical model](#).

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

- Direct inference with implicit distribution can be difficult because of the need to estimate the density ratio $\frac{P(X, \mathbf{z})}{Q_\phi(\mathbf{z})}$

Hierarchical variational family

Capturing the latent variable dependencies plays the key role to accurately estimate the uncertainty.

- 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} which we choose as variational family

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

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

Semi-implicit variational inference (SIVI)

- We call the hierarchical model semi-implicit because it requires $q(\mathbf{z} | \psi)$ to be explicit while allows $q_\phi(\psi)$ to be implicit, and $h_\phi(\mathbf{z}) = \mathbb{E}_{q_\phi(\psi)} q(\mathbf{z} | \psi)$ and ELBO is generally 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 $\bar{\mathcal{L}}$

Degeneracy of $\underline{\mathcal{L}}$

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

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

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

Asymptotically exact ELBO

- 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 Jensen gap can also be narrowed from upper side by $\bar{\mathcal{L}}_k = \bar{\mathcal{L}} - A_k$

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

The regularized lower bound $\underline{\mathcal{L}}_K$ is an asymptotically exact ELBO that satisfies $\underline{\mathcal{L}}_0 = \underline{\mathcal{L}}$ and $\lim_{K \rightarrow \infty} \underline{\mathcal{L}}_K = \mathcal{L}$. The regularized upper bound satisfies $\bar{\mathcal{L}}_1 = \bar{\mathcal{L}}$, $\bar{\mathcal{L}}_{K+1} \leq \bar{\mathcal{L}}_K$, 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

$t = t + 1$

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

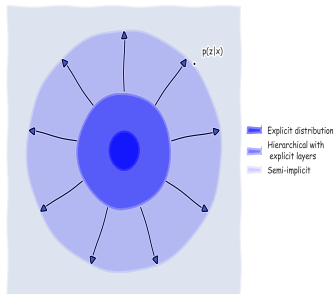
$\phi = \phi + \eta_t \nabla_\phi \underline{L}_{K_t}(\{\psi^{(k)}\}_{1,K_t}, \{\psi_j\}_{1,J}, \{\mathbf{z}_j\}_{1,J})$

end

Methods to expand variational distribution family

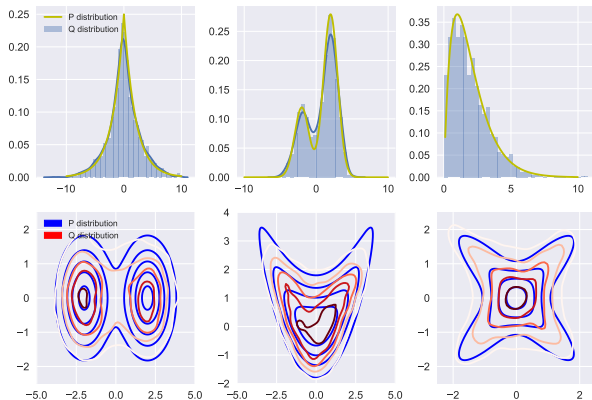
Expand variational family via stochastic and/or deterministic method

- Hierarchical models: eg. Negative Binomial \leftrightarrow Poisson-Gamma hierarchy; Hierarchical variational model (Ranganath et al., 2016)
- Normalizing Flow: transfer simple distribution with a chain of simple invertible mapping $\mathbf{z}_t = f_t \circ \dots \circ f_0(\mathbf{z}_0)$ (Rezende and Mohamed, 2015)
- Modeling the dependencies between univariate marginals with copula (Tran et al., 2015)
- Implicit distribution $\mathbf{z} = f(\epsilon)$, where f is not invertible; (Tran et al., 2017)
- Our approach: hierarchy with explicit conditional layer, implicit mixing layers (semi-implicit)



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($z; 2, 1$)
$z \sim \mathcal{N}\left(\psi, \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix}\right),$ $\psi \sim q(\psi)$	$0.5\mathcal{N}(z; -2, I) + 0.5\mathcal{N}(z; 2, I)$ $\mathcal{N}(z_1; z_1^2/4, 1)\mathcal{N}(z_2; 0, 4)$ $0.5\mathcal{N}\left(z; 0, \begin{bmatrix} 2 & 1.8 \\ 1.8 & 2 \end{bmatrix}\right) + 0.5\mathcal{N}\left(z; 0, \begin{bmatrix} 2 & -1.8 \\ -1.8 & 2 \end{bmatrix}\right)$



Model:

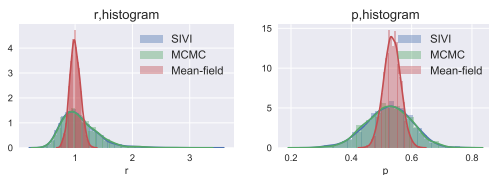
$$x_i \stackrel{i.i.d.}{\sim} \text{NB}(r, p), \quad r \sim \text{Gamma}(a, 1/b), \quad p \sim \text{Beta}(\alpha, \beta),$$

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) :

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



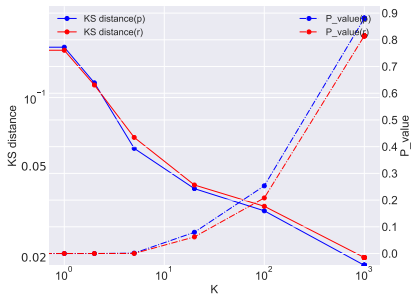
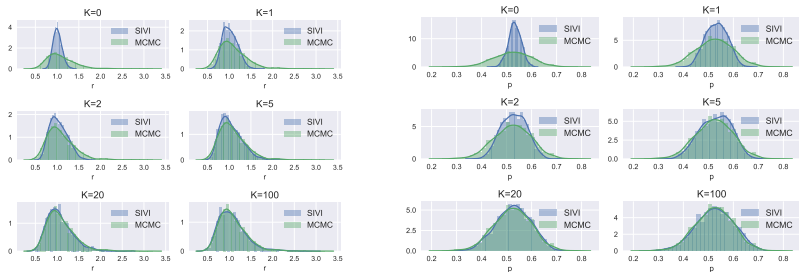


Figure: Kolmogorov-Smirnov (KS) distance and its corresponding p -value between the marginal posteriors of r and p inferred by SIVI and MCMC. SIVI rapidly improves as K increases.

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

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

- 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:

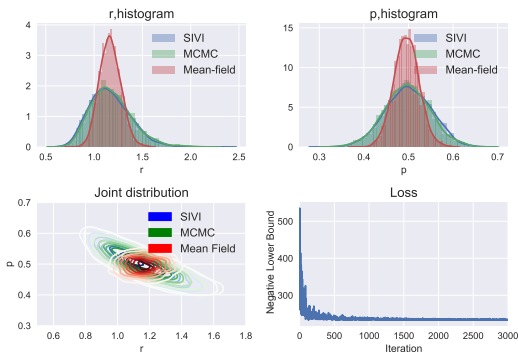
$$p(n_i, l_i | r, p) = r^{l_i} p^{n_i} (1-p)^r / Z_i, \quad r \sim \text{Gamma}(a, 1/b), \quad p \sim \text{Beta}(\alpha, \beta)$$

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):

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

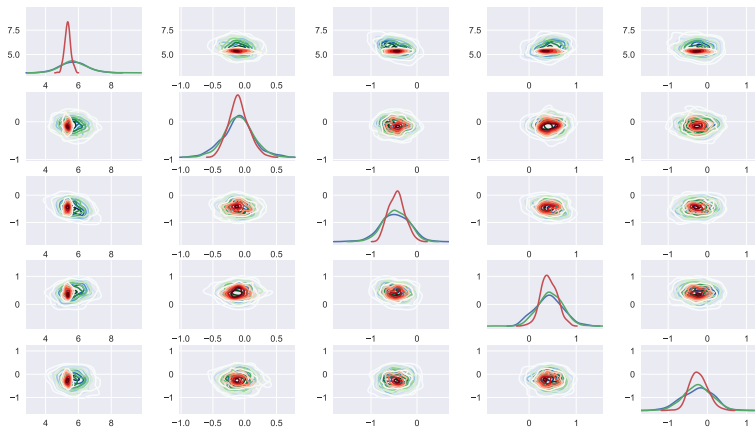


Bayesian logistic regression (pairwise joint distributions)

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

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

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



Bayesian logistic regression

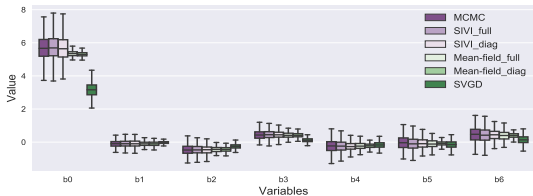


Figure: Comparing univariate marginals

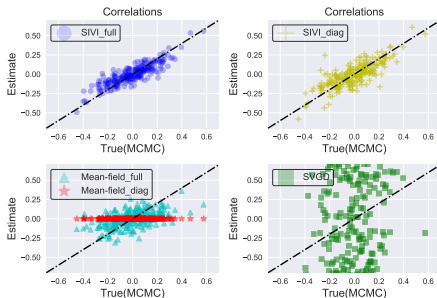


Figure: Comparing posterior covariance matrix

Bayesian logistic regression (predictive uncertainty)

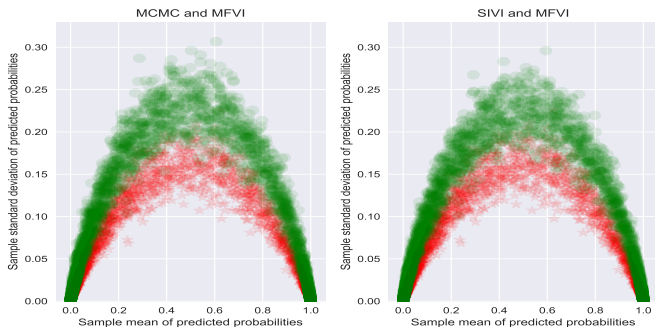


Figure: Comparison of MFVI (red) with a full covariance matrix, MCMC (green on left), and SIVI (green on right) with a full covariance matrix on quantifying predictive uncertainty for Bayesian logistic regression on *waveform*

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

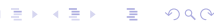
Methods	$-\log p(\mathbf{x})$
<i>Results below form Burda et al. (2015)</i>	
VAE + IWAE	= 86.76
IWAE + IWAE	= 84.78
<i>Results below form Salimans et al. (2015)</i>	
DLGM + HVI (1 leapfrog step)	= 88.08
DLGM + HVI (4 leapfrog step)	= 86.40
DLGM + HVI (8 leapfrog steps)	= 85.51
<i>Results below form Rezende & Mohamed (2015)</i>	
DLGM+NICE (Dinh et al., 2014) (k = 80)	\leq 87.2
DLGM+NF (k = 40)	\leq 85.7
DLGM+NF (k = 80)	\leq 85.1
<i>Results below form Gregor et al. (2015)</i>	
DLGM	\approx 86.60
NADE	= 88.33
DBM 2hl	\approx 84.62
DBN 2hl	\approx 84.55
EoNADE-5 2hl (128 orderings)	= 84.68
DARN 1hl	\approx 84.13
<i>Results below form Maaløe et al. (2016)</i>	
Auxiliary VAE (L=1, IW=1)	\leq 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

- Uncertainty estimation is difficult but important in Variational Inference
- One key to get an accurate uncertainty estimation is to construct a flexible variational distribution that can capture the dependencies between latent variables
- Balancing the expressiveness and tractability, 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 fast via the inferred stochastic variational inference network.

Thank you!

Welcome to our poster at Hall B # 177

¹Reproducible code is at <https://github.com/mingzhang-yin/SIVI> 

Related inference methods

- VAE: Changing the empirical data distribution leads to degenerated \mathcal{L}

$$\mathcal{L}_{VAE} = \mathbb{E}_{\mathbf{x} \sim D(\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})}$$

$$\underline{\mathcal{L}}_{SIVI} = \mathbb{E}_{\epsilon \sim q(\epsilon)} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}(\epsilon))} \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x}(\epsilon))}$$

- Data augmentation: iteratively sample from $p(\mathbf{z}|\psi)$ and $p(\psi|\mathbf{z})$ with

$$p(\mathbf{z}) = \int p(\mathbf{z}, \psi) d\psi$$

- Auxiliary Deep Generative Models (Maaløe et al., 2016): optimize on a less tighter bound

$$\log p(\mathbf{x}) \geq \mathbb{E}_{h_\phi(\mathbf{z}|\mathbf{x})} \log \frac{p(\mathbf{x}, \mathbf{z})}{h_\phi(\mathbf{z}|\mathbf{x})} \geq \mathbb{E}_{q_\phi(\mathbf{z}, \mathbf{a}|\mathbf{x})} \log \frac{p(\mathbf{x}, \mathbf{z}, \mathbf{a})}{q_\phi(\mathbf{z}, \mathbf{a}|\mathbf{x})}$$