

Overdispersed Black-Box Variational Inference

Francisco J. R. Ruiz, Michalis K. Titsias, David M. Blei

Columbia University
Athens University of Economics and Business

June 27th, 2016



Overdispersed Black-Box Variational Inference

- ▶ General variational inference for any probabilistic model
- ▶ Builds on black-box variational inference (BBVI)
- ▶ Reduces the variance of the estimator (\implies faster convergence)
- ▶ Requires a variational distribution in the exponential family
- ▶ Key idea: analyze the optimal *importance sampling proposal*

Notation

- ▶ Probabilistic model $p(\mathbf{x}, \mathbf{z})$
 - ▶ \mathbf{x} : Data
 - ▶ \mathbf{z} : Latent variables
- ▶ Assume the posterior $p(\mathbf{z} | \mathbf{x})$ is intractable:

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{\int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}}$$

- ▶ We wish to *approximate* the posterior using **variational inference**

Variational Inference

- ▶ Approximate the posterior with a simpler distribution $q(\mathbf{z}; \boldsymbol{\lambda})$
- ▶ Minimize the KL divergence w.r.t. $\boldsymbol{\lambda}$

$$\boldsymbol{\lambda}^* = \arg \min_{\boldsymbol{\lambda}} D_{\text{KL}}(q(\mathbf{z}; \boldsymbol{\lambda}) || p(\mathbf{z} | \mathbf{x}))$$

- ▶ Evidence lower bound (ELBO):

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\lambda})} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \boldsymbol{\lambda})]$$

Variational Inference

- ▶ Approximate the posterior with a simpler distribution $q(\mathbf{z}; \boldsymbol{\lambda})$
- ▶ Minimize the KL divergence w.r.t. $\boldsymbol{\lambda}$

$$\boldsymbol{\lambda}^* = \arg \min_{\boldsymbol{\lambda}} D_{\text{KL}}(q(\mathbf{z}; \boldsymbol{\lambda}) || p(\mathbf{z} | \mathbf{x}))$$

- ▶ Evidence lower bound (ELBO):

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\lambda})} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \boldsymbol{\lambda})]$$

- ▶ Optimization problem:
 - ▶ Conditionally conjugate models: *coordinate ascent*
 - ▶ Non-conjugate models: one recent approach is **BBVI**

Examples of Conditionally Non-Conjugate Models

- ▶ Time series models
- ▶ Probabilistic matrix factorization
- ▶ Deep probabilistic models
- ▶ Correlated topic models
- ▶ ...

Black-Box Variational Inference¹

- ▶ Stochastic optimization
- ▶ Builds Monte Carlo estimates of the gradient $\nabla_{\lambda}\mathcal{L}$
- ▶ Relies on the score function method:

$$\nabla_{\lambda}\mathcal{L} = \mathbb{E}_{q(\mathbf{z};\lambda)} [f(\mathbf{z})],$$

where

$$f(\mathbf{z}) \triangleq \nabla_{\lambda} \log q(\mathbf{z}; \lambda) (\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \lambda))$$

¹Ranganath et al. (2014)

Black-Box Variational Inference¹

- ▶ Stochastic optimization
- ▶ Builds Monte Carlo estimates of the gradient $\nabla_{\lambda}\mathcal{L}$
- ▶ Relies on the score function method:

$$\nabla_{\lambda}\mathcal{L} = \mathbb{E}_{q(\mathbf{z};\lambda)} [f(\mathbf{z})],$$

where

$$f(\mathbf{z}) \triangleq \nabla_{\lambda} \log q(\mathbf{z}; \lambda) (\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \lambda))$$

- ▶ Algorithm:
 1. Sample $\mathbf{z}^{(s)}$ iid from $q(\mathbf{z}; \lambda)$
 2. Evaluate $f(\mathbf{z}^{(s)})$ for each sample s
 3. Obtain a Monte Carlo estimate of the gradient
 4. Take a gradient step for λ

¹Ranganath et al. (2014)

Controlling the Variance

- ▶ The estimator of the gradient may suffer from **high variance**
- ▶ This leads to slow convergence
- ▶ Methods to reduce the variance:
 - ▶ Rao-Blackwellization²
 - ▶ Control variates³
 - ▶ Reparameterization trick⁴
 - ▶ Local expectations⁵

²Casella and Robert (1996); Ranganath et al. (2014)

³Ross (2002); Paisley et al. (2012); Ranganath et al. (2014); Gu et al. (2016)

⁴Price (1958); Bonnet (1964); Salimans and Knowles (2013); Kingma and Welling (2014); Rezende et al. (2014); Kucukelbir et al. (2015)

⁵Titsias and Lázaro-Gredilla (2015)

Controlling the Variance

- ▶ The estimator of the gradient may suffer from **high variance**
- ▶ This leads to slow convergence
- ▶ Methods to reduce the variance:
 - ▶ Rao-Blackwellization²
 - ▶ Control variates³
 - ▶ Reparameterization trick⁴
 - ▶ Local expectations⁵
- ▶ New method: Overdispersed BBVI

²Casella and Robert (1996); Ranganath et al. (2014)

³Ross (2002); Paisley et al. (2012); Ranganath et al. (2014); Gu et al. (2016)

⁴Price (1958); Bonnet (1964); Salimans and Knowles (2013); Kingma and Welling (2014); Rezende et al. (2014); Kucukelbir et al. (2015)

⁵Titsias and Lázaro-Gredilla (2015)

Overdispersed BBVI

- ▶ Builds on BBVI
- ▶ Samples from another distribution $r(\mathbf{z}) \neq q(\mathbf{z}; \lambda)$

$$\nabla_{\lambda} \mathcal{L} = \mathbb{E}_{q(\mathbf{z}; \lambda)} [f(\mathbf{z})] = \mathbb{E}_{r(\mathbf{z})} \left[f(\mathbf{z}) \frac{q(\mathbf{z}; \lambda)}{r(\mathbf{z})} \right]$$

- ▶ The optimal importance sampling proposal⁶ is

$$r_n^*(\mathbf{z}) \propto q(\mathbf{z}; \lambda) |f_n(\mathbf{z})|$$

- ▶ The optimal proposal is intractable
- ▶ O-BBVI searches for another proposal $r(\mathbf{z})$

⁶Robert and Casella (2005); Owen (2013)

Overdispersed BBVI

- ▶ Assume an exponential family variational distribution

$$q(\mathbf{z}; \boldsymbol{\lambda}) \propto \exp\{\boldsymbol{\lambda}^\top t(\mathbf{z}) - A(\boldsymbol{\lambda})\}$$

- ▶ Recall the optimal proposal:

$$r_n^*(\mathbf{z}) \propto q(\mathbf{z}; \boldsymbol{\lambda}) |f_n(\mathbf{z})|$$
$$f_n(\mathbf{z}) = \nabla_{\lambda_n} \log q(\mathbf{z}; \boldsymbol{\lambda}) (\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \boldsymbol{\lambda}))$$

Overdispersed BBVI

- ▶ Assume an exponential family variational distribution

$$q(\mathbf{z}; \boldsymbol{\lambda}) \propto \exp\{\boldsymbol{\lambda}^\top t(\mathbf{z}) - A(\boldsymbol{\lambda})\}$$

- ▶ Recall the optimal proposal:

$$r_n^*(\mathbf{z}) \propto q(\mathbf{z}; \boldsymbol{\lambda}) |f_n(\mathbf{z})|$$
$$f_n(\mathbf{z}) = \nabla_{\lambda_n} \log q(\mathbf{z}; \boldsymbol{\lambda}) (\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \boldsymbol{\lambda}))$$

- ▶ The optimal proposal assigns higher mass in the *tails* of $q(\mathbf{z}; \boldsymbol{\lambda})$

Overdispersed BBVI

- ▶ Assume an exponential family variational distribution

$$q(\mathbf{z}; \boldsymbol{\lambda}) \propto \exp\{\boldsymbol{\lambda}^\top t(\mathbf{z}) - A(\boldsymbol{\lambda})\}$$

- ▶ Recall the optimal proposal:

$$r_n^*(\mathbf{z}) \propto q(\mathbf{z}; \boldsymbol{\lambda}) |f_n(\mathbf{z})|$$
$$f_n(\mathbf{z}) = \nabla_{\lambda_n} \log q(\mathbf{z}; \boldsymbol{\lambda}) (\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \boldsymbol{\lambda}))$$

- ▶ The optimal proposal assigns higher mass in the *tails* of $q(\mathbf{z}; \boldsymbol{\lambda})$
- ▶ We use an **overdispersed distribution**⁷

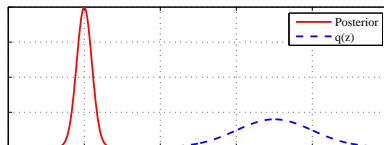
$$r(\mathbf{z}; \boldsymbol{\lambda}, \tau) \propto \exp\left\{\frac{\boldsymbol{\lambda}^\top t(\mathbf{z}) - A(\boldsymbol{\lambda})}{\tau}\right\}$$

⁷Jørgensen (1987)

Heavier Tails

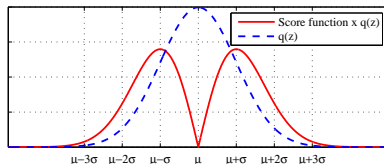
$$r_n^*(\mathbf{z}) \propto q(\mathbf{z}; \lambda) |f_n(\mathbf{z})|$$
$$f_n(\mathbf{z}) = \nabla_{\lambda_n} \log q(\mathbf{z}; \lambda) (\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \lambda))$$

- ▶ Through the model $p(\mathbf{x}, \mathbf{z})$



- ▶ Through the score function

$$\nabla_{\lambda_n} \log q(\mathbf{z}; \lambda) = t_n(\mathbf{z}) - \mathbb{E}_{q(\mathbf{z}; \lambda)} [t_n(\mathbf{z})]$$



Implementation

- ▶ Importance sampling fails in high dimensionality settings
 - We use local expectations⁸
 - A proposal distribution per latent variable

$$\begin{aligned}\nabla_{\lambda_n} \mathcal{L} &= \mathbb{E}_{q(\mathbf{z}_n; \lambda_n)} \left[\mathbb{E}_{q(\mathbf{z}_{-n}; \lambda_{-n})} [f_n(\mathbf{z})] \right] \\ &= \mathbb{E}_{r(\mathbf{z}_n; \lambda_n, \tau_n)} \left[\frac{q(\mathbf{z}_n; \lambda_n)}{r(\mathbf{z}_n; \lambda_n, \tau_n)} \mathbb{E}_{q(\mathbf{z}_{-n}; \lambda_{-n})} [f_n(\mathbf{z})] \right]\end{aligned}$$

⁸Titsias and Lázaro-Gredilla (2015)

Implementation

- ▶ Importance sampling fails in high dimensionality settings
 - We use local expectations⁸
 - A proposal distribution per latent variable

$$\begin{aligned}\nabla_{\lambda_n} \mathcal{L} &= \mathbb{E}_{q(\mathbf{z}_n; \lambda_n)} \left[\mathbb{E}_{q(\mathbf{z}_{-n}; \lambda_{-n})} [f_n(\mathbf{z})] \right] \\ &= \mathbb{E}_{r(\mathbf{z}_n; \lambda_n, \tau_n)} \left[\frac{q(\mathbf{z}_n; \lambda_n)}{r(\mathbf{z}_n; \lambda_n, \tau_n)} \mathbb{E}_{q(\mathbf{z}_{-n}; \lambda_{-n})} [f_n(\mathbf{z})] \right]\end{aligned}$$

- ▶ Algorithm:
 1. Sample $\mathbf{z}^{(0)} \sim q(\mathbf{z}; \lambda)$
 2. For each n , sample $\mathbf{z}_n^s \sim r(\mathbf{z}_n; \lambda_n, \tau_n)$, for $s = 1, \dots, S$
 3. For each n , obtain a Monte Carlo estimate of $\nabla_{\lambda_n} \mathcal{L}$
 4. Take a gradient step for λ

⁸Titsias and Lázaro-Gredilla (2015)

Implementation

- ▶ Need to choose the dispersion coefficients τ_n
 - Gradient steps for τ_n to minimize the variance
 - Monte Carlo estimator with little extra overhead

⁹Veach and Guibas (1995)

Implementation

- ▶ Need to choose the dispersion coefficients τ_n
 - Gradient steps for τ_n to minimize the variance
 - Monte Carlo estimator with little extra overhead
- ▶ High variance of the importance weights
 - Multiple importance sampling⁹
 - The proposal $r(z_n; \lambda_n, \tau_{n1}, \tau_{n2})$ can be a mixture

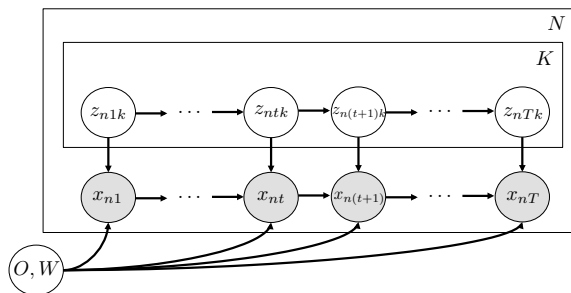
⁹Veach and Guibas (1995)

Full Algorithm

- ▶ Control variates
- ▶ Rao-Blackwellization
- ▶ O-BBVI with
 - ▶ Local expectations
 - ▶ Adaptation of the dispersion coefficients
 - ▶ Multiple proposals

Experiments: GN-TS Model

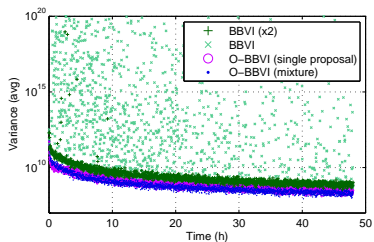
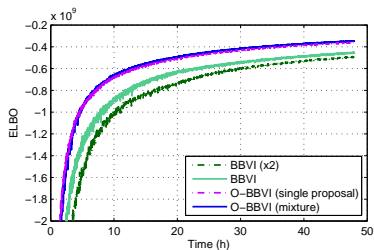
Gamma-Normal Time Series¹⁰ Model



¹⁰Ranganath et al. (2014)

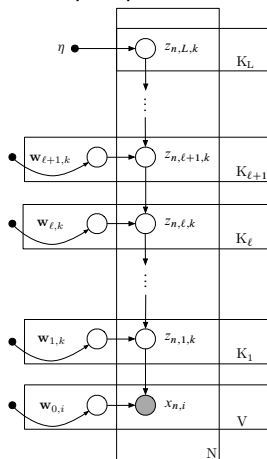
Experiments: GN-TS Model

Dataset: Synthetic



Experiments: DEF

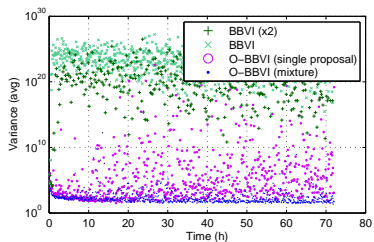
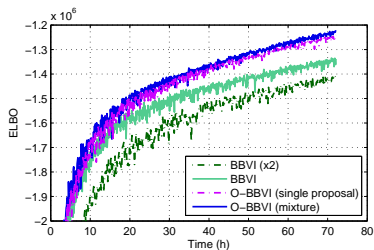
Poisson Deep Exponential Family¹¹



¹¹Ranganath et al. (2015)

Experiments: DEF

Dataset: Papers in NIPS'11 conference



Summary: O-BBVI

- ▶ Unconventional application of importance sampling to general VI
- ▶ Reduce the variance of the gradient estimator
- ▶ Lower variance than BBVI with $2\times$ Monte Carlo samples
- ▶ Faster convergence

Thank you for your attention!

