# [ Paper review 44 ]

# Variational Inference ; A Review for Statisticians

## ( Blei, et.al , 2018 )

## [ Contents ]

# 1. Abstract

core problem : **difficult-to-compute pdf**

This paper : review VARIATIONAL INFERENCE

- approximates via optimization

- step 1) posit a family of densities

    step 2) find the member of that family, closest to the target

# 2. Introduction

VI :faster & easier to scale to LARGE data

$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{z})p(\mathbf{x} \mid \mathbf{z}).$$

- latent variables : $\mathbf{z} = z_{1:m}$

  ( help govern the distn of data )
- observation : $x = x_{1:n}$

Draw latent variable from prior $p(z)$

Relate them to observation, via likelihood $p(\mathbf{z} \mid \mathbf{x})$

MCMC

- step 1) construct ergodic Markov chain on $z$

  ( whose stationary distn is posterior $p(\mathbf{z} \mid \mathbf{x})$ )
- step 2) sample from the chain ( = collect samples from stationary distn )
- step 3) approximate posterior with collected samples

VI is needed, when we need faster speed than MCMC

- when data sets are **large**
- when models are **complex**

Rather than sampling, use optimization!

Key point

- **KEY POINT 1** choose approximating distn to be **flexible**
- **KEY POINT 2** simple enough for efficient optimization

## VI vs MCMC

MCMC :

- computationally intensive, but (asymptotically) exact samples
- suited to smaller dataset

VI :

- faster than MCMC ( can use stochastic optimization ), but do not guarantee exact samples
- suited to largerdataset
- when we want to quickly explore many models

Not only data size, but also **geometry of the posterior**

- multi-mode : VI > MCMC

## Modern research in VI

- problem which involve massive data
- using improved optimization method
- easy to apply to a wide class of models
- increase the accuracy of VI ( by stretching the boundaries of approx distn )

## This paper...

- [Section 2] MVFI, CAVI

- [Section 3] Bayesian Mixture of Gaussians

- [Section 4-1&2] When joint density of $z$ and $x$ are exponential family

  [Section 4-3] SVI

# 3. Variational Inference

Goal of VI : approximate conditional density of latent variables, given observed variables ( = $P(Z \mid X)$ )

key : solve using **optimization**

( use family of densities over latent variables, parameterized by free **variational parameters** )

## 3-1. Problem of Approximate Inference

$p(\mathbf{z} \mid \mathbf{x}) = \frac{p(\mathbf{z},\mathbf{x})}{p(\mathbf{x})}$

- evidence : $p(\mathbf{x}) = \int p(\mathbf{z},\mathbf{x})\mathrm{d}\mathbf{z}$ ( intractable )

## Bayesian Mixture of Gaussians

- $K$ mixture components

- mean params : $\mu = \{\mu_1, \ldots, \mu_K\}$
  - drawn from $p(\mu_k) = \mathcal{N}(0, \sigma^2)$
- how to generate $x_i$ ?
  - step 1) choose cluster assignment $c_i$
  - step 2) draw $x_i$ from $\mathcal{N}(c_i^\top \mu, 1)$
- Full hierarchical model :
$$\mu_k \sim \mathcal{N}(0, \sigma^2), \qquad k = 1, \ldots, K,$$
$$c_i \sim \text{Categorical}(1/K, \ldots, 1/K), \qquad i = 1, \ldots, n.$$
$$x_i \mid c_i, \mu \sim \mathcal{N}(c_i^\top \mu, 1) \qquad i = 1, \ldots, n$$

- Joint pdf : (latent var : $\mathbf{z} = \{\mu, \mathbf{c}\}$ )

  $p(\mu, \mathbf{c}, \mathbf{x}) = p(\mu) \prod_{i=1}^{n} p(c_i) p(x_i \mid c_i, \mu).$

- Evidence :

$$p(\mathbf{x}) = \int p(\mu) \prod_{i=1}^{n} \sum_{c_i} p(c_i) \, p(x_i \mid c_i, \mu) \, \mathrm{d}\mu$$
$$= \sum_{\mathbf{c}} p(\mathbf{c}) \int p(\mu) \prod_{i=1}^{n} p(x_i \mid c_i, \mu) \, \mathrm{d}\mu$$

- time complexity of $K$-dim : $O(K^n)$

## 3-2. ELBO

Optimization problem ( $q^*(\mathbf{z}) = \underset{q(\mathbf{z}) \in \mathcal{Q}}{\arg\min} \, \mathrm{kL}(q(\mathbf{z}) \| p(\mathbf{z} \mid \mathbf{x}))$ )

Minimize KL Divergence :

$$\mathrm{KL}(q(\mathbf{z}) \| p(\mathbf{z} \mid \mathbf{x})) = \mathbb{E}[\log q(\mathbf{z})] - \mathbb{E}[\log p(\mathbf{z} \mid \mathbf{x})]$$
$$= \mathbb{E}[\log q(\mathbf{z})] - \mathbb{E}[\log p(\mathbf{z}, \mathbf{x})] + \log p(\mathbf{x})$$

Maximize ELBO:

$$\mathrm{ELBO}(q) = \mathbb{E}[\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}[\log q(\mathbf{z})]$$

Interpretation of ELBO :

$$\mathrm{ELBO}(q) = \mathbb{E}[\log p(\mathbf{z})] + \mathbb{E}[\log p(\mathbf{x} \mid \mathbf{z})] - \mathbb{E}[\log q(\mathbf{z})]$$
$$= \mathbb{E}[\log p(\mathbf{x} \mid \mathbf{z})] - \mathrm{kL}(q(\mathbf{z}) \| p(\mathbf{z}))$$

- first term ) fit well
- second term ) regularize well

Relationship between ELBO & $\log p(x)$

- used for **model selection** criterion

EM vs VI

- EM assumes, expectation under $p(z \mid x)$ is computable
- VI does not estimate fixed-model parameters ( classical params are treated as latent variables )

## 3-3. MFVI

Assumption : independency between latent variables

$$\to q(\mathbf{z}) = \prod_{j=1}^{m} q_j(z_j)$$

Each latent variable $z_j$ is governed by its own variational factor, $q_j(z_j)$

( these variational factors are chosen to MAXIMIZE ELBO )

- ex) choose as Gaussian factor, or categorical factor

Researchers have also studied more complex families

- 1) Structured VI
- 2) Mixture of variational densities

$\rightarrow$ both improve the fidelity of the approximation

( trade-off : dificult to solve variational optimization problem )

## Bayesian Mixture of Gaussians (cont)

MFVI : $q(\mu, \mathbf{c}) = \prod_{k=1}^{K} q\left(\mu_k; m_k, s_k^2\right) \prod_{i=1}^{n} q\left(c_i; \varphi_i\right).$

- $q\left(\mu_k; m_k, s_k^2\right)$ : Gaussian distn
- $q\left(c_i; \varphi_i\right)$ : its assignment probabilities are a $K$ -vector $\varphi_i$

Summary : ELBO is defined by..

- model definition : $p(\mu, \mathbf{c}, \mathbf{x}) = p(\mu) \prod_{i=1}^{n} p\left(c_i\right) p\left(x_i \mid c_i, \mu\right)$
- MFVI : $q(\mu, \mathbf{c}) = \prod_{k=1}^{K} q\left(\mu_k; m_k, s_k^2\right) \prod_{i=1}^{n} q\left(c_i; \varphi_i\right)$

# 3-4. CAVI ( Coordinate ascent MFVI )

CAVI : **iteratively** optimizes each factor of the MF variational density

optimal solution

- conditional : $q_j^*\left(z_j\right) \propto \exp\{\mathbb{E}_{-j}\left[\log p\left(z_j \mid \mathbf{z}_{-j}, \mathbf{x}\right)\right]\}.$

- joint : $q_j^*\left(z_j\right) \propto \exp\{\mathbb{E}_{-j}\left[\log p\left(z_j, \mathbf{z}_{-j}, \mathbf{x}\right)\right]\}.$

  ( expectation on RHS do not involve $j^{th}$ variational factor $\rightarrow$ valid coordinate update )

---

**Algorithm 1:** Coordinate ascent variational inference (CAVI)

---

**Input:** A model $p(\mathbf{x}, \mathbf{z})$, a data set $\mathbf{x}$

**Output:** A variational density $q(\mathbf{z}) = \prod_{j=1}^{m} q_j(z_j)$

**Initialize:** Variational factors $q_j(z_j)$

**while** *the* CAVI *has not converged* **do**
    **for** $j \in \{1, \dots, m\}$ **do**
        | Set $q_j(z_j) \propto \exp\{\mathbb{E}_{-j}[\log p(z_j \mid \mathbf{z}_{-j}, \mathbf{x})]\}$
    **end**
    Compute ELBO$(q) = \mathbb{E}\left[\log p(\mathbf{z}, \mathbf{x})\right] - \mathbb{E}\left[\log q(\mathbf{z})\right]$
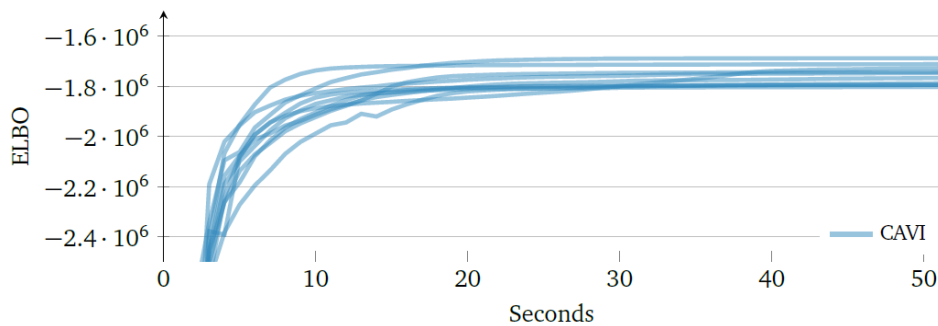**end**
**return** $q(\mathbf{z})$

---

( can see that it is closely related to Gibbs sampling )

# 3-5. Practicalities

## Initialization

- ELBO : (usually) non-convex objectiv function
- CAVI only guarantees local optimum ( sensitive to initialization )



/

- 10 random initialization, reaching different values!

  ( many local optima in ELBO )

- Not always bad!
    - ex) Mixture of Gaussian : many posterior modes
    - exploring latent clusters, predicting new observation


## Assessing convergence

- computing the ELBO of full dataset may be undesirable
- proxy! average log predictive of a small held-out dataset


## Numerical stability

- probabilities should be between 0~1
- use **log-sum-exp** trick

  $\log[\sum_i \exp(x_i)] = \alpha + \log[\sum_i \exp(x_i - \alpha)]$.


# 4. A complete example : Bayesian Mixture of Gaussians

notation

- $K$ real valued mean params : $\mu = \mu_{1:K}$
- $n$ latent class assignments : $\mathbf{c} = c_{1:n}$


ELBO for mixture of Gaussians

- variational parameters : $\mathbf{m}, \mathbf{s}^2, \varphi$

$$\mathrm{ELBO}\left(\mathbf{m}, \mathbf{s}^2, \varphi\right) = \sum_{k=1}^{K} \mathbb{E}\left[\log p\left(\mu_k\right); m_k, s_k^2\right]$$

$$+ \sum_{i=1}^{n}\left(\mathbb{E}\left[\log p\left(c_i\right); \varphi_i\right] + \mathbb{E}\left[\log p\left(x_i \mid c_i, \mu\right); \varphi_i, \mathbf{m}, \mathbf{s}^2\right]\right)$$

$$- \sum_{i=1}^{n} \mathbb{E}\left[\log q\left(c_i; \varphi_i\right)\right] - \sum_{k=1}^{K} \mathbb{E}\left[\log q\left(\mu_k; m_k, s_k^2\right)\right]$$

CAVI updates each variational parameters in turn

# 4-1. (step 1)The variational density of the "mixture assignments"

(review) optimal solution : $q_j^*\left(z_j\right) \propto \exp\{\mathbb{E}_{-j}\left[\log p\left(z_j, \mathbf{z}_{-j}, \mathbf{x}\right)\right]\}$

(1) derive variational update for $c_i$ ( cluster assignment )

- $q^*\left(c_i; \varphi_i\right) \propto \exp\{\log p\left(c_i\right) + \mathbb{E}\left[\log p\left(x_i \mid c_i, \mu\right); \mathbf{m}, \mathbf{s}^2\right]\}$.
    - 1st term) log prior of $c_i$ : $\log p\left(c_i\right) = -\log K$
    - 2nd term) expected log of the $c_i$th Gaussian density
        - $p\left(x_i \mid c_i, \mu\right) = \prod_{k=1}^{K} p\left(x_i \mid \mu_k\right)^{c_{ik}}$
        - Thus,

        $$\mathbb{E}\left[\log p\left(x_i \mid c_i, \mu\right)\right] = \sum_{k} c_{ik} \mathbb{E}\left[\log p\left(x_i \mid \mu_k\right); m_k, s_k^2\right]$$
        $$= \sum_{k} c_{ik} \mathbb{E}\left[-(x_i - \mu_k)^2/2; m_k, s_k^2\right] + \text{const.}$$
        $$= \sum_{k} c_{ik}\left(\mathbb{E}\left[\mu_k; m_k, s_k^2\right] x_i - \mathbb{E}\left[\mu_k^2; m_k, s_k^2\right]/2\right) + \text{const.}$$

- result : $\varphi_{ik} \propto \exp\{\mathbb{E}\left[\mu_k; m_k, s_k^2\right] x_i - \mathbb{E}\left[\mu_k^2; m_k, s_k^2\right]/2\}$

# 4-2. (step 2)The variational density of the "mixture-component means"

(2) variational density $q\left(\mu_k; m_k, s_k^2\right)$ of $k^{th}$ mixture components

$q\left(\mu_k\right) \propto \exp\{\log p\left(\mu_k\right) + \sum_{i=1}^{n} \mathbb{E}\left[\log p\left(x_i \mid c_i, \mu\right); \varphi_i, \mathbf{m}_{-k}, \mathbf{s}_{-k}^2\right]\}$.

Unnormallized log of $q\left(\mu_k\right)$ :

$$\log q(\mu_k) = \log p(\mu_k) + \sum_i \mathbb{E}\left[\log p(x_i \mid c_i, \mu); \varphi_i, \mathbf{m}_{-k}, \mathbf{s}^2_{-k}\right] + \text{const.}$$

$$= \log p(\mu_k) + \sum_i \mathbb{E}\left[c_{ik} \log p(x_i \mid \mu_k); \varphi_i\right] + \text{const.}$$

$$= -\mu_k^2/2\sigma^2 + \sum_i \mathbb{E}\left[c_{ik}; \varphi_i\right] \log p(x_i \mid \mu_k) + \text{const.}$$

$$= -\mu_k^2/2\sigma^2 + \sum_i \varphi_{ik}\left(-(x_i - \mu_k)^2/2\right) + \text{const.}$$

$$= -\mu_k^2/2\sigma^2 + \sum_i \varphi_{ik} x_i \mu_k - \varphi_{ik}\mu_k^2/2 + \text{const.}$$

$$= \left(\sum_i \varphi_{ik} x_i\right)\mu_k - \left(1/2\sigma^2 + \sum_i \varphi_{ik}/2\right)\mu_k^2 + \text{const.}$$

where $\varphi_{ik} = \mathbb{E}\left[c_{ik}; \varphi_i\right]$.

Thus, $m_k = \dfrac{\sum_i \varphi_{ik} x_i}{1/\sigma^2 + \sum_i \varphi_{ik}}, \quad s_k^2 = \dfrac{1}{1/\sigma^2 + \sum_i \varphi_{ik}}$

## 4-3. CAVI for the mixture of Gaussians

---
**Algorithm 2:** CAVI for a Gaussian mixture model

---
**Input:** Data $x_{1:n}$, number of components $K$, prior variance of component means $\sigma^2$

**Output:** Variational densities $q(\mu_k; m_k, s_k^2)$ (Gaussian) and $q(c_i; \varphi_i)$ ($K$-categorical)

**Initialize:** Variational parameters $\mathbf{m} = m_{1:K}$, $\mathbf{s}^2 = s^2_{1:K}$, and $\boldsymbol{\varphi} = \varphi_{1:n}$

**while** *the* ELBO *has not converged* **do**

    **for** $i \in \{1, \ldots, n\}$ **do**

        Set $\varphi_{ik} \propto \exp\{\mathbb{E}\left[\mu_k; m_k, s_k^2\right] x_i - \mathbb{E}\left[\mu_k^2; m_k, s_k^2\right]/2\}$

    **end**

    **for** $k \in \{1, \ldots, K\}$ **do**

        Set $m_k \longleftarrow \dfrac{\sum_i \varphi_{ik} x_i}{1/\sigma^2 + \sum_i \varphi_{ik}}$

        Set $s_k^2 \longleftarrow \dfrac{1}{1/\sigma^2 + \sum_i \varphi_{ik}}$

    **end**

    Compute ELBO$(\mathbf{m}, \mathbf{s}^2, \boldsymbol{\varphi})$

**end**

**return** $q(\mathbf{m}, \mathbf{s}^2, \boldsymbol{\varphi})$

---

Approximate predictive : (mixture of Gaussians)

$$p(x_{\text{new}} \mid x_{1:n}) \approx \frac{1}{K} \sum_{k=1}^{K} p(x_{\text{new}} \mid m_k)$$

# 5. VI with exponential families

(Until now....)

- **MFVI**

- **CAVI** ( coordinate ascent algorithm for optimizing ELBO )

- demonstration using **simple mixture of Gaussians**

  ( available in closed-form )


Now, will work with **exponential family**

- working with this simplifies VI
- easier to derive CAVI
- section 5-1) general case

  section 5-2) conditionally conjugate models

  section 5-3) SVI


# 5-1. Complete conditionals in the exponential family

suppose **complete conditional** is "exponential family" :

$$p\left(z_j \mid \mathbf{z}_{-j}, \mathbf{x}\right) = h\left(z_j\right) \exp\left\{\eta_j(\mathbf{z}_{-j}, \mathbf{x})^\top z_j - a\left(\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)\right)\right\}$$


CAVI with MFVI

- coordinate update of $q_j^*\left(z_j\right) \propto \exp\{\mathbb{E}_{-j}\left[\log p\left(z_j \mid \mathbf{z}_{-j}, \mathbf{x}\right)\right]\}$ =

  $$\begin{aligned} q\left(z_j\right) &\propto \exp\{\mathbb{E}\left[\log p\left(z_j \mid \mathbf{z}_{-j}, \mathbf{x}\right)\right]\} \\ &= \exp\left\{\mathbb{E}\left[\log h\left(z_j\right)\exp\left\{\eta_j(\mathbf{z}_{-j}, \mathbf{x})^\top z_j - a\left(\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)\right)\right\}\right]\right\} \\ &= \exp\left\{\log h\left(z_j\right) + \mathbb{E}[\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)]^\top z_j - \mathbb{E}\left[a\left(\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)\right)\right]\right\} \\ &\propto h\left(z_j\right)\exp\left\{\mathbb{E}[\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)]^\top z_j\right\} \end{aligned}$$

- set $v_j = \mathbb{E}\left[\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)\right]$


# 5-2. Conditional conjugacy and Bayesian models

Special case of exponential family : "conditional conjugate models" with local & global variables


## Conditionally conjugate models

- $\beta$ : global latent variables
- $\mathbf{z}$ : local latent variables
- joint pdf : $p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta)\prod_{i=1}^n p\left(z_i, x_i \mid \beta\right)$


Assumption 1) joint density of each $\left(x_i, z_i\right)$ pair, conditional on $\beta$ = exponential family

$$p\left(z_i, x_i \mid \beta\right) = h\left(z_i, x_i\right)\exp\left\{\beta^\top t\left(z_i, x_i\right) - a(\beta)\right\}\ldots\ldots\ldots\ldots\ldots\ldots(a)$$


Assumption 2) prior (on global variables) to be conjugate prior

$$p(\beta) = h(\beta) \exp\left\{\alpha^\top [\beta, -a(\beta)] - a(\alpha)\right\} \dots\dots\dots\dots\dots\dots\dots\dots\dots \text{(b)}$$

- natural (hyper) parameter $\alpha = [\alpha_1, \alpha_2]^\top$

(a) and (b) : conjugate $\rightarrow \hat{\alpha} = \left[\alpha_1 + \sum_{i=1}^n t(z_i, x_i), \alpha_2 + n\right]^\top$

Complete conditional of local variable $z_i$ :

- (given $\beta$ and $x_i$) $z_i$ is conditionally independent!

  $\rightarrow \quad p(z_i \mid x_i, \beta, \mathbf{z}_{-i}, \mathbf{x}_{-i}) = p(z_i \mid x_i, \beta)$

- assumption ) exponential family

  $\rightarrow \quad p(z_i \mid x_i, \beta) = h(z_i) \exp\left\{\eta(\beta, x_i)^\top z_i - a(\eta(\beta, x_i))\right\}$

# VI in conditionally conjugate models

describe CAVI for general class of models

notation

- $\lambda$ : **global** variational parameter

  $q(\beta \mid \lambda)$ : variational posterior approximation on $\beta$

- $\phi$ : **local** variational parameter

  $q(z_i \mid \phi)$ : variational posterior on each local variable $z_i$

**Local** variational update : $\varphi_i = \mathbb{E}_\lambda [\eta(\beta, x_i)]$ .... by

- (1) $v_j = \mathbb{E}[\eta_j(\mathbf{z}_{-j}, \mathbf{x})]$
- (2) $p(z_i \mid x_i, \beta, \mathbf{z}_{-i}, \mathbf{x}_{-i}) = p(z_i \mid x_i, \beta)$

**Global** variational update : $\lambda = \left[\alpha_1 + \sum_{i=1}^n \mathbb{E}_{\varphi_i}[t(z_i, x_i)], \alpha_2 + n\right]^\top$

- expectation of $\hat{\alpha} = \left[\alpha_1 + \sum_{i=1}^n t(z_i, x_i), \alpha_2 + n\right]^\top$

CAVI optimizes the ELBO by iterating "local updates" and "global updates"

To assess convergence, compute ELBO at each iteration!

ELBO :

- $\mathrm{ELBO}(q) = \mathbb{E}[\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}[\log q(\mathbf{z})]$

  $p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i \mid \beta)$

- Therefore, $\mathrm{ELBO} = \left(\alpha_1 + \sum_{i=1}^n \mathbb{E}_{\varphi_i}[t(z_i, x_i)]\right)^\top \mathbb{E}_\lambda[\beta] - (\alpha_2 + n)\mathbb{E}_\lambda[a(\beta)] - \mathbb{E}[\log q(\beta, \mathbf{z})]$

  - $\mathbb{E}[\log q(\beta, \mathbf{z})] = \lambda^\top \mathbb{E}_\lambda[t(\beta)] - a(\lambda) + \sum_{i=1}^n \varphi_i^\top \mathbb{E}_{\varphi_i}[z_i] - a(\varphi_i)$

# 5-3. SVI ( Stochastic Variational Inference )

Modern problems : require analyizing massive data

but, most do not easily scale ( ex. CAVI )


CAVI, not scalable!

- requires iteration through entire data at each iteration
- alternative : **gradient-based optimization**

  ( = Key to SVI )


SVI focuses on optimizing the global variational params $\lambda$ of conditionally conjugate model

- step 1) subsample data
- step 2) use current global param to compute **optimal local params** for the subsampled data
- step 3) adjust the current **global params**


## Natural gradient of ELBO

SVI focuses on optimizing the global variational params $\lambda$


Euclidan gradient of ELBO :

- $\nabla_\lambda \mathrm{ELBO} = a''(\lambda)\left(\mathbb{E}_\varphi[\hat\alpha] - \lambda\right)$ ( Hoffman et al. 2013)

Natural gradient : (premultiply by inverse Fisher info )

- $g(\lambda) = \mathbb{E}_\varphi[\hat\alpha] - \lambda$

Update param, using **natural gradient** in gradient-based optimization method

- at each iteration, update $\lambda_t = \lambda_{t-1} + \epsilon_t g\left(\lambda_{t-1}\right)$

  ( $= \lambda_t = (1 - \epsilon_t)\lambda_{t-1} + \epsilon_t \mathbb{E}_\varphi[\hat\alpha]$ )


## Stochastic Optimization of the ELBO

goal : construct a cheaply computed, noisy, unbiiased natural gradient

- $g(\lambda) = \mathbb{E}_\varphi[\hat\alpha] - \lambda$
- $\hat\alpha = \left[\alpha_1 + \sum_{i=1}^n t\left(z_i, x_i\right), \alpha_2 + n\right]^\top$

$\to g(\lambda) = \alpha + \left[\sum_{i=1}^n \mathbb{E}_{\varphi_i^*}\left[t\left(z_i, x_i\right)\right], n\right]^\top - \lambda$

  ( $\varphi_i^*$ : optimized local variational param )


Construct noisy natural gradient, by **sampling an index from the data** & rescaling

- $t \sim \mathrm{Unif}(1, \ldots, n)$

$$\hat{g}(\lambda) = \alpha + n \left[ \mathbb{E}_{\varphi_t^*} \left[ t\left(z_t, x_t\right) \right], 1 \right]^\top - \lambda$$

- $\hat{g}(\lambda)$ is unbiased ( $\mathbb{E}_t[\hat{g}(\lambda)] = g(\lambda)$ ) & ceheap to compute

  ( not only sample 1 data, but can also use mini-batch )

Step-size sequence

- should follow conditions of Robbins and Monro,

  $\sum_t \epsilon_t = \infty; \sum_t \epsilon_t^2 < \infty$

---

**Algorithm 3:** SVI for conditionally conjugate models

---

**Input:** Model $p(\mathbf{x}, \mathbf{z})$, data $\mathbf{x}$, and step size sequence $\epsilon_t$

**Output:** Global variational densities $q_\lambda(\beta)$

**Initialize:** Variational parameters $\lambda_0$

**while** *TRUE* **do**

    Choose a data point uniformly at random, $t \sim \text{Unif}(1, \ldots, n)$

    Optimize its local variational parameters $\varphi_t^* = \mathbb{E}_\lambda \left[ \eta(\beta, x_t) \right]$

    Compute the coordinate update as though $x_t$ was repeated $n$ times,

$$\hat{\lambda} = \alpha + n \mathbb{E}_{\varphi_t^*} \left[ f(z_t, x_t) \right]$$

    Update the global variational parameter, $\lambda_t = (1 - \epsilon_t) \lambda_t + \epsilon_t \hat{\lambda}_t$

**end**

**return** $\lambda$

---

# 6. Discussion

Summary

- MFVI
- CAVI
- Bayesian Mixture of Gaussians
- special case of exponential families & conditional conjugacy
- SVI