

# Variational Inference: an introduction

Massil Achab

May 18, 2017

This presentation is inspired by the very insightful review from Blei et al. [2017].

## 1 Inference of probabilistic models

A probabilistic model asserts how observations arise from a natural phenomenon.

We design the model via a *joint distribution*

$$p(\mathbf{x}, \mathbf{z})$$

of observed variables  $\mathbf{x}$  corresponding to data, and latent variables  $\mathbf{z}$  that provide the hidden structure to generate from  $\mathbf{x}$ .

The *likelihood*

$$p(\mathbf{x}|\mathbf{z})$$

is a probability that describes how any data  $\mathbf{x}$  is likely given a particular hidden pattern described by  $\mathbf{z}$ .

The *prior*

$$p(\mathbf{z})$$

posits a generating process of the hidden structure.

Inference amounts to conditioning on data and computing the *posterior*

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

This last distribution is difficult to compute because of the normalizing constant. Such computation often requires *approximate inference*. For decades, MCMC was the method of choice to do approximate inference. MCMC methods construct an ergodic Markov chain on  $\mathbf{z}$  with  $p(\mathbf{z}|\mathbf{x})$  as stationary distribution. Another approach, Variational Inference (VI) approximates the posterior  $p(\mathbf{z}|\mathbf{x})$  through optimization and tends to be faster than MCMC sampling.

## 2 Variational Inference

### 2.1 Idea behind Variational Inference

Core idea of VI:

- posit a family of distributions  $q(\mathbf{z}) \in \mathcal{D}$

- match  $q(\mathbf{z})$  to the posterior  $p(\mathbf{z}|\mathbf{x})$

This strategy converts the problem of computing the posterior  $p(\mathbf{z}|\mathbf{x})$  into an optimization problem of the form:

$$q^*(\mathbf{z}) = \arg \min_{q(\mathbf{z}) \in \mathcal{D}} \text{divergence}(q(\mathbf{z}), p(\mathbf{z}|\mathbf{x})). \quad (1)$$

The optimized  $q^*(\mathbf{z})$  is then used as a proxy of  $p(\mathbf{z}|\mathbf{x})$ .

MCMC and VI are different approaches for solving the same problem. Comparison:

- MCMC asymptotically provides exact samples from the posterior distribution.
- VI faster and can be parallelized.
- VI performs well on mixture models.

## 2.2 Evidence Lower Bound

The usual criterion to match  $q(\mathbf{z})$  to  $p(\mathbf{z}|\mathbf{x})$  uses the Kullback-Leibler divergence. We now minimize:

$$\begin{aligned} \text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) &= \mathbb{E}^q[\log q(\mathbf{z})] - \mathbb{E}^q[\log p(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}^q[\log q(\mathbf{z})] - \mathbb{E}^q[\log p(\mathbf{x}, \mathbf{z})] + \log p(\mathbf{x}) \end{aligned}$$

The *evidence* is defined as  $p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z})d\mathbf{z}$ . Since the evidence does not depend on  $q(\mathbf{z})$ , we optimize an alternative objective:

$$\text{ELBO}(q) = \mathbb{E}^q[\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}^q[\log q(\mathbf{z})]. \quad (2)$$

Maximizing the ELBO function is equivalent to minimize the Kullback-Leibler divergence above. The name **ELBO** stands for **Evidence Lower Bound** because of the following inequality:

$$\begin{aligned} \log p(\mathbf{x}) &= \text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) + \text{ELBO}(q) \\ &\geq \text{ELBO}(q), \end{aligned}$$

since Kullback-Leibler divergence always takes positive values.

## 2.3 Mean-field variational family

A popular choice of family is the *mean-field variational family* where the latent variables are mutually independent. A generic member writes:

$$q(\mathbf{z}) = \prod_{j=1}^m q_j(z_j) \quad (3)$$

Let emphasize that the variational family is not a model of the observed data.

The mutual independence of the latent variables triggers the separability of the ELBO's second term:

$$\mathbb{E}^q[\log q(\mathbf{z})] = \sum_{j=1}^m \mathbb{E}^q[\log q_j(z_j)]$$

## 2.4 CAVI: Coordinate Ascent mean-field Variational Inference

We have now cast the approximate conditional inference as an optimization problem. The previous remark on the separability incites the use of a coordinate ascent algorithm. Rewriting ELBO as a function of  $q_j$ , we prove that the optimal density  $q_j^*$  satisfies

$$q_j^*(z_j) \propto \exp\{\mathbb{E}_{-j}^q[\log p(z_j | \mathbf{z}_{-j}, \mathbf{x})]\} \quad (4)$$

## 3 Application: Bayesian Mixture of Gaussians

We consider  $K$  mixture components with means  $\boldsymbol{\mu} = \{\mu_1, \dots, \mu_K\}$  drawn from  $\mathcal{N}(0, \sigma^2)$ . The full hierarchical model writes

$$\begin{aligned} \mu_k &\sim \mathcal{N}(0, \sigma^2) & k &= 1, \dots, K, \\ c_i &\sim \text{Cat}(1/K, \dots, 1/K) & i &= 1, \dots, n, \\ x_i | c_i, \boldsymbol{\mu} &\sim \mathcal{N}(c_i^\top \boldsymbol{\mu}, 1) & i &= 1, \dots, n, \end{aligned}$$

where  $c_i \in \mathbb{R}^K$  is a one-hot encoded vector that assigns the latent class to  $x_i$ .

In this application,  $\mathbf{z} = (\mathbf{c}, \boldsymbol{\mu})$ . Following the previously detailed steps, we efficiently obtain an approximation of the posterior density  $p(\mathbf{z} | \mathbf{x})$ . See Section 3 in Blei et al. [2017] for more details.

## References

David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, (just-accepted), 2017.