

Introduction to Variational Inference

Dmitry Kropotov
Lomonosov Moscow State University



Contents

Variational Inference

Mean Field Approximation

Example: Latent Dirichlet Allocation

Variational Parametric Approximation

Example: Bayesian Logistic Regression

The problem: finding posterior distribution

In Bayesian Inference we are interested in finding posterior distributions.

Bayesian Ensembling:

Probabilistic model: $p(\mathbf{x}, \theta) = p(\mathbf{x}|\theta)p(\theta)$, $\mathbf{x} = [\mathbf{x}_{train}, \mathbf{x}_{test}]$.

Prediction: $p(\mathbf{x}_{test}|\mathbf{x}_{train}) = \int p(\mathbf{x}_{test}|\theta)p(\theta|\mathbf{x}_{train})d\theta$.

EM-algorithm:

Latent variable model: $p(\mathbf{x}, \mathbf{z}|\theta)$.

Training: $\log p(\mathbf{x}|\theta) \geq \underbrace{\mathbb{E}_{q(\mathbf{z})} \log p(\mathbf{x}, \mathbf{z}|\theta) - \mathbb{E}_{q(\mathbf{z})} \log q(\mathbf{z})}_{ELBO(q, \theta)} \rightarrow \max_{q(\mathbf{z}), \theta}.$

E-step: $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}, \theta)$.

The problem: finding posterior distribution

In Bayesian Inference we are interested in finding posterior distributions.

Bayesian Ensembling:

Probabilistic model: $p(\mathbf{x}, \theta) = p(\mathbf{x}|\theta)p(\theta)$, $\mathbf{x} = [\mathbf{x}_{train}, \mathbf{x}_{test}]$.

Prediction: $p(\mathbf{x}_{test}|\mathbf{x}_{train}) = \int p(\mathbf{x}_{test}|\theta)p(\theta|\mathbf{x}_{train})d\theta$.

EM-algorithm:

Latent variable model: $p(\mathbf{x}, \mathbf{z}|\theta)$.

Training: $\log p(\mathbf{x}|\theta) \geq \underbrace{\mathbb{E}_{q(\mathbf{z})} \log p(\mathbf{x}, \mathbf{z}|\theta) - \mathbb{E}_{q(\mathbf{z})} \log q(\mathbf{z})}_{ELBO(q, \theta)} \rightarrow \max_{q(\mathbf{z}), \theta}.$

E-step: $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}, \theta)$.

Posterior distributions can be calculated analytically only for simple conjugate models!

Variational Inference

Probabilistic model: $p(\mathbf{x}, \boldsymbol{\theta}) = p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})$.

Main idea: Find posterior approximation $p(\boldsymbol{\theta}|\mathbf{x}) \approx q(\boldsymbol{\theta}) \in \mathcal{Q}$, using the following criterion function:

$$F(q) := \text{KL}(q(\boldsymbol{\theta}) || p(\boldsymbol{\theta}|\mathbf{x})) \rightarrow \min_{q(\boldsymbol{\theta}) \in \mathcal{Q}}. \quad (1)$$

Variational Inference

Probabilistic model: $p(\mathbf{x}, \boldsymbol{\theta}) = p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})$.

Main idea: Find posterior approximation $p(\boldsymbol{\theta}|\mathbf{x}) \approx q(\boldsymbol{\theta}) \in \mathcal{Q}$, using the following criterion function:

$$F(q) := \text{KL}(q(\boldsymbol{\theta}) || p(\boldsymbol{\theta}|\mathbf{x})) \rightarrow \min_{q(\boldsymbol{\theta}) \in \mathcal{Q}}. \quad (1)$$

$$\begin{aligned} F(q) &= \int q(\boldsymbol{\theta}) \log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{x})} d\boldsymbol{\theta} = \int q(\boldsymbol{\theta}) \log \frac{q(\boldsymbol{\theta})p(\mathbf{x})}{p(\mathbf{x}, \boldsymbol{\theta})} d\boldsymbol{\theta} = \\ &= \int q(\boldsymbol{\theta}) \log q(\boldsymbol{\theta}) d\boldsymbol{\theta} + \underbrace{\int q(\boldsymbol{\theta}) \log p(\mathbf{x}) d\boldsymbol{\theta}}_{\log p(\mathbf{x})} - \int q(\boldsymbol{\theta}) \log p(\mathbf{x}, \boldsymbol{\theta}) d\boldsymbol{\theta} = \\ &= \mathbb{E}_{q(\boldsymbol{\theta})}[\log q(\boldsymbol{\theta}) - \log p(\mathbf{x}, \boldsymbol{\theta})] + \text{const} \rightarrow \min_{q \in \mathcal{Q}} \quad (2) \end{aligned}$$

Variational Inference

Probabilistic model: $p(\mathbf{x}, \boldsymbol{\theta}) = p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})$.

Main idea: Find posterior approximation $p(\boldsymbol{\theta}|\mathbf{x}) \approx q(\boldsymbol{\theta}) \in \mathcal{Q}$, using the following criterion function:

$$F(q) := \text{KL}(q(\boldsymbol{\theta}) || p(\boldsymbol{\theta}|\mathbf{x})) \rightarrow \min_{q(\boldsymbol{\theta}) \in \mathcal{Q}}. \quad (1)$$

$$\begin{aligned} F(q) &= \int q(\boldsymbol{\theta}) \log \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{x})} d\boldsymbol{\theta} = \int q(\boldsymbol{\theta}) \log \frac{q(\boldsymbol{\theta})p(\mathbf{x})}{p(\mathbf{x}, \boldsymbol{\theta})} d\boldsymbol{\theta} = \\ &= \int q(\boldsymbol{\theta}) \log q(\boldsymbol{\theta}) d\boldsymbol{\theta} + \underbrace{\int q(\boldsymbol{\theta}) \log p(\mathbf{x}) d\boldsymbol{\theta}}_{\log p(\mathbf{x})} - \int q(\boldsymbol{\theta}) \log p(\mathbf{x}, \boldsymbol{\theta}) d\boldsymbol{\theta} = \\ &= \mathbb{E}_{q(\boldsymbol{\theta})}[\log q(\boldsymbol{\theta}) - \log p(\mathbf{x}, \boldsymbol{\theta})] + \text{const} \rightarrow \min_{q \in \mathcal{Q}} \quad (2) \end{aligned}$$

The problem (2) is equivalent to maximizing ELBO:

$$\log p(\mathbf{x}) \geq \text{ELBO}(q) = \mathbb{E}_{q(\boldsymbol{\theta})} \log p(\mathbf{x}, \boldsymbol{\theta}) - \mathbb{E}_{q(\boldsymbol{\theta})} \log q(\boldsymbol{\theta}) \rightarrow \max.$$

Variational EM-algorithm

Latent variable model: $p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})$.

Conventional EM-algorithm:

$$\log p(\mathbf{x}|\boldsymbol{\theta}) \geq \underbrace{\mathbb{E}_{q(\mathbf{z})} \log p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) - \mathbb{E}_{q(\mathbf{z})} \log q(\mathbf{z})}_{ELBO(q, \boldsymbol{\theta})} \rightarrow \max_{q(\mathbf{z}), \boldsymbol{\theta}},$$

$$q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}).$$

Variational EM-algorithm:

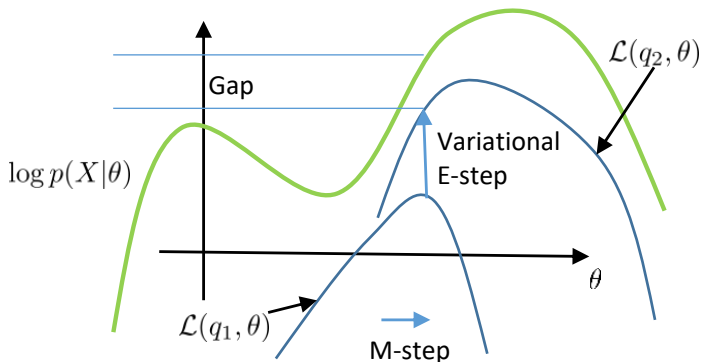
$$\log p(\mathbf{x}|\boldsymbol{\theta}) \geq \underbrace{\mathbb{E}_{q(\mathbf{z})} \log p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) - \mathbb{E}_{q(\mathbf{z})} \log q(\mathbf{z})}_{ELBO(q, \boldsymbol{\theta})} \rightarrow \max_{q(\mathbf{z}) \in \mathcal{Q}, \boldsymbol{\theta}},$$

$q(\mathbf{z}) \in \mathcal{Q}$ – variational approximation to $p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$.

Variational EM-algorithm

$$\log p(\mathbf{x}|\boldsymbol{\theta}) \geq \mathcal{L}(q, \boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{z})} \log p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) - \mathbb{E}_{q(\mathbf{z})} \log q(\mathbf{z}) \rightarrow \max_{q \in \mathcal{Q}, \boldsymbol{\theta}}.$$

Even in the case of inexact E-step in many situations we are able to find good model parameters.



ELBO interpretation

Probabilistic model: $p(\mathbf{x}, \boldsymbol{\theta}) = p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})$.

Finding variational approximation $p(\boldsymbol{\theta}|\mathbf{x}) \approx q(\boldsymbol{\theta}) \in \mathcal{Q}$:

$$\begin{aligned}\text{ELBO}(q) &= \mathbb{E}_{q(\boldsymbol{\theta})} \log p(\mathbf{x}, \boldsymbol{\theta}) - \mathbb{E}_{q(\boldsymbol{\theta})} \log q(\boldsymbol{\theta}) = \mathbb{E}_{q(\boldsymbol{\theta})} [\log p(\mathbf{x}|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) - \log q(\boldsymbol{\theta})] = \\ &= \mathbb{E}_{q(\boldsymbol{\theta})} \log p(\mathbf{x}|\boldsymbol{\theta}) + \mathbb{E}_{q(\boldsymbol{\theta})} \log \frac{p(\boldsymbol{\theta})}{q(\boldsymbol{\theta})} = \underbrace{\mathbb{E}_{q(\boldsymbol{\theta})} \log p(\mathbf{x}|\boldsymbol{\theta})}_{\text{data term}} - \underbrace{KL(q(\boldsymbol{\theta}) || p(\boldsymbol{\theta}))}_{\text{KL term}} \rightarrow \max_{q \in \mathcal{Q}}.\end{aligned}$$

Maximum likelihood approach:

$$\mathbb{E}_{q(\boldsymbol{\theta})} \log p(\mathbf{x}|\boldsymbol{\theta}) \rightarrow \max_q \Leftrightarrow q(\boldsymbol{\theta}) = \delta(\boldsymbol{\theta} - \boldsymbol{\theta}_{ML}), \quad \boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{x}|\boldsymbol{\theta})$$

Mean Field Approximation

Factorized family of variational distributions:

$$q(\boldsymbol{\theta}) = \prod_{j=1}^m q_j(\boldsymbol{\theta}_j), \quad \boldsymbol{\theta} = [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_m]. \quad (3)$$

Let's fix all factors $\{q_j(\boldsymbol{\theta}_j)\}_{j \neq i}$ except one in the factorized family (3) and find best variational approximation for $q_i(\boldsymbol{\theta}_i)$:

$$\begin{aligned} \text{ELBO}(q) &= \mathbb{E}_{q(\boldsymbol{\theta})}[\log p(\mathbf{x}, \boldsymbol{\theta}) - \log q(\boldsymbol{\theta})] = \\ &= \int \prod_{j=1}^m q_j(\boldsymbol{\theta}_j) \log p(\mathbf{x}, \boldsymbol{\theta}) d\boldsymbol{\theta} - \int \prod_{j=1}^m q_j(\boldsymbol{\theta}_j) \left(\sum_{k=1}^m \log q_k(\boldsymbol{\theta}_k) \right) d\boldsymbol{\theta} = \\ &= \int \prod_{j=1}^m q_j(\boldsymbol{\theta}_j) \log p(\mathbf{x}, \boldsymbol{\theta}) d\boldsymbol{\theta} - \sum_{k=1}^m \int \prod_{j=1}^m q_j(\boldsymbol{\theta}_j) \log q_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta} = \\ &= \int \prod_{j=1}^m q_j(\boldsymbol{\theta}_j) \log p(\mathbf{x}, \boldsymbol{\theta}) d\boldsymbol{\theta} - \sum_{k=1}^m \int q_k(\boldsymbol{\theta}_k) \log q_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k = \dots \end{aligned}$$

Mean Field Approximation

$$\begin{aligned} \dots &= \int q_i(\boldsymbol{\theta}_i) \left[\underbrace{\int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j}_{r_i(\boldsymbol{\theta}_i, \mathbf{x})} \right] d\boldsymbol{\theta}_i - \int q_i(\boldsymbol{\theta}_i) \log q_i(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i + \text{const} = \\ &= \int q_i(\boldsymbol{\theta}_i) \log \frac{\exp(r_i(\boldsymbol{\theta}_i, \mathbf{x}))}{q_i(\boldsymbol{\theta}_i)} d\boldsymbol{\theta}_i + \text{const} = [Z_i - \text{normalizing constant for distribution } r_i] = \\ &= \int q_i(\boldsymbol{\theta}_i) \log \underbrace{\frac{\exp(r_i(\boldsymbol{\theta}_i, \mathbf{x}))}{Z_i}}_{\hat{r}_i(\boldsymbol{\theta}_i, \mathbf{x})} \cdot \frac{Z_i}{q_i(\boldsymbol{\theta}_i)} d\boldsymbol{\theta}_i + \text{const} = \\ &= \int q_i(\boldsymbol{\theta}_i) \log \frac{\hat{r}_i(\boldsymbol{\theta}_i, \mathbf{x})}{q_i(\boldsymbol{\theta}_i)} d\boldsymbol{\theta}_i + \log Z_i + \text{const} = \text{const} - KL(q_i(\boldsymbol{\theta}_i) \parallel \hat{r}_i(\boldsymbol{\theta}_i, \mathbf{x})) \rightarrow \max_{q_i(\boldsymbol{\theta}_i)}. \end{aligned}$$

Solution: $q_i(\boldsymbol{\theta}_i) = \hat{r}_i(\boldsymbol{\theta}_i, \mathbf{x}) = \frac{1}{Z_i} \exp(\int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j)$

Mean Field Approximation

General scheme for Mean Field Variational Inference:

1. Initialize $q(\boldsymbol{\theta}) = \prod_{j=1}^m q_j(\boldsymbol{\theta}_j)$;
2. For each factor $q_i(\boldsymbol{\theta}_i)$, $i = 1, \dots, m$ do:
 - Calculate $q_i(\boldsymbol{\theta}_i) = \frac{1}{Z_i} \exp(\int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j)$, where $Z_i = \int \exp(\int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j) d\boldsymbol{\theta}_i$;
3. Calculate $\text{ELBO}(q) = \int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j=1}^m q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta} - \sum_{j=1}^m \int q_j(\boldsymbol{\theta}_j) \log q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j$;
4. Repeat until convergence of $\text{ELBO}(q)$.

Mean Field Approximation

General scheme for Mean Field Variational Inference:

1. Initialize $q(\boldsymbol{\theta}) = \prod_{j=1}^m q_j(\boldsymbol{\theta}_j)$;
2. For each factor $q_i(\boldsymbol{\theta}_i)$, $i = 1, \dots, m$ do:
 - Calculate $q_i(\boldsymbol{\theta}_i) = \frac{1}{Z_i} \exp(\int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j)$, where $Z_i = \int \exp(\int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j) d\boldsymbol{\theta}_i$;
3. Calculate $\text{ELBO}(q) = \int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j=1}^m q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta} - \sum_{j=1}^m \int q_j(\boldsymbol{\theta}_j) \log q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j$;
4. Repeat until convergence of $\text{ELBO}(q)$.

Main assumption: we are able to calculate $\int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j$ and Z_i analytically.

Mean Field Approximation

General scheme for Mean Field Variational Inference:

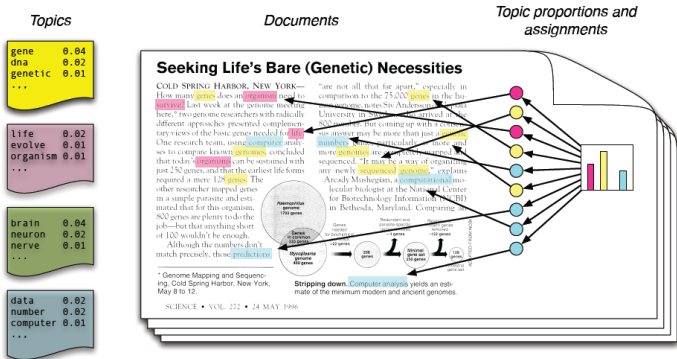
1. Initialize $q(\boldsymbol{\theta}) = \prod_{j=1}^m q_j(\boldsymbol{\theta}_j)$;
2. For each factor $q_i(\boldsymbol{\theta}_i)$, $i = 1, \dots, m$ do:
 - Calculate $q_i(\boldsymbol{\theta}_i) = \frac{1}{Z_i} \exp(\int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j)$, where $Z_i = \int \exp(\int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j) d\boldsymbol{\theta}_i$;
3. Calculate $\text{ELBO}(q) = \int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j=1}^m q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta} - \sum_{j=1}^m \int q_j(\boldsymbol{\theta}_j) \log q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j$;
4. Repeat until convergence of $\text{ELBO}(q)$.

Main assumption: we are able to calculate $\int \log p(\mathbf{x}, \boldsymbol{\theta}) \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j$ and Z_i analytically. For many probabilistic models this calculation is much simpler than $p(\mathbf{x}) = \int p(\mathbf{x}, \boldsymbol{\theta}) d\boldsymbol{\theta}$.

Example: Latent Dirichlet Allocation (LDA) [Blei, Ng & Jordan, 2003]

LDA is a tool for topic modelling. LDA assumptions:

- ▶ Each document is a collection of words regardless their order;
- ▶ Each topic is a probability distribution over set of words;
- ▶ Each document is a discrete probability mixture of several topics.



Example: Latent Dirichlet Allocation

Generative process for one document with n words:

1. Choose topic probabilities θ from prior distribution;
2. For each position $1, \dots, n$ in the document:
 - ▶ Choose topic for current position z_i using topic probabilities θ ;
 - ▶ Choose current word w_i using word probabilities Φ for current topic z_i .

Example: Latent Dirichlet Allocation

Probabilistic model for one document:

$$p(\mathbf{w}, \mathbf{z}, \boldsymbol{\theta} | \Phi, \alpha) = p(\boldsymbol{\theta} | \alpha) \prod_{i=1}^n p(z_i | \boldsymbol{\theta}) p(w_i | z_i, \Phi),$$

$$p(\boldsymbol{\theta} | \alpha) = \text{Dir}(\boldsymbol{\theta} | \alpha) \propto \prod_{t=1}^T \theta_t^{\alpha-1},$$

$$p(z_i | \boldsymbol{\theta}) = \prod_{t=1}^T \theta_t^{[z_i=t]},$$

$$p(w_i | z_i, \Phi) = \prod_{t=1}^T \prod_{w=1}^W \phi_{t,w}^{[w_i=w][z_i=t]}.$$

Variational EM-algorithm for LDA

LDA training: $\log p(\mathbf{w}|\Phi, \alpha) \geq \text{ELBO}(q, \Phi, \alpha) = \mathbb{E}_{q(\mathbf{z}, \boldsymbol{\theta})} \log \frac{p(\mathbf{w}, \mathbf{z}, \boldsymbol{\theta}|\Phi, \alpha)}{q(\mathbf{z}, \boldsymbol{\theta})} \rightarrow \max_{q(\mathbf{z}, \boldsymbol{\theta}), \Phi, \alpha}.$

Mean field approximation: $q(\mathbf{z}, \boldsymbol{\theta}) = q(\mathbf{z})q(\boldsymbol{\theta}).$

Calculation of $q(\boldsymbol{\theta})$:

$$\begin{aligned}\log q(\boldsymbol{\theta}) &= \int \log p(\mathbf{w}, \mathbf{z}, \boldsymbol{\theta}|\Phi, \alpha) q(\mathbf{z}) d\mathbf{z} + \text{const} = \\ &= \log \text{Dir}(\boldsymbol{\theta}|\alpha) + \sum_{i=1}^n \sum_{z_i} \log p(z_i|\boldsymbol{\theta}) q(z_i) + \text{const} = \\ &= \sum_{t=1}^T (\alpha - 1) \log \theta_t + \sum_{i=1}^n \sum_{z_i} [z_i = t] \log \theta_t q(z_i) + \text{const} = \\ &= \sum_{t=1}^T \log \theta_t \underbrace{\left[\alpha - 1 + \sum_{i=1}^n q(z_i = t) \right]}_{\gamma_t - 1} + \text{const} = \log \text{Dir}(\boldsymbol{\theta}|\boldsymbol{\gamma}).\end{aligned}$$

Variational Parametric Approximation

Alternative to mean field variational approximation is to use **parametric variational approximation**: $q(\boldsymbol{\theta}) = q(\boldsymbol{\theta}|\boldsymbol{\lambda})$, where $\boldsymbol{\lambda}$ are some parameters.

In this case variational inference transforms to parametric optimization problem:

$$\text{ELBO}(q, \boldsymbol{\theta}) = \text{ELBO}(\boldsymbol{\lambda}, \boldsymbol{\theta}) = \mathbb{E}_{q(\boldsymbol{\theta}|\boldsymbol{\lambda})} \log \frac{p(\mathbf{x}, \boldsymbol{\theta})}{q(\boldsymbol{\theta}|\boldsymbol{\lambda})} \rightarrow \max_{\boldsymbol{\lambda}, \boldsymbol{\theta}}. \quad (4)$$

If we're able to calculate derivatives of ELBO w.r.t. $\boldsymbol{\lambda}$ and $\boldsymbol{\theta}$, we can solve problem (4) using some numerical optimization solver.

Example: Bayesian Logistic Regression

Consider a 2-class classification problem. We have a dataset $(\mathbf{y}, X) = \{y_i, \mathbf{x}_i\}_{i=1}^n$, where $\mathbf{x}_i \in \mathbb{R}^d$ – feature vectors and $y_i \in \{-1, +1\}$ – class labels and want to train a linear classifier:

$$\hat{y}(\mathbf{x}) = \text{sign}\left(\sum_{j=1}^d w_j x_j\right) = \text{sign}(\mathbf{w}^T \mathbf{x}).$$

Probabilistic model:

$$p(\mathbf{y}, \mathbf{w} | X, \alpha) = p(\mathbf{y} | \mathbf{w}, X) p(\mathbf{w} | \alpha) = p(\mathbf{w} | \alpha) \prod_{i=1}^n p(y_i | \mathbf{w}, \mathbf{x}_i),$$

$$p(y_i | \mathbf{w}, \mathbf{x}_i) = \sigma(y_i \mathbf{w}^T \mathbf{x}_i) = \frac{1}{1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)},$$

$$p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha I).$$

Use variational EM-algorithm for training:

$$\log p(\mathbf{y} | X, \alpha) \geq \text{ELBO}(q(\mathbf{w}), \alpha) = \mathbb{E}_{q(\mathbf{w})} \log \frac{p(\mathbf{y}, \mathbf{w} | X, \alpha)}{q(\mathbf{w})} \rightarrow \max_{q(\mathbf{w}) \in \mathcal{Q}, \alpha}.$$

Example: Bayesian Logistic Regression

Parametric variational family: $q(\mathbf{w}|\boldsymbol{\mu}, \Sigma) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \Sigma)$.

$$\begin{aligned}\log p(\mathbf{y}|X, \alpha) &\geq \text{ELBO}(\boldsymbol{\mu}, \Sigma, \alpha) = \mathbb{E}_{q(\mathbf{w}|\boldsymbol{\mu}, \Sigma)} \log p(\mathbf{y}|\mathbf{w}, X) - \underbrace{KL(q(\mathbf{w}|\boldsymbol{\mu}, \Sigma) \parallel p(\mathbf{w}|\alpha))}_{\text{analytical expression}} = \\ &= \mathbb{E}_{q(\mathbf{w}|\boldsymbol{\mu}, \Sigma)} \sum_{i=1}^n \log p(y_i|\mathbf{w}, \mathbf{x}_i) - \underbrace{KL(\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \Sigma) \parallel \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha I))}_{\text{analytical expression}} = \\ &= \sum_{i=1}^n \mathbb{E}_{\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \Sigma)} \log \sigma(y_i \mathbf{w}^T \mathbf{x}_i) - KL(\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \Sigma) \parallel \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha I)) \rightarrow \max_{\boldsymbol{\mu}, \Sigma, \alpha}.\end{aligned}$$

If $\mathbf{w} \sim \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \Sigma)$, then $u_i = \mathbf{w}^T \mathbf{x}_i \sim \mathcal{N}(u_i|m_i, s_i^2)$, where $m_i = \boldsymbol{\mu}^T \mathbf{x}_i$, $s_i^2 = \mathbf{x}_i^T \Sigma \mathbf{x}_i$.
Hence, $\mathbb{E}_{\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \Sigma)} \log \sigma(y_i \mathbf{w}^T \mathbf{x}_i) = \mathbb{E}_{\mathcal{N}(u_i|m_i, s_i^2)} \log \sigma(y_i u_i)$.

Finally, if $u_i \sim \mathcal{N}(u_i|m_i, s_i^2)$, then $\xi_i \sim \mathcal{N}(\xi_i|0, 1)$ and $u_i = \xi_i s_i + m_i$. Hence,
 $\mathbb{E}_{\mathcal{N}(u_i|m_i, s_i^2)} \log \sigma(y_i u_i) = \mathbb{E}_{\mathcal{N}(\xi_i|0, 1)} \log \sigma(y_i (\xi_i s_i + m_i))$.

Example: Bayesian Logistic Regression

Final expression for ELBO:

$$\text{ELBO}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \alpha) = \sum_{i=1}^n \mathbb{E}_{\mathcal{N}(\xi_i|0,1)} \log \sigma(y_i(\xi_i s_i + m_i)) - KL(\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) || \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha I)) \rightarrow \max_{\boldsymbol{\mu}, \boldsymbol{\Sigma}, \alpha}$$
$$m_i = \boldsymbol{\mu}^T \mathbf{x}_i, s_i^2 = \mathbf{x}_i^T \boldsymbol{\Sigma} \mathbf{x}_i.$$

For solving optimization problem we need to know derivatives of ELBO w.r.t. parameters $\boldsymbol{\mu}, \boldsymbol{\Sigma}, \alpha$. These expressions can be calculated in the following way:

$$\nabla_{\boldsymbol{\mu}} \text{ELBO}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \alpha) = \sum_{i=1}^n \mathbb{E}_{\mathcal{N}(\xi_i|0,1)} \nabla_{\boldsymbol{\mu}} \log \sigma(y_i(\xi_i s_i + m_i)) - \nabla_{\boldsymbol{\mu}} KL(\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) || \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha I))$$

Expectation w.r.t. one-dimensional standard normal distribution can be calculated using Gauss-Hermite quadrature.

Conclusions

- ▶ Variational Inference approach transforms Bayesian inference problem to a certain type of optimization problem. In this sense it is usually much faster than MCMC sampling approach;
- ▶ Variational Inference requires a careful choice of variational approximation family. Depending of properties of probabilistic model this could be mean field or parametric approximation or mixture of them;
- ▶ Variational Inference can give poor posterior approximation due to restrictions in variational approximation family. Variational EM-algorithm usually gives good results for model parameters;
- ▶ Scalability issues of Variational Inference are not covered. See the next lecture!