

STK4051/9051 Computational statistics

Geir Storvik

Variational inference

- Assume observations \mathbf{x} , variables of interest \mathbf{z}
 - \mathbf{z} may include parameters and/or latent variables
- Core: Posterior distribution

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

- Inference:

$$E[h(\mathbf{z})|\mathbf{x}] = \int_{\mathbf{z}} h(\mathbf{z})p(\mathbf{z}|\mathbf{x})d\mathbf{z}$$

- Main problem: $p(\mathbf{z}|\mathbf{x})$ can be very complex
- "Standard" methods such as MCMC do not scale well with big data

- Main source: Blei et al. (2017)
- Idea:
 - Approximate $p(\mathbf{z}|\mathbf{x})$ by a simpler $q^*(\mathbf{z})$
 - Perform inference by

$$E[h(\mathbf{z})|\mathbf{x}] \approx \int_{\mathbf{z}} h(\mathbf{z})q^*(\mathbf{z})d\mathbf{z} \quad (*)$$

- Question: How to find $q^*(\mathbf{z})$?
- Approach: Define $q^*(\mathbf{z})$ as

$$q^*(\mathbf{z}) = \arg \min_{q(\mathbf{z}) \in \mathcal{Q}} \mathcal{D}(q(\mathbf{z}), p(\mathbf{z}|\mathbf{x}))$$

where

- \mathcal{Q} is some class of distributions; **variational family**
- $\mathcal{D}(\cdot, \cdot)$ is some distance measure between distributions
- Note: q^* will depend on \mathbf{x} although not explicitly shown in the notation!
- **Integration problem** now mainly transformed to an **optimization problem**
 - Integration of (*) assumed to be simple

- Most common: **Kullback-Leibler divergence**:

$$\begin{aligned}\mathcal{D}(q(\mathbf{z}), p(\mathbf{z}|\mathbf{x})) &= KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) \\ &= \int_{\mathbf{z}} q(\mathbf{z}) \log \left(\frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \right) d\mathbf{z} \\ &= E^q[\log q(\mathbf{z})] - E^q[\log p(\mathbf{z}|\mathbf{x})]\end{aligned}$$

- **Properties:**

- $KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) \geq 0$
- $KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) = 0 \Rightarrow q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x})$
- $KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) \neq KL(p(\mathbf{z}|\mathbf{x})||q(\mathbf{z}))$

- Rewriting of the KL divergence:

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

- Minimizing KL is equivalent maximizing

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

- ELBO = **Evidence lower bound**:

$$\text{ELBO}(q) = \log p(\mathbf{x}) - KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) \leq \log p(\mathbf{x})$$

- Practical formulation:

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

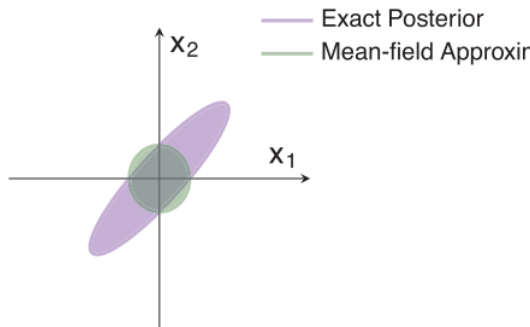
- First term: Fit to data
- Second part: Penalty term

- Common choice: **Mean-field variational family**

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

that is **independence** structure

- Also need to specify $q_j(z_j)$.
 - For continuous variables, Gaussian is possible
 - For categorical variables, any discrete distribution.
- Typically, point estimates quite accurate, uncertainty measures too low



Example - Gaussian mixtures

- Model

$$\mu_k \sim \mathcal{N}(0, \sigma^2),$$

$$k = 1, \dots, K$$

$$c_i \sim \text{Categorical}(1/K, \dots, 1/K),$$

$$i = 1, \dots, n$$

$$x_i | c_i, \boldsymbol{\mu} \sim \mathcal{N}(\mu_{c_i}, 1)$$

Of interest

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

- Variational approximation:

$$q(\boldsymbol{\mu}, \mathbf{c}) = \prod_{k=1}^K q(\mu_k; m_k, s_k^2) \prod_{i=1}^n q(c_i; \boldsymbol{\phi}_i)$$

where $\boldsymbol{\phi}_i = (\phi_{i,1}, \dots, \phi_{i,K})$, $\sum_{j=1}^K \phi_{i,j} = 1$ and

$$q(\mu_k; m_k, s_k^2) = \mathcal{N}(\mu_k; m_k, s_k^2)$$

$$q(c_i; \boldsymbol{\phi}_i) = \phi_{i,c_i}$$

- CAVI = Coordinate ascent variational inference
- Given all other variables, the optimal $q_j^*(z_j)$ is

$$\begin{aligned}q_j^*(z_j) &\propto \exp\{E_{z_j^q}^q[\log p(z_j | \mathbf{z}_{-j}, \mathbf{x})]\} \\ &\propto \exp\{E_{z_j^q}^q[\log p(z_j, \mathbf{z}_{-j}, \mathbf{x})]\}\end{aligned}$$

- Algorithm

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

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

Initialize: Variational factors $q_j(z_j)$

while the ELBO has not converged **do**

for $j \in \{1, \dots, m\}$ **do**

 Set $q_j(z_j) \propto \exp\{E_{z_j^q}^q[\log p(z_j, \mathbf{z}_{-j}, \mathbf{x})]\}$

end

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

end

- Note: Similar structure as the Gibbs sampler!

$$\begin{aligned}
 q^*(c_i; \phi_i) &\propto \exp\{E_{\mathbf{c}_{-i}, \boldsymbol{\mu}}[\log p(c_i, \mathbf{c}_{-i}, \boldsymbol{\mu}, \mathbf{x})]\} \\
 &\propto \exp\{E_{\mathbf{c}_{-i}, \boldsymbol{\mu}}[\log(p(c_i)) + \log(p(\mathbf{c}_{-i})) + \log(\boldsymbol{\mu}) + \sum_{j=1}^n \log p(x_j|c_j, \boldsymbol{\mu})]\} \\
 &\propto \exp\{E_{\mathbf{c}_{-i}, \boldsymbol{\mu}}[\log(p(c_i)) + \log p(x_i|c_i, \boldsymbol{\mu})]\} \\
 &\propto \exp\{E_{\mathbf{c}_{-i}, \boldsymbol{\mu}}[\frac{1}{K} + \log p(x_i|\mu_{c_i})]\} \\
 &\propto \exp\{E_{\boldsymbol{\mu}}[[-\frac{1}{2}(x_i - \mu_{c_i})^2]]\} \\
 &\propto \exp\{E_{\boldsymbol{\mu}}[-\frac{1}{2}(x_i^2 - 2x_i\mu_{c_i} + \mu_{c_i}^2)]\} \\
 &\propto \exp\{E_{\boldsymbol{\mu}}[x_i m_{c_i} - \frac{1}{2}s_{c_i}^2 - \frac{1}{2}m_{c_i}^2]\}
 \end{aligned}$$

$$\begin{aligned}
 q^*(\mu_k; m_k, s_k^2) &\propto \exp\{E_{\mathbf{c}, \mu_{-k}}[\log p(\mathbf{c}, \boldsymbol{\mu}, \mathbf{x})]\} \\
 &\propto \exp\{E_{\mathbf{c}, \mu_{-k}}[\log p(\mu_k) + \sum_{i=1}^n \log p(x_i | c_i, \boldsymbol{\mu})]\} \\
 &\propto \exp\{E_{\mathbf{c}, \mu_{-k}}[-\frac{1}{2\sigma^2} \mu_k^2 + \sum_{i=1}^n I(c_i = k) \log p(x_i | \mu_k)]\} \\
 &\propto \exp\{E_{\mathbf{c}, \mu_{-k}}[-\frac{1}{2\sigma^2} \mu_k^2 + \sum_{i=1}^n \phi_{ik} [-\frac{1}{2}(x_i - \mu_k)^2]]\} \\
 &\propto \exp\{E_{\mathbf{c}, \mu_{-k}}[-\frac{1}{2\sigma^2} \mu_k^2 + -\frac{1}{2} \sum_{i=1}^n \phi_{ik} [x_i^2 - 2x_i \mu_k + \mu_k^2]]\} \\
 &\propto \exp\{-\frac{1}{2}[(\frac{1}{\sigma^2} + \sum_{i=1}^n \phi_{ik})\mu_k^2 - 2 \sum_{i=1}^n \phi_{ik} x_i \mu_k]\} \\
 &\propto \exp\{-\frac{1}{2}(\frac{1}{\sigma^2} + \sum_{i=1}^n \phi_{ik})[\mu_k - \frac{\sum_{i=1}^n \phi_{ik} x_i}{\frac{1}{\sigma^2} + \sum_{i=1}^n \phi_{ik}}]^2\}
 \end{aligned}$$

Example - calculating ELBO

$$\begin{aligned} \text{ELBO}(q) &= E^q[\log p(\mathbf{z}, \mathbf{x})] - E^q[\log q(\mathbf{z})] \\ &= E^q\left[\sum_{i=1}^n \log p(c_i) + \sum_{k=1}^K \log p(\mu_k) + \sum_{i=1}^n \log p(x_i | c_i, \boldsymbol{\mu})\right] - \\ &\quad E^q\left[\sum_{i=1}^n \log q(c_i) + \sum_{k=1}^K \log q(\mu_k)\right] \\ &= E^q\left[\sum_{i=1}^n \log\left(\frac{1}{K}\right) - \frac{1}{2\sigma^2} \sum_{k=1}^K \mu_k^2 - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K I(c_i = k)(x_i - \mu_k)^2\right] - \\ &\quad E^q\left[\sum_{i=1}^n \sum_{k=1}^K \phi_{ik} \log(\phi_{ik}) + \sum_{k=1}^K \left[-\frac{1}{2} \log(s_k^2) - \frac{1}{2s_k^2} (\mu_k - m_k)^2\right]\right] \\ &= E^q\left[-\frac{1}{2\sigma^2} \sum_{k=1}^K [s_k^2 + m_j^2] - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K \phi_{ik} [x_i^2 - 2x_i\mu_k + \mu_k^2] - \right. \\ &\quad \left. E^q\left[\sum_{i=1}^n \sum_{k=1}^K \phi_{ik} \log(\phi_{ik}) + \sum_{k=1}^K \left[-\frac{1}{2} \log(s_k^2) - \frac{1}{2}\right]\right]\right] \\ &= -\frac{1}{2\sigma^2} \sum_{k=1}^K [s_k^2 + m_j^2] + \sum_{i=1}^n \sum_{k=1}^K \phi_{ik} [x_i m_k - \frac{1}{2} s_k^2 - \frac{1}{2} m_k^2] - \\ &\quad \sum_{i=1}^n \sum_{k=1}^K \phi_{ik} \log(\phi_{ik}) + \frac{1}{2} \sum_{k=1}^K \log(s_k^2) \end{aligned}$$

- D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.