

EM & Variational Bayes

Hanxiao Liu

September 9, 2014

Outline

1. EM Algorithm
 - 1.1 Introduction
 - 1.2 Example: Mixture of vMFs
2. Variational Bayes
 - 2.1 Introduction
 - 2.2 Example: Bayesian Mixture of Gaussians

MLE by Gradient Ascent

Goal: maximize $\mathcal{L}(\boldsymbol{\theta}; \mathbf{X}) = \log p(\mathbf{X}|\boldsymbol{\theta})$ w.r.t $\boldsymbol{\theta}$

Gradient Ascent (GA)

▶ One-step view: $\boldsymbol{\theta}^{t+1} \leftarrow \nabla \mathcal{L}(\boldsymbol{\theta}^t; \mathbf{X}) + \boldsymbol{\theta}^t$

▶ Two-step view:

1. $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^t) = \mathcal{L}(\boldsymbol{\theta}^t; \mathbf{X}) + (\boldsymbol{\theta}^t - \boldsymbol{\theta}) \nabla \mathcal{L}(\boldsymbol{\theta}^t; \mathbf{X}) - \frac{1}{2} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}\|_2^2$
2. $\boldsymbol{\theta}^{t+1} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^t)$

Drawbacks

1. $\nabla \mathcal{L}$ can be too complicated to work with
2. Too general to be efficient for structured problems

MLE by EM

Expectation-maximization (EM)

1. Expectation: $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^t) = \mathbb{E}_{\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^t} \mathcal{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z})$
2. Maximization: $\boldsymbol{\theta}^{t+1} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^t)$

- ▶ Replace $\underbrace{\mathcal{L}(\boldsymbol{\theta}; \mathbf{X})}_{\text{log-likelihood}}$ by $\underbrace{\mathcal{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z})}_{\text{complete log-likelihood}}$
- ▶ $\mathcal{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z})$ is a random function w.r.t \mathbf{Z}
—use the expected function as a surrogate

why EM is superior

A comparison between $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$, i.e., the local concave model

1. EM

$$\begin{aligned} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^t) &= \mathbb{E}_{\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^t} \mathcal{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z}) \\ &= \mathcal{L}(\boldsymbol{\theta}; \mathbf{X}) - D_{KL}(p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^t) \| p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})) + C \end{aligned}$$

2. GA

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^t) = \mathcal{L}(\boldsymbol{\theta}^t; \mathbf{X}) + (\boldsymbol{\theta}^t - \boldsymbol{\theta}) \nabla \mathcal{L}(\boldsymbol{\theta}^t; \mathbf{X}) - \frac{1}{2} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}\|_2^2$$

Example: vMF mixture

Notations

- ▶ $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$, $\boldsymbol{\theta} = \left\{ \boldsymbol{\pi} \in \Delta^{k-1}, \{(\boldsymbol{\mu}_i, \kappa_i)\}_{i=1}^k \right\}$
- ▶ $\mathbf{Z} = \{z_{ij} \in \{0, 1\}\}$
 - ▶ $z_{ij} = 1 \implies \mathbf{x}_i \sim$ the j -th mixture component

Log-likelihood

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{X}) = \sum_{i=1}^n \log p(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{i=1}^n \underbrace{\log \sum_{j=1}^k \pi_j \text{vMF}(\mathbf{x}_i | \boldsymbol{\mu}_j, \kappa_j)}_{\text{log sum coupling}}$$

Complete log-likelihood

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z}) = \sum_{i=1}^n \log p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{j=1}^k z_{ij} \log \left(\pi_j \text{vMF}(\mathbf{x}_i | \boldsymbol{\mu}_j, \kappa_j) \right)$$

E-step

Compute $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^t) \triangleq \mathbb{E}_{\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^t} \mathcal{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z})$

$$\begin{aligned} Q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\kappa}; \boldsymbol{\pi}^t, \boldsymbol{\mu}^t, \boldsymbol{\kappa}^t) &= \mathbb{E}_{\mathbf{Z}|\mathbf{X}, \boldsymbol{\pi}^t, \boldsymbol{\mu}^t, \boldsymbol{\kappa}^t} \sum_{i=1}^n \sum_{j=1}^k z_{ij} \log \left(\pi_j \text{vMF}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\kappa}_j) \right) \\ &= \sum_{i=1}^n \sum_{j=1}^k w_{ij}^t \log \text{vMF}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\kappa}_j) + w_{ij}^t \log \pi_j \end{aligned}$$

where

$$\begin{aligned} w_{ij}^t &= \mathbb{E}_{z_{ij}|\mathbf{X}, \boldsymbol{\pi}^t, \boldsymbol{\mu}^t, \boldsymbol{\kappa}^t} [z_{ij}] = p(z_{ij} = 1 | \mathbf{x}_i, \boldsymbol{\pi}^t, \boldsymbol{\mu}^t, \boldsymbol{\kappa}^t) \\ &= \frac{\pi_j^t \cdot \text{vMF}(\mathbf{x}_i | \boldsymbol{\mu}_j^t, \boldsymbol{\kappa}_j^t)}{\sum_{u=1}^k \pi_u^t \cdot \text{vMF}(\mathbf{x}_i | \boldsymbol{\mu}_u^t, \boldsymbol{\kappa}_u^t)} \end{aligned}$$

M-step

Maximize

$$Q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\kappa}; \boldsymbol{\pi}^t, \boldsymbol{\mu}^t, \boldsymbol{\kappa}^t) = \sum_{i=1}^n \sum_{j=1}^k w_{ij}^t \log \text{vMF}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\kappa}_j) + w_{ij}^t \log \pi_j$$

w.r.t $\boldsymbol{\pi}$, $\boldsymbol{\mu}$ and $\boldsymbol{\kappa}$ s.t. $|\boldsymbol{\pi}|_1 = 1$ and $\|\boldsymbol{\mu}_j\|_2 = 1, \forall j \in [k]$

To impose constraints, maximize

$$\tilde{Q} \triangleq Q + \lambda (1 - \boldsymbol{\pi}^\top \mathbf{1}) + \sum_{j=1}^k \nu_j (1 - \boldsymbol{\mu}_j^\top \boldsymbol{\mu}_j)$$

M-step

$$\begin{aligned}\tilde{Q}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\kappa}; \boldsymbol{\pi}^t, \boldsymbol{\mu}^t, \boldsymbol{\kappa}^t) &= \sum_{i=1}^n \sum_{j=1}^k w_{ij}^t \log \text{vMF}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\kappa}_j) + w_{ij}^t \log \pi_j \\ &\quad + \lambda (\mathbf{1} - \boldsymbol{\pi}^\top \mathbf{1}) + \sum_{j=1}^k \nu_j (1 - \boldsymbol{\mu}_j^\top \boldsymbol{\mu}_j)\end{aligned}$$

Updating π_j^t

Combining $\sum_j^k \pi_j = \sum_j^k w_{ij}^t = 1$ with

$$\partial_{\pi_j} \tilde{Q} = \frac{\sum_{i=1}^n w_{ij}^t}{\pi_j} - \lambda = 0$$

$$\implies \pi_j^{t+1} \leftarrow \frac{\sum_{i=1}^n w_{ij}^t}{n}$$

M-step

$$\begin{aligned}\tilde{Q}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\kappa}; \boldsymbol{\pi}^t, \boldsymbol{\mu}^t, \boldsymbol{\kappa}^t) &= \sum_{i=1}^n \sum_{j=1}^k w_{ij}^t \log \text{vMF}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\kappa}_j) + w_{ij}^t \log \pi_j \\ &+ \lambda (1 - \boldsymbol{\pi}^\top \mathbf{1}) + \sum_{j=1}^k \nu_j (1 - \boldsymbol{\mu}_j^\top \boldsymbol{\mu}_j)\end{aligned}$$

Updating $\boldsymbol{\mu}_j^t$

$$\log \text{vMF}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\kappa}_j) = \kappa_j \boldsymbol{\mu}_j^\top \mathbf{x}_i + C \quad (\text{w.r.t } \boldsymbol{\mu}_j)$$

$$\partial_{\boldsymbol{\mu}_j} \tilde{Q} = \kappa_j \sum_{i=1}^n w_{ij}^t \mathbf{x}_i - \nu_j \boldsymbol{\mu}_j = 0$$

$$\implies \boldsymbol{\mu}_j^{t+1} \leftarrow \frac{\mathbf{r}_j}{\|\mathbf{r}_j\|_2} \quad \text{where } \mathbf{r}_j = \sum_{i=1}^n w_{ij}^t \mathbf{x}_i$$

M-step

Updating κ_j^t

$$\blacktriangleright C_p(\kappa_j) = \frac{\kappa_j^{\frac{p}{2}-1}}{(2\pi)^{\frac{p}{2}} I_{\frac{p}{2}-1}(\kappa_j)}$$

- \blacktriangleright the recurrence property of modified Bessel function ¹

$$\partial_{\kappa_j} \log I_{\frac{p}{2}-1}(\kappa_j) = \frac{\frac{p}{2}-1}{\kappa_j} + \frac{I_{\frac{p}{2}}(\kappa_j)}{I_{\frac{p}{2}-1}(\kappa_j)}$$

$$\partial_{\kappa_j} \tilde{Q} = \sum_{i=1}^n w_{ij}^t \left(-\frac{I_{\frac{p}{2}}(\kappa_j)}{I_{\frac{p}{2}-1}(\kappa_j)} + \boldsymbol{\mu}_j^\top \mathbf{x}_i \right) = 0$$

$$\implies \frac{I_{\frac{p}{2}}(\kappa_j)}{I_{\frac{p}{2}-1}(\kappa_j)} = \bar{r}_j \implies \kappa_j^{t+1} \approx \frac{\bar{r}_j p - \bar{r}_j^3}{1 - \bar{r}_j^2} \quad [?]$$

where $\bar{r}_j = \frac{\sum_{i=1}^n w_{ij}^t \boldsymbol{\mu}_j^\top \mathbf{x}_i}{\sum_{i=1}^n w_{ij}^t}$

¹<http://functions.wolfram.com/Bessel-TypeFunctions/BesselK/introductions/Bessels/05/>

An alternative view of EM

EM - original definition

1. Expectation: $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^t) = \mathbb{E}_{\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^t} \mathcal{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z})$ why?
2. Maximization: $\boldsymbol{\theta}^{t+1} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^t)$

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}; \mathbf{X}) &= \mathbb{E}_q \log p(\mathbf{X}|\boldsymbol{\theta}) \\ &= \underbrace{\mathbb{E}_q \left[\log \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right]}_{\text{VLB}(q, \boldsymbol{\theta})} + \underbrace{\mathbb{E}_q \left[\log \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} \right]}_{D_{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}))} \end{aligned}$$

EM - coordinate ascent

1. $q^{t+1} = \operatorname{argmax}_q \text{VLB}(q, \boldsymbol{\theta}^t)$
2. $\boldsymbol{\theta}^{t+1} = \operatorname{argmax}_{\boldsymbol{\theta}} \text{VLB}(q^{t+1}, \boldsymbol{\theta})$

Show the equivalence?

Bayes Inference

Notations

- ▶ θ : hyper parameters
- ▶ \mathbf{Z} : hidden variables + random parameters

Goals

1. find a good posterior $q(\mathbf{Z}) \approx p(\mathbf{Z}|\mathbf{X}; \theta)$
2. estimate θ by Empirical Bayes, i.e., maximize $\mathcal{L}(\theta; \mathbf{X})$ w.r.t θ

$$\mathcal{L}(\theta; \mathbf{X}) = \underbrace{\mathbb{E}_q \left[\log \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \right]}_{\text{VLB}(q, \theta)} + \underbrace{\mathbb{E}_q \left[\log \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X}, \theta)} \right]}_{D_{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \theta))}$$

both goals can be achieved via the same procedure as EM

Variational Bayes Inference

One should have $q \rightarrow p(\mathbf{Z}; \mathbf{X}, \boldsymbol{\theta}^*)$ by alternating between

1. $q^{t+1} = \operatorname{argmax}_q \text{VLB}(q, \boldsymbol{\theta}^t)$
2. $\boldsymbol{\theta}^{t+1} = \operatorname{argmax}_{\boldsymbol{\theta}} \text{VLB}(q^{t+1}, \boldsymbol{\theta})$

However, we do not want q to be too complicated

- ▶ e.g., $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^t) = \mathbb{E}_q \mathcal{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z})$ can be intractable

Solution: modify the first step as

$$q^{t+1} = \operatorname{argmax}_{q \in \mathcal{Q}} \text{VLB}(q, \boldsymbol{\theta}^t)$$

\mathcal{Q} - some tractable distribution families

- ▶ Recall: without \mathcal{Q} , $q^{t+1} \equiv p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^t)$

Variational Bayes Inference

Goal: solve $\operatorname{argmax}_{q \in \mathcal{Q}} \text{VLB}(q, \theta^t)$

usually, $\mathcal{Q} = \left\{ q \mid q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i) \stackrel{\Delta}{=} \prod_{i=1}^M q_i \right\}$

Coordinate ascent

$$\begin{aligned} \text{VLB}(q_j; q_{-j}, \theta^t) &= \mathbb{E}_q \left[\log \frac{p(\mathbf{X}, \mathbf{Z}; \theta^t)}{q(\mathbf{Z})} \right] \\ &= \mathbb{E}_q \log p(\mathbf{X}, \mathbf{Z}; \theta^t) - \sum_{i=1}^M \mathbb{E}_q \log q_i \\ &= \mathbb{E}_{q_j} \left(\mathbb{E}_{q_{-j}} \log p(\mathbf{X}, \mathbf{Z}; \theta^t) \right) - \mathbb{E}_{q_j} \log q_j + C \\ &= -D_{KL} \left(\log q_j \parallel \mathbb{E}_{q_{-j}} \log p(\mathbf{X}, \mathbf{Z}; \theta^t) \right) + C \end{aligned}$$

$$\log q_j^* = \mathbb{E}_{q_{-j}} \log p(\mathbf{X}, \mathbf{Z}; \theta^t)$$

Example: Bayes Mixture of Gaussians

Consider putting a prior over the means in GM ²

- ▶ For $k = 1, 2 \dots K$, $\mu_k \sim \mathcal{N}(0, \tau^2)$
- ▶ For $i = 1, 2 \dots N$
 1. $z_i \sim \text{Mult}(\boldsymbol{\pi})$
 2. $x_i \sim \mathcal{N}(\mu_{z_i}, \sigma^2)$

$$\begin{aligned} p(\mathbf{z}, \boldsymbol{\mu} | \mathbf{X}) &= \frac{p(\mathbf{X} | \mathbf{z}, \boldsymbol{\mu}) p(\mathbf{z}) p(\boldsymbol{\mu})}{p(\mathbf{X})} \\ &= \frac{\prod_{i=1}^N p(z_i) p(x_i | z_i, \boldsymbol{\mu}) \prod_{k=1}^K p(\mu_k)}{\int \sum_{\mathbf{z}} \prod_{i=1}^N p(z_i) p(x_i | z_i, \boldsymbol{\mu}) \prod_{k=1}^K p(\mu_k) d\boldsymbol{\mu}} \\ q(\mathbf{z}, \boldsymbol{\mu}) &= \prod_{i=1}^N q(z_i; \boldsymbol{\phi}_i) \prod_{k=1}^K q(\mu_k; \tilde{\mu}_k, \tilde{\sigma}_k^2) \end{aligned}$$

²<https://www.cs.princeton.edu/courses/archive/fall11/cos597C/lectures/variational-inference-i.pdf>

Example: Bayes Mixture of Gaussians

$$\begin{aligned}\log q^*(z_j) &= \mathbb{E}_{q \setminus z_j} \log p(\mathbf{z}, \boldsymbol{\mu}, \mathbf{X}) \\ &= \mathbb{E}_{q \setminus z_j} \left(\sum_{i=1}^N \log p(z_i) + \log p(x_i | z_i, \boldsymbol{\mu}) + \sum_{k=1}^K \log p(\mu_k) \right) \\ &= \log p(z_j) + \mathbb{E}_{q(\mu_{z_j})} \log p(x_j | z_j, \mu_{z_j}) + C \\ &= \log \pi_{z_j} + x_j \underbrace{\mathbb{E}_{q(\mu_{z_j})} [\mu_{z_j}]}_{\tilde{\mu}_{z_j}} - \frac{1}{2} \underbrace{\mathbb{E}_{q(\mu_{z_j})} [\mu_{z_j}^2]}_{\tilde{\mu}_{z_j}^2 + \tilde{\sigma}_{z_j}^2} + C\end{aligned}$$

By observation $q^*(z_j) \sim \text{Mult}$, we can update ϕ_j accordingly

Example: Bayes Mixture of Gaussians

$$\begin{aligned}\log q^*(\mu_j) &= \mathbb{E}_{q \setminus \mu_j} \log p(\mathbf{z}, \boldsymbol{\mu}, \mathbf{X}) \\ &= \mathbb{E}_{q \setminus \mu_j} \left(\sum_{i=1}^N \log p(z_i) + \log p(x_i | z_i, \mu_{z_i}) + \sum_{k=1}^K \log p(\mu_k) \right) \\ &= \mathbb{E}_{q \setminus \mu_j} \sum_{i=1}^N \sum_{k=1}^K \delta_{z_i=k} \log \mathcal{N}(x_i | \mu_k) + \log p(\mu_j) + C \\ &= \sum_{i=1}^N \underbrace{\mathbb{E}_{z_i} [\delta_{z_i=j}]}_{\phi_i^j} \log \mathcal{N}(x_i | \mu_j) + \log p(\mu_j) + C\end{aligned}$$

Observing that $q^*(\mu_j) \sim \mathcal{N}$, $\tilde{\mu}_j$ and $\tilde{\sigma}_j^2$ can be updated accordingly

Stay tuned

Next topics

- ▶ LDA (Wanli)
- ▶ Bayes vMF