

Expectation Propagation

Thomas P Minka

Kuhwan Jeong¹

¹Department of Statistics, Seoul National University, South Korea

December, 2016

Example 1

- Given a parameter θ , observations $\mathbf{x}_1, \dots, \mathbf{x}_n$ are independently sampled from

$$p(\mathbf{x}|\theta) = (1 - w)\mathcal{N}(\mathbf{x}; \theta, \mathbf{I}) + w\mathcal{N}(\mathbf{x}; \mathbf{0}, 10\mathbf{I}).$$

- We use a Gaussian prior for θ :

$$\theta \sim \mathcal{N}(\theta; \mathbf{0}, 100\mathbf{I}_d).$$

- The joint distribution of θ and $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is

$$p(D, \theta) = p(\theta) \prod_i p(\mathbf{x}_i|\theta).$$

- We want to know $p(\theta|D)$.

Assumed-density filtering

$$\begin{aligned} p(\theta|D) &\propto p(\theta)p(x_1|\theta)p(x_2|\theta)p(x_3|\theta)\cdots p(x_n|\theta) \\ &\propto p(\theta|x_1)p(x_2|\theta)p(x_3|\theta)\cdots p(x_n|\theta) \\ &\propto p(\theta|x_1, x_2)p(x_3|\theta)\cdots p(x_n|\theta) \\ &\vdots \\ &\propto p(\theta|x_1, x_2, \dots, x_{n-1})p(x_n|\theta) \end{aligned}$$

$$\begin{aligned} p(\theta|D) &\approx q_1(\theta)p(x_2|\theta)p(x_3|\theta)\cdots p(x_n|\theta) \\ &\approx q_2(\theta)p(x_3|\theta)\cdots p(x_n|\theta) \\ &\vdots \\ &\approx q_{n-1}(\theta)p(x_n|\theta) \end{aligned}$$

Assumed-density filtering

1. Initialize with $q(\theta) = p(\theta)$.
2. For $i = 1, \dots, n$,
 - ① Define $\hat{q}(\theta) = \frac{p(\mathbf{x}_i|\theta)q(\theta)}{\int p(\mathbf{x}_i|\theta)q(\theta)d\theta}$.
 - ② Find $q^{new}(\theta) \in \mathcal{H}$ minimizing the KL-divergence

$$D(\hat{q}(\theta)||q^{new}(\theta)).$$

- Here \mathcal{H} is a class of distributions.
- After the i^{th} iteration, $q^{new}(\theta)$ approximates the posterior distribution given data $\{\mathbf{x}_1, \dots, \mathbf{x}_i\}$ with a prior $p(\theta)$.

Assumed-density filtering

- In the example, let

$$q(\theta) = \mathcal{N}(\theta; \mathbf{m}_\theta, v_\theta \mathbf{I}),$$
$$q^{new}(\theta) = \mathcal{N}(\theta; \mathbf{m}_\theta^{new}, v_\theta^{new} \mathbf{I}).$$

- Zeroing the gradient of $D(\hat{p}(\theta) || q^{new}(\theta))$ gives the conditions

$$\mathbf{m}_\theta^{new} = \int \hat{p}(\theta) \theta d\theta,$$
$$v_\theta^{new} d + (\mathbf{m}_\theta^{new})^T (\mathbf{m}_\theta^{new}) = \int \hat{p}(\theta) \theta^T \theta d\theta,$$

or in other words, expectation constraints:

$$\mathbb{E}_{q^{new}}[\theta] = \mathbb{E}_{\hat{p}}[\theta],$$
$$\mathbb{E}_{q^{new}}[\theta^T \theta] = \mathbb{E}_{\hat{p}}[\theta^T \theta].$$

Assumed-density filtering

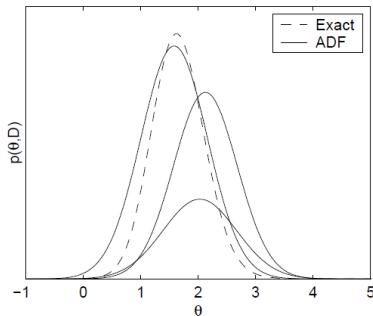
1. Initialize $\mathbf{m}_\theta = \mathbf{0}$, $v_\theta = 100$ (the prior).
2. For $i = 1, \dots, n$,

$$r_i = \frac{(1-w)\mathcal{N}(\mathbf{x}_i; \mathbf{m}_\theta, (v_\theta+1)\mathbf{I})}{(1-w)\mathcal{N}(\mathbf{x}_i; \mathbf{m}_\theta, (v_\theta+1)\mathbf{I}) + w\mathcal{N}(\mathbf{x}_i; \mathbf{0}, 10\mathbf{I})},$$

$$\mathbf{m}_\theta^{new} = \mathbf{m}_\theta + v_\theta r_i \frac{\mathbf{x}_i - \mathbf{m}_\theta}{v_\theta + 1},$$

$$v_\theta^{new} = v_\theta - r_i \frac{v_\theta^2}{v_\theta + 1} + r_i(1-r_i) \frac{v_\theta^2 (\mathbf{x}_i - \mathbf{m}_\theta)^T (\mathbf{x}_i - \mathbf{m}_\theta)}{d(v_\theta + 1)^2}.$$

Assumed-density filtering



- ADF depends on the order in which data is processed.
- No theory is available for how ADF varies with ordering.
- The error increases whenever similar data points are processed together.
- Processing the data in sorted order is especially bad.

Assumed-density filtering

1. Initialize with $q(\theta) = p(\theta)$.
2. For $i = 1, \dots, n$,
 - ① Define $\hat{q}(\theta) = \frac{p(\mathbf{x}_i|\theta)q(\theta)}{\int p(\mathbf{x}_i|\theta)q(\theta)d\theta}$.
 - ② Find $q^{new}(\theta) \in \mathcal{H}$ minimizing the KL-divergence

$$D(\hat{q}(\theta)||q^{new}(\theta)).$$

- ③ Define

$$\tilde{p}(\mathbf{x}_i|\theta) \propto \frac{q^{new}(\theta)}{q(\theta)}.$$

3. Then

$$q(\theta) \propto p(\theta) \prod_i \tilde{p}(\mathbf{x}_i|\theta).$$

Expectation Propagation

1. Initialize with $q(\theta) = p(\theta)$ and $\tilde{p}(\mathbf{x}_i|\theta) = 1$.

2. Until all $\tilde{p}(\mathbf{x}_i|\theta)$ converge :

① Choose a $\tilde{p}(\mathbf{x}_i|\theta)$ to refine.

② Define

$$q^{-i}(\theta) \propto q(\theta)/\tilde{p}(\mathbf{x}_i|\theta)$$

and

$$\hat{q}(\theta) = \frac{p(\mathbf{x}_i|\theta)q^{-i}(\theta)}{\int p(\mathbf{x}_i|\theta)q^{-i}(\theta)d\theta}.$$

③ Find $q^{new}(\theta) \in \mathcal{H}$ minimizing the KL-divergence

$$D(\hat{q}(\theta)||q^{new}(\theta)).$$

④ Set

$$\tilde{p}(\mathbf{x}_i|\theta) \propto \frac{q^{new}(\theta)}{q^{-i}(\theta)}.$$

- Here $q^{-i}(\theta)$ approximates the posterior distribution given data $D \setminus \{\mathbf{x}_i\}$ with a prior $p(\theta)$.

Example 2

- There are n documents.
- There are W distinct words $(1, \dots, W)$.
- n_{dw} is the number of word w in a document d .
- There are K topics ϕ_1, \dots, ϕ_K where $\phi_k = (\phi_{k1}, \dots, \phi_{kW})$.
- $\theta_d = (\theta_{d1}, \dots, \theta_{dK})$ is a topic proportion of a document d .

Example 2

- The probability of a document d is

$$\begin{aligned} p(d|\alpha, \phi) &= \int \mathcal{D}(\theta|\alpha) \prod_{w=1}^W \left(\sum_k \theta_k \phi_{kw} \right)^{n_{dw}} d\theta, \\ &= \int \mathcal{D}(\theta|\alpha) \prod_{w=1}^W p(w|\theta)^{n_{dw}} d\theta. \end{aligned}$$

- EP approximates each term $p(w|\theta)$ by a simpler term

$$\tilde{p}(w|\theta) = \prod_k \theta_k^{\beta_{dwk}},$$

giving

$$q_d(\theta) = \mathcal{D}(\theta|\alpha) \prod_{w=1}^W \tilde{p}(w|\theta)^{n_{dw}} = \mathcal{D}(\theta|\gamma_d)$$

where $\gamma_{dk} = \alpha_k + \sum_w n_{dw} \beta_{dwk}$.

Expectation Propagation

1. Initialize with $\gamma_d = \alpha$ and $\beta_{dwk} = 0$
($q_d(\theta) = \mathcal{D}(\theta|\alpha), \tilde{p}(w|\theta) = 1$).

2. Until all $\tilde{p}(w|\theta)$ converge :

- ① Choose a $\tilde{p}(w|\theta)$ to refine.
- ② Define

$$\gamma_{dk}^{-w} = \gamma_{dk} - n_{dw}\beta_{dwk}.$$

and

$$\hat{q}(\theta) = \frac{(\sum_k \theta_k \phi_{kw})^{n_{dw}} \mathcal{D}(\theta|\gamma_d^{-w})}{\int (\sum_k \theta_k \phi_{kw})^{n_{dw}} \mathcal{D}(\theta|\gamma_d^{-w}) d\theta}.$$

- ③ Find $q_d^{new}(\theta) = \mathcal{D}(\theta|\gamma_d)$ by matching the mean and variance of $q_d^{new}(\theta)$ against those of $\hat{q}(\theta)$.
- ④ Set

$$\beta_{dwk} = \frac{\gamma_{dk} - \gamma_{dk}^{-w}}{n_{dw}}.$$

Estimation of ϕ

- Maximize the following lower bound to the log-likelihood:

$$\begin{aligned}\sum_{d=1}^n \log p(d|\alpha, \phi) &\geq \sum_d \int q_d(\boldsymbol{\theta}) \log \left\{ \mathcal{D}(\boldsymbol{\theta}|\alpha) \prod_w \left(\sum_k \theta_k \phi_{kw} \right)^{n_{dw}} \right\} d\boldsymbol{\theta} \\ &\quad - \sum_d \int q_d(\boldsymbol{\theta}) \log q_d(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \sum_{dw} n_{dw} \int q_d(\boldsymbol{\theta}) \log \left(\sum_k \theta_k \phi_{kw} \right) d\boldsymbol{\theta} + \text{const.}\end{aligned}$$

- By zeroing the derivative with respect to ϕ_{kw} , we obtain

$$\phi_{kw}^{\text{new}} \propto \sum_d n_{dw} \int q_d(\boldsymbol{\theta}) \frac{\theta_k \phi_{kw}}{\sum_k \theta_k \phi_{kw}} d\boldsymbol{\theta}.$$

Experimental Results

- The models are compared quantitatively using test perplexity

$$\exp \left(- \frac{\sum_d \log p(d)}{\sum_d n_d} \right).$$

