

PRML chapter10

1. Variational Inference

- aim : evaluate $p(\mathbf{Z}|\mathbf{X})$
- optimize $q(\mathbf{Z})$, means maximize $L(q)$ or minimize $KL(q||p)$
- decompose the log marginal probability using

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + KL(q||p),$$

under

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right\} d\mathbf{Z}$$

$$KL(q||p) = - \int q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} \right\} d\mathbf{Z}$$

- the parameters are now stochastic variables
- consider a restricted family of distributions $q(\mathbf{Z})$ and then seek the member of this family for which the KL divergence is minimized

(a) Factorized distributions

- Suppose that $q(\mathbf{Z})$ can be decomposed as below

$$q(\mathbf{Z}) = \prod_{i=1}^M q_i(Z_i)$$

- maximize $L(q)$

$$\begin{aligned} \mathcal{L}(q) &= \int \prod_i q_i \left\{ \ln p(\mathbf{X}, \mathbf{Z}) - \sum_i \ln q_i \right\} d\mathbf{Z} = \int q_j \left\{ \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i d\mathbf{Z}_i \right\} d\mathbf{Z}_j - \int q_j \ln q_j d\mathbf{Z}_j + \text{const} \\ &= \int q_j \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j - \int q_j \ln q_j d\mathbf{Z}_j + \text{const} \end{aligned}$$

- it is Kullback-Leibler divergence so we could gain a general expression for the optimal solution

$$\ln q_j^*(\mathbf{Z}_j) = E_{i \neq j}[\ln p(\mathbf{X}, \mathbf{Z})] + \text{const.}$$

(b) Properties of factorized approximations

- consider the problem of approximating a general distribution by a factorized distribution.
- show the difference between $KL(p||q)$ and $KL(q||p)$
- The former avoids the region that $p(\mathbf{Z})$ is low, otherwise the latter tries to cover the region that $p(\mathbf{Z})$ is not zero.

(c) Example: The univariate Gaussian

- illustrate the factorized variational approximation using a Gaussian distribution over a single variable x
- likelihood is

$$p(D|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{\frac{N}{2}} \exp\left\{-\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2\right\}$$

- prior distribution is

$$p(\mu|\tau) = N(\mu|\mu_0, (\lambda_0\tau)^{-1})p(\tau) = \text{Gam}(\tau|a_0, b_0)$$

- we gain

$$\ln q_\mu^*(\mu) = E_\tau[\ln p(D|\mu, \tau) + \ln p(\mu|\tau)] + \text{const} = -\frac{E[\tau]}{2} \left\{ \lambda_0 (\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right\} + \text{const.}$$

$$\begin{aligned} \ln q_\tau^*(\tau) &= E_\mu[\ln p(D|\mu, \tau) + \ln p(\mu|\tau)] + \ln p(\tau) + \text{const} \\ &= (a_0 - 1) \ln \tau - b_0 \tau + \frac{N}{2} \ln \tau - \frac{\tau}{2} E_\mu \left[\sum_{n=1}^N (x_n - \mu)^2 + \lambda_0 (\mu - \mu_0)^2 \right] + \text{const} \end{aligned}$$

2. Illustration: Variational Mixture of Gaussians

- later come back this session

3. Variational Linear Regression

- α 's prior distribution is gamma distribution and

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi_n, \beta^{-1}), p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$$

(a) Variational distribution

- posterior distribution is expressed by the factorized expression

$$q(\mathbf{w}, \alpha) = q(\mathbf{w})q(\alpha)$$

- of course we soon gain

$$\begin{aligned} \ln q^*(\alpha) &= \ln p(\alpha) + E_{\mathbf{w}}[\ln p(\mathbf{w}|\alpha)] + \text{const} \\ &= (a_0 - 1) \ln \alpha - b_0 \alpha + \frac{M}{2} \ln \alpha - \frac{\alpha}{2} E[\mathbf{w}^T \mathbf{w}] + \text{const.} \end{aligned}$$

- and

$$\begin{aligned} \ln q^*(\mathbf{w}) &= \ln p(\mathbf{t}|\mathbf{w}) + E_\alpha[\ln p(\mathbf{w}|\alpha)] + \text{const} \\ &= -\frac{\beta}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi_n - t_n\}^2 - \frac{1}{2} E[\alpha] \mathbf{w}^T \mathbf{w} + \text{const} \\ &= -\frac{1}{2} \mathbf{w}^T (E[\alpha] \mathbf{I} + \beta \Phi^T \Phi) \mathbf{w} + \beta \mathbf{w}^T \Phi^T \mathbf{t} + \text{const.} \end{aligned}$$

- The evaluation of the variational posterior distribution begins by initializing the parameters of one of the distributions $q(\mathbf{w})$ or $q(\alpha)$, and then alternately re-estimates these factors in turn until a suitable convergence criterion is satisfied

(b) Predictive distribution

- predictive distribution is calculated easily by

$$p(t|\mathbf{x}, \mathbf{t}) = \int p(t|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathbf{t})d\mathbf{w} \simeq \int p(t|\mathbf{x}, \mathbf{w})q(\mathbf{w})d\mathbf{w}$$

(c) lower bound

- another important quantity

$$\begin{aligned} \mathcal{L}(q) &= E[\ln p(\mathbf{w}, \alpha, \mathbf{t})] - E[\ln q(\mathbf{w}, \alpha)] \\ &= E_{\mathbf{w}}[\ln p(\mathbf{t}|\mathbf{w})] + E_{\mathbf{w}, \alpha}[\ln p(\mathbf{w}|\alpha)] + E_{\alpha}[\ln p(\alpha)] - E_{\alpha}[\ln q(\mathbf{w})]_{\mathbf{w}} - E[\ln q(\alpha)] \end{aligned}$$

4. Exponential Family Distributions

- make a further distinction between latent variables and parameters
- let joint distribution be

$$p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\eta}) = \prod_{n=1}^N h(\mathbf{x}_n, \mathbf{z}_n) g(\boldsymbol{\eta}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}_n, \mathbf{z}_n) \}$$

- prior distribution of $\boldsymbol{\eta}$ is

$$p(\boldsymbol{\eta}|\nu_0, \boldsymbol{\nu}_0) = f(\nu_0, \chi_0) g(\boldsymbol{\eta})^{\nu_0} \exp \{ \nu_0 \boldsymbol{\eta}^T \chi_0 \}$$

- we gain solution as

$$\begin{aligned} \ln q^*(\boldsymbol{\eta}) &= \ln p(\boldsymbol{\eta}|\nu_0, \chi_0) + E_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\eta})] + \text{const} \\ \ln q^*(\mathbf{Z}) &= E_{\boldsymbol{\eta}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\eta})] + \text{const} \end{aligned}$$

- they are dependent on each other so use EM algorithms

5. Local Variational Methods

- approximate the convex function $f(x)$ by a linear function $y(x, \lambda)$

$$f(x) = \max_{\lambda} \{ \lambda x - \lambda + \lambda \ln(-\lambda) \}$$

- more generally

$$f(x) = \max_{\lambda} \{ \lambda x - g(\lambda) \}$$

$$g(\lambda) = \max_x \{ \lambda x - f(x) \}$$

6. Variational Logistic Regression

(a) Variational posterior distribution

- In Bayesian logistic regression model,

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{w})p(\mathbf{w})d\mathbf{w} = \int \left[\prod_{n=1}^N p(t_n|\mathbf{w}) \right] p(\mathbf{w})d\mathbf{w}$$

- in last session, we talked about the variational lower bound on the logistic sigmoid function.

- use it for $p(t|\mathbf{w})$, then

$$p(t|\mathbf{w}) = e^{at} \sigma(-a) \geq e^{at} \sigma(\xi) \exp \left\{ -(a + \xi)/2 - \lambda(\xi) (a^2 - \xi^2) \right\}$$

- let ξ be a variational parameter,

$$p(\mathbf{t}, \mathbf{w}) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w})h(\mathbf{w}, \boldsymbol{\xi})p(\mathbf{w})$$

where

$$h(\mathbf{w}, \boldsymbol{\xi}) = \prod_{n=1}^N \sigma(\xi_n) \exp \left\{ \mathbf{w}^T \boldsymbol{\phi}_n t_n - (\mathbf{w}^T \boldsymbol{\phi}_n + \xi_n) / 2 - \lambda(\xi_n) \left([\mathbf{w}^T \boldsymbol{\phi}_n]^2 - \xi_n^2 \right) \right\}$$

- Since log function is monotonically increasing

$$\ln \{ p(\mathbf{t}|\mathbf{w})p(\mathbf{w}) \} \geq \ln p(\mathbf{w}) + \sum_{n=1}^N \left\{ \ln \sigma(\xi_n) + \mathbf{w}^T \boldsymbol{\phi}_n t_n - (\mathbf{w}^T \boldsymbol{\phi}_n + \xi_n) / 2 - \lambda(\xi_n) \left([\mathbf{w}^T \boldsymbol{\phi}_n]^2 - \xi_n^2 \right) \right\}$$

- substitute $p(\mathbf{w})$, right side of the inequality becomes

$$-\frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N \left\{ \mathbf{w}^T \boldsymbol{\phi}_n (t_n - 1/2) - \lambda(\xi_n) \mathbf{w}^T (\boldsymbol{\phi}_n \boldsymbol{\phi}_n^T) \mathbf{w} \right\} + \text{const.}$$

- this is quadratic function of \mathbf{w} ,

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

- we shall use it shortly to evaluate the predictive distribution for new data points

(b) Optimizing the variational parameters

- how to estimate the variational parameter ξ
- Two ways to achieve this goal, one is to use EM algorithm the other is to integrate analytically and perform a direct maximization

7. Expectation Propagation

- we have a joint distribution

$$p(\mathcal{D}, \boldsymbol{\theta}) = \prod_i f_i(\boldsymbol{\theta})$$

, and want to approximate the posterior distribution $p(\boldsymbol{\theta}|\mathcal{D})$ by

$$q(\boldsymbol{\theta}) = \frac{1}{Z} \prod_i \tilde{f}_i(\boldsymbol{\theta})$$

- we also want to approximate the model evidence $p(\mathcal{D})$
- First, we initialize the approximating factor $\tilde{f}_i(\boldsymbol{\theta})$
- and initialize posterior approximation $q(\boldsymbol{\theta}) \propto \prod_i \tilde{f}_i(\boldsymbol{\theta})$
- Then, choose a factor $\tilde{f}_j(\boldsymbol{\theta})$ that we want to refine
- remove it from the posterior distribution

$$q^{\setminus j}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_j(\boldsymbol{\theta})}$$

- calculate $q^{new}(\boldsymbol{\theta})$ by equaling to $q^{\setminus j}(\boldsymbol{\theta})\tilde{f}_j(\boldsymbol{\theta})$

- normalization constant is

$$Z_j = \int q^{\setminus j}(\boldsymbol{\theta}) f_j(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

- evaluate the new factor

$$\tilde{f}_j(\boldsymbol{\theta}) = Z_j \frac{q^{new}(\boldsymbol{\theta})}{q^{\setminus j}(\boldsymbol{\theta})}$$

- evaluate the model evidence

$$p(\mathcal{D}) \simeq \int \prod_i \tilde{f}_i(\boldsymbol{\theta}) d\boldsymbol{\theta}$$