MAST90125: Bayesian Statistical learning

Lecture 21: Expectation propagation

Feng Liu and Guoqi Qian



What have we covered so far

- In the last lecture, we discussed one method for approximate Bayesian inference, Variational Bayes. This method was based on partitioning the parameter vector $\boldsymbol{\theta}$ into sub-vectors $\boldsymbol{\theta}_1, \dots \boldsymbol{\theta}_K$ and finding approximate posteriors $Q(\boldsymbol{\theta}_i)$; $i=1,\dots,K$ for each sub-vector.
- ▶ In this lecture, we will discuss another method for approximate Bayesian inference, Expectation propagation. This method is based on partitioning the data, rather than the parameters.

Partitioning the posterior

In expectation propagation, it is assumed that $p(\theta|y)$ can be factorised as,

$$p(\boldsymbol{\theta}|y_1,\ldots y_n)=\prod_{i=0}^n f_i(\boldsymbol{\theta}),$$

where *i* represents a datapoint.

- ▶ Question: How would you choose $f_i(\theta)$?
 - In this course, we have typically assumed that conditional on θ , observations are *i.i.d.* This means that.

$$p(\theta|y_1,\ldots y_n) = \frac{p(y_1,\ldots y_n|\theta)p(\theta)}{p(y_1,\ldots y_n)} = \frac{p(\theta)\prod_{i=1}^n p(y_i|\theta)}{p(y_1,\ldots y_n)}.$$

More importantly, for Bayesian computing it is usually sufficient to work with the known un-normalised density. Therefore a natural choice for $f_i(\theta)$ is $p(\theta)$ if i = 0 and $p(y_i|\theta)$ if i = 1, ..., n.

Approximating the posterior

- As we have been discovering since lecture 8, it is often not practical to determine the exact posterior analytically. So, does working with $f_i(\theta)$ as defined on the previous slide, and nothing else, allow us to determine the exact posteriors?
 - ► In general, no.
- ightharpoonup To get around this, in expectation propagation we define g(heta) such that

$$g(\theta) = \prod_{i=0}^n g_i(\theta).$$

▶ We then estimate $g_i(\theta)$ such that $g_i(\theta) \approx f_i(\theta) \Rightarrow g(\theta) \approx p(\theta|y_1,...,n)$.

How do we learn $g_i(\theta)$?

▶ To determine $g_i(\theta)$, we need to define a *cavity distribution*

$$g_{-i}(oldsymbol{ heta}) \propto rac{g(oldsymbol{ heta})}{g_i(oldsymbol{ heta})}$$

and tilted distribution assumed proportional to

$$g_{-i}(\theta)f_i(\theta)$$
.

- ▶ Having defined these distributions, the approximation of $g_{-i}(\theta)f_i(\theta)$ is the new estimate of $g(\theta)$. Then the new estimate of $g_i(\theta)$ is defined as $g(\theta)/g_{-i}(\theta)$.
- ► This leaves us with unanswered questions:
 - ▶ How do we measure that the distribution $g(\theta)$ approximates $g_{-i}(\theta)f_i(\theta)$?
 - ▶ How many choices do we have over the distributional form of $g(\theta)$?

How do we define distribution similarity?

- ▶ We have already encountered the Kullback-Leibler divergence many times.
- ▶ In expectation propagation, the divergence we are considering is

$$D_{\mathsf{KL}}\{c\mathsf{g}_{-i}(\boldsymbol{\theta})\mathsf{f}_i(\boldsymbol{\theta})||\mathsf{g}(\boldsymbol{\theta})\} = \int c\mathsf{g}_{-i}(\boldsymbol{\theta})\mathsf{f}_i(\boldsymbol{\theta})\{\log(c\mathsf{g}_{-i}(\boldsymbol{\theta})\mathsf{f}_i(\boldsymbol{\theta})) - \log(\mathsf{g}(\boldsymbol{\theta}))\}d\boldsymbol{\theta}.$$

- As D_{KL} is a strictly non-negative measure such that $D_{KL}(g||g) = 0$, we want to minimise this function.
 - Note the order of approximate and exact distributions has been reversed to that of D_{KL} in variational Bayes.
 - D_{KL} as used in expectation propagation is data-point specific. This means the final result of expectation propagation is not necessarily a global (i.e. over all data-points/prior) approximation to the posterior.

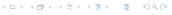
Choosing the distributional form of $g(\theta)$

- In expectation propagation, we need to choose the distributional form of the approximating distribution, $g(\theta)$.
 - Note: This is unlike variational Bayes, where the distributional form of $Q(\theta_i)$ was dictated by the distributional form of the joint distribution.
- We will assume that $g(\theta)$ is a member of the exponential family of distributions. Hence we can write,

$$g(\theta|\eta) = h(\theta)a(\eta)e^{\eta'u(\theta)} = h(\theta)e^{\eta'u(\theta)-A(\eta)},$$

where $A(\eta) = -\log(a(\eta))$, η is a vector of natural parameters, $u(\theta)$ is a vector of sufficient statistics with a first derivative satisfying

$$\frac{dg(\theta|\eta)}{d\eta} = (u(\theta) - A'(\eta))h(\theta)e^{\eta'u(\theta) - A(\eta)} = (u(\theta) - A'(\eta))g(\theta|\eta).$$



Minimising $D_{KL}\{g_{-i}(\boldsymbol{\theta})f_i(\boldsymbol{\theta})||g(\boldsymbol{\theta})\}$

▶ To minimise the KL divergence, we need to solve $dD_{KL}/d\eta = 0$.

$$0 = \frac{dD_{KL}\{g_{-i}(\theta)f_i(\theta)||g(\theta|\eta)\}}{d\eta} = \frac{d}{d\eta} \int cg_{-i}(\theta)f_i(\theta)\{\log(cg_{-i}(\theta)f_i(\theta)) - \log(g(\theta|\eta))\}d\theta$$
$$= c \int g_{-i}(\theta)f_i(\theta) \frac{d\{\log(cg_{-i}(\theta)f_i(\theta)) - \log(g(\theta|\eta))\}}{d\eta}d\theta$$
$$= -c \int g_{-i}(\theta)f_i(\theta)(u(\theta) - A'(\eta))d\theta$$

$$A'(\eta) = E_{cg_{-i}(\theta)f_i(\theta)}(u(\theta))$$

► In addition, we know that

$$\int \frac{dg(\theta|\eta)}{d\eta}d\theta = \frac{d}{d\eta}\int g(\theta|\eta)d\theta = \frac{d}{d\eta}1 = 0.$$



Minimising $D_{KI}\{g_{-i}(\theta)f_i(\theta)||g(\theta)\}$

At the same time, we know that

$$\int \frac{dg(\theta|\eta)}{d\eta}d\theta = \int (u(\theta) - A'(\eta))g(\theta|\eta)d\theta = E_{\theta|\eta}(u(\theta)) - A'(\eta).$$

► Combining the two results for $\int \frac{dg(\theta|\eta)}{d\eta} d\theta$, and the result for D_{KL} we find that

$$E_{\theta|\eta}(u(\theta)) = A'(\eta) \tag{1}$$

$$E_{c\sigma_{-i}(\theta)f_{i}(\theta)}(u(\theta)) = A'(\eta) \tag{2}$$

$$f_{cg_{-i}(\boldsymbol{\theta})f_i(\boldsymbol{\theta})}(u(\boldsymbol{\theta})) = A'(\boldsymbol{\eta})$$
 (2)

Aims of the algorithm

► Combining (1) and (2) we get the equality

$$E_{cg_{-i}(\theta)f_i(\theta)}(u(\theta)) = E_{\theta|\eta}(u(\theta)).$$

Thus determining $g(\theta)$ requires iterative moment matching until convergence.

Before writing up the algorithm, let's choose a specific distribution for $g(\theta)$ with desirable properties. Usually, $g(\theta)$ is assumed multivariate normal, $\mathcal{N}(\mu, \Sigma)$, with natural parameters $\Sigma^{-1}, \Sigma^{-1}\mu$. This implies that $g_i(\theta) = \mathcal{N}(\mu_i, \Sigma_i)$. If we match the kernels of $g(\theta)$ and $\prod_{i=0}^n g_i(\theta)$,

$$e^{-\frac{\theta'\Sigma^{-1}\theta}{2}}e^{\theta'\Sigma^{-1}\mu}=\prod_{i=0}^{n}e^{-\frac{\theta'\Sigma_{i}^{-1}\theta}{2}}e^{\theta'\Sigma_{i}^{-1}\mu_{i}}=e^{-\frac{\theta'\sum_{i=0}^{n}\Sigma_{i}^{-1}\theta}{2}}e^{\theta'\sum_{i=0}^{n}\Sigma_{i}^{-1}\mu_{i}},$$

we find
$$\Sigma^{-1}=\sum_{i=0}^n \Sigma_i^{-1}$$
 and $\Sigma^{-1}\mu=\sum_{i=0}^n \Sigma_i^{-1}\mu_i$.



Expectation propagation for a regression with known error variance.

- ► This algorithm assumes
 - $f_i(\beta) = p(y_i|\theta_i) = p(y_i|\mathbf{x}_i'\beta) = d(y_i, \eta^{-1}(\mathbf{x}_i'\beta))$, where d is some arbitrary distribution and η is a link function.
 - $ightharpoonup g(oldsymbol{eta}) = \mathcal{N}(oldsymbol{\mu}, oldsymbol{\Sigma}), \ g_i(oldsymbol{eta}) = \mathcal{N}(oldsymbol{\mu}_i, oldsymbol{\Sigma}_i).$
 - $g_0(\beta)$ is fixed to the prior $p(\beta)$, so does not need updating.
- ightharpoonup For $t = 1, \ldots,$
 - ightharpoonup For $i = 1, \ldots, n$
 - 1 Compute the (natural) parameters of the cavity distribution $g_{-i}(\beta)$, $\Sigma_{-i}^{-1} = \Sigma^{-1} \Sigma_{i}^{-1}$, $\Sigma_{-i}^{-1} \mu_{-i} = \Sigma^{-1} \mu \Sigma_{i}^{-1} \mu_{i} \Rightarrow \mu_{-i} = \Sigma_{-i} (\Sigma^{-1} \mu \Sigma_{i}^{-1} \mu_{i})$.
 - 2 Compute the parameters of $g_{-i}(\theta_i)$ $M_{-i} = \mathbf{x}'_i u_{-i}, V_{-i} = \mathbf{x}'_i \Sigma_{-i} \mathbf{x}_i$

Expectation propagation for a regression with known error variance.

3 Construct the un-normalised tilted distribution and calculate,

$$E_k = \int_{-\infty}^{\infty} \theta_i^k g_{-i}(\theta_i) f_i(\theta_i) d\theta_i \quad \text{for } k = 0, 1, 2,$$

where E_0 is the normalising constant for tilted distribution. Then moment match and set $M=\frac{E_1}{E_0}$ and $V=\frac{E_2}{E_0}-M^2$ (M and V are two moments that $g(\theta)$ should have, i.e., the moment matching). Note in practice, we will need to specify finite bounds on the integral. The simplest choice would be $M_{-i}\pm\delta\sqrt{V_{-i}}$ for suitably large δ .

Expectation propagation for a regression with known error variance.

3 Construct the un-normalised tilted distribution and calculate,

$$E_k = \int_{-\infty}^{\infty} \theta_i^k g_{-i}(\theta_i) f_i(\theta_i) d\theta_i$$
 for $k = 0, 1, 2,$

where E_0 is the normalising constant for tilted distribution. Then moment match and set $M=\frac{E_1}{E_0}$ and $V=\frac{E_2}{E_0}-M^2$ (M and V are two moments that $g(\theta)$ should have, i.e., the moment matching). Note in practice, we will need to specify finite bounds on the integral. The simplest choice would be $M_{-i}\pm\delta\sqrt{V_{-i}}$ for suitably large δ .

- 4 Determine the natural parameters of $g_i(\theta_i)$, $M_i/V_i = M/V M_{-i}/V_{-i}$. $1/V_i = 1/V 1/V_{-i}$.
- 5 Transform the natural parameters found in step 4 to those of $g_i(\beta)$, $\Sigma_i^{-1} \mu_i = \mathbf{x}_i M_i / V_i, \Sigma_i^{-1} = \mathbf{x}_i (1/V_i) \mathbf{x}_i'$
- 6 Update the natural parameters of $g(\beta)$, $\Sigma^{-1} = \Sigma_i^{-1} + \Sigma_{-i}^{-1}$ and $\Sigma^{-1}\mu = \Sigma_i^{-1}\mu_i + \Sigma_{-i}^{-1}\mu_{-i}$.

Stop once estimates have converged.



Comments

The differencing method used in the expectation propagation algorithm for finding the natural parameters of $g_{-i}()$, would hold for any g in the exponential family. This is because the definition $g(\theta) = \prod_{i=0}^{n} g_i(\theta)$ implies

$$e^{\boldsymbol{\eta}'u(\boldsymbol{\theta})}=e^{\sum_{i=0}^n \boldsymbol{\eta}_i'u(\boldsymbol{\theta})}.$$

- In previous lectures, we considered at hierarchical models and normal based regression. To fit such models in a fully Bayesian way, we needed to estimate parameters additional to β , such as the error variance, σ_e^2 .
 - How do you think we could do this using expectation propagation?
 - In the case of a normal regression, we could do this by assuming $g(\beta, \sigma_e^2) = \prod_{i=0}^n g_{i,\beta}(\beta) g_{i,\sigma_e^2}(\sigma_e^2)$. This would double the number of factors g_i we need to estimate in the algorithm.

Extending expectation propagation to more complicated models

- In the case of a normal regression, we could do this by assuming $g(\beta, \sigma_e^2) = \prod_{i=0}^n g_{i,\beta}(\beta) g_{i,\sigma_e^2}(\sigma_e^2)$. This would double the number of factors g_i we need to estimate in the algorithm.
- In the write-up of the expectation propagation algorithm, we assumed g_i was normal. Would you necessarily want $g_{i,\sigma_e^2}(\sigma_e^2)$ to be normally distributed?

Extending expectation propagation to more complicated models

- In the case of a normal regression, we could do this by assuming $g(\beta, \sigma_e^2) = \prod_{i=0}^n g_{i,\beta}(\beta) g_{i,\sigma_e^2}(\sigma_e^2)$. This would double the number of factors g_i we need to estimate in the algorithm.
- In the write-up of the expectation propagation algorithm, we assumed g_i was normal. Would you necessarily want $g_{i,\sigma_e^2}(\sigma_e^2)$ to be normally distributed?
 - ▶ Given that $\sigma^2 \in (0, \infty)$ and a normal random variable can take any value from $(-\infty, \infty)$, a normal distribution would be inappropriate. A log-normal might be a good choice though.
- ▶ While we will not look into such examples, there is nothing to stop you in practice from considering more complicated models. In so doing, you may want to
 - ightharpoonup Allow factor blocks to differ by distribution \Rightarrow do not feel restricted to normals.



Example of expectation propagation

- ► To demonstrate expectation propagation, we will return to the logistic regression example.
- As a reminder, the basic logistic regression from a Bayesian perspective corresponds to assuming a flat prior $p(\beta) \propto 1$.
- ► For the purposes of comparison, we will compare the approximate inference obtained using expectation propagation to
 - ▶ fitting a standard glm ⇔ Normal approximation at posterior mode.
 - ▶ Bayesian estimation using Metropolis-Hasting algorithm.
 - ▶ Where possible, exact posterior inference.
- ► The code required for this example is contained in a separate R markdown document.