

A Variational Approach for Bayesian Density Regression

Eric Chu

ERICCHUU@STAT.TAMU.EDU

Department of Statistics

Texas A&M University

College Station, TX 77840, USA

Abstract

In the Bayesian density regression problem, mixture of expert models are often used because of their flexibility in estimating conditional densities. In this paper, we discuss the case when covariate dependent weights are used in the approximating mixture density. Under this framework, however, traditional Bayesian methods results in computational difficulties when the dimension of the covariates is large. In order to remedy this problem and to provide a method for faster inference, we propose using a variational approximation to estimate the conditional density. We also discuss different alternative for approximating quantities that lack a closed form so that a coordinate ascent algorithm is viable.

Keywords: Bayesian Density Regression, Variational Bayes, Mixture Models

1. Introduction

In the Bayesian density regression problem, we observe data $(y_n, x_n)_{n=1}^N$, and the goal is the estimate the conditional density of $y | x$. A common approach for doing this is to model the density using a mixture of gaussians, such as the following,

$$f(y | x) = \sum_k \pi_k \mathcal{N}(y | \mu_k(x), \tau_k^{-1}) \quad (1)$$

While the representation of the density using predictor-independent weights yields less expensive computation, it often lacks flexibility to make it useful in practice and results in a reliance on have too many mixture components. As a result, there have been many proposed models that consider predictor-dependent weights using a kernel stick-breaking process (**dunsonpark:08**) or logit stick-breaking prior (**durante:17**) to generate the weights. In the former method, the increased flexibility comes at heavy computational cost, and in the later method, the process from which the weights are generated does not allow for intuitive inference on the covariates. In our proposed model, the covariates enter through a logistic link function so that we can naturally perform inference on the coefficients. More specifically, we can model

$$f(y | x) = \sum_k^K \pi_k(x) \mathcal{N}(y | \mu_k(x), \tau_k^{-1}) \quad (2)$$

where $\mu_k(x) = x^\top \beta_k$ and $\pi_k \propto \exp(x^\top \gamma_k)$.

2. Notation and Prior Specification

Appendix A.

Taking the expectation with respect to the other variational parameters, we can derive the following variational distribution for \mathbf{Z} ,

$$\begin{aligned}\ln q^*(\mathbf{Z}) &= \sum_n \sum_k z_{nk} \left\{ -\frac{1}{2} \ln(2\pi) + \frac{1}{2} E_{q(\tau)}[\ln \tau_k] - \frac{1}{2} E_{q(\beta, \tau)}[\tau_k (y_n - x_n^\top \beta_k)^2] \right. \\ &\quad \left. + x_n^\top E_{q(\gamma)}[\gamma_k] - E_{q(\gamma)} \left[\ln \left(\sum_j \exp\{x_n^\top \gamma_j\} \right) \right] \right\} \\ &= \sum_n \sum_k z_{nk} \ln \rho_{nk}\end{aligned}$$

where we have defined

$$\begin{aligned}\ln \rho_{nk} &= -\frac{1}{2} \ln(2\pi) + \frac{1}{2} E_{q(\tau)}[\ln \tau_k] - \frac{1}{2} E_{q(\beta, \tau)}[\tau_k (y_n - x_n^\top \beta_k)^2] \\ &\quad + x_n^\top E_{q(\gamma)}[\gamma_k] - E_{q(\gamma)} \left[\ln \left(\sum_j \exp\{x_n^\top \gamma_j\} \right) \right]\end{aligned}\tag{3}$$

Exponentiating and normalizing, we have

$$q^*(\mathbf{Z}) = \prod_n \prod_k r_{nk}^{z_{nk}}, \quad r_{nk} = \frac{\rho_{nk}}{\sum_j \rho_{nj}}\tag{4}$$

For the discrete distribution $q^*(\mathbf{Z})$ given in (4) above, we have $E[z_{nk}] = r_{nk}$. Note, however, that in order to compute the expectation in closed form, we need an expression for the four expectations involved in the quantity $\ln \rho_{nk}$, as defined in (3).

From the results derived in Appendix C, we know that $q^*(\tau_k) = \text{Ga}(\tau_k \mid a_k, b_k)$. We can then compute the following expectation with respect to $q^*(\tau)$.

$$E_{q(\tau)}[\ln \tau_k] = \psi(a_k) - \psi(b_k)\tag{5}$$

Again from Appendix C, we can then compute the following expectation with respect to $q^*(\beta_k, \tau_k)$.

$$\begin{aligned}E_{q(\beta, \tau)}[\tau_k (y_n - x_n^\top \beta_k)^2] &= E \left[\tau_k \left(y_n - m_k^\top x_n x_n^\top m_k + \text{tr} \left(x_n x_n^\top (\tau_k V_k)^{-1} \right) - 2y_n x_n^\top m_k \right) \right] \\ &= \frac{a_k}{b_k} (y_n^2 + m_k^\top x_n x_n^\top m_k) + \text{tr} (x_n x_n^\top V_k^{-1}) \\ &= \frac{a_k}{b_k} (y_n + m_k^\top x_n)^2 + x_n^\top V_k^{-1} x_n\end{aligned}\tag{6}$$

From the expression derived in (11) of Appendix B, we have $q^*(\gamma_k) = \mathcal{N}(\gamma_k \mid \mu_k, \mathbf{Q}_k^{-1})$, then we have

$$E_{q(\gamma_k)}[\gamma_k] = \mu_k\tag{7}$$

Using the bound discussed in Appendix B, equation (10), we can then compute the following expectation with respect to $q^*(\gamma)$.

$$\begin{aligned}
& \mathbb{E}_{q(\gamma)} \left[\ln \left(\sum_j^K \exp\{x_n^\top \gamma_j\} \right) \right] \\
& \approx \mathbb{E}_{q(\gamma)} \left[\alpha_n + \sum_{j=1}^K \frac{x_n^\top \gamma_j - \alpha_n + \xi_{nj}}{2} + \lambda(\xi_{nj}) \left((x_n^\top \gamma_j - \alpha_n)^2 - \xi_{nj}^2 \right) + \log \left(1 + e^{\xi_{nj}} \right) \right] \\
& = \alpha_n + \sum_j^K \frac{1}{2} (x_n^\top \mu_j - \alpha_n + \xi_{nj}) + \lambda(\xi_{nj}) \left((x_n^\top \mu_j - \alpha_k)^2 - \xi_{nj}^2 + x_j^\top Q_k^{-1} x_j \right) + \log(1 + e^{\xi_{nj}})
\end{aligned} \tag{8}$$

Gathering the results in (5), (6), (7), and (8), and substituting these into (3), we can compute $E[z_{nk}] = r_{nk}$ in closed form.

Appendix B.

For the variational distribution for $\gamma_k, k = 1, \dots, K$, we first note the following bound given by **bouchard:07**, $\sum_{j=1}^K e^{t_j} \leq \prod_{j=1}^K (1 + e^{t_j})$. Setting $t_j = x_n^\top \gamma_j - \alpha_n$ and then taking log, we have the following bound:

$$\log \left(\sum_{j=1}^K \exp\{x_n^\top \gamma_j\} \right) \leq \alpha_n + \sum_{j=1}^K \log(1 + \exp\{x_n^\top \gamma_j - \alpha_n\}) \tag{9}$$

If we then use the bound from **jj:2001**,

$$\log(1 + e^x) \leq \frac{x - t}{2} + \frac{1}{4t} \tanh\left(\frac{t}{2}\right) (x^2 - t^2) + \log(1 + e^t)$$

then we arrive at the following bound:

$$\log \left(\sum_{j=1}^K \exp\{x_n^\top \gamma_j\} \right) \leq \alpha_n + \sum_{j=1}^K \frac{x_n^\top \gamma_j - \alpha_n + \xi_{nj}}{2} + \lambda(\xi_{nj}) \left((x_n^\top \gamma_j - \alpha_n)^2 - \xi_{nj}^2 \right) + \log(1 + e^{\xi_{nj}}) \tag{10}$$

where $\lambda(\xi) = \frac{1}{4\xi} \tanh\left(\frac{\xi}{2}\right)$. Then we can substitute this back into $\ln q^*(\gamma_k)$ to obtain an approximation for the left hand side of (9), thus allowing us to obtain a closed form for the

variational distribution. Note that all of the equalities above are written up to constants.

$$\begin{aligned}
\ln q^*(\gamma_k) &= -\frac{1}{2}\gamma_k^\top \gamma_k + \sum_n r_{nk} x_n^\top \gamma_k - \sum_n r_{nk} \ln \left(\sum_j \exp\{x_n^\top \gamma_j\} \right) \\
&\approx -\frac{1}{2}\gamma_k^\top \gamma_k + \gamma_k^\top \sum_n r_{nk} x_n \\
&\quad - \sum_n r_{nk} \left\{ \alpha_n + \sum_{j=1}^K \frac{x_n^\top \gamma_j - \alpha_n + \xi_{nj}}{2} + \lambda(\xi_{nj}) ((x_n^\top \gamma_j - \alpha_n)^2 - \xi_{nj}^2) + \log(1 + e^{\xi_{nj}}) \right\} \\
&= -\frac{1}{2}\gamma_k^\top \gamma_k + \gamma_k^\top \sum_n r_{nk} x_n - \sum_n r_{nk} \left\{ \frac{1}{2}\gamma_k^\top x_n + \lambda(\xi_{nj}) (\gamma_j^\top x_n x_n^\top \gamma_j - 2\alpha_n \gamma_j^\top x_n) \right\} \\
&= -\frac{1}{2}\gamma_k^\top \left(\mathbf{I}_D + 2 \sum_n r_{nk} \lambda(\xi_{nk}) x_n x_n^\top \right) \gamma_k + \gamma_k^\top \left(\sum_n r_{nk} \left(\frac{1}{2} + 2\lambda(\xi_{nk}) \alpha_n x_n \right) \right)
\end{aligned}$$

Exponentiating, we can recover $q^*(\gamma_k) = \mathcal{N}(\mu_k, \mathbf{Q}_k^{-1})$, where

$$\begin{aligned}
\mu_k &= \mathbf{Q}_k^{-1} \eta_k \\
\eta_k &= \sum_n r_{nk} \left(\frac{1}{2} + 2\lambda(\xi_{nj}) \alpha_n \right) x_n \\
\mathbf{Q}_k &= \mathbf{I}_D + 2 \sum_n r_{nk} \lambda(\xi_{nk}) x_n x_n^\top
\end{aligned} \tag{11}$$

The additional parameters introduced in the two upper bounds can be updated using the following equations

$$\begin{aligned}
\xi_{nk} &= \sqrt{(\mu_k^\top x_n - \alpha_n)^2 + x_n^\top \mathbf{Q}_k^{-1} x_n} \quad \forall k, n \\
\alpha_n &= \frac{\frac{1}{2} \left(\frac{K}{2} - 1 \right) + \sum_{j=1}^K \lambda(\xi_{nj}) \mu_j^\top x_n}{\sum_{j=1}^K \lambda(\xi_{nj})} \quad \forall n
\end{aligned}$$

Appendix C.

Using results from Appendix A, we can write the following expression for the joint variational distribution of (β_k, τ_k) ,

$$\begin{aligned}
\ln q^*(\beta_k, \tau_k) &= \sum_n -\frac{1}{2} r_{nk} \tau_k (y_n^2 + \beta_k^\top x_n x_n^\top \beta_k - 2y_n \beta_k^\top x_n) + \frac{r_{nk}}{2} \ln \tau_k + \frac{D}{2} \ln \tau_k \\
&\quad - \frac{\tau_k}{2} (\beta_k^\top \Lambda_0 \beta_k + m_0^\top \Lambda_0 m_0 - 2\beta_k^\top \Lambda_0 m_0) + (a_0 - 1) \ln \tau_k - b_0 \tau_k
\end{aligned} \tag{12}$$

We first consider terms on the right hand side of (12) that depend on β_k to find $\ln q^*(\beta_k | \tau_k)$, giving

$$\ln q^*(\beta_k | \tau_k) = -\frac{\tau_k}{2} \beta_k^\top \left[\sum_n r_{nk} x_n x_n^\top + \Lambda_0 \right] \beta_k + \tau_k \beta_k^\top \left[\sum_n r_{nk} y_n x_n + \Lambda_0 m_0 \right] \tag{13}$$

$$q^*(\beta_k | \tau_k) = \mathcal{N}(m_k, (\tau_k V_k)^{-1}) \quad (14)$$

$$\begin{aligned} m_k &= V_k^{-1} b_k \\ V_k &= \sum_n r_{nk} x_n x_n^\top + \Lambda_0 \\ b_k &= \sum_n r_{nk} y_n x_n + \Lambda_0 m_0 \end{aligned} \quad (15)$$

Then we can make use of the relation $\ln q^*(\tau_k) = \ln q^*(\beta_k, \tau_k) - \ln q^*(\beta_k | \tau_k)$, where the quantities on the right hand side come from (12) and (14). Note that equality below is written up to constants, keeping only terms involving τ_k .

$$\begin{aligned} \ln q^*(\tau_k) &= (a_0 + N_k - 1) \ln \tau_k - \tau_k \left\{ b_0 + \frac{1}{2} \left(\sum_n r_{nk} y_n^2 + m_0^\top \Lambda_0 m_0 - m_k^\top V_k m_k \right) \right. \\ &\quad + \frac{1}{2} \beta_k^\top \left(\sum_n r_{nk} x_n x_n^\top + \Lambda_0 - V_k \right) \beta_k \\ &\quad \left. - 2 \beta_k^\top \left(\sum_n r_{nk} y_n x_n + \Lambda_0 m_0 - V_k m_k \right) \right\} \end{aligned} \quad (16)$$

Exponentiating, we arrive at the following distribution

$$q^*(\tau_k) = \text{Ga}(\tau_k | a_k, b_k) \quad (17)$$

where we have defined

$$\begin{aligned} a_k &= a_0 + N_k \\ b_k &= b_0 + \frac{1}{2} \sum_n r_{nk} y_n^2 + m_0^\top \Lambda_0 m_0 - b_k^\top V_k^{-1} b_k \end{aligned} \quad (18)$$

The expression for b_k arises by noting that the three following simplifications for the summation terms in the coefficient of τ_k in (16),

$$\begin{aligned} \sum_n r_{nk} y_n^2 + m_0^\top \Lambda_0 m_0 - m_k^\top V_k m_k &= \sum_n r_{nk} y_n^2 + m_0^\top \Lambda_0 m_0 - b_k^\top V_k^{-1} b_k \\ \sum_n r_{nk} x_n x_n^\top + \Lambda_0 - V_k &= 0 \\ \sum_n r_{nk} y_n x_n + \Lambda_0 m_0 - V_k m_k &= 0 \end{aligned}$$

where the first equality holds by expanding $m_k^\top V_k m_k = b_k^\top (V_k^{-1})^\top V_k V_k^{-1} b_k = b_k^\top V_k^{-1} b_k$. The second equality holds by recalling the definition of V_k in (15), and the third equality holds by observing from (15) that $V_k m_k = b_k$.

Appendix D.

sample