

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: May 14, 2021

This is Bayes Rule!

Given:

$$\arg \min_{q(\theta)} - \sum_{n=1}^N \left[\int q(\theta) \log p(\mathbf{x}_n | \theta) d\theta \right] + KL(q(\theta) \| p(\theta)) \quad (1)$$

$$= \arg \max_{q(\theta)} \int q(\theta) \log p(\mathbf{X} | \theta) d\theta + \int q(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta \quad (2)$$

$$= \arg \max_{q(\theta)} \underbrace{\int q(\theta) \log \frac{p(\mathbf{X}, \theta)}{q(\theta)} d\theta}_{\mathcal{L}(q)} \quad (3)$$

The identity we know

$$\log p(\mathbf{X}) = \mathcal{L}(q) + KL(q(\theta) \| p(\theta | X)) \quad (4)$$

As the $\log p(\mathbf{X})$ is constant wrt to θ

$$\arg \max_{q(\theta)} \mathcal{L}(q) = \arg \min_q KL(q(\theta) \| p(\theta | X)) \quad (5)$$

$KL(q(\theta) \| p(\theta | X))$ will attain its minimum value when $q(\theta) = p(\theta | X)$

Therefore by minimizing the given equation we can get the posterior of parameters.

Intuition of given equation

1. The first term is one such a way that best explains the data
2. Second term is regularizer makes sure that $q(\theta)$ close to prior $p(\theta)$

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: May 14, 2021

Mean-Field VI for Sparse Bayesian Linear Regression

Given

$$y_n \sim \mathcal{N}(y_n | \omega^T \mathbf{x}_n, \beta^{-1}) \quad (6)$$

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\boldsymbol{\omega}, \beta^{-1}\mathbf{I}) \quad (7)$$

$$p(\boldsymbol{\omega}) = \mathcal{N}(\boldsymbol{\omega} | 0, \boldsymbol{\alpha}^{-1}) \text{ where } \boldsymbol{\alpha}^{-1} = \text{diag}(\alpha_1^{-1}, \dots, \alpha_D^{-1}) \quad (8)$$

$$\omega_d \sim \mathcal{N}(\omega_d | 0, \alpha_d) \quad (9)$$

$$\beta \sim \text{Gamma}(\beta | a_0, b_0) \quad (10)$$

$$\alpha_d \sim \text{Gamma}(\alpha_d | e_0, f_0) \quad (11)$$

The joint probability can be written as

$$\log p(\mathbf{y}, \boldsymbol{\omega}, \beta, \alpha_1, \dots, \alpha_D | \mathbf{X}) = \log p(\mathbf{y} | \beta, \boldsymbol{\omega}, \mathbf{X}) + \log p(\boldsymbol{\omega} | \alpha_1, \dots, \alpha_D) + \log p(\beta) + \sum_{d=1}^D \log p(\alpha_d) \quad (12)$$

Now update for $\boldsymbol{\omega}$

$$\log q_{\boldsymbol{\omega}}^*(\boldsymbol{\omega}) = \mathbb{E}_{\beta, \alpha_1, \dots, \alpha_D} [\log p(\mathbf{y}, \boldsymbol{\omega}, \beta, \alpha_1, \dots, \alpha_D | \mathbf{X})] \quad (13)$$

$$= \mathbb{E}_{\beta, \alpha_1, \dots, \alpha_D} [\log p(\mathbf{y} | \mathbf{X}\boldsymbol{\omega}, \beta) + \log p(\boldsymbol{\omega} | \alpha_1, \dots, \alpha_D)] + \text{const} \quad (14)$$

$$= \mathbb{E}_{\beta, \alpha_1, \dots, \alpha_D} [\log \mathcal{N}(\mathbf{y} | \mathbf{X}\boldsymbol{\omega}, \beta^{-1}\mathbf{I}) + \log \mathcal{N}(\boldsymbol{\omega} | 0, \boldsymbol{\alpha})] \quad (15)$$

$$\text{using gaussian properties} \quad (16)$$

$$q_{\boldsymbol{\omega}}^*(\boldsymbol{\omega}) = \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\omega}}, \boldsymbol{\Sigma}_{\boldsymbol{\omega}}) \quad (17)$$

$$\boldsymbol{\Sigma}_{\boldsymbol{\omega}} = (\mathbb{E}[\boldsymbol{\alpha}] + \mathbb{E}[\beta] \mathbf{X}^T \mathbf{X})^{-1} \quad (18)$$

$$\boldsymbol{\mu}_{\boldsymbol{\omega}} = \mathbb{E}[\beta] \boldsymbol{\Sigma}_{\boldsymbol{\omega}} \mathbf{X}^T \mathbf{y} \quad (19)$$

$$\mathbb{E}[\boldsymbol{\alpha}] = \text{diag}(\mathbb{E}[\alpha_1], \dots, \mathbb{E}[\alpha_D]) \quad (20)$$

update for β

$$\log q_{\beta}^*(\beta) = \mathbb{E}_{\boldsymbol{\omega}, \alpha_1, \dots, \alpha_D} [\log p(\mathbf{y} | \mathbf{X}\boldsymbol{\omega}, \beta) + \log p(\beta | a_0, b_0)] + \text{const} \quad (21)$$

$$= \mathbb{E}_{\boldsymbol{\omega}, \alpha_1, \dots, \alpha_D} [\log \mathcal{N}(\mathbf{y} | \mathbf{X}\boldsymbol{\omega}, \beta^{-1}\mathbf{I}) + \log \text{Gamma}(\beta | a_0, b_0)] \quad (22)$$

$$\text{by conjugacy} \quad (23)$$

$$q_{\beta}^*(\beta) = \text{Gamma}(\beta | a_{\beta}, b_{\beta}) \quad (24)$$

$$a_{\beta} = \frac{N}{2} + a_0 \quad (25)$$

$$b_{\beta} = \mathbb{E}[\frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\omega}\|^2] + b_0 \quad (26)$$

$$= \frac{\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbb{E}[\boldsymbol{\omega}] + \text{Tr}(\mathbf{X}^T \mathbf{X} \mathbb{E}[\boldsymbol{\omega} \boldsymbol{\omega}^T])}{2} + b_0 \quad (27)$$

Now update for α_d

$$\log q_{\alpha_d}^*(\alpha_d) = \mathbb{E}_{\boldsymbol{\omega}, \beta, \alpha_{i \neq d}} [\log p(\omega_d | 0, \alpha_d) + \log p(\alpha_d | e_0, f_0)] + \text{const} \quad (28)$$

$$= \mathbb{E}_{\boldsymbol{\omega}, \beta, \alpha_{i \neq d}} [\log \mathcal{N}(\omega_d | 0, \alpha_d) + \log \text{Gamma}(\alpha_d | e_0, f_0)] \quad (29)$$

$$\text{by conjugacy} \quad (30)$$

$$q_{\alpha_d}^*(\alpha_d) = \text{Gamma}(\alpha_d | e_{\alpha_d}, f_{\alpha_d}) \quad (31)$$

$$e_{\alpha_d} = \frac{1}{2} + e_0 \quad (32)$$

$$f_{\alpha_d} = \frac{\mathbb{E}[\omega_d^2]}{2} + f_0 \quad (33)$$

We following expectations are required for above calculations

$$\mathbb{E}[\boldsymbol{\omega}] = \boldsymbol{\mu}_{\boldsymbol{\omega}} \quad (34)$$

$$\mathbb{E}[\boldsymbol{\omega} \boldsymbol{\omega}^T] = \boldsymbol{\Sigma}_{\boldsymbol{\omega}} + \boldsymbol{\mu}_{\boldsymbol{\omega}} \boldsymbol{\mu}_{\boldsymbol{\omega}}^T \quad (35)$$

$$\mathbb{E}[\omega_d^2] = [\boldsymbol{\Sigma}_{\boldsymbol{\omega}}]_{dd} + [\boldsymbol{\mu}_{\boldsymbol{\omega}} \boldsymbol{\mu}_{\boldsymbol{\omega}}^T]_{dd} \quad \forall d \quad (36)$$

$$\mathbb{E}[\alpha_d] = \frac{e_{\alpha_d}}{f_{\alpha_d}} \quad \forall d \quad (37)$$

$$\mathbb{E}[\beta] = \frac{a_{\beta}}{b_{\beta}} \quad (38)$$

Mean Field VI Algorithm

1. Set $a_{\beta} = a_0 + \frac{N}{2}$ and $e_{\alpha_d} = e_0 + \frac{1}{2}$
2. Set $t=0$. Initialize, $\boldsymbol{\Sigma}_{\boldsymbol{w}}^{(0)} = \mathbf{I}_D$ and $\boldsymbol{\mu}_{\boldsymbol{w}}^{(0)} = 0$ so,

$$\begin{aligned} f_{\alpha_d}^{(0)} &= f_0 + \frac{1}{2} \\ b_{\beta}^{(0)} &= b_0 + \frac{\mathbf{y}^T \mathbf{y} + \text{Tr}(\mathbf{X}^T \mathbf{X})}{2} \\ \mathbb{E}[\beta]^{(0)} &= \frac{a_{\beta}}{b_{\beta}^{(0)}} \\ \mathbb{E}[\alpha_d]^{(0)} &= \frac{e_{\alpha_d}}{f_{\alpha_d}^{(0)}}, \forall d \end{aligned}$$

3. Set $t=t+1$. Until convergence, repeat.

$$\begin{aligned} \boldsymbol{\Sigma}_{\boldsymbol{\omega}}^{(t)} &= (\mathbb{E}[\boldsymbol{\alpha}]^{(t-1)} + \mathbb{E}[\beta]^{(t-1)} \mathbf{X}^T \mathbf{X})^{-1} \\ \boldsymbol{\mu}_{\boldsymbol{\omega}}^{(t)} &= \mathbb{E}[\beta]^{(t-1)} \boldsymbol{\Sigma}_{\boldsymbol{\omega}}^{(t)} \mathbf{X}^T \mathbf{y} \\ b_{\beta}^{(t)} &= \frac{\mathbf{y}^T \mathbf{y} - 2 \mathbf{y}^T \mathbf{X} \mathbb{E}[\boldsymbol{\omega}]^{(t)} + \text{Tr}(\mathbf{X}^T \mathbf{X} \mathbb{E}[\boldsymbol{\omega} \boldsymbol{\omega}^T]^{(t)})}{2} + b_0 \\ f_{\alpha_d}^{(t)} &= \frac{\mathbb{E}[\omega_d^2]^{(t)}}{2} + f_0, \forall d \\ \mathbb{E}[\beta]^{(t)} &= \frac{a_{\beta}}{b_{\beta}^{(t)}} \\ \mathbb{E}[\alpha_d]^{(t)} &= \frac{e_{\alpha_d}}{f_{\alpha_d}^{(t)}}, \forall d \end{aligned}$$

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: May 14, 2021

Gibbs Sampling

Given

$$x_1, \dots, x_N \sim \text{Poisson}(x_n | \lambda_n) \quad (39)$$

$$\lambda_n \sim \text{Gamma}(\lambda_n | \alpha, \beta) \quad (40)$$

$$\alpha \sim \text{Gamma}(\alpha | a, b) \quad (41)$$

$$\beta \sim \text{Gamma}(\beta | c, d) \quad (42)$$

To perform the gibbs sampling. of posteriror we need CP's. Finding the conditional posteriors (CP)

wrt $\lambda_n \forall n$

$$p(\lambda_n | \lambda_{-n}, \alpha, \beta, X) \propto p(x_n | \lambda_n) p(\lambda_n | \alpha, \beta) \quad (43)$$

$$= \text{Poisson}(x_n | \lambda_n) \text{Gamma}(\lambda_n | \alpha, \beta) \quad (44)$$

$$\text{by conjugacy} \quad (45)$$

$$p(\lambda_n | \lambda_{-n}, \alpha, \beta, x_n) = \text{Gamma}(\lambda_n | \alpha_{\lambda_n}, \beta_{\lambda_n}) \quad (46)$$

$$\alpha_{\lambda_n} = x_n + \alpha \quad (47)$$

$$\beta_{\lambda_n} = \beta + 1 \quad (48)$$

wrt α

$$p(\alpha | \lambda_1, \dots, \lambda_N, \beta, X) \propto \prod_{n=1}^N p(\lambda_n | \alpha, \beta) p(\alpha | a, b) \quad (49)$$

$$\propto \prod_{n=1}^N \text{Gamma}(\lambda_n | \alpha, \beta) \text{Gamma}(\alpha | a, b) \quad (50)$$

$$\propto \left(\frac{\beta^a}{\Gamma(\alpha)} \right)^N \left[\prod_{n=1}^N \lambda_n \right]^{\alpha-1} \exp(-\beta \sum_{n=1}^N \lambda_n - b\alpha) \alpha^{a-1} \quad (51)$$

The above equation is not in the closed form wrt α

$$p(\beta | \lambda_1, \dots, \lambda_N, \alpha, X) \propto \prod_{n=1}^N p(\lambda_n | \alpha, \beta) p(\beta | c, d) \quad (52)$$

$$\propto \prod_{n=1}^N \text{Gamma}(\lambda_n | \alpha, \beta) \text{Gamma}(\beta | c, d) \quad (53)$$

$$\text{by conjugacy} \quad (54)$$

$$p(\beta | \lambda_1, \dots, \lambda_N, \alpha, X) = \text{Gamma}(c + N\alpha, d + \sum_{n=1}^N \lambda_n) \quad (55)$$

Student Name: Musale Krushna Pavan
Roll Number: 20111268
Date: May 14, 2021

Using Samples for Prediction

Consider matrix factorization model $N \times M$

$$\text{likelihood : } p(r_{ij}|\mathbf{u}_i, \mathbf{v}_j) = \mathcal{N}(r_{ij}|\mathbf{u}_i^T \mathbf{v}_j, \beta^{-1}) \quad (56)$$

$$\text{ppd : } p(r_{ij}|R) = \int p(r_{ij}|\mathbf{u}_i, \mathbf{v}_j)p(\mathbf{u}_i, \mathbf{v}_j|R)d\mathbf{u}_i d\mathbf{v}_j \quad (57)$$

Given the samples $\{\mathbf{U}^{(s)}, \mathbf{V}^{(s)}\}_{s=1}^S$ from gibbs sampler. PPD will can written as:

$$p(r_{ij}|R) \approx \frac{1}{S} \sum_{s=1}^S p(r_{ij}|\mathbf{u}_i^{(s)}, \mathbf{v}_j^{(s)}) \quad (58)$$

Now calculating expectation of r_{ij}

$$\mathbb{E}_{r_{ij}|R}[r_{ij}] = \int r_{ij} p(r_{ij}|R) dr_{ij} \quad (59)$$

$$= \frac{1}{S} \sum_{s=1}^S \int r_{ij} \mathcal{N}(r_{ij}|\mathbf{u}_i^{T(s)} \mathbf{v}_j^{(s)}, \beta^{-1}) dr_{ij} \quad (60)$$

$$= \frac{1}{S} \sum_{s=1}^S \mathbb{E}[\mathcal{N}(r_{ij}|\mathbf{u}_i^{T(s)} \mathbf{v}_j^{(s)}, \beta^{-1})] \quad (61)$$

$$= \frac{1}{S} \sum_{s=1}^S \mathbf{u}_i^{T(s)} \mathbf{v}_j^{(s)} \quad (62)$$

Now calculating expectation of r_{ij}

$$\text{var}_{r_{ij}|R}(r_{ij}) = \mathbb{E}[r_{ij}^2] - [\mathbb{E}[r_{ij}]]^2 \quad (63)$$

$$\mathbb{E}_{r_{ij}|R}[r_{ij}^2] = \int r_{ij}^2 p(r_{ij}|R) dr_{ij} - \left(\frac{1}{S} \sum_{s=1}^S \mathbf{u}_i^{T(s)} \mathbf{v}_j^{(s)}\right)^2 \quad (64)$$

$$= \frac{1}{S} \sum_{s=1}^S \int r_{ij}^2 \mathcal{N}(r_{ij}|\mathbf{u}_i^{T(s)} \mathbf{v}_j^{(s)}, \beta^{-1}) dr_{ij} - \left(\frac{1}{S} \sum_{s=1}^S \mathbf{u}_i^{T(s)} \mathbf{v}_j^{(s)}\right)^2 \quad (65)$$

$$= \frac{1}{S} \sum_{s=1}^S \left(\beta^{-1} + (\mathbf{u}_i^{T(s)} \mathbf{v}_j^{(s)})^2\right) - \left(\frac{1}{S} \sum_{s=1}^S \mathbf{u}_i^{T(s)} \mathbf{v}_j^{(s)}\right)^2 \quad (66)$$

$$= \beta^{-1} + \frac{1}{S} \sum_{s=1}^S (\mathbf{u}_i^{T(s)} \mathbf{v}_j^{(s)})^2 - \left(\frac{1}{S} \sum_{s=1}^S \mathbf{u}_i^{T(s)} \mathbf{v}_j^{(s)}\right)^2 \quad (67)$$

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: May 14, 2021

Rejection Sampling

Given

$$p(x) \propto \exp(\sin(x)) \quad -\pi \leq x \leq \pi \quad (68)$$

$$\text{poposal distribution for rejection sampling} \quad (69)$$

$$q(x) = \mathcal{N}(x|0, \sigma^2) \quad (70)$$

For doing rejection sampling we need the following condition

$$Mq(x) \geq \tilde{p}(x) \quad \forall -\pi \leq x \leq \pi \quad (71)$$

$$M \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{x^2}{2\sigma^2}\right) \geq \exp(\sin(x)) \quad (72)$$

$$M \geq \frac{\sigma\sqrt{2\pi} \exp(\sin(x))}{\exp\left(-\frac{x^2}{2\sigma^2}\right)} \quad (73)$$

$$\text{Approximating max value of M for } \forall -\pi \leq x \leq \pi \quad (74)$$

$$\max \exp(\sin(x)) = e \quad (75)$$

$$\min \exp\left(-\frac{x^2}{2\sigma^2}\right) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{\pi^2}{2\sigma^2}\right) \quad (76)$$

$$M \geq \sigma\sqrt{2\pi} \exp\left(\frac{\pi^2}{2\sigma^2} + 1\right) \quad (77)$$

considering the $\sigma^2 = 1$ we get $M = \sqrt{2\pi} \exp\left(\frac{\pi^2}{2} + 1\right) = 947.4183259814814$

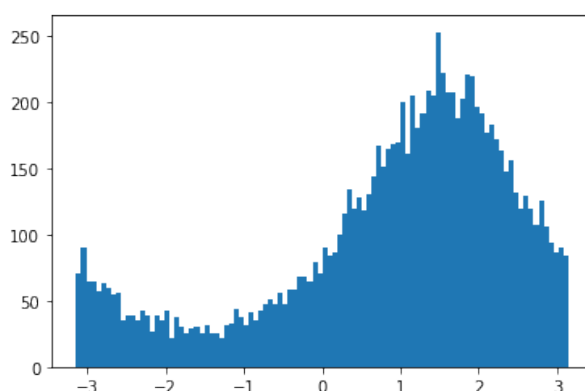


Figure 1: Hist of samples from $p(x)$