

Student Name: Ajita Shree

Roll Number: 20111262

Date: May 14, 2021

This is Bayes' rule!

Given observations x_1, x_2, \dots, x_N , drawn i.i.d from the observations model $p(x/\theta)$ and prior distribution $p(\theta)$



To Prove: solving the below equation is equivalent to finding the posterior distribution of **Eqn 1:** $\theta \operatorname{argmin}_{q(\theta)} - \sum_{n \in N} [\int q(\theta) \log(x_n/\theta) d\theta] + KL(q(\theta)||p(\theta))$

- **Explanation:** The above objective function is ELBO expression, i.e. expected lower bound optimization.
- Maximizing ELBO will give an approximating distribution $q(\theta)$ which explains data well i.e. give it large probability i.e. large exp log-likelihood i.e. $E_q[\log(p(X/\theta))]$
- KL term however, takes care of that $q(\theta)$ is close to the prior distribution $p(\theta)$, this will act as simple regularizer so that it doesn't over-fit the data much.
- MAP, posterior estimation of $\theta, p(\theta/x) \propto p(x/\theta)p(\theta)$
- Here also, there are two major objectives
 - Maximizing the likelihood of the data
 - Regularizer $p(\theta)$ will ensure there is not much overfitting.
- As we have seen, both approaches are for optimizing same objectives and hence, minimizing the -ve of ELBO is equivalent to maximizing the posterior prediction distribution.
- **Formal proof is as follows:** Let us optimize Eqn 1 wrt $q(\theta)$ and find out the optimal $q(\theta)$ value; Differentiating and equating it to 0 will give
 - $d/d(q(\theta)) \left\{ - \sum_{n \in N} \left[\int q(\theta) \log(x_n/\theta) d\theta \right] + \int q(\theta) \log(q(\theta)/p(\theta/X)) d\theta \right\}$
 - $\sum_{n \in N} \int \log(x_n/\theta) d\theta + \int \log(q(\theta)/p(\theta)) = 0$
 - $\int \log(q(\theta)/p(\theta)) = \sum_{n \in N} \int \log(x_n/\theta) d\theta$
 - Rearranging terms and removing integrals, we have
 - $\log(q(\theta)) = \sum_{n \in N} \log(x_n/\theta) + \log(p(\theta))$
 - $q(\theta) = \prod_{n \in N} p(x_n/\theta) \cdot p(\theta)$
 - **Hence proved, best possible value of $q(\theta)$ is nothing but the posterior distribution $p(\theta/X)$**

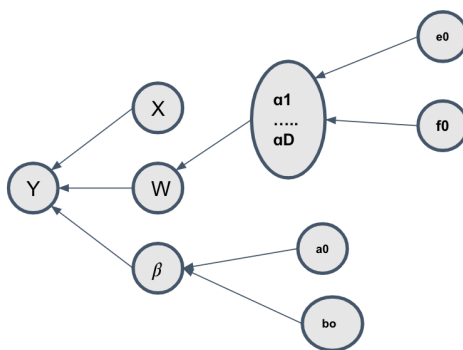
Student Name: Ajita Shree

Roll Number: 20111262

Date: May 14, 2021

Mean Field VI for Sparse Bayesian Linear Regression

The probability distributions of all rvs and their respective graphical model is as follows.



- **Given** $y_n = N(y_n/w^T x_n, \beta^{-1})$
- $p(w) = N(w/0, \text{diag}(\alpha_1^{-1}, \dots, \alpha_D^{-1}))$
- $p(\beta) = \text{Gamma}(\beta/a_0, b_0)$
- $\alpha_d = \text{Gamma}(\alpha_d/e_0, f_0)$
- To Derive: $q(w, \beta, \alpha_1 \dots \alpha_D) = q(w)q(\beta)q(\alpha_1) \dots q(\alpha_D) = p(w, \beta, \alpha_1, \alpha_2, \dots, \alpha_D / Y, X)$
- Note : $\text{Gamma}(\eta/\tau_1, \tau_2) = \frac{\tau_2^{\tau_1}}{\Gamma(\tau_1)} \eta^{\tau_1-1} \exp(-\tau_2 \eta)$

• **Solution:**

- **Joint distribution** $p(y, w, \beta, \alpha_1, \alpha_2, \dots, \alpha_D / X) = \prod_{n \in N} p(y_n/w^T x_n, \beta) * p(w/\alpha_1, \dots, \alpha_D) * p(\beta) * \prod_{d \in D} p(\alpha_d)$
- Taking log of joint distribution will give
- $\log(p(y, w, \beta, \alpha_1, \dots, \alpha_D) / X) = \sum_{n \in N} \log(p(y_n/w, x_n, \beta)) + \log(p(w/\alpha_1 \dots \alpha_D)) + \log(\beta) + \sum_{d \in D} \log(p(\alpha_d))$
- $\log(p(y, w, \beta, \alpha_1, \dots, \alpha_D) / X) = \sum_{n \in N} \log N(y_n/w^T x_n, \beta^{-1}) + \log N(w/0, \text{diag}(\alpha_1^{-1}, \dots, \alpha_D^{-1})) + \log \text{Gamma}(\beta/a_0, b_0) + \sum_{d \in D} \log \text{Gamma}(\alpha_d/e_0, f_0)$
- Expanding the probability distributions, we will get
- $\sum_{n \in N} \log(\sqrt{\frac{\beta}{2 * \pi}} \exp(\frac{-\beta}{2} * (y_n - w^T x_n)^2)) + \log(\sqrt{(\frac{\alpha_1 \dots \alpha_D}{(2\pi)^D}}) \exp(\frac{-w^T \Sigma w}{2}))$

- $+ \log(\frac{b_0^{a_0} \beta^{a_0-1}}{T(a_0)} \exp(-b_0 \beta)) + \sum_{d \in D} \log(\frac{f_0^{e_0} \alpha_d^{e_0-1}}{T(e_0)} \exp(-f_0 \alpha_d))$ where Σ is diagonal matrix with α_d at d th row, col.
- **Eq1: Joint dist** $\log(p(y, w, \beta, \alpha_1, \dots, \alpha_D)/X) \propto (N/2) \log \beta - (\beta/2) \sum_{n \in N} (y_n - w^T x_n)^2 + (1/2) \sum_{d \in D} \log \alpha_d - (1/2) w^T \Sigma w + (a_0 - 1) \log \beta - b_0 \beta + (e_0 - 1) \sum_{d \in D} \log \alpha_d - f_0 \sum_{d \in D} \alpha_d$
- **Note 1:** We know that standard mean field VI algorithm, the approximate distribution for an rv $z_j, q_j(z_j) \propto [E_{i \neq j}(\log p(X, Z))]$, where X is the true data and z 's are the rvs on which data is dependent.
- Using the fact in **Note 1**, we can write $\log q_w^*(w), \log q_\beta^*(\beta)$ and $\log q_{\alpha_d}^*(\alpha_d)$ as follows:
- **Eq2: For w**
 1. $\log(q_w^*(w)) = E_{\beta, \alpha_1, \dots, \alpha_D} [\log p(y, w, \beta, \alpha_1, \dots, \alpha_D)/X]$
 2. Taking terms containing w from equation 1, we have
 3. $E_{q, \beta, \alpha_1, \dots, \alpha_D} [\frac{-\beta}{2} \sum_{n \in N} (y_n - w^T x_n)^2 - \frac{1}{2} w^T \Sigma w] + \text{const.}$
 4. $(-1/2) [E[\beta] w^T (\sum_{n \in N} x_n x_n^T) w - E[\beta] 2 w^T \sum_{n \in N} y_n x_n + w^T \text{diag}(E[\alpha_1] \dots E[\alpha_D] w)]$
 5. **For w 's update, $E[\beta], E[\alpha_1] \dots E[\alpha_D]$ will be needed**
- **Eq3: For β**
 1. $\log(q_\beta^*(\beta)) = E_{w, \alpha_1, \dots, \alpha_D} [\log p(y, w, \beta, \alpha_1, \dots, \alpha_D)/X]$
 2. $\propto E[(N/2) \log \beta - (\beta/2) \sum_{n \in N} (y_n - w^T x_n)^2 + (a_0 - 1) \log \beta - b_0 \beta]$
 3. On further solving, we have, $(N/2 + a_0 - 1) \log \beta - \beta (\sum_{n \in N} (1/2) E(w^T (\sum_{n \in N} x_n x_n^T) w - 2 w^T \sum_{n \in N} y_n x_n) + b_0)$
 4. The above form is similar to Gamma form with hyperparameters $N/2 + a_0$ and $(1/2) E(w^T (\sum_{n \in N} x_n x_n^T) w - 2 w^T \sum_{n \in N} y_n x_n) + b_0$
 5. **For β 's update, $E[w^T (\sum_{n \in N} x_n x_n^T) w]$ and $E[w]$ will be needed**
- **Eq4: For $\alpha_d, d = 1, 2, \dots, d$**
 - $\log(q_{\alpha_d}^*(\alpha_d)) = E_{w, \beta} [\log p(y, w, \beta, \alpha_1, \dots, \alpha_D)/X]$
 - $\propto E[\log(\alpha_d/2 - w_d^2 \alpha_d/2 + (e_0 - 1) \log(\alpha_d) - f_0 \alpha_d)]$
 - $\propto (1/2 + e_0 - 1) \log \alpha_d - \alpha_d (f_0 + E[w_d^2]/2)$
 - The above form is the Gamma distribution with hyper parameters $1/2 + e_0$ and $f_0 + E[w_d^2]/2$; **dependent on $E[w_d^2]$ term**

Mean Field VI Algorithm

1. **Input:** Model in the form of priors and likelihood
2. **Output:** A variational distribution $q_w, q_\beta, q_{\alpha_1, \dots, \alpha_D}$
3. **ELBO Expression** = $E_q[\log P(Y, w^T X, \beta, \alpha_1, \dots, \alpha_D)] - E_q[\log q(w, \beta, \alpha_1, \dots, \alpha_D)]$
4. While the ELBO is not converged
 - $q_w, q_\beta, q_{\alpha_1, \dots, \alpha_D}$ are updated in an alternating fashion using ALT-OPT because as shown in above equations, update for $q_w, q_\beta, q_{\alpha_1, \dots, \alpha_D}$ are dependent on each other.

Student Name: Ajita Shree

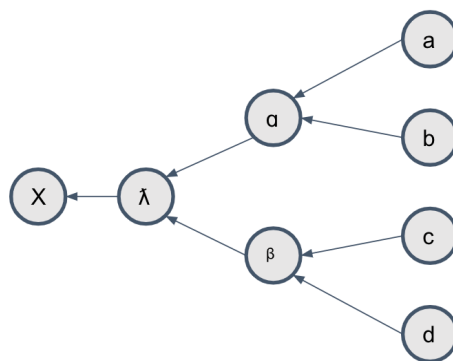
Roll Number: 20111262

Date: May 14, 2021

Gibbs Sampling

Below is the probability distribution for all the random variables and the graphical network.

- $p(x_n/\lambda_n) = \text{Poisson}(x_n/\lambda_n)$
- $p(\lambda_n/\alpha, \beta) = \text{Gamma}(\lambda_n/\alpha, \beta); n = 1, 2 \dots N$
- $p(\alpha/a, b) = \text{Gamma}(\alpha/a, b)$
- $p(\beta/c, d) = \text{Gamma}(\beta/c, d)$



Solution

Todo: derive the CP of all variable, if the CP are in closed form or not?

- Joint distribution, $p(X, \lambda, \alpha, \beta/a, b, c, d) = \prod_{n \in N} [p(x_n/\lambda_n)p(\lambda_n/\alpha, \beta)] * p(\alpha/a, b)p(\beta/c, d)$
- $p(\lambda_n/x_n, \alpha, \beta) = \text{Poisson}(x_n/\lambda_n) * \text{Gamma}(\lambda_n/\alpha, \beta)$
- Using the property of Conjugate priors, $p(\lambda_n/x_n, \alpha, \beta)$ will be a Negative binomial distribution with parameters $\alpha + x_n$ and $\beta + 1$, Hence closed form.
- $p(\beta/\lambda_1 \dots \lambda_n, \alpha,) = \prod_{n \in N} \text{Gamma}(\lambda_n/\alpha, \beta) * \text{Gamma}(\beta/c, d)$
- Using the property of Conjugate priors, $p(\beta/\lambda_1 \dots \lambda_n, \alpha, \beta)$ will be a CG distribution with parameters $c + na, d + \sum_{n \in N} \lambda_n$ Hence closed form.
- $p(\alpha/\lambda_1 \dots \lambda_n, \beta) = \prod_{n \in N} \text{Gamma}(\lambda_n/\alpha, \beta) * \text{Gamma}(\alpha/a, b)$
- The Closed form for CP of α is not available because the terms are not conjugate priors of each other.
- Ref: Wikipedia, Conjugate Priors

Student Name: Ajita Shree

Roll Number: 20111262

Date: May 14, 2021

Using Samples for Prediction

Given a matrix factorization $N * M$, R matrix where $p(r_{ij}/u_i, v_j) = N(r_{ij}/u_i^T v_j, \beta^{-1})$, where u_i and v_j are latent factors of i th and j th column. The PPD of r_{ij} is $p(r_{ij}/R) = p(r_{ij}/u_i, v_j)p(u_i, v_j/R)du_i dv_j$. Note we can write each $r_{ij} = u_i^T v_j + \epsilon_{ij}$ where $\epsilon_{ij} = N(\epsilon_{ij}/0, \beta^{-1})$. Also, we are given set of S samples $u^s, v^s_{s \in S}$ generated by Gibbs Sampler.

Solution

- **Derivation of sample based approximation of the mean**
- We know, for rv x with probability distribution $p(x)$, $E(x) = \int xp(x)dx$
- $E(r_{ij}) = \int r_{ij}p(r_{ij}/R)dr$, Computing it is intractable, but we have samples $U^s, V^s_{s \in S}$ generated by Gibbs Sampler for $p(u_i, v_j/R)$ and we also know that $r_{ij} = u_i^T v_j + \epsilon_{ij}$
- Expectation, $E[r_{ij}] = E[u_i^T v_j] + E[\epsilon_{ij}] = \int u_i^T v_j p(u_i, v_j/R)du_i dv_j + 0$; ($E[\epsilon_{ij}] = 0$)
- $E[r_{ij}] = \int u_i^T v_j \{u_i^s, v_j^s\}_{s \in S} du_i dv_j = \frac{\sum_{s \in S} u_i^{sT} v_j^s}{S}$
- Hence done
- **Derivation of sample based approximation of variance**
- $Var[r_{ij}] = E[r_{ij}^2] - E[r_{ij}]^2$
- Based on LOTUS, for a rv x with probability distribution $p(x)$ and its function $g(x)$, $E(g(x)) = \int g(x)p(x)$
- $E[r_{ij}^2] = E[(u_i^T v_j + \epsilon_{ij})^2] = E[(u_i^T v_j)^2] + E[\epsilon_{ij}^2] + E[2\epsilon_{ij}(u_i^T v_j)]$
- Using the fact $E[\epsilon_{ij}^2] = Var[\epsilon_{ij}] + E[\epsilon_{ij}]$, above eqn can be written as
- $E[r_{ij}^2] = \int (u_i^T v_j)^2 \{u_i^s, v_j^s\}_{s \in S} du_i dv_j + \beta^{-1} + 2E[\epsilon_{ij}u_i^T v_j]$
- Term 3: $E[\epsilon_{ij}u_i^T v_j] = \int \epsilon_{ij}u_i^T v_j \{u_i^s, v_j^s\}_{s \in S} du_i dv_j N(\epsilon_{ij}, 0, \beta^{-1})d\epsilon_{ij}$
- Above exp can be solved with identity, $\int x e^{x^2} = e^{x^2}/2 + C$
- Finally, $E[r_{ij}^2] = \frac{(\sum_{s \in S} u_i^{sT} v_j^s)^2}{S} + \beta^{-1} - 2\beta^{-1} \frac{\sum_{s \in S} u_i^{sT} v_j^s}{S}$
- $Var[\epsilon_{ij}] = \frac{(\sum_{s \in S} u_i^{sT} v_j^s)^2}{S} + \beta^{-1} - 2\beta^{-1} \frac{\sum_{s \in S} u_i^{sT} v_j^s}{S} - \left(\frac{(\sum_{s \in S} u_i^{sT} v_j^s)^2}{S} \right)^2$

Student Name: Ajita Shree

Roll Number: 20111262

Date: May 14, 2021

Rejection Sampling

Consider a distribution $p(x) \propto \exp(\sin(x))$ for $-\pi \leq x \leq \pi$, use proposal distribution $q(x) = N(x/0, \sigma^2)$

- The expression for M w.r.t. $Mq(x) \geq p'(x)$ is as follows:
- $M.N(x/0, \sigma^2) \geq \exp(\sin(x))$
- $M \frac{\exp(-0.5(\frac{x}{\sigma})^2)}{\sigma\sqrt{(2\pi)}} \geq \exp(\sin(x))$
- Rearranging terms, we will have $M \geq \sigma\sqrt{2\pi}\exp(\sin(x) + 0.5 * (\frac{x}{\sigma})^2)$
- Considering range of $x \in [-\pi, \pi]$, $M \geq \sigma\sqrt{2\pi}\exp(1 + 0.5 * (\frac{\pi}{\sigma})^2)$
- **Rejection Sampling Algorithm :**
 - sample x^* from $q(x)$
 - sample a uniform rv u $[0, Mq(x^*)]$
 - if $u \leq p(z^*)$, accept else reject
 - All accepted z^* 's are random samples from $p(x)$

Below plot is corresponding to the above algorithm and the M expression and sigma value = 1.5, 10000 samples are drawn from the distribution $p(x)$ and plotted below in the histogram with bin size 500. The M value calculated to be 91.36.

