

# DD2434/FDD3434 Machine Learning, Advanced Course

## Module 3 Exercise Solutions

November 2023

## Contents

<b>3 Variational Inference – Exercises</b>	<b>1</b>
3.1 KL-Divergence between 1D-Gaussians . . . . .	1
3.1.1 Solution . . . . .	1
3.2 ELBO derivation and motivation . . . . .	3
3.2.1 Solution . . . . .	3
3.3 Beta and Binomial model . . . . .	6
3.3.1 Solution . . . . .	6
3.4 Gaussian Mixture Model - light . . . . .	9
3.4.1 Solution . . . . .	9
3.5 Mixture Model with Bernoulli observations . . . . .	14
3.5.1 Solution . . . . .	14
3.6 Cartesian Matrix Model (from assignment 1B, 2017) . . . . .	14
3.6.1 Solution . . . . .	14
3.7 Troll factories (from assignment, 1B 2022) . . . . .	17
3.7.1 Solution . . . . .	17

## 3 Variational Inference – Exercises

### 3.1 KL-Divergence between 1D-Gaussians

Let  $p_1(x) = \mathcal{N}(\mu_1, \sigma_1^2)$  and  $p_2(x) = \mathcal{N}(\mu_2, \sigma_2^2)$ .

- a) Derive a closed form expression for  $\mathcal{D}_{KL}(p_1(x) \| p_2(x))$ .
- b) Implement a function for calculating the KL between two Gaussians based on your derivation. Implement a function that calculates the Euclidean norm, i.e.:  $D_{Euc}(p_1(x) \| p_2(x)) = \sqrt{\int (p_1(x) - p_2(x))^2 dx}$  (for stability reasons, you may want to integrate on a much smaller interval than  $[-\infty, \infty]$ ).

#### 3.1.1 Solution

# KL 1D-Gaussians

den 13 november 2024 11:23

$$D_{KL}(p_1 \parallel p_2) = \int_X p_1(x) \log \frac{p_1(x)}{p_2(x)} dx$$

$$= E_{p_1(x)} [\log p_1(x) - \log p_2(x)] = E_{p_1(x)} [\log p_1(x)] - E_{p_1(x)} [\log p_2(x)]$$

$$= \left\{ \log N(\mu, \sigma^2) = \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} (x-\mu)^2 \right\}$$

$$= E_{p_1(x)} \left[ -\frac{1}{2} \log 2\pi - \log \sigma_1 - \frac{1}{2\sigma_1^2} (x-\mu_1)^2 \right]$$

$$- E_{p_1(x)} \left[ -\frac{1}{2} \log 2\pi - \log \sigma_2 - \frac{1}{2\sigma_2^2} (x-\mu_2)^2 \right]$$

$$= -\log \sigma_1 + \log \sigma_2 - \frac{1}{2\sigma_1^2} \underbrace{E_{p_1(x)} [(x-\mu_1)^2]}_{= V(x) = \sigma_1^2} + \frac{1}{2\sigma_2^2} \underbrace{E_{p_1(x)} [(x-\mu_2)^2]}_{\equiv ①}$$

$$= \left\{ ① = E_{p_1(x)} [x^2 - 2x\mu_2 + \mu_2^2] \approx \underbrace{E_{p_1}[x^2]}_{= V(x) + E(x)^2} - 2[E_{p_1}[x]\mu_2 + \mu_2^2] = \sigma_1^2 + \mu_1^2 - 2\mu_1\mu_2 + \mu_2^2 \right\}$$

$$\approx \sigma_1^2 + \mu_1^2$$

$$= \log \frac{\sigma_2}{\sigma_1} - \frac{1}{2} + \frac{1}{2\sigma_2^2} \left( \sigma_1^2 + (\mu_1 - \mu_2)^2 \right)$$

### 3.2 ELBO derivation and motivation

Given some posterior  $p(Z|X)$  and variational distribution  $q(Z)$  (here X is the data and Z the set of latent variables and parameters):

a) show that  $\mathcal{D}_{KL}(q(Z)||p(Z|X))$  can be decomposed into the ELBO and log-marginal likelihood  $\log p(x)$ . *Hint: use equations (3)-(5) but with more detailed steps so that you understand the derivation better.*

b) Motivate why maximizing the ELBO implicitly minimizes the KL-divergence.

c) Alternative motivation for the ELBO:

Start with  $\log p(X)$  and use Jensen's inequality to show that  $\log p(X)$  is lower bounded by the ELBO.

d) In one sentence, describe qualitatively why it is easier to maximize the ELBO rather than minimizing the KL-divergence explicitly.

#### 3.2.1 Solution

# ELBO derivation and motivation

den 13 november 2024 13:41

$$a) D_{KL}(q(z) \| p(z|x)) = \int_z q(z) \log\left(\frac{q(z)}{p(z|x)}\right) dz$$

↑  
definition  
of KL-div.

$$= \int_z q(z) \log\left(\frac{q(z)}{\left(\frac{p(x,z)}{p(x)}\right)}\right) dz = \int_z q(z) \log\left(\frac{p(x)q(z)}{p(x,z)}\right) dz$$

↑ Bayes thm

$$= \int_z q(z) \log p(x) + q(z) \log\left(\frac{q(z)}{p(x,z)}\right) dz =$$

↑  
doesn't  
depend on z  
⇒ can be moved out of ∫      ↓  
=  $-\log\left(\frac{p(x,z)}{q(z)}\right)$

$$= \log p(x) \int_z q(z) dz - \int_z q(z) \log \frac{p(x,z)}{q(z)} dz =$$

↓  
integral over  
pdf = 1

$$= \log p(x) - \mathbb{E}_{q(z)} \left[ \log \frac{p(x,z)}{q(z)} \right]$$

↓  
 $= \text{ELBO}(q, p; x)$

b) Rearranging our equation above gives:

$$\log p(x) = D_{KL}(q(z) \parallel p(z|x)) + ELBO(x) \quad (*)$$

If we fix the parameters of  $p$ , and only maximize the parameters of  $q(z)$ , then  $\log p(x)$  is a constant. Therefore, as  $ELBO(x)$  increases,  $D_{KL}(q \parallel p)$  must decrease for  $(*)$  to hold.

c)

$$\begin{aligned} \log p(x) &= \log \int_z p(x,z) dz = \log \int_z \frac{q(z)}{q(z)} p(x,z) dz \\ &= \log E_{q(z)} \left[ \frac{p(x,z)}{q(z)} \right] \geq E_{q(z)} \left[ \log \frac{p(x,z)}{q(z)} \right] \end{aligned}$$

QED

d) The ELBO doesn't contain  $p(x)$ , which is usually intractable for complex models, as opposed to  $D_{KL}(q(z) \parallel p(z|x))$

### 3.3 Beta and Binomial model

Let  $X = (X_1, \dots, X_N)$  be i.i.d. where  $X_n|m, \theta \sim \text{Binomial}(m, \theta)$  and  $\theta \sim \text{Beta}(\alpha, \beta)$ .

- a) Derive the CAVI updates for  $q(\theta)$  using equation ??.
- b) How does this compare to the posterior in exercise 1.1 of Module 1? Describe qualitatively in one sentence why this is the case.

#### 3.3.1 Solution

# Exercises Module 3

Sunday, 15 October 2023

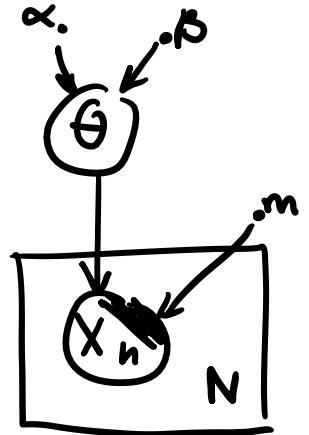
15:49

$$X = (X_1, \dots, X_N), \quad X_n | m, \theta \sim \text{Bin}(m, \theta)$$
$$\theta \sim \text{Beta}(\alpha, \beta)$$

a) CAVI update

Solution:

Can use fact:  $E_{\theta} [\log p(\theta, x)]$



1. Google  $p(X_n | m, \theta)$  and  $p(\theta)$  (Wikipedia page of Beta distribution and Binomial distribution usually works)

$$p(X_n = x_n | m, \theta) = \binom{m}{x_n} \theta^{x_n} (1 - \theta)^{m - x_n}$$

$$p(\theta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)}$$

2. Log form of  $p(X, \theta)$ :

$$\log p(X, \theta) = \log \prod_{n=1}^N \binom{m}{x_n} \theta^{x_n} (1 - \theta)^{m - x_n} \cdot \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)}$$

Since this is the expected value w.r.t. all variables except theta,  
But we have no other variables.

$$3. q^*(\theta) = \cancel{E_{-\theta} \left[ \log \prod_{n=1}^N \binom{m}{x_n} \theta^{x_n} \cdot (1-\theta)^{m-x_n} \cdot \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{\beta(\alpha, \beta)} \right]}$$

Only interested in the expression up to additive constant

$$\stackrel{+}{=} \sum_{n=1}^N \log \binom{m}{x_n} \theta^{x_n} (1-\theta)^{m-x_n} \cdot \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\begin{aligned} &= \sum_{n=1}^N \log \binom{m}{x_n} + x_n \log \theta + (m-x_n) \log (1-\theta) + (\alpha-1) \log \theta \\ &\quad + (\beta-1) \log (1-\theta) \stackrel{+}{=} \log \theta \cdot \left( \sum_{n=1}^N x_n + \alpha-1 \right) \\ &\quad + \log (1-\theta) \cdot (N-m - \sum x_n + \beta-1) \end{aligned}$$

Note Same as for the true posterior  
in Module 1. exercise. Why?

Two explanations:

1. since  $q^*(\theta) = \cancel{E_{-\theta} [\log p(X, \theta)]}$  simplifies to  $\log p(X, \theta)$ ,  
which is the same as when showing the conjugate-prior  
relation in Module 1.
2. Our mean-field assumption is only over one variable,  
hence, we are not making any simplifying assumption at all.

### 3.4 Gaussian Mixture Model - light

Here we will examine an simpler version of the Gaussian Mixture model. Still  $p(X_n|Z_n = k, \mu_k, \tau_k) = \text{Normal}(\mu_k, \frac{1}{\tau_k})$ ,  $p(Z_n|\pi) = \text{Categorical}(\pi)$ , but we assume  $\pi$  and  $\tau_k$  are given and let  $p(\mu_k) = \text{Normal}(\nu_k, \sigma_k)$ .

- a) Write the DGM/Bayes net for the model.
- b) Write out  $\log p(X, Z, \mu)$ .
- c) Apply and state the mean-field approximation for  $Z$  and  $\mu$ .
- d) Derive the associated CAVI updates using ??.
- e) Implement the CAVI algorithm ?? and apply it to simulated data using the generative model (If you are unfamiliar with this, it will be shown in the Exercise session of module 3). Try simulating data for different  $K$ ,  $N$ ,  $\nu_k$ ,  $\tau$  and  $\pi$  - under what circumstances does it have trouble finding all clusters?

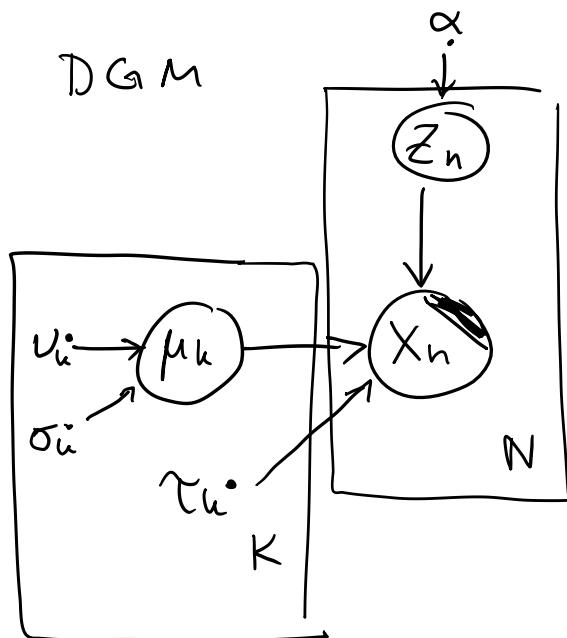
#### 3.4.1 Solution

# GMM - light

den 13 november 2023

11:34

a) DGM



$$b) \log p(x, z, \mu) = \sum_{n,k}^{N,K} \log p(x_n | z_n, \mu_k, \tau_k) \mathbb{1}_{\{z_n=k\}}$$

$$+ \log p(z | \pi) + \log p(\mu)$$

$$\log p(x | z, \mu) = \sum_{n=1}^{N,K} \mathbb{1}_{\{z_n=k\}} \left( \log \left( \frac{\tau_k}{\sqrt{2\pi}} \right) - \frac{\tau_k}{2} (x_n - \mu_k)^2 \right)$$

$$\log p(z | \pi) = \sum_{n=1, k=1}^{N, K} z_n^k \log \pi_k$$

$$\log p(\mu) = \log \prod_{k=1}^K p(\mu_k) = \sum_{k=1}^K \log \frac{1}{\sigma_k \sqrt{2\pi}} - \frac{1}{2\sigma_k^2} (\mu_k - \nu_k)^2$$

c) Mean-field approximation:  $q(z, \mu) = q(z)q(\mu)$

eq. 7.

$$\log q^*(z) = \mathbb{E}_z [\log p(x, z, \mu)] + \mathbb{E}_\mu [\log p(x|z, \mu)] \quad \textcircled{1}$$

$$+ \underbrace{\mathbb{E}_\mu [\log p(z|\pi)]}_{\textcircled{2}}$$

$$= \log p(z|\pi) \text{ since } p(z|\pi) \text{ not dep. on } \mu$$

$$\textcircled{1} = \mathbb{E}_\mu \left[ \sum_{n=1, k=1}^{N, K} \mathbb{1}_{\{z_n=k\}} \left( \log \frac{\tau_n}{\sqrt{2\pi}} - \frac{\tau_n}{2} (x_n - \mu_n)^2 \right) \right]$$

$$= \sum_{n, k}^{N, K} \mathbb{1}_{\{z_n=k\}} \left( \log \frac{\tau_n}{\sqrt{2\pi}} - \frac{\tau_n}{2} \mathbb{E}_\mu [(x_n - \mu_n)^2] \right)$$

$$\Rightarrow \textcircled{1} + \textcircled{2} =$$

$$\sum_{n, k}^{N, K} \mathbb{1}_{\{z_n=k\}} \left( \log \frac{\tau_n}{\sqrt{2\pi}} - \frac{\tau_n}{2} (x_n^2 - 2x_n \mathbb{E}_\mu [\mu_n] + \mathbb{E}_\mu [\mu_n^2] + \log \tau_n) \right)$$

$\equiv \log p_{nk}$

after we derive  
 $q(\mu)$  update we can get  
 expressions for these

$$= \sum_{n, k}^{N, K} \mathbb{1}_{\{z_n=k\}} \log p_{nk}$$

$$\Rightarrow q^*(z) \propto \prod_{n, k}^{N, K} p_{nk}^{\mathbb{1}_{\{z_n=k\}}} \Rightarrow q^*(z) = \prod_{n, k}^{N, K} r_{nk}^{\mathbb{1}_{\{z_n=k\}}},$$

i.e.  $q(z) = \prod_{n=1}^N q(z_n)$ , where  $r_{nk} = \frac{p_{nk}}{\sum_{k=1}^K p_{nk}}$

where each  $q(z_n) = \text{Cat}(r_{n1}, \dots, r_{nK})$

$$\log q(\mu)^* = E_{\mu} [\log p(x, z, \mu)] + E_z [\log p(x|z, \mu)] + \log p(\mu)$$

$$\begin{aligned}
 &= \sum_{n=1, k=1}^{N, K} E_z \left[ \mathbb{1}_{\{z_n=k\}} \left( \log \frac{q_{nk}}{\sqrt{2\pi}} - \frac{\tau_k}{2} (x_n - \mu_k)^2 \right) \right] + \sum_{k=1}^K \log \frac{1}{\sigma_k \sqrt{2\pi}} - \frac{1}{2\sigma_k^2} (\mu_k - v_k)^2 \\
 &\stackrel{+}{=} \sum_{k=1}^K \sum_{n=1}^N \underbrace{-E_z [\mathbb{1}_{\{z_n=k\}}]}_{=q(z_n=k)} \frac{\tau_k}{2} \left( x_n^2 - 2\mu_k x_n + \mu_k^2 \right) - \frac{1}{2\sigma_k^2} (\mu_k^2 - 2\mu_k v_k + v_k^2) \\
 &= \sum_{k=1}^K -\frac{1}{2} \left( \underbrace{\left( \sum_{n=1}^N q(z_n=k) \cdot x_n \right) \tau_k}_{\equiv A_k^N} \cdot 2 \cdot \mu_k + \left( \sum_{n=1}^N q(z_n=k) \right) \tau_k \mu_k^2 \right) - \frac{1}{2} \left( \underbrace{\frac{1}{\sigma_k^2} \cdot 2 \cdot \mu_k v_k + \frac{1}{\sigma_k^2} \mu_k^2}_{\equiv B_k^N} \right) \\
 &= \sum_{k=1}^K -\frac{1}{2} \left( -2 \cdot \mu_k \left( A_k^N \cdot \tau_k + \frac{v_k}{\sigma_k^2} \right) + \mu_k^2 \left( B_k^N \cdot \tau_k + \frac{1}{\sigma_k^2} \right) \right) \\
 &= \sum_{k=1}^K -\frac{1}{2} \left( \underbrace{-2 \mu_k}_{\sigma_k^{2*}} \underbrace{\left( A_k^N \cdot \tau_k + \frac{v_k}{\sigma_k^2} \right)}_{B_k^N \cdot \tau_k + \frac{1}{\sigma_k^2}} + \mu_k^2 \underbrace{\left( B_k^N \cdot \tau_k + \frac{1}{\sigma_k^2} \right)}_{v_k^{*}} \right)
 \end{aligned}$$

$$\Rightarrow q(\mu) = \prod_{k=1}^K q(\mu_k), \quad q(\mu_k) = N(v_k^*, \sigma_k^{2*})$$

$$\Rightarrow \begin{cases} \mathbb{E}_\mu[\mu_n] = \nu_k^*, \\ \mathbb{E}_{\mu_n}[\mu_n^2] = \text{Var}(\mu_n) - \mathbb{E}[\mu_n]^2 = \sigma_k^{*2} - \nu_k^{*2} \end{cases}$$

d) See Bernoulli Mixture Model solution  
 for how to implement CAVI updates in  
 VI-algorithm.

### 3.5 Mixture Model with Bernoulli observations

In the video lectures the CAVI updates for a Mixture model with Gaussian observation model is introduced, i.e.,  $p(X_n|Z_n = k, \mu_k, \tau_k) = \text{Normal}(\mu_k, \frac{1}{\tau_k})$ ,  $p(Z_n|\pi) = \text{Categorical}(\pi)$ ,  $p(\pi) = \text{Dirichlet}(\alpha)$ .

In this exercise we examine the Bernoulli Mixture Model, with same priors for  $Z_n|\pi$  and  $\pi$ , but with  $X_n = \{X_{n1}, \dots, X_{nD}\}$  with observational model  $p(X_{nd}|Z_n = k, \theta_{kd}) = \text{Bernoulli}(\theta_{kd})$  with prior  $p(\theta) = \prod_k \prod_d p(\theta_{nd})$ , where  $p(\theta_{ka}) = \text{Beta}(a, b)$ .

- a) Write the DGM/Bayes net for the model.
- b) Write out  $\log p(X, Z, \pi, \theta)$ .
- c) Apply and state the mean-field approximation for  $Z$ ,  $\pi$  and  $\theta$ .
- d) Derive the associated CAVI updates using ??.
- e) Derive a closed form expression for the ELBO.
- f) Implement the CAVI algorithm ?? and apply it to simulated data using the generative model. Try simulating data for different K, N,  $\theta_k$  and  $\pi$  - under what circumstances does it have trouble finding all clusters?

#### 3.5.1 Solution

See pdf in Canvas.

### 3.6 Cartesian Matrix Model (from assignment 1B, 2017)

The Cartesian Matrix Model (CMM) is defined as follows. There are  $R$  row distributions  $\{N(\mu_r, \lambda_r^{-1}) : 1 \leq r \leq R\}$ , each variance  $\lambda_r^{-1}$  is known and each  $\mu_r$  has prior distribution  $N(\mu, \lambda^{-1})$ . There are also  $C$  column distributions  $\{N(\xi_c, \tau_c^{-1}) : 1 \leq c \leq C\}$ , each variance  $\tau_c^{-1}$  is known and each  $\xi_c$  has prior distribution  $N(\xi, \tau^{-1})$ . All hyper-parameters are known. A matrix  $S$  is generated by, for each row  $1 \leq r \leq R$  and each column  $1 \leq c \leq C$ , setting  $S_{rc} = X_r + Y_c$  where  $X_r$  is sampled from  $N(\mu_r, \lambda_r^{-1})$  and  $Y_c$  from  $N(\xi_c, \tau_c^{-1})$ . Use Variational Inference in order to obtain a variational distribution

$$q(\mu_1, \dots, \mu_R, \xi_1, \dots, \xi_C) = \prod_r q(\mu_r) \prod_c q(\xi_c)$$

that approximates  $p(\mu_1, \dots, \mu_R, \xi_1, \dots, \xi_C | S)$ . Tip: what distribution do you get from the sum of two Gaussian random variables? What is the relation between the means?

**Question 15:** Present the algorithm written down in a formal manner (using both text and mathematical notation, but not pseudo code).

Figure 1: From Assignment 2, 2017

#### 3.6.1 Solution

##### Solution:

Since each cell of the matrix denoted as  $S_{rc}$  is a sum of two Gaussians then as we know we can get the Gaussian result where the mean is the sum of the mean of  $X$  and the mean of  $Y$ . The variance is also summed. So by having  $\mu_r$  and  $\xi_c$ , the distribution of each cell can be

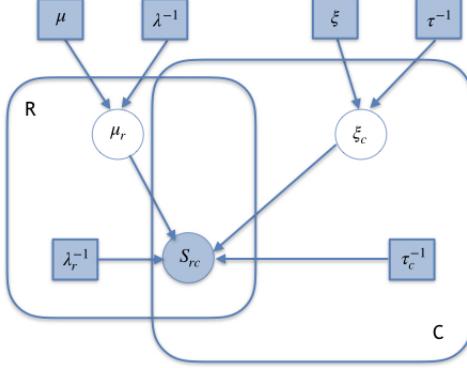


Figure 2: The graphical model (solution).

expressed as  $\mathcal{N}(\mu_r + \xi_c, \lambda_r^{-1} + \tau_c^{-1})$ . In Fig. 2, the graphical model is illustrated. Below are the calculations of  $q_l(\mu_l)$  and  $q_t(\xi_t)$  according to variational approximation ( $l$  and  $t$  represent any relevant index for calculating all of the possible  $q_l(\mu_l)$  and  $q_t(\xi_t)$ ). The known parameters are represented with  $\Theta = \{\lambda_{1:R}^{-1}, \lambda^{-1}, \mu, \tau_{1:C}^{-1}, \tau^{-1}, \xi\}$ .

$$\begin{aligned}\log q(\mu_l) &= E_{\substack{\xi_{1:C} \\ \mu_{1:R-l}}} [\log p(S, \mu_{1:R}, \xi_{1:C} | \Theta)] + \text{const} \\ \log q(\xi_t) &= E_{\substack{\mu_{1:R} \\ \xi_{1:C-t}}} [\log p(S, \mu_{1:R}, \xi_{1:C} | \Theta)] + \text{const}\end{aligned}$$

To calculate the above we first calculate  $\log q_l(\mu_l)$  which needs in return to first calculate  $p(S, \mu_{1:R}, \xi_{1:C})$  (Note that we choose indexes  $l$  and later  $t$  to denote a specific index so that it is not mixed with the indexes in the sum and product terms):

$$\begin{aligned}p(S, \mu_{1:R}, \xi_{1:C} | \Theta) &= p(S | \mu_{1:R}, \xi_{1:C}, \Theta) p(\mu_{1:R}, \xi_{1:C} | \Theta) \\ &= \prod_{c=1}^C \prod_{r=1}^R p(S_{rc} | \mu_r + \xi_c, \lambda_r^{-1} + \tau_c^{-1}) \prod_{r=1}^R p(\mu_r | \mu, \lambda^{-1}) \prod_{c=1}^C p(\xi_c | \xi, \tau^{-1})\end{aligned}$$

Notice that each  $\mu_r$  and each  $\xi_c$  are independent in the above expression, resulting in factors. Now we calculate  $\log q_l(\mu_l)$  by taking the expectation over the log of the expression above:

$$\log q(\mu_l) = -\frac{1}{2} E_{-\mu_l} \left[ \sum_{c=1}^C \sum_{r=1}^R \frac{(S_{rc} - (\mu_r + \xi_c))^2}{\lambda_r^{-1} + \tau_c^{-1}} + \frac{1}{\lambda^{-1}} \sum_{r=1}^R (\mu_r - \mu)^2 + \frac{1}{\tau^{-1}} \sum_{c=1}^C (\xi_c - \xi)^2 \right] + \text{const}$$

All the components of  $r$  in  $\sum_{r=1}^R$  are independent on  $\mu_l$  except  $r = l$ ; adding those independent

terms in the above into constant results in:

$$-\frac{1}{2}E_{\xi_{1:C}} \left[ \sum_{c=1}^C \frac{S_{lc}^2 - 2S_{lc}(\mu_l + \xi_c) + (\mu_l^2 + \xi_c^2 + 2\mu_l\xi_c)}{\lambda_l^{-1} + \tau_c^{-1}} + \lambda(\mu_l^2 + \mu^2 - 2\mu\mu_l) + \frac{1}{\tau^{-1}} \sum_{c=1}^C (\xi_c - \xi)^2 \right] \\ + const$$

The  $\mu_l$ -independent terms in the above expression can be pushed into constant resulting in (the expectation goes through the terms as below):

$$\begin{aligned} \log q(\mu_l) &= -\frac{1}{2}E_{\xi_{1:C}} \left[ \sum_{c=1}^C \frac{-2S_{lc}(\mu_l) + (\mu_l^2 + 2\mu_l\xi_c)}{\lambda_l^{-1} + \tau_c^{-1}} + \lambda(\mu_l^2 - 2\mu\mu_l) \right] + const \\ &= -\frac{1}{2}E_{\xi_{1:C}} \left[ \left( \lambda + \sum_{c=1}^C \frac{1}{\lambda_l^{-1} + \tau_c^{-1}} \right) \mu_l^2 - 2\mu_l \left( \mu\lambda + \sum_{c=1}^C E_{\xi_c} \left[ \frac{S_{lc} - \xi_c}{\lambda_l^{-1} + \tau_c^{-1}} \right] \right) \right] + const \\ &= -\frac{1}{2} \left[ \mu_l^2 \left( \lambda + \sum_{c=1}^C \frac{1}{\lambda_l^{-1} + \tau_c^{-1}} \right) - 2\mu_l \left( \mu\lambda + \sum_{c=1}^C E_{\xi_c} \left[ \frac{S_{lc} - \xi_c}{\lambda_l^{-1} + \tau_c^{-1}} \right] \right) \right] + const \\ &= -\frac{1}{2} \left[ \mu_l^2 \left( \lambda + \sum_{c=1}^C \frac{1}{\lambda_l^{-1} + \tau_c^{-1}} \right) - 2\mu_l \left( \mu\lambda + \sum_{c=1}^C \frac{S_{lc} - E_{\xi_c}[\xi_c]}{\lambda_l^{-1} + \tau_c^{-1}} \right) \right] + const \end{aligned}$$

By completing the square, the above expression becomes in a form of Gaussian with mean  $m_l$  and precision  $p_l$ , where

$$\begin{aligned} p_l &= \lambda + \sum_{c=1}^C \frac{1}{\lambda_l^{-1} + \tau_c^{-1}} \\ m_l &= \frac{\mu\lambda + \sum_{c=1}^C \frac{S_{lc} - E_{\xi_c}[\xi_c]}{\lambda_l^{-1} + \tau_c^{-1}}}{p_l} \\ q(\mu_l) &= \mathcal{N}(\mu_l | m_l, p_l) \end{aligned}$$

Similarly we can calculate  $\log q_t(\xi_t)$  and so it results in a Gaussian with mean  $m_t$  and precision  $p_t$ , where

$$\begin{aligned} p_t &= \tau + \sum_{r=1}^R \frac{1}{\tau_t^{-1} + \lambda_r^{-1}} \\ m_t &= \frac{\xi\tau + \sum_{r=1}^R \frac{S_{rt} - E_{\mu_r}[\mu_r]}{\tau_t^{-1} + \lambda_r^{-1}}}{p_t} \\ q(\xi_t) &= \mathcal{N}(\xi_t | m_t, p_t) \end{aligned}$$

Since we have approximated  $q(\mu_l)$ , we have R number of estimations each having similar result as in  $m_l$  and  $p_l$ . Having  $q(\xi_t)$  form, we have C number of estimations each having similar result as in  $m_t$  and  $p_t$ .

The only remaining parts to compute are  $E_{\xi_c}[\xi_c]$  and  $E_{\mu_r}[\mu_r]$ . Since  $\xi_c$  and  $\mu_r$  are Gaussian random variables, we can substitute the expression of the means

$$m_l = \frac{\mu\lambda + \sum_{c=1}^C \frac{S_{lc} - m_c}{\lambda_l^{-1} + \tau_c^{-1}}}{p_l}$$

$$m_t = \frac{\xi\tau + \sum_{r=1}^R \frac{S_{rt} - m_r}{\tau_t^{-1} + \lambda_r^{-1}}}{p_t}$$

We found the distributions of  $q(\mu_l)$  and  $q(\xi_t)$ , therefore we can calculate  $\prod_{r=1}^R q(\mu_r) \prod_{c=1}^C q(\xi_c)$ .

### 3.7 Troll factories (from assignment, 1B 2022)

On a social media platform, K troll factories have posted N comments on a live news report from an ongoing war. A security agency wants to extract information on the troll factories, but due to integrity protection policies, the platform can only provide metadata of the posts, such as comment length  $X_n$  of each post as well as response time  $T_n$ . Together with the security agency's disinformation team, the newly employed ML expert develops a model which infers comment to factory assignment,  $Z_n$ , factory post volume fraction,  $\pi$ , troll factory specific response rate,  $\lambda_k$ , and average comment length,  $\mu_k$ , and precision  $\tau_k$  with the following distributions:

- $X_n|\mu_k, \tau_k, Z_n = k \sim Lognormal(\mu_k, \tau_k^{-1})$  - based on the assumptions that each troll factory has its own strategy for comment length and variation in length and that comments are always of positive length.
- $\mu_k, \tau_k | \nu, \kappa, \alpha, \beta \sim NormalGamma(\nu, \kappa, \alpha, \beta)$
- $T_n|\lambda_k, Z_n = k \sim Exp(\lambda_k)$  - Comments are written as reactions to events with a factory specific response rate.
- $\lambda_k | a, b \sim Gamma(a, b)$  - The factory specific response rate is unknown, but the domain experts provide reasonable values for  $a$  and  $b$ .
- $Z_n | \pi \sim Categorical(\pi)$  - Each post is associated with K different factories.
- $\pi | \delta \sim Dirichlet(\delta)$

Note that comment length is a discrete entity, but we approximate the likelihood of observations with a continuous distribution in the model.

- Provide a graphical model for the model described above.
- Derive the CAVI update equations of each variational distribution.

#### 3.7.1 Solution

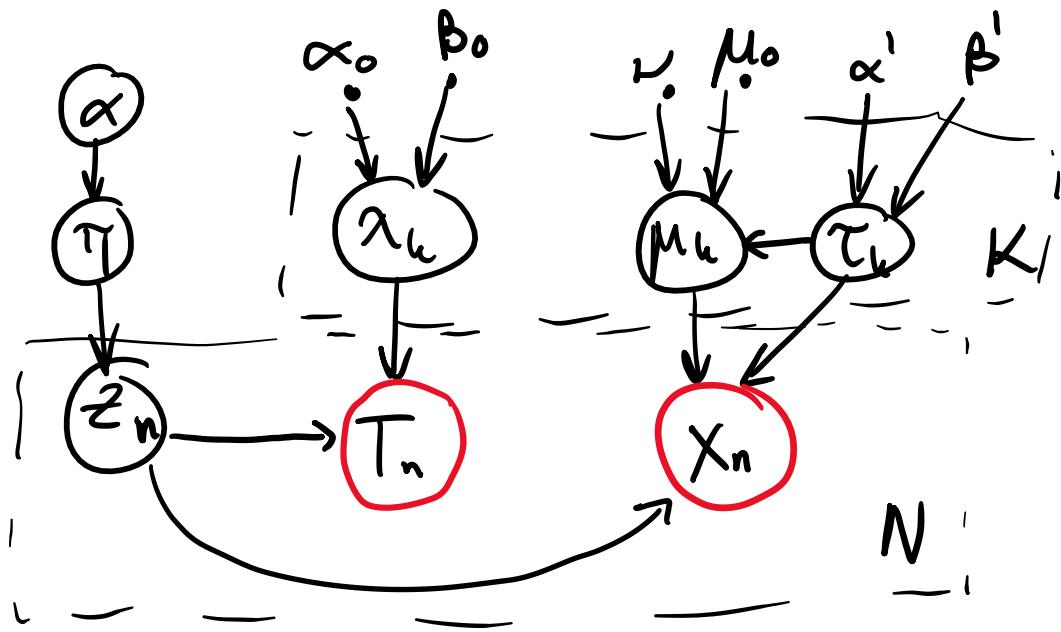
# Troll factories

Friday, 10 November 2023 12:09

This is a problem from last years 1B.

All updates are not shown as they give too much information on this assignment.

A) Draw the DGM



B) Derive the CAVI updates

$$\ln \text{LogNormal}(x_n | \mu_k) = -\ln x_n + \frac{\ln \tau_k}{2} - \frac{\tau_k (\ln x_n - \mu_k)^2}{2} - \frac{\ln 2\pi}{2}$$

$$\ln \text{Exp}(t_n | \lambda_k) = \ln \lambda_k - \lambda_k t_n$$

$$\ln P(\pi | \alpha) = \sum_{k=1}^K (\alpha - 1) \ln(\pi_k) + \text{Const.}$$

$$\ln P(\lambda_k | \alpha_0, \beta_0) = (\alpha_0 - 1) \ln \lambda_k - \beta_0 \lambda_k + \text{Const.}$$

$$\ln P(\mu_k, \tau_k | \nu, \mu'_0, \alpha', \beta') = (\alpha' - \frac{1}{2}) \ln \tau_k - \beta' \tau_k - \frac{\nu \tau_k (\mu_k - \mu')^2}{2} + \text{Const.}$$

Denoting joint probability  $\phi = P(\mathbf{X}, \mathbf{S}, \mathbf{Z}, \pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T})$ ,

$$\begin{aligned}\phi &= P(\mathbf{X}|\mathbf{Z}, \mathbf{M}, \mathbf{T})P(\mathbb{T}|\mathbf{Z}, \boldsymbol{\Lambda})P(\mathbf{Z}|\pi)P(\pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T}) \\ \ln \phi &= \ln P(\mathbf{X}|\mathbf{Z}, \mathbf{M}, \mathbf{T}) + \ln P(\mathbb{T}|\mathbf{Z}, \boldsymbol{\Lambda}) + \ln P(\mathbf{Z}|\pi) + \ln P(\pi) + \ln P(\boldsymbol{\Lambda}) + \ln P(\mathbf{M}, \mathbf{T}) \\ &= (\sum_{n=1}^N \sum_{k=1}^K z_{nk} [\ln \text{LogNormal}(\mathbf{x}_n | \mu_k, \tau_k^{-1}) + \ln \text{Exp}(\mathbf{t}_n | \lambda_k) + \ln \pi_k]) + \\ &\quad \ln P(\pi | \alpha) + \sum_{k=1}^K (\ln P(\lambda_k | \alpha_0, \beta_0) + \ln P(\mu_k, \tau_k | \nu, \mu', \alpha', \beta'))\end{aligned}$$

$$\begin{aligned}\ln q_{\mathbf{Z}}(\mathbf{Z}) &= \mathbb{E}_{\pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T}}[\ln \phi] \\ &= \mathbb{E}_{\pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T}}[(\sum_{n=1}^N \sum_{k=1}^K z_{nk} [\ln \text{LogNormal}(\mathbf{x}_n | \mu_k, \tau_k^{-1}) + \ln \text{Exp}(\mathbf{t}_n | \lambda_k) + \ln \pi_k]) \\ &\quad + \ln P(\pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T})] \\ &= \sum_{n=1}^N \sum_{k=1}^K z_{nk} [\mathbb{E}_{\mu_k, \tau_k}(\ln \text{LogNormal}(\mathbf{x}_n | \mu_k, \tau_k^{-1})) + \mathbb{E}_{\lambda_k}(\ln \text{Exp}(\mathbf{t}_n | \lambda_k)) + \mathbb{E}_{\pi_k}(\ln \pi_k)] + \text{Const}\end{aligned}$$

Since the term  $\ln P(\pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T})$  is not a function of  $\mathbf{Z}$ , it is a constant. Moreover  $q_{\mathbf{Z}}(\mathbf{Z})$  can be further broken down to  $\prod_{n=1}^N q_{\mathbf{z}_n}(\mathbf{z}_n)$ . Thus,

$$q_{\mathbf{z}_n}(\mathbf{z}_n) = \sum_{k=1}^K z_{nk} [\mathbb{E}_{\mu_k, \tau_k}(\ln \text{LogNormal}(\mathbf{x}_n | \mu_k, \tau_k^{-1})) + \mathbb{E}_{\lambda_k}(\ln \text{Exp}(\mathbf{t}_n | \lambda_k)) + \mathbb{E}_{\pi_k}(\ln \pi_k)] + \text{Const}$$

Now evaluating each of the expectation terms using the above equations,

$$\begin{aligned}\mathbb{E}_{\mu_k, \tau_k}(\ln \text{LogNormal}(\mathbf{x}_n | \mu_k, \tau_k^{-1})) &= \frac{1}{2} \mathbb{E}_{\tau_k}[\ln \tau_k] - \frac{1}{2} \mathbb{E}_{\mu_k, \tau_k}[\tau_k (\ln x_n - \mu_k)^2] + \text{Const.} \\ \mathbb{E}_{\lambda_k}(\ln \text{Exp}(\mathbf{t}_n | \lambda_k)) &= \mathbb{E}_{\lambda_k}[\ln \lambda_k] - \mathbb{E}_{\lambda_k}[\lambda_k] t_n\end{aligned}$$

Now for the term  $q_{\pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T}}(\pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T})$

$$\begin{aligned}\ln q_{\pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T}}(\pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T}) &= \mathbb{E}_{\mathbf{Z}}[\ln \phi] \\ &= \mathbb{E}_{\mathbf{Z}}[(\sum_{n=1}^N \sum_{k=1}^K z_{nk} [\ln \text{LogNormal}(\mathbf{x}_n | \mu_k, \tau_k^{-1}) + \ln \text{Exp}(\mathbf{t}_n | \lambda_k) + \ln \pi_k]) \\ &\quad + \ln P(\pi, \boldsymbol{\Lambda}, \mathbf{M}, \mathbf{T})] \\ &= \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}_{z_{nk}}[z_{nk} [\ln \text{LogNormal}(\mathbf{x}_n | \mu_k, \tau_k^{-1}) + \ln \text{Exp}(\mathbf{t}_n | \lambda_k) + \ln \pi_k] \\ &\quad + \ln P(\pi) + \ln P(\boldsymbol{\Lambda}) + \ln P(\mathbf{M}, \mathbf{T})] \\ &= \sum_{n=1}^N \sum_{k=1}^K r_{nk} [\ln \text{LogNormal}(\mathbf{x}_n | \mu_k, \tau_k^{-1}) + \ln \text{Exp}(\mathbf{t}_n | \lambda_k) + \ln \pi_k] \\ &\quad + \sum_{k=1}^K (\alpha - 1) \ln \pi_k + \sum_{k=1}^K [(\alpha_0 - 1) \ln \lambda_k - \beta_0 \lambda_k] \\ &\quad + \sum_{k=1}^K [(\alpha' - \frac{1}{2}) \ln \tau_k - \beta' \tau_k - \frac{\nu \tau_k (\mu_k - \mu')^2}{2}] + \text{Const}\end{aligned}$$

The above expression brings out the underlying independence, which is  $\Pi$ ,  $\Lambda$  and  $(\mathbf{M}, \mathbf{T})$ , don't appear together or in pairs, and the fact that for  $k \neq k'$ , the terms corresponding to  $k$  and  $k'$  don't appear together. Therefore it can be further factorised as  $q_{\pi, \Lambda, \mathbf{M}, \mathbf{T}}(\pi, \Lambda, \mathbf{M}, \mathbf{T}) = q_\pi(\pi)q_\Lambda(\Lambda)q_{\mathbf{M}, \mathbf{T}}(\mathbf{M}, \mathbf{T}) = q_\pi(\pi)\prod_{k=1}^K q_{\lambda_k}(\lambda_k)q_{\mu_{\mathbf{k}}, \tau_{\mathbf{k}}}(\mu_{\mathbf{k}}, \tau_{\mathbf{k}})$ . We can use the above equations to get each of the factorised units.

$$\begin{aligned}\ln q_\pi(\pi) &= \sum_{k=1}^K (\alpha - 1) \ln \pi_k + \sum_{n=1}^N \sum_{k=1}^K r_{nk} \ln \pi_k \\ &= \sum_{k=1}^K [(\alpha - 1) + N_k] \ln \pi_k \\ q_\pi(\pi) &= \text{Dir}(\pi | \alpha + N_k)\end{aligned}$$

where  $N_k = \sum_{n=1}^N r_{nk}$ . For  $q_{\lambda_k}(\lambda_k)$ ,

$$\begin{aligned}\ln q_{\lambda_k}(\lambda_k) &= (\alpha_0 - 1) \ln \lambda_k - \beta_0 \lambda_k + \sum_{n=1}^N r_{nk} [\ln \lambda_k - t_n \lambda_k] \\ &= \ln \lambda_k [\alpha_0 - 1 + N_k] - (\beta_0 + \sum_{n=1}^N r_{nk} t_n) \lambda_k + \text{Const} \\ q_{\lambda_k}(\lambda_k) &= \text{Gamma}(\alpha_0 + N_k, \beta_0 + \sum_{n=1}^N t_n r_{nk})\end{aligned}$$