



KTH ROYAL INSTITUTE OF TECHNOLOGY

DD2434  
MACHINE LEARNING ADVANCED COURSE

YEAR 2017 - 2018

---

## Assignment 2

---

INDIVIDUAL REPORT

PICÓ ORISTRELL Sandra

940531 – 2308 sandrapo@kth.se

*Professors :*

HERMAN Pawel

KJELLSTRÖM Hedvig

LAGERGREN Jens

January 31, 2018

# Contents

<b>1</b>	<b>Dependences in a Directed Graphical Model</b>	<b>2</b>
<b>2</b>	<b>The Sum-HMM</b>	<b>4</b>
2.0.1	Non-biased dices . . . . .	4
2.0.2	Medium biased dice . . . . .	5
2.0.3	Total biased dice . . . . .	5
<b>3</b>	<b>Simple VI</b>	<b>5</b>
<b>4</b>	<b>Variational Inference</b>	<b>9</b>
	<b>Appendices</b>	<b>11</b>
<b>A</b>	<b>Code question 3</b>	<b>11</b>
<b>B</b>	<b>Code question 10</b>	<b>12</b>
	<b>References</b>	<b>13</b>

# 1 Dependences in a Directed Graphical Model

**Important note:** In order to get a **E** in this assignment I decided to complete Task 2.1,2.3 and 2.6.

Considering the following Directed Acyclic Graph (DAG) shown below, answer the next two questions:

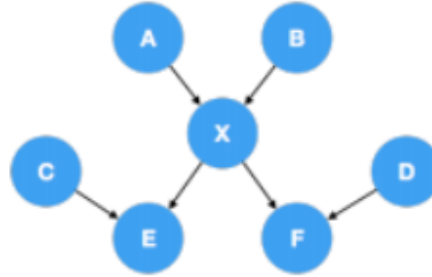


Figure 1: The DAG

**Question 1:** Which pairs of variables, not including  $X$ , are dependent conditioned on  $X$ ?

In the figure attached above we can see a Directed Acyclic Graphical Model or also called Bayesian Network. In this kind of models, the arrow represents that one node has direct influence over the other. For example, in the graph of this exercise we can state that the node  $D$  has direct influence over  $F$ .

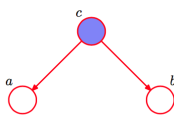
Another term that is important to define is the *conditional independence* term. As Bishop say [1], conditional independence is an important concept for probability distributions over multiple variables. An example of two variables that are conditionally independent is the following:

If:

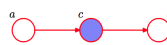
$$p(a|b, c) = p(a|c) \quad (1)$$

then, we state that  $a$  is conditionally independent of  $b$  given  $c$ .

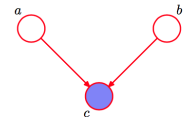
At the same time, there are different graph structures that verify the property described above. However, as explained in [1], there are three specific configurations that can represent the property mentioned. That ones are the following ones:



(a)  $C$  as a tail-to-tail node



(b)  $C$  as a head-to-tail node



(c)  $C$  as a head-to-head node

Figure 2: Conditional independence property in  $C$  node

It is also important to state that the mathematical notation used to define conditional independence property is:

$$a \perp\!\!\!\perp b | c \quad (2)$$

For this reason, based on all explained above, the pairs of variables that appear in the graph that are dependent conditioned on X ( not including X) is only one: [A,B].

$$A \not\perp B|X$$

However, whether conditioned on X or not, we can also find two other pairs of variables that are dependent:

$$E \not\perp C|X \quad F \not\perp D|X$$

To sum up the question, we can then state that the pairs of variables that dependent conditioned on X (not including X) are:

- (A,B)
- (E,C)
- (F,D)

**Question 2:** Which pairs of variables, not including X, are dependent, not conditioned on X?

In order to analyze the graph in a proper way, we are going to study all the different pairs of variables by steps.

First of all, based on all the information already explained in *question 1*, the node X is *tail-to-tail* respect to E and F because the node is connected to the tails of the two arrows. For this reason, if we write the conditional distribution of E and F, given X, we are going to obtain the conditional independence property. However, this question is related to be *not conditioned on X*. So, if we are in the situation where X is not observed, we can test to see if E and F are independent by marginalizing over X. This condition will imply the following:

$$E \not\perp F|\emptyset \tag{3}$$

where  $\emptyset$  means not being observed anything. Based on the justification above, we can state that the pair [E,F] are dependent, not conditioned on X.

To continue, now we are going to analyze the *head-to-tail* nodes. In the graph, we find the four pairs of variables that have this configuration related to X: [A,F] , [B,E], [A,E] and [B,F]. Therefore, based on the same reasoning that above, we can state:

$$A \not\perp F|\emptyset \quad B \not\perp E|\emptyset \quad A \not\perp E|\emptyset \quad B \not\perp F|\emptyset$$

And, finally, as we did in the previous question, we can also state that E and F are dependent to C and D respectively, either conditioned or not conditioned on X.

$$E \not\perp C|\emptyset \quad F \not\perp D|\emptyset$$

To sum up, based on all explained above, the pairs of variables that appear in the graph that are dependent not conditioned on X(not including X) are:

- (E,F)
- (A,F)
- (B,E)
- (B,F)
- (A,E)
- (E,C)
- (F,D)

## 2 The Sum-HMM

**Question 3:** Implement the Sum-HMM, i.e., write your own code for it.

I attached the code in the appendix section.

**Question 4:** Provide data generated using at least three different sets of categorical dice distributions that provide reasonable tests for the correctness of your program.

**Question 5:** Motivate your test and why the result of it indicates correctness.

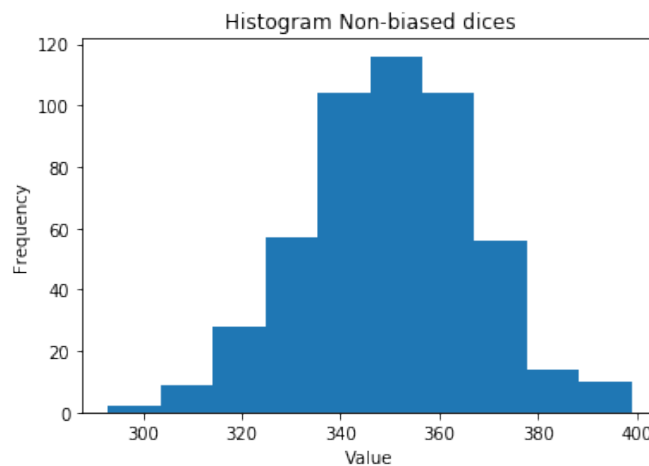
I will answer question 4 and question 5 together. In these questions, in order to check if the program is working properly, we are going to test different cases. Specifically, the tests that we want to check are the following ones:

- Non-biased dice.
- Medium biased dice.
- Total biased dice.

We have to take into account that the original problem was simplified. We are going to consider that the two groups of tables are going to have the same categorical distribution, i.e. behave in the same manner. Also, the number of players will be 500 and the number of tables will be 100 for all the cases tested.

### 2.0.1 Non-biased dices

In the non-biased dices, as the name indicates, we will have the same probability per each number. In other words,  $p(\text{number}) = \frac{1}{6}$ .



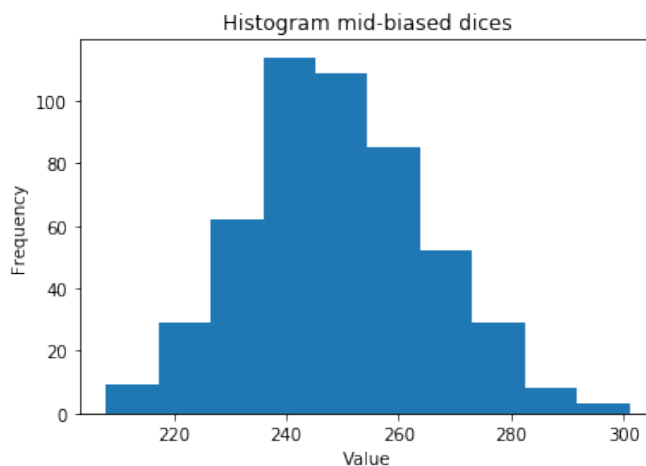
The previous figure reveals the result obtained in the case described. Concretely, we can state that it is similar to the Normal distribution. The reason behind that is the central limit theorem. As the number of players is high enough, the random variable tends to behave as the Normal distribution.

Moreover, if we analyze the mean of the distribution we realize that it is in the "middle" by the simple fact that we are studying the case where the two groups of tables have uniform distribution.

## 2.0.2 Medium biased dice

This case, we will have different distributions per each group of tables. The first group will behave as a uniform distribution and the second group will act with a probability of 0.5 for the numbers 1 and 2.

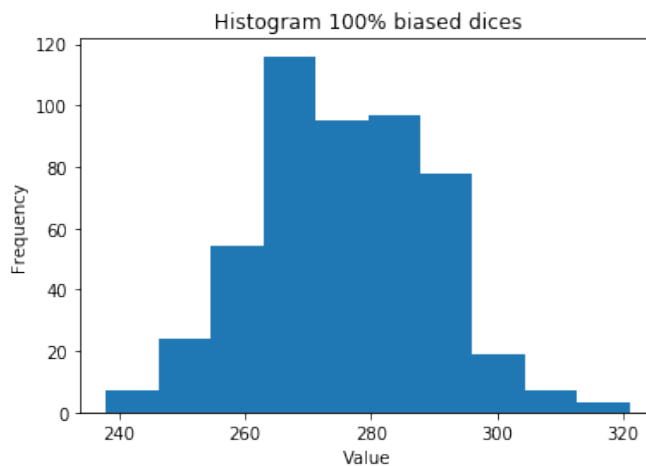
The result of this case is:



If we compare the result obtained in the previous section (2.0.1) with the one obtained now we are going to realize that the mean is displaced. The reason is the distribution for the second group of tables.

## 2.0.3 Total biased dice

In this case, the second group of tables will only obtain the number 2. In other words, the probability of having a 2 will be  $\frac{6}{6}$ . However, regarding the first group of tables we will have a uniform distribution. That is the reason why we obtain the following histogram moved to the left side:



## 3 Simple VI

Consider the model defined by the following equations:

Likelihood function given by:

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{\frac{N}{2}} \exp - \frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 \quad (4)$$

Conjugate prior distributions for  $\mu$  and  $\tau$  given by:

$$p(\mu|\tau) = \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1}) \quad (5)$$

$$p(\tau) = \text{Gam}(\tau|a_0, b_0) \quad (6)$$

where Gam is the gamma distribution defined by:

$$\text{Gam}(\lambda|a, b) = \frac{1}{(a)} b^a \lambda^{a-1} \exp(-b\lambda). \quad (7)$$

**Question 8:** Implement the VI algorithm for the variational distribution in Equation (10.24) in Bishop.

In order to implement the VI algorithm for the variational distribution we need to understand how the algorithm works. In other words, understand the theory in a proper way. For this reason, I will explain a little bit the formulas used. However, the whole python implementation is attached in the appendix section.

The final goal is to infer the posterior  $p(\mu, \tau|D)$  where  $\tau = \frac{1}{\sigma^2}$  represents the precision. At the same time, we need to take into consideration the conjugate prior distributions described previously in the description of this section.

In order to infer the posterior, we are going to use a factorized variational approximation. This approximation is followed by the Equation (10.24) stated in Bishop.

$$q(\mu, \tau) = q_\mu(\mu) q_\tau(\tau) \quad (8)$$

Then, in order to obtain the approximation we will need to compute two differentiated steps: updating  $q_\mu(\mu)$  and updating  $q_\tau(\tau)$ .

In order to find the optimal form for  $q_\mu(\mu)$  we need to compute the following equations:

$$\log q_\mu(\mu) = -\frac{\mathbb{E}_{q_\tau}[\tau]}{2} (\lambda_0(\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu_0)^2) + \text{const}$$

After completing the square we can also understand that  $q_\mu(\mu)$  is a Gaussian. All the steps followed can be stated through the code.

On the other hand, to achieve the optimal form of  $q_\tau(\tau)$  we will need to compute the next equation. Here, we are going to recognize the log of a Gamma distribution.

$$\log q_\tau(\tau) = (a_0 - 1) \log \tau - b_0 \tau + \frac{N}{2} \log \tau - \frac{\tau}{2} \mathbb{E}_{q_\mu}(\lambda_0(\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu_0)^2) + \text{const}$$

After finishing all the previous steps and computing the parameters  $a_N, b_N, \mu_N, \lambda_N$  using the equations also found in Bishop, we are going to have the optimal distributions for  $q_\tau(\tau)$  and  $q_\mu(\mu)$ . Then, after doing an initial guess for  $\mathbb{E}_{q_\tau}[\tau]$  and computing again  $q_\tau(\tau)$  we are going to extract  $\mathbb{E}_{q_\mu}[\mu]$  and  $\mathbb{E}_{q_\mu}[\mu^2]$  and compute the distribution  $q_\lambda(\lambda)$ .

**Question 9:** Describe the exact posterior

To describe the exact posterior we are going to start with the likelihood stated in the description of this task.

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{\frac{N}{2}} \exp\left(-\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2\right). \quad (9)$$

We can re-write this equation through adding the sample mean ( $\bar{x}$ ) and the sample covariance ( $s$ ). Applying that, we are going to express the likelihood as:

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{\frac{N}{2}} \exp\left(-\frac{\tau}{2} (N(\mu - \bar{x})^2 + Ns)\right) \quad (10)$$

Where  $\bar{x}$  and  $s$  are defined as:

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n \quad (11)$$

$$s = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2 \quad (12)$$

Therefore, as the description also indicates, the conjugate prior distributions follow Gaussian-Gamma distribution. That is why the joint probability will be:

$$p(\mu, \tau) = p(\mu|\tau)p(\tau) = \frac{b^a \sqrt{\lambda_0}}{\Gamma(a) \sqrt{2\pi}} \tau^{a_0 - \frac{1}{2}} \exp(-b_0 \tau) \exp\left(-\frac{\lambda_0 \tau (\mu - \mu_0)^2}{2}\right) \quad (13)$$

The fact of having a Gaussian likelihood and also, as defined before, a prior of Gaussian-Gamma implies that the posterior distribution will be also a Normal-Gamma distribution. We can state that combining the likelihood described and the priors mentioned:

$$p(\mu, \tau|D) \propto \tau^{N/2 + a_0 - 1/2} \exp\left(-\tau(b_0 - \frac{1}{2}Ns)\right) \exp\left\{-\frac{\tau(\lambda_0(\mu - \mu_0)^2 + N(\bar{x} - \mu)^2)}{2}\right\}$$

Then, if we complete the squares in the last exponential, we are going to obtain:

$$\begin{aligned} \lambda_0(\mu - \mu_0)^2 + N(\bar{x} - \mu)^2 &= \lambda_0\mu^2 + \lambda_0\mu_0^2 - 2\lambda_0\mu\mu_0 + N\bar{x} + N\mu^2 - 2N\bar{x}\mu = \dots = \\ &= (\lambda_0 + N) \left(\mu - \frac{\lambda_0\mu_0 + N\bar{x}}{\lambda_0 + N}\right)^2 + \frac{\lambda_0 N(\bar{x} - \mu_0)^2}{\lambda_0 + N} \end{aligned}$$

The next step is using the previous result and substitute it in the posterior formula.

$$\begin{aligned} p(\mu, \tau|D) &\propto \tau^{N/2 + a_0 - 1/2} \exp\left\{-\tau\left(\frac{1}{2}Ns + b_0 + \frac{\lambda_0 N(\bar{x} - \mu_0)^2}{\lambda_0 + N}\right)\right\} \\ &\quad \exp\left\{-\frac{\tau}{2}(\lambda_0 + N) \left(\mu - \frac{\lambda_0\mu_0 + N\bar{x}}{\lambda_0 + N}\right)^2\right\} \end{aligned}$$

Moreover, from the previous equation we can also discover the parameters of the posterior distribution. As said before, the posterior will be also a Normal-Gamma distribution.

$$\mu_{posterior} = \frac{\lambda_0\mu_0 + N\bar{x}}{\lambda_0 + N} \quad (14)$$

$$\lambda_{posterior} = \lambda_0 + N \quad (15)$$

$$a_{posterior} = a_0 + \frac{N}{2} \quad (16)$$



$$b_{posterior} = b_0 + \frac{1}{2} \left( Ns + \frac{\lambda_0 N(\bar{x} - \mu_0)^2}{\lambda_0 + N} \right) \quad (17)$$

**Question 10:** Compare the variational distribution with the exact posterior. Run the inference for a couple of interesting cases and describe the difference.

The main goal of this question is to compare the variational distribution with the exact posterior. However, it is also important to identify a couple of interesting cases to test and extract conclusions from them. Specifically, I will comment three cases/plots.

The main intention for plotting three different cases is to realize about the importance of the number of data points. We are going to see how the variational distribution is improving as the number of samples is increasing. Concretely, the first case will test using a  $N = 20$ , the second case using a  $N = 100$  and the other case, using a  $N = 900$ .

The variational distribution improves with the number of samples by the simple fact that, as the prior is not good, we need more datapoints to achieve a good variational distribution.

However, in order to understand how the variational inference algorithm work, I also plot four images in each case. Each of the images represent different moments in the execution of the algorithm (a,b,c,d). The first image is representing the initial factorized approximation. The second and the third image are referring to the moment after updating  $q_\mu(\mu)$  and  $q_\tau(\tau)$  respectively. To conclude, the forth image is the one obtained at the end of the algorithm with 20 iterations.

It has to be also take into account that the plots attached below are for a simple example with the following parameters:

$$\tau = 2.2 \quad \mu = 1 \quad \lambda_0 \text{ and } \mu_0 = 1 \quad a_0 \text{ and } b_0 = 0.01$$

With this initialization we are close to have non-informative priors. Obviously, if we try values closer to the real posterior we are going to have better results in the VI. The reason behind that is because it will imply easier convergence.

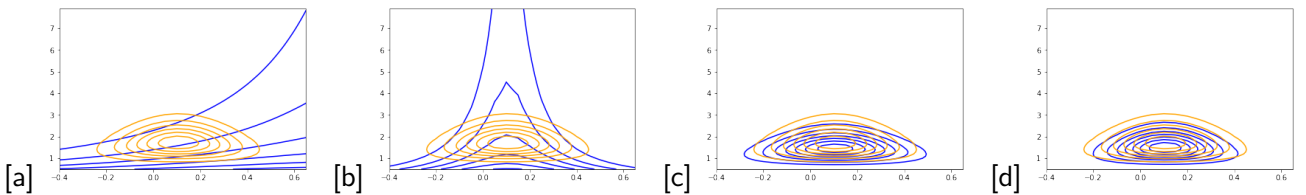


Figure 3: Case 1: Four plots(a,b,c,d) corresponding to the evolution of the VI algorithm using  $N = 20$ . (X axis =  $\mu$ , Y axis =  $\tau$ ). The estimated posterior is represented by the color blue and computed through the VI algorithm, while the exact posterior is plotted in orange. The number of iterations used in the algorithm is 20. (a) Initial factorized approximation, (b) Iteration 1 after updating  $q_\mu(\mu)$ , (c) Iteration 1 after updating  $q_\tau(\tau)$  and (d) Last iteration.

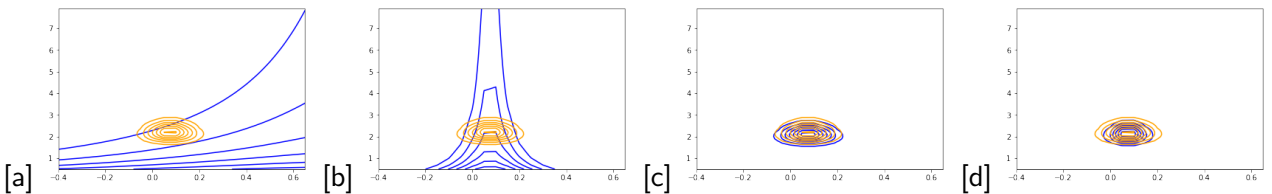


Figure 4: Case 2: Four plots corresponding to the evolution of the VI algorithm using  $N = 100$ . (X axis =  $\mu$ , Y axis =  $\tau$ ). The estimated posterior is represented by the color blue and computed through the VI algorithm, while the exact posterior is plotted in orange. The number of iterations used in the algorithm is 20. (a) Initial factorized approximation, (b) Iteration 1 after updating  $q_\mu(\mu)$ , (c) Iteration 1 after updating  $q_\tau(\tau)$  and (d) Last iteration.

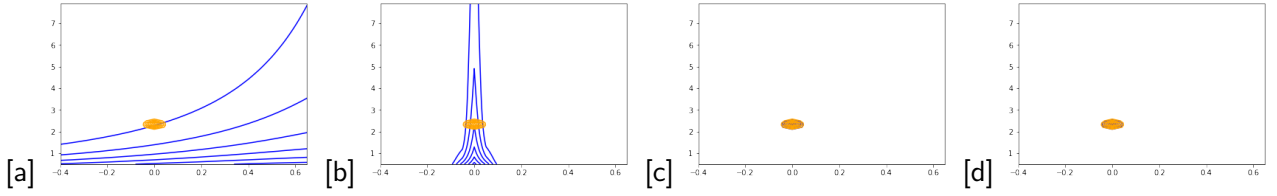


Figure 5: Case 3: Four plots corresponding to the evolution of the VI algorithm using  $N = 900$ . (X axis =  $\mu$ , Y axis =  $\tau$ ). The estimated posterior is represented by the color blue and computed through the VI algorithm, while the exact posterior is plotted in orange. The number of iterations used in the algorithm is 20. (a) Initial factorized approximation, (b) Iteration 1 after updating  $q_\mu(\mu)$ , (c) Iteration 1 after updating  $q_\tau(\tau)$  and (d) Last iteration.

## 4 Variational Inference

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

To sum up all the information described in the description of this task we can state that the Cartesian Matrix Model (CMM) is defined as:

- $R$  rows distribution  $\mathcal{N}(\mu_r, \lambda_r^{-1}) : 1 \leq r \leq R$ , each variance  $\lambda_r^{-1}$  is known and each  $\mu_r$  has a prior distribution  $\mathcal{N}(\mu, \lambda^{-1})$
- $C$  column distribution  $\mathcal{N}(\xi_c, \tau_c^{-1}) : 1 \leq c \leq C$  each variance  $\tau_c^{-1}$  is known and each  $\xi_c$  has a prior distribution  $\mathcal{N}(\xi, \tau^{-1})$
- Hyper-parameters known.
- Matrix  $S$  is generated by  $S_{rc} = X_r + Y_c$
- Use Variational Inference to obtain variational distribution :

$$q(\mu_1, \dots, \mu_R, \xi_1, \dots, \xi_C) = \prod_r q(\mu_r) \prod_c q(\xi_c) \quad (18)$$

To continue, we can define the distribution  $S_{rc}$  through the sum of the two Gaussian variables. The reason behind is because we know that from a sum of gaussians we are going to achieve another gaussian.

$$S_{rc} \sim \mathcal{N}(\mu_r + \xi_c, \lambda_r^{-1} + \tau_c^{-1}). \quad (19)$$

Then, we can define the joint distribution as:

$$p(S, \mu_1, \dots, \mu_R, \xi_1, \dots, \xi_C) = p(S | \mu_1, \dots, \mu_R, \xi_1, \dots, \xi_C) p(\mu_1, \dots, \mu_R) p(\xi_1, \dots, \xi_C) = \dots \quad (20)$$

$$\prod_{r=1}^R \prod_{c=1}^C \mathcal{N}(S_{rc} | \mu_r + \xi_c, \lambda_r^{-1} + \tau_c^{-1}) \prod_{r=1}^R \mathcal{N}(\mu_r | \mu, \lambda_r^{-1}) \prod_{c=1}^C \mathcal{N}(\xi_c | \xi, \tau^{-1}). \quad (21)$$

Also, we need to state that in order to obtain the previous joint distribution the assumption of rows and column distributions being independent has been made.

Then, in order to obtain the optimized factors, we are going to compute the log form of the join probability.

$$\ln p(S, \mu_1, \dots, \mu_R, \xi_1, \dots, \xi_C) \quad (22)$$

Moreover, we are going to define the overall variance as  $\beta_{rc}^{-1} = \lambda_r^{-1} + \tau_c^{-1}$

Based on what we said, the log form of the join probability is equal to:

$$\begin{aligned} & \sum_{c=1}^C \sum_{r=1}^R \ln \mathcal{N}(S_{rc} | \mu_r + \xi_c, \lambda_r^{-1} + \tau_c^{-1}) + \sum_{r=1}^R \ln \mathcal{N}(\mu_r | \mu, \lambda^{-1}) + \sum_{c=1}^C \ln \mathcal{N}(\xi_c | \xi, \tau^{-1}) = \\ & \sum_{c=1}^C \sum_{r=1}^R \left( \frac{1}{2} \ln \beta_{rc} - \frac{1}{2} \ln(2\pi) - \frac{\beta_{rc}}{2} (S_{rc} - \mu_r - \xi_c)^2 \right) + \sum_{r=1}^R \left( \frac{1}{2} \ln \lambda - \frac{1}{2} \ln(2\pi) - \frac{\lambda}{2} (\mu_r - \mu)^2 \right) \\ & + \sum_{c=1}^C \left( \frac{1}{2} \ln \tau - \frac{1}{2} \ln(2\pi) - \frac{\tau}{2} (\xi_c - \xi)^2 \right) \end{aligned}$$

And then, we are able to define  $\ln q^*(\mu_r)$  and  $\ln q^*(\xi_c)$ . First of all we are going to calculate  $\ln q^*(\mu)$  and then the  $\ln q^*(\xi_c)$ .

In order to calculate the first optimized factor we are going to compute the expectation of the log joint probability over all parameters (except of  $\mu_r$ ).

$$\begin{aligned} \ln q_r^*(\mu_r) &= E_{i \neq r} [\ln p(\mu_1, \dots, \mu_R, \xi_1, \dots, \xi_C, S)] + \text{const} = \\ &= E_{i \neq r} \left[ \ln \left( \prod_{r=1}^R \prod_{c=1}^C \mathcal{N}(S_{rc} | \mu_r + \xi_c, \lambda_r^{-1} + \tau_c^{-1}) \prod_{r=1}^R \mathcal{N}(\mu_r | \mu, \lambda^{-1}) \prod_{c=1}^C \mathcal{N}(\xi_c | \xi, \tau^{-1}) \right) \right] + \text{const} = \\ &= E_{i \neq r} \left[ \sum_{r=1}^R \sum_{c=1}^C \ln \mathcal{N}(S_{rc} | \mu_r + \xi_c, \lambda_r^{-1} + \tau_c^{-1}) + \sum_{r=1}^R \ln \mathcal{N}(\mu_r | \mu, \lambda^{-1}) + \sum_{c=1}^C \ln \mathcal{N}(\xi_c | \xi, \tau^{-1}) \right] + \text{const} = \end{aligned} \quad (23)$$

And then, calculating the normal distributions:

$$\begin{aligned} &= E_{i \neq r} \left[ \sum_{r=1}^R \sum_{c=1}^C -\frac{(S_{rc} - \mu_r - \xi_c)^2}{2(\lambda_r^{-1} + \tau_c^{-1})} + \sum_{r=1}^R -\frac{(\mu_r - \mu)^2}{2\lambda^{-1}} + \sum_{c=1}^C -\frac{(\xi_c - \xi)^2}{2\tau^{-1}} \right] + \text{const} \\ &= -\frac{1}{2} E_{i \neq r} \left[ \sum_{r=1}^R \sum_{c=1}^C \frac{(S_{rc} - \mu_r - \xi_c)^2}{\lambda_r^{-1} + \tau_c^{-1}} + \sum_{r=1}^R \lambda (\mu_r - \mu)^2 + \sum_{c=1}^C \tau (\xi_c - \xi)^2 \right] + \text{const} \\ &= -\frac{1}{2} E_{i \neq r} \left[ \sum_{c=1}^C \left( \frac{(S_{rc} - \mu_r - \xi_c)^2}{\lambda_r^{-1} + \tau_c^{-1}} + \lambda (\mu_r - \mu)^2 \right) \right] + \text{const} \end{aligned} \quad (24)$$

To continue, we are only keeping the  $\mu_r$ -dependent terms and taking into account  $a_{rc} = \frac{1}{\lambda_r^{-1} + \tau_c^{-1}}$  we can continue with the calculus of the squares:

$$\begin{aligned} &= -\frac{1}{2} E_{i \neq r} \left[ \sum_{c=1}^C a_{rc} (S_{rc}^2 - 2S_{rc}(\mu_r + \xi_c) + \mu_r^2 + 2\mu_r \xi_c + \xi_c^2) + \lambda (\mu_r^2 - 2\mu \mu_r + \mu^2) \right] + \text{const} \\ &= -\frac{1}{2} E_{i \neq r} \left[ \sum_{c=1}^C (a_{rc} + \lambda) \mu_r^2 + 2(a_{rc}(\xi_c - S_{rc}) - \mu \lambda) \mu_r \right] + \text{const} \\ &= -\frac{1}{2} \left[ \left( \sum_{c=1}^C (a_{rc} + \lambda) \right) \mu_r^2 + 2 \left( \sum_{c=1}^C (a_{rc}(E[\xi_c] - S_{rc}) - \mu \lambda) \right) \mu_r \right] + \text{const} \\ &= -\frac{1}{2} \left( \sum_{c=1}^C a_{rc} + C\lambda \right) \mu_r^2 - \left( \sum_{c=1}^C a_{rc}(E[\xi_c] - S_{rc}) - C\mu \lambda \right) \mu_r + \text{const} \end{aligned} \quad (25)$$

With all the expressions stated above, we can conclude:

$$\mu_r \sim N(a_r, b_r) \quad (26)$$

$$b_r = \left( \sum_{c=1}^C a_{rc} + C\lambda \right)^{-1} \quad (27)$$

$$a_r = b_r \left( \sum_{c=1}^C a_{rc}(S_{rc} - E[\xi_c]) + C\mu\lambda \right) \quad (28)$$

$$\xi_c \sim N(a_c, b_c) \quad (29)$$

$$b_c = \left( \sum_{r=1}^R a_{rc} + R\tau \right)^{-1} \quad (30)$$

$$a_c = b_c \left( \sum_{r=1}^R a_{rc}(S_{rc} - E[\mu_r]) + R\xi\tau \right) \quad (31)$$

## Appendices

### A Code question 3

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def Draw(total_sum):
5
6     plt.title("Histogram mid-biased dices")
7     plt.hist(total_sum)
8     plt.xlabel("Value")
9     plt.ylabel("Frequency")
10    plt.show()
11
12    #Number of players and number of tables for all the experiments.
13    num_players = 500
14    num_tables = 100
15
16
17    #Define all the possible distributions.
18    uniform_distribution = [1./6,1./6,1./6,1./6,1./6,1./6]
19    mid_bias_distribution = [0.5,0.5,0,0,0,0]
20    biased_distribution = [0,1,0,0,0,0] #In this case the probability of the number 2 will ...
        be 6/6.
21
22    #HMM
23
24    # Create transition matrix ( probabilities of going to different groups of tables (1/4 ...
        and 3/4))
25    A = [[0.25, 0.75], [0.75, 0.25]]
26    #Define the emission matrix.
27    B = [uniform_distribution, mid_bias_distribution]
28    #Initial conditions. We can be either in one group of tables or in the other.

```

```

29 pi = [0.5, 0.5]
30
31
32 # Alpha-pass algorithm
33 total_sum = np.zeros(num_players)
34 for player in range(num_players):
35     pi = [0.5, 0.5]
36     output_distribution = np.dot(pi, B)
37     total_sum[player] = total_sum[player] + np.random.choice(a=[1, 2, 3, 4, 5, 6], ...
38         p=output_distribution)
39     for timestep in range(1,num_tables):
40         pi = np.dot(pi,A)
41         output_distribution = np.dot(pi,B)
42         total_sum[player] = total_sum[player] + np.random.choice(a=[1, 2, 3, 4, 5, 6], ...
43             p = output_distribution)
44
45 Draw(total_sum)

```

## B Code question 10

```

1 #Import libraries
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import scipy.stats as stats
5 import sys
6
7 #Create the data
8 def GenerateData(mu,tau,N):
9     x = np.random.normal(loc=mu, scale=tau**(-0.5), size=N)
10    return x
11
12 #Prior distribution of tau will be the Gaussian-Gamma distribution.
13 def normal_gamma(mean, precision, shape, rate):
14     normal = stats.norm(loc=mean, scale=(precision*range_tau)**(-.5))
15     gamma = stats.gamma(a=shape, scale=1.0/rate)
16     return [[normal.pdf(u) [t]*gamma.pdf(range_tau[t])
17         for u in range_mu] for t in range(len(range_tau))]
18
19 #VI algorithm. Following Bishop equations.
20 def VIalgorithm(a_init,b_init,mu_init,lambda_init,N,data,its):
21     sample_mean = np.mean(x)
22     sample_var = np.var(x)
23     a_N = a_init
24     b_N = b_init
25     vi_posterior= []
26     for i in range(20):
27         mean_tau = a_N/b_N
28         mu_N = (lambda_init*mu_init+N*sample_mean)/(lambda_init+N)
29         lambda_N = (N+lambda_init)*mean_tau
30         vi_posterior1 = normal_gamma(mu_N,lambda_N,a_N,b_N)
31         mean_mu = mu_N
32         mean_mu2 = mu_N**2 + 1/lambda_N
33         a_N = a_init+(N+1)/2
34         b_N = ...
35         b_init+.5*(mean_mu2*(N+lambda_init)-2*mean_mu*(N*sample_mean+lambda_init*mu_init))+lambda
36         vi_posterior2 = normal_gamma(mu_N,lambda_N,a_N,b_N)
37         if (i < 1):
38             vi_posterior.append(vi_posterior1)
39             vi_posterior.append(vi_posterior2)
40         vi_posterior.append(vi_posterior2)
41     return vi_posterior

```

```

41
42 #Initialize the prior parameters.
43 #In that case, non-informative priors.
44 def priorInit():
45     mu_init = 1
46     lambda_init = 1
47     a_init = 0.01
48     b_init = 0.01
49     return a_init,mu_init,lambda_init,b_init
50
51 #Calculate the exact posterior.
52 def ExactPosterior(mu_init,lambda_init,a_init,b_init,N,data):
53     sample_mean = np.mean(data)
54     sample_var = np.var(data)
55     mu_posterior = (lambda_init*mu_init+N*sample_mean)/(lambda_init + N)
56     lambda_posterior = lambda_init+N
57     a_posterior = a_init + N/2
58     b_posterior = b_init + 0.5*(N*sample_var + ...
        (lambda_init*N*(sample_mean-mu_init))/(lambda_init+N))
59     real_posterior = normal_gamma(mu_posterior,lambda_posterior,a_posterior,b_posterior)
60     return real_posterior
61
62 def DrawGraphs(prior,real_posterior,vi_posterior):
63     fig = plt.figure()
64     plt.contour(prior,colors = 'blue')
65     plt.contour(real_posterior,colors = 'orange')
66     plt.show()
67     for j in range(0,3):
68         fig = plt.figure()
69         plt.contour(vi_posterior[j],colors='blue')
70         plt.contour(real_posterior,colors='orange')
71         plt.show()
72
73 def Main():
74
75     #Parameters
76     mu = 0
77     tau = 2.2
78     N_data = 100
79     iterations = 20
80     data = GenerateData(mu,tau,N_data)
81     #To compute the graphs
82     range_mu = np.arange(-.4,0.7,0.05)
83     range_tau = np.arange(0.5,8,0.1)
84     sample_mean = np.mean(data)
85     sample_var = np.var(data)
86     a_init,mu_init,lambda_init,b_init = priorInit()
87     prior = normal_gamma(mu_init,lambda_init,a_init,b_init)
88     exact_posterior = ExactPosterior(mu_init,lambda_init,a_init,b_init,N_data,data)
89     infer_posterior = ...
        VIALgorithm(a_init,b_init,mu_init,lambda_init,N_data,data,iterations)
90     DrawGraphs(prior,exact_posterior,infer_posterior)
91
92 if __name__ == "__main__":
93     Main()

```

## References

- [1] C. M.Bishop, "Pattern recognition and machine learning," 2009.