

Machine Learning Advanced Course: Assignment 1A

Michel Le Dez, Corentin Flandre

December 05, 2023

Contents

1	Assignment 1A	2
1.1	Exponential Family	2
1.2	Dependencies in a Directed Graphical Model	4
1.3	CAVI	5
1.4	SVI - LDA	11
1.5	BBVI	12

1 Assignment 1A

1.1 Exponential Family

An exponential-family distribution with natural parameters is in the following form:

$$p(x|\theta) = h(x) \exp(\eta(\theta) \cdot T(x) - A(\eta))$$

Question 1.1.1:

We have:

- $\theta = \lambda$
- $\eta(\theta) = \log(\theta) = \log(\lambda)$
- $h(x) = \frac{1}{x!}$
- $T(x) = x$
- $A(\eta) = e^\eta = e^{\log \lambda} = \lambda$

Therefore:

$$\begin{aligned} p(x|\theta) &= h(x) \exp(\eta(\theta) \cdot T(x) - A(\eta)) \\ \iff p(x|\lambda) &= \frac{1}{x!} \exp(x \log(\lambda) - \lambda) \\ \boxed{p(x|\lambda) &= \frac{\lambda^x}{x!} e^{-\lambda}} \end{aligned}$$

We recognize the probability mass function of the **Poisson distribution**.

Question 1.1.2:

We have:

- $\theta = [\alpha, \beta]$
- $\eta(\theta) = [\theta_1 - 1, -\theta_2] = [\alpha - 1, -\beta]$
- $h(x) = 1$
- $T(x) = [\log x, x]$
- $A(\eta) = \log \Gamma(\eta_1 + 1) - (\eta_1 + 1) \log(-\eta_2) = \log \Gamma(\alpha) - \alpha \log(\beta) = \log\left(\frac{\Gamma(\alpha)}{\beta^\alpha}\right)$

Therefore:

$$\begin{aligned} p(x|\theta) &= h(x) \exp(\eta(\theta) \cdot T(x) - A(\eta)) \\ \iff p(x|\alpha, \beta) &= \exp((\alpha - 1) \log x - \beta x - \log\left(\frac{\Gamma(\alpha)}{\beta^\alpha}\right)) \\ \boxed{p(x|\alpha, \beta) &= \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}} \end{aligned}$$

We recognize the probability density function of the **Gamma distribution**.

Question 1.1.3:

We have:

- $\theta = [\mu, \sigma^2]$

- $\eta(\theta) = [\frac{\theta_1}{\theta_2}, -\frac{1}{2\theta_2}] = [\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}]$
- $h(x) = \frac{1}{\sqrt{2\pi}}$
- $T(x) = [x, x^2]$
- $A(\eta) = -\frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \log(-2\eta_2) = \frac{\mu^2}{2\sigma^2} - \log(\frac{1}{\sigma})$

Therefore:

$$\begin{aligned}
 p(x|\theta) &= h(x) \exp(\eta(\theta) \cdot T(x) - A(\eta)) \\
 \iff p(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi}} \exp\left(\frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2 - \frac{\mu^2}{2\sigma^2} + \log\left(\frac{1}{\sigma}\right)\right) \\
 \boxed{p(x|\mu, \sigma^2) &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)}
 \end{aligned}$$

We recognize the probability density function of the **Normal distribution**.

Question 1.1.4:

We have:

- $\theta = \lambda$
- $\eta(\theta) = -\theta = -\lambda$
- $h(x) = 2$
- $T(x) = x$
- $A(\eta) = -\log(-\frac{\eta}{2}) = -\log(\frac{\lambda}{2})$

Therefore:

$$\begin{aligned}
 p(x|\theta) &= h(x) \exp(\eta(\theta) \cdot T(x) - A(\eta)) \\
 \iff p(x|\lambda) &= 2 \exp(-\lambda x + \log(\frac{\lambda}{2})) \\
 \boxed{p(x|\lambda) &= \lambda e^{-\lambda x}}
 \end{aligned}$$

We recognize the probability density function of the **Exponential distribution**.

Question 1.1.5:

We have:

- $\theta = [\psi_1, \psi_2]$
- $\eta(\theta) = [\theta_1 - 1, \theta_2 - 2] = [\psi_1 - 1, \psi_2 - 2]$
- $h(x) = 1$
- $T(x) = [\log x, \log(1 - x)]$
- $A(\eta) = \log \Gamma(\eta_1 + 1) + \log \Gamma(\eta_2 + 1) - \log \Gamma(\eta_1 + \eta_2 + 2) = -\log \frac{\Gamma(\psi_1 + \psi_2)}{\Gamma(\psi_1)\Gamma(\psi_2)}$

Therefore:

$$\begin{aligned}
 p(x|\theta) &= h(x) \exp(\eta(\theta) \cdot T(x) - A(\eta)) \\
 \iff p(x|\psi_1, \psi_2) &= \exp((\psi_1 - 1) \log x + (\psi_2 - 1) \log(1 - x) + \log \frac{\Gamma(\psi_1 + \psi_2)}{\Gamma(\psi_1)\Gamma(\psi_2)}) \\
 \boxed{p(x|\psi_1, \psi_2) &= \frac{\Gamma(\psi_1 + \psi_2)}{\Gamma(\psi_1)\Gamma(\psi_2)} x^{\psi_1-1} (1 - x)^{\psi_2-1}}
 \end{aligned}$$

We recognize the probability density function of the **Beta distribution**.

1.2 Dependencies in a Directed Graphical Model

Question 1.2.6: Yes.

Question 1.2.7: No.

Question 1.2.8: Yes.

Question 1.2.9: No.

Question 1.2.10: No.

Question 1.2.11: No.

1.3 CAVI

We have the following distribution:

$$p(\tau) = \text{Gam}(\tau|a_0, b_0) = \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} e^{-b_0\tau} \quad (1)$$

$$p(\mu|\tau) = \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1}) = \frac{\sqrt{\lambda_0\tau}}{\sqrt{2\pi}} \exp(-\frac{\lambda_0\tau}{2}(\mu - \mu_0)^2) \quad (2)$$

$$p(D|\mu, \tau) = \prod_{n=1}^N \frac{\sqrt{\tau}}{\sqrt{2\pi}} \exp(-\frac{\tau}{2}(x_n - \mu)^2) = (\frac{\tau}{2\pi})^{\frac{N}{2}} \exp(-\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2) \quad (3)$$

Question 1.3.12: The function implementation for generating data points, as well as the code for displaying the histogram is in the annex. Here are the the histograms we obtained:

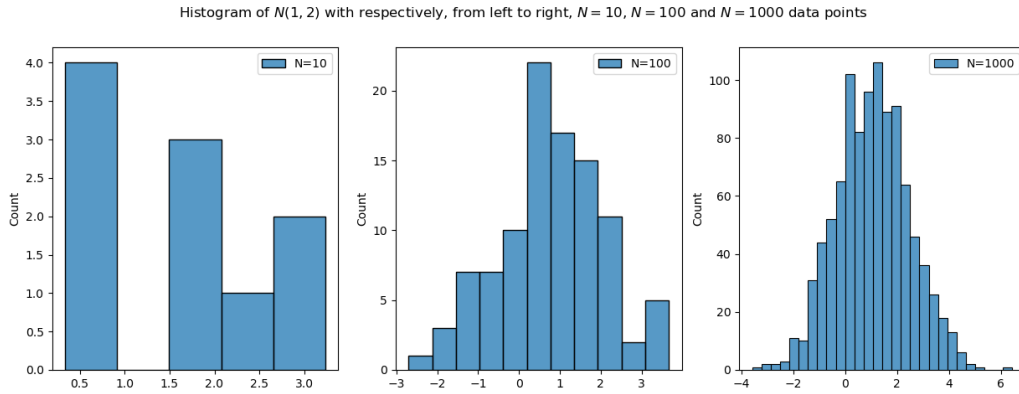


Figure 1: Histogram of $\mathcal{N}(\mu, \frac{1}{\tau})$ with $\mu = 1$ and $\tau = 0.5$ with respectively, from left to right, $N = 10$, $N = 100$, $N = 1000$ data points.

We observe that the more data we have, the closer the histogram is to the normal distribution that generated it.

Question 1.3.13: The likelihood of the data points $D = x_{1:N}$ given the parameter μ, τ is as follows:

$$l(\mu, \tau) := p(D|\mu, \tau) = (\frac{\tau}{2\pi})^{\frac{N}{2}} \exp(-\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2)$$

$$\iff \log(l(\mu, \tau)) = \frac{N}{2} \log \tau - \frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 + \text{const}$$

We are looking for the parameters μ and τ which maximise the likelihood $l(\mu, \tau)$, which is equivalent to maximising the log-likelihood given that the log function is a monotonically increasing one. The constant term above gathers all the terms that do not depend on μ or τ . Deriving the gradient of $l(\mu, \tau)$ and setting it to 0 at (μ_{MLE}, τ_{MLE}) yields the following system of equations:

$$\begin{cases} \tau_{MLE} \sum_{n=1}^N x_n - \tau_{MLE} N \mu_{MLE} = 0 \\ \frac{N}{2\tau_{MLE}} - \frac{1}{2} \sum_{n=1}^N (x_n - \mu_{MLE})^2 = 0 \end{cases}$$

which yields to the following solution:

$$\begin{cases} \mu_{MLE} = \bar{x} = \frac{1}{N} \sum_{n=1}^N x_n \\ \tau_{MLE} = \frac{1}{\frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2} \end{cases}$$

To verify that the point (μ_{MLE}, τ_{MLE}) definitely maximises the likelihood we can compute the hessian of the log-likelihood at this point, which yields to:

$$\begin{bmatrix} -\frac{N^2}{\sum_{n=1}^N (x_n - \bar{x})^2} & 0 \\ 0 & -\frac{1}{2N} (\sum_{n=1}^N (x_n - \bar{x})^2)^2 \end{bmatrix}$$

We can see the eigenvalues of the hessian are strictly negative, therefore (μ_{MLE}, τ_{MLE}) definitely maximises the likelihood.

Question 1.3.14: To compute the posterior, we are going to use the Bayes' theorem and gather in the constant term all the terms that do not depend on μ or τ . Here is the posterior:

$$\begin{aligned} p(\mu, \tau|D) &= p(D|\mu, \tau)p(\mu|\tau)p(\tau)/p(D) \\ \iff \log p(\mu, \tau|D) &= \log p(D|\mu, \tau) + \log p(\mu|\tau) + \log p(\tau) + \text{const} \\ &= \frac{N}{2} \log \tau - \frac{\tau}{2} \sum_{n=1}^N (x_n^2 + \mu^2 - 2x_n\mu) + \frac{1}{2} \log \tau - \frac{\lambda_0 \tau}{2} (\mu^2 + \mu_0^2 - 2\mu\mu_0) \\ &\quad + (a_0 - 1) \log \tau - b_0 \tau + \text{const} \\ &= (a_0 + \frac{N}{2} - \frac{1}{2}) \log \tau - (b_0 + \frac{1}{2} \sum_{n=1}^N x_n^2 + \frac{\lambda_0 \mu_0^2}{2}) \tau + (\sum_{n=1}^N x_n + \lambda_0 \mu_0) \tau \mu \\ &\quad - \frac{\tau}{2} (\lambda_0 + N) \mu^2 + \text{const} \end{aligned}$$

However, we know that for $\mu, \tau \sim \text{NormalGamma}(\mu_0^*, \lambda_0^*, a_0^*, b_0^*)$, the logarithm of the probability density function is as follows:

$$\begin{aligned} \log p(\mu, \tau|\mu_0^*, \lambda_0^*, a_0^*, b_0^*) &= (a_0^* - \frac{1}{2}) \log \tau - b_0^* \tau - \frac{\lambda_0^* \mu_0^{*2}}{2} + \lambda_0^* \mu_0^* \tau \mu - \frac{\tau}{2} \lambda_0^* \mu^2 + \text{const} \\ &= (a_0^* - \frac{1}{2}) \log \tau - b_0^* \tau - \frac{\tau \lambda_0^*}{2} (\mu - \mu_0^*)^2 + \text{const} \end{aligned}$$

By identification, we have $a_0^* = a_0 + \frac{N}{2}$, $\lambda_0^* = \lambda_0 + N$ and $\mu_0^* = \frac{\sum_{n=1}^N x_n + \lambda_0 \mu_0}{\lambda_0 + N}$. For b_0^* , let's rewrite the log posterior in the form of the second equality above by adding and subtracting the missing term $\frac{1}{2} \frac{(\sum_{n=1}^N x_n + \lambda_0 \mu_0)^2}{\lambda_0 + N} \tau$ for completing the square, which yields to:

$$\begin{aligned} \log p(\mu, \tau|D) &= (a_0 + \frac{N}{2} - \frac{1}{2}) \log \tau - (b_0 + \frac{1}{2} \sum_{n=1}^N x_n^2 + \frac{\lambda_0 \mu_0^2}{2} - \frac{1}{2} \frac{(\sum_{n=1}^N x_n + \lambda_0 \mu_0)^2}{\lambda_0 + N}) \tau \\ &\quad - \frac{\tau (\lambda_0 + N)}{2} (\mu - \frac{\sum_{n=1}^N x_n + \lambda_0 \mu_0}{\lambda_0 + N})^2 + \text{const} \end{aligned}$$

Here, we can easily identify b_0^* and we summarise the results below:

$\mu, \tau|D \sim \text{NormalGamma}(\mu_0^*, \lambda_0^*, a_0^*, b_0^*)$ with the following parameters

$$\begin{aligned} \mu_0^* &= \frac{\sum_{n=1}^N x_n + \lambda_0 \mu_0}{\lambda_0 + N} \\ \lambda_0^* &= \lambda_0 + N \\ a_0^* &= a_0 + \frac{N}{2} \\ b_0^* &= b_0 + \frac{1}{2} \sum_{n=1}^N x_n^2 + \frac{\lambda_0 \mu_0^2}{2} - \frac{1}{2} \frac{(\sum_{n=1}^N x_n + \lambda_0 \mu_0)^2}{\lambda_0 + N} \end{aligned}$$

Question 1.3.15: The mean field approximation for the variational distribution is the following:

$$q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau).$$

The log of the joint distribution can be written as follows:

$$\log p(x, \mu, \tau) = \log p(x|\mu, \tau) + \log p(\mu|\tau) + \log p(\tau),$$

with:

$$\begin{aligned}\log p(x|\mu, \tau) &= \frac{N}{2} \log \tau - \frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 + \text{const} \\ \log p(\mu|\tau) &= \frac{1}{2} \log \tau - \frac{\lambda_0 \tau}{2} (\mu - \mu_0)^2 + \text{const} \\ \log p(\tau) &= (a_0 - 1) \log \tau - b_0 \tau + \text{const}\end{aligned}$$

where the constant terms include terms that do not depend on μ or τ . Let's derive now the coordinate ascent update for μ by including terms that do not depend on μ in the constant term (i.e $\log p(\tau)$, $\frac{N}{2} \log \tau$, $\frac{1}{2} \log \tau$, $-\frac{1}{2} \mathbf{E}_{q(\tau)}[\tau] \sum_{n=1}^N x_n^2$ and $-\frac{1}{2} \mathbf{E}_{q(\tau)}[\tau] \lambda_0 \mu_0^2$):

$$\begin{aligned}\log q^*(\mu) &= \mathbf{E}_{q(\tau)}[\log p(x, \mu, \tau)] \\ &= -\mathbf{E}_{q(\tau)} \left[\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 + \frac{\lambda_0 \tau}{2} (\mu - \mu_0)^2 \right] + \text{const} \\ &= -\frac{1}{2} \mathbf{E}_{q(\tau)}[\tau] \left(\sum_{n=1}^N (x_n - \mu)^2 + \lambda_0 (\mu - \mu_0)^2 \right) + \text{const} \\ &= \mathbf{E}_{q(\tau)}[\tau] \left(\sum_{n=1}^N x_n + \lambda_0 \mu_0 \right) \mu - \frac{1}{2} \mathbf{E}_{q(\tau)}[\tau] (\lambda_0 + N) \mu^2 + \text{const}\end{aligned}$$

However, we know that for $\mu \sim \text{Normal}(\tilde{\mu}_0, \tilde{\lambda}_0^{-1})$, the log of the probability density function is as follows:

$$\log p(\mu|\tilde{\mu}_0, \tilde{\lambda}_0^{-1}) = \tilde{\lambda}_0 \tilde{\mu}_0 \mu - \frac{\tilde{\lambda}_0}{2} \mu^2 + \text{const}$$

By identification, we have:

$$\begin{aligned}q^*(\mu) &= \text{Normal}(\mu|\tilde{\mu}_0, \tilde{\lambda}_0^{-1}) \text{ with the following parameters} \\ \tilde{\mu}_0 &= \frac{\sum_{n=1}^N x_n + \lambda_0 \mu_0}{\lambda_0 + N} \\ \tilde{\lambda}_0 &= \mathbf{E}_{q(\tau)}[\tau] (\lambda_0 + N) \\ \text{with} \\ \mathbf{E}_{q(\tau)}[\tau] &= \frac{\tilde{a}_0}{\tilde{b}_0}\end{aligned}$$

Let's derive now the coordinate ascent update for τ by including terms that do not depend on τ in the constant term:

$$\begin{aligned}\log q^*(\tau) &= \mathbf{E}_{q(\mu)}[\log p(x, \mu, \tau)] \\ &= \frac{N}{2} \log \tau - \frac{\tau}{2} \sum_{n=1}^N \mathbf{E}_{q(\mu)}[(x_n - \mu)^2] + \frac{1}{2} \log \tau - \frac{\lambda_0 \tau}{2} \mathbf{E}_{q(\mu)}[(\mu - \mu_0)^2] + (a_0 - 1) \log \tau - b_0 \tau \\ &= (a_0 + \frac{N+1}{2} - 1) \log \tau - (b_0 + \frac{1}{2} \mathbf{E}_{q(\mu)}[\sum_{n=1}^N (x_n - \mu)^2 + \lambda_0 (\mu - \mu_0)^2]) \tau\end{aligned}$$

However, we know that for $\tau \sim \text{Gamma}(\tilde{a}_0, \tilde{b}_0)$, the log of the probability density function is as follows:

$$\log p(\tau|\tilde{a}_0, \tilde{b}_0) = (\tilde{a}_0 - 1) \log \tau - \tilde{b}_0 \tau + \text{const}$$

By identification, we have:

$$q^*(\tau) = \text{Gamma}(\tau|\tilde{a}_0, \tilde{b}_0) \text{ with the following parameters}$$

$$\tilde{a}_0 = a_0 + \frac{N+1}{2}$$

$$\tilde{b}_0 = b_0 + \frac{1}{2} \left[\sum_{n=1}^N x_n^2 + \lambda_0 \mu_0^2 - 2(\lambda_0 \mu_0 + \sum_{n=1}^N x_n) \mathbf{E}_{q(\mu)}[\mu] + (\lambda_0 + \sum_{n=1}^N x_n^2) \mathbf{E}_{q(\mu)}[\mu^2] \right]$$

with

$$\mathbf{E}_{q(\mu)}[\mu] = \tilde{\mu}_0$$

$$\mathbf{E}_{q(\mu)}[\mu^2] = \frac{1}{\tilde{\lambda}_0} + \tilde{\mu}_0^2$$

Now that we have derived the variational distributions, we can compute the ELBO \mathcal{L} to monitor its convergence during the CAVI algorithm:

$$\mathcal{L} = \mathbf{E}_{q(\mu, \tau)} \left[\log \frac{p(x, \mu, \tau)}{q(\mu, \tau)} \right] = \mathbf{E}_{q(\mu)q(\tau)} \left[\log \frac{p(x|\mu, \tau)p(\mu|\tau)p(\tau)}{q(\mu)q(\tau)} \right]$$

which finally gives:

$$\mathcal{L} = \mathbf{E}_{q(\mu)q(\tau)} [\log p(x|\mu, \tau)] + \mathbf{E}_{q(\mu)q(\tau)} [\log p(\mu|\tau)] + \mathbf{E}_{q(\tau)} [\log p(\tau)]$$

$$- \mathbf{E}_{q(\mu)} [\log q(\mu)] - \mathbf{E}_{q(\tau)} [\log q(\tau)]$$

with:

$$\mathbf{E}_{q(\mu)q(\tau)} [\log p(x|\mu, \tau)] = -\frac{N}{2} \log(2\pi) + \frac{N}{2} \mathbf{E}_{q(\tau)} [\log q(\tau)] - \frac{1}{2} \mathbf{E}_{q(\tau)} [\tau] \sum_{n=1}^N (x_n^2 + \mathbf{E}_{q(\mu)}[\mu^2] - 2x_n \mathbf{E}_{q(\mu)}[\mu])$$

$$\mathbf{E}_{q(\mu)q(\tau)} [\log p(\mu|\tau)] = \frac{1}{2} \log\left(\frac{\lambda_0}{2\pi}\right) + \frac{1}{2} \mathbf{E}_{q(\tau)} [\log q(\tau)] - \frac{\lambda_0}{2} \mathbf{E}_{q(\tau)} [\tau] (\mathbf{E}_{q(\mu)}[\mu^2] + \mu_0^2 - 2\mu_0 \mathbf{E}_{q(\mu)}[\mu])$$

$$\mathbf{E}_{q(\tau)} [\log p(\tau)] = (a_0 - 1) \mathbf{E}_{q(\tau)} [\log q(\tau)] + a_0 \log b_0 - b_0 \mathbf{E}_{q(\tau)} [\tau] - \log \Gamma(a_0)$$

$$\mathbf{E}_{q(\mu)} [\log q(\mu)] = -H(\mu) = -\frac{1}{2} \log\left(\frac{2\pi}{\tilde{\lambda}_0}\right) - \frac{1}{2}$$

$$\mathbf{E}_{q(\tau)} [\log q(\tau)] = -(\tilde{a}_0 - \log \tilde{b}_0 + \log \Gamma(\tilde{a}_0) + (1 - \tilde{a}_0) \psi(\tilde{a}_0))$$

and:

$$\mathbf{E}_{q(\tau)} [\log q(\tau)] = \psi(\tilde{a}_0) - \log \tilde{b}_0$$

$$\mathbf{E}_{q(\tau)} [\tau] = \frac{\tilde{a}_0}{\tilde{b}_0}$$

$$\mathbf{E}_{q(\mu)} [\mu] = \tilde{\mu}_0$$

$$\mathbf{E}_{q(\mu)} [\mu^2] = \frac{1}{\tilde{\lambda}_0} + \tilde{\mu}_0^2$$

We implemented the CAVI algorithm using the variational distributions found above, and we monitored the convergence of the ELBO. Finally, we compared the variational distribution $q(\mu, \tau)$ obtained at the end of the CAVI algorithm with the true posterior that we computed at **Question 1.3.14**. Here are the results we obtained for the 3 data sets generated at **Question 1.13.12**. We also display the ML estimate computed at **Question 1.13.13**.

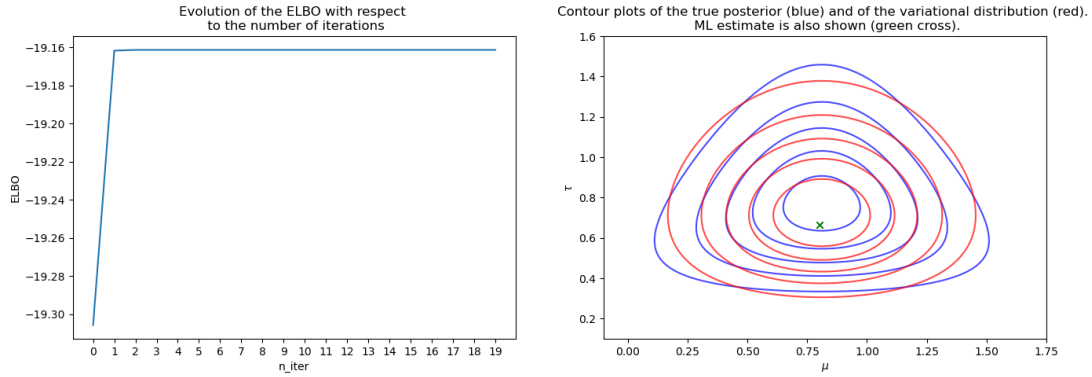
First data set:

Figure 2: Left: Evolution of the ELBO with respect to the number of iterations. Right: Comparison of the true posterior (blue) with the variational distribution (red) obtained with the CAVI algorithm. The ML estimate (green) is also displayed.

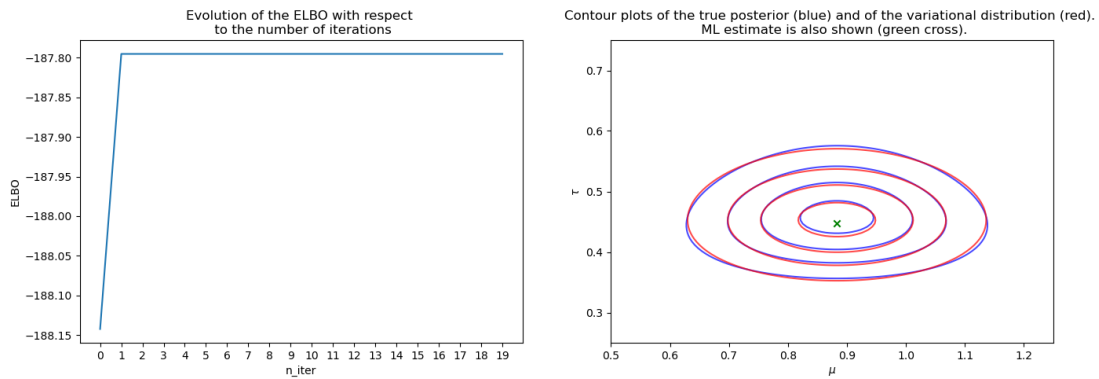
Second data set:

Figure 3: Left: Evolution of the ELBO with respect to the number of iterations. Right: Comparison of the true posterior (blue) with the variational distribution (red) obtained with the CAVI algorithm. The ML estimate (green) is also displayed.

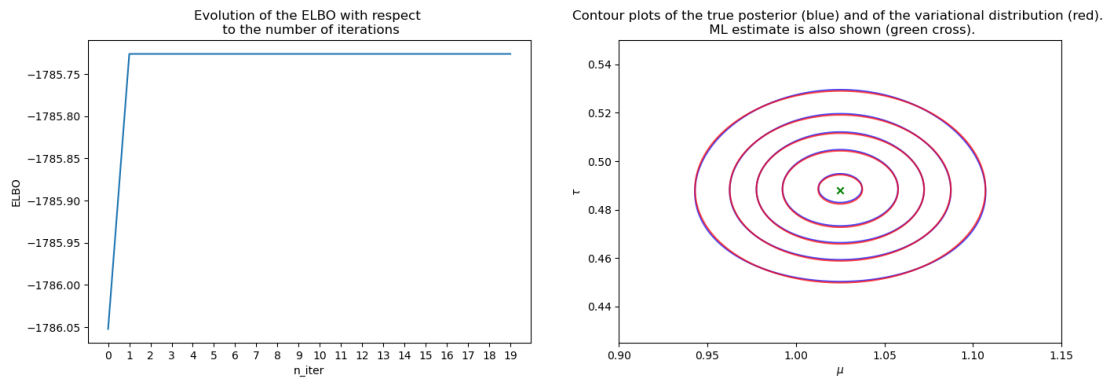
Third data set:

Figure 4: Left: Evolution of the ELBO with respect to the number of iterations. Right: Comparison of the true posterior (blue) with the variational distribution (red) obtained with the CAVI algorithm. The ML estimate (green) is also displayed.

Discussion of the results:

*Concerning the ELBO, we notice that its convergence is very fast, a small number of iterations is enough to converge. Concerning the comparison of the variational distribution obtained at the end of the CAVI algorithm with the true posterior distribution, we notice that **the more data points we have in our data set, the more accurate will be the variational distribution to the true posterior** as the right plot of Figure 5 shows where they almost overlap perfectly. It is the same observation concerning the ML estimate, the more data points we have in our data set, the more accurate it will be as we can see with the right plot of Figure 5 where the ML estimate is effectively at the center of the true posterior.*

1.4 SVI - LDA

Question 1.4.16: According to the Hoffman paper, local hidden variables $z_{n,j}$ are defined by the fact that their complete conditional probability distribution are determined by β the global hidden variables, α the fixed parameters and the other local variables x_n , $z_{n,-j}$ in the n^{th} context. Mathematically, the definition is the following:

$$p(z_{n,j}|x_n, x_{-n}, z_{-n}, z_{n,-j}, \beta, \alpha) = p(z_{n,j}|x_n, z_{n,-j}, \beta, \alpha)$$

Question 1.4.17: In the LDA model of the Hoffman paper, the global hidden variables are the $\beta_k, k = 1 \cdots K$ and the local hidden variables are the $\theta_d, d = 1 \cdots D$ and $z_{d,n}, d = 1 \cdots D, n = 1 \cdots N$.

Question 1.4.18: To compute the ELBO, we used this source [1].

$$\begin{aligned} ELBO = & \sum_{d=1}^D \sum_{n=1}^N \sum_{k=1}^K \sum_{w=1}^W \phi_{dnk} (\psi(\lambda_{kw}) - \psi(\sum_{w'} \lambda_{kw'})) w_{dnw} + \sum_{d=1}^D \sum_{n=1}^N \sum_{k=1}^K \phi_{dnk} (\psi(\gamma_{dk}) - \psi(\sum_{k'} \gamma_{dk'})) \\ & - D \log B(\alpha) + \sum_{k=1}^K \sum_{d=1}^D (\alpha_k - 1) (\psi(\gamma_{dk}) - \psi(\sum_{k'} \gamma_{dk'})) \\ & - K \log B(\eta) + \sum_{k=1}^K \sum_{w=1}^W (\eta_w - 1) (\psi(\lambda_{kw}) - \psi(\sum_{w'} \lambda_{kw'})) \\ & + \sum_{d=1}^D \left[\log B(\gamma_d) + \left(\sum_{k=1}^K \gamma_{dk} - K \right) \psi\left(\sum_{k=1}^K \gamma_{dk}\right) - \sum_{k=1}^K (\gamma_{dk} - 1) \psi(\gamma_{dk}) \right] \\ & + \sum_{k=1}^K \left[\log B(\lambda_k) + \left(\sum_{w=1}^W \lambda_{kw} - W \right) \psi\left(\sum_{w=1}^W \lambda_{kw}\right) - \sum_{w=1}^W (\lambda_{kw} - 1) \psi(\lambda_{kw}) \right] \\ & - \sum_{d=1}^D \sum_{n=1}^N \sum_{k=1}^K \phi_{dnk} \log \phi_{dnk} \end{aligned}$$

Question 1.4.19: Implementing the SVI updates in the notebook and running the two algorithms on the three data sets yielded the following results.

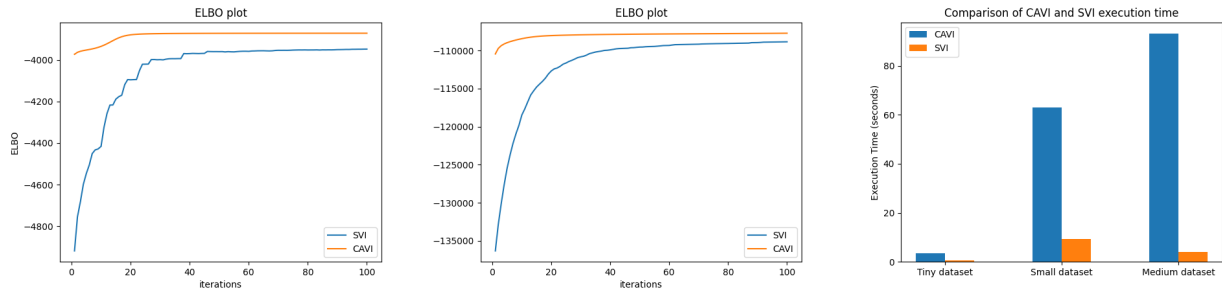


Figure 5: Respectively from left to right: (a) Evolution of the ELBO with respect to the number of iterations for both the CAVI algorithm and the SVI algorithm for the tiny dataset; (b) Same as (a) but for the small dataset; (c) Comparison of the execution time of the CAVI and SVI algorithms for all three datasets.

Discussion of the results:

Concerning the ELBO, we notice that both the one of the SVI and CAVI algorithm converge approximately around the same value (in terms of order of magnitude). The one of the CAVI algorithm seems to be faster in term of iterations to converge. However, what is of interest is the execution time and if we look at the bar plot, we can see that **the SVI algorithm is definitely faster than the CAVI algorithm for all three datasets.**

1.5 BBVI

Question 1.5.20: First, let's derive the gradient of the ELBO with respect to ν :

$$\begin{aligned}\nabla_\nu \mathcal{L} &= \nabla_\nu \int q(\theta|\nu, \epsilon^2) [\log p(x, \theta) - \log q(\theta|\nu, \epsilon^2)] d\theta \\ &= \int \nabla_\nu q(\theta|\nu, \epsilon^2) [\log p(x, \theta) - \log q(\theta|\nu, \epsilon^2)] d\theta - \int q(\theta|\nu, \epsilon^2) \nabla_\nu \log q(\theta|\nu, \epsilon^2) d\theta \\ &= \int q(\theta|\nu, \epsilon^2) \nabla_\nu \log q(\theta|\nu, \epsilon^2) [\log p(x, \theta) - \log q(\theta|\nu, \epsilon^2)] d\theta\end{aligned}$$

where we used $\nabla_\nu \log q(\theta|\nu, \epsilon^2) = \frac{\nabla_\nu q(\theta|\nu, \epsilon^2)}{q(\theta|\nu, \epsilon^2)}$ and $\int q(\theta|\nu, \epsilon^2) \nabla_\nu \log q(\theta|\nu, \epsilon^2) d\theta = \int \nabla_\nu q(\theta|\nu, \epsilon^2) d\theta$ which equals 0 by commuting the gradient and the integral and noticing that the integral equals 1. Therefore, an estimate of the gradient of the ELBO using one sample $z \sim q(\theta|\nu, \epsilon^2)$ is obtained as follows:

$$\begin{aligned}\nabla_\nu \mathcal{L} &\approx \nabla_\nu \log q(z|\nu, \epsilon^2) [\log p(x, z) - \log q(z|\nu, \epsilon^2)] \\ \text{with} \\ \nabla_\nu \log q(z|\nu, \epsilon^2) &= \frac{1}{\epsilon^2} (\log z - \nu) \\ \log q(z|\nu, \epsilon^2) &= -\log(z\epsilon\sqrt{2\pi}) - \frac{1}{2\epsilon^2} (\log z - \nu)^2 \\ \log p(x, z) &= -\log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} (x - z)^2 + \alpha \log \beta - \log \Gamma(\alpha) + (\alpha - 1) \log z - \beta z\end{aligned}$$

Question 1.5.21: Control variates is a variance reduction method used to compute a noisy estimate of the components of the gradient of the ELBO (obtained by the Rao-Blackwellization method) with low variance. Mathematically, if we want to estimate the expected value of f , it introduces a family $\hat{f}(z) = f(z) - a(h(z) - \mathbf{E}[h(z)])$ such that $\mathbf{E}[\hat{f}] = \mathbf{E}[f]$ but $\text{Var}[\hat{f}] < \text{Var}[f]$.

References

1. Blei, D., Ng, A., and Jordan, M. (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022. <https://www.cs.columbia.edu/~blei/papers/BleiLafferty2009.pdf>.


```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import scipy.special as sp_spec
import pandas as pd
np.random.seed(3)
```

Assignment 1.3 - CAVI

Consider the model defined by Equation (10.21)-(10.23) in Bishop, for which DGM is presented below:



Question 1.3.12:

Implement a function that generates data points for the given model.

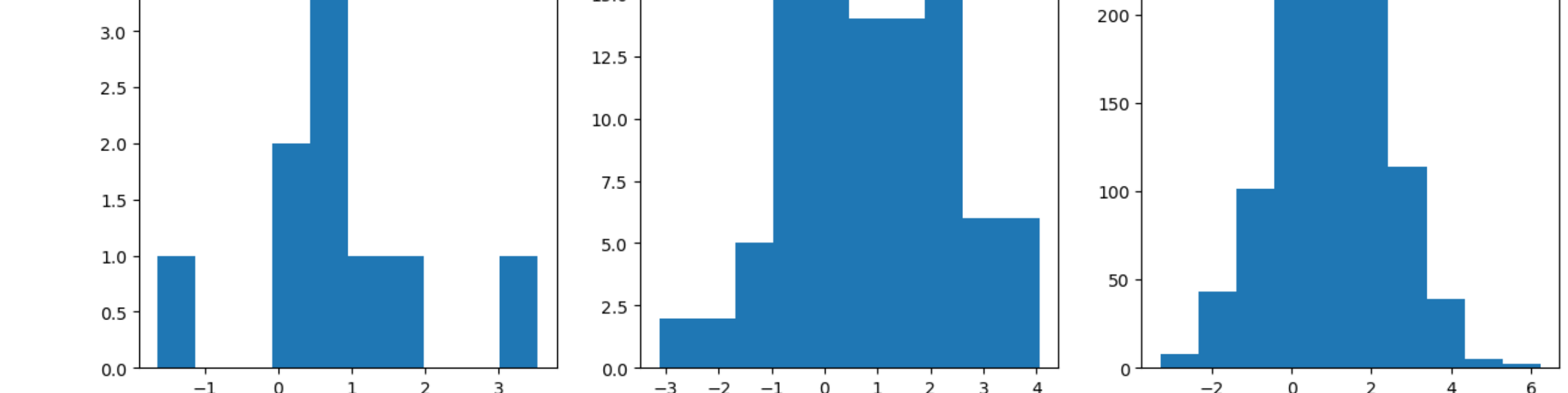
```
In [2]: def generate_data(mu, tau, N):
# Insert your code here
D = np.random.normal(mu, np.sqrt(1/tau), N)
return D
```

Set $\mu = 1$, $\tau = 0.5$ and generate datasets with size $N=10, 100, 1000$. Plot the histogram for each of 3 datasets you generated.

```
In [3]: mu = 1
tau = 0.5

dataset_1 = generate_data(mu, tau, 10)
dataset_2 = generate_data(mu, tau, 100)
dataset_3 = generate_data(mu, tau, 1000)

# Visualize the datasets via histograms
# Insert your code here
list_dataset = [dataset_1, dataset_2, dataset_3]
fig, axes = plt.subplots(1, 3, figsize=(15, 5))
for i in range(len(list_dataset)):
    dataset = list_dataset[i]
    axes[i].hist(dataset, label=f'N={len(dataset)}')
# sns.histplot(dataset, label=f'N={len(dataset)}', ax=axes[i])
axes[i].legend()
plt.suptitle('Histogram of $N(1, 2)$ with respectively, from left to right, $N=10$, $N=100$ and $N=1000$ data points')
plt.savefig('histogram.png')
plt.show()
```



Question 1.3.13:

Find ML estimates of the variables μ and τ

```
In [4]: def ML_est(data):
# insert your code
mu_ml = np.mean(data)
tau_ml = 1/np.var(data, ddof=0)
return mu_ml, tau_ml
```

```
In [5]: print('For N=10, the Maximum Likelihood estimates for the parameters mu and tau are respectively', np.round(ML_est(dataset_1),2))
print('For N=100, the Maximum Likelihood estimates for the parameters mu and tau are respectively', np.round(ML_est(dataset_2),2))
print('For N=1000, the Maximum Likelihood estimates for the parameters mu and tau are respectively', np.round(ML_est(dataset_3),2))

For N=10, the Maximum Likelihood estimates for the parameters mu and tau are respectively [0.8 0.67]
For N=100, the Maximum Likelihood estimates for the parameters mu and tau are respectively [0.88 0.45]
For N=1000, the Maximum Likelihood estimates for the parameters mu and tau are respectively [1.03 0.49]
```

We observe that the more data we have, the closer the parameters are to the true values: [1.00, 0.50].

Question 1.3.14:

You will implement the VI algorithm for the variational distribution in Equation (10.24) in Bishop. Start with introducing the prior parameters:

```
In [6]: # prior parameters
mu_0 = 1
lambda_0 = 0.5
a_0 = 2
b_0 = 1
```

Continue with a helper function that computes ELBO:

```
In [7]: def compute_elbo(D, a_0, b_0, mu_0, lambda_0, a_N, b_N, mu_N, lambda_N):
# given the prior and posterior parameters together with the data,
# compute ELBO here
N = len(D)
# eps = 0.0001
E_log_tau = sp_spec.digamma(a_N) - np.log(b_N)
CR_likelihood = N/2*np.log(2*np.pi) + N/2*E_log_tau - 0.5*(a_N/b_N)*np.sum(D**2 + 1/lambda_N + mu_N**2 - 2*D*mu_N)
CR_mu = 0.5*np.log(lambda_0/(2*np.pi)) + 0.5*E_log_tau - lambda_0/2*(a_N/b_N)*(1/lambda_N + mu_N**2 + mu_0**2 - 2*mu_0*mu_N)
CR_tau = (a_0 - 1)*E_log_tau + a_0*np.log(b_0) + 0.5
H_mu = 0.5*np.log(2*np.pi/lambda_N) + 0.5
H_tau = a_N - np.log(b_N) + sp_spec.gammaln(a_N) + (1 - a_N)*sp_spec.digamma(a_N)
elbo = CR_likelihood + CR_mu + CR_tau + H_mu + H_tau
return elbo
```

Now, implement the CAVI algorithm:

```
In [8]: def update_q_mu(D, mu_0, lambda_0, a_N, b_N):
N = len(D)
mu_N = (lambda_0*mu_0 + np.sum(D)) / (lambda_0 + N)
lambda_N = (lambda_0 + N) * (a_N / b_N)
return mu_N, lambda_N
```

```
In [9]: def update_q_tau(D, a_0, b_0, mu_0, lambda_0, mu_N, lambda_N):
N = len(D)
a_N = a_0 + (N+1)/2
b_N = b_0 + 0.5*(np.sum(D**2) + lambda_0*(mu_0**2 - 2*(np.sum(D) + lambda_0*mu_0)*mu_N + (lambda_0 + N)*(1/lambda_N + mu_N**2)))
return a_N, b_N
```

```
In [10]: def CAVI(D, a_0, b_0, mu_0, lambda_0, n_iter):
# make an initial guess for the expected value of tau
initial_guess_exp_tau = 5
a_N = 20
b_N = 10

# CAVI iterations ...
# save ELBO for each iteration, plot them afterwards to show convergence
elbos = np.zeros(n_iter)
for i in range(0, n_iter):
    mu_N, lambda_N = update_q_mu(D, mu_0, lambda_0, a_N, b_N)
    a_N, b_N = update_q_tau(D, a_0, b_0, mu_0, lambda_0, mu_N, lambda_N)
    elbos[i] = compute_elbo(D, a_0, b_0, mu_0, lambda_0, a_N, b_N, mu_N, lambda_N)

return a_N, b_N, mu_N, lambda_N, elbos
```

Question 1.3.15:

What is the exact posterior? First derive it in closed form, and then implement a function that computes it for the given parameters:

```
In [11]: def compute_exact_posterior(D, a_0, b_0, mu_0, lambda_0):
# your implementation
N = len(D)
mu_0_star = (np.sum(D)+lambda_0*mu_0)/(lambda_0+N)
lambda_0_star = lambda_0 + N
a_0_star = a_0 + N/2
b_0_star = b_0 + np.sum(D**2)/2 + lambda_0*(mu_0**2)/2 - 0.5*((np.sum(D) + lambda_0*mu_0)**2)/(lambda_0 + N)

return mu_0_star, lambda_0_star, a_0_star, b_0_star
```

Question 1.3.16:

Run the VI algorithm on the datasets. Compare the inferred variational distribution with the exact posterior and the ML estimate. Visualize the results and discuss your findings.

```
In [12]: def log_pdf_exact_post(mu, tau, mu_0_star, lambda_0_star, a_0_star, b_0_star):
return a_0_star * np.log(b_0_star) - sp_spec.gammaln(a_0_star) + 0.5 * np.log(lambda_0_star / (2 * np.pi)) + (a_0_star - 1/2) * np.log(tau) - 1/2 * tau * (mu - mu_0_star)**2

def log_pdf_approx_post(mu, tau, mu_N, lambda_N, a_N, b_N):
return a_N * np.log(b_N) - sp_spec.gammaln(a_N) + 0.5 * np.log(lambda_N / (2 * np.pi)) + (a_N - 1) * np.log(tau) - b_N * tau - 0.5 * lambda_N * (mu - mu_N)**2
```

```
First dataset

In [13]: dataset = dataset_1
n_iter = 20

# ML estimates, VI estimates, true parameters
mu_ml, tau_ml = ML_est(dataset)
a_N, b_N, mu_N, lambda_N, elbos = CAVI(dataset, a_0, b_0, mu_0, lambda_0, n_iter)
mu_0_star, lambda_0_star, a_0_star, b_0_star = compute_exact_posterior(dataset, a_0, b_0, mu_0, lambda_0)

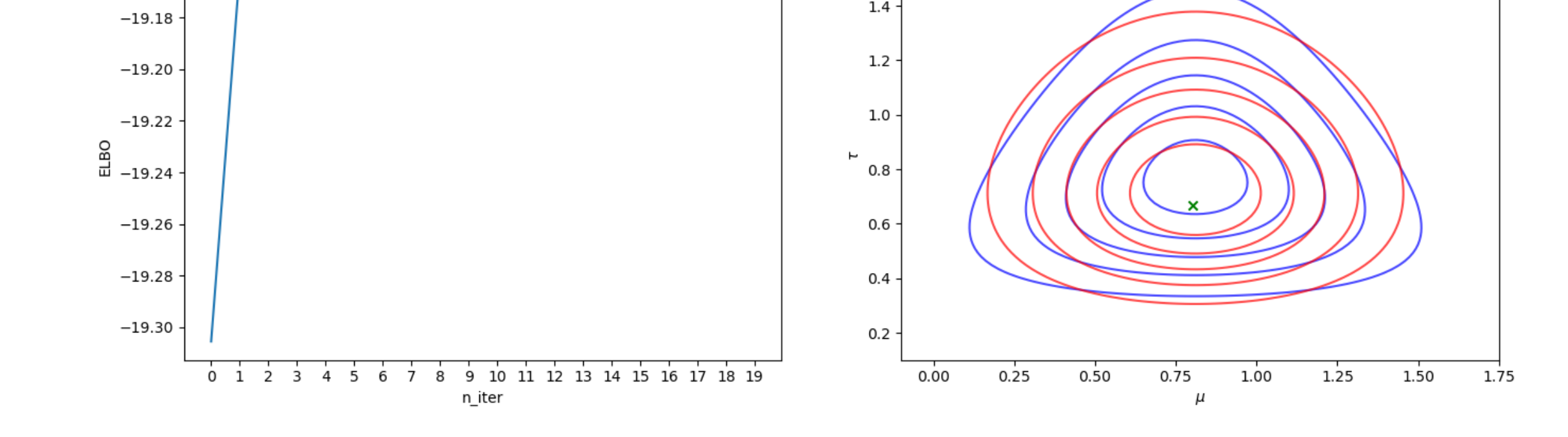
# plot elbos, show convergence
fig, ax = plt.subplots(1, 2, figsize=(16, 5))
ax[0].plot(elbos)
ax[0].set_title('Evolution of the ELBO with respect \n to the number of iterations')
ax[0].set_ylabel('ELBO')
ax[0].set_xlabel('n_iter')
ax[0].set_xticks(np.arange(0,n_iter))

# compare exact_post_dist with the CAVI result ( = q(a_N, b_N, mu_N, lambda_N) ) using for ex. contour plots, show also ML estimate on this plot
mu = np.linspace(-0.1, 1.75, 100)
tau = np.linspace(0.1, 1.6, 100)
muv, tauv = np.meshgrid(mu, tau)

# compute exact posterior and approx posterior on the grid defined above
lg_pdf_ex_post = log_pdf_exact_post(muv, tauv, mu_0_star, lambda_0_star, a_0_star, b_0_star)
pdf_ex_post = np.exp(lg_pdf_ex_post - sp_spec.logsumexp(lg_pdf_ex_post))

lg_pdf_aprx_post = log_pdf_approx_post(muv, tauv, mu_N, lambda_N, a_N, b_N)
pdf_aprx_post = np.exp(lg_pdf_aprx_post - sp_spec.logsumexp(lg_pdf_aprx_post))

# graph contour
levels = 5
ax[1].contour(muv, tauv, pdf_ex_post, colors='blue', levels=levels, alpha=0.7)
ax[1].contour(muv, tauv, pdf_aprx_post, colors='red', levels=levels, alpha=0.7)
ax[1].scatter(mu_ml, tau_ml, color='green', label='ML Estimates', marker='x')
ax[1].set_title('Contour plots of the true posterior (blue) and of the variational distribution (red). \n ML estimate is also shown (green cross).')
ax[1].set_ylabel(r'$\tau$')
ax[1].set_xlabel(r'$\mu$')
# ax.set_title('Contours des Distributions Réelle et Approximative')
# plt.savefig('ELBO_and_Contour_dataset_1.png')
plt.show()
```



```
Second dataset

In [14]: dataset = dataset_2
n_iter = 20

# ML estimates, VI estimates, true parameters
mu_ml, tau_ml = ML_est(dataset)
a_N, b_N, mu_N, lambda_N, elbos = CAVI(dataset, a_0, b_0, mu_0, lambda_0, n_iter)
mu_0_star, lambda_0_star, a_0_star, b_0_star = compute_exact_posterior(dataset, a_0, b_0, mu_0, lambda_0)

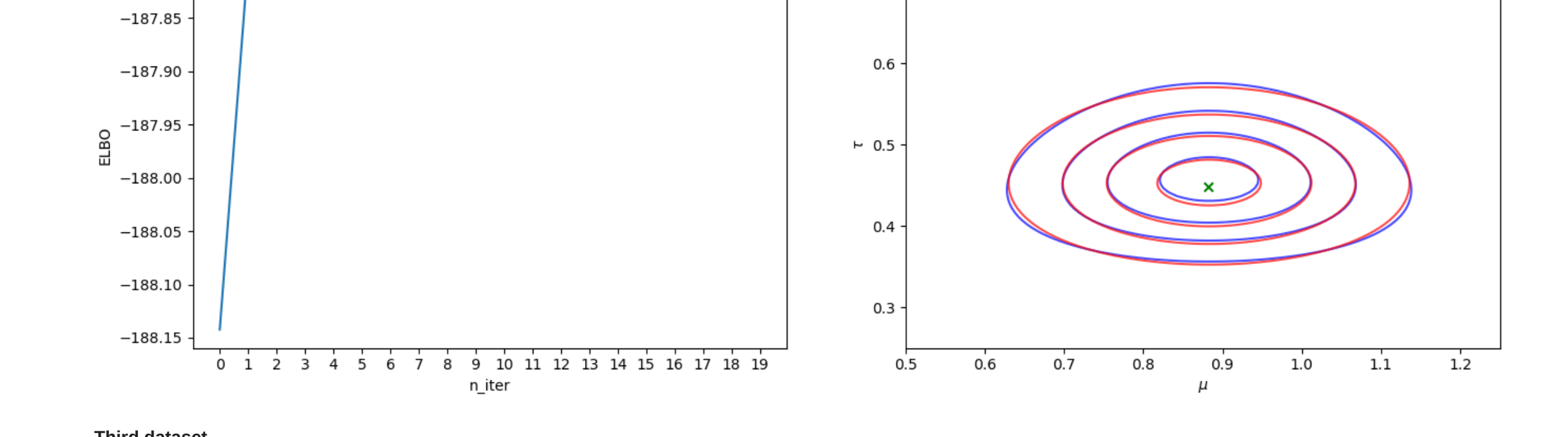
# plot elbos, show convergence
fig, ax = plt.subplots(1, 2, figsize=(16, 5))
ax[0].plot(elbos)
ax[0].set_title('Evolution of the ELBO with respect \n to the number of iterations')
ax[0].set_ylabel('ELBO')
ax[0].set_xlabel('n_iter')
ax[0].set_xticks(np.arange(0,n_iter))

# compare exact_post_dist with the CAVI result ( = q(a_N, b_N, mu_N, lambda_N) ) using for ex. contour plots, show also ML estimate on this plot
mu = np.linspace(0.5, 1.25, 100)
tau = np.linspace(0.25, 0.75, 100)
muv, tauv = np.meshgrid(mu, tau)

# compute exact posterior and approx posterior on the grid defined above
lg_pdf_ex_post = log_pdf_exact_post(muv, tauv, mu_0_star, lambda_0_star, a_0_star, b_0_star)
pdf_ex_post = np.exp(lg_pdf_ex_post - sp_spec.logsumexp(lg_pdf_ex_post))

lg_pdf_aprx_post = log_pdf_approx_post(muv, tauv, mu_N, lambda_N, a_N, b_N)
pdf_aprx_post = np.exp(lg_pdf_aprx_post - sp_spec.logsumexp(lg_pdf_aprx_post))

# graph contour
levels = 5
ax[1].contour(muv, tauv, pdf_ex_post, colors='blue', levels=levels, alpha=0.7)
ax[1].contour(muv, tauv, pdf_aprx_post, colors='red', levels=levels, alpha=0.7)
ax[1].scatter(mu_ml, tau_ml, color='green', label='ML Estimates', marker='x')
ax[1].set_title('Contour plots of the true posterior (blue) and of the variational distribution (red). \n ML estimate is also shown (green cross).')
ax[1].set_ylabel(r'$\tau$')
ax[1].set_xlabel(r'$\mu$')
# ax.set_title('Contours des Distributions Réelle et Approximative')
# plt.savefig('ELBO_and_Contour_dataset_2.png')
plt.show()
```



```
Third dataset

In [15]: dataset = dataset_3
n_iter = 20

# ML estimates, VI estimates, true parameters
mu_ml, tau_ml = ML_est(dataset)
a_N, b_N, mu_N, lambda_N, elbos = CAVI(dataset, a_0, b_0, mu_0, lambda_0, n_iter)
mu_0_star, lambda_0_star, a_0_star, b_0_star = compute_exact_posterior(dataset, a_0, b_0, mu_0, lambda_0)

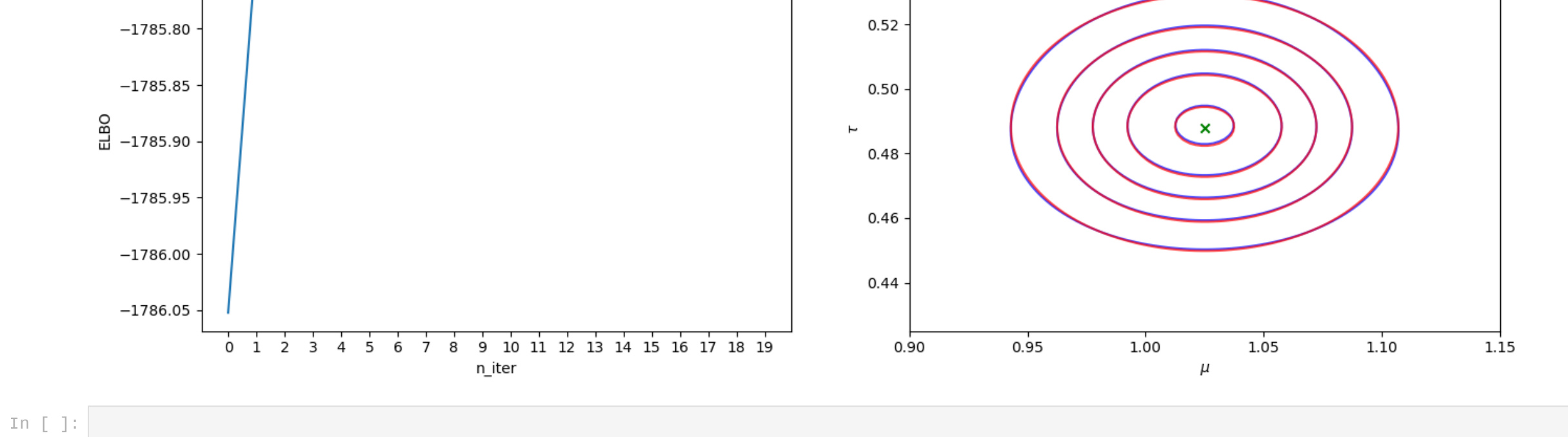
# plot elbos, show convergence
fig, ax = plt.subplots(1, 2, figsize=(16, 5))
ax[0].plot(elbos)
ax[0].set_title('Evolution of the ELBO with respect \n to the number of iterations')
ax[0].set_ylabel('ELBO')
ax[0].set_xlabel('n_iter')
ax[0].set_xticks(np.arange(0,n_iter))

# compare exact_post_dist with the CAVI result ( = q(a_N, b_N, mu_N, lambda_N) ) using for ex. contour plots, show also ML estimate on this plot
mu = np.linspace(0.9, 1.15, 100)
tau = np.linspace(0.425, 0.550, 100)
muv, tauv = np.meshgrid(mu, tau)

# compute exact posterior and approx posterior on the grid defined above
lg_pdf_ex_post = log_pdf_exact_post(muv, tauv, mu_0_star, lambda_0_star, a_0_star, b_0_star)
pdf_ex_post = np.exp(lg_pdf_ex_post - sp_spec.logsumexp(lg_pdf_ex_post))

lg_pdf_aprx_post = log_pdf_approx_post(muv, tauv, mu_N, lambda_N, a_N, b_N)
pdf_aprx_post = np.exp(lg_pdf_aprx_post - sp_spec.logsumexp(lg_pdf_aprx_post))

# graph contour
levels = 5
ax[1].contour(muv, tauv, pdf_ex_post, colors='blue', levels=levels, alpha=0.7)
ax[1].contour(muv, tauv, pdf_aprx_post, colors='red', levels=levels, alpha=0.7)
ax[1].scatter(mu_ml, tau_ml, color='green', label='ML Estimates', marker='x')
ax[1].set_title('Contour plots of the true posterior (blue) and of the variational distribution (red). \n ML estimate is also shown (green cross).')
ax[1].set_ylabel(r'$\tau$')
ax[1].set_xlabel(r'$\mu$')
# ax.set_title('Contours des Distributions Réelle et Approximative')
# plt.savefig('ELBO_and_Contour_dataset_3.png')
plt.show()
```



In []:

