



KTH ROYAL INSTITUTE OF TECHNOLOGY

Machine Learning Advanced Course

ASSIGNMENT 2

Castellano Giovanni

October 3, 2022

2.1 Dependencies in a Directed Graphical Model

Question 2.1.1

Yes

Question 2.1.2

No

Question 2.1.3

Yes

Question 2.1.4

No

Question 2.1.5

No

Question 2.1.6

No

2.2 Likelihood of a Tree Graphical Model

The algorithm is based on the algorithm studied in module 7 lecture 6 where we studied how to compute the probability that a node X_u has a certain value after observing the values of the leaves. In particular, during the lecture we defined two probabilities; one concerning the nodes above X_u , and the second one concerning the nodes below X_u . To compute the probability of the leaves I computed the value of the second function on the root node for all possible values that the root can assume. Subsequently, I multiplied the results for the probability that the root assumes those values and summed them up.

- Small Tree:
0.009786075668602368
0.015371111945909397
0.02429470256988136
0.005921848333806081
0.016186321212555956

- Medium Tree:
1.7611947348905713e-18
2.996933026124685e-18
2.891411201505415e-18
4.6788419411270006e-18
5.664006737201378e-18
- Large Tree:
3.63097513075208e-69
3.9405421986921234e-67
5.549061147187144e-67
9.89990102807915e-67
3.11420969368965e-72

2.3 Simple Variational Inference

Exact posterior

Suppose to have some data $\mathbf{X} = x_1, \dots, x_n \sim \text{Normal}(\mu, \lambda^{-1})$, where $\lambda = 1/\sigma^2$. Suppose moreover, that we do not know the exact values of the parameters μ and λ but we know that $(\mu, \lambda) \sim \text{NormalGamma}(m, c, a, b)$, and therefore:

$$\mu|\lambda \sim \text{Normal}(m, (c\lambda)^{-1}) \quad (1)$$

$$\lambda \sim \text{Gamma}(a, b) \quad (2)$$

$$P(\mu, \lambda) = P(\mu|\lambda)P(\lambda) = \text{Normal}(\mu|m, (c\lambda)^{-1}) * \text{Gamma}(\lambda|a, b) \quad (3)$$

What we want to prove here, is that the NormalGamma is a conjugate prior for the Normal distribution, i.e.:

$$P(\mu, \lambda|\mathbf{X}) \sim \text{NormalGamma}(M, C, A, B) \quad (4)$$

We have:

$$\text{NormalGamma}(\mu, \lambda|m, c, a, b) = \text{Normal}(\mu|m, (c\lambda)^{-1})\text{Gamma}(\lambda|a, b) \quad (5)$$

$$\text{NormalGamma}(\mu, \lambda|m, c, a, b) = \sqrt{\frac{c\lambda}{2\pi}} e^{-0.5c\lambda(\mu-m)^2} \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda} \quad (6)$$

$$\sqrt{\frac{c\lambda}{2\pi}} e^{-0.5c\lambda(\mu-m)^2} \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda} \propto \lambda^{a-1/2} e^{-0.5\lambda(c\mu^2 - 2cm\mu + cm^2 + 2b)} \quad (7)$$

Where we developed the square and joined the exponentials; moreover for any x :

$$Normal(x|\mu, \lambda^{-1}) = \sqrt{\frac{\lambda}{2\pi}} e^{-0.5\lambda(x-\mu)^2} \quad (8)$$

$$\sqrt{\frac{\lambda}{2\pi}} e^{-0.5\lambda(x-\mu)^2} \propto \lambda^{1/2} e^{-0.5\lambda(\mu^2 - 2x\mu + x^2)} \quad (9)$$

thus:

$$P(\mu, \lambda|\mathbf{X}) \propto P(\mu, \lambda) * P(\mathbf{X}|\mu, \lambda) \quad (10)$$

using equation 7 and equation 9 we get:

$$P(\mu, \lambda)P(\mathbf{X}|\mu, \lambda) \propto \lambda^{a-1/2} e^{-0.5\lambda(c\mu^2 - 2cm\mu + cm^2 + 2b)} \lambda^{n/2} e^{-0.5\lambda(n\mu^2 - 2(\sum_i^n x_i)\mu + \sum_i^n x_i^2)} \quad (11)$$

where we are supposing that the data points are drawn independently.

$$P(\mu, \lambda)P(\mathbf{X}|\mu, \lambda) \propto \lambda^{a+n/2-1/2} e^{-0.5\lambda((c+n)\mu^2 - 2(cm + \sum_i^n x_i)\mu + cm^2 + 2b + \sum_i^n x_i^2)} \quad (12)$$

therefore:

$$P(\mu, \lambda)P(\mathbf{X}|\mu, \lambda) \propto \lambda^{A-1/2} e^{-0.5\lambda(C\mu^2 - 2(CM)\mu + CM^2 + 2B)} \quad (13)$$

that is up to a constant a NormalGamma distribution, and we have defined:

$$A = a + n/2 \quad (14)$$

$$C = c + n \quad (15)$$

$$CM = cm + \sum_i^n x_i \quad (16)$$

$$CM^2 + 2B = cm^2 + 2b + \sum_i^n x_i^2 \quad (17)$$

Notice that equations 15, 16 and 17 contains 3 unknown values, and since they are exactly 3 equations, it is possible to join them and show that they hold for these values of M and B:

$$M = \frac{cm + \sum_i^n x_i}{c + n} \quad (18)$$

$$B = b + 0.5(cm^2 - CM^2 + \sum_i^n x_i^2) \quad (19)$$

substituting A, B, C and M in equation 13 we would get again 12. This proof is based on the proof provided by Rebecca C. Steorts and Lei Qian [1].

Approximation of the posterior with the VI algorithm

The following plots illustrate the prior probability of the parameters (green), the exact posterior (blue), and the approximated posterior (red). Moreover, we can observe in blue the parameters of the gaussian from which the points are drawn, and in orange the same parameters estimated from the generated data points.

To begin, figure 1 shows how the distributions appear when the posterior is approximated only using 2 data points. Moreover, in this case the prior is centered on the true parameters. In this case, the algorithm converged after 3 iterations. We notice that the center of the posterior is close to the center of the prior, and far from the estimated parameters (orange point). When we increase the number of data points (figure 2) from 2 to 50 however, we observe that the center of the posterior moves closer to the orange point; and If we increase the number of points from 50 to 500 (figure 3), we can notice that now the posterior is perfectly centered on the orange point, moreover, the two points are very close (this is due to the law of large numbers). It is also interesting that the variance of the posterior decreased as the number of points increased. In all the three cases, we observed a good approximation of the posterior. However, the approximation seems to get better as the number of data points increase.

In figure 4 we show an example where the approximation of the posterior is not very good. As we can observe, with this configuration of parameters, the prior distribution assumes a skewed shape; in particular, we notice some tails on the bottom. Such tails also appear in the exact posterior, however, we do not observe them in the approximated posterior; the shape of the approximation remains very simple. This is of course due to the assumption of factorization. We notice

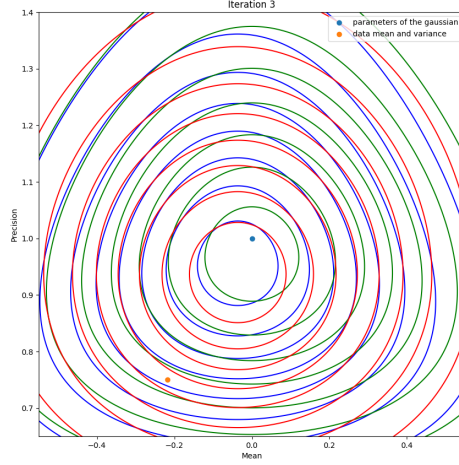


Figure 1: In green the prior; in blue the exact posterior; in red the approximated posterior. $\text{mean} = 0$, $\text{precision} = 1$, $a_0 = 20$, $b_0 = 20$, $\mu_0 = 0$, $\lambda_0 = 10$, $N_{\text{points}} = 2$

however, that such tails disappear from the posterior as we increase the number of points (figure 5), in such case we get again a good approximation.

Lastly, in figure 6, we can observe a case where the posterior, the prior, and the true parameters are far from each other. Also in this case we notice that the exact posterior is a bit skewed and the approximation is not perfect. Moreover, the posterior is squeezed against the bottom axis. This is due to the particular values of a_0 and b_0 .

In general the algorithm needs a very small number of iterations to converge.

The implementation of the algorithm is partially based on the code provided by MICHAEL P. J. CAMILLERI [2]

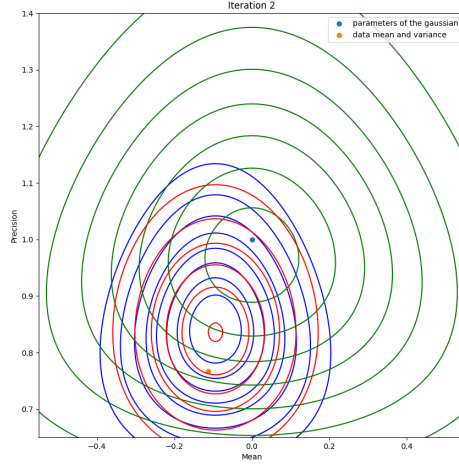


Figure 2: In green the prior; in blue the exact posterior; in red the approximated posterior. $\text{mean} = 0$, $\text{precision} = 1$, $a_0 = 20$, $b_0 = 20$, $\mu_0 = 0$, $\lambda_0 = 10$, $N_{\text{points}} = 50$.

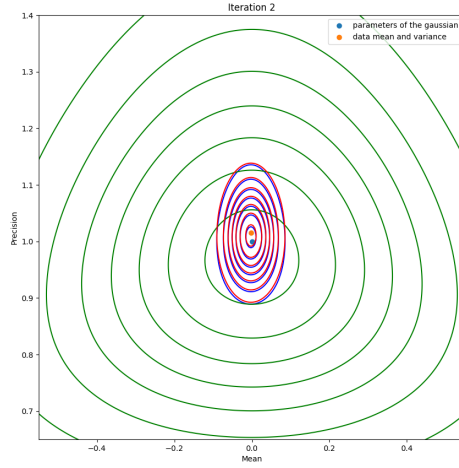


Figure 3: In green the prior; in blue the exact posterior; in red the approximated posterior. $\text{mean} = 0$, $\text{precision} = 1$, $a_0 = 20$, $b_0 = 20$, $\mu_0 = 0$, $\lambda_0 = 10$, $N_{\text{points}} = 500$.

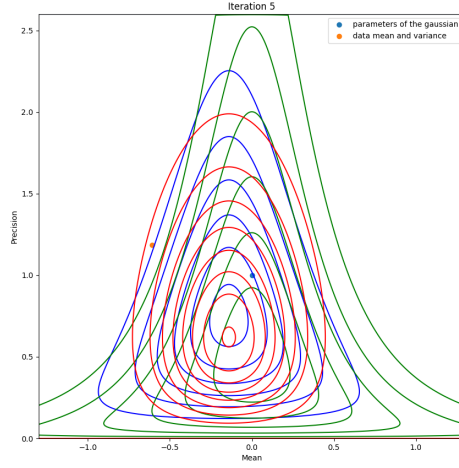


Figure 4: In green the prior; in blue the exact posterior; in red the approximated posterior. $\text{mean} = 0$, $\text{precision} = 1$, $a_0 = 1$, $b_0 = 1$, $\mu_0 = 0$, $\lambda_0 = 10$, $N_{\text{points}} = 3$.

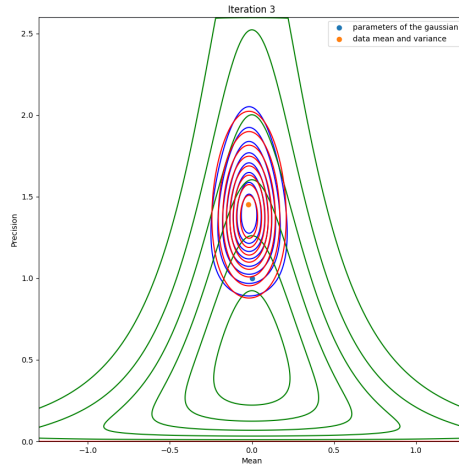


Figure 5: In green the prior; in blue the exact posterior; in red the approximated posterior. $\text{mean} = 0$, $\text{precision} = 1$, $a_0 = 1$, $b_0 = 1$, $\mu_0 = 0$, $\lambda_0 = 10$, $N_{\text{points}} = 50$.

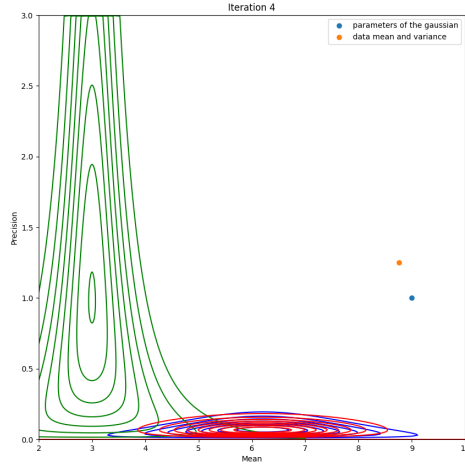


Figure 6: In green the prior; in blue the exact posterior; in red the approximated posterior. $\text{mean} = 9$, $\text{precision} = 1$, $a_0 = 1$, $b_0 = 5$, $\mu_0 = 3$, $\lambda_0 = 4$, $N_{\text{points}} = 5$.

2.4 Super Epicentra - Expectation-Maximization

Derivation of the algorithm

To derive the algorithm we will start from the general formulation of the EM algorithm.

Expectation step:

compute the responsibilities:

$$\gamma(Z_{ij}) = P(Z_j = 1 | D_i, \vec{\Theta}^{old}) \quad (20)$$

for each $j \in 1, \dots, K$ and $i \in 1, \dots, N$; where K is the number of distributions in the mixture and N is the number of data points; Z_j is the j th component of the vector \mathbf{Z} that is equal to one when a data point $D_i \in \mathbf{D}$ is generated by the j th distribution of the mixture.

Maximization step:

Maximize the log likelihood:

$$\log(P(\mathbf{D}|\vec{\Theta})) = \sum_i^N \sum_j^K \gamma(Z_{ij}) \log(P(D_i, Z_j | \vec{\Theta})) \quad (21)$$

In our case, $\mathbf{D} = (\mathbf{X}, \mathbf{S})$, where \vec{x}_i and s_i are drawn respectively from a Gaussian and a Poisson distribution. Therefore in the expectation step we get:

$$\begin{aligned} \gamma(Z_{ij}) &= P(Z_j = 1 | s_i, \vec{x}_i, \vec{\Theta}^{old}) \\ &= \frac{P(Z_j = 1)P(s_i, \vec{x}_i | Z_j = 1, \vec{\Theta}^{old})}{\sum_l^k P(Z_l = 1)P(\vec{x}_i, s_i | Z_l = 1, \vec{\Theta}^{old})} \\ &= \frac{P(Z_j = 1)P(\vec{x}_i | Z_j = 1, \vec{\Theta}^{old})P(s_i | Z_j = 1, \vec{\Theta}^{old})}{\sum_l^k P(Z_l = 1)P(\vec{x}_i | Z_l = 1, \vec{\Theta}^{old})P(s_i | Z_l = 1, \vec{\Theta}^{old})} \\ &= \frac{\pi_j N(\vec{x}_i | \vec{u}_j, \Sigma_j) Pois(s_i | \lambda_j)}{\sum_l^k \pi_l N(\vec{x}_i | \vec{u}_l, \Sigma_l) Pois(s_i | \lambda_l)} \end{aligned} \quad (22)$$

Where we used the Bayes formula, and we splitted the probability using the conditional independence given Z_j . Moreover, we kept the vector $\vec{\Theta}^{old}$ on the right side of the probability when we used the Bayes formula since it is just a constant. Lastly, we named π_j the probability that $Z_j = 1$, therefore, the vector π represents the weights of the mixture.

Let's proceed now with maximization step:

$$\begin{aligned}
\log(P(\mathbf{D}|\vec{\Theta})) &= \log(P(\mathbf{X}, \mathbf{S}|\vec{\Theta})) \\
&= \sum_i^N \sum_j^K \gamma(Z_{ij}) \log(P(\vec{x}_i, s_i, Z_j = 1|\vec{\Theta})) \\
&= \sum_i^N \sum_j^K \gamma(Z_{ij}) \log(P(\vec{x}_i, s_i|Z_j = 1, \vec{\Theta})P(Z_j = 1|\vec{\Theta})) \\
&= \sum_i^N \sum_j^K \gamma(Z_{ij}) \log(P(\vec{x}_i|Z_j = 1, \vec{\Theta})P(s_i|Z_j = 1, \vec{\Theta})P(Z_j = 1|\vec{\Theta})) \\
&= \sum_i^N \sum_j^K \gamma(Z_{ij}) [\log(P(\vec{x}_i|Z_j = 1, \vec{\Theta})) + \log(P(s_i|Z_j = 1, \vec{\Theta})) + \log(P(Z_j = 1|\vec{\Theta}))]
\end{aligned} \tag{23}$$

Where again we used the conditional independence to split the probabilities. To maximize this quantity we have to take the derivative with respect to each component of $\vec{\Theta}$ and set them equal to zero. Remember that the vector $\vec{\Theta}$ contains the parameters $\vec{\mu}_l$, Σ_l , λ_l , and π_l for each $l \in 1, \dots, K$. Moreover, only the l th term of the sum over K depends on μ_l , Σ_l , λ_l , and π_l . this helps us to simplify the derivatives.

Let's start from the derivative with respect to $\vec{\mu}_l$:

$$\begin{aligned}
\frac{\partial \log(P(\mathbf{X}, \mathbf{S}))}{\partial \vec{\mu}_l} &= \sum_i^N \gamma(z_{il}) \frac{\partial \log(P(\vec{x}_i|Z_l = 1, \vec{\mu}_l, \Sigma_l))}{\partial \vec{\mu}_l} \\
&= \sum_i^N \gamma(z_{il}) \frac{\partial [\log(\frac{1}{\sqrt{(2\pi)^2 |\Sigma_l|}}) - 0.5(\vec{x}_i - \vec{u}_l)^T \Sigma_l^{-1} (\vec{x}_i - \vec{u}_l)]}{\partial \vec{\mu}_l} \\
&= \sum_i^N \gamma(z_{il}) (\vec{x}_i - \vec{u}_l)^T \Sigma_l^{-1}
\end{aligned} \tag{24}$$

Where we used the the formula of the derivative of a quadratic form with respect to the vector, and the chain rule. Setting equation 24 equal to zero:

$$\sum_i^N \gamma(z_{il}) (\vec{x}_i - \vec{u}_l)^T \Sigma_l^{-1} = 0 \tag{25}$$

$$\sum_i^N \gamma(z_{il}) (\vec{x}_i - \vec{u}_l)^T = 0 \tag{26}$$

$$\sum_i^N \gamma(z_{il}) \vec{x}_i - \sum_i^N \gamma(z_{il}) \vec{u}_l = 0 \quad (27)$$

$$\vec{u}_l = \frac{\sum_i^N \gamma(z_{il}) \vec{x}_i}{\sum_i^N \gamma(z_{il})} \quad (28)$$

Now we can take the derivative with respect to Σ_l , we are supposing that such matrix is diagonal $\Sigma_l = \text{diag}(\sigma_{11}^l, \sigma_{22}^l)$, therefore, to avoid working with matrices we can take the derivatives with respect to the two elements of the diagonal; if we consider that $\vec{\mu}_l = (\mu_{l1}, \mu_{l2})$ and $\vec{x}_i = (x_{i1}, x_{i2})$, we have:

$$\begin{aligned} \frac{\partial \log(P(\mathbf{X}, \mathbf{S}))}{\partial \sigma_{11}^l} &= \sum_i^N \gamma(z_{il}) \frac{\partial \log(P(x_i | Z_l = 1, \vec{\mu}_l, \Sigma_l))}{\partial \sigma_{11}^l} \\ &= \sum_i^N \gamma(z_{il}) \left[\frac{\partial \log\left(\frac{1}{\sqrt{(2\pi)^2 \sigma_{11}^l \sigma_{22}^l}}\right)}{\partial \sigma_{11}^l} - \frac{\partial 0.5((x_{i1} - u_{l1})^2 \frac{1}{\sigma_{11}^l} + (x_{i2} - u_{l2})^2 \frac{1}{\sigma_{22}^l})}{\partial \sigma_{11}^l} \right] \end{aligned} \quad (29)$$

where we calculated the determinant of the matrix, the inverse, and we developed the quadratic form. Considering that:

$$\begin{aligned} \frac{\partial \log\left(\frac{1}{\sqrt{(2\pi)^2 \sigma_{11}^l \sigma_{22}^l}}\right)}{\partial \sigma_{11}^l} &= \frac{\partial [-\log(\sqrt{(2\pi)^2 \sigma_{11}^l \sigma_{22}^l})]}{\partial \sigma_{11}^l} \\ &= -\frac{1}{\sqrt{(2\pi)^2 \sigma_{11}^l \sigma_{22}^l}} \frac{1}{2\sqrt{(2\pi)^2 \sigma_{11}^l \sigma_{22}^l}} (2\pi)^2 \sigma_{22}^l \\ &= -0.5 \frac{1}{\sigma_{11}^l} \end{aligned} \quad (30)$$

moreover,

$$\frac{\partial 0.5((x_{i1} - u_{l1})^2 \frac{1}{\sigma_{11}^l} + (x_{i2} - u_{l2})^2 \frac{1}{\sigma_{22}^l})}{\partial \sigma_{11}^l} = -0.5(x_{i1} - u_{l1})^2 \frac{1}{(\sigma_{11}^l)^2} \quad (31)$$

joining equations 29, 30, 31 and setting the derivative equal to zero:

$$\frac{\partial \log(P(\mathbf{X}, \mathbf{S}))}{\partial \sigma_{11}^l} = \sum_i^N -0.5\gamma(z_{il}) \left[\frac{1}{\sigma_{11}^l} - \frac{1}{(\sigma_{11}^l)^2} (x_{i1} - u_{l1})^2 \right] = 0 \quad (32)$$

multiplying by σ_{11}^l :

$$\sum_i^N \gamma(z_{il}) \left[1 - \frac{1}{\sigma_{11}^l} (x_{i1} - u_{l1})^2 \right] = 0 \quad (33)$$

and we can conclude that:

$$\sigma_{11}^l = \frac{\sum_i^N \gamma(z_{il}) (x_{i1} - u_{l1})^2}{\sum_i^N \gamma(z_{il})} \quad (34)$$

following the same steps we can also find that:

$$\sigma_{22}^l = \frac{\sum_i^N \gamma(z_{il}) (x_{i2} - u_{l2})^2}{\sum_i^N \gamma(z_{il})} \quad (35)$$

Let's proceed now with the derivative with respect to λ_l

$$\begin{aligned} \frac{\partial \log(P(\mathbf{X}, \mathbf{S}))}{\partial \lambda_l} &= \sum_i^N \gamma(z_{il}) \frac{\partial \log(P(s_i | Z_l = 1, \lambda_l))}{\partial \lambda_l} \\ &= \sum_i^N \gamma(z_{il}) \left[\frac{\partial \log(\lambda_l^{s_i}) + \log(e^{-\lambda_l}) - \log(s_i!)}{\partial \lambda_l} \right] \\ &= \sum_i^N \gamma(z_{il}) \left[\frac{\partial s_i \log(\lambda_l) - \lambda_l - \log(s_i!)}{\partial \lambda_l} \right] \\ &= \sum_i^N \gamma(z_{il}) \left[\frac{s_i}{\lambda_l} - 1 \right] \end{aligned} \quad (36)$$

setting it equal to zero:

$$\sum_i^N \gamma(z_{il}) \left[\frac{s_i}{\lambda_l} - 1 \right] = 0 \quad (37)$$

$$\lambda_l = \frac{\sum_i^N \gamma(z_{il}) s_i}{\sum_i^N \gamma(z_{il})} \quad (38)$$

Lastly, we compute the derivative with respect to π_l ; here, we have to take in consideration that $\sum_j^K \pi_j = 1$; we will set this constraint using a Lagrange multiplier β :

$$\begin{aligned}\frac{\partial \log(P(\mathbf{X}, \mathbf{S}))}{\partial \pi_l} &= \sum_i^N \gamma(z_{il}) \frac{\partial \log(P(Z_l = 1 | \pi_l)) + \beta(\sum_j^K \pi_j - 1)}{\partial \pi_l} \\ &= \sum_i^N \gamma(z_{il}) \frac{1}{\pi_l} + \beta\end{aligned}\quad (39)$$

$$\sum_i^N \gamma(z_{il}) \frac{1}{\pi_l} + \beta = 0 \quad (40)$$

$$\pi_l = -\frac{\sum_i^N \gamma(z_{il})}{\beta} \quad (41)$$

using the constraint we can calculate the value of the multiplier:

$$\sum_j^K \pi_j = 1 = \sum_j^K -\frac{\sum_i^N \gamma(z_{ij})}{\beta} \quad (42)$$

$$\beta = -\sum_j^K \sum_i^N \gamma(z_{ij}) = -N \quad (43)$$

substituting equation 43 in 41:

$$\pi_l = \frac{\sum_i^N \gamma(z_{il})}{N} \quad (44)$$

Implementation and results

To implement the algorithm, we first initialized the vector $\vec{\pi}$ uniformly, moreover, we randomly initialized the parameters λ_l and $\vec{\mu}_l$; lastly, we set the diagonal elements of the covariance matrices equal to the variances of the data in the two axis.

Plotting the data, we clearly observe three clusters; therefore, we would expect to get a good result running the algorithm with the parameter $k = 3$. Figure 7 and figure 8 show the position of the distributions at the first and at the last iteration respectively. As we can observe we got a very good result. However, we have to keep in mind that the EM algorithm converges to a local maximum, and therefore, depends on the initial positions of the distributions. Figure 9 illustrates the final position of the distributions if we start with a different initial condition. In this case, we did not get a good result. In general for a gaussian mixture, a good idea is to initialize the means of the gaussians equal to the centroid returned by the k-means algorithm. Moreover, we could run the algorithm many times any keep the solution that maximizes the log likelihood.

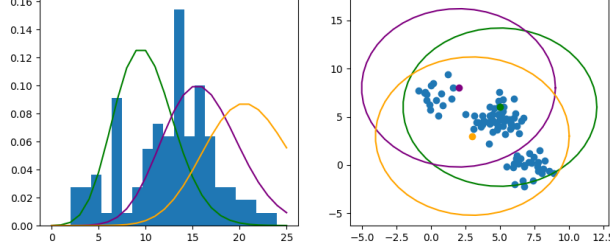


Figure 7: On the left the strengths of the earthquakes plotted using an histogram and the Poisson mixture; on the right, the locations of the earthquakes and the gaussian mixture, the ellipses represent the confidence intervals (3 times the standard deviation). Iteration number = 1; $k = 3$.

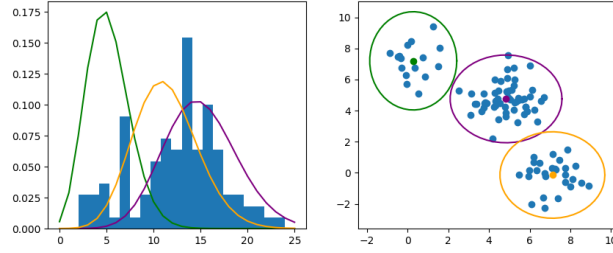


Figure 8: On the left the strengths of the earthquakes plotted using an histogram and the Poisson mixture; on the right, the locations of the earthquakes and the gaussian mixture, the ellipses represent the confidence intervals (3 times the standard deviation). Iteration number = 12; $k = 3$.

In figure 10 we plot the final position of the distributions when $k = 2$. As expected, one of the gaussians is covering two clusters.

Let's see now what happens if we try to set $k = 5$. In figure 11 we can notice a very interesting result; two of the Gaussians have a very low variance and they are concentrated over a very small number of points. In general, we will observe this kind of behavior when we use more gaussians than we need. This happens because when the mean of a gaussian is placed very close to a data point, the log likelihood becomes indefinitely high as the variance goes to zero. This is a very easy way to get a large likelihood, however, it does not mean that we got a good result. For instance, if we have a number of gaussians equal to the number of data points, we will observe that each gaussian will collapse on one of the data points and we would get a very high likelihood, this result however, is not useful. Moreover, this behavior leads to many complications in the implementation of the algorithm. First, if the mean of a gaussian is placed on a data point, we will

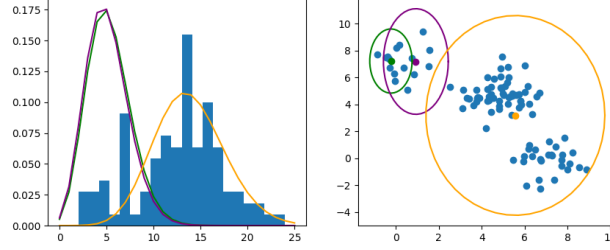


Figure 9: On the left the strengths of the earthquakes plotted using an histogram and the Poisson mixture; on the right, the locations of the earthquakes and the gaussian mixture, the ellipses represent the confidence intervals (3 times the standard deviation). Iteration number = 15; $k = 3$.

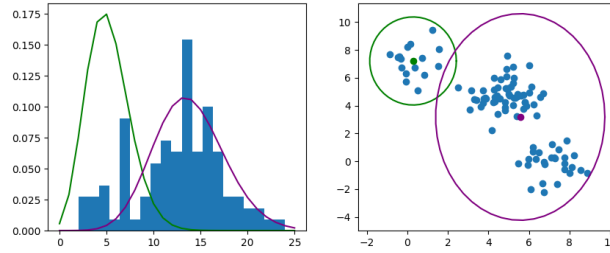


Figure 10: On the left the strengths of the earthquakes plotted using an histogram and the Poisson mixture; on the right, the locations of the earthquakes and the gaussian mixture, the ellipses represent the confidence intervals (3 times the standard deviation). Iteration number = 8; $k = 2$.

get a $\mathbf{0}$ covariance matrix, and since we have the determinant of this matrix in the denominator of the gaussian pdf, we cannot compute the likelihood. Second, even if the mean is not in the same position of a data point, when the variances of the gaussian are very low, Python will approximate the likelihood that a data point was generated by that gaussian to zero when such data point is far from the mean. This approximation does not allow us to compute the log likelihood to check for the convergence of the algorithm. In fact, in the formula of the log likelihood of the data we will have a $\log(0)$. To solve these problems we can check if the covariance matrix becomes equal to the zero matrix, and in such case, we can reinitialize the mean and the covariance matrix to some new values. This in general could be done whenever the variances get very low, checking if the likelihood of the points are approximated to zero. What I observed however, is that in most of cases even if we reinitialize a gaussian, it will tend to collapse on the same data points again. We have to remember, that in our problem

we are fitting some gaussians and some poissons at the same time, therefore, we could observe different behaviours if we work just with gaussians. Notice that we never observe a rotation of the ellipses, this is because we are using a diagonal covariance matrix.

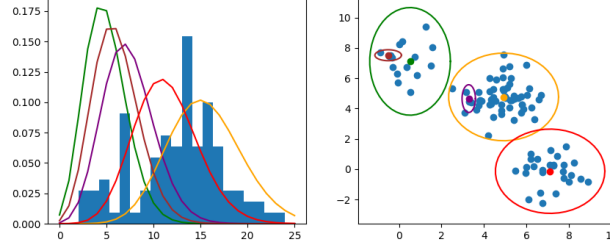


Figure 11: On the left the strengths of the earthquakes plotted using an histogram and the Poisson mixture; on the right, the locations of the earthquakes and the gaussian mixture, the ellipses represent the confidence intervals (3 times the standard deviation). Iteration number = 8; $k = 5$.

Figure 12 illustrates the log likelihood at each iteration for different values of k . As expected, we get the maximum value when $k = 3$. We know that the EM algorithm increases the log likelihood at each step, however, in the plot we can observe that sometimes it decreases, this due to the reinitialization of the gaussian; every time a gaussian is reinitialized, the log likelihood slightly decreases, moreover, after the reinitialization we never observe a significant improvement in the log likelihood.

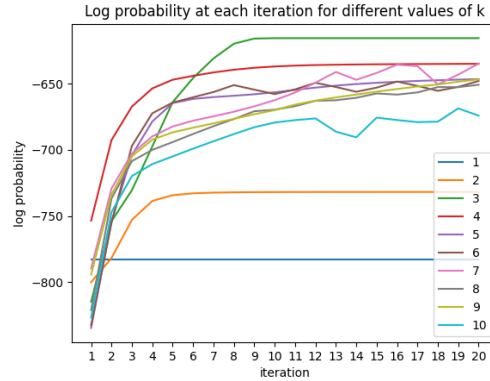


Figure 12: The log likelihood for different values of K at each iteration. The likelihood decreases when a gaussian is reinitialized.

It is also worth to mention that the problem of the collapsing gaussian could be

solved using a bayesian approach; i.e, we consider a prior distribution for the parameters where we have a low probability to get low variances. Using MAP rather than MLE, it is less likely that we will observe the phenomenon of the collapsing gaussian.

To conclude, we provide the values of the π vector when $k = 3$:

0.2700935758292452

0.16000900078909946

0.5698974233816554

We can observe that they sum up to one, and one has a significantly higher value .

References

- [1] R. C. Steorts and L. Qian, “Introduction to the normal gamma model.”
- [2] M. P. J. CAMILLERI, “Variational inference for bayesian univariate gaussian.”

Dynamic algorithm

```

1  import numpy as np
2  from Tree import Tree
3
4  map = {}
5
6  def S(node_u, value_k, tree_topology, theta, beta, values = [0,1,2,3,4]):
7      # print(k)
8      children = []
9      for i in range(len(tree_topology)):
10         if tree_topology[i] == node_u:
11             children.append(i)
12     if not children:
13         return -1
14     product = 1
15     for child in children:
16         sum = 0
17         for value in values:
18             key = str(child)+" "+str(value)
19             s = map.get(key)
20             if s is None:
21                 s = S(child, value, tree_topology, theta, beta)
22                 map[key] = s
23             if s == -1:
24
25                 sum += theta[child][value_k][int(beta[child])]
26                 break
27             else:
28                 sum += s*theta[child][value_k][value]
29     product *= sum
30     return product
31
32
33 def calculate_likelihood(tree_topology, theta, beta):
34     """
35     This function calculates the likelihood of a sample of leaves.
36     :param: tree_topology: A tree topology. Type: numpy array. Dimensions: (num_nodes, )
37     :param: theta: CPD of the tree. Type: list of numpy arrays. Dimensions (approximately):
38     :param: beta: A list of node assignments. Type: numpy array. Dimensions: (num_nodes, )
39             Note: Inner nodes are assigned to np.nan. The leaves have values in [K]
40     :return: likelihood: The likelihood of beta. Type: float.
41     This is a suggested template. You don't have to use it.
42     """
43     sum = 0
44     values = [0, 1, 2, 3, 4]

```

```

45     for value in values:
46         sum += S(0,value,tree_topology,theta,beta)*theta[0][value]
47
48     # Start: Example Code Segment. Delete this segment completely before you implement the
49     print("\tCalculating the likelihood...")
50     likelihood = sum
51     # End: Example Code Segment
52     for k in list(map):
53         del map[k]
54     return likelihood
55
56
57 def main():
58     print("Hello World!")
59     print("This file is the solution template for question 2.2.")
60
61     print("\n1. Load tree data from file and print it\n")
62
63     filename = "data/q2_2_small_tree.pkl" # "data/q2_2_medium_tree.pkl", "data/q2_2_large_
64     print("filename: ", filename)
65
66     t = Tree()
67     t.load_tree(filename)
68     t.print()
69     print("K of the tree: ", t.k, "\talphabet: ", np.arange(t.k))
70
71     print("\n2. Calculate likelihood of each FILTERED sample\n")
72     # These filtered samples already available in the tree object.
73     # Alternatively, if you want, you can load them from corresponding .txt or .npy files
74
75     for sample_idx in range(t.num_samples):
76         beta = t.filtered_samples[sample_idx]
77         print("\n\tSample: ", sample_idx, "\tBeta: ", beta)
78         sample_likelihood = calculate_likelihood(t.get_topology_array(), t.get_theta_array())
79         print("\tLikelihood: ", sample_likelihood)
80
81
82 if __name__ == "__main__":
83     main()

```

Variational Inference algorithm

```
1 from matplotlib import pyplot as plt
2 from scipy.stats import norm, gamma
3 import numpy as np
4
5 tolerance = 1e-4
6 #parameters of the normal that generates the points
7 mean = 9
8 precision = 1
9 #parameters of the prior distribution
10 #gamma
11 a_0 = 1
12 b_0 = 0.5
13 #gaussian
14 mu_0 = 3
15 lambda_0 = 4
16 #number of points
17 N_points = 5
18 #compute the prior for some values of the mean and the precision
19 mean_axis = np.linspace(2, 10, 400)
20 tau_axis = np.linspace(0, 3, 400)
21
22 prior_values = np.empty([len(mean_axis), len(tau_axis)])
23 p_tau = gamma(a=a_0, loc=0, scale=1/b_0)
24 for i, tau in enumerate(tau_axis[1:]):
25     p_mu = norm(loc=mu_0, scale=1/np.sqrt(lambda_0 * tau))
26     prior_values[i,:] = p_tau.pdf(tau) * p_mu.pdf(mean_axis)
27
28 # generate points
29 data = np.random.normal(loc=mean, scale=np.sqrt(1 / precision), size=N_points)
30
31 #compute exact posterior
32 #it is interesting to note that the parameters depend on the number of points
33 mu_N = (lambda_0 * mu_0 + np.sum(data))/(N_points + lambda_0)
34 a_N = a_0 + N_points*0.5
35 b_N = b_0 + 0.5 *
36     lambda_0 * (mu_0**2) + 0.5 * np.square(data).sum() - 0.5
37     * (N_points + lambda_0) * mu_N**2
38
39 exact_posterior_values = np.empty([len(mean_axis), len(tau_axis)])
40 p_tau = gamma(a=a_N, loc=0, scale=1/b_N)
41 for i, tau in enumerate(tau_axis[1:]):
42     p_mu = norm(loc=mu_N, scale=1/np.sqrt((N_points+lambda_0) * tau))
43     exact_posterior_values[i,:] = p_mu.pdf(mean_axis) * p_tau.pdf(tau)
44
```

```

45
46 E_tau = 10
47 E_tau_prev = np.NaN
48
49 iteration = 0
50 while np.isnan(E_tau_prev) or (np.abs(E_tau - E_tau_prev) >
51 tolerance and iteration < 10):
52     mu_N = (lambda_0 * mu_0 + N_points * data.mean()) / (lambda_0 + N_points)
53     lambda_N = (lambda_0 + N_points) * E_tau
54     E_mu = mu_N
55     E_mu2 = lambda_N ** -1 + E_mu ** 2
56     a_N = a_0 + (N_points + 1) / 2
57     b_N = b_0 + 0.5 * (np.square(data) - 2
58         * E_mu * data + E_mu2).sum() + 0.5 * lambda_0 *
59         (E_mu2 - 2 * E_mu * mu_0 + mu_0 ** 2)
60     E_tau_prev = E_tau
61     E_tau = a_N / b_N
62
63     pTau = gamma(a=a_N, loc=0, scale=1 / b_N)
64     pMu = norm(loc=mu_N, scale=1 / np.sqrt(lambda_N))
65
66     values = np.outer(pMu.pdf(mean_axis), pTau.pdf(tau_axis))
67
68     # Now Plot
69     plt.figure(figsize=[10, 10])
70     plt.plot(mean, precision, 'o', label='parameters of the gaussian')
71     plt.plot(data.mean(), 1 / data.var(), 'o', label='data mean and variance')
72     plt.contour(*np.meshgrid(mean_axis, tau_axis), exact_posterior_values,
73         levels=8, colors=['blue'])
74     plt.contour(*np.meshgrid(mean_axis, tau_axis), prior_values, levels=8, colors=['green'])
75     plt.contour(*np.meshgrid(mean_axis, tau_axis), values.T, levels=8, colors=['red'])
76
77     plt.title('Iteration {}'.format(iteration))
78     plt.xlabel('Mean')
79     plt.ylabel('Precision')
80     plt.legend(loc="upper right")
81     plt.show()
82     iteration += 1

```

EM algorithm

```
1 import numpy as np, numpy.random
2 from scipy.stats import multivariate_normal as normal
3 from scipy.stats import poisson
4 import matplotlib.pyplot as plt
5 import random
6 import matplotlib as mpl
7 from matplotlib.patches import Ellipse
8
9 def column(matrix, i):
10     return [row[i] for row in matrix]
11
12
13 def EM (S,X, n, k):
14
15     X0 = column(X, 0)
16     X1 = column(X, 1)
17     tolerance = 0.001
18     prob = np.zeros(n)
19     pi = [1/k for i in range(k)]
20     responsibility = np.zeros((k,len(S)))
21     lambda_ = [random.randrange(int(np.min(S)),int(np.max(S))) for i in range(k)]
22     mu = [[random.randrange(int(np.min(X0)),int(np.max(X0))),
23            random.randrange(int(np.min(X1)),int(np.max(X1)))] for i in range(k)]
24     sigma = np.ones((k,2,2))
25     for i in range(k):
26         sigma[i] = [[np.var(X0),0],[0,np.var(X1)]]
27
28
29
30     fig = plt.figure(figsize=(10,4))
31     ax = fig.add_subplot(121)
32     ax2 = fig.add_subplot(122)
33     #plot poisson
34     colors = ["green","purple","orange","brown","red","yellow","pink"]
35     ax.hist(S, bins=20, density=True)
36     for j in range(k):
37         x = [i for i in range(0, 26)]
38         y = [poisson.pmf(i, lambda_[j]) for i in range(0, 26)]
39         ax.plot(x, y, color=colors[j])
40     #plot gaussians
41     ax2.scatter(X0, X1)
42     for i in range(k):
43         ax2.scatter(mu[i][0], mu[i][1], color=colors[i])
44         u = mu[i][0] # x-position of the center
```



```

45     v = mu[i][1] # y-position of the center
46     a = np.sqrt(sigma[i][0][0]) * 3 # radius on the x-axis
47     b = np.sqrt(sigma[i][1][1]) * 3 # radius on the y-axis
48     t = np.linspace(0, 2 * np.pi, 50)
49     ellipse_x = [u + a * np.cos(j) for j in t]
50     ellipse_y = [v + b * np.sin(j) for j in t]
51     ax2.plot(ellipse_x, ellipse_y, '-', color=colors[i])
52 plt.savefig('images/test6/{}.png'.format(0))
53 plt.show()
54
55
56
57 iteration = 0
58 while iteration < n and (iteration < 2 or
59     np.abs(prob[iteration-1]-prob[iteration-2]) >= tolerance):
60     #expectation
61     for i in range(k):
62         for j in range(len(X)):
63             responsibility[i][j] = pi[i]*normal.pdf(X[j],mean=mu[i], cov=sigma[i] )
64             *poisson.pmf(S[j], mu=lambda_[i])
65             sum = 0
66             for l in range(k):
67                 sum +=pi[l]*normal.pdf(X[j],mean=mu[l],cov=sigma[l])*
68                 poisson.pmf(S[j],lambda_[l])
69             responsibility[i][j] /= sum
70     #maximization
71     for i in range(k):
72         tot_responsibility = np.sum(responsibility[i])
73         mu[i] = np.zeros(2)
74         sigma[i] = np.zeros((2,2))
75         lambda_[i] = 0
76         for j in range(len(X)):
77             mu[i] += X[j]*responsibility[i][j]
78         mu[i] /= tot_responsibility
79         for j in range(len(X)):
80             lambda_[i] += S[j]*responsibility[i][j]
81             sigma[i][0][0] += responsibility[i][j]* (X[j][0] - mu[i][0])**2
82             sigma[i][1][1] += responsibility[i][j]* (X[j][1] - mu[i][1])**2
83         sigma[i] /= tot_responsibility
84         lambda_[i] /= tot_responsibility
85         pi[i] = tot_responsibility/len(X)
86
87     #check if the matrix is semidefinite positive
88     if np.linalg.det(sigma[i]) == 0:
89         sigma[i] = [[np.var(X0), 0], [0, np.var(X1)]]
90         mu[i] = [random.randrange(int(np.min(X0)),

```

```

91         int(np.max(X0))), random.randrange(int(np.min(X1)), int(np.max(X1)))])
92         # sigma[i] += 1.0e-8*numpy.identity(2)
93
94     #compute log probability
95     for j in range(len(S)):
96         for i in range(k):
97             poisson_p = poisson.pmf(S[j], lambda_[i])
98             normal_p = normal.pdf(X[j], mean=mu[i], cov=sigma[i])
99             if normal_p == 0:
100                 sigma[i] = [[np.var(X0), 0], [0, np.var(X1)]]
101                 mu[i] = [random.randrange(int(np.min(X0)),
102                                     int(np.max(X0))),
103                           random.randrange(int(np.min(X1)),
104                                     int(np.max(X1)))]
105                 normal_p = normal.pdf(X[j], mean=mu[i], cov=sigma[i])
106             prob[iteration] += responsibility[i][j] * (np.log(normal_p)
107                                                         + np.log(poisson_p) + np.log(pi[i]))
108
109     iteration += 1
110
111     fig = plt.figure(figsize=(10, 4))
112     ax = fig.add_subplot(121)
113     ax2 = fig.add_subplot(122)
114     ax.hist(S, bins=20, density=True)
115     for j in range(k):
116         x = [i for i in range(0, 26)]
117         y = [poisson.pmf(i, lambda_[j]) for i in range(0, 26)]
118         ax.plot(x, y, color=colors[j])
119
120     # plot gaussians
121     ax2.scatter(X0, X1)
122     for i in range(k):
123         ax2.scatter(mu[i][0], mu[i][1], color=colors[i])
124         u = mu[i][0] # x-position of the center
125         v = mu[i][1] # y-position of the center
126         a = np.sqrt(sigma[i][0][0]) * 3 # radius on the x-axis
127         b = np.sqrt(sigma[i][1][1]) * 3 # radius on the y-axis
128         t = np.linspace(0, 2 * np.pi, 50)
129         ellipse_x = [u + a * np.cos(j) for j in t]
130         ellipse_y = [v + b * np.sin(j) for j in t]
131         ax2.plot(ellipse_x, ellipse_y, '-', color=colors[i])
132
133     plt.savefig('images/test6/{}.png'.format(iteration))
134     plt.show()
135
136     #plot logs
137     # x = [i for i in range(1,n+1)]
138     # plt.xticks(x)
139     # plt.plot(x, prob[:,], label=(k))

```

```
137         return pi, lambda_, mu, sigma, np.max(prob)
138
139     X = np.loadtxt("X.txt")
140     S = np.loadtxt("S.txt")
141
142
143     pi, lambda_, mu, sigma, prob = EM(S, X, 50, 2)
```

I discussed the assignment with my teammate Alice De Schutter.