# Assignment 2 - DD2434 Machine Learning, Advanced Course

August Regnell 970712-9491

19 December, 2020

Code found in Appendix and in the following **Git repo**

## Contents

# 1 Dependencies in a Directed Graphical Model

## 1.1 Question 2.1.1

No.

## 1.2 Question 2.1.2

Yes.

## 1.3 Question 2.1.3

$A = \{\mu_{r,c} \mid r \in [R], c \in [C]\}$

## 1.4 Question 2.1.4

No.

## 1.5 Question 2.1.5

No.

## 1.6 Question 2.1.6
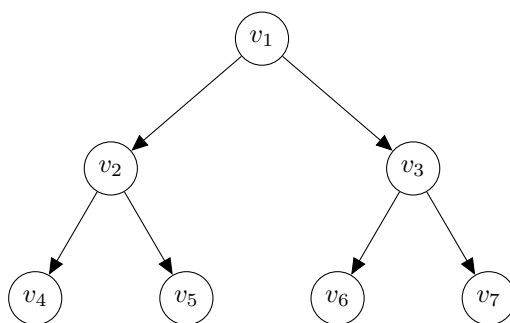
$B = \{Z_m^n \mid n \in [N], m \in [M]\} \cup \{C^n \mid n \in [N]\}$

# 2   Likelihood of a Tree Graphical Model

## 2.1   Question 2.2.7

Let $T$ be a binary tree, with vertex set $V(T)$ and leaf set $L(T)$, and consider the graphical model $T, \Theta$ described as follows. For each vertex $v \in V(T)$ there is an associated random variable $X_v$ that assumes values in $[K]$. Moreover, for each $v \in V(T)$, the CPD $\theta_v = p(X_v \mid x_{pa(v)})$ is a categorical distribution. Let $\beta = \{x_l : l \in L(T)\}$ be an assignment of values to all the leaves of T.

   Consider the likelihood of the following small tree consisting of seven vertices $\{v_1, \ldots, v_7\}$ numbered from top to bottom and left to right. (**Note:** the trees provided in the data are not full trees like this example)



Its likelihood is given by:

$$p(T \mid \Theta) = p(v_1)p(v_2 \mid v_1)p(v_3 \mid v_1)p(v_4 \mid v_2)p(v_5 \mid v_2)p(v_6 \mid v_3)p(v_6 \mid v_3)$$

For this tree we have that $\beta = \{x_l : l \in L(T)\} = \{v_4, v_5, v_6, v_7\}$, that is:

Calculating $p(\beta \mid T, \Theta)$ can be done by marginalizing on $V(T) \setminus L(T)$.

$$p(\beta \mid T, \Theta)$$

$$= \sum_{v \in V(T) \setminus L(T)} p(x_{v_1}) p(x_{v_2} \mid x_{v_1}) p(x_{v_3} \mid x_{v_1}) p(X_{v_1} \mid x_{v_2}) p(X_{v_1} \mid x_{v_2}) p(X_{v_6} \mid x_{v_3}) p(X_{v_7} \mid x_{v_3})$$

$$= \sum_{x_1} \sum_{x_2} \sum_{x_3} p(x_{v_1}) p(x_{v_2} \mid x_{v_1}) p(x_{v_3} \mid x_{v_1}) p(X_{v_1} \mid x_{v_2}) p(X_{v_1} \mid x_{v_2}) p(X_{v_6} \mid x_{v_3}) p(X_{v_7} \mid x_{v_3})$$

Generalizing for a larger tree gives:

$$p(\beta \mid T, \Theta) = \sum_{v \in V(T) \setminus L(T)} p(x_{v_1}) p(x_{v_2} \mid x_{v_1}) p(x_{v_3} \mid x_{v_1}) \cdots \prod_{v \in L(T)} p(X_v \mid x_{pa(v)})$$

$$= \sum_{v_1} \sum_{v_2} \cdots \sum_{v_N} p(x_{v_1}) p(x_{v_2} \mid x_{v_1}) p(x_{v_3} \mid x_{v_1}) \cdots \prod_{v \in L(T)} p(X_v \mid x_{pa(v)})$$

Calculating the likelihood this way is extremely inefficient. It would in fact be $\mathcal{O}(K^N)$ where $N = \mid V(T) \setminus L(T) \mid$. To make the calculation more efficient we can use dynamic programming and conditional independence to bring the sum inside the product and break the problem into smaller sub-problems.

$$p(\beta \mid T, \Theta) = \sum_{x_{v_1}} p(x_{v_1}) \mu(x_{v_1})$$

where

$$\mu(x_v) = \prod_{c \in de(v)} \sum_{x_c} p(X_c = x_c \mid X_v = x_v) \mu(x_c)$$

where $de(v)$ denotes the descendants of $v$ and $\mu(x_v) = \prod_{c \in de(v)} p(X_c = x_c^* \mid X_v = x_v)$ if $de(v) \subseteq \beta$. The $*$ denotes that its value has been observed.

This method is much more efficient than the previous one. It is in fact $\mathcal{O}(NK^2)$. It has been implemented in matrix form in python.

## 2.2   Question 2.2.8

This implementation was run on the data provided. The results on the different tree sizes and their five respective samples are shown in table 1.

| Tree size | 1 | 2 | 3 | 4 | 5 |
|:---:|:---:|:---:|:---:|:---:|:---:|
| Small | 0.016179 | 0.01541 | 0.01137 | 0.0086 | 0.0409 |
| Medium | $4.336 \cdot 10^{-18}$ | $3.094 \cdot 10^{-20}$ | $1.050 \cdot 10^{-16}$ | $6.585 \cdot 10^{-16}$ | $1.488 \cdot 10^{-18}$ |
| Large | $3.288 \cdot 10^{-69}$ | $1.109 \cdot 10^{-66}$ | $2.522 \cdot 10^{-68}$ | $1.242 \cdot 10^{-66}$ | $3.535 \cdot 10^{-69}$ |

Table 1: $p(\beta \mid T, \Theta)$ for the three different tree sizes and their five respective samples.

To ponder the plausibility of these results one can do some rough calculations. Every vertex can take five different values, so if they would be equally likely they would have a probability of 0.2. Using this approximation one can use the respective trees number of leaves (3, 26 and 101) to get a rough approximation of the likelihood. These would be:

$$\left(\frac{1}{5}\right)^3 = 0.008$$

$$\left(\frac{1}{5}\right)^{26} \approx 6.71 \cdot 10^{-19}$$

$$\left(\frac{1}{5}\right)^{101} \approx 2.54 \cdot 10^{-71}$$

We see that these approximations are close to the results.

# 3   Simple Variational Inference

## 3.1   Question 2.3.9

The VI algorithm was implemented using python.

## 3.2   Question 2.3.10

We are given the following likelihood function for the data

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

and the following conjugate priors distribution for $\mu$ and $\tau$

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

$$p(\tau) = \mathrm{Gam}(\tau \mid a_0, b_0)$$

We want to find the exact posterior, $p(\mu, \tau \mid \mathcal{D})$.

$$p(\mu, \tau \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mu, \tau)p(\mu, \tau)}{p(\mathcal{D})} \propto p(\mathcal{D} \mid \mu, \tau)p(\mu, \tau) = p(\mathcal{D} \mid \mu, \tau)p(\mu \mid \tau)p(\tau)$$

Taking the logarithm of both sides we get

$$\log(p(\mu, \tau \mid \mathcal{D})) \overset{+}{=} \log(p(\mathcal{D} \mid \mu, \tau)) + \log(p(\mu \mid \tau)) + \log(p(\tau))$$

$$\overset{+}{=} \frac{N}{2}\log(\tau) - \frac{\tau}{2}\sum_{n=1}^{N}(x_n - \mu)^2 + \frac{1}{2}\log(\tau) - \frac{\lambda_0 \tau}{2}(\mu - \mu_0)^2 + (a_0 - 1)\log(\tau) - b_0\tau$$

$$= (a_0 + \frac{N}{2} - 1)\log(\tau) - b_0\tau + \frac{1}{2}\log(\tau) - \frac{\tau}{2}\left(N\mu^2 - 2N\mu\bar{x} + \sum_{n=1}^{N}x_n^2\right) - \frac{\lambda_0 \tau}{2}(\mu^2 - 2\mu\mu_0 + \mu_0^2)$$

$$= (a_0 + \frac{N}{2} - 1)\log(\tau) - (b_0 + \frac{1}{2}\lambda_0\mu_0^2 + \frac{1}{2}\sum_{n=1}^{N}x_n^2)\tau + \frac{1}{2}\log(\tau) - \frac{\tau}{2}\left((N + \lambda_0)\mu^2 - 2\mu(N\bar{x} + \lambda_0\mu_0)\right)$$

$$= \left(a_0 + \frac{N}{2} - 1\right)\log(\tau) - \left(b_0 + \frac{1}{2}(\lambda_0\mu_0^2 - \frac{(N\bar{x} + \lambda_0\mu_0)^2}{N + \lambda_0} + \sum_{n=1}^{N}x_n^2)\right)\tau + \frac{1}{2}\log(\tau) - \frac{\tau(N + \lambda_0)}{2}\left(\mu - \frac{N\bar{x} + \lambda_0\mu_0}{N + \lambda_0}\right)^2$$

$$= (\hat{a} - 1)\log(\tau) - \hat{b}\tau + \frac{1}{2}\log(\tau) - \frac{\tau\hat{\beta}}{2}(\mu - \hat{\mu})^2$$

Notice that from the third to second last expression the square was completed for the last term, leading to adding $\frac{\tau}{2}\frac{(N\bar{x}+\lambda_0\mu_0)^2}{N+\lambda_0}$.

From the last two expressions we see that the exact posterior is *normal-gamma* distributed. That is

$$p(\mu,\tau) = \mathcal{N}(\mu \mid \hat{\mu}, (\hat{\beta}\tau)^{-1})\mathrm{Gam}(\lambda \mid \hat{a},\hat{b})$$

where

$$\hat{\beta} = N + \lambda_0$$

$$\hat{\mu} = \frac{N\bar{x} + \lambda_0\mu_0}{N+\lambda_0} = \frac{N\bar{x} + \lambda_0\mu_0}{\hat{\beta}}$$

$$\hat{a} = a_0 + \frac{N}{2}$$

$$\hat{b} = b_0 + \frac{1}{2}\left(\lambda_0\mu_0^2 - \frac{(N\bar{x}+\lambda_0\mu_0)^2}{N+\lambda_0} + \sum_{n=1}^{N}x_n^2\right) = b_0 + \frac{1}{2}\left(\lambda_0\mu_0^2 - \hat{\beta}\hat{\mu}^2 + \sum_{n=1}^{N}x_n^2\right)$$

### 3.3   Question 2.3.11

**Case 1**

The first case is the one used in Figure 10.4 in Bishop. That is $\mu$ and $\tau$ are the parameters for a univariate standard normal distribution. Thus, the data was simulated from a univariate standard normal distribution with 100 data points. Assuming the conjugate priors were relatively well known the parameters shown in table 2 were chosen.

| $a_0$ | $b_0$ | $\mu_0$ | $\lambda_0$ |
|---|---|---|---|
| 10 | 10 | 0 | 10 |

Table 2: Chosen parameters of priors for Case 1

These results can be seen in figure 1 and table 3.

Figure 1: Contour plots of posterior calculated by VI (in red) and by exact posterior (in blue) for Case 1. Horizontal axis $\mu$, vertical axis $\tau$. **Left:** Zoomed out. **Right:** Zoomed in.

| Method | $a$ | $b$ | $\mu$ | $\lambda/\beta$ |
|--------|-----|-----|-------|-----------------|
| VI | 60 | 57.55 | -0.09 | 114.68 |
| Exact | 60 | 57.07 | -0.09 | 110 |

Table 3: Parameters calculated by VI and exact posterior for Case 2.

We see that the exact posterior is much better at capturing the variance of $\mu$ when the precision, $\tau$, is low. This can be seen in the blue contours being wider for smaller $\tau$. However, we see almost no difference in the calculated parameters, which means that the difference is captured in the differing distributions.

Moreover, we also see that the outer contours are wide, which is due to the small amount of data points.

**Case 2**

We will now increase the number of data points to $n = 10,000$ but keep the distribution we simulate the data from the same, as well as the prior parameters (they can thus be seen in table 2).
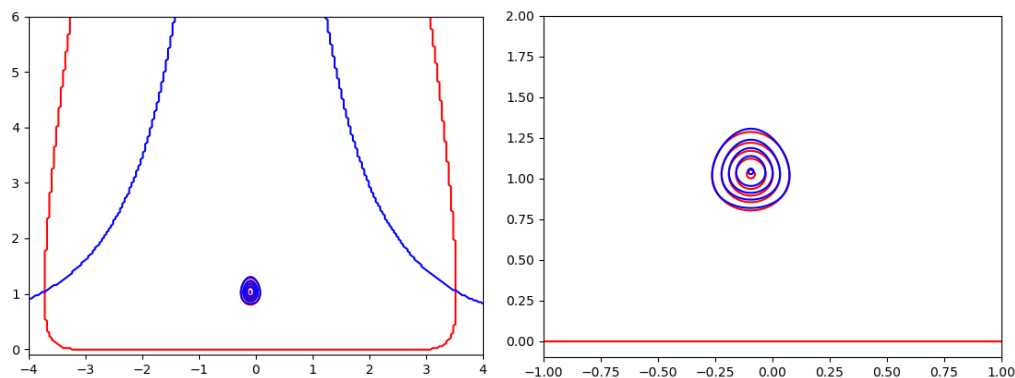
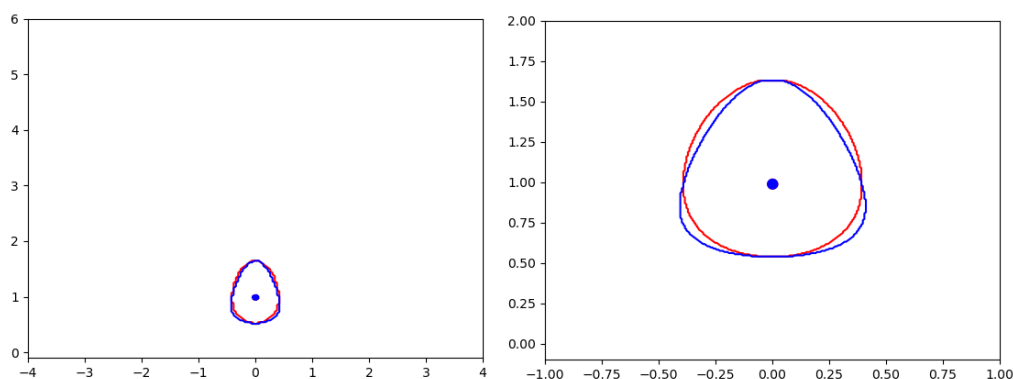The results can be seen in figure 2 and table 4.

Figure 2: Contour plots of posterior calculated by VI (in red) and by exact posterior (in blue) for Case 2. Horizontal axis $\mu$, vertical axis $\tau$. **Left:** Zoomed out. **Right:** Zoomed in.

| Method | $a$ | $b$ | $\mu$ | $\lambda/\beta$ |
|--------|------|---------|--------|---------|
| VI | 5010 | 5073.67 | 0.0011 | 9884.39 |
| Exact | 5010 | 5073.16 | 0.0011 | 10010 |

Table 4: Parameters calculated by VI and exact posterior for Case 2.

We see that the contour plots now become much tighter, indicating that the increased number of data points significantly has increased the accuracy. In the zoomed picture we also see that most of the probability is contained in the neighborhood of the true parameters, $(0, 1)$. We still see the difference in the shape between the exact and the posterior estimated by the VI algorithm.

**Case 3**

For the last case we choose a different scenario, were $\mu = 10$ and $\tau = 0.5$. We simultaneously reduce the number of data points to $n = 1,000$ and make a bad guess for the prior parameters (assuming they are not well known). The chosen parameters can be seen in table 5.

| $a_0$ | $b_0$ | $\mu_0$ | $\lambda_0$ |
|-------|-------|---------|-------------|
| 2 | 0.1 | 5 | 1 |

Table 5: Chosen parameters of priors for Case 3

The results can be seen in figure 3 and table 6.

Figure 3: Contour plots of posterior calculated by VI (in red) and by exact posterior (in blue) for Case 3. Horizontal axis $\mu$, vertical axis $\tau$. **Left:** Zoomed out. **Right:** Zoomed in.

| Method | $a$ | $b$ | $\mu$ | $\lambda/\beta$ |
|--------|-----|--------|-------|---------|
| VI | 502 | 149.42 | 9.987 | 3363.07 |
| Exact | 502 | 149.27 | 9.987 | 1001 |

Table 6: Parameters calculated by VI and exact posterior for Case 3.

Contrary to the other cases, in Case 3 neither of the methods are able to get a good estimate of the true values of $\tau$ (they still do well for $\mu$). The contour plots should be centered around $(10, 1/0.5^2) = (10, 4)$ if they had estimated the parameters correctly. This is due to the bad choice of prior parameters (a better choice of prior parameters gives the distribution shown in figure 4). This choice would be equal to believing $\mu$ is centered around 5 and $\tau$ around $2/0.1 = 20$ (with a variance $\text{Var}(\tau) = 2/0.1^2 = 200$). This shows us that these methods have a hard time when little is known about the prior distributions.

We still see the familiar difference in shape in the zoomed out figure, but the contours start to coincide in the zoomed in figure.

Figure 4: Case 3 but with a better choice of prior parameters.

# 4   Mixture of trees with observable variables

## 4.1   Question 2.4.12

The EM algorithm was implemented using python. Sieving was applied in the following way:

1. Choose 100 random seeds and simulate one tree from every seed

2. Run the algorithm for 10 iterations for these 100 trees and choose the 10 seeds that gave the highest likelihood

3. Re-simulate the trees from the 10 best seeds and run the algorithm for 100 iterations

4. Choose the best seed and re-run it to convergence (relative error $< 10^{-6}$) or a maximum of 100 iterations

## 4.2   Question 2.4.13

The algorithm was run on the provided data. The graphs of the log- and regular likelihood versus the iteration of the algorithm can be seen in figure 5.



Figure 5: The likelihood of the data given the current (in terms of iteration) tree mixture.

Further, the inferred trees were also compared with the real trees using the unweighted Robinson-Foulds metric and the likelihood. The results can be seen in table 7 and 8 respectively.

|  | Real Tree 0 | Real Tree 1 | Real Tree 2 |
|---|---|---|---|
| **Inferred Tree 0** | 4 | 5 | 4 |
| **Inferred Tree 1** | 0 | 3 | 4 |
| **Inferred Tree 2** | 2 | 3 | 4 |

Table 7: Robin-Foulds metric between the inferred and real trees.

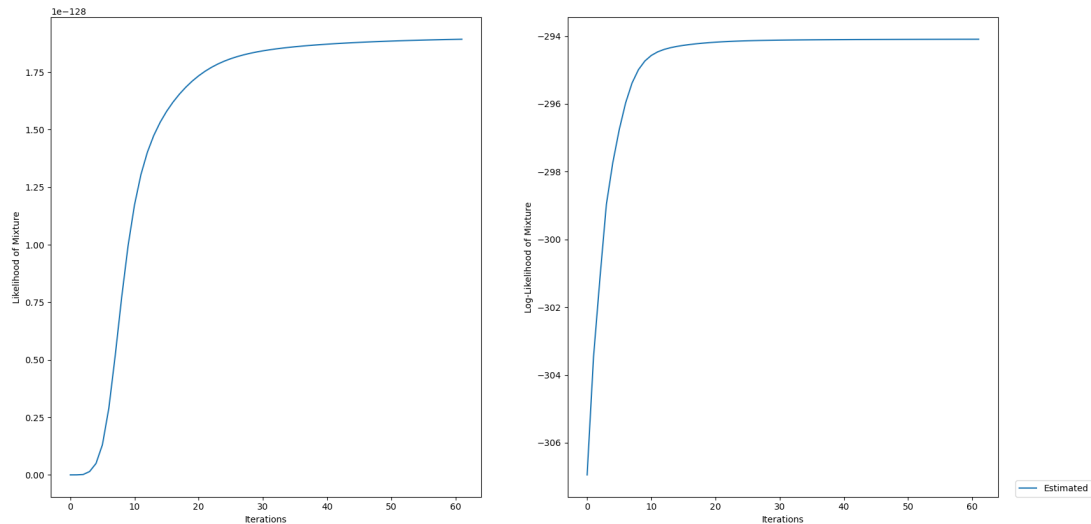|  | Log-likelihood | Likelihood |
|---|---|---|
| **EM** | -294.093 | 1.89298e-128 |
| **Real** | -311.449 | 5.48542e-136 |

Table 8: Values of likelihood for the inferred and real tree mixture.

From table 7 we see that the inferred and real trees are not so similar in terms of structure, with an exception of inferred tree 1 and the real tree 0 which have distance of 0. Note that a value of 5 is large for these trees since the only contain five nodes. Matching the inferred and real trees to get the minimal total distance we get a total distance of 7. Furthermore, looking at the likelihood of the samples given a certain tree mixture we see that the EM algorithm produces a much higher likelihood than the real tree mixture. This is probably due to the low amount of samples, which could lead to overfitting.

### 4.3   Question 2.4.14

**Scenario 1**

As a first test, we use the same tree mixture as the previous question but increase the number of samples to $n = 1000$. This, as expected, made the algorithm take much longer to converge. The results can be seen in figure 6. We see that the log-likelihood graph has a similar shape to the previous example with $n = 100$. However, we do get a horizontal line for the likelihood due to python not being able to handle so small numbers (around $10^{-1300}$).

Figure 6: The likelihood of the data given the current (in terms of iteration) tree mixture.

|  | Real Tree 0 | Real Tree 1 | Real Tree 2 |
|---|---|---|---|
| **Inferred Tree 0** | 0 | 3 | 4 |
| **Inferred Tree 1** | 4 | 3 | 0 |
| **Inferred Tree 2** | 4 | 5 | 4 |

Table 9: Robin-Foulds metric between the inferred and real trees.

|  | Log-likelihood | Likelihood |
|---|---|---|
| **EM** | -3031.45 | 3.4834e-1316 |
| **Real** | -3052.08 | 3.1732e-1325 |

Table 10: Values of likelihood for the inferred and real tree mixture.

We know see that the inferred trees better matches the real tree in structure. Matching in the same way as the previous question gives a minimal total distance of 5, where the pair real tree 1 and inferred tree 2 is the only contributor. This indicates that the higher number of samples leads to a better matching, more precisely since this resulted in two perfect matches (0 distance).

Furthermore, looking at the difference in likelihood we see that the gap in log-likelihood has increased, from $\approx 17$ to $\approx 21$. However, the number of samples is also ten times more, so this means that the gap has shrunk relatively. This supports the hypothesis that more samples decreased the overfitting.

**Scenario 2**

For the second scenario, the number of samples was once again reduced to $n = 100$. However, the number of nodes were increased to 7 while the number of clusters reduced to 2. What one first notices is that the algorithm converged after only 10 iterations. This could be due to the lower number of clusters.



Figure 7: The likelihood of the data given the current (in terms of iteration) tree mixture.

|  | Real Tree 0 | Real Tree 1 |
|---|---|---|
| **Inferred Tree 0** | 3 | 4 |
| **Inferred Tree 1** | 3 | 4 |

Table 11: Robin-Foulds metric between the inferred and real trees.

|  | Log-likelihood | Likelihood |
|---|---|---|
| **EM** | -410.593 | 4.80378e-179 |
| **Real** | -424.39 | 4.89653e-185 |

Table 12: Values of likelihood for the inferred and real tree mixture.

We first of all see now perfect match in terms of structure between the inferred and real trees. As in question 2.4.13 we see that the minimal total distance is 7. Compared to the same question, this tree mixture has 14 nodes compared to 15. This indicates that the

algorithm does about as good for this tree mixture as the tree mixture with 5 nodes and 2 clusters.

Looking at the likelihoods, we first of all see that they are much lower, even though it is the same number of samples as in 2.4.13. This is due to the likelihood of every node in the respective trees are multiplied with each other and the likelihoods of the different trees are added (after multiplying by their $\pi_k$). We also see that the difference in likelihood between the inferred and real tree mixture is smaller. This is probably due to less overfitting, due to this tree mixture having less variables to optimize ($\pi$ and $2 \cdot 7$ nodes).

**Scenario 3**

For the third and final scenario, the number of samples was once again put at $n = 100$. However, the number of nodes were decreased to 4 while the number of clusters increased to 5. The threshold for convergence was decreased to $10^{-8}$ to not make the algorithm



Figure 8: The likelihood of the data given the current (in terms of iteration) tree mixture.

|        | RT 0 | RT 1 | RT 2 | RT 3 | RT 4 |
|--------|------|------|------|------|------|
| **IT 0** | 2    | 2    | 3    | 2    | 0    |
| **IT 1** | 4    | 4    | 3    | 4    | 2    |
| **IT 2** | 4    | 4    | 3    | 4    | 2    |
| **IT 3** | 2    | 2    | 3    | 2    | 0    |
| **IT 4** | 4    | 4    | 3    | 4    | 2    |

Table 13: Robin-Foulds metric between the inferred (IT) and real trees (RT).

|        | Log-likelihood | Likelihood    |
|--------|----------------|---------------|
| **EM** | -260.678       | 6.15237e-114  |
| **Real** | -267.582     | 6.17358e-117  |

Table 14: Values of likelihood for the inferred and real tree mixture.

At first we see that both IT 0 and IT 3 is a perfect match, structure wise, of RT 4. However, choosing one of these when calculating the minimal distance is not the only solution, there are in fact several combinations without these that will grant the minimal value of 13. Moreover, this total distance is higher to the other scenarios, even in the relative case (13 distance with $4 \cdot 5 = 20$ different nodes). We thus conclude that the algorithm had a harder time with these particular structure of tree mixture. It would probably do better with more samples, as was shown with the first scenario.[1]

We also notice that the likelihoods are quite high compared to the other scenarios, they are in fact the highest. This is due to the same reasoning explained in the previous scenario, namely that the number of nodes per tree has the largest influence on the likelihood (if the number of samples is kept constant).

---

[1]When $n$ was increased to 300 the minimum total distance was reduced to 9

# A   Appendix: Code

## A.1   2.2 - Likelihood of a Tree Graphical Model implementation

```python
import numpy as np
from Tree import Tree
from Tree import Node


def calculate_likelihood(root, beta, k):
    print("Calculating the likelihood...")

    likelihood = recursive(root, beta, k)

    return likelihood


def recursive(node, beta, k):
    if len(node.descendants) == 0: # If the vertex is a leaf
        obs_index = int(beta[int(node.name)])
        theta = np.array(node.cat)
        theta = theta[:,obs_index] # Extract probabilities of the observed vertex given all
        return theta
    else:
        theta = np.array(node.cat)
        descendant_theta = np.ones(k)
        for descendant in node.descendants:
            descendant_theta *= recursive(descendant, beta, k) # Multiply probabilities of
        theta = np.dot(theta, descendant_theta)
        return theta


def main(treeSize):
    print("\nRunning algorithm for " + treeSize + " tree...")
    if treeSize == "small":
        filename = "data/q2_2/q2_2_small_tree.pkl"
    elif treeSize == "medium":
        filename = "data/q2_2/q2_2_medium_tree.pkl"
    else:
        filename = "data/q2_2/q2_2_large_tree.pkl"
```

```python
38        t = Tree()
39        t.load_tree(filename)
40        # t.print()
41
42        for sample_idx in range(t.num_samples):
43            beta = t.filtered_samples[sample_idx]
44            sample_likelihood = calculate_likelihood(t.root, beta, t.k)
45            print("\tLikelihood for sample " + str(sample_idx) + ": ", sample_likelihood)
46
47
48    def ownTree():
49        topology_array = np.array([float('nan'), 0., 0., 1., 1.])
50        theta_array = [
51            np.array([0.2, 0.8]),
52            np.array([[0.9, 0.1], [0.9, 0.1]]),
53            np.array([[0.05, 0.95], [0.1, 0.9]]),
54            np.array([[0.9, 0.1], [0.9, 0.1]]),
55            np.array([[0.1, 0.9], [0.1, 0.9]])
56        ]
57        t = Tree()
58        t.load_tree_from_direct_arrays(topology_array, theta_array)
59        t.print()
60
61        t.sample_tree(1)
62
63
64        beta = t.filtered_samples[0]
65        print(beta)
66        sample_likelihood = calculate_likelihood(t.root, beta, t.k)
67        print("\tLikelihood: ", sample_likelihood)
68
69    # ownTree()
70
71    if __name__ == "__main__":
72        main("small")
73        main("medium")
74        main("large")
```

## A.2   2.3 - Simple Variational Inference implementation

```python
1   import numpy as np
2   from scipy.stats import norm
3   from scipy.stats import gamma
4   import matplotlib.pyplot as plt
5
6
7   def calcMu(a_N, b_N, x, mu0, lambda0):
8     N = len(x)                    # N
9     x_mean = np.mean(x)           # Mean of x
10    tau_expected = a_N / b_N      # E[tau]
11
12    mu_N = (lambda0 * mu0 + N * x_mean) / (lambda0 + N)
13    lambda_N = (lambda0 + N) * tau_expected
14
15    return mu_N, lambda_N
16
17
18  def calcTau(mu_N, lambda_N, x, a0, b0, mu0, lambda0):
19    N = len(x)                                    # N
20    x_mean = np.mean(x)                           # Mean of x
21    x_square = sum(x ** 2)                        # sum(x^2)
22    mu_expected = mu_N                            # E[mu]
23    mu_square_expected = 1 / lambda_N + mu_N ** 2   # E[mu^2]
24
25    a_N = a0 + N/2
26    b_N = b0 + (1/2) * (x_square + lambda0 * mu0 ** 2) - (mu0 * lambda0 + N * x_mean) * mu_ex
27
28    return a_N, b_N
29
30
31  def getData(n, mu, std_2, seed = 100):
32    np.random.seed(seed)
33    return np.random.normal(mu, std_2, n) # Drawing from gaussian distrbution
34
35
36  def threshold(a_N, b_N, mu_N, lambda_N, a_N_next, b_N_next, mu_N_next, lambda_N_next, thres
37    # Checking if the change in variables are within the threshold
38    if abs(a_N_next / a_N - 1) < threshold:
```

```python
39          if abs(b_N_next / b_N - 1) < threshold:
40            if abs(mu_N_next / mu_N - 1) < threshold:
41              if abs(lambda_N_next / lambda_N - 1) < threshold:
42                return False
43      return True


46  def printVIResults(data, iteration):
47      labels = ["a", "b", "mu", "lambda"]
48      for i in range(len(labels)):
49        print(labels[i] + " | From " + str(data[i]) + " to " + str(data[i+4]))
50      print("It took " + str(iteration) + " iterations.")


53  def printTrueResults(data):
54      labels = ["a", "b", "mu", "beta"]
55      for i in range(len(labels)):
56        print(labels[i] + " | " + str(data[i]) )


59  def q_mu(mus, mu, lambd):
60      return norm.pdf(mus, mu, np.sqrt(1 / lambd)) # Gaussian pdf


63  def q_tau(taus, a, b):
64      return gamma.pdf(taus, a, loc = 0, scale = 1 / b) # Gamma pdf


67  def calcNormGamParam(x, mu_0, lambda_0, a0, b0):
68      N = len(x)                # N
69      x_mean = np.mean(x)       # Mean of x
70      x_square = sum(x ** 2)    # sum(x^2)
71
72      beta = N + lambda_0       # Without tau
73      mu = (N * x_mean + lambda_0 * mu_0) / beta
74
75      a = a0 + N / 2
76      b = b0 + (lambda_0 * (mu_0 ** 2) - beta * (mu ** 2) + x_square) / 2
77
78      return beta, mu, a, b
79
```

```
80
81  def relativeErrors(vi_results, true_results):
82    labels = ["a", "b", "mu"]
83    for i in range(len(labels)):
84      print(labels[i] + " | " + str(round(100*abs(vi_results[i]/true_results[i]-1),2)) + "%"
85
86
87  def vi_algorithm(a_N, b_N, mu_N, lambda_N, x, mu0, lambda0, a0, b0):
88    thresh = 1e-7
89
90    running = True
91    iteration = 0
92
93    # Running algorithm
94    while running:
95      #Calculating new parameters
96      mu_N_next, lambda_N_next = calcMu(a_N, b_N, x, mu0, lambda0)
97      a_N_next, b_N_next = calcTau(mu_N, lambda_N, x, a0, b0, mu0, lambda0)
98
99      #Checking if threshold is satisfied
100     running = threshold(a_N, b_N, mu_N, lambda_N, a_N_next, b_N_next, mu_N_next, lambda_N_n
101
102     # Saving parameters
103     a_N, b_N, mu_N, lambda_N = a_N_next, b_N_next, mu_N_next, lambda_N_next
104     iteration += 1
105
106   return a_N, b_N, mu_N, lambda_N, iteration
107
108
109 def pTrue(x, y, beta, mu, a, b):
110   return norm.pdf(x, mu, np.sqrt(1 / (beta * y))) * gamma.pdf(y, a, loc = 0, scale = 1 / b)
111
112
113 def plotResults(a, b, mu, precision, center, std_2, exact = False):
114   # Getting interval that fits data
115   mus = np.linspace(-4,4,300) * (std_2 ** (0.7)) + center
116   taus = np.linspace(-0.1,6,200) * std_2 ** (-2)
117   # mus = np.linspace(8.5,11.5,300)
118   # taus = np.linspace(-0.1,15,300)
119   # mus = np.linspace(9.9,10.05,300)
120   # taus = np.linspace(3,4,300)
```

```python
121
122       # For the exact posterior
123       if exact:
124         color = "blue"
125         Ms, Ts = np.meshgrid(mus, taus, indexing="ij")
126         Z = np.zeros_like(Ms)
127
128         for i in range(Z.shape[0]):
129           for j in range(Z.shape[1]):
130               Z[i][j] = pTrue(mus[i], taus[j], precision, mu, a, b)
131
132       # The posterior calculated by the VI algorithm
133       else:
134         color = "red"
135         q_mus = q_mu(mus, mu, precision)
136         q_taus = q_tau(taus, a, b)
137         Ms, Ts = np.meshgrid(mus, taus, indexing="ij")
138
139         Z = np.outer(q_mus, q_taus)
140
141       # Plotting the contour
142       plt.contour(Ms, Ts, Z, 10, colors = color)
143
144
145   def main():
146       # Data attributes
147       n = 10000
148       center = -2
149       std = 30
150
151       print("1. Simulating guassian data with\nmu =", center, "| sigma =", std, "| n =", n)
152
153       x = getData(n, center, std, 1021)
154
155       # --------------------------------------------------------------------------------
156
157       a0, b0 = 0.01, 9
158       mu0, lambda0 = -2, 10
159
160       print("\n2. Setting prior parameters\na0 =", a0, "| b0 =", b0, "| mu0 =", mu0, "| lambda0
161
```

```python
162      # ----------------------------------------------------------------------------------
163
164      a_start, b_start = 1e-9, 1e-9
165      mu_start, lambda_start = 1e-9, 1e-9
166
167      print("\n3. Setting start values\na_start =", a_start, "| b_start =", b_start, "| mu_star
168
169      a_N, b_N = a_start, b_start
170      mu_N, lambda_N = mu_start, lambda_start
171
172      # ----------------------------------------------------------------------------------
173
174      print("\n4. Running VI algorithm...")
175
176      a_N, b_N, mu_N, lambda_N, iteration = vi_algorithm(a_start, b_start, mu_start, lambda_sta
177
178      printVIResults([a_start, b_start, mu_start, lambda_start, a_N, b_N, mu_N, lambda_N], iter
179
180      # ----------------------------------------------------------------------------------
181
182      print("\n5. Calculating true posterior...")
183
184      beta, mu, a, b = calcNormGamParam(x, mu0, lambda0, a0, b0)
185
186      printTrueResults([a, b, mu, beta])
187
188      # ----------------------------------------------------------------------------------
189
190      print("\n6. Calculate relative errors")
191
192      relativeErrors([a_N, b_N, mu_N], [a, b, mu])
193
194      # ----------------------------------------------------------------------------------
195
196      print("\n7. Plot results")
197
198      plotResults(a_N, b_N, mu_N, lambda_N, center, std)
199
200      plotResults(a, b, mu, beta, center, std, True)
201
202      plt.show()
```

```
203
204
205   main()
```

## A.3   2.4 - Mixture of trees with observable variables implementation

```python
1   import numpy as np
2   import matplotlib.pyplot as plt
3   import Kruskal_v1 as kr
4   import Kruskal_v2 as kr2
5   import queue
6   from Tree import TreeMixture
7   from Tree import Tree
8   import sys
9   import random
10  from Phylogeny import tree_to_newick_rec
11  from tabulate import tabulate
12  import dendropy
13
14
15  def save_results(loglikelihood, topology_array, theta_array, filename):
16
17      likelihood_filename = filename + "_em_loglikelihood.npy"
18      topology_array_filename = filename + "_em_topology.npy"
19      theta_array_filename = filename + "_em_theta.npy"
20      print("Saving log-likelihood to ", likelihood_filename, ", topology_array to: ", topolo
21            ", theta_array to: ", theta_array_filename, "...")
22      np.save(likelihood_filename, loglikelihood)
23      np.save(topology_array_filename, topology_array)
24      np.save(theta_array_filename, theta_array)
25
26
27  def em_algorithm(seed_val, samples, num_clusters, max_num_iter=100, debugging = False):
28
29      # Set the seeds
30      np.random.seed(seed_val)
31      num_start_tms = 100
32      seeds = np.array(np.random.ranf((num_start_tms,)) * 1e5).astype(int)
33
34
35      # Run the algorithm
36      print("Running EM algorithm...")
37
38
```

```python
39      # Creating intital tms
40      print("\nCreating initial", num_start_tms, "tree mixtures...")
41      tree_mixtures = create_initial_tms(seeds, num_clusters, samples)
42
43      if debugging:
44          em_one_dimension(tree_mixtures[0], samples, num_clusters, 10)
45
46          return
47
48      # First step of sieving
49      print("\nRunning 10 iterations for the first", num_start_tms, "tree mixtures...\n")
50      loglikelihoods = []
51      j = 0
52      for tm in tree_mixtures:
53          loglikelihoods.append(em_one_dimension(tm, samples, num_clusters, 10))
54          j += 1
55          if j % 10 == 0:
56              print(str(j) + " of the first", num_start_tms, "tree mixtures done")
57
58
59      # Choosing 10 best tree mixtures
60      num_best_trees = 10
61      print("\nChoosing the", num_best_trees, "best tree mixtures")
62      indices = np.argsort(-np.array(loglikelihoods))
63      sieved_indices = indices[:num_best_trees]
64      best_seeds = seeds[sieved_indices]
65
66
67      # Recreating best tms
68      print("\nRecreating the", num_best_trees, "best tree mixtures and running them for", ma
69      tree_mixtures = create_initial_tms(best_seeds, num_clusters, samples)
70
71
72      # Running 100 iterations for the best trees
73      loglikelihoods = []
74      j = 0
75      for tm in tree_mixtures:
76          loglikelihoods.append(em_one_dimension(tm, samples, num_clusters, 100))
77          j += 1
78          print(str(j) + " of the first", num_best_trees, "tree mixtures done")
79
```

```python
80        # Choosing the best tree mixture
81        print("\nChoosing the best tree mixtures and running it until convergence\n")
82        indices = np.argsort(-np.array(loglikelihoods))
83        best_seed = best_seeds[indices[0]]
84        tree_mixtures = create_initial_tms([best_seed], num_clusters, samples)
85        threshold = 1e-6
86
87        tm, loglikelihoods = em_one_dimension_final(tree_mixtures[0], samples, num_clusters, ma
88
89        print("\n The final loglikelihood was", loglikelihoods[-1])
90
91        topology_list, theta_list = get_arrays(tm, num_clusters)
92
93        return loglikelihoods, topology_list, theta_list, tm.pi
94
95
96    def get_arrays(tm, num_clusters):
97        topology_list = []
98        theta_list = []
99        for i in range(num_clusters):
100           topology_list.append(tm.clusters[i].get_topology_array())
101           theta_list.append(tm.clusters[i].get_theta_array())
102
103       topology_list = np.array(topology_list)
104       theta_list = np.array(theta_list)
105
106       return topology_list, theta_list
107
108
109   def create_initial_tms(seeds, num_clusters, samples):
110       tree_mixtures = []
111       for seed in seeds:
112           tm = TreeMixture(num_clusters=num_clusters, num_nodes=samples.shape[1])
113           tm.simulate_pi(seed_val=seed)
114           tm.simulate_trees(seed_val=seed)
115           tree_mixtures.append(tm)
116
117       return tree_mixtures
118
119
120   def em_one_dimension_final(tm, samples, num_clusters, max_num_iterations, threshold):
```

```python
121        d = tm.clusters[0].k
122        loglikelihoods = []
123        iteration = 0
124
125        for _ in range(max_num_iterations):
126            # Step 0 - Create new tree mixture
127            tm_new = TreeMixture(num_clusters=num_clusters, num_nodes=samples.shape[1])
128
129            # Step 1 - Calculate responsibilities
130            r_new = responsibilities(tm, samples)
131
132            # Step 2 - Calculate pi'
133            tm_new.pi = new_pi(r_new)
134
135            # Step 3 - Create G_k's
136            tm_new.clusters = create_new_tms2(r_new, samples, num_clusters, d)
137
138            iteration += 1
139            loglikelihoods.append(log_likelihood(tm_new, samples)[0])  # Calculating loglikelih
140
141            tm = tm_new
142
143            if iteration > 1:
144                if loglikelihoods[-2] > loglikelihoods[-1]:
145                    print("!! The likelihood decreased !!") # Debugging
146                if abs(loglikelihoods[-1] / loglikelihoods[-2] - 1) < threshold:
147                    break
148
149        if iteration < max_num_iterations:
150            print("The algorithm converged after", iteration, "iterations")
151        else:
152            print("The algorithm did not converge after", iteration, "iterations")
153
154        return tm, loglikelihoods
155
156
157    def em_one_dimension(tm, samples, num_clusters, max_num_iterations):
158        d = tm.clusters[0].k
159
160        for _ in range(max_num_iterations):
161            # Step 0 - Create new tree mixture
```

```
162            tm_new = TreeMixture(num_clusters=num_clusters, num_nodes=samples.shape[1])
163
164            # Step 1 - Calculate responsibilities
165            r_new = responsibilities(tm, samples)
166
167            # Step 2 - Calculate pi'
168            tm_new.pi = new_pi(r_new)
169
170            # Step 3 - Create G_k's
171            tm_new.clusters = create_new_tms2(r_new, samples, num_clusters, d) # Adding the tre
172
173            tm = tm_new
174
175        loglikelihood = log_likelihood(tm_new, samples) # Calculating loglikelihood
176
177        return loglikelihood[0]
178
179
180    def create_new_tms(r_new, samples, num_clusters, d):
181        clusters_new = []
182        num_vertices = samples.shape[1]
183        G_ks = [kr.Graph(num_vertices) for i in range(num_clusters)] # Creating a graph for eve
184        for (k, G_k) in enumerate(G_ks):
185            root_cdf = np.array([q_root(k, a, r_new, samples) for a in range(d)]) # Calculating
186
187            for s in range(num_vertices):
188                for t in range(s):
189                    G_k.addEdge(t, s, mutual_information(k, t, s, r_new, samples, d)) # Assign
190
191            # print(G_k.graph)
192            max_G = G_k.maximum_spanning_tree() # Calculating the maximum spanning tree
193
194            clusters_new.append(create_tree(max_G, k, r_new, samples, d, root_cdf)) # Adding th
195
196        return clusters_new
197
198
199    def create_new_tms2(r_new, samples, num_clusters, d):
200        clusters_new = []
201        num_vertices = samples.shape[1]
202        G_ks = [set() for i in range(num_clusters)] # Creating a graph for every cluster
```

```python
203     for (k, G_k) in enumerate(G_ks):
204         root_cdf = np.array([q_root(k, a, r_new, samples) for a in range(d)]) # Calculating
205
206         for s in range(num_vertices):
207             for t in range(s):
208                 G_k.add((t, s, mutual_information(k, t, s, r_new, samples, d))) # Assign va
209
210         vertices = list(range(num_vertices))
211         graph = {
212             'vertices': vertices,
213             'edges': G_k
214         }
215
216         result = kr2.maximum_spanning_tree(graph) # Calculating the maximum spanning tree
217         max_G = np.array([np.array([edge[0], edge[1]]) for edge in result])
218
219         clusters_new.append(create_tree(max_G, k, r_new, samples, d, root_cdf)) # Adding th
220
221     return clusters_new
222
223
224 def q_root(k, a, r, samples):
225     return q_single(k, 0, a, r, samples)
226
227
228 def log_likelihood(tm, samples):
229     return sum(np.log(responsibilities(tm, samples, True)))
230
231
232 def nan_function(x):
233     if x == 0:
234         return float('nan')
235     else:
236         return x
237
238
239 def create_tree(G, k, r, samples, d, root_cdf):
240     topology_array, index_array = topology_index(G, d)
241
242     theta_array = t_array(k, topology_array, index_array, r, samples, d)
243
```

```
244        theta_array.insert(0, root_cdf)
245
246        t = Tree()
247        t.load_tree_from_direct_arrays(np.array(topology_array), theta_array)
248
249        return t
250
251
252    def t_array(k, topology_array, index_array, r, samples, d):
253        theta_array = []
254        for s in range(1,len(topology_array)):
255            t = topology_array[s] # Index of parent
256            theta = np.array([[q_conditional(k, s, t, a, b, r, samples) for a in range(d)] for
257            theta_array.append(theta)
258
259        return theta_array
260
261
262    def topology_index(G, d):
263        topology_array = [float('nan')]
264        index_array = []
265        visited_rows = []
266
267        q = queue.Queue()
268        q.put(0)
269        i = 0
270
271        while not q.empty():
272            x = q.get()
273
274            index_1, index_2 = np.where(G == x)
275            index_array.append(x)
276            for j in range(len(index_1)):
277                row = index_1[j]
278                descendant = int(G[row, (index_2[j] + 1) % 2])
279                if row not in visited_rows:
280                    visited_rows.append(row)
281                    topology_array.append(i)
282                    q.put(descendant)
283            i += 1
284
```

```python
285         return topology_array, index_array
286
287
288    def mutual_information(k, s, t, r, samples, d):
289        information = 0
290
291        for a in range(d):
292            q_a = q_single(k, s, a, r, samples)
293            for b in range(d):
294                q_b = q_single(k, t, b, r, samples)
295                q_ab = q_joint(k, s, t, a, b, r, samples)
296                if q_ab != 0:
297                    information += (q_ab + sys.float_info.epsilon) * np.log(q_ab / (q_a * q_b )
298
299        return information
300
301
302    def q_conditional(k, s, t, a, b, r, samples):
303        q = q_joint(k, s, t, a, b, r, samples) / q_single(k, t, b, r, samples)
304
305        return q
306
307
308    def q_single(k, s, a, r, samples):
309        denominator = sum(r[:,k]) # Sum of r_n,k over N
310
311        column = samples[:,s]
312
313        indices = np.where(column == a)[0]
314
315        numerator = sum(r[indices,k])
316
317        q = numerator / denominator
318
319        return q
320
321
322    def q_joint(k, s, t, a, b, r, samples):
323        denominator = sum(r[:,k]) # Sum of r_n,k over N
324
325
```

```python
326         columns = samples[:,[s,t]]
327
328         indices = np.unique(np.where((columns == [a,b]).all(axis=1))[0])
329
330
331         numerator = sum(r[indices,k])
332
333         q = numerator / denominator
334
335         return q
336
337
338     def responsibilities(tm, samples, log_prob = False):
339         n_clusters = len(tm.clusters)
340         n_samples = len(samples)
341
342         p_matrix = np.zeros((n_samples, n_clusters))
343
344         for i in range(n_samples): # Calculate probabilities of every sample for every cluster
345             p_matrix[i,:] = [responsibility(tree.root, samples[i]) for tree in tm.clusters]
346
347         r_unscaled = p_matrix * tm.pi # Multiply by categorical
348
349         probabilities = r_unscaled.sum(axis=1).reshape((n_samples,1)) # Get normalizing factors
350
351         if log_prob:
352             return probabilities
353
354         r = r_unscaled * probabilities ** (-1) # Normalize r
355
356         return r
357
358
359     def new_pi(r):
360         N = r.shape[0]
361         return r.sum(axis=0) / N
362
363
364     def responsibility(node, sample):
365         if node.ancestor == None:
366             p = node.cat[sample[int(node.name)]] # Value of node
```

```python
367        else:
368            p =  node.cat[sample[int(node.ancestor.name)]][sample[int(node.name)]] # Value of a
369        if len(node.descendants) > 0:
370            for descendant in node.descendants:
371                p *= responsibility(descendant, sample)
372        return p


375 def rf_analysis(real_values_filename, output_filename, num_clusters):
376     print("\n4.1.1 Loading ground truth trees from Newick files:\n")
377
378     # If you want to compare two trees, make sure you specify the same Taxon Namespace!
379     tns = dendropy.TaxonNamespace()
380
381     realTrees = []
382     for i in range(num_clusters):
383         filename = real_values_filename + "_tree_" + str(i) + "_newick.txt"
384         with open(filename, 'r') as input_file:
385             newick_str = input_file.read()
386         realTrees.append(dendropy.Tree.get(data=newick_str, schema="newick", taxon_namespac
387
388
389     print("\n4.1.2 Loading inferred trees")
390     filename = output_filename + "_em_topology.npy"  # This is the result you have.
391     topology_list = np.load(filename)
392
393     inferredTrees = []
394     for i in range(num_clusters):
395         rt = Tree()
396         rt.load_tree_from_direct_arrays(topology_list[i])
397         inferredTrees.append(dendropy.Tree.get(data=rt.newick, schema="newick", taxon_names
398
399     print("\n4.1.3 Compare trees and print Robinson-Foulds (RF) distance:\n")
400     for i in range(num_clusters):
401         print("\tt" +str(i) + " vs inferred trees")
402         for j in range(num_clusters):
403             print("\tRF distance: \t", dendropy.calculate.treecompare.symmetric_difference(
404
405
406 def true_log_likelihood(real_values_filename, sample_filename, num_clusters):
407     samples = np.loadtxt(sample_filename, delimiter="\t", dtype=np.int32)
```

```python
408
409     tm = TreeMixture(num_clusters=num_clusters, num_nodes=samples.shape[1])
410
411     tm.load_mixture(real_values_filename)
412
413     topology_array, theta_array = get_arrays(tm, num_clusters)
414
415     print("\nStructure of the true trees:")
416     for i in range(num_clusters):
417         print("\n\tCluster: ", i)
418         print("Pi: ", tm.pi[i])
419         print("\tTopology: \t", topology_array[i])
420         print("\tTheta: \t", theta_array[i])
421
422     return log_likelihood(tm, samples)
423
424
425 def create_real_trees1():
426     # Create tree for scenario 1
427     filename = "data/q2_4/more_samples"
428
429     tm = TreeMixture(num_clusters=3, num_nodes=5)
430     tm.load_mixture("data/q2_4/q2_4_tree_mixture.pkl")
431     tm.samples = list()
432     tm.sample_assignments = list()
433
434     tm.sample_mixtures(1000, 1337)
435
436     tm.save_mixture(filename)
437
438
439 def create_real_trees2():
440     # Create tree for scenario 2
441     filename = "data/q2_4/more_nodes_less_clusters.pkl"
442
443     tm = TreeMixture(num_clusters=2, num_nodes=7)
444
445     tm.simulate_pi(seed_val=1337)
446     tm.simulate_trees(seed_val=1337)
447
448     tm.sample_mixtures(100, 1337)
```

```
449
450        tm.save_mixture(filename, True)
451
452
453   def create_real_trees3():
454        # Create tree for scenario 3
455        filename = "data/q2_4/less_nodes_more_clusters.pkl"
456
457        tm = TreeMixture(num_clusters=5, num_nodes=4)
458
459        tm.simulate_pi(seed_val=1337)
460        tm.simulate_trees(seed_val=1337)
461
462        tm.sample_mixtures(100, 1337)
463
464        tm.save_mixture(filename, True)
465
466
467   def main(scenario = "normal"):
468        seed_val = 123412567
469
470        if scenario == "normal":
471            sample_filename = "data/q2_4/q2_4_tree_mixture.pkl_samples.txt"
472            output_filename = "q2_4_results.txt"
473            real_values_filename = "data/q2_4/q2_4_tree_mixture.pkl"
474            num_clusters = 3
475        elif scenario == "more samples":
476            sample_filename = "data/q2_4/more_samples_samples.txt"
477            output_filename = "q2_4_1_results.txt"
478            real_values_filename = "data/q2_4/q2_4_tree_mixture.pkl"
479            num_clusters = 3
480        elif scenario == "ml":
481            sample_filename = "data/q2_4/more_nodes_less_clusters.pkl_samples.txt"
482            output_filename = "q2_4_2_results.txt"
483            real_values_filename = "data/q2_4/more_nodes_less_clusters.pkl"
484            num_clusters = 2
485        elif scenario == "lm":
486            sample_filename = "data/q2_4/less_nodes_more_clusters.pkl_samples.txt"
487            output_filename = "q2_4_3_results.txt"
488            real_values_filename = "data/q2_4/less_nodes_more_clusters.pkl"
489            num_clusters = 5
```

```
490
491     print("\n1. Load samples from txt file.\n")
492
493     samples = np.loadtxt(sample_filename, delimiter="\t", dtype=np.int32)
494     num_samples, num_nodes = samples.shape
495     print("\tnum_samples: ", num_samples, "\tnum_nodes: ", num_nodes)
496     print("\tSamples: \n", samples)
497
498     print("\n2. Run EM Algorithm.\n")
499
500     loglikelihood, topology_array, theta_array, pi = em_algorithm(seed_val, samples, num_cl
501
502     print("\n3. Save, print and plot the results.\n")
503
504     save_results(loglikelihood, topology_array, theta_array, output_filename)
505
506     for i in range(num_clusters):
507         print("\n\tCluster: ", i)
508         print("Pi: ", pi[i])
509         print("\tTopology: \t", topology_array[i])
510         print("\tTheta: \t", theta_array[i])
511
512     plt.figure(figsize=(8, 3))
513     plt.subplot(121)
514     plt.plot(np.exp(loglikelihood), label='Estimated')
515     plt.ylabel("Likelihood of Mixture")
516     plt.xlabel("Iterations")
517     plt.subplot(122)
518     plt.plot(loglikelihood, label='Estimated')
519     plt.ylabel("Log-Likelihood of Mixture")
520     plt.xlabel("Iterations")
521     plt.legend(loc=(1.04, 0))
522     plt.show()
523
524     if real_values_filename != "":
525         print("\n4. Retrieve real results and compare.\n")
526         print("\tComparing the results with real values...")
527
528         print("\t4.1. Make the Robinson-Foulds distance analysis.\n")
529         rf_analysis(real_values_filename, output_filename, num_clusters)
530
```

```python
531          print("\n\t4.2. Make the likelihood comparison.\n")
532          true_likelihood = true_log_likelihood(real_values_filename, sample_filename, num_cl
533
534          data = [("EM", loglikelihood[-1], np.exp(loglikelihood[-1])),
535                  ("True", true_likelihood, np.exp(true_likelihood))]
536
537          headers = ["","Loglikelihood","Likelihood"]
538
539          print(tabulate(data, headers=headers))
540
541
542  # create_real_trees1()
543
544
545  # create_real_trees2()
546
547
548  # create_real_trees3()
549
550
551  if __name__ == "__main__":
552      main()
```