

COMP0086

# Probabilistic and Unsupervised Learning Coursework

October 2024

University College London

# Contents

<b>1</b>	<b>Models for binary vectors</b>	<b>3</b>
1.1	Question 1a . . . . .	3
1.2	Question 1b . . . . .	3
1.3	Question 1c . . . . .	4
1.4	Question 1d . . . . .	6
1.5	Question 1e . . . . .	7
<b>2</b>	<b>Model selection</b>	<b>10</b>
2.1	Question 2a . . . . .	10
2.2	Question 2b . . . . .	10
2.3	Question 2c . . . . .	11
<b>3</b>	<b>EM for Binary Data</b>	<b>15</b>
3.1	Question 3a . . . . .	15
3.2	Question 3b . . . . .	15
3.3	Question 3c . . . . .	16
3.4	Question 3d . . . . .	18
3.5	Question 3e . . . . .	27
3.6	Bonus: Question 3f . . . . .	30
3.7	Bonus: Question 3g . . . . .	31
<b>4</b>	<b>Bonus: LGSSMs, EM and SSID</b>	<b>32</b>
4.1	Question 4a . . . . .	32
<b>5</b>	<b>Decrypting Messages with MCMC</b>	<b>37</b>
5.1	Question 5a . . . . .	37
5.2	Question 5b . . . . .	40
5.3	Question 5c . . . . .	40
5.4	Question 5d . . . . .	41
5.5	Question 5e . . . . .	46
5.6	Question 5f . . . . .	46
<b>6</b>	<b>Bonus: Implementing Gibbs sampling for LDA</b>	<b>49</b>
<b>7</b>	<b>Optimization</b>	<b>50</b>

7.1	Question 7a . . . . .	50
7.2	Question 7b . . . . .	51
7.2.1	Question 7b i . . . . .	51
7.2.2	Question 7b ii . . . . .	52
<b>8</b>	<b>Bonus: Eigenvalues as solutions of an optimization problem</b>	<b>53</b>
8.1	Question 8a . . . . .	53
8.2	Question 8b . . . . .	54
8.3	Question 8c . . . . .	55

# 1 Models for binary vectors

## 1.1 Question 1a

A multivariate Gaussian would not be an appropriate model because of the nature of the data. The dataset contains discrete data, describing the pixels with binary values (0 or 1). In contrast, the multivariate normal distribution is designed to model continuous data, with values that can be any real number. The Gaussian distribution would predict a range of values between 0 and 1, which aren't valid pixel values for the images encoded in binary. On the other hand, the Bernoulli distribution has the same sample space as the dataset, and thus would be a better choice of model to describe it.

## 1.2 Question 1b

For a D-dimensional multivariate Bernoulli distribution, the likelihood of observing a single image  $\mathbf{x}^{(n)}$  given the parameter vector  $\mathbf{p}$  is given by:

$$P(\mathbf{x}^{(n)}|\mathbf{p}) = \prod_{d=1}^D p_d^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}}$$

Knowing that the N images are modeled as i.i.d. samples from the distribution, the likelihood function can be written as:

$$L(\mathbf{p}) = \prod_{n=1}^N P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}|\mathbf{p}) = \prod_{n=1}^N \prod_{d=1}^D p_d^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}}$$

We take the logarithm of the likelihood function to obtain the log-likelihood:

$$\log L(\mathbf{p}) = \sum_{n=1}^N \sum_{d=1}^D \left[ x_d^{(n)} \log(p_d) + (1 - x_d^{(n)}) \log(1 - p_d) \right]$$

This transforms the products into sums, and doesn't move the extrema as the logarithm is a monotonically increasing function. To find the maximum likelihood estimate, we take the derivative of the log-likelihood with respect to each parameter  $p_d$  and set it equal to zero:

$$\frac{\partial}{\partial p_d} \log L(\mathbf{p}) = \sum_{n=1}^N \left[ \frac{x_d^{(n)}}{p_d} - \frac{1 - x_d^{(n)}}{1 - p_d} \right] = 0$$

Rearranging the above equation gives us:

$$\sum_{n=1}^N x_d^{(n)} \cdot (1 - p_d) = \sum_{n=1}^N (1 - x_d^{(n)}) \cdot p_d$$

Which we can simplify to get:

$$\sum_{n=1}^N x_d^{(n)} = p_d \cdot N$$

Thus, the maximum likelihood estimate for  $p_d$  is:

$$p_d = \frac{1}{N} \sum_{n=1}^N x_d^{(n)}$$

Or in vector form:

$$\mathbf{p}_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)}$$

### 1.3 Question 1c

To find the Maximum a posteriori (MAP) estimate, we use Bayes theorem to incorporate prior beliefs. In our case, Bayes theorem states that:

$$P(\mathbf{p} | \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}) = \frac{P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} | \mathbf{p}) \cdot P(\mathbf{p})}{P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\})}$$

To find the MAP estimate, we maximise the posterior:

$$\begin{aligned}
\mathbf{p}_{\text{MAP}} &= \arg \max_{\mathbf{p}} P(\mathbf{p} \mid \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}) \\
&= \arg \max_{\mathbf{p}} P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} \mid \mathbf{p}) \cdot P(\mathbf{p}) \\
&= \arg \max_{\mathbf{p}} (\log P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} \mid \mathbf{p}) + \log P(\mathbf{p}))
\end{aligned}$$

As before, the likelihood of observing  $N$  independent images is:

$$P(\mathbf{x}^{(n)} \mid \mathbf{p}) = \prod_{n=1}^N \prod_{d=1}^D p_d^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}}$$

And the corresponding log-likelihood:

$$\log P(\mathbf{x}^{(n)} \mid \mathbf{p}) = \sum_{n=1}^N \sum_{d=1}^D \left[ x_d^{(n)} \log(p_d) + (1 - x_d^{(n)}) \log(1 - p_d) \right]$$

We know that the prior for each  $p_d$  is given as a Beta distribution by:

$$P(p_d) = \frac{1}{B(\alpha, \beta)} p_d^{\alpha-1} (1 - p_d)^{\beta-1}$$

Since the priors on each  $p_d$  are independent, the prior for  $\mathbf{p}$  is:

$$P(\mathbf{p}) = \prod_{d=1}^D P(p_d) = \prod_{d=1}^D \frac{1}{B(\alpha, \beta)} p_d^{\alpha-1} (1 - p_d)^{\beta-1}$$

We get the log-prior by taking the logarithm, which yields:

$$\log P(\mathbf{p}) = \sum_{d=1}^D [(\alpha - 1) \log(p_d) + (\beta - 1) \log(1 - p_d)] + \text{constant}$$

The constants won't change the position of the extrema so they aren't very important to the calculation. The log of the posterior is:

$$\log P(\mathbf{p} \mid \mathbf{x}^{(n)}) \propto \log P(\mathbf{x}^{(n)} \mid \mathbf{p}) + \log P(\mathbf{p})$$

We can then substitute the expressions for  $\log P(\mathbf{x}^{(n)} \mid \mathbf{p})$  and  $\log P(\mathbf{p})$  to get:

$$\log P(\mathbf{p} | \mathbf{x}^{(n)}) \propto \sum_{d=1}^D \left[ \sum_{n=1}^N \left( x_d^{(n)} \log(p_d) + (1 - x_d^{(n)}) \log(1 - p_d) \right) + (\alpha - 1) \log(p_d) + (\beta - 1) \log(1 - p_d) \right]$$

To maximize this, we take the derivative with respect to  $p_d$  and set it to zero:

$$\frac{\partial}{\partial p_d} \log P(\mathbf{p} | \mathbf{x}^{(n)}) = \frac{\sum_{n=1}^N x_d^{(n)}}{p_d} - \frac{\sum_{n=1}^N (1 - x_d^{(n)})}{1 - p_d} + \frac{\alpha - 1}{p_d} - \frac{\beta - 1}{1 - p_d} = 0$$

Rearranging the terms yields:

$$\sum_{n=1}^N x_d^{(n)} + (\alpha - 1) = p_d (N + \alpha + \beta - 2)$$

Solving for  $p_d$  gives the MAP estimate for  $p_d$ :

$$p_d = \frac{\sum_{n=1}^N x_d^{(n)} + \alpha - 1}{N - 2 + \alpha + \beta}$$

We can write the MAP estimate in vectorised form:

$$\mathbf{p}_{MAP} = \frac{\sum_{n=1}^N \mathbf{x}^{(n)} + \alpha - 1}{N - 2 + \alpha + \beta}$$

because the images were modelled as independently and identically distributed samples.

## 1.4 Question 1d

Listing 1 shows the code that was written to find the Maximum Likelihood estimate of the parameters, as well as display Figure 1.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # We load the data into a numpy array
5 x = np.loadtxt('binarydigits.txt')
6

```

```

7 # We compute the maximum likelihood estimate of the parameters of the multivariate Bernoulli
8 maximum_likelihood = np.mean(x, axis=0)
9
10 # We plot the learnt parameters with matplotlib
11 plt.figure()
12 plt.title("Maximum Likelihood estimate of the parameters")
13 plt.imshow(np.reshape(maximum_likelihood, (8, 8)), cmap='inferno')
14 plt.colorbar()
15 plt.show()

```

Listing 1: Python code used to compute ML estimates and display figure 1

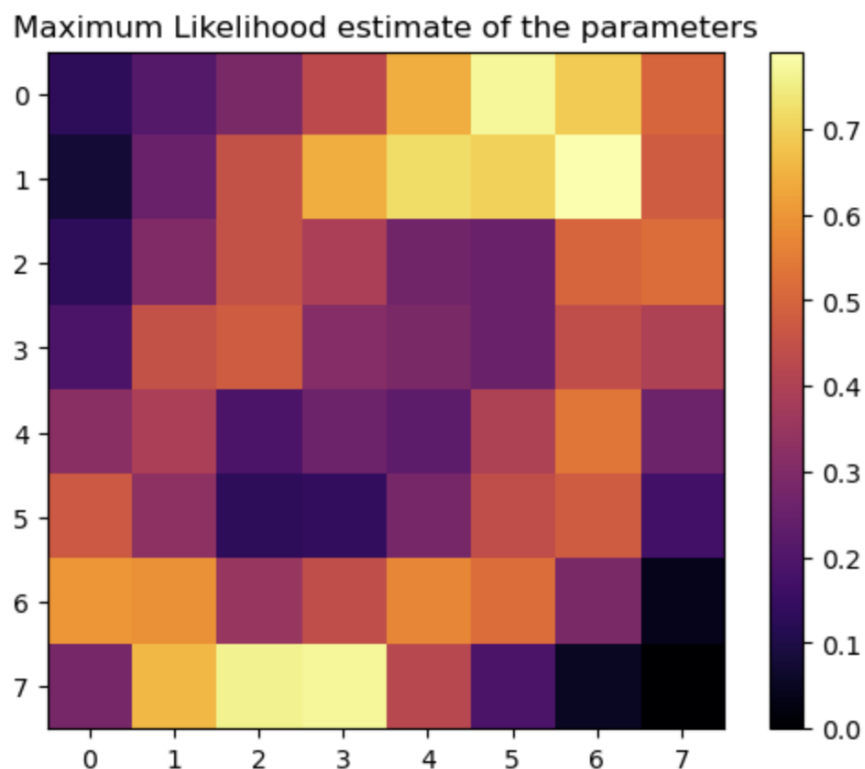


Figure 1: Maximum Likelihood estimate of the parameters

## 1.5 Question 1e

Listing 2 shows the code that was written to find the MAP estimate of the parameters, as well as display Figure 2.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # We load the data into a numpy array
5 x = np.loadtxt('binarydigits.txt')

```



```

6
7 # We define values of alpha and beta
8 alpha = 3
9 beta = 3
10
11 # We calculate the MAP estimate of the parameters
12 maximum_a_posteriori = (alpha - 1 + np.sum(x, axis=0)) / (x.shape[0] + alpha + beta - 2)
13
14 # We plot the learned parameters
15 plt.figure()
16 plt.title("Maximum a posteriori estimate of the parameters")
17 plt.imshow(np.reshape(maximum_a_posteriori, (8, 8)), cmap='inferno')
18 plt.colorbar()

```

Listing 2: Python code used to compute MAP estimates and display figure 2

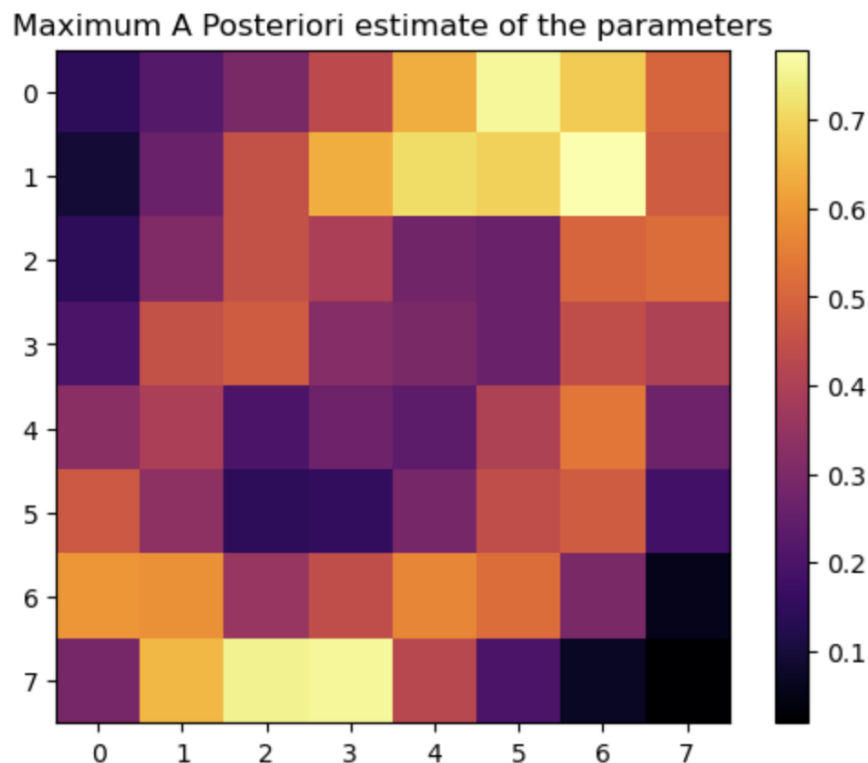


Figure 2: Maximum A Posteriori estimate of the parameters

From figures 1 and 2, we see that both methods find similar estimates of the parameters. However, introducing a prior makes the MAP estimate less certain about its predictions. This is not obvious from figures 1 and 2, so an easy way to show this is by plotting the difference between the MAP and ML estimates:

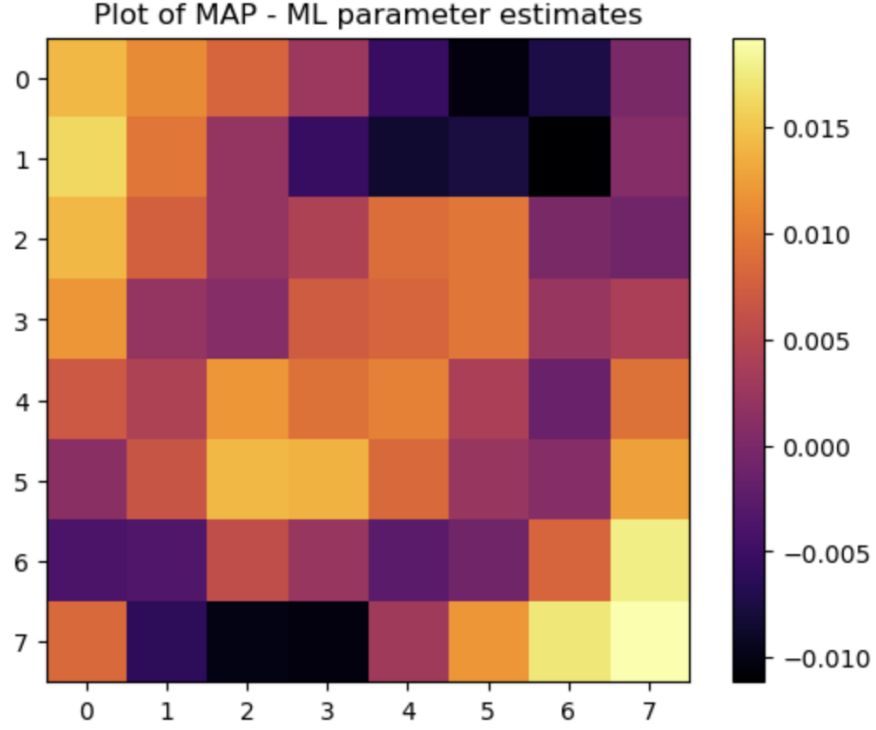


Figure 3: Plot of MAP-ML parameter estimates

We see that the difference is negative when parameters estimates are higher, meaning the MAP is smaller than the ML, and it is positive for lower parameter estimates, meaning the MAP is bigger than the ML. The MAP solution thus gives higher probability to cases it hasn't seen before, and this uncertainty will slowly decrease as more examples are given, since the prior will have less and less impact. The MAP estimate is better because it allows us to have some uncertainty which reflects the quantity of data the model has seen.

## 2 Model selection

### 2.1 Question 2a

In this model, all components are generated from a Bernoulli distribution with  $p_d = 0.5$ . Plugging this into the Bernoulli likelihood equation:

$$P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} | \mathbf{p}^{(n)} = (0.5, 0.5, \dots, 0.5)) = \prod_{n=1}^N \prod_{d=1}^D (0.5)^{x_d^{(n)}} (0.5)^{1-x_d^{(n)}}$$

Since this is a Bernoulli distribution,  $x_d$  can only take the values 0 or 1. If  $x_d = 0$ , the the likelihood simplifies to:

$$P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} | \mathbf{p}^{(n)} = (0.5, 0.5, \dots, 0.5)) = \prod_{n=1}^N \prod_{d=1}^D (0.5)^0 (0.5)^{1-0} = \prod_{n=1}^N \prod_{d=1}^D (0.5)$$

On the other hand if  $x_d = 1$ , the the likelihood simplifies to:

$$P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} | \mathbf{p}^{(n)} = (0.5, 0.5, \dots, 0.5)) = \prod_{n=1}^N \prod_{d=1}^D (0.5)^1 (0.5)^{1-1} = \prod_{n=1}^N \prod_{d=1}^D (0.5)$$

This means that no matter the value of  $x_d$  we have the same likelihood function. Therefore the likelihood function for this model is:

$$P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} | \mathbf{p}^{(n)}) = \prod_{n=1}^N \prod_{d=1}^D (0.5) = (0.5)^{N \cdot D}$$

### 2.2 Question 2b

In this model, all components are generated from Bernoulli distributions with unknown, but identical  $p_d$ . Plugging this into the Bernoulli likelihood equation:

$$P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} | \mathbf{p}^{(n)} = (p_d, p_d, \dots, p_d)) = \prod_{n=1}^N \prod_{d=1}^D (p_d)^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}}$$

To find the relative probability, we have to integrate over all possible values of  $p_d$ , which on our case is just one variable (since all values of  $p_d$  are all identical). This yields:

$$P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}) = \int_0^1 \prod_{n=1}^N \prod_{d=1}^D (p_d)^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}} dp_d$$

By using power rules, we can change the products into sums, which will make a known integral form appear:

$$\begin{aligned} P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}) &= \int_0^1 (p_d)^{\sum_{n=1}^N \sum_{d=1}^D x_d^{(n)}} (1 - p_d)^{\sum_{n=1}^N \sum_{d=1}^D (1-x_d^{(n)})} dp_d \\ &= \int_0^1 (p_d)^{\sum_{n=1}^N \sum_{d=1}^D x_d^{(n)}} (1 - p_d)^{\sum_{n=1}^N \sum_{d=1}^D 1 - \sum_{n=1}^N \sum_{d=1}^D x_d^{(n)}} dp_d \\ &= \int_0^1 (p_d)^{y-1} (1 - p_d)^{z-1} dp_d \\ &= \frac{\Gamma(\sum_{n=1}^N \sum_{d=1}^D x_d^{(n)} + 1) \cdot \Gamma(\sum_{n=1}^N \sum_{d=1}^D 1 - \sum_{n=1}^N \sum_{d=1}^D x_d^{(n)} + 1)}{\Gamma(\sum_{n=1}^N \sum_{d=1}^D 1 + 2)} \\ &= \frac{\Gamma(\sum_{n=1}^N \sum_{d=1}^D x_d^{(n)} + 1) \cdot \Gamma(N \cdot D - \sum_{n=1}^N \sum_{d=1}^D x_d^{(n)} + 1)}{\Gamma(N \cdot D + 2)} \end{aligned}$$

In the final few steps of the computation above, we have used the Beta function integral property:

$$B(y, z) = \int_0^1 t^{y-1} (1 - t)^{z-1} dt$$

as well as the following substitutions:

$$\begin{aligned} y &= \sum_{n=1}^N \sum_{d=1}^D x_d^{(n)} + 1 \\ z &= \sum_{n=1}^N \sum_{d=1}^D 1 - \sum_{n=1}^N \sum_{d=1}^D x_d^{(n)} + 1 \end{aligned}$$

## 2.3 Question 2c

In this model, each component is Bernoulli distributed with separate, unknown  $p_d$ . The likelihood function thus becomes:

$$P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} | \mathbf{p}^{(n)} = (p_1, p_2, p_3, \dots, p_D)) = \prod_{n=1}^N \prod_{d=1}^D (p_d)^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}}$$

To find the relative probability, we have to integrate over all possible values of  $p_d$ , which gives:

$$P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}) = \int_0^1 \dots \int_0^1 \prod_{n=1}^N \prod_{d=1}^D (p_d)^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}} dp_1 dp_2 \dots dp_D$$

We proceed as we did previously:

$$\begin{aligned} P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}) &= \int_0^1 \dots \int_0^1 \prod_{n=1}^N \prod_{d=1}^D (p_d)^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}} dp_1 dp_2 \dots dp_D \\ &= \prod_{d=1}^D \int_0^1 \prod_{n=1}^N (p_d)^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}} dp_d \\ &= \prod_{d=1}^D \int_0^1 (p_d)^{\sum_{n=1}^N x_d^{(n)}} (1 - p_d)^{\sum_{n=1}^N (1-x_d^{(n)})} dp_d \\ &= \prod_{d=1}^D \int_0^1 (p_d)^{y-1} (1 - p_d)^{z-1} dp_d \\ &= \prod_{d=1}^D \frac{\Gamma(\sum_{n=1}^N x_d^{(n)} + 1) \cdot \Gamma(\sum_{n=1}^N (1 - x_d^{(n)}) + 1)}{\Gamma(\sum_{n=1}^N + 2)} \end{aligned}$$

Where in the second step, we have taken out the product over  $d$  so that we could evaluate each integral with one value of  $p_d$ . Then we used the trick of putting the product over  $n$  in the exponent and using the following substitutions:

$$\begin{aligned} y &= \sum_{n=1}^N x_d^{(n)} + 1 \\ z &= \sum_{n=1}^N (1 - x_d^{(n)}) + 1 \end{aligned}$$

to make the Beta function integral appear.

Now we can find the posterior probabilities of each of the three models having generated the data in "binarydigits.txt". Bayes rule states that:

$$P(M_i|\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}) = \frac{P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}|M_i)P(M_i)}{P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\})}$$

The evidence can be written as a sum over our three models, and rewriting the denominator as:  $P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}) = P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}|M_i) \cdot P(M_i)$ , we have:

$$P(M_i|\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}) = \frac{P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}|M_i)}{\sum_{i=1}^3 P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}|M_i)}$$

The posterior probabilities of each of the three models having generated the data in binary-digits.txt is registered in the table below:

Model	Posterior Probability
Model 1	$9.14 \times 10^{-255}$
Model 2	$1.43 \times 10^{-188}$
Model 3	$1.00 \times 10^0$

Table 1: Posterior Probabilities for Each Model

These number (while they don't seem to add up to 1 because of numerical precision) are sensible, since models 1 and 2 would predict that all the pixels have the same probability of being black or white. For example, if all  $p_d = 0.5$ , this means that the colours are just random and that there is no underlying pattern to the data, which is obviously false. The third model's probability is very close to one since we assume that one of the models did generate the data, so its probability is just 1 minus the probability that models 1 and 2 generated the data.

The code used to calculate these results can be found in Listing 3.

```

1 import numpy as np
2 from scipy.special import betaln, logsumexp
3
4 # We load the binary digit data
5 x = np.loadtxt('binarydigits.txt')
6 n, d = x.shape
7
8 # We calculate the log likelihoods for each model
9 # Model 1: p_d = 0.5
10 log_pd_m1 = n * d * np.log(0.5)
11
12 # Model 2: unknown identical p_d
13 # We count the total number of 1s
14 k = np.sum(x).astype(int)
15 log_pd_m2 = betaln(k + 1, n * d - k + 1)
16
17 # Model 3: separate, unknown p_d

```

```

18 k_d = np.sum(x, axis=0).astype(int)
19 log_pd_m3 = np.sum(betaln(k_d + 1, n - k_d + 1))
20
21 # We combine log likelihoods for each model
22 log_pd_total = np.array([log_pd_m1, log_pd_m2, log_pd_m3])
23
24 # We calculate posterior log probabilities for each model
25 log_pm_data = log_pd_total - logsumexp(log_pd_total)
26
27 # We go back to linear space by taking the exponential
28 posterior_probs = np.exp(log_pm_data)
29
30 for i, prob in enumerate(posterior_probs, start=1):
31     print(f"Model {i}: Posterior Probability = {prob:.2e}")

```

Listing 3: Python code used to compute the posterior probabilities of each model having generated the data in binarydigits.txt

## 3 EM for Binary Data

### 3.1 Question 3a

We are looking to model the binary data using a mixture of  $K$  multivariate Bernoulli distributions. The likelihood function for this type of distribution is:

$$P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} | \pi, \mathbf{P}) = \prod_{d=1}^D \prod_{n=1}^N p_{kd}^{x_d^{(n)}} \cdot (1 - p_{kd})^{1-x_d^{(n)}}$$

where the parameters are:

1.  $\pi = (\pi_1, \pi_2, \dots, \pi_K)$ , representing the mixing proportion for the  $k$ -th component. We have  $0 \leq \pi_k \leq 1$  and  $\sum_{k=1}^K \pi_k = 1$ .
2.  $\mathbf{P}$ , which is a matrix containing the Bernoulli parameters  $p_{kd}$ . They represent the probability that a pixel  $d$  has the value 1 with a mixture component  $k$ .

For a single image, the equation for the likelihood function becomes:

$$P(\mathbf{x}^{(n)} | \pi, \mathbf{P}) = \sum_{k=1}^K \pi_k \prod_{d=1}^D p_{kd}^{x_d^{(n)}} \cdot (1 - p_{kd})^{1-x_d^{(n)}}$$

This is the likelihood function for a binary image.

### 3.2 Question 3b

The responsibility is defined by the following expression:

$$r_{nk} = P(s^{(n)} = k | \mathbf{x}^{(n)}, \pi, \mathbf{P})$$

Bayes theorem states that:

$$r_{nk} = \frac{P(\mathbf{x}^{(n)} | s^{(n)} = k, \pi, \mathbf{P}) P(s^{(n)} = k | \pi, \mathbf{P})}{P(\mathbf{x}^{(n)} | \pi, \mathbf{P})}$$

Since the latent variable is not dependent on  $\mathbf{P}$ , we have:  $P(s^{(n)} = k | \pi, \mathbf{P}) = P(s^{(n)} = k | \pi)$ . Additionally, we are told that  $P(s^{(n)} = k | \pi) = \pi_k$ , which lets us write the numerator as:

$$P(\mathbf{x}^{(n)} | s^{(n)} = k, \mathbf{P}) \cdot \pi_k$$



To evaluate the other expression, we take the likelihood function for a mixture component  $k$  and a corresponding probability of choosing component  $k$ ,  $\pi_k$ :

$$P(\mathbf{x}^{(n)} | s^{(n)} = k, \pi, \mathbf{P}) = \prod_{d=1}^D (p_{kd}^{x_d^{(n)}} (1 - p_{kd})^{1-x_d^{(n)}})$$

We know an expression for the denominator  $P(\mathbf{x}^{(n)} | \pi, \mathbf{P})$  from question 3a:

$$P(\mathbf{x}^{(n)} | \pi, \mathbf{P}) = \sum_{j=1}^K \pi_j \prod_{d=1}^D \left( p_{jd}^{x_d^{(n)}} (1 - p_{jd})^{1-x_d^{(n)}} \right)$$

Therefore the responsibility is equal to:

$$r_{nk} = \frac{\pi_k \prod_{d=1}^D \left( p_{kd}^{x_d^{(n)}} (1 - p_{kd})^{1-x_d^{(n)}} \right)}{\sum_{j=1}^K \pi_j \prod_{d=1}^D \left( p_{jd}^{x_d^{(n)}} (1 - p_{jd})^{1-x_d^{(n)}} \right)}$$

### 3.3 Question 3c

We aim to maximize the expected log-joint probability:

$$\arg \max_{\pi, \mathbf{P}} \left\langle \sum_{n=1}^N \log P(\mathbf{x}^{(n)}, s^{(n)} | \pi, \mathbf{P}) \right\rangle_{q(\{s^{(n)}\})}$$

where  $q(\{s^{(n)}\})$  is the distribution over the latent variables given in the previous question. Firstly, we expand the expected log-joint distribution for each data point when  $s^{(n)} = k$ :

$$\log P(\mathbf{x}^{(n)}, s^{(n)} = k | \pi, \mathbf{P}) = \log (P(\mathbf{x}^{(n)} | s^{(n)} = k, \mathbf{P}) \pi_k)$$

where we used  $P(s^{(n)} = k | \pi) = \pi_k$  as before. Writing out the distribution and using the log properties gives:

$$\log P(\mathbf{x}^{(n)}, s^{(n)} = k | \pi, \mathbf{P}) = \log \pi_k + \sum_{d=1}^D \left[ x_d^{(n)} \log p_{kd} + (1 - x_d^{(n)}) \log(1 - p_{kd}) \right]$$

We now calculate the expectation value of this for all data points. This means we sum over  $n$ , and the expectation computation will be a sum over  $k$  since the latent variable is discrete. This gives:

$$\sum_{n=1}^N \sum_{k=1}^K r_{nk} \left[ \log \pi_k + \sum_{d=1}^D \left( x_d^{(n)} \log p_{kd} + (1 - x_d^{(n)}) \log(1 - p_{kd}) \right) \right]$$

The goal now is to take the derivative with respect  $\pi$  and  $\mathbf{P}$ , set them to zero, and solve for  $\pi$  and  $\mathbf{P}$ . These will be the parameters that maximise the expectation value. We will start by maximising with respect to  $\pi$ . Because we are subject to the constraint  $\sum_{k=1}^K \pi_k = 1$ , we can maximise using a Lagrange multiplier  $\lambda$  and a Lagrange function defined as follows:

$$\mathcal{L} = \sum_{k=1}^K \left( \sum_{n=1}^N r_{nk} \right) \log \pi_k + \lambda \left( 1 - \sum_{k=1}^K \pi_k \right)$$

We now take the derivative of  $\mathcal{L}$  with respect to  $\pi_k$ , set it to zero and solve for  $\pi_k$ :

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \pi_k} &= \frac{\sum_{n=1}^N r_{nk}}{\pi_k} - \lambda = 0 \\ \pi_k &= \frac{\sum_{n=1}^N r_{nk}}{\lambda} \end{aligned}$$

Using the constraint  $\sum_{k=1}^K \pi_k = 1$ :

$$\sum_{k=1}^K \pi_k = \sum_{k=1}^K \frac{\sum_{n=1}^N r_{nk}}{\lambda} = \frac{N}{\lambda} = 1$$

We find that  $\lambda = N$ , and the updated mixing coefficients are:

$$\pi_k = \frac{\sum_{n=1}^N r_{nk}}{N}$$

Now to maximise with respect to  $p_{kd}$ . Since constants won't affect where the maximum is, we only take into account terms containing  $p_{kd}$ :

$$\sum_{n=1}^N r_{nk} \left[ x_d^{(n)} \log p_{kd} + (1 - x_d^{(n)}) \log(1 - p_{kd}) \right]$$

Like we did for  $\pi_k$ , we take the derivative with respect to  $p_{kd}$ , set it to zero and solve for  $p_{kd}$ :

$$\begin{aligned}
& \frac{\partial}{\partial p_{kd}} \left( \sum_{n=1}^N r_{nk} \left[ x_d^{(n)} \log p_{kd} + (1 - x_d^{(n)}) \log(1 - p_{kd}) \right] \right) = 0 \\
& \Rightarrow \sum_{n=1}^N r_{nk} \left( \frac{x_d^{(n)}}{p_{kd}} - \frac{1 - x_d^{(n)}}{1 - p_{kd}} \right) = 0 \\
& \Rightarrow \frac{1}{p_{kd}(1 - p_{kd})} \left( \sum_{n=1}^N r_{nk} x_d^{(n)} (1 - p_{kd}) - \sum_{n=1}^N r_{nk} (1 - x_d^{(n)}) p_{kd} \right) = 0 \\
& \Rightarrow \sum_{n=1}^N r_{nk} x_d^{(n)} (1 - p_{kd}) = \sum_{n=1}^N r_{nk} (1 - x_d^{(n)}) p_{kd} \\
& \Rightarrow (1 - p_{kd}) \sum_{n=1}^N r_{nk} x_d^{(n)} = p_{kd} \sum_{n=1}^N r_{nk} (1 - x_d^{(n)}) \\
& \Rightarrow \left( \sum_{n=1}^N r_{nk} x_d^{(n)} \right) = p_{kd} \left( \sum_{n=1}^N r_{nk} x_d^{(n)} + \sum_{n=1}^N r_{nk} (1 - x_d^{(n)}) \right) \\
& \Rightarrow \sum_{n=1}^N r_{nk} x_d^{(n)} = p_{kd} \sum_{n=1}^N r_{nk}
\end{aligned}$$

Finally, solving for  $p_{kd}$ :

$$p_{kd} = \frac{\sum_{n=1}^N r_{nk} x_d^{(n)}}{\sum_{n=1}^N r_{nk}}$$

### 3.4 Question 3d

As suggested by the hint, I've introduced weak priors on  $\mathbf{P}$  and  $\pi$ . I've chosen a uniform Dirichlet distribution for the prior of  $\pi$ , meaning all components are assumed to be equally likely before observing any data. I've chosen a Beta distribution  $\mathbf{P}$ , with both  $\alpha$  and  $\beta$  parameters set to 1 (making it a uniform distribution) to indicate no prior knowledge. I've implemented the EM algorithm for a mixture of  $K$  multivariate Bernoullis and run it for  $K = \{2, 3, 4, 7, 10\}$ . The plots of the log posteriors as a function of the iteration number for each value of  $K$  are shown below.

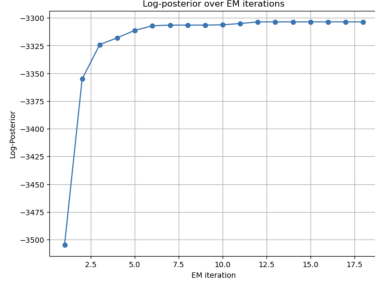


Figure 4:  $K = 2$ , Convergence step = 17

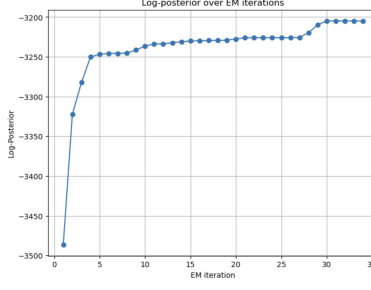


Figure 5:  $K = 3$ , Convergence step = 33

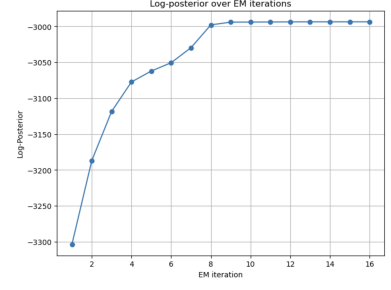


Figure 6:  $K = 4$ , Convergence step = 15

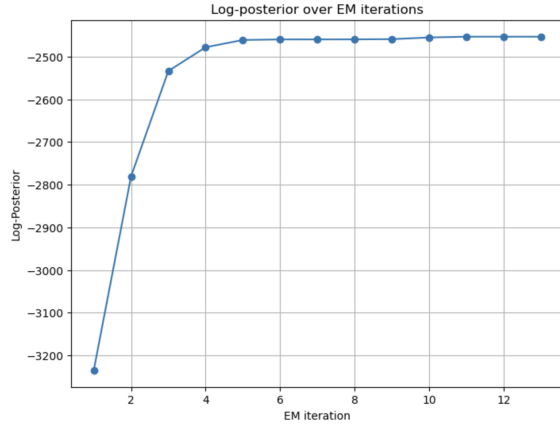


Figure 7:  $K = 7$ , Convergence step = 12

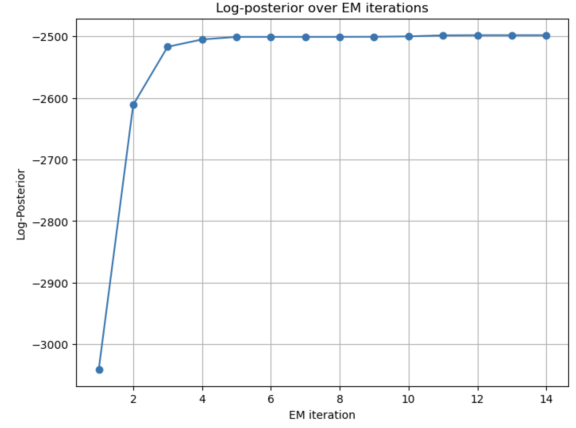


Figure 8:  $K = 10$ , Convergence step = 13

The parameters  $\pi$  and  $\mathbf{P}$  are shown in tables 2 and 3 respectively:

$K$	Values of $\pi_k$
2	0.5804, 0.4196
3	0.2700, 0.1300, 0.6000
4	0.4402, 0.2810, 0.1090, 0.1698
7	0.1600, 0.0900, 0.1600, 0.2300, 0.0500, 0.0600, 0.2500
10	0.1400, 0.0500, 0.1800, 0.1600, 0.0300, 0.0300, 0.0300, 0.0800, 0.2000, 0.1000

Table 2: Table showing the values of  $\pi_k$  for different  $K$

K	Optimised Bernoulli parameters	
2		
3		
4		
7		
10		

Table 3: Table showing the optimised Bernoulli parameters for different K values

The code used to run the EM algorithm and plot the graphs can be found in Listing 4, and the code used to display the parameters can be found in Listing 5:

```

1 import torch
2 import torch.nn.functional as F
3 from torch.distributions import Dirichlet
4 from typing import List, Tuple
5 from scipy.special import betaln
6 import pandas as pd
7 import matplotlib.pyplot as plt
8 import numpy as np
9
10 def init_params(k: int, d: int, epsilon: float = 1e-6):
11     """
12     This function randomly initialises the model parameters (log_pi and log_p_matrix)
13
14     INPUT:
15     k      : Number of mixture components
16     d      : Dimensionality of the data (number of features)
17     epsilon : Small constant used to avoid numerical instability
18
19     OUTPUT:
20     log_pi      : Log mixing proportions, Shape: (1, k)
21     log_p_matrix : Log probability of the Bernoulli parameters, Shape: (d, k)
22     """
23     # We initialise the mixing proportions using a Dirichlet distribution

```

```

25     # The .unsqueeze method adds a dimension to the pi vector (goes from shape (k) to (1, k)
26     )
27     pi = Dirichlet(torch.ones(k)).sample().unsqueeze(0)
28     log_pi = torch.log(pi)
29
30     # We initialize the Bernoulli parameters
31     # I've encountered numerical errors due to taking the log of 0 later on in the code
32     # The clamping with the epsilon variable is here to avoid that
33     p_initial = torch.rand(d, k)
34     p_initial = torch.clamp(p_initial, min=epsilon, max=1 - epsilon)
35     log_p_matrix = torch.log(p_initial)
36
37     return log_pi, log_p_matrix
38
39 def compute_log_one_minus_p_matrix(log_p_matrix: torch.Tensor, epsilon: float = 1e-6) ->
40     torch.Tensor:
41     """
42     This function computes log(1 - p) from log(p)
43
44     INPUT:
45     log_p_matrix      : Log probabilities of p, Shape: (d, k)
46     epsilon           : Small constant used to avoid numerical instability
47
48     OUTPUT:
49     log_one_minus_p_matrix: Log probabilities of (1 - p), Shape: (d, k)
50     """
51     # The epsilon is just here to avoid numerical instability (taking the log of 0)
52     log_one_minus_p_matrix = torch.log1p(-torch.exp(log_p_matrix) + epsilon)
53     return log_one_minus_p_matrix
54
55 def compute_log_pi_repeated(log_pi: torch.Tensor, n: int) -> torch.Tensor:
56     """
57     This function repeats log_pi n times along the first dimension
58
59     INPUT:
60     log_pi      : Log mixing proportions, Shape: (1, k)
61     n           : Number of repetitions
62
63     OUTPUT:
64     Repeated log_pi, Shape: (n, k)
65     """
66     return log_pi.expand(n, -1)
67
68 def compute_log_e_step(x: torch.Tensor, log_p_matrix: torch.Tensor, log_one_minus_p_matrix:
69     torch.Tensor, log_pi: torch.Tensor) -> torch.Tensor:
70     """
71     This function computes the log responsibilities for each data point and mixture
72     component
73
74     INPUT:
75     x                : Data matrix, Shape: (n, d)
76     log_p_matrix      : Log probabilities of Bernoulli parameters, Shape: (d, k)
77     log_one_minus_p_matrix : Log probabilities of (1 - p), Shape: (d, k)
78     log_pi            : Log mixing proportions, Shape: (1, k)
79
80     OUTPUT:
81     log_r_responsibilities : Log responsibilities (shape: (n, k)
82     """

```

```

80 # We compute the log-probability of each data point belonging to each component
81 log_component_p = x @ log_p_matrix + (1 - x) @ log_one_minus_p_matrix
82
83 # We add the log of mixing proportions to get the unnormalized log-responsibilities
84 log_r_unnormalised = log_component_p + log_pi
85
86 # We normalise the log-responsibilities across components to ensure they sum to 1
87 log_r_normaliser = torch.logsumexp(log_r_unnormalised, dim=1, keepdim=True)
88
89 # We subtract the normaliser to get the final log responsibilities
90 log_r_responsibilities = log_r_unnormalised - log_r_normaliser
91 return log_r_responsibilities
92
93 def compute_log_pi_hat(log_responsibility: torch.Tensor) -> torch.Tensor:
94     """
95     This function updates the log mixing proportions.
96
97     INPUT:
98     log_responsibility: Log responsibilities, Shape: (n, k)
99
100     OUTPUT:
101     log_pi_hat = Updated log_pi, Shape: (1, k)
102     """
103     # We compute the sum of the log-responsibilities across all data points for each
104     # component
105     # and normalise the log mixing proportions by subtracting log(n)
106     log_pi_hat = torch.logsumexp(log_responsibility, dim=0, keepdim=True) - torch.log(torch.
107         tensor([n], dtype=log_responsibility.dtype))
108     return log_pi_hat
109
110 def compute_log_p_matrix_hat(x: torch.Tensor, log_responsibility: torch.Tensor,
111     alpha_parameter: float, beta_parameter: float, epsilon: float = 1e-6,) -> torch.Tensor:
112     """
113     This function updates the log Bernoulli parameters.
114
115     INPUT:
116     x : Data matrix, Shape: (n, d)
117     log_responsibility : Log responsibilities, Shape: (n, k)
118     alpha_parameter : Alpha parameter of the Beta prior
119     beta_parameter : Beta parameter of the Beta prior
120     epsilon : Small constant used to avoid numerical instability
121
122     OUTPUT:
123     log_p_matrix_hat = Updated log_p_matrix, Shape: (d, k)
124     """
125     # We go from log space back to linear space
126     responsibility = torch.exp(log_responsibility)
127
128     # We compute the numerator
129     numerator = x.T @ responsibility + (alpha_parameter - 1)
130
131     # We compute the denominator
132     denominator = torch.sum(responsibility, dim=0) + (alpha_parameter + beta_parameter - 2)
133     denominator = denominator.unsqueeze(0)
134
135     # we compute the updated probability matrix
136     p_matrix_hat = numerator / denominator
137     # Again, we use epsilon to prevent taking log(0)
138     p_matrix_hat = torch.clamp(p_matrix_hat, min=epsilon, max=1 - epsilon)

```

```

136
137     log_p_matrix_hat = torch.log(p_matrix_hat)
138     return log_p_matrix_hat
139
140 def compute_log_likelihood(x: torch.Tensor, log_p_matrix: torch.Tensor,
141     log_one_minus_p_matrix: torch.Tensor, log_pi: torch.Tensor) -> float:
142     """
143     This function computes the total log-likelihood of the data given the model parameters
144
145     INPUT:
146     x                : Data matrix, Shape: (n, d)
147     log_p_matrix      : Log probabilities of p, Shape: (d, k)
148     log_one_minus_p_matrix : Log probabilities of (1 - p), Shape: (d, k)
149     log_pi            : Log mixing proportions, Shape: (1, k)
150
151     OUTPUT:
152     log_likelihood     : Total log-likelihood, Shape: (scalar)
153     """
154     # We calculate the log probabilities
155     log_component_probabilities = x @ log_p_matrix + (1 - x) @ log_one_minus_p_matrix
156     log_pi_repeated = compute_log_pi_repeated(log_pi, n)
157     log_probabilities = log_component_probabilities + log_pi_repeated
158
159     # We compute the log-likelihood for each data point by summing over components
160     log_likelihood = torch.sum(torch.logsumexp(log_probabilities, dim=1))
161     return log_likelihood.item()
162
163 def compute_log_prior(log_p_matrix: torch.Tensor, log_one_minus_p_matrix: torch.Tensor,
164     alpha_parameter: float, beta_parameter: float) -> float:
165     """
166     This function omputes the log-prior probability of the model parameters
167
168     INPUT:
169     log_p_matrix      : Log probabilities of p, Shape: (d, k)
170     log_one_minus_p_matrix : Log probabilities of (1 - p), Shape: (d, k)
171     alpha_parameter    : Alpha parameter of the Beta prior
172     beta_parameter     : Beta parameter of the Beta prior
173
174     OUTPUT:
175     log_prior         : Total log-prior, Shape: (scalar)
176     """
177     # We find the log of the beta function normalisation constant
178     beta_ln = betaln(alpha_parameter, beta_parameter)
179     # We calculate the log-prior
180     log_prior = beta_ln * log_p_matrix.numel() + torch.sum((alpha_parameter - 1) *
181         log_p_matrix + (beta_parameter - 1) * log_one_minus_p_matrix)
182     return log_prior.item()
183
184 def compute_unnormalised_log_posterior_likelihood(x: torch.Tensor, log_p_matrix: torch.
185     Tensor, log_one_minus_p_matrix: torch.Tensor, log_pi: torch.Tensor, alpha_parameter:
186     float, beta_parameter: float) -> float:
187     """
188     This function computes the unnormalised log-posterior probability of the model
189     parameters
190
191     INPUT:
192     x                : Data matrix, Shape: (n, d))
193     log_p_matrix      : Log probabilities of p, Shape: (d, k))
194     log_one_minus_p_matrix : Log probabilities of (1 - p), Shape: (d, k))

```



```

189     log_pi                : Log mixing proportions, Shape: (1, k))
190     alpha_parameter       : Alpha parameter of the Beta prior
191     beta_parameter        : Beta parameter of the Beta prior
192
193     OUTPUT:
194     log_posterior         : Unnormalized log-posterior, Shape: (scalar)
195     """
196     log_likelihood = compute_log_likelihood(x, log_p_matrix, log_one_minus_p_matrix, log_pi)
197     log_prior = compute_log_prior(log_p_matrix, log_one_minus_p_matrix, alpha_parameter,
198                                   beta_parameter)
199     log_posterior = log_likelihood + log_prior
200     return log_posterior
201
202 def run_expectation_maximisation(x: torch.Tensor, log_pi: torch.Tensor, log_p_matrix: torch.
203     Tensor, alpha_parameter: float, beta_parameter: float, max_number_of_steps: int, epsilon
204     : float) -> Tuple[torch.Tensor, torch.Tensor, torch.Tensor, List[float]]:
205     """
206     This is the main function that runs the EM algorithm
207
208     INPUT:
209     x                : Data matrix, Shape: (n, d))
210     log_pi           : Initial log mixing proportions, Shape: (1, k))
211     log_p_matrix      : Initial log Bernoulli parameters, Shape: (d, k))
212     alpha_parameter   : Alpha parameter of the Beta prior
213     beta_parameter    : Beta parameter of the Beta prior
214     max_number_of_steps : Maximum number of EM iterations
215     epsilon           : Small constant representing the convergence threshold
216
217     OUTPUT:
218     Updated log_pi, log_p_matrix, log_responsibility, and log_posteriors over iterations
219     """
220     log_posteriors = []
221     for step in range(max_number_of_steps):
222         log_one_minus_p_matrix = compute_log_one_minus_p_matrix(log_p_matrix)
223
224         # We perform the E-step
225         log_responsibility = compute_log_e_step(x, log_p_matrix, log_one_minus_p_matrix,
226                                                  log_pi)
227
228         # We perform the M-step
229         log_pi = compute_log_pi_hat(log_responsibility)
230         log_p_matrix = compute_log_p_matrix_hat(x, log_responsibility, alpha_parameter,
231                                                  beta_parameter)
232
233         # Recompute log_one_minus_p_matrix after updating log_p_matrix
234         log_one_minus_p_matrix = compute_log_one_minus_p_matrix(log_p_matrix)
235
236         # Compute log-posterior using the UPDATED log_one_minus_p_matrix
237         log_posterior = compute_unnormalised_log_posterior_likelihood(
238             x, log_p_matrix, log_one_minus_p_matrix, log_pi, alpha_parameter, beta_parameter
239         )
240         log_posteriors.append(log_posterior)
241
242         # As required, we use another constant epsilon to check for convergence and end the
243         # loop early
244         if step > 0 and abs(log_posteriors[-1] - log_posteriors[-2]) < epsilon:
245             print(f"Converged at iteration {step}")
246             break

```

```

242     return log_pi, log_p_matrix, log_responsibility, log_posteriors
243
244 def plot_log_posteriors(log_posteriors: List[float], title: str = "Log-posterior over EM
245     iterations") -> None:
246     """
247     This function plots the log-posteriors over EM iterations
248
249     INPUT:
250     log_posteriors : List of log-posterior values.
251     title          : Title of the plot.
252
253     OUTPUT:
254     Plot of the log-posterior over EM iterations
255     """
256     plt.figure(figsize=(8, 6))
257     plt.plot(range(1, len(log_posteriors) + 1), log_posteriors, marker='o')
258     plt.title(title)
259     plt.xlabel('EM iteration')
260     plt.ylabel('Log-Posterior')
261     plt.grid(True)
262     plt.show()
263
264 # We choose the inputs to the function
265 k = 3
266 max_number_of_steps = 100
267 data_np = np.loadtxt('binarydigits.txt')
268 x = torch.tensor(data_np, dtype=torch.float32)
269 n, d = x.shape
270 epsilon = 1e-4
271
272 # We initialise the parameters
273 log_pi, log_p_matrix = init_params(k, d)
274
275 # We set the prior parameters for the Beta distribution to 1 (this choice is explained in
276 # the submitted pdf file)
277 alpha_parameter = 1.0
278 beta_parameter = 1.0
279
280 # We run EM algorithm
281 log_pi, log_p_matrix, log_responsibility, log_posteriors = run_expectation_maximisation(
282     x,
283     log_pi,
284     log_p_matrix,
285     alpha_parameter,
286     beta_parameter,
287     max_number_of_steps,
288     epsilon,
289 )
290
291 plot_log_posteriors(log_posteriors)
292
293 #### Code to display the Log-posterior against the EM iteration for different values of K and
294 # the parameters ####
295
296 # We initialise the values of K to test and display
297 K_values = [2, 3, 4, 7, 10]
298
299 # Store log-likelihoods and parameters for each K
300 log_likelihoods_per_k = {}

```

```

298 parameters_per_k = {}
299
300 # Run the EM algorithm for each K
301 for k in K_values:
302     print(f"K = {k}")
303     log_pi, log_p_matrix = init_params(k, d)
304
305     # We run the EM algorithm
306     log_pi, log_p_matrix, log_responsibility, log_posteriors = run_expectation_maximisation(
307         x, log_pi, log_p_matrix, alpha_parameter, beta_parameter, max_number_of_steps,
308         epsilon=epsilon)
309
310     # Store the log-likelihoods and parameters
311     log_likelihoods_per_k[k] = log_posteriors
312     parameters_per_k[k] = {'log_pi': log_pi.detach(), 'log_p_matrix': log_p_matrix.detach()}
313
314     # Plot the log-likelihood over iterations
315     plot_log_posteriors(log_posteriors)
316
317     # Display the parameters found
318     print(f"Parameters for K = {k}:")
319     print("Mixing proportions (pi):")
320     print(torch.exp(log_pi))
321     print("Bernoulli parameters (p_matrix):")
322     print(torch.exp(log_p_matrix))

```

Listing 4: Python code running the EM algorithm for a mixture of K multivariate Bernoullis

```

1 def visualise_p_matrices(log_p_matrices: List[np.ndarray], log_pis: List[np.ndarray], ks:
2     List[int], figure_title: str, figure_path: str) -> None:
3     """
4     This function lets us visualise the P matrices contained in the parameters_per_k
5     variable
6
7     INPUT:
8     log_p_matrices : List of log probability matrice, Shape: (d, k)
9     log_pis        : List of log mixing proportions, each array has shape (1, k)
10    ks              : List of k values for each P matrix (2, 3, 4, 7 10)
11    figure_title    : Title for the figures
12    figure_path     : File path to save the figures
13
14    OUTPUT:
15    The visual representation of the P matrices
16    """
17    # We find the number of P matrices to show and the maximum k value
18    n = len(ks)
19    m = np.max(ks)
20
21    # We set up the plot dimensions
22    fig = plt.figure()
23    fig.set_figwidth(15)
24    fig.set_figheight(10)
25
26    for i, (log_p_matrix, log_pi, k) in enumerate(zip(log_p_matrices, log_pis, ks)):
27        # We calculate P matrix and mixing proportions
28        p_matrix = np.exp(log_p_matrix).reshape(8, 8, k)
29        pi_values = np.exp(log_pi)
30
31        for j in range(k):

```

```

30         ax = plt.subplot(n, m, m * i + j + 1)
31         ax.imshow(p_matrix[:, :, j], cmap="inferno", interpolation="none")
32         ax.tick_params(axis="both", which="both", bottom=False, left=False)
33         ax.xaxis.set_ticklabels([])
34         ax.yaxis.set_ticklabels([])
35         ax.set_title(f"pi_{j}: {np.round(pi_values[0, j], 2)}")
36         if j == 0:
37             ax.set_ylabel(f"k={k}")
38
39     fig.suptitle(figure_title)
40     plt.savefig(figure_path)
41     plt.show()
42
43 ks = [2, 3, 4, 7, 10]
44
45 # We display as 8x8 images the Bernoulli parameters
46 for i in ks:
47     log_pi = np.array(parameters_per_k[i][ 'log_pi' ])
48     log_p_matrix = np.array(parameters_per_k[i][ 'log_p_matrix' ])
49     ks = [i]
50     log_p_matrices = [log_p_matrix]
51     log_pis = [log_pi]
52     figure_title = "test"
53     figure_path = f"path_{i}.png"
54     visualise_p_matrices(log_p_matrices, log_pis, ks, figure_title, figure_path)
55
56 ks = [2, 3, 4, 7, 10]
57 # We display the pi_k in the linear space
58 for i in ks:
59     print(torch.exp(parameters_per_k[i][ 'log_pi' ]))

```

Listing 5: Python code used to display the parameters

### 3.5 Question 3e

I ran the algorithm 3 times (test runs 1, 2 and 3) with randomly chosen initial conditions, and the resulting parameters are shown below in tables 4 to 7.

<b>K</b>	$\pi_k$ on run 1	$\pi_k$ on run 2	$\pi_k$ on run 3
2	0.5804, 0.4196	0.4099, 0.5901	0.7600, 0.2400
3	0.3399, 0.2200, 0.4401	0.4002, 0.2998, 0.2999	0.6900, 0.1800, 0.1300
4	0.1500, 0.4103, 0.1300, 0.3097	0.1615, 0.3193, 0.2200, 0.2992	0.1000, 0.3700, 0.2900, 0.2400
7	0.2200, 0.1600, 0.0800, 0.0300, 0.1400, 0.2500, 0.1200	0.0501, 0.1600, 0.1299, 0.2600, 0.1500, 0.0300, 0.2200	0.0200, 0.0700, 0.2600, 0.1806, 0.1100, 0.3100, 0.0494
10	0.0100, 0.0700, 0.2001, 0.1190, 0.1600, 0.2099, 0.0300, 0.0100, 0.1510, 0.0400	0.0700, 0.0400, 0.0900, 0.1685, 0.0900, 0.1800, 0.0900, 0.0800, 0.1109, 0.0806	0.1200, 0.0300, 0.0200, 0.0600, 0.1300, 0.0400, 0.0600, 0.2700, 0.1100, 0.1600

Table 4: Table showing the values of  $\pi_k$  for different K on the 3 test runs

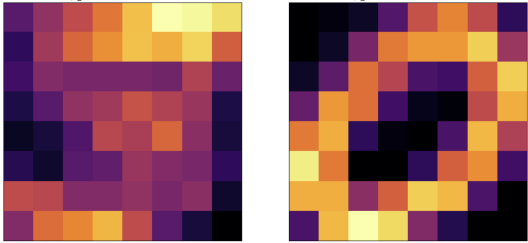

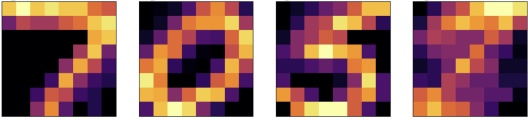

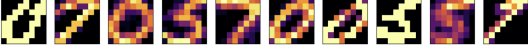
<b>K</b>	<b>Optimised Bernoulli parameters</b>	
2		
3		
4		
7		
10		

Table 5: Table showing the optimised Bernoulli parameters for the first test run

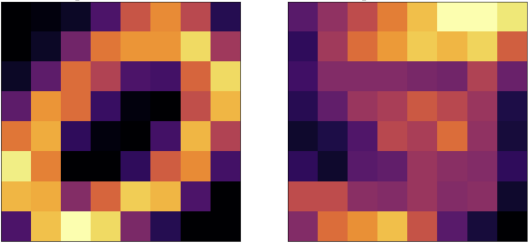

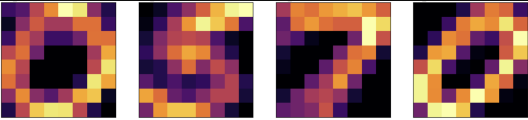


K	Optimised Bernoulli parameters	
2		
3		
4		
7		
10		

Table 6: Table showing the optimised Bernoulli parameters for the second test run

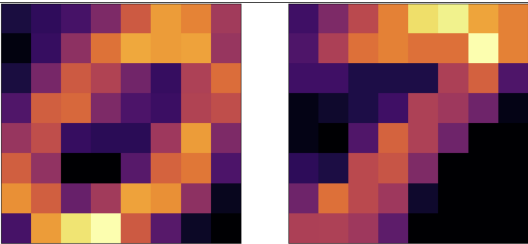
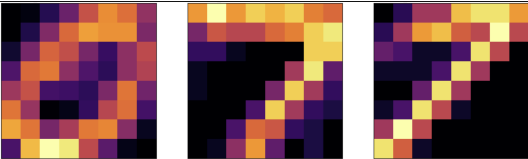
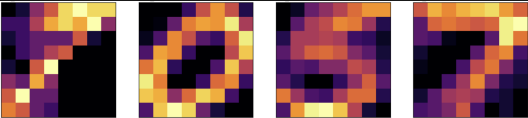


K	Optimised Bernoulli parameters	
2		
3		
4		
7		
10		

Table 7: Table showing the optimised Bernoulli parameters for the third test run

Looking at the results above, we see that there is a similarity to the solutions found. For example, in the solution obtained in question 3d and the first two test runs, we find the same patterns. When  $K = 3$ , all three runs show the numbers 5, 7, 0 in the same order. A similar pattern can be found for  $K = 2$ , where the numbers 5, 0 appear every time, up to a permutation since test run 2 shows 0, 5 instead. We see that as the value of  $K$  gets higher, we lose this similarity. For example, comparing the Bernoulli parameters for test runs 1 and 3 when  $K = 10$ , they have different numbers of 0, 5 and 7s and the numbers are in very different orders. The similarity in the solutions does seem to depend on  $K$ .

The algorithm works well in general. From tables 5 to 7, we see that it correctly finds clusters of the numbers 0, 5 and 7. However, its performance is very dependent on the value of  $K$ . If we look at the  $K = 2$  case of our test runs, we see that the algorithm finds two clusters, one which is clearly identifiable as the 0, and the other seems to be a mix of the 5 and 7, resulting in the image seeming blurry. This makes sense, as the three numbers would require one cluster each, so we would need at least  $K = 3$ . Now looking at the other end of the spectrum, when  $K = 10$  we see that some of the images look to be only 2 colours. This is because the probability for this cluster is really small (e.g.  $\pi_1 = 0.01$  for the test run 1), which results in loosely defined clusters. The probability is must better distributed among the clusters for  $K = 3$ , where each  $\pi_k$  is closer to 30%. With this in mind, one way to improve the model is to change our belief on what the ideal number of clusters is, which we've seen empirically is at least 3. If we could incorporate prior knowledge of the data into the Dirichlet and Beta distributions we defined (such as the fact that we don't like clusters with very low probability), we could improve the model.

### 3.6 Bonus: Question 3f

Our log likelihood are calculated using the natural logarithm, which is of base  $e$ . To encode the log-likelihood in bits, we need to convert from (base  $e$ ) to base 2, which we can do this by dividing the log-likelihood by  $\log(2)$ :

$$\text{Log-likelihood in bits} = \frac{\text{log-likelihood}}{\log(2)}$$

The naive way to encode the binary data is to use the 0s and 1s that correspond to the pixels as the encoding. Each binary image has  $D$  pixel, so for a dataset with  $N$  binary images the length of the encoding would be:

$$\text{Naïve encoding length} = N \times D \text{ bits}$$

Gzip is a compression algorithm that finds patterns in data to reduce its storage size. If the log likelihood in bits is lower than the gzip encoding length, this shows that the EM algorithm capture the structure of the data well and is a better compression algorithm. If gzip works better than the EM model, it could be because the EM algorithm doesn't capture all the patterns contained the binarydigits.txt.

A difference between the two could come about because of the difference in approach by both methods. The EM model tries to understand the underlying probability distribution, while gzip compresses regardless of the probabilistic background. This can make for better compression if the model that we've made can't represent all the patterns. This would be the case for low values of  $K$ . For example, if we compare the  $K = 2$  and  $K = 4$  Bernoulli parameters for the second test run (in table 6), we see that the  $K = 4$  parameters capture the difference between a round 0 and a slanted one, which we don't see in  $K = 2$ .

### 3.7 Bonus: Question 3g

We must consider both the costs of the data and the parameters, which makes the total cost:

$$\text{Total cost} = \text{Cost of the model parameters} + \text{Cost of the data given a certain model}$$

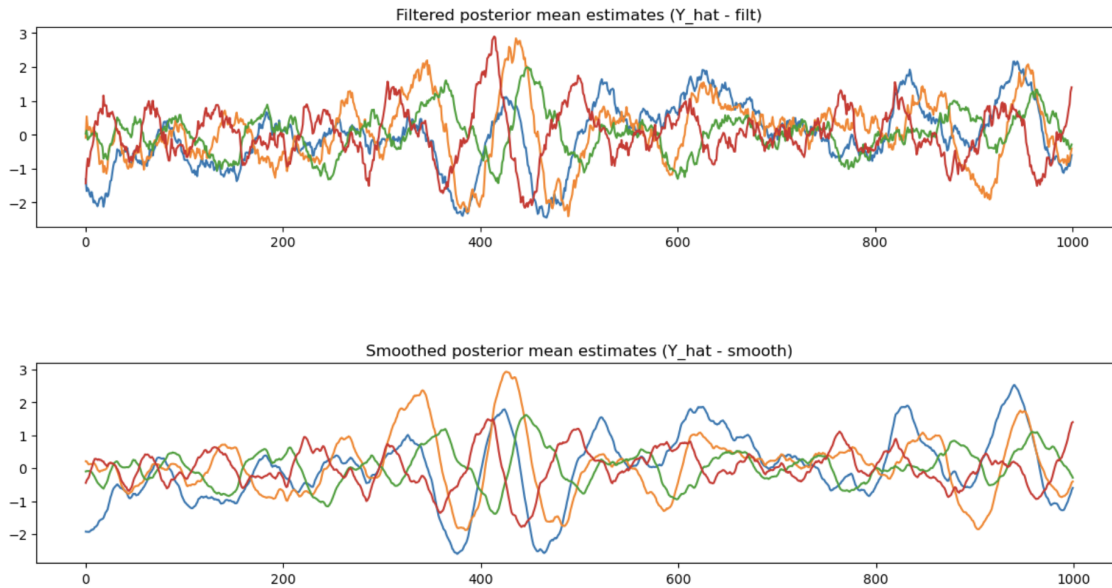
This sets our model apart from the gzip algorithm, which doesn't need to take into account parameter cost since it compresses regardless of the probabilistic background. As  $K$  increases, our model becomes more complex and requires more parameters, which would increase the total cost of storing the model. However, if  $K$  is small the model might not understand all patterns in the data as discussed in the previous question, which would result in a higher cost of the data. Therefore they must be an optimal  $K$  that minimises the total cost.



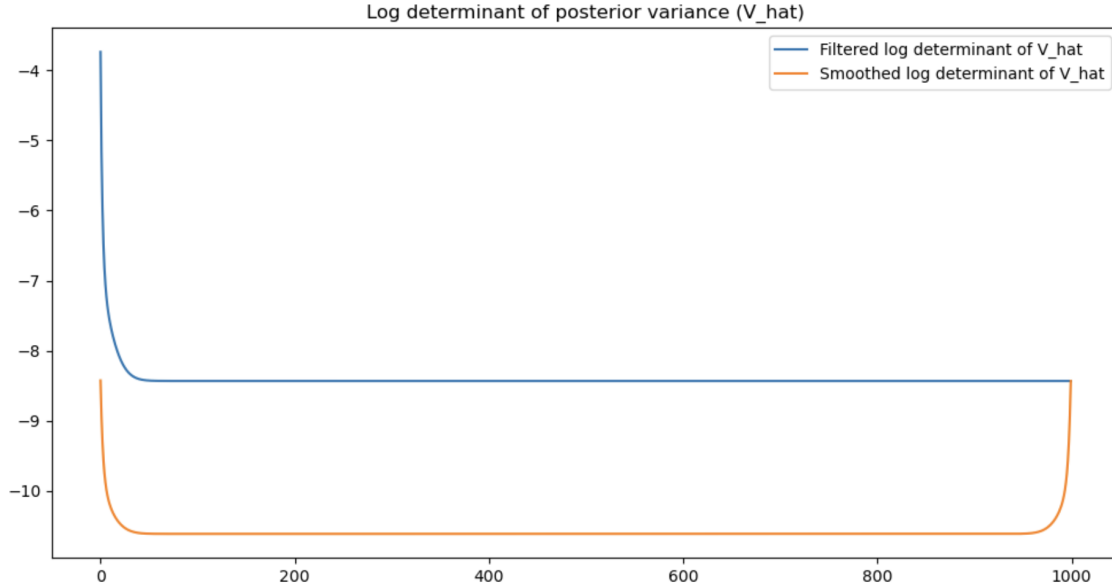
## 4 Bonus: LGSSMs, EM and SSID

### 4.1 Question 4a

The figures below show the Kalman filter and Kalman smoother respectively.



The behaviour of  $Y$  on both plots is similar, as we can see that both plots have similar shapes, however they differ in their uncertainty. This is linked to how both were computed: the filtered estimates were generated using only the data up to each time step. Each estimate depends only on the information available up to that time, resulting in higher uncertainty. On the other hand, the smoothed estimates use data from the entire time series, which makes more accurate state estimates. The log determinant of posterior variances for the Kalman filter and smoother were also calculated and displayed in the figure below.



The log determinant of the posterior variance is a measure of the uncertainty in the state estimates. We see from the figure above that the log determinant for the filtered estimates is higher than the one for the smoothed estimates, which is what we would expect. They have similar shapes near the start of the time series as the algorithm has little information to work with. The Kalman filter only has the initial state estimate and the initial observation, and the Kalman smoother starts lower because it has access to future observations. The two plots differ at the end, where the Kalman smoother uncertainty increases as it can no longer look at the future observations, which makes it lose confidence in its estimates, and join the Kalman filter (as both can only see the past observations now). The code used to compute and display the figures in this question can be found in Listing 6:

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def run_ssm_kalman(X, y_init, Q_init, A, Q, C, R, mode='smooth'):
5     """
6     Calculates kalman-smoother estimates of SSM state posterior.
7     :param X:      data, [d, t_max] numpy array
8     :param y_init:  initial latent state, [k,] numpy array
9     :param Q_init:  initial variance, [k, k] numpy array
10    :param A:       latent dynamics matrix, [k, k] numpy array
11    :param Q:       innovations covariance matrix, [k, k] numpy array
12    :param C:       output loading matrix, [d, k] numpy array
13    :param R:       output noise matrix, [d, d] numpy array
14    :param mode:    'forw' or 'filt' for forward filtering, 'smooth' for also backward
                    filtering
15    :return:
16    y_hat:         posterior mean estimates, [k, t_max] numpy array
17    V_hat:         posterior variances on y_t, [t_max, k, k] numpy array
18    V_joint:       posterior covariances between y_{t+1}, y_t, [t_max, k, k] numpy array
19    likelihood:    conditional log-likelihoods log(p(x_t|x_{1:t-1})), [t_max,] numpy array

```

```

20     """
21     d, k = C.shape
22     t_max = X.shape[1]
23
24     # dimension checks
25     assert np.all(X.shape == (d, t_max)), "Shape of X must be (%d, %d), %s provided" % (d,
26         t_max, X.shape)
27     assert np.all(y_init.shape == (k,)), "Shape of y_init must be (%d,), %s provided" % (k,
28         y_init.shape)
29     assert np.all(Q_init.shape == (k, k)), "Shape of Q_init must be (%d, %d), %s provided" %
30         (k, k, Q_init.shape)
31     assert np.all(A.shape == (k, k)), "Shape of A must be (%d, %d), %s provided" % (k, k, A.
32         shape)
33     assert np.all(Q.shape == (k, k)), "Shape of Q must be (%d, %d), %s provided" % (k, k, Q.
34         shape)
35     assert np.all(C.shape == (d, k)), "Shape of C must be (%d, %d), %s provided" % (d, k, C.
36         shape)
37     assert np.all(R.shape == (d, d)), "Shape of R must be (%d, %d), %s provided" % (d, k, R.
38         shape)
39
40     y_filt = np.zeros((k, t_max)) # filtering estimate:  $\hat{y}_t$ 
41     V_filt = np.zeros((t_max, k, k)) # filtering variance:  $\hat{V}_t$ 
42     y_hat = np.zeros((k, t_max)) # smoothing estimate:  $\hat{y}_T$ 
43     V_hat = np.zeros((t_max, k, k)) # smoothing variance:  $\hat{V}_T$ 
44     K = np.zeros((t_max, k, X.shape[0])) # Kalman gain
45     J = np.zeros((t_max, k, k)) # smoothing gain
46     likelihood = np.zeros(t_max) # conditional log-likelihood:  $p(x_t|x_{1:t-1})$ 
47
48     I_k = np.eye(k)
49
50     # forward pass
51
52     V_pred = Q_init
53     y_pred = y_init
54
55     for t in range(t_max):
56         x_pred_err = X[:, t] - C.dot(y_pred)
57         V_x_pred = C.dot(V_pred.dot(C.T)) + R
58         V_x_pred_inv = np.linalg.inv(V_x_pred)
59         likelihood[t] = -0.5 * (np.linalg.slogdet(2 * np.pi * (V_x_pred))[1] +
60             x_pred_err.T.dot(V_x_pred_inv).dot(x_pred_err))
61
62         K[t] = V_pred.dot(C.T).dot(V_x_pred_inv)
63
64         y_filt[:, t] = y_pred + K[t].dot(x_pred_err)
65         V_filt[t] = V_pred - K[t].dot(C).dot(V_pred)
66
67         # symmetrise the variance to avoid numerical drift
68         V_filt[t] = (V_filt[t] + V_filt[t].T) / 2.0
69
70         y_pred = A.dot(y_filt[:, t])
71         V_pred = A.dot(V_filt[t]).dot(A.T) + Q
72
73     # backward pass
74
75     if mode == 'filt' or mode == 'forw':
76         # skip if filtering/forward pass only
77         y_hat = y_filt
78         V_hat = V_filt

```

```

72     V_joint = None
73     else:
74         V_joint = np.zeros_like(V_filt)
75         y_hat[:, -1] = y_filt[:, -1]
76         V_hat[-1] = V_filt[-1]
77
78         for t in range(t_max - 2, -1, -1):
79             J[t] = V_filt[t].dot(A.T).dot(np.linalg.inv(A.dot(V_filt[t]).dot(A.T) + Q))
80             y_hat[:, t] = y_filt[:, t] + J[t].dot((y_hat[:, t + 1] - A.dot(y_filt[:, t])))
81             V_hat[t] = V_filt[t] + J[t].dot(V_hat[t + 1] - A.dot(V_filt[t]).dot(A.T) - Q).
                dot(J[t].T)
82
83         V_joint[-2] = (I_k - K[-1].dot(C)).dot(A).dot(V_filt[-2])
84
85         for t in range(t_max - 3, -1, -1):
86             V_joint[t] = V_filt[t + 1].dot(J[t].T) + J[t + 1].dot(V_joint[t + 1] - A.dot(
                V_filt[t + 1])).dot(J[t].T)
87
88     return y_hat, V_hat, V_joint, likelihood
89
90 # We start by defining the matrices as given in the question
91 A = 0.99 * np.array([
92     [np.cos(2 * np.pi / 180), -np.sin(2 * np.pi / 180), 0, 0],
93     [np.sin(2 * np.pi / 180), np.cos(2 * np.pi / 180), 0, 0],
94     [0, 0, np.cos(2 * np.pi / 90), -np.sin(2 * np.pi / 90)],
95     [0, 0, np.sin(2 * np.pi / 90), np.cos(2 * np.pi / 90)]
96 ])
97
98 Q = np.eye(4) - A.dot(A.T)
99 C = np.array([
100     [1, 0, 1, 0],
101     [0, 1, 0, 1],
102     [1, 0, 0, 1],
103     [0, 1, 0, 1],
104     [0.5, 0.5, 0.5, 0.5]
105 ])
106 R = np.eye(5)
107
108 # We load data and transpose X such that we have Shape: (d, T)
109 X = np.loadtxt('ssm_spins.txt').T
110
111 # We initialise the state and covariance
112 y_init = np.zeros(4)
113 Q_init = np.eye(4)
114
115 # We call the function contained in the ssm_kalman.py for filtering
116 y_hat_filt, V_hat_filt, V_joint_filt, likelihood_filt = run_ssm_kalman(X, y_init, Q_init, A,
    Q, C, R, mode='filt')
117
118 # We call the function contained in the ssm_kalman.py for smoothing
119 y_hat_smooth, V_hat_smooth, V_joint_smooth, likelihood_smooth = run_ssm_kalman(X, y_init,
    Q_init, A, Q, C, R, mode='smooth')
120
121 # We plotting y_hat and estimates for filtering and smoothing
122 plt.figure(figsize=(12, 8))
123 plt.subplot(3, 1, 2)
124 plt.plot(y_hat_filt.T)
125 plt.title("Filtered posterior mean estimates (Y_hat - filt)")
126

```

```

127 plt.subplot(3, 1, 3)
128 plt.plot(y_hat_smooth.T)
129 plt.title("Smoothed posterior mean estimates (Y_hat - smooth)")
130 plt.tight_layout()
131 plt.show()
132
133 # We plot log determinant of V_hat to show the uncertainty
134 logdet_filt = [np.linalg.slogdet(V)[1] for V in V_hat_filt]
135 logdet_smooth = [np.linalg.slogdet(V)[1] for V in V_hat_smooth]
136
137 plt.figure(figsize=(12, 6))
138 plt.plot(logdet_filt, label='Filtered log determinant of V_hat')
139 plt.plot(logdet_smooth, label='Smoothed log determinant of V_hat')
140 plt.title("Log determinant of posterior variance (V_hat)")
141 plt.legend()
142 plt.show()

```

Listing 6: Python code used to answer question 4a

# 5 Decrypting Messages with MCMC

## 5.1 Question 5a

The formula for the ML estimates of the transition probabilities  $p(s_i = \alpha \mid s_{i-1} = \beta) \equiv \psi(\alpha, \beta)$  are simply the number of times the pair  $(\alpha, \beta)$  occurs over the total number occurrences of  $\beta$ :

$$\psi(\alpha, \beta) = \frac{N_{(\alpha, \beta)}}{N_{\beta, \text{total}}}$$

Below is the table showing the transition probabilities  $\psi(\alpha, \beta)$  for all the characters. The table is too big for the text to be readable on the page. I suggest either zooming in on the computer or running the code given at the end of the question.

	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	z
1	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
2	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
3	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
4	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
5	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
6	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
7	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
8	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
9	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
10	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
11	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
12	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
13	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
14	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
15	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
16	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
17	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
18	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
19	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
20	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	0.0300	
21	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	0.0300	
22	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	0.0300	
23	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	0.0300	
24	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	0.0300	
25	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0400	
26	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
27	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
28	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
29	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
30	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
31	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
32	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
33	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
34	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
35	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
36	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
37	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
38	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
39	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
40	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	
41	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.0300	0.					

Below is the table showing the invariant probability distribution. The probabilities are given in descending order.

Symbol	Probability	Symbol	Probability
space	0.165211	x	0.001303
e	0.100472	!	0.001260
t	0.072515	?	0.001006
a	0.064910	j	0.000825
o	0.060963	-	0.000586
n	0.059082	;	0.000367
i	0.055083	:	0.000315
h	0.053595	z	0.000766
s	0.052255	q	0.000748
r	0.047498	(	0.000214
d	0.037947	)	0.000214
l	0.030960	1	0.000114
u	0.020655	*	0.000090
c	0.019650	8	0.000056
m	0.019776	2	0.000042
w	0.018999	0	0.000049
f	0.017609	5	0.000015
g	0.016464	6	0.000016
y	0.014834	3	0.000018
p	0.014485	9	0.000010
,	0.012797	7	0.000012
b	0.011110	,	0.000002
.	0.009906	"	0.000007
v	0.008630	=	0.000001
k	0.006550	/	0.000003

Table 8: Table showing the invariant probability distribution

The code use to load the text, learn the transition statistics, estimate the transition probabilities, the stationary distribution, and display both tables can be found in Listing 7:

```

1 import numpy as np
2 from collections import Counter, defaultdict
3 import pandas as pd
4
5 # We start by loading War and Peace, the symbols list and the encrypted message using the
  # functions below:
6
7 def load_symbols(file_path):

```

```

8     with open(file_path, 'r', encoding='utf-8') as file:
9         symbols = []
10        for line in file:
11            symbol = line.rstrip() # Use rstrip() to remove newline but not internal spaces
12            # I've had to add the line below because I couldn't get the space to appear when
13            # loading the symbols
14            if symbol == " ":
15                symbols.append(" ")
16            else:
17                symbols.append(symbol)
18        return symbols
19
20 def load_encrypted_message(file_path):
21     with open(file_path, 'r', encoding='utf-8') as file:
22         encrypted_message = file.read().strip()
23     return encrypted_message
24
25 def load_war_and_peace(file_path):
26     with open(file_path, 'r', encoding='utf-8') as file:
27         text = file.read()
28         # I've truncated the text such that the text starts at the start of the book
29         # I've made sure only symbols that exist in the symbols.txt file remained in the
30         text
31         text = text[7473:].lower()
32         filtered_text = ''.join(char.lower() for char in text if char.lower() in symbols)
33     return filtered_text
34
35 def calculate_transition_statistics(text):
36     # We count occurrences of each symbol and each symbol pair
37     single_counts = Counter(text)
38     pair_counts = Counter(zip(text, text[1:]))
39
40     # We sum to find the total symbols count
41     total_symbols = sum(single_counts.values())
42
43     # We calculate the stationary distribution
44     stationary_dist = {char: count / total_symbols for char, count in single_counts.items()}
45
46     # We calculate transition probabilities
47     transition_probs = defaultdict(dict)
48     for (char1, char2), count in pair_counts.items():
49         transition_probs[char1][char2] = count / single_counts[char1]
50
51     return stationary_dist, transition_probs
52
53 # We load the three files
54 symbols = load_symbols('symbols.txt')
55 encrypted_message = load_encrypted_message('message.txt')
56 war_and_peace_text = load_war_and_peace('war_and_peace.txt')
57
58 # We get the invariant distribution and the transition probabilities
59 probabilities = calculate_transition_statistics(war_and_peace_text)
60 invariant_distrib = probabilities[0]
61 transition_prob = probabilities[1]
62
63 # We creating stationary distribution table
64 invariant_df = pd.DataFrame(list(invariant_distrib.items()), columns=['Symbol', 'Probability'])

```



```

64 # We create an NxN transition probability matrix DataFrame from the transition_prob
    dictionary
65 transition_df = pd.DataFrame(transition_prob).fillna(0)
66
67 # Get the table into Latex
68 transition_df.to_latex("transition_table.tex", index=True, float_format="%.2e")

```

Listing 7: Python code used to answer question 5a

## 5.2 Question 5b

Since we are assuming that we have a uniform prior distribution over all possible permutations, the first choice in the permutation is chosen independently of the others. However, as soon as a particular permutation  $\sigma$  is selected, each mapping  $\sigma(s)$  is dependent on the other mappings in that permutation (since each symbol (original text) must map to a unique symbol (encrypted text)). So the latent variables are not independent.

The joint probability of observing an encrypted sequence  $e_1 e_2 \cdots e_n$  given a specific permutation  $\sigma$  is determined by the probability of the next character given the previous one. We know that the first one we choose can be modelled as sampled from the stationary distribution, and that the rest will be the conditional probabilities of each character based on the preceding character. We can express this as:

$$p(e_1 e_2 \cdots e_n | \sigma) = p(e_1 | \sigma) \prod_{i=2}^n p(e_i | e_{i-1}, \sigma)$$

Now using the transition probability  $\psi(\alpha, \beta)$  for transitions from symbol  $\beta$  to  $\alpha$  (given by the Markov model), we can rewrite this as:

$$p(e_1 e_2 \cdots e_n | \sigma) = p(\sigma^{-1}(e_1)) \prod_{i=2}^n \psi(\sigma^{-1}(e_i), \sigma^{-1}(e_{i-1}))$$

where we have used the  $\sigma^{-1}$  to denote the act of decrypting.

## 5.3 Question 5c

The proposal probability  $S(\sigma \rightarrow \sigma')$  is the probability of proposing a move from the current permutation  $\sigma$  to a new permutation  $\sigma'$ . We would do this by choosing two symbols  $s$  and  $s'$  randomly and swapping their encrypted symbols  $\sigma(s)$  and  $\sigma(s')$ . Since they are chosen randomly from 53 different symbols, each pair  $(s, s')$  has an equal chance of being selected and there are  $\binom{53}{2} = \frac{53(53-1)}{2}$  possible ways to choose them. Thus, the probability  $\frac{2}{53(53-1)} = 7.26 \times$

$10^{-4}$ . Therefore, the proposal probability  $S(\sigma \rightarrow \sigma')$  doesn't depend on the permutations  $\sigma$  and  $\sigma'$ , only on the symbols  $(s, s')$  being chosen randomly.

In the Metropolis-Hastings algorithm, the acceptance probability kernel is given by:

$$A(x_{i+1}|x_i) := \min \left\{ 1, \frac{q(x_i|x_{i+1})p(x_{i+1})}{q(x_{i+1}|x_i)p(x_i)} \right\}$$

Now in our case, for a proposed move from  $\sigma$  to  $\sigma'$  it becomes:

$$A(\sigma \rightarrow \sigma') = \min \left( 1, \frac{p(e_1 e_2 \cdots e_n | \sigma') S(\sigma' \rightarrow \sigma)}{p(e_1 e_2 \cdots e_n | \sigma) S(\sigma \rightarrow \sigma')} \right)$$

Since  $S(\sigma \rightarrow \sigma') = S(\sigma' \rightarrow \sigma)$  in this symmetric proposal setup, the acceptance probability simplifies to:

$$A(\sigma \rightarrow \sigma') = \min \left( 1, \frac{p(e_1 e_2 \cdots e_n | \sigma')}{p(e_1 e_2 \cdots e_n | \sigma)} \right)$$

## 5.4 Question 5d

The implemented Metropolis-Hastings sampler is shown in Listing 8.

```

1 from collections import Counter, defaultdict
2 import random
3 import numpy as np
4
5 def load_symbols(file_path):
6     """
7     This function reads the symbols.txt file and loads it in Python as a list
8
9     INPUT:
10    file_path    : File path of the symbols.txt document
11
12    OUTPUT:
13    symbols      : A list of strings containing all the symbols in the symbols.txt file
14    """
15    with open(file_path, 'r', encoding='utf-8') as file:
16        symbols = []
17        for line in file:
18            symbol = line.rstrip()
19            # I've had to add the line below because I couldn't get the space to appear when
20            # loading the symbols
21            if symbol == " ":
22                symbols.append(" ")
23            else:
24                symbols.append(symbol)
25    return symbols
26
27 def load_encrypted_message(file_path):
28     """
29     This function loads the encrypted message into Python as a str

```

```

29
30 INPUT:
31 file_path          : File path of the message.txt document
32
33 OUTPUT:
34 encrypted_message   : A string containing the encrypted message
35 """
36 with open(file_path, 'r', encoding='utf-8') as file:
37     encrypted_message = file.read().strip()
38 return encrypted_message
39
40 def load_war_and_peace(file_path):
41     """
42     This function reads the war_and_peace.txt file and outputs
43
44     INPUT:
45     file_path          : File path of the war_and_peace.txt document
46
47     OUTPUT:
48     filtered_text      : The War and Peace processed text, to include only lower case and
49                          symbols contained in the symbols.txt file
50     """
51     with open(file_path, 'r', encoding='utf-8') as file:
52         text = file.read()
53         # I've truncated the text to start at the beginning of the book
54         text = text[7473:].lower()
55         # We filter out characters not in the symbols list
56         filtered_text = ''.join(char.lower() for char in text if char.lower() in symbols)
57     return filtered_text
58
59 def calculate_transition_statistics(text):
60     """
61     This function calculates the stationary distribution and the transition probabilities
62
63     INPUT:
64     text              : The large corpus of English text (in our case, the
65                        war_and_peace_text variable)
66
67     OUTPUT:
68     stationary_dist    : A dictionary containing the probabilities of the stationary
69                          distribution
69     transition_probs    : A dictionary containing all the transition probabilities
70     """
71     # We count occurrences of each symbol and each symbol pair
72     single_counts = Counter(text)
73     pair_counts = Counter(zip(text, text[1:]))
74
75     # We sum to find the total symbols count
76     total_symbols = sum(single_counts.values())
77
78     # We calculate the stationary distribution
79     stationary_dist = {char: count / total_symbols for char, count in single_counts.items()}
80
81     # We calculate transition probabilities
82     transition_probs = defaultdict(dict)
83     for (char1, char2), count in pair_counts.items():
84         transition_probs[char1][char2] = count / single_counts[char1]
85
86     return stationary_dist, transition_probs

```

```

85
86
87 def generate_initial_permutation_based_on_frequencies(encrypted_message):
88     """
89     This function generates an initial decryption key
90
91     INPUT:
92     encrypted_message      : The encrypted message contained in message.txt
93
94     OUTPUT:
95     initial_permutation    : The initial permutation
96     """
97     # We generate the initial permutation by matching frequencies between the encrypted and
98     # training texts
99     # We count the frequency of each symbol in each text
100    encrypted_counts = Counter(encrypted_message)
101    encrypted_symbols = list(encrypted_counts.keys())
102
103    war_and_peace_counts = Counter(war_and_peace_text)
104    war_and_peace_symbols = list(war_and_peace_counts.keys())
105
106    # We sort the symbols by frequency
107    encrypted_symbols_sorted = [symbol for symbol, count in encrypted_counts.most_common()]
108    war_and_peace_symbols_sorted = [symbol for symbol, count in war_and_peace_counts.
109                                   most_common()]
110
111    # We map the encrypted symbols to most frequent symbols in the training text for our
112    # initial permutation
113    initial_permutation = {}
114    for enc_sym, dec_sym in zip(encrypted_symbols_sorted, war_and_peace_symbols_sorted):
115        initial_permutation[enc_sym] = dec_sym
116
117    # The method above might not map every symbol, since they might not be present in the
118    # encrypted text
119    # The code below handles any remaining symbols by assigning them arbitrarily from the
120    # remaining symbols
121    all_symbols = set(symbols)
122    unmapped_symbols = all_symbols - set(initial_permutation.values())
123    unmapped_encrypted_symbols = set(encrypted_message) - set(initial_permutation.keys())
124    for enc_sym, dec_sym in zip(unmapped_encrypted_symbols, unmapped_symbols):
125        initial_permutation[enc_sym] = dec_sym
126
127    return initial_permutation
128
129 def decrypt_text(encrypted_message, permutation):
130     """
131     This function decrypts the encrypted message using the provided permutation
132
133     INPUT:
134     encrypted_message      : The encrypted message contained in message.txt
135     permutation            : The current permutation
136
137     OUTPUT:
138     decrypted_message      : The decrypted message
139     """
140    decrypted_message = ''.join([permutation.get(char, char) for char in encrypted_message])
141    return decrypted_message
142
143 def calculate_posterior_probability(encrypted_message, permutation, transition_probs,

```

```

stationary_dist):
139     """
140     This function calculates the log posterior probability of the decrypted message given
        the current permutation
141
142     INPUT:
143     encrypted_message    : The encrypted message contained in message.txt
144     permutation          : The current permutation
145     transition_probs     : The transition probability dictionary
146     stationary_dist      : The stationary distribution probability dictionary
147
148     OUTPUT:
149     log_prob             : The log posterior probability of the decrypted message
150     """
151     # We decrypt the message using the function defined previously
152     decrypted_message = decrypt_text(encrypted_message, permutation)
153     log_prob = 0.0
154
155     # As suggested in the problem, we add the log probability of the first character from
        the stationary distribution
156     first_char = decrypted_message[0]
157     if first_char in stationary_dist and stationary_dist[first_char] > 0:
158         log_prob += np.log(stationary_dist[first_char])
159     else:
160         # We add a small probability for any unseen characters
161         log_prob += np.log(1e-10)
162
163     # We compute the log probabilities of transitions
164     for i in range(1, len(decrypted_message)):
165         prev_char = decrypted_message[i - 1]
166         curr_char = decrypted_message[i]
167         if prev_char in transition_probs and curr_char in transition_probs[prev_char]:
168             prob = transition_probs[prev_char][curr_char]
169             if prob > 0:
170                 log_prob += np.log(prob)
171             else:
172                 log_prob += np.log(1e-10)
173         else:
174             # Similarly to before, we add a small probability for unseen transitions
175             log_prob += np.log(1e-10)
176     return log_prob
177
178 # We load the three files
179 symbols = load_symbols('symbols.txt')
180 encrypted_message = load_encrypted_message('message.txt')
181 war_and_peace_text = load_war_and_peace('war_and_peace.txt')
182
183 # We calculate the transition statistics
184 stationary_dist, transition_probs = calculate_transition_statistics(war_and_peace_text)
185
186 # We initialise the Metropolis-Hastings parameters
187 N = len(encrypted_message)
188 iterations = 10000
189 report_every = 100
190 initial_permutation = generate_initial_permutation_based_on_frequencies(encrypted_message)
191 best_permutation = initial_permutation.copy()
192 best_prob = -np.inf
193
194 # The following code runs the un the Metropolis-Hastings sampling

```

```

195 current_permutation = initial_permutation.copy()
196 for iteration in range(iterations):
197     # We propose a new permutation by swapping two symbols
198     s, s_prime = random.sample(current_permutation.keys(), 2)
199     proposed_permutation = current_permutation.copy()
200     proposed_permutation[s], proposed_permutation[s_prime] = proposed_permutation[s_prime],
        proposed_permutation[s]
201
202     # We calculate the posterior probabilities of the current and proposed permutations
203     current_prob = calculate_posterior_probability(encrypted_message, current_permutation,
        transition_probs, stationary_dist)
204     proposed_prob = calculate_posterior_probability(encrypted_message, proposed_permutation,
        transition_probs, stationary_dist)
205
206     # We calculate acceptance probability
207     acceptance_probability = min(1, np.exp(proposed_prob - current_prob))
208
209     # We decide whether to accept the proposed permutation
210     if random.uniform(0, 1) < acceptance_probability:
211         current_permutation = proposed_permutation
212         current_prob = proposed_prob
213
214     # We update the best permutation found so far
215     if current_prob > best_prob:
216         best_permutation = current_permutation.copy()
217         best_prob = current_prob
218
219     # We report the decryption every 'report_every' iterations
220     if iteration % report_every == 0:
221         decrypted_text = decrypt_text(encrypted_message[:60], current_permutation)
222         print(f"Iteration {iteration}: {decrypted_text}")

```

Listing 8: Python code showing implemented MH sampler

When run, the code above yields the following decryption of the first 60 symbols after every 100 iterations:

Iteration	Text
0	on lw whunges tnm lhse furnestpre wetsi lw ytades gtfe le ih
100	on d, iungew tna diwe hurnewtpre ,etws d, mtlyew gthe de si
200	an d, iongew uns diwe yolnewuple ,euwr d, muthew guye de ri
300	in d, oungep ans dope yulnepable ,eapr d, mathep gaye de ro
400	in m. .oungep ans mope yulnepable .eapr m. dathep gaye me ro
500	in m. .oungep ans mope yulnepable .eapr m. dathep gaye me ro
600	in my youngep ans mope .ulnepable yeapr my dathep ga.e me ro
700	in my youngep and mope .ulnepable yeaps my rathep ga.e me so
800	in my younger and more .ulnerable years my pather ga.e me so
900	in my younger and more zulnerable years my father gaze me so
1000	in my younger and more zulnerable years my father gaze me so
1100	in my younger and more vulnerable years my father gave me so

1200	in my younger and more vulnerable years my father gave me so
1300	in my younger and more vulnerable years my father gave me so
1400	in my younger and more vulnerable years my father gave me so
1500	in my younger and more vulnerable years my father gave me so
1600	in my younger and more vulnerable years my father gave me so
1700	in my younger and more vulnerable years my father gave me so
1800	in my younger and more vulnerable years my father gave me so
1900	in my younger and more vulnerable years my father gave me so
2000	in my younger and more vulnerable years my father gave me so

I've displayed no more than 2000 iterations in the table above as the first 60 characters of the message didn't change after this. The chain converged reasonably quickly to a sensible message, and this could be due to the way we initialised the permutation. Following the hint, the permutation wasn't initialised at random and was instead obtained by looking at the frequency of every character in the encrypted message and the war and peace text. By matching the most frequent characters from both, we used the natural frequency distribution of English letters to make an informed guess at the mapping of encrypted symbols. The advantage of this is obvious for our case, however, an issue can arise when the encrypt is too small to determine meaning statistics on the frequency of the characters. Even then, it is better than a random guess as it increases the initial likelihood of correct decryption. Another possible ways to smartly initialise the permutation is to run it many times with random parameters, compute the log-likelihood for each, and keep the best one.

## 5.5 Question 5e

A Markov chain is called ergodic if it is possible to reach any state from any other state, given enough time, and it has a unique stationary distribution. If some transition probabilities  $\psi(\alpha, \beta)$  are zero, it means that certain transitions between states are not allowed under the current model. This would stop the chain from being ergodic as it would prevent certain states reaching others. To restore ergodicity, we can add a small probability to every transition, which will make sure that all states can be reached. This is done and clearly commented in Listing 8.

## 5.6 Question 5f

Using only symbol probabilities would significantly weaken the algorithm for two main reason. The first is the lack of context. By mapping the symbols in both the English and

encrypted text based on their overall frequency, we are assuming that the characters are independent of each other, which we've seen to not be true in the transition probability matrix. Composite sounds such as "th", "kn" or "oo" for example make certain sequences of letters more probable. This method would work if you assumed that all the characters have distinct frequencies and the both the English and encrypted texts are large (and diverse) enough to show the correct frequency of the characters. Needless to say, this is impractical and very hard to verify. The second reason arises from cases where we don't have representative frequencies for the symbols, as many have similar frequencies, which would make deciphering smaller phrases impossible. High-frequency symbols might be mapped correctly, but without transitions giving information on the structure of the language, low-frequency symbols are difficult to decrypt accurately.

Second-order Markov chains would introduce the dependence on the last two symbols (instead of one) to predict the next. In theory, this would be even better than the first-order Markov chain, as the model would have an even better description of the language's structure. The issue arises when considering the implementation of this method. The transition matrix would go from having  $53^2$  to  $53^3$  entries, which makes it a lot heavier to work with computationally. On top of that, we would need a much larger corpus of English text to have accurate transition probabilities for rare sequences of characters. In essence, it is theoretically a good idea to use second-order Markov chains, but in practice it will make difficult to obtain a good transition probability matrix, will make the convergence of the sampler slower (because of the increase in complexity), and will increase computation costs and data requirements.

Having two different symbols map to the same encrypted one would create issues when decrypting. The first of the loss of information, as there would be no way for the Markov chain to distinguish and separate the two original symbols, leading to a decrypted message which won't read like correct English. Additionally, the log probability function relies on modelling the decryption process as a permutation of symbols. If we didn't have one-to-one mapping and the function could not uniquely determine the English from the encrypted text. There will then be ambiguity for the sampler, as multiple plaintext sequences could correspond to the same encrypted text, which might stop it from converging quickly to a solution.

Similarly to what we said for second-order Markov chains, using this method to decrypt Chinese would lead to significant issues computationally. If Chinese has more than 1000 characters, then the transition probability matrix would have  $1000^2$  entries, which is very difficult to store and use in the sampling. Moreover, this would require an enormous amount of text to get accurate transition probabilities, as some characters would be very rare. The set of all possible permutations is so large that it would take the MC sampler a very long time



to converge too. In conclusion, I don't think this method would work to decrypt Chinese symbols.

## 6 Bonus: Implementing Gibbs sampling for LDA

Not attempted

## 7 Optimization

### 7.1 Question 7a

We want to solve for the local extrema of  $f(x, y) = x + 2y$  subject to the constraint  $y^2 + xy = 1$  using Lagrange multipliers. This method finds the extrema of  $f(x, y)$  with the constraint  $g(x, y) = 0$ , so first, we will rewrite the constraint as:

$$g(x, y) = y^2 + xy - 1$$

with  $g(x, y) = 0$ , to match the notation. The Lagrange multiplier method uses the fact that, at any point where the function  $f$  reaches a local extremum on the constraint curve, the gradients of  $f(x, y)$  and  $g(x, y)$  are aligned. Since we're constrained to the curve defined by  $g(x, y) = 0$ , any movement along the constraint doesn't change  $g(x, y)$ . So, at the extrema, the gradients of  $f(x, y)$  and  $g(x, y)$  must be proportional, meaning they point in the same or opposite directions. Mathematically this translates to:

$$f(x, y) = \lambda g(x, y)$$

$$f(x, y) - \lambda g(x, y) = 0$$

Intuitively from the expression above, we define the Lagrange function as:

$$\mathcal{L}(x, y, \lambda) = f(x, y) - \lambda g(x, y) = x + 2y - \lambda(y^2 + xy - 1) = 0$$

It follows that the partial derivatives with respect to  $x$ ,  $y$  and  $\lambda$  of the Lagrange function would be equal to 0. This gives:

$$\frac{\partial \mathcal{L}}{\partial x} = 1 - \lambda y = 0.$$

$$\frac{\partial \mathcal{L}}{\partial y} = 2 - \lambda(2y + x) = 0.$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = -(y^2 + xy - 1) = 0.$$

Which results in the following system of equations:

$$\begin{cases} 1 - \lambda y = 0, \\ 2 - \lambda(2y + x) = 0, \\ y^2 + xy = 1. \end{cases}$$

We can now solve this system of equations. From the first equation we have:

$$\lambda = \frac{1}{y}$$

Substituting this into the second equation gives:

$$2 - \frac{1}{y}(2y + x) = 0 \Rightarrow 2 - 2 - \frac{x}{y} = 0 \Rightarrow x = 0.$$

Substituting  $x = 0$  into the constraint equation:

$$y^2 + 0 \cdot y = 1 \Rightarrow y^2 = 1 \Rightarrow y = \pm 1$$

The solution to the system of equation is:

$$\begin{cases} \lambda = \mp 1 \\ x = 0 \\ y = \pm 1 \end{cases}$$

This tells us that there are two local extrema:

1. When  $\lambda = 1$ , the extrema is at  $(x = 0, y = -1)$
2. When  $\lambda = -1$ , the extrema is at  $(x = 0, y = 1)$

## 7.2 Question 7b

### 7.2.1 Question 7b i

Newton's method is used as an iterative technique to find a root of a function. We want to find  $x$  such that  $x = \ln(a)$ . Using the definition of the natural logarithm, we have:

$$\exp(x) = a$$

We can rewrite this equation as a root-finding problem:

$$f(x, a) = \exp(x) - a = 0$$

Therefore the function  $f(x, a)$  we are looking for is:

$$f(x, a) = \exp(x) - a$$

### 7.2.2 Question 7b ii

The update equation for Newton's algorithm is given by:

$$x_{n+1} = x_n - \frac{f(x_n, a)}{f'(x_n, a)}$$

where  $f'(x_n, a)$  is the derivative of  $f(x_n, a)$  with respect to  $x$ . For our case:

$$\frac{\partial f(x, a)}{\partial x} = \frac{\partial}{\partial x}(\exp(x) - a) = \exp(x)$$

Substituting in the function we found in part (i) into the update equation, we get:

$$x_{n+1} = x_n - \frac{\exp(x_n) - a}{\exp(x_n)}$$

## 8 Bonus: Eigenvalues as solutions of an optimization problem

### 8.1 Question 8a

The Extreme Value Theorem states that a continuous function reaches its maximum and minimum on a compact set. We are trying to solve the optimisation problem:

$$x^* := \arg \max_{x \in \mathbb{R}^n} R_A(x)$$

But since  $\mathbb{R}^n$  isn't compact, we need to search a compact subset. We use the hint that is given, which states we should at the unit sphere:

$$S = \{x \in \mathbb{R}^n \mid \|x\| = 1\},$$

which is a compact set. We can do this because  $R_A(x)$  doesn't change if we scale  $x$ , as is shown below:

$$R_A(cx) = \frac{(cx)^T A(cx)}{(cx)^T (cx)} = \frac{c^2 x^T A x}{c^2 x^T x} = \frac{x^T A x}{x^T x} = R_A(x)$$

In this compact set,  $R_A(x)$  becomes:

$$R_A(x) = \frac{q_A(x)}{\|x\|^2} = q_A(x)$$

as  $\|x\|^2 = 1$  in the unit sphere. Thanks to this, we see that  $R_A(x)$  is continuous, and thus satisfies all the requirements of the Extreme Value Theorem.  $\sup_{x \in S} R_A(x)$  is attained, we now need to show that this is also true for  $x \in \mathbb{R}^n$ . Let  $z = \frac{x}{\|x\|}$ . Substituting this into our equation for  $R_A(x)$  yields:

$$\begin{aligned} R_A(x) &= \frac{x^T A x}{\|x\|^2} \\ &= \frac{x^T A x}{x^T x} \\ &= \left( \frac{x}{\|x\|} \right)^T A \left( \frac{x}{\|x\|} \right) \\ &= z^T A z \\ &= R_A(z) \end{aligned}$$

This equality shows that the superior in the compact set  $S$  is the same as the one in  $\mathbb{R}^n$ , so we have shown that  $\sup_{x \in \mathbb{R}^n} R_A(x)$  is attained.

## 8.2 Question 8b

We are told that eigenvectors  $\xi_1, \dots, \xi_n$  form an orthonormal basis and we any vector  $x$  can be represented through that orthonormal basis as:

$$x = \sum_{i=1}^n (\xi_i^T x) \xi_i$$

With this in mind, we expand our expression using the eigenbasis, looking at the nominator first then the denominator. We substitute  $x = \sum_{i=1}^n (\xi_i^T x) \xi_i$  into  $x^T Ax$ :

$$x^T Ax = \left( \sum_{i=1}^n (\xi_i^T x) \xi_i \right)^T A \left( \sum_{j=1}^n (\xi_j^T x) \xi_j \right)$$

Using linearity of inner product:

$$x^T Ax = \sum_{i=1}^n \sum_{j=1}^n (\xi_i^T x) (\xi_j^T x) \xi_i^T A \xi_j$$

As  $\xi_i$  are eigenvectors of  $A$  with  $A\xi_j = \lambda_j \xi_j$ , we have  $\xi_i^T A \xi_j = \lambda_j \delta_{ij}$ . Therefore we have:

$$x^T Ax = \sum_{i=1}^n (\xi_i^T x)^2 \lambda_i$$

We proceed similarly when looking at the denominator:

$$x^T x = \left( \sum_{i=1}^n (\xi_i^T x) \xi_i \right)^T \left( \sum_{j=1}^n (\xi_j^T x) \xi_j \right) = \sum_{i=1}^n (\xi_i^T x)^2$$

Putting everything back together gives:

$$R_A(x) = \frac{x^T Ax}{x^T x} = \frac{\sum_{i=1}^n (\xi_i^T x)^2 \lambda_i}{\sum_{i=1}^n (\xi_i^T x)^2}$$

The question tells us that the eigenvalues of  $A$   $\lambda_1 \geq \dots \geq \lambda_n$  are enumerated in decreasing size, meaning  $\lambda_1$  is the largest. We can write:

$$R_A(x) = \frac{x^T Ax}{x^T x} \leq \frac{\sum_{i=1}^n (\xi_i^T x)^2 \lambda_1}{\sum_{i=1}^n (\xi_i^T x)^2} \leq \lambda_1$$

as required.

### 8.3 Question 8c

No other eigenvectors have the eigenvalue  $\lambda_1$ , since we have made sure that  $x$  isn't one of the eigenvectors with degenerate eigenvalue  $\lambda_1$ . This means the inequality  $\lambda_1 \geq \dots \geq \lambda_n$  becomes strict  $\lambda_1 > \dots > \lambda_n$ . Using the final step of our calculation in the previous question, we see that the solution arises:

$$\begin{aligned} R_A(x) &= \frac{x^T Ax}{x^T x} \\ &= \frac{\sum_{i=1}^n (\xi_i^T x)^2 \lambda_i}{\sum_{i=1}^n (\xi_i^T x)^2} \\ &< \frac{\lambda_1 \sum_{i=1}^n (\xi_i^T x)^2}{\sum_{i=1}^n (\xi_i^T x)^2} \\ &< \lambda_1 \end{aligned}$$