

# COMP0086: Probabilistic and Unsupervised Learning

January 24, 2025

## Contents

|                     |           |
|---------------------|-----------|
| <b>Question 1</b>   | <b>3</b>  |
| 1-a . . . . .       | 3         |
| 1-b . . . . .       | 3         |
| 1-c . . . . .       | 3         |
| 1-d and e . . . . . | 4         |
| <b>Question 2</b>   | <b>6</b>  |
| <b>Question 3</b>   | <b>8</b>  |
| 3-a . . . . .       | 8         |
| 3-b . . . . .       | 8         |
| 3-c . . . . .       | 8         |
| 3-d . . . . .       | 9         |
| 3-e . . . . .       | 10        |
| <b>Question 5</b>   | <b>13</b> |
| 5-a . . . . .       | 13        |
| 5-b . . . . .       | 14        |
| 5-c . . . . .       | 15        |
| 5-d . . . . .       | 16        |
| 5-e . . . . .       | 19        |
| 5-f . . . . .       | 20        |
| <b>Question 7</b>   | <b>22</b> |
| 7-a . . . . .       | 22        |
| 7-b . . . . .       | 22        |
| <b>Question 4</b>   | <b>23</b> |
| 4-a . . . . .       | 23        |
| <b>Question 8</b>   | <b>25</b> |
| 8-a . . . . .       | 25        |

|               |    |
|---------------|----|
| 8-b . . . . . | 25 |
|---------------|----|

## Question 1

### Question 1.a

The dataset consists of binary images with each pixel as 0 or 1, meaning the data is discrete. A multivariate Gaussian model assumes continuous values, which would allow non-binary values between 0 and 1, making it unsuitable. A better choice would be a Bernoulli distribution, which models binary outcomes directly.

### Question 1.b

The maximum likelihood estimator (MLE) is defined as

$$\hat{\theta}_{\text{ML}} := \arg \max_{\theta \in \mathcal{T}} p(x_1, \dots, x_n \mid \theta), \quad (1)$$

which can be found by solving maximum likelihood equation:

$$\nabla_{\theta} p(\mathcal{D} \mid \theta) = \nabla_{\theta} \left( \prod_{i=1}^n p(x_i \mid \theta) \right) = 0, \quad (2)$$

where we will estimate parameter  $p$  in further discussion. For convenience, we use the log-likelihood to convert the product into a sum, allowing us to rewrite the expression as:

$$\sum_{i=1}^n \nabla_p \log p(x_i \mid \theta) = 0. \quad (3)$$

Substituting with our multivariate Bernoulli distribution, this becomes:

$$\sum_{i=1}^n \nabla_p \log \prod_{d=1}^D p_d^{x_d} (1 - p_d)^{(1-x_d)} = 0. \quad (4)$$

We will first solve the expression inside the summation over  $n$  and consider the summation afterward. By taking the log of the product, we convert it into a sum, leading to:

$$\nabla_p \sum_{d=1}^D (x_d \log p_d + (1 - x_d) \log (1 - p_d)) = 0. \quad (5)$$

Taking the derivative yields:

$$\frac{x_d}{\hat{p}_d} - \frac{(1 - x_d)}{1 - \hat{p}_d} = 0, \quad (6)$$

where  $\hat{p}_d$  represents the estimated probability for the  $d$ -th element in the final MLE solution.

Now, we incorporate the summation across all  $N$  observations:

$$\frac{\sum_{n=1}^N x_d^{(n)}}{\hat{p}_d} - \frac{\sum_{n=1}^N (1 - x_d^{(n)})}{1 - \hat{p}_d} = 0, \quad (7)$$

where  $N$  is the total number of pixels in our binary images. Since each  $x$  value is either 1 or 0, let  $N_1$  be the count of ones and  $N_0$  the count of zeros. In this case, we have:

$$\frac{N_1}{\hat{p}_d} - \frac{N_0}{1 - \hat{p}_d} = 0 \Rightarrow \hat{p}_d = \frac{N_1}{N_1 + N_0}, \quad (8)$$

which gives the  $d$ -th element of the MLE solution. Thus, the overall solution can be expressed as:

$$\hat{\mathbf{p}}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)}. \quad (9)$$

### Question 1.c

The maximum a posteriori (MAP) estimator is defined as:

$$\hat{\theta}_{\text{MAP}} := \arg \max_{\theta \in \mathcal{T}} p(\theta \mid x_1, \dots, x_n), \quad (10)$$

which can be rewritten using Bayes' theorem as:

$$\hat{\theta}_{\text{MAP}} := \arg \max_{\theta \in \mathcal{T}} \frac{p(x_1, \dots, x_n \mid \theta) p(\theta)}{p(x_1, \dots, x_n)}. \quad (11)$$

Following a similar approach to 1(b), we take the log-likelihood and then differentiate with respect to  $p$ . Since the evidence term  $p(x_1, \dots, x_n)$  is independent of  $p$ , we can focus on maximizing  $p(x_1, \dots, x_n \mid \theta) p(\theta)$ :

$$\nabla_{\theta} p(\theta \mid \mathcal{D}) = \nabla_{\theta} \left( \prod_{i=1}^n p(x_i \mid \theta) p(\theta) \right) = 0. \quad (12)$$

Using the log-likelihood to convert the product into a sum, we can rewrite this as:

$$\sum_{i=1}^n \nabla_p \log p(x_i \mid \theta) p(\theta) = 0. \quad (13)$$

Substituting the Beta prior and the multivariate Bernoulli distribution, we obtain:

$$\sum_{i=1}^n \nabla_p \log \prod_{d=1}^D p_d^{x_d} (1 - p_d)^{(1-x_d)} \prod_{d=1}^D \frac{1}{B(\alpha, \beta)} p_d^{\alpha-1} (1 - p_d)^{\beta-1} = 0. \quad (14)$$

Taking the log of each term, we express this as:

$$-\sum_{d=1}^D \log(B(\alpha, \beta)) + \sum_{d=1}^D x_d \log(p_d) + (1 - x_d) \log(1 - p_d) + \sum_{d=1}^D (\alpha - 1) \log(p_d) + (\beta - 1) \log(1 - p_d). \quad (15)$$

To simplify, we focus on the  $d$ -th element and find its derivative as:

$$\left( \frac{x_d}{\hat{p}_d} - \frac{1 - x_d}{1 - \hat{p}_d} \right) + \left( \frac{\alpha - 1}{\hat{p}_d} - \frac{\beta - 1}{1 - \hat{p}_d} \right) = 0. \quad (16)$$

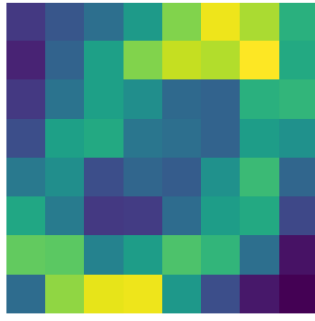
Since  $B(\alpha, \beta)$  is independent of  $p_d$ , it no longer appears in the derivative.

Applying the summation as in 1(b), we get:

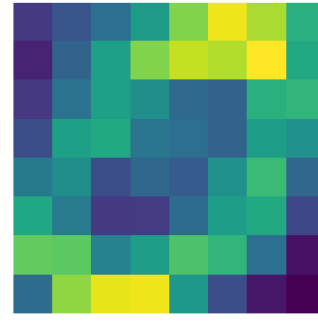
$$\begin{aligned} & \left( \frac{N_1}{\hat{p}_d} - \frac{N_0}{1 - \hat{p}_d} \right) + \left( \frac{\alpha - 1}{\hat{p}_d} - \frac{\beta - 1}{1 - \hat{p}_d} \right) = 0 \\ \Rightarrow & \frac{N_1 + \alpha - 1}{\hat{p}_d} = \frac{N_0 + \beta - 1}{1 - \hat{p}_d} \\ \Rightarrow & \hat{p}_d = \frac{N_1 + \alpha - 1}{N + \alpha + \beta - 2}. \end{aligned} \quad (17)$$

The final MAP estimator can thus be expressed as:

$$\hat{\mathbf{p}}_{\text{MAP}} = \frac{\alpha - 1 + \sum_{n=1}^N \mathbf{x}^{(n)}}{N + \alpha + \beta - 2}. \quad (18)$$



(a) MLE



(b) MAP

Figure 1: ML and MAP parameters of a multivariate Bernoulli

## Question 1.d and e

Here, we use the code to plot both ML parameters of a multivariate Bernoulli and MAP parameters.

```

1 def MLE(x):
2     return np.mean(x, axis=0)
3 def MAP(x, alpha, beta):
4     N, _ = x.shape
5     return (alpha - 1 + np.sum(x, axis=0)) / (N + alpha + beta - 2)

```

Listing 1: Question 1 code

Maximum Likelihood (ML) and Maximum A Posteriori (MAP) estimation offer different approaches to parameter estimation, particularly in their handling of prior information. ML estimation relies solely on the observed data, maximizing the likelihood to provide the best estimate under the assumption that no prior knowledge about the parameters exists. In contrast, MAP incorporates prior beliefs or information through a prior distribution, allowing for potentially better estimates, especially when the data is limited or noisy. However, MAP's reliance on the prior can be a drawback if the prior is poorly chosen, as it could introduce bias into the estimate. Therefore, when reliable prior knowledge is available, MAP can improve the estimation, particularly if the prior is well-aligned with the likelihood function, such as when a conjugate prior is used. In situations where prior information is unavailable or uncertain, ML is often the more robust choice, avoiding the risk of introducing bias from an inappropriate prior.

## Question 2

Let us label these three model as  $M_1$ ,  $M_2$  and  $M_3$ . Then, the equally likely priori for different models can be written as:

$$P(M_1) = P(M_2) = P(M_3) = \frac{1}{3}. \quad (19)$$

Now, We will derive the likelihood  $P(x | p)$ .

For (a), we substitute  $p_d = 0.5$  into likelihood equation:

$$P(\mathbf{x}^1 \dots \mathbf{x}^N | \mathbf{p}^1) = \prod_{n=1}^N \prod_{d=1}^D 0.5^{x_d^{(n)}} (1 - 0.5)^{(1-x_d^{(n)})} = 0.5^{ND}, \quad (20)$$

where we can apply log:

$$\log P(\mathbf{x}^1 \dots \mathbf{x}^N | \mathbf{p}^1) = -ND \log(2). \quad (21)$$

For (b), we have  $\mathbf{p}$  in forms of  $(p_d, \dots, p_d)$  where  $p_d$  can be any random number between 0 and 1.

$$\begin{aligned} P(\mathbf{x}^1 \dots \mathbf{x}^N | \mathbf{p}^2) &= \int_0^1 \prod_{n=1}^N \prod_{d=1}^D (p_d)^{x_d^{(n)}} (1 - p_d)^{(1-x_d^{(n)})} dp_d, \\ \Rightarrow P(\mathbf{x}^1 \dots \mathbf{x}^N | \mathbf{p}^2) &= \int_0^1 p_d^{\sum_1^N \sum_d^D x_d^{(n)}} (1 - p_d)^{\sum_1^N \sum_d^D (1-x_d^{(n)})} dp_d, \end{aligned} \quad (22)$$

which reminds us the form of Beta distribution:

$$B(\alpha, \beta) = \int_0^1 u^{\alpha-1} (1-u)^{\beta-1} du. \quad (23)$$

Thus, we can find this likelihood as:

$$P(\mathbf{x}^1 \dots \mathbf{x}^N | \mathbf{p}^2) = B \left( \sum_1^N \sum_d^D x_d^{(n)} + 1, \sum_1^N \sum_d^D (1 - x_d^{(n)}) + 1 \right). \quad (24)$$

For (c), we repeat the similar process as (b) with different entries of  $p_d$ . The likelihood function can be expressed as:

$$\begin{aligned} P(\mathbf{x}^1 \dots \mathbf{x}^N | \mathbf{p}^3) &= \int_0^1 \dots \int_0^1 \prod_{n=1}^N (p_d)^{x_d^{(n)}} (1 - p_d)^{(1-x_d^{(n)})} dp_1 \dots dp_D, \\ \Rightarrow P(\mathbf{x}^1 \dots \mathbf{x}^N | \mathbf{p}^3) &= \prod_{d=1}^D \left( \int_0^1 p_d^{\sum_1^N x_d^{(n)}} (1 - p_d)^{\sum_1^N (1-x_d^{(n)})} dp_d \right), \end{aligned} \quad (25)$$

where we shall integrate over different  $p_d$ . Luckily, we still have the form of Beta distribution and our likelihood can be expressed as:

$$P(\mathbf{x}^1 \dots \mathbf{x}^N | \mathbf{p}^3) = \prod_{d=1}^D B \left( \sum_1^N x_d^{(n)} + 1, \sum_1^N (1 - x_d^{(n)}) + 1 \right). \quad (26)$$

Following the Bayes rule, we can find the posterior probabilities as:

$$P(M_i | \mathbf{x}) = \frac{P(\mathbf{x} | M_i) P(M_i)}{P(\mathbf{x})}, \quad (27)$$

where the evidence term  $P(\mathbf{x})$  is constant. For simplicity, we can compare the logarithm of the posterior:

$$\begin{aligned} \log P(M_1 | \mathbf{x}) &\propto -ND \log(2) - \log(3), \\ \log P(M_2 | \mathbf{x}) &\propto \text{Beta} \ln \left( \sum_1^N \sum_d^D x_d^{(n)} + 1, \sum_1^N \sum_d^D (1 - x_d^{(n)}) + 1 \right) - \log(3), \\ \log P(M_3 | \mathbf{x}) &\propto \sum_{d=1}^D \text{Beta} \ln \left( \sum_1^N x_d^{(n)} + 1, \sum_1^N (1 - x_d^{(n)}) + 1 \right) - \log(3), \end{aligned} \quad (28)$$

where  $\text{Beta} \ln$  is the logarithm of Beta distribution. Here, we can calculate these with the following code.

```

1 #Load Data
2 X = np.loadtxt('binarydigits.txt')
3 N, D = X.shape
4 # Model 1
5 log_P_M1 = - N*D*np.log(2) - np.log(3)
6 # Model 2
7 log_P_M2 = betaln(np.sum(X)+1,np.sum(1-X)+1) - np.log(3)
8 # Model 3
9 log_P_M3 = 0
10 for d in range(D):
11     X_d = X[d,:]
12     log_P_M3 += betaln(np.sum(X_d)+1,N-np.sum(X_d)+1)
13 log_P_M3 -= np.log(3)
14 print(log_P_M1,log_P_M2,log_P_M3)

```

Listing 2: Question 2 code

Thus, we get the final results as:

$$\begin{aligned}
 \log P(M_1 | \mathbf{x}) &\propto -4437.24, \\
 \log P(M_2 | \mathbf{x}) &\propto -4284.82, \\
 \log P(M_3 | \mathbf{x}) &\propto -3727.40.
 \end{aligned}
 \tag{29}$$

### Question 3

#### Question 3.a

Now, we want to add a new set of parameters  $\pi_1 \dots \pi_K$  as a proportions to our original multivariate Bernoulli distributions:

$$P(\mathbf{x}^{(n)} | \mathbf{P}) = \prod_{d=1}^D p_d^{x_d} (1 - p_d)^{(1-x_d)}, \quad (30)$$

which can be expressed as:

$$P(\mathbf{x}^{(n)} | \boldsymbol{\pi}, \mathbf{P}) = \sum_{k=1}^K \pi_k \prod_{d=1}^D p_{kd}^{x_d} (1 - p_{kd})^{(1-x_d)}. \quad (31)$$

For N images, the likelihood can be written as:

$$P(\mathbf{x}^1 \dots \mathbf{x}^N | \boldsymbol{\pi}, \mathbf{P}) = \prod_{n=1}^N \sum_{k=1}^K \pi_k \prod_{d=1}^D p_{kd}^{x_d} (1 - p_{kd})^{(1-x_d)}. \quad (32)$$

#### Question 3.b

Using the expression of responsibility, we can get:

$$\begin{aligned} r_{nk} &= P(s^{(n)} = k | \mathbf{x}^{(n)}, \boldsymbol{\pi}, \mathbf{P}), \\ \Rightarrow r_{nk} &= \frac{P(\mathbf{x}^{(n)} | s^{(n)} = k, \boldsymbol{\pi}, \mathbf{P}) P(s^{(n)} = k | \boldsymbol{\pi}, \mathbf{P})}{\sum_{k'=1}^K P(\mathbf{x}^{(n)} | s^{(n)} = k', \boldsymbol{\pi}, \mathbf{P}) P(s^{(n)} = k' | \boldsymbol{\pi}, \mathbf{P})}, \\ \Rightarrow r_{nk} &= \frac{\pi_k \prod_{d=1}^D p_{kd}^{x_d} (1 - p_{kd})^{(1-x_d)}}{\sum_{k'=1}^K \pi_{k'} \prod_{d=1}^D p_{k'd}^{x_d} (1 - p_{k'd})^{(1-x_d)}}. \end{aligned} \quad (33)$$

#### Question 3.c

Let us calculate the expected log-joint first.

$$\begin{aligned} &\left\langle \sum_n \log P(\mathbf{x}^{(n)}, s^{(n)} | \boldsymbol{\pi}, \mathbf{P}) \right\rangle_{q(\{s^{(n)}\})}, \\ &= \sum_{k=1}^K \sum_{n=1}^N q(s^{(n)} = k) \log P(\mathbf{x}^{(n)}, s^{(n)} | \boldsymbol{\pi}, \mathbf{P}), \\ &= \sum_{k,n} r_{nk} \log P(s^{(n)} | \mathbf{x}^{(n)}, \boldsymbol{\pi}, \mathbf{P}) P(\mathbf{x}^{(n)} | s^{(n)}, \boldsymbol{\pi}, \mathbf{P}), \\ &= \sum_{k,n} r_{nk} \log \left( \pi_k \prod_{d=1}^D p_{kd}^{x_d} (1 - p_{kd})^{(1-x_d)} \right), \\ &= \sum_{k,n} r_{nk} \left( \log \pi_k + \sum_{d=1}^D x_d^{(n)} \log(p_{kd}) + (1 - x_d^{(n)}) \log(1 - p_{kd}) \right). \end{aligned} \quad (34)$$

We denote the expectation above as  $E$ . To maximize  $E$ , we first take the partial derivative with respect to  $p_{kd}$ :

$$\begin{aligned} &\frac{\partial E}{\partial p_{kd}} = 0, \\ \Rightarrow &\sum_{n=1}^D r_{nk} \left( \sum_{d=1}^D x_d^{(n)} \frac{1}{p_{kd}} - (1 - x_d^{(n)}) \frac{1}{1 - p_{kd}} \right) = 0, \\ \Rightarrow &\sum_{n=1}^D r_{nk} \left( x_d^{(n)} \frac{1}{p_{kd}} - (1 - x_d^{(n)}) \frac{1}{1 - p_{kd}} \right) = 0, \\ \Rightarrow &(1 - p_{kd}) \sum_{n=1}^D (r_{nk} x_d^{(n)}) = p_{kd} \sum_{n=1}^D r_{nk} (1 - x_d^{(n)}), \end{aligned} \quad (35)$$



which leads to:

$$\hat{p}_{kd} = \frac{\sum_{n=1}^D \left( r_{nk} x_d^{(n)} \right)}{\sum_{n=1}^D r_{nk}}. \quad (36)$$

Now, we move to the derivative with respect to  $\pi_k$ . We need to apply a Lagrangian multiplier as a constraint  $\sum_{k=1}^K \pi_k = 1$  is given. Thus, we have:

$$\begin{aligned} \frac{\partial E}{\partial \pi_k} + \lambda &= 0, \\ \Rightarrow \sum_{k,n} \frac{r_{nk}}{\pi_k} + \lambda &= 0, \\ \Rightarrow \sum_{k,n} r_{nk} + \sum_k \pi_k \lambda &= 0, \\ \Rightarrow N + \lambda &= 0, \end{aligned} \quad (37)$$

where we notice that  $r_{nk}$  is normalized and the constraint of  $\pi_k$  is used. In that case, we can find  $\hat{\pi}_k$ :

$$\hat{\pi}_k = \frac{\sum_n r_{nk}}{N}. \quad (38)$$

### Question 3.d

(d) Implement the EM algorithm for a mixture of K multivariate Bernoulli. Let us first define the likelihood function.

```

1 def likelihood_K_Mixing_Bernoulli(X,K,Pi,P):
2     N, D = X.shape
3     likelihood = 0
4     for n in range(N): # sum over N
5         likelihood_N = []
6         for k in range(K): # sum over K
7             likelihood_NK = Pi[k]*np.prod((P[k,:]**X[n,:])*((1-P[k,:])**(1-X[n,:])))
8             likelihood_N.append(likelihood_NK)
9         likelihood += np.log(np.sum(likelihood_N))
10    return likelihood

```

Question (3) Log-likelihood

Then, the E-step and M-step are coded as:

```

1 def E_step(X,K,Pi,P):
2     N, D = X.shape
3     # create an empty matrix for responsibility
4     Responsibility = np.empty((N,K))
5     # sum over N
6     for n in range(N):
7         Numerator = np.empty(K)
8         # sum over K
9         for k in range(K):
10            Numerator[k] = Pi[k]*np.prod((P[k,:]**X[n,:])*((1-P[k,:])**(1-X[n,:])))
11            Denominator = np.sum(Numerator)
12            # find responsibility
13            Responsibility[n,:] = Numerator/Denominator
14    return Responsibility
15
16 def M_step(X,K,Responsibility):
17     N, D = X.shape
18     # create hat_p_kd and hat_pi_k
19     Hat_P = np.empty((K,D))
20     Hat_Pi = np.empty(K)
21     for k in range(K):
22         # Follow the expressions in Q3(c)
23         Hat_P[k,:] = np.sum(Responsibility[:,k]*X.T,1)/np.sum(Responsibility[:,k])
24         Hat_Pi[k] = np.sum(Responsibility[:,k])/N
25    return Hat_P, Hat_Pi

```

Question (3) E-step and M-step

Finally, we combine above code together to iterate our EM steps:

| K values | $\pi_k$ values   |
|----------|--|
| k=2      | 0.59014409, 0.40985591   |
| k=3      | 0.42000426, 0.17999998, 0.39999575   |
| k=4      | 0.3199043 , 0.16981744, 0.19026642, 0.32001185   |
| k=7      | 0.09, 0.09969762, 0.08037638, 0.18030827, 0.06, 0.08962362, 0.39999412                         |
| k=10     | 0.14, 0.04, 0.15000267, 0.10999683, 0.12999862, 0.06000787, 0.04, 0.03999998, 0.13999403, 0.15 |

Table 1:  $\pi_k$  values

```

1 def EM_Start(X,K):
2     N, D = X.shape
3     max_iterations = 100 # set iterations as 100
4     convergence = 0.0001 # less than convergence, the loop breaks
5     # initials Pi and P
6     Pi_initial = np.random.rand(K)
7     Pi__initial_normal = Pi_initial/np.sum(Pi_initial)
8     P_initial = np.random.rand(K,D)
9     Pi = Pi__initial_normal
10    P = P_initial
11    log_likelihoods = []
12    log_likelihoods.append(likelihood_K_Mixing_Bernoulli(X,K,Pi,P))
13    # start our EM
14    for i in range(max_iterations):
15        # E-step
16        Responsibility = E_step(X,K,Pi,P)
17        # M-step
18        P, Pi = M_step(X,K,Responsibility)
19        log_likelihoods.append(likelihood_K_Mixing_Bernoulli(X,K,Pi,P))
20        difference = np.abs(log_likelihoods[i]-log_likelihoods[i+1])
21        if difference < convergence:
22            break
23    return Responsibility, P, Pi, log_likelihoods

```

Question (3) EM Alogrithm

Following the questions, we run EM algorithm for different K values.

```

1 # load data
2 X = np.loadtxt('binarydigits.txt')
3 # K list:
4 K_list = [2,3,4,7,10]
5 # create the list for collection
6 Responsibility_list = []
7 P_list = []
8 Pi_list = []
9
10 plt.figure(figsize=(10,8))
11 # EM algorithm for different K
12 for K in K_list:
13     Responsibility, P, Pi, log_likelihoods= EM_Start(X, K)
14     # collect parameters
15     Responsibility_list.append(Responsibility)
16     P_list.append(P)
17     Pi_list.append(Pi)
18
19     # plot
20     num_iterations = len(log_likelihoods)
21     iterations = np.linspace(0, num_iterations, num_iterations)
22     plt.plot(iterations, log_likelihoods, label = "K = " + str(K))
23
24 print(Responsibility_list, P_list, Pi_list)
25 plt.legend(loc='lower right')
26 plt.xlabel('Iterations', fontsize = 15)
27 plt.ylabel('Log-likelihood', fontsize = 15)
28 plt.grid()
29 plt.title("Log-likelihood VS Iterations",fontsize = 15)
30 plt.savefig('Q3(d).png')

```

Question (3) Vary K values

Here we have the plot and parameters:

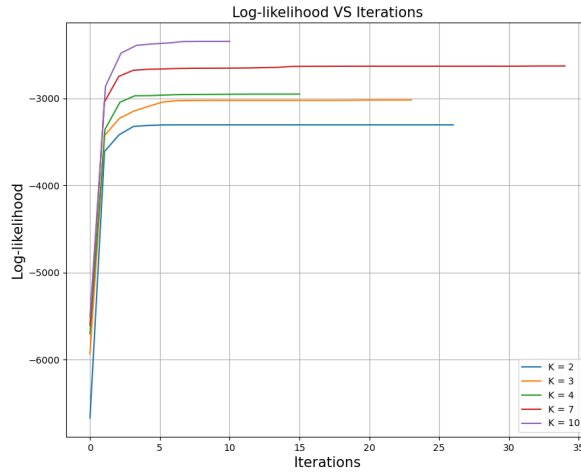


Figure 2: Log-likelihood vs Iterations

### Question 3.e

(e) Run the algorithm a few times starting from randomly chosen initial conditions. We get different log-likelihoods and parameters each time and the log-likelihoods will converge to a fix value after several iterations. We can deduce that the value of log-likelihood depends on value of K from Figure 2.

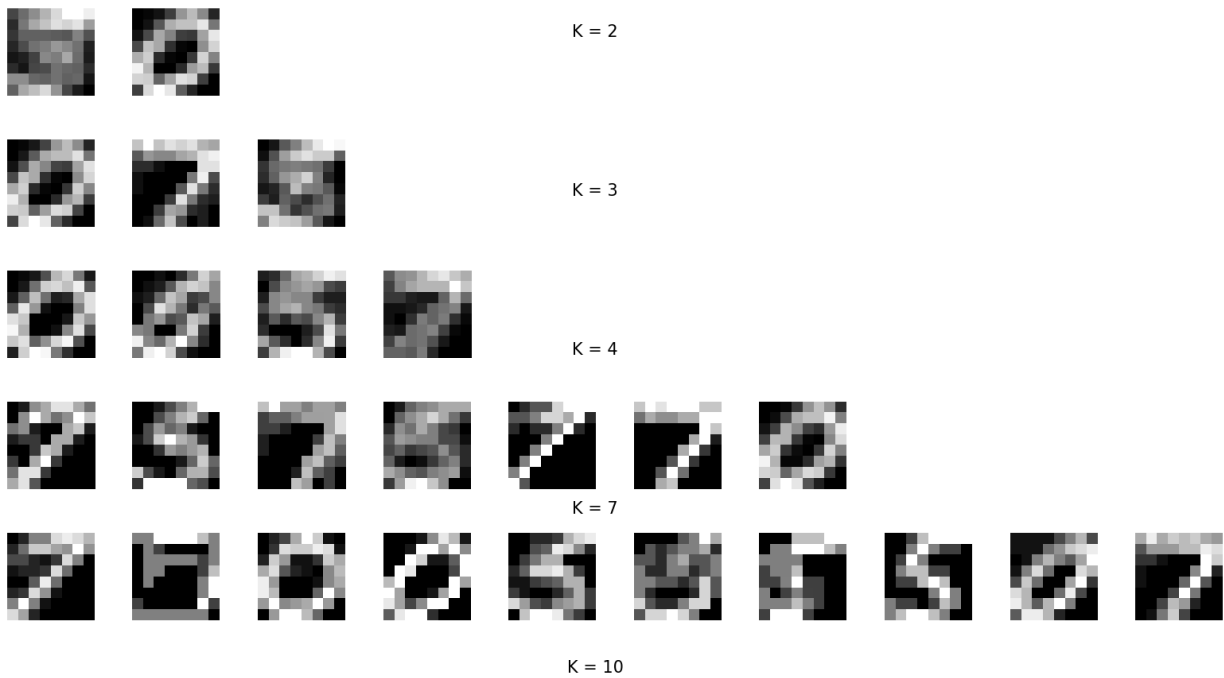


Figure 3: Learned probability vectors

The plot shows the number 0, 5 and 7. We get different mixtures of the numbers that depends on the choices of K-value. In  $K = 2$ , we have roughly the same value of  $\pi_k$  which gives a mixture of 0, 5 and 7 into two different group. For higher K values, we see more groups. It seems that we can get a clear version for each number at  $K = 3$ .

Here is the comment on the algorithm. Overall, it works well as we can find when we check the responsibilities. At the end of the iterations, we get an array containing only one element non-zero. This means our algorithm leaves only one cluster at the end. If we add more K values i.e in Figure 4, we can find that the algorithm works not well at large K. It is not a good choice to have large value K since we only three types of number. A desired K value should be larger or equal to 3. The assumption is made as we know the data set compositions. In this case, one can

improve our model by using prior information of the data. For example, we have three numbers in the data set which means  $K = 3$  could be a good value. At the same time, we notice that handwritings are influenced by personal style which gives different types of writing number 0, 5 and 7. Once we have a preferred K-value, our algorithm could be more efficient.

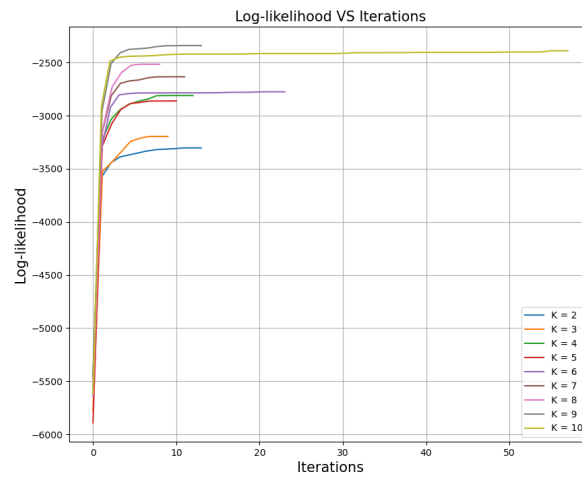


Figure 4: Log-likelihood vs Iterations

## Question 5

### Question 5.a

Let  $N(\alpha, \beta)$  be the counts of numbers of occurrences of symbols and pairs of symbols. We have the transition function in form of  $\psi(\alpha, \beta)$  which represents a change from  $\beta$  to  $\alpha$ . In that case, we can assume that

$$\sum_{\beta} \psi(\alpha, \beta) = 1, \quad (39)$$

as sum of the probability from  $\beta$  to  $\alpha$  must be 1. This can be expressed as a constraint with Lagrangian Multiplier  $\lambda$ . In that case, we can find ML by:

$$\frac{\partial \left( \log p(s_1 s_2 \cdots s_n) + \sum \lambda \left( 1 - \sum_{\beta} \psi(\alpha, \beta) \right) \right)}{\partial \psi(\alpha, \beta)} = 0, \quad (40)$$

where the log-likelihood is :

$$\log p(s_1 s_2 \cdots s_n) = \log p(s_1) + \sum_{i=2}^n \log(\psi(\alpha, \beta)). \quad (41)$$

Then, we can write equation 40 as:

$$\begin{aligned} \frac{N(\alpha, \beta)}{\psi(\alpha, \beta)} - \lambda &= 0, \\ \Rightarrow N(\alpha, \beta) &= \psi(\alpha, \beta) \lambda. \end{aligned} \quad (42)$$

Recall that  $\sum_{\beta} \psi(\alpha, \beta) = 1$ . Thus, we get:

$$\begin{aligned} \sum_{\beta} N(\alpha, \beta) &= \sum_{\beta} \psi(\alpha, \beta) \lambda, \\ \Rightarrow \sum_{\beta} N(\alpha, \beta) &= \lambda, \\ \Rightarrow \psi(\alpha, \beta) &= \frac{N(\alpha, \beta)}{\sum_{\beta} N(\alpha, \beta)}. \end{aligned} \quad (43)$$

Here is the code for estimated probabilities.

```
1 # load data and preprocess them
2 message = open('message.txt', 'r').read().replace('\n', '')
3 symbols = open('symbols.txt', 'r').read().replace('\n', '')
4 war_and_peace = open('2600-0.txt', 'r').read().replace('\n', '').lower()
```

Question 5(a) Load data

```
1 def Psi(text, symbols):
2     num_text = len(text)
3     # numbers of symbols
4     num_sym = len(symbols)
5     # initial the counts
6     transition_counts = np.zeros((num_sym, num_sym))
7     # go through text
8     for i in range(num_text):
9         # pick the character
10        current_sym = text[i-1] # beta
11        next_sym = text[i] # alpha
12        # check whether these character in the symbols.
13        # if character within symbols, we get id. Otherwise, return -1
14        current_sym_id = symbols.find(current_sym)
15        next_sym_id = symbols.find(next_sym)
16        if (current_sym_id != -1) and (next_sym_id != -1):
17            transition_counts[current_sym_id, next_sym_id] += 1
18        else:
19            continue
20    # normalize
21    row_sums = np.sum(transition_counts, axis=1).reshape(-1, 1)
22    Psi = np.divide(transition_counts, row_sums, where=row_sums != 0)
23    return Psi
```

Question 5(a) Transition Matrix

```

1 transition_matrix = Transition_Matrix(war_and_peace, symbols)
2 plt.figure(figsize=(25,25))
3 sns.heatmap(np.round(transition_matrix,2),annot=True, xticklabels=True,yticklabels=True, square=
4             True).get_figure
5 plt.savefig('Q5(a).png')

```

Question 5(a) Estimated probabilities in heatmap

Here is the output plot.

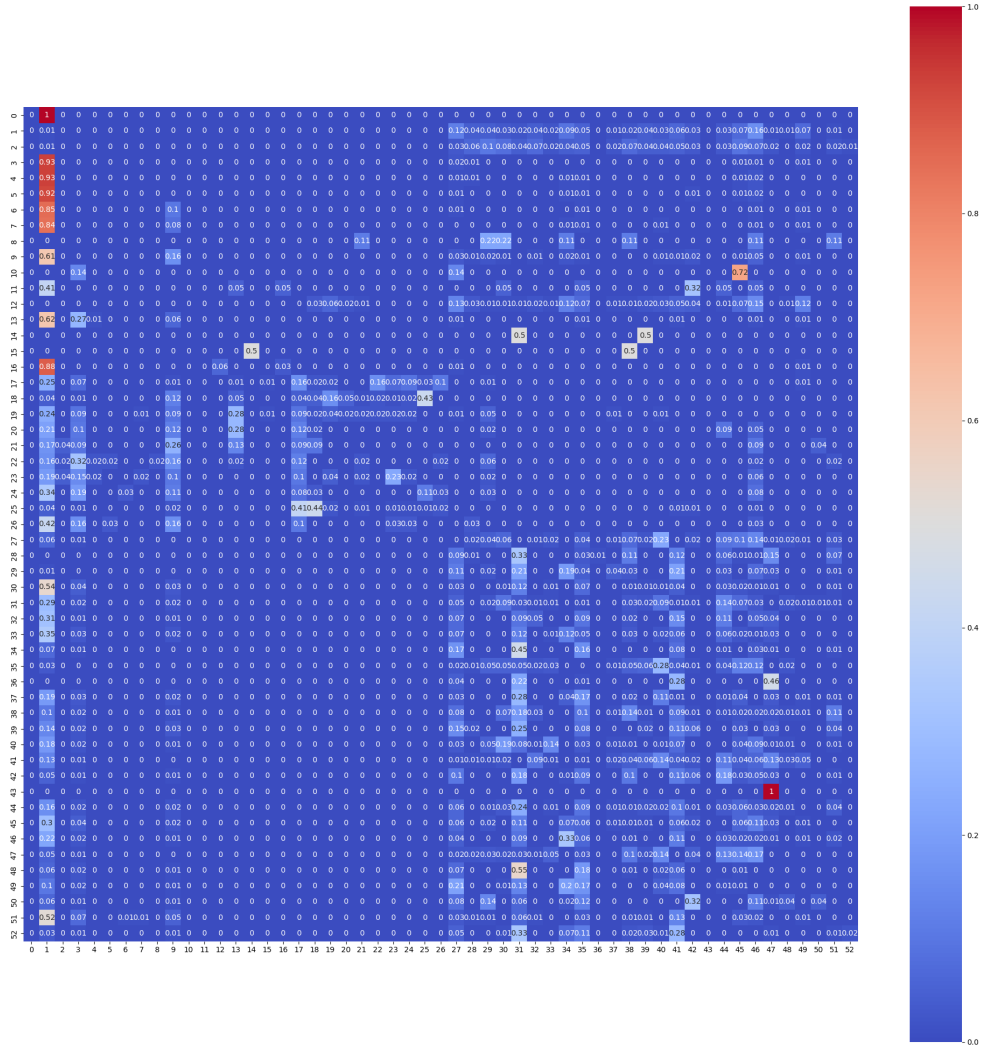


Figure 5: Heatmap for estimated probability

## Question 5.b

In the question, we assume that the mapping between symbols is one-to-one. This means one cannot choose a selected symbol. For example, if the previous  $\sigma(a) = !$ , the following  $\sigma \neq !$ . Thus, the latent variables are not independent.

The joint probability can be expressed as:

$$p(e_1 e_2 \cdots s_n \mid \sigma) = p(e_1 \mid \sigma) \prod_{i=2}^n p(e_i \mid e_{i-1}, \sigma). \quad (44)$$

### Question 5.c

Choosing 2 random symbols from  $n$ , there is  $\binom{n}{2}$  possible choices. In this case, the proposal probability can be written as:

$$S(\sigma \rightarrow \sigma') = \frac{1}{\binom{n}{2}} = \frac{2}{53 \times 52}, \quad (45)$$

where we let  $n = 53$ . At the same time, we find the proposal probability is symmetric. Then, we can express MH acceptance probability as:

$$A(x_{i+1} \mid x_i) = \min \left\{ 1, \frac{q(x_i \mid x_{i+1}) p(x_{i+1})}{q(x_{i+1} \mid x_i) p(x_i)} \right\}, \quad (46)$$

where  $q(x_{i+1} \mid x_i) = q(x_i \mid x_{i+1})$  for symmetric proposal probability. Substituting the probability in our question, we get:

$$\begin{aligned} A(\sigma \mid \sigma') &= \min \left\{ 1, \frac{p(\sigma')}{p(\sigma)} \right\}, \\ \Rightarrow A(\sigma \mid \sigma') &= \min \left\{ 1, \frac{p(e_1 \mid \sigma') \prod_{i=2}^n p(e_i \mid e_{i-1}, \sigma')}{p(e_1 \mid \sigma) \prod_{i=2}^n p(e_i \mid e_{i-1}, \sigma)} \right\}, \\ \Rightarrow A(\sigma \mid \sigma') &= \min \left\{ 1, \frac{p((\sigma')^{-1}(e_1)) \prod_{i=2}^n \psi((\sigma')^{-1}(e_i), \sigma'(e_{i-1}))}{\prod_{i=2}^n \psi(\sigma^{-1}(e_i), \sigma^{-1}(e_{i-1}))} \right\}, \end{aligned} \quad (47)$$

where we write  $\sigma^{-1}$  as a function of  $e_i$  that decrypt the symbols and apply the transition probability defined in (a).

## Question 5.d

Here is the Metropolis-Hastings (MH) chain in code. We first define a function to decrypt the message which helps us to use the key from message. Then, we follow the instruction in 5(c) and define the proposal for the key. After that, we compute the log-likelihood function.

```
1 def decrypt_message(text, symbols, key):
2     # sigma inverse (e)
3     # decode the key to symbols
4     decoded_message = str()
5     for i in text:
6         index = list(key).index(i)
7         letter = symbols[index]
8         decoded_message += letter
9     return decoded_message
10
11 def Proposal(key):
12     # Select two distinct indices to swap
13     key_new = key.copy()
14     idx1, idx2 = random.sample(range(len(key)), 2)
15
16     # Swap the characters at these indices
17     key_new[idx1], key_new[idx2] = key[idx2], key[idx1]
18
19     return key_new
20
21 def log_likelihood_function(text, symbols, key, matrix):
22     decoded_message = decrypt_message(text, symbols, key)
23     log_likelihood = 0
24     for i in range(1, len(text)):
25         beta = decoded_message[i-1]
26         alpha = decoded_message[i]
27         beta_id = symbols.index(beta)
28         alpha_id = symbols.index(alpha)
29         log_likelihood += np.log(matrix[beta_id, alpha_id])
30     return log_likelihood
31
32
33 def MCMC_MH(symbols, text, matrix, intelligent_guess, iterations=10000):
34     key = intelligent_guess.copy()
35     log_likelihoods = []
36     for i in range(iterations):
37         proposal_key = Proposal(key).copy()
38         current_prob = log_likelihood_function(text, symbols, key, matrix)
39         proposal_prob = log_likelihood_function(text, symbols, proposal_key, matrix)
40         log_accept_ratio = proposal_prob - current_prob
41         u = np.log(np.random.rand())
42         if u < log_accept_ratio:
43             key = proposal_key.copy()
44         # print every 100 iteration
45         if i % 100 == 0:
46             print('Iteration', i, ': \n', decrypt_message(text, symbols, key)[:60])
47             print(current_prob)
48             log_likelihoods.append(current_prob)
49     return key, log_likelihoods
```

Question 5(d) Metropolis-Hastings (MH) chain

It is also important to have good initial guess. As we investigate 'War and Peace' in previous question, it is a good idea to find some hints from this book. We reorder the characters based on occurrence in both 'War and Peace' and 'message'. We guess that the occurrence are related. For example, we find 'space' is the most character in 'War and Peace' and 'p' is the most character in 'message'. Then, 'space' might be encoded as 'p'. We follow this guess and repeat for all characters. For those are not shown in the 'War and Peace', we randomly map them to rest character. Then, we reorder again to follow the symbols. This gives us an initial guess for key and here is the code.

```
1 # Order text by counts
2 def count_and_rank_characters(message, symbols):
3     # count each character in the message
4     char_count = Counter(message)
5     # rank characters by frequency, with the most common characters first
6     ranked_chars = [char for char, _ in char_count.most_common()]
7     # create an ordered list of symbols based on the ranked characters
8     ordered_symbols = []
9     seen_symbols = set()
```



```

10 # First, add symbols that match the ranked characters
11 for char in ranked_chars:
12     if char in symbols:
13         ordered_symbols.append(char)
14         seen_symbols.add(char)
15 # get the remaining symbols that were not in the ranked characters
16 remaining_symbols = [char for char in symbols if char not in seen_symbols]
17 # shuffle the remaining symbols randomly and add to the ordered list
18 random.shuffle(remaining_symbols)
19 ordered_symbols.extend(remaining_symbols)
20 return ordered_symbols
21 # mapping text to find the key
22 def combined_mapping_functions(list1, list2, original):
23     # map list1 elements to unique numbers (their index)
24     number_mapping = {element: index for index, element in enumerate(list1)}
25     # map list1 to list2 if they have the same length
26     if len(list1) != len(list2):
27         raise ValueError("list1 and list2 must have the same length.")
28     shuffled_map = {list1[i]: list2[i] for i in range(len(list1))}
29     # reorder shuffled_map to match the order of original
30     reordered_map = {key: shuffled_map[key] for key in original if key in shuffled_map}
31     # return the reorder values
32     return reordered_map.values()

```

Question 5(d) Initial guess

To start the chain, we use

```

1 r1 = count_and_rank_characters(war_and_peace, symbols)
2 r2 = count_and_rank_characters(message, symbols)
3 transition_matrix = Psi(war_and_peace, symbols) + 0.001
4 intelligent_guess = combined_mapping_functions(r1, r2, list(symbols))
5 decrypt_key, log_values = MCMC_MH(symbols, message, transition_matrix, list(intelligent_guess),
    iterations=10000)

```

Question 5(d) Start to decrypt messages

Here is the list of iterations. After 6000 iteration, the text is fixed.

```

1 Iteration 0 :
2 on lw whunges tnm lhse furnestpre wetsi lw ytades gtfe le ih
3 -4913.005879300815
4 Iteration 100 :
5 on cd dhungem tns chme lurnemtire detmp cd ytawem gtfe ce ph
6 -4816.650894335458
7 Iteration 200 :
8 on cd dtungel hns ctfe murnelhire dehlf cd yhawel ghme ce ft
9 -4648.682392184368
10 Iteration 300 :
11 on cd dtungel hns ctfe murnelhire dehlf cd yhawel ghme ce ft
12 -4624.31119562341
13 Iteration 400 :
14 on cd dtungel hns ctfe murnelhire dehlf cd yhawel ghme ce ft
15 -4624.521313081218
16 Iteration 500 :
17 on cd diungel tns cile murnelthre detlf cd ytawel gtme ce fi
18 -4567.078251264491
19 Iteration 600 :
20 on cd diungel tns cile murnelthre detlf cd ytawel gtme ce fi
21 -4563.6137610479645
22 Iteration 700 :
23 on cy yiungel tnf cile murnelthre yetls cy dtawel gtme ce si
24 -4504.451488550505
25 Iteration 800 :
26 on cy yiungel tn, cile zurnelthre yetls cy dtawel gtze ce si
27 -4490.852170848223
28 Iteration 900 :
29 on cy yiungel tn, cile zurnelthre yetls cy dtawel gtze ce si
30 -4458.92164593478
31 Iteration 1000 :
32 on cy yiungel tnd cile zurnelthre yetls cy ,twael gtze ce si
33 -4448.969449694636
34 Iteration 1100 :
35 on cy yiungel tnd cile zurnelthre yetls cy ,twael gtze ce si
36 -4444.745708121956
37 Iteration 1200 :

```

```

38 in cy youngel tnd cole vurnelthre yetls cy ,twael gtve ce so
39 -4320.917470787748
40 Iteration 1300 :
41 in cy youngel tnd cole vurnelthre yetls cy ,twael gtve ce so
42 -4320.917470787748
43 Iteration 1400 :
44 an cy ytunger ond ctile vurnelohre yeols cy mowiel gove ce st
45 -4194.391689769387
46 Iteration 1500 :
47 an cy ytunger ond ctire vulnerohle yeors cy mowier gove ce st
48 -4144.736226056704
49 Iteration 1600 :
50 an hy ytunger ond htre vulnerocle yeors hy mowier gove he st
51 -4128.0836534097825
52 Iteration 1700 :
53 an hy ytunger ond htre vulnerocle yeors hy mowier gove he st
54 -4129.434031119794
55 Iteration 1800 :
56 an hy ytunger ond htre vulnerocle yeors hy mowier gove he st
57 -4125.538626283407
58 Iteration 1900 :
59 an hy ytunger ond htre vulnerofle yeors hy mowier gove he st
60 -4125.500471863725
61 Iteration 2000 :
62 an hy ytunger ond htre vulnerofle yeors hy mowier gove he st
63 -4111.046671905882
64 Iteration 2100 :
65 an hy ytunger ond htre vulnerofle yeors hy mowier gove he st
66 -4110.360219082787
67 Iteration 2200 :
68 an hy ytunger ond htre vulnerofle yeors hy mowier gove he st
69 -4104.101692632136
70 Iteration 2300 :
71 an iy ytunger ond itre vulnerofle yeors iy mowher gove ie st
72 -4079.3138075944585
73 Iteration 2400 :
74 an ty yiunger ond tire vulnerofle yeors ty mowher gove te si
75 -3936.715608710116
76 Iteration 2500 :
77 an ty yiunger ond tire vulnerofle yeors ty mowher gove te si
78 -3934.365399182161
79 Iteration 2600 :
80 an ty yiunger ond tire vulnerofle yeors ty mowher gove te si
81 -3934.365399182161
82 Iteration 2700 :
83 an ty yiunger ond tire vulneromle yeors ty fowher gove te si
84 -3931.989558103351
85 Iteration 2800 :
86 an ty yiunger ond tire vulneromle yeors ty fowher gove te si
87 -3921.829691740311
88 Iteration 2900 :
89 an ty yiunger ond tire vulneromle yeors ty fowher gove te si
90 -3921.0980737400137
91 Iteration 3000 :
92 an ty yiunger ond tire vulneromle yeors ty fowher gove te si
93 -3921.0980737400137
94 Iteration 3100 :
95 an ty yiunger ond tire vulneromle yeors ty fowher gove te si
96 -3921.0980737400137
97 Iteration 3200 :
98 an ty yiunger ond tire vulneromle yeors ty fowher gove te si
99 -3923.457207966006
100 Iteration 3300 :
101 an ty yiunger ond tire vulneromle yeors ty fowher gove te si
102 -3913.6758807342153
103 Iteration 3400 :
104 an ty yiunger ond tire vulneromle yeors ty fowher gove te si
105 -3904.838654218881
106 Iteration 3500 :
107 an ty yiunger ond tire vulneromle yeors ty fowher gove te si
108 -3904.838654218881
109 Iteration 3600 :
110 an ty yiunger ond tire vulneromle yeors ty fowher gove te si
111 -3906.044784164552
112 Iteration 3700 :
113 an ty yiunger ond tire vulneromle yeors ty fowher gove te si

```

```

114 -3903.915903707606
115 Iteration 3800 :
116 an wy yiunger ond wire vulneromle yeors wy fother gove we si
117 -3789.8702741883667
118 Iteration 3900 :
119 an wy yiunger ond wire vulneromle yeors wy fother gove we si
120 -3789.8702741883667
121 Iteration 4000 :
122 on wy yiunger and wire vulneramle years wy father gave we si
123 -3747.334993839508
124 Iteration 4100 :
125 on wy yiunger and wire vulneramle years wy father gave we si
126 -3748.0666118398053
127 Iteration 4200 :
128 on my yiunger and mire vulnerawle years my father gave me si
129 -3737.330370823384
130 Iteration 4300 :
131 on my yiunger and mire vulnerawle years my father gave me si
132 -3737.330370823384
133 Iteration 4400 :
134 on my yiunger and mire vulnerawle years my father gave me si
135 -3737.330370823384
136 Iteration 4500 :
137 on my yiunger and mire vulnerawle years my father gave me si
138 -3737.330370823384
139 Iteration 4600 :
140 on my yiunger and mire vulnerawle years my father gave me si
141 -3739.8215556333244
142 Iteration 4700 :
143 on my yiunger and mire vulnerawle years my father gave me si
144 -3739.8215556333244
145 Iteration 4800 :
146 on my yiunger and mire vulnerawle years my father gave me si
147 -3738.6947738045683
148 Iteration 4900 :
149 in my younger and more vulnerawle years my father gave me so
150 -3583.4950298706535
151 Iteration 5000 :
152 in my younger and more vulnerawle years my father gave me so
153 -3583.4950298706535
154 Iteration 5100 :
155 in my younger and more vulnerawle years my father gave me so
156 -3583.4950298706535
157 Iteration 5200 :
158 in my younger and more vulnerawle years my father gave me so
159 -3583.4950298706535
160 Iteration 5300 :
161 in my younger and more vulnerawle years my father gave me so
162 -3583.4950298706535
163 Iteration 5400 :
164 in my younger and more vulnerawle years my father gave me so
165 -3582.7634118703563
166 Iteration 5500 :
167 in my younger and more vulnerawle years my father gave me so
168 -3580.7622846745494
169 Iteration 5600 :
170 in my younger and more vulnerawle years my father gave me so
171 -3580.7622846745494
172 Iteration 5700 :
173 in my younger and more vulnerawle years my father gave me so
174 -3580.7622846745494
175 Iteration 5800 :
176 in my younger and more vulnerable years my father gave me so
177 -3526.706700483252
178 Iteration 5900 :
179 in my younger and more vulnerable years my father gave me so
180 -3526.706700483252
181 Iteration 6000 :
182 in my younger and more vulnerable years my father gave me so
183 -3526.706700483252

```

We also track the change of log-likelihood in Figure 6

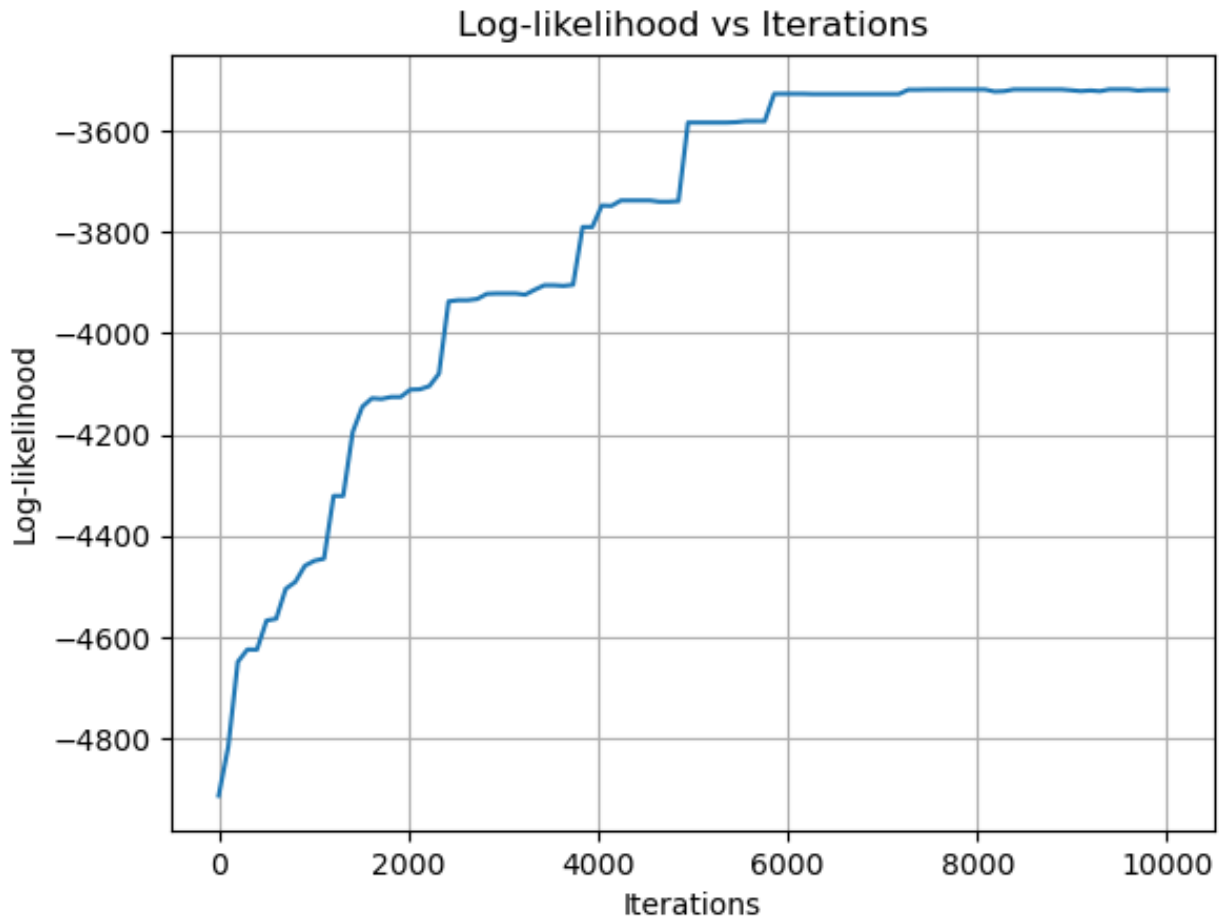


Figure 6: log-likelihood value for 10000 iterations

### Question 5.e

A Markov chain is called an ergodic chain if some power of the transition matrix has only positive elements. In other words, it can transit from any state to any state. In our transition matrix, we find zero values which means that we no longer hold the ergodic property. To solve this, we add a small value like 0.001 to our transition matrix before running MCMC. Thus, we restore the ergodicity. The code has already shown in 5(d) and it works good as we get some good starting points. Instead of negative infinite log-likelihood, the none-zero transition matrix avoid a lot of issues.

### Question 5.f

If we consider symbol probabilities alone, it would be insufficient for the English language. The structure of the language is influenced by sequences of letters, not just independent symbols. The interpretative ability is reduced when we rely on single-symbol probabilities. This approach would likely produce an output with little semantic coherence, as it disregards how symbols combine within words and phrases.

Using a second-order Markov chain can improve accuracy in text decoding. However, this increases the complexity, as it requires a three-dimensional transition matrix or tensor. Additionally, since we have successfully decoded messages with a first-order Markov chain, the benefit of using a second-order model may only result in a minor improvement in output.

If the encryption scheme permits multiple symbols to map to the same encrypted value, it introduces ambiguity in decoding. Our probability function is defined under the assumption of bijectivity, as are the algorithms that follow. Allowing non-bijective mappings leads to coding errors, causing the entire algorithm to fail.

For Chinese, with over 10,000 characters, the transition matrix would be vast and highly sparse, with most entries near zero. This extreme sparsity and high dimensionality would make many states nearly unreachable within finite

MCMC iterations, significantly limiting the model's efficiency and making decoding impractical.

## Question 7

### Question 7.a

Find the local extrema of the function  $f(x, y) := x + 2y$  subject to the constraint  $y^2 + xy = 1$ . To apply Lagrangian multiplier, one can write:

$$\nabla f = \lambda \nabla g. \quad (48)$$

Differentiating over  $x$  and  $y$  separately, we get:

$$\begin{aligned} 1 &= \lambda y, \\ 2 &= \lambda(2y + x). \end{aligned} \quad (49)$$

Together with the constraint, we have three equation for three unknown value. Starting with  $y = \frac{1}{\lambda}$  with  $\lambda \neq 0$ , we find:

$$\begin{aligned} 2 &= \lambda\left(\frac{2}{\lambda} + x\right) \Rightarrow x = 0, \\ y^2 &= 1 \Rightarrow y = \pm 1, \end{aligned} \quad (50)$$

where the local extrema are points  $(0, 1)$  and  $(0, -1)$ . Thus, the Lagrangian multiplier is  $\pm 1$  and the function value at local extrema is  $\pm 2$ .

### Question 7.b

Suppose we have a numerical routine to evaluate the exponential function  $\exp(x)$ . Using Newton's method, we can compute  $\ln(a)$ . Let  $x = \ln(a)$ , we the define

$$f(x, a) = e^x - a, \quad (51)$$

where the root of function is  $\ln(a)$ . Applying Newton's method, we have:

$$\begin{aligned} x_{n+1} &= x_n - \frac{f(x_n)}{f'(x_n)}, \\ \Rightarrow x_{n+1} &= x_n - \frac{e^{x_n} - a}{e^{x_n}}. \end{aligned} \quad (52)$$

The algorithm always converges as  $f'' > 0$ .

## Bonus Question 4

### Question 4.a

We follow the instruction of the question. Here are the plots for filter mode and smooth mode:

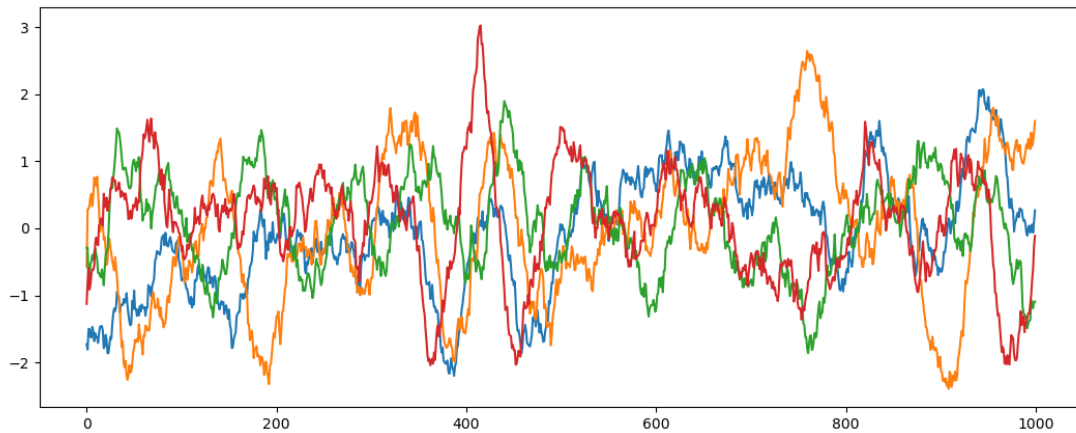


Figure 7: Kalman Filter

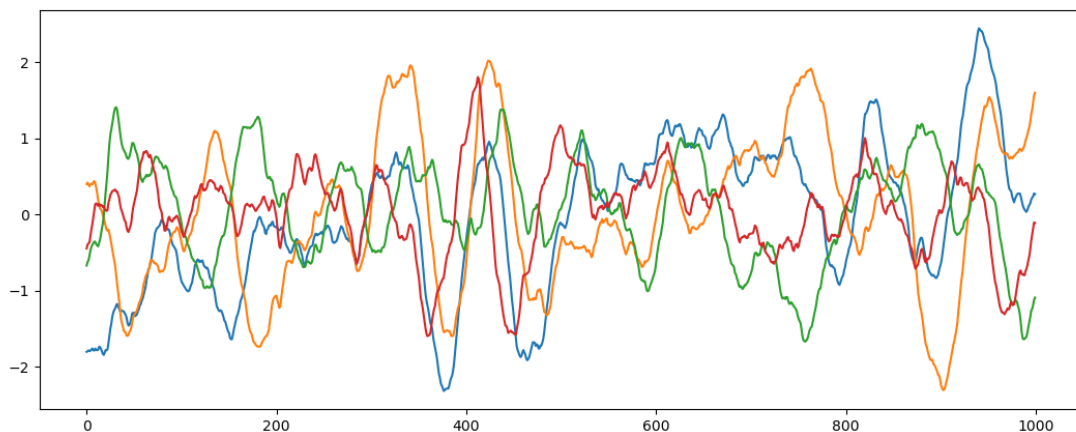


Figure 8: Kalman Smooth

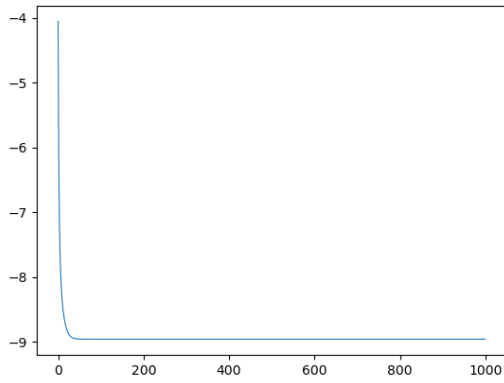
Figure 7 and Figure 8 have a similar trends. The filter mode has more details or zigzag in the plot and the smooth mode looks smooth.

For log determinant, we can find that they all converge quickly. We have the filter log-determinant of covariance at about -9 and smooth on at about -11.

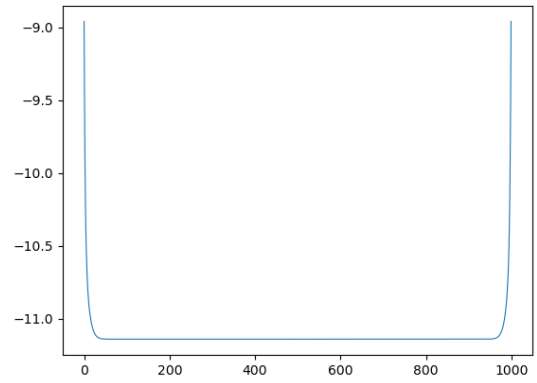
The difference can be explained if we look at the code. For filter mode, we stop at  $t$ . However, we have the whole dataset for smooth mode which yields a relatively smaller covariance and a smooth estimate.

The code is here

```
1 def logdet(A):
2     return 2 * np.sum(np.log(np.diag(np.linalg.cholesky(A))))
3 # %%
4 # load data
5 train = np.loadtxt('ssm_spins.txt')
6 # %%
7 # def parameter
8 theta1 = 2 * np.pi / 180
```



(a) Log-det(V) for Kalman Filter



(b) Log-det(V) for Kalman Smooth

Figure 9: Plot of log determinant for Covariances

```

9  theta2 = 2 * np.pi / 90
10 A = 0.99 * np.array([
11     [np.cos(theta1), -np.sin(theta1), 0, 0],
12     [np.sin(theta1), np.cos(theta1), 0, 0],
13     [0, 0, np.cos(theta2), -np.sin(theta2)],
14     [0, 0, np.sin(theta2), np.cos(theta2)]
15 ])
16 C = np.array([
17     [1, 0, 1, 0],
18     [0, 1, 0, 1],
19     [1, 0, 0, 1],
20     [0, 1, 1, 1],
21     [0.5, 0.5, 0.5, 0.5]
22 ])
23 R = np.eye(5)
24 Q = np.eye(4) - A @ A.T
25 # %%
26 x = train.T
27 Y0 = np.random.multivariate_normal(np.zeros(4), np.eye(4))
28 Q0 = np.eye(4)
29 Y, V, _, L = run_ssm_kalman(x, Y0, Q0, A, Q, C, R, "filt")
30 Ys, Vs, Vj, Ls = run_ssm_kalman(x, Y0, Q0, A, Q, C, R, mode='smooth')

```

Question 4(a)



## Bonus Question 6

## Bonus Question 8

### Question 8.a

Starting with extreme value theorem of calculus, a continuous function on a compact domain attains its maximum and minimum. In that case, we need to transform the original not-compact space to a compact space. From the hints, we know that unit sphere is good choice. Let us first focus on  $R_A$ :

$$\begin{aligned} R_A(x) &= \frac{x^\top A x}{x^\top x} = \frac{q_A(x)}{\|x\|^2}, \\ \Rightarrow R_A(x) &= \frac{x^\top}{\|x\|} A \frac{x}{\|x\|}. \end{aligned} \quad (53)$$

One can define  $y = \frac{x}{\|x\|}$ , and its normal can be write as:

$$\|y\| = \sqrt{y^\top y} = \sqrt{\frac{x^\top}{\|x\|} \frac{x}{\|x\|}} = 1, \quad (54)$$

which satisfies a unit sphere. Thus, we can rewrite  $R_A$  in terms of  $y$ :

$$R_A(y) = y^\top A y \quad \text{for } y \in \mathbb{S}, \quad (55)$$

where its domain is compact. Finally, we attain  $\sup_{x \in \mathbb{R}^n} R_A(x)$ .

### Question 8.b

Now, we substitute  $x = \sum_{i=1}^n (\xi_i^\top x) \xi_i$  into  $R_A$ :

$$\begin{aligned} R_A(x) &= \frac{(\sum_{i=1}^n (\xi_i^\top x) \xi_i)^\top A (\sum_{i=1}^n (\xi_i^\top x) \xi_i)}{(\sum_{i=1}^n (\xi_i^\top x) \xi_i)^\top (\sum_{i=1}^n (\xi_i^\top x) \xi_i)}, \\ \Rightarrow R_A(x) &= \frac{(\sum_{i=1}^n (\xi_i^\top x) \xi_i)^\top (\sum_{i=1}^n (\xi_i^\top x) A \xi_i)}{(\sum_{i=1}^n (\xi_i^\top x) \xi_i)^\top (\sum_{i=1}^n (\xi_i^\top x) \xi_i)}. \end{aligned} \quad (56)$$

Using  $A \xi_i = \lambda_i \xi_i$ , we can get:

$$R_A(x) = \frac{(\sum_{i=1}^n (\xi_i^\top x) \xi_i)^\top (\sum_{i=1}^n (\xi_i^\top x) \lambda_i \xi_i)}{(\sum_{i=1}^n (\xi_i^\top x) \xi_i)^\top (\sum_{i=1}^n (\xi_i^\top x) \xi_i)}, \quad (57)$$

where we can use ONB property to reduce the summations. Since  $\lambda_1$  is the biggest eigenvalues, we can deduce:

$$R_A(x) = \frac{\sum_{i=1}^n \lambda_i ((\xi_i^\top x) \xi_i)^\top ((\xi_i^\top x) \xi_i)}{\sum_{i=1}^n ((\xi_i^\top x) \xi_i)^\top ((\xi_i^\top x) \xi_i)} \leq \lambda_1 \times \frac{\sum_{i=1}^n ((\xi_i^\top x) \xi_i)^\top ((\xi_i^\top x) \xi_i)}{\sum_{i=1}^n ((\xi_i^\top x) \xi_i)^\top ((\xi_i^\top x) \xi_i)}, \quad (58)$$

and this is

$$R_A(x) \leq \lambda_1. \quad (59)$$

(c) We can use the similar process in (b). Instead of summing from  $i$  to  $n$ , we will remove the span  $\{\xi_1, \dots, \xi_k\}$ .

$$R_A(x) = \frac{\sum_{i=k+1}^n \lambda_i ((\xi_i^\top x) \xi_i)^\top ((\xi_i^\top x) \xi_i)}{\sum_{i=k+1}^n ((\xi_i^\top x) \xi_i)^\top ((\xi_i^\top x) \xi_i)}. \quad (60)$$

This equation must be equal or less than  $\lambda_{k+1}$ . From the question, we know that  $\lambda_1$  is strictly larger than  $\lambda_i$  for  $i = k+1, \dots, n$ . Thus, we have

$$R_A(x) \leq \lambda_{k+1} < \lambda_1. \quad (61)$$