

Probabilistic and Unsupervised Learning Summative Assignment  
COMP0086

Richard Huang

Nov 13 2024

# Question 1

(a) The dataset contains binary images where each pixel is either 0 or 1, representing black or white, respectively. A multivariate Gaussian distribution, however, assumes continuous values for each variable. Specifically, each element in a multivariate Gaussian vector can take any real value, which doesn't match the binary nature of this dataset.

Additionally, a multivariate Gaussian is defined by mean and covariance parameters and is used for data that tends to follow a continuous, bell-shaped distribution in each dimension. Since binary data (0 or 1) does not follow a continuous distribution, a Gaussian model would not capture the true distribution of the binary image data effectively. Therefore, a multivariate Gaussian model is not suitable for modeling binary vectors like these images.

(b) To find the ML estimate of  $\mathbf{p}$ , we maximize the likelihood of the data given  $\mathbf{p}$ . The likelihood function for a single sample  $\mathbf{x}$  is given by:

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

For  $N$  images, the likelihood of the entire dataset becomes:

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

Taking the logarithm to simplify, we get the log-likelihood:

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

Define  $s_d = \sum_{n=1}^N x_d^{(n)}$ , the total count of 1s at pixel  $d$  across all images. Then, the log-likelihood simplifies to:

$$\log P(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}|\mathbf{p}) = \sum_{d=1}^D (s_d \log p_d + (N - s_d) \log(1 - p_d))$$

To maximize with respect to  $p_d$ , take the derivative of the log-likelihood with respect to  $p_d$  and set it to zero:

$$\frac{\partial}{\partial p_d} \log P = \frac{s_d}{p_d} - \frac{N - s_d}{1 - p_d} = 0$$

Solving this equation for  $p_d$ , we find that:

$$p_{\text{MLE}} = \frac{1}{N} \sum_{n=1}^N x_d^{(n)}$$

Thus, the maximum likelihood estimate for each pixel parameter  $p_{\text{MLE}}$  is given by:

$$p_{\text{MLE}} = \frac{1}{N} \sum_{n=1}^N x_d^{(n)}$$

(c) To find the maximum a posteriori (MAP) estimator for  $\mathbf{p}$ , we assume independent Beta priors for each  $p_d$  with parameters  $\alpha$  and  $\beta$ . The prior distribution for each  $p_d$  is given by:

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

where  $B(\alpha, \beta)$  is the Beta function. The prior distribution for  $\mathbf{p}$  is:

$$P(\mathbf{p}) = \prod_{d=1}^D P(p_d)$$

The likelihood function for the observed data is:

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

The posterior distribution  $P(\mathbf{p}|\mathbf{x})$  is proportional to the product of the likelihood and the prior:

$$P(\mathbf{p}|\mathbf{x}) \propto P(\mathbf{x}|\mathbf{p})P(\mathbf{p})$$

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

Simplifying for each  $p_d$ , we have:

$$P(p_d|\mathbf{x}) \propto p_d^{\sum_{n=1}^N x_d^{(n)} + \alpha - 1} (1-p_d)^{N - \sum_{n=1}^N x_d^{(n)} + \beta - 1}$$

To find the MAP estimate, we maximize  $P(p_d|\mathbf{x})$  with respect to  $p_d$ , which is equivalent to maximizing the logarithm of  $P(p_d|\mathbf{x})$ :

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

Taking the derivative with respect to  $p_d$  and setting it to zero:

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

Solving for  $p_d$ , we get:

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

Expanding and simplifying:

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

Thus, the MAP estimate for  $p_d$  is:

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

Therefore, the MAP estimator for  $\mathbf{p}$  is:

$$p_d^{\text{MAP}} = \frac{\sum_{n=1}^N x_d^{(n)} + \alpha - 1}{N + \alpha + \beta - 2} \quad \text{for each } d = 1, \dots, D.$$

(d) The visualization of the learned parameter vector for question(d) is shown in Figure 1

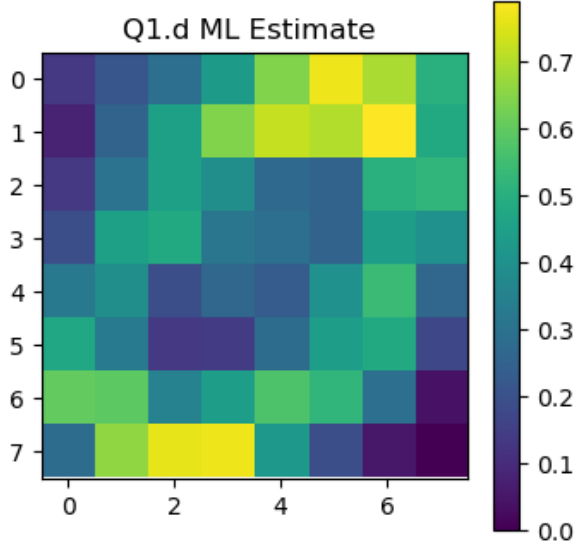


Figure 1: ML Parameters

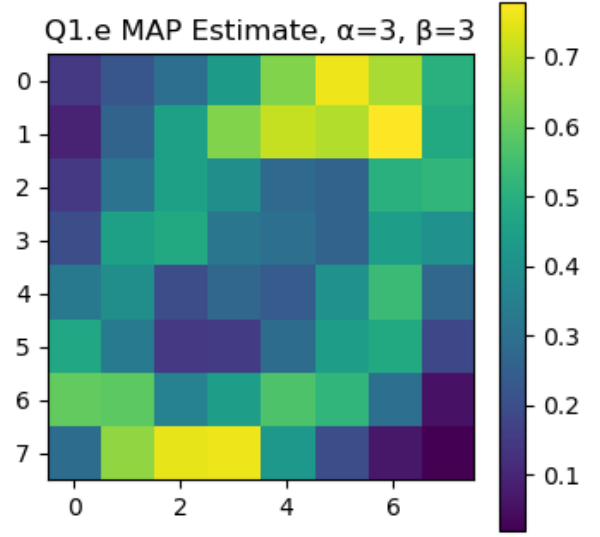


Figure 2: MAP Parameters

(e) The visualization of the learned parameter vector for question(d) is shown in Figure 2

The ML estimate for the multivariate Bernoulli model calculates each parameter  $p_d$  as the proportion of times the pixel is 1 across all images in the dataset. This approach directly reflects the observed data, giving an estimate that is unbiased by any prior beliefs. However, this estimate can lead to extreme values (close to 0 or 1) for pixels that are consistently off or on, particularly in small datasets. These extreme values might result in overfitting, where the model captures noise rather than general patterns.

In contrast, the MAP estimate incorporates a prior, specifically a Beta distribution with parameters  $\alpha = 3$  and  $\beta = 3$ , which introduces regularization. This prior pulls each parameter  $p_d$  toward 0.5, a value indicating no strong preference for on or off. The MAP estimate is thus less influenced by rare occurrences or noise in the data, providing more conservative probabilities. This regularization effect is particularly useful in cases where there is limited data for certain pixels, as the prior stabilizes the estimate and prevents extreme values.

While MAP regularization helps avoid overfitting and can improve generalization to new data, it introduces bias. The prior's pull toward 0.5 means that the MAP estimate may understate genuine patterns in the data, especially when certain pixels are consistently on or off. In large datasets, where data patterns are clear and representative of the true distribution, the ML estimate is typically more accurate, as it directly reflects the observed data without the influence of a prior. However, in smaller or noisier datasets, the MAP estimate is often preferable, as it offers a balanced estimate that mitigates overfitting risks. Thus, while the MAP estimate may be more robust in uncertain data conditions, it can underrepresent true patterns compared to the ML estimate when ample data is available.

## Python code for d&e:

```
1
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 #Load Data
6 data = np.loadtxt('/Users/richardhuang/Documents/GitHub/Projects/Probabilistic
7 Summative/binarydigits.txt')
8
9 def learn_ml_params(data):
10     """
11     Learn the Maximum Likelihood (ML) parameters for a multivariate Bernoulli distribution.
12     This function computes the ML estimate and reshapes it to an 8x8 image format.
13     """
14     p_mle = np.mean(data, axis=0)
15     return p_mle.reshape(8, 8)
16
17 def learn_map_params(data, alpha=3, beta=3):
18     """
19     Learn the Maximum A Posteriori (MAP) parameters for a multivariate Bernoulli distribution
20     with Beta prior parameters alpha and beta. This function computes the MAP estimate and
21     reshapes it to an 8x8 image format.
22     """
23     N = data.shape[0]
24     pixel_sums = np.sum(data, axis=0)
25     p_map = (pixel_sums + alpha - 1) / (N + alpha + beta - 2)
26     return p_map.reshape(8, 8)
27
28 #Calculate ML and MAP parameters
29 p_mle_image = learn_ml_params(data)
30 p_map_image = learn_map_params(data, alpha=3, beta=3)
31
32 plt.figure(figsize=(4, 4))
33 plt.imshow(p_mle_image, cmap='viridis', interpolation='nearest')
34 plt.colorbar()
35 plt.title("Q1.d ML Estimate ")
36 plt.show()
37
38 plt.figure(figsize=(4, 4))
39 plt.imshow(p_map_image, cmap='viridis', interpolation='nearest')
40 plt.colorbar()
41 plt.title("Q1.e MAP Estimate, \alpha = 3, \beta=3")
42 plt.show()
43
44 plt.tight_layout()
45 plt.show()
```

Listing 1: solutions/q1.py

## Question 2

(a) In model (a), we assume that every pixel in the image is generated by a Bernoulli distribution with the same parameter  $p_d = 0.5$  for all pixels. The probability mass function for a Bernoulli distribution is:

$$P(x_d|p_d) = p_d^{x_d} (1 - p_d)^{1-x_d}$$

where  $x_d \in \{0, 1\}$  represents the value of the  $d$ -th pixel in the image and  $p_d = 0.5$  is the probability of the pixel being 1.

Since  $p_d = 0.5$  for all pixels, the likelihood of observing a specific image  $\mathbf{x}^{(n)}$  with  $D$  pixels is:

$$P(\mathbf{x}^{(n)}|p_d = 0.5) = \prod_{d=1}^D 0.5 = 0.5^D$$

The total likelihood for the entire dataset, which consists of  $N$  independent images, is:

$$P(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}|p_d = 0.5) = \prod_{n=1}^N P(\mathbf{x}^{(n)}|p_d = 0.5) = (0.5^D)^N = 0.5^{N \times D}$$

Thus, the likelihood of the dataset for model (a) is:

$$P(\mathbf{x}|p) = 0.5^{N \times D}$$

(b) In model (b), we assume that each pixel in the image is generated by a Bernoulli distribution with an unknown parameter  $p_d$ , which is the same for all pixels in the image. The likelihood of observing an image  $\mathbf{x}^{(n)}$  with  $D$  pixels is:

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

Since the images are independent, the total likelihood for the entire dataset of  $N$  images is:

$$P(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}|p_d) = \prod_{n=1}^N \prod_{d=1}^D p_d^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}}$$

Thus, the likelihood of the dataset for model (b) is:

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

(c) In model (c), we assume that each pixel has its own separate Bernoulli distribution with its own unknown parameter  $p_d$ . The likelihood of observing an image  $\mathbf{x}^{(n)}$  is:

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

The total likelihood for the entire dataset of  $N$  images is:

$$P(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)} | \mathbf{p}) = \prod_{n=1}^N \prod_{d=1}^D p_d^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}}$$

Thus, the likelihood of the dataset for model (c) is:

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

**Posterior probabilities** Since the priors for the models are uniform, we have:

$$P(M_a) = P(M_b) = P(M_c) = \frac{1}{3}$$

The marginal likelihood  $P(\mathbf{x})$  is computed as:

$$P(\mathbf{x}) = P(\mathbf{x} | M_a) \cdot P(M_a) + P(\mathbf{x} | M_b) \cdot P(M_b) + P(\mathbf{x} | M_c) \cdot P(M_c)$$

Given the three models  $M_a$ ,  $M_b$ , and  $M_c$ , the posterior probabilities can be computed using Bayes' Theorem:

$$P(M_i | \mathbf{x}) = \frac{P(\mathbf{x} | M_i) \cdot P(M_i)}{P(\mathbf{x})}$$

Where the evidence  $P(\mathbf{x})$  is:

$$P(\mathbf{x}) = P(\mathbf{x} | M_a) \cdot P(M_a) + P(\mathbf{x} | M_b) \cdot P(M_b) + P(\mathbf{x} | M_c) \cdot P(M_c)$$

Thus, the posterior probabilities for each model are:

$$P(M_a | \mathbf{x}) = \frac{P(\mathbf{x} | M_a) \cdot P(M_a)}{P(\mathbf{x})}$$

$$P(M_b | \mathbf{x}) = \frac{P(\mathbf{x} | M_b) \cdot P(M_b)}{P(\mathbf{x})}$$

$$P(M_c | \mathbf{x}) = \frac{P(\mathbf{x} | M_c) \cdot P(M_c)}{P(\mathbf{x})}$$

Where the likelihoods for each model are as follows: - For model  $M_a$ ,  $P(\mathbf{x} | M_a) = 0.5^{N \times D}$ , - For model  $M_b$ ,  $P(\mathbf{x} | M_b) = \prod_{d=1}^D p_d^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}}$ , - For model  $M_c$ ,  $P(\mathbf{x} | M_c) = \prod_{d=1}^D p_d^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}}$ .

Where  $p_d$  is the mean of the data for each pixel in the case of models (b) and (c).

The posterior probabilities for each model, calculated using Bayes' theorem, are as follows:

$$P(M_a | \mathbf{x}) = 0.37434997449465635$$

$$P(M_b | \mathbf{x}) = 0.3128250127526719$$

$$P(M_c | \mathbf{x}) = 0.3128250127526717$$

Model	Posterior Probability
Model (a)	0.37434997449465635
Model (b)	0.3128250127526719
Model (c)	0.3128250127526717

Table 1: Posterior probabilities for each model

## Python code for Q2:

```
1 N, D = data.shape
2
3 # Model (a) - Log likelihood
4 P_x_given_Ma = N * D * np.log(0.5)
5
6
7 # Model (b) - Log likelihood
8 P_x_given_Mb = np.sum(np.log(np.power(p_d_b, np.sum(data, axis=0)) *
9 np.power(1 - p_d_b, N - np.sum(data, axis=0))))
10
11 # Model (c) - Log likelihood
12 P_x_given_Mc = np.sum(np.sum(np.log(np.power(p_d_b, data) * np.power(1 - p_d_b, 1 - data)),
13 axis=0))
14
15 # Priors
16 P_Ma = P_Mb = P_Mc = 1/3
17
18 # Evidence (marginal likelihood)
19 P_x = P_x_given_Ma * P_Ma + P_x_given_Mb * P_Mb + P_x_given_Mc * P_Mc
20
21 if P_x == 0:
22     print("Warning: The evidence (P_x) is zero!")
23
24 # Posterior probabilities using Bayes' theorem
25 P_Ma_given_x = (P_x_given_Ma * P_Ma) / P_x if P_x != 0 else 0
26 P_Mb_given_x = (P_x_given_Mb * P_Mb) / P_x if P_x != 0 else 0
27 P_Mc_given_x = (P_x_given_Mc * P_Mc) / P_x if P_x != 0 else 0
28
29
30 print(f"Posterior probability for Model (a): {P_Ma_given_x}")
31 print(f"Posterior probability for Model (b): {P_Mb_given_x}")
32 print(f"Posterior probability for Model (c): {P_Mc_given_x}")
```

Listing 2: solutions/q1.py



# Question 3

(a) To model a mixture of  $K$  multivariate Bernoulli distributions for the binary image dataset, we define the following parameters:

$\pi_1, \pi_2, \dots, \pi_K$ : the mixing proportions for the  $K$  components, where  $0 \leq \pi_k \leq 1$  and  $\sum_{k=1}^K \pi_k = 1$ . -  $\mathbf{P}$ : a  $K \times D$  matrix where each element  $p_{kd}$  represents the probability that pixel  $d$  takes the value 1 under mixture component  $k$ .

The likelihood for a single image  $\mathbf{x} = (x_1, x_2, \dots, x_D)$  under this mixture model is given by:

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

Since the images are independent and identically distributed (i.i.d.) under the model, the likelihood of the entire dataset  $X = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$  is:

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

(b) To compute the responsibility of mixture component  $k$  for data vector  $\mathbf{x}^{(n)}$ , we define the responsibility  $r_{nk}$  as the posterior probability that the hidden variable  $s^{(n)} = k$  given  $\mathbf{x}^{(n)}, \pi, \mathbf{P}$ .

This is expressed as:

$$r_{nk} \equiv P(s^{(n)} = k|\mathbf{x}^{(n)}, \pi, \mathbf{P}) = \frac{P(s^{(n)} = k|\pi)P(\mathbf{x}^{(n)}|s^{(n)} = k, \mathbf{P})}{\sum_{j=1}^K P(s^{(n)} = j|\pi)P(\mathbf{x}^{(n)}|s^{(n)} = j, \mathbf{P})}$$

Since  $P(s^{(n)} = k|\pi) = \pi_k$  and  $P(\mathbf{x}^{(n)}|s^{(n)} = k, \mathbf{P}) = \prod_{d=1}^D p_{kd}^{x_d^{(n)}} (1 - p_{kd})^{1-x_d^{(n)}}$ , we can rewrite  $r_{nk}$  as:

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

This expression represents the responsibility of component  $k$  for data point  $\mathbf{x}^{(n)}$ , and it forms the E-step in the Expectation-Maximization (EM) algorithm. During this E-step, we compute the responsibilities  $r_{nk}$  for each data point and each component  $k$  given the current estimates of the parameters  $\pi$  and  $\mathbf{P}$ .

(c) To find the maximizing parameters for the expected log-joint likelihood with respect to the mixing proportions  $\pi$  and the Bernoulli parameters  $\mathbf{P}$ , we can proceed as follows:

The goal is to maximize the expected log-joint likelihood, which is given by:

$$\mathbb{E}_{q(\{s^{(n)}\})} \left[ \sum_n \log P(\mathbf{x}^{(n)}, s^{(n)}|\pi, \mathbf{P}) \right] = \sum_n \sum_{k=1}^K r_{nk} \log P(\mathbf{x}^{(n)}, s^{(n)} = k|\pi, \mathbf{P})$$

where  $r_{nk} = P(s^{(n)} = k|\mathbf{x}^{(n)}, \pi, \mathbf{P})$  is the responsibility assigned to component  $k$  for data point  $\mathbf{x}^{(n)}$ .

Maximizing with Respect to  $\pi$  (Mixing Proportions): The term involving  $\pi_k$  is:

$$\sum_n \sum_{k=1}^K r_{nk} \log \pi_k$$

To maximize with respect to  $\pi_k$  (subject to  $\sum_{k=1}^K \pi_k = 1$ ), we use a Lagrange multiplier  $\lambda$  to enforce this constraint. The Lagrange function is:

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

Taking the partial derivative of  $\mathcal{L}$  with respect to  $\pi_k$  and setting it to zero gives:

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

Rearranging, we find:

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

To satisfy the constraint  $\sum_{k=1}^K \pi_k = 1$ , we substitute  $\pi_k = \frac{\sum_{n=1}^N r_{nk}}{\lambda}$  into the constraint, which gives:

$$\lambda = N$$

Therefore, the update rule for  $\pi_k$  becomes:

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

Maximizing with Respect to  $\mathbf{P}$  (Bernoulli Parameters): For each Bernoulli parameter  $p_{kd}$ , the relevant term in the expected log-joint likelihood is:

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

To maximize with respect to  $p_{kd}$ , take the partial derivative with respect to  $p_{kd}$  and set it to zero:

$$\frac{\partial}{\partial p_{kd}} \sum_n \sum_{k=1}^K r_{nk} \left( x_d^{(n)} \log p_{kd} + (1 - x_d^{(n)}) \log(1 - p_{kd}) \right) = 0$$

Solving this equation gives the update rule:

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

This expression updates  $p_{kd}$  as the weighted average of the pixel values  $x_d^{(n)}$  for all data points, where the weights are given by the responsibilities  $r_{nk}$ .

The M-step of the EM algorithm for this model involves updating the parameters  $\pi$  and  $\mathbf{P}$  using the following rules:

$$\begin{aligned} \pi_k &= \frac{\sum_{n=1}^N r_{nk}}{N} \\ p_{kd} &= \frac{\sum_{n=1}^N r_{nk} x_d^{(n)}}{\sum_{n=1}^N r_{nk}} \end{aligned}$$

These updates maximize the expected log-joint likelihood given the responsibilities computed in the E-step.

(d) Figure 3 displays the convergence of the log-likelihood over iterations for different values of  $K$  in the EM algorithm. Each subplot corresponds to a distinct value of  $K$ , specifically  $K = 2, 3, 4, 7$ , and  $10$ , as labeled. For each  $K$  value, the log-likelihood is observed to increase as the number of iterations progresses, eventually reaching a plateau, which indicates that the EM algorithm has found a stable set of parameters and converged. This behavior is consistent across all values of  $K$ , with convergence generally achieved within approximately 15–20 iterations.

As  $K$  increases, the final log-likelihood at convergence improves. For example, the log-likelihood converges around -3300 for  $K = 2$ , whereas it reaches approximately -2400 for  $K = 10$ . This increase in the log-likelihood with higher values of  $K$  is expected, as models with more components possess greater flexibility to fit the data distribution more closely. Therefore, higher values of  $K$  allow the model to capture more patterns within the data, resulting in a better fit.

However, while the increasing log-likelihood values with higher  $K$  suggest an improved fit, this does not imply that the model with the highest  $K$  (here  $K = 10$ ) is necessarily the best. Increasing  $K$  introduces more parameters into the model, which can lead to overfitting, where the model captures noise rather than genuine underlying patterns in the data.

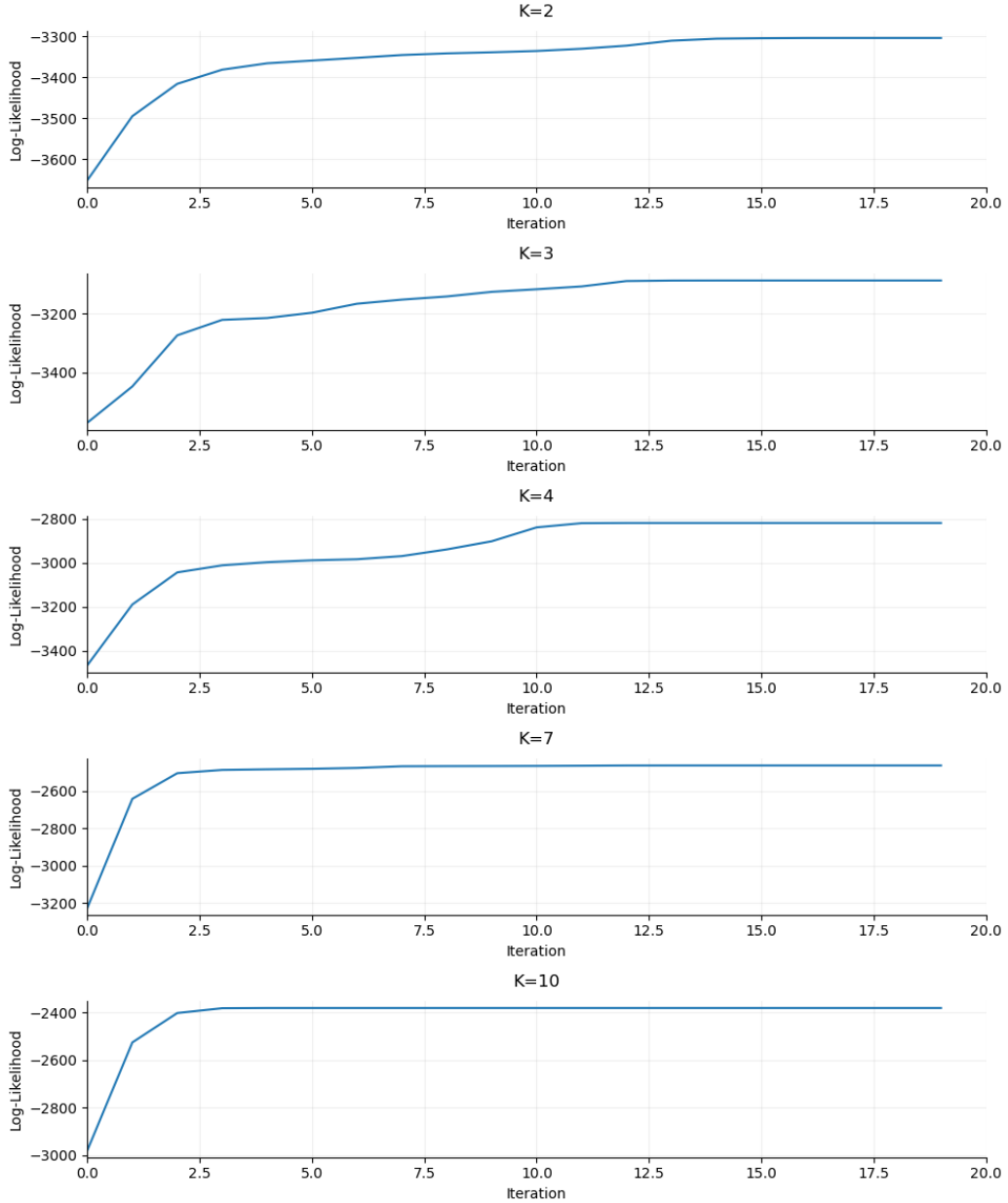


Figure 3: Log likelihood Function

The Figure 4 shows the learned parameters of a mixture model with varying values of  $K$  (the number of components). Each row corresponds to a specific  $K$  value, with each component represented by an  $8 \times 8$  grid displaying pixel probabilities. Higher values (in yellow) indicate a greater likelihood of the pixel being "on." Each component's mixing proportion,  $\pi$ , is shown above the grid. As  $K$  increases, the model captures more diverse patterns, with finer structural details becoming apparent. However, for higher  $K$ , some components have low mixing proportions, indicating they contribute minimally to the mixture.

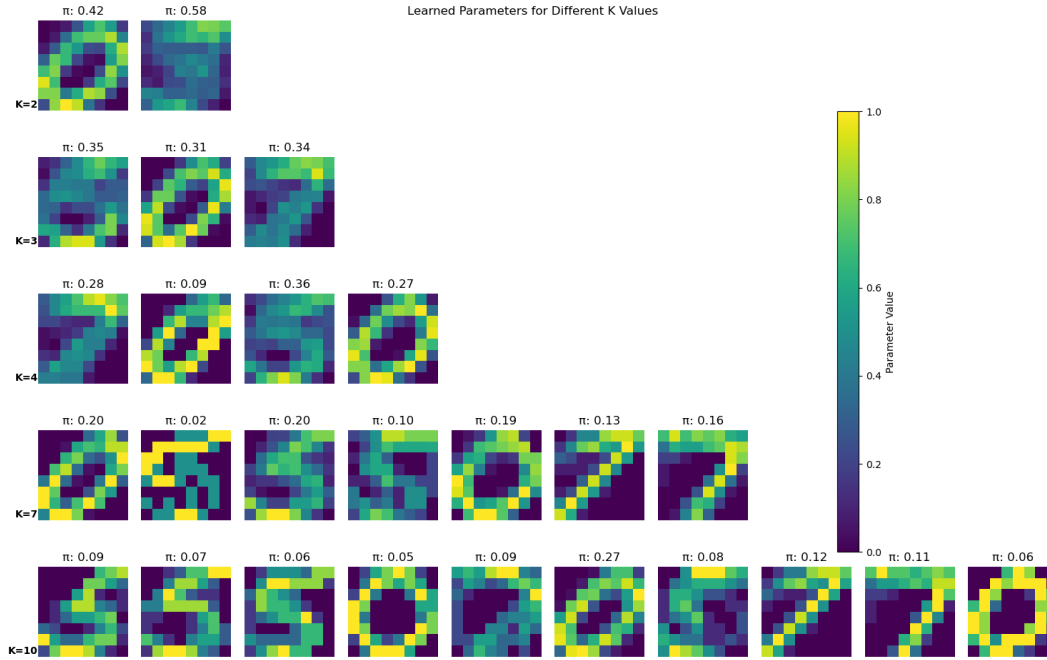


Figure 4: EM Learned Parameters for Different K Values

Python code for (d) (Log likelihood vs iteration number with different  $K$ ):

```

1 # Initialize the parameters for EM algorithm
2 def initialize_parameters(K, D):
3     # Randomly initialize mixing proportions and ensure they sum to 1
4     pi = np.random.dirichlet(alpha=np.ones(K), size=1)[0]
5
6     P = 0.2 + 0.6 * np.random.rand(K, D)
7     return pi, P
8
9 # Calculate the log-likelihood of the data
10 def compute_log_likelihood(X, pi, P):
11     N, D = X.shape
12     K = pi.shape[0]
13     log_likelihood = 0
14     for n in range(N):
15         log_prob_sum = np.log(np.sum([
16             pi[k] * np.prod(P[k]**X[n] * (1 - P[k])** (1 - X[n])) for k in range(K)
17         ]))
18         log_likelihood += log_prob_sum
19     return log_likelihood
20
21 # E-step: calculate responsibilities
22 def e_step(X, pi, P):
23     N, D = X.shape
24     K = pi.shape[0]
25     responsibilities = np.zeros((N, K))
26
27     for n in range(N):
28         log_resps = np.array([

```

```

29         np.log(pi[k] + 1e-10) + np.sum(X[n] * np.log(P[k] + 1e-10) + (1 - X[n])
30         * np.log(1 - P[k] + 1e-10))
31         for k in range(K)
32     ])
33     max_log_resp = np.max(log_resps)
34     log_resps -= max_log_resp
35     responsibilities[n, :] = np.exp(log_resps)
36     responsibilities[n, :] /= np.sum(responsibilities[n, :])
37
38     return responsibilities
39
40 # M-step: update parameters based on responsibilities
41 def m_step(X, responsibilities):
42     N, D = X.shape
43     K = responsibilities.shape[1]
44     pi_new = np.sum(responsibilities, axis=0) / N
45     P_new = np.dot(responsibilities.T, X) / np.sum(responsibilities, axis=0)[:, None]
46     return pi_new, P_new
47
48 # EM algorithm
49 def plot_log_likelihoods_subplots(data, K_values):
50     """Plot log-likelihoods in the specified format"""
51     fig, axes = plt.subplots(len(K_values), 1, figsize=(10, 12))
52
53
54
55     for idx, K in enumerate(K_values):
56         # Run EM
57         pi, P, log_likelihoods = em_algorithm(data, K, max_iter=20)
58
59         # Plot
60         ax = axes[idx]
61         ax.plot(np.arange(len(log_likelihoods)), log_likelihoods, '-', linewidth=1.5,
62                 color='#1f77b4')
63
64         ax.set_title(f'K={K}', pad=10)
65         ax.set_xlabel('Iteration')
66         ax.set_ylabel('Log-Likelihood')
67
68         ax.set_xlim(0, 50)
69
70         ax.grid(True, linestyle='-', alpha=0.2)
71         ax.spines['top'].set_visible(False)
72         ax.spines['right'].set_visible(False)
73
74         ax.tick_params(direction='out')
75
76     plt.tight_layout()
77     plt.show()
78
79 # Modified EM algorithm to ensure proper convergence behavior
80 def em_algorithm(X, K, max_iter=20):
81     """EM algorithm with fixed number of iterations"""
82     pi, P = initialize_parameters(K, X.shape[1])
83     log_likelihoods = []
84
85     for _ in range(max_iter):
86         # E-step
87         resp = e_step(X, pi, P)
88
89         # M-step
90         pi, P = m_step(X, resp)
91
92         # Compute and store log likelihood
93         log_likelihood = compute_log_likelihood(X, pi, P)
94         log_likelihoods.append(log_likelihood)
95
96     return pi, P, log_likelihoods
97
98 # Run with the specified K values
99 data = np.loadtxt('binarydigits.txt')

```

```

99 plot_log_likelihoods_subplots(data, [2, 3, 4, 7, 10])
100

```

Listing 3: src/solutions/q1.py

## Python code for (d) (Display parameters found):

```

1 plt.figure(figsize=(15, 10))
2
3
4 max_K = 10
5 subplot_idx = 1
6
7 for K in [2, 3, 4, 7, 10]:
8     # Run EM algorithm
9     pi, P, log_likelihoods = em_algorithm(data, K, 20)
10
11     # Plot each component for current K
12     for k in range(K):
13         subplot_pos = (subplot_idx - 1) * max_K + k + 1
14         plt.subplot(5, max_K, subplot_pos)
15
16         # Display the image
17         component_params = P[k].reshape(8, 8)
18         im = plt.imshow(component_params,
19                         cmap='viridis',
20                         vmin=0,
21                         vmax=1,
22                         interpolation='nearest')
23         plt.axis('off')
24
25         # Add pi value above each component
26         plt.title(f' : {pi[k]:.2f}')
27
28         # Add K value to bottom-left corner of first component in each row
29         if k == 0:
30             plt.text(-0.5, 7, f'K={K}',
31                     color='black',
32                     fontweight='bold',
33                     fontsize=10,
34                     ha='right',
35                     va='center')
36
37         subplot_idx += 1
38
39 plt.colorbar(im, ax=plt.gcf().axes, shrink=0.8, label='Parameter Value')
40
41 plt.suptitle('Learned Parameters for Different K Values', y=0.95)
42 plt.tight_layout()
43 plt.show()

```

Listing 4: src/solutions/q1.py

(e) Running the algorithm a few times starting from randomly chosen initial conditions and the parameters are shown below:

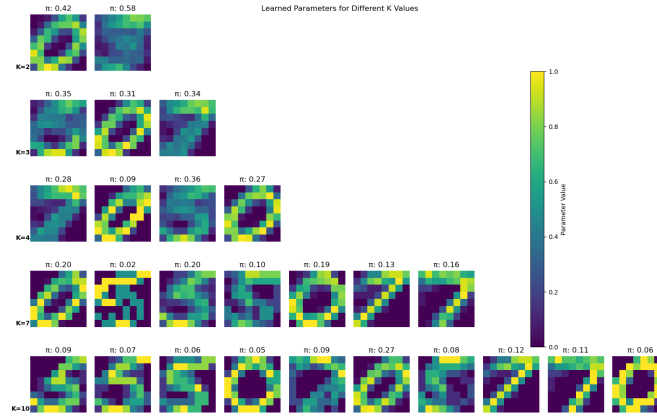


Figure 5: EM Learned Parameters Trial 0

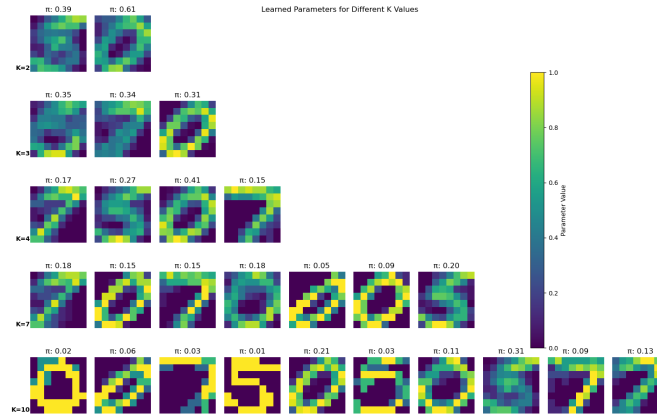


Figure 6: EM Learned Parameters Trial 1

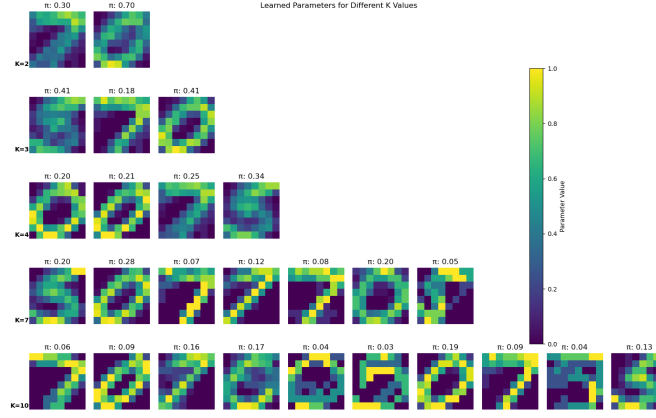


Figure 7: EM Learned Parameters Trial 2

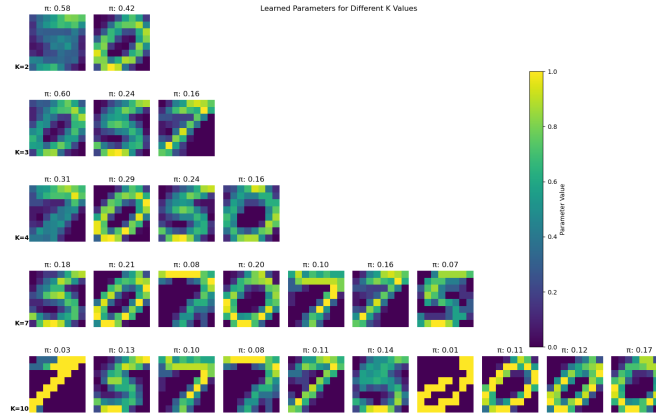


Figure 8: EM Learned Parameters Trial 3

The learned parameters from the EM algorithm indicate varying degrees of distinctiveness across different values of  $K$ . For lower values of  $K$  (e.g.,  $K = 2$  and  $K = 3$ ), the clusters exhibit recognizable patterns, suggesting that the algorithm effectively identifies distinct groups in the data. The images for these clusters show clear features, indicating successful clustering.

As  $K$  increases (e.g., from  $K = 4$  to  $K = 10$ ), the clusters become more complex and less distinct. This may indicate that the algorithm is attempting to fit noise in the data rather than capturing meaningful patterns, leading to less interpretable results.

The mixing proportions provide insight into the relative size of each cluster:

For  $K = 2$ , one cluster has a significantly higher proportion ( $\pi \approx 0.66$ ) compared to the other ( $\pi \approx 0.34$ ). This suggests that the algorithm has identified a dominant pattern in the data.

As  $K$  increases, the proportions become more evenly distributed, which can be a sign of overfitting, where the model attempts to create too many clusters for the available data.

The responsibilities indicate how likely each data point belongs to each cluster. If the responsibilities are concentrated (i.e., a data point has a high probability of belonging to one cluster), it suggests that the clusters are well-defined. Conversely, if responsibilities are spread out, it may indicate overlapping clusters or poor separation, which can complicate the interpretation of the results.

Each learned parameter matrix can be interpreted as a representation of specific features or patterns in the data. For example, if the data consists of binary images of digits, certain clusters may represent specific digits or variations of digits.

The patterns in the learned parameters can be analyzed to understand what features are being captured by each cluster. Clusters that show circular patterns may correspond to digits like '0' or '6'.

To enhance the performance of the EM algorithm, the following improvements can be considered:

- Better Initialization: Using k-means clustering to initialize the parameters can lead to better convergence and



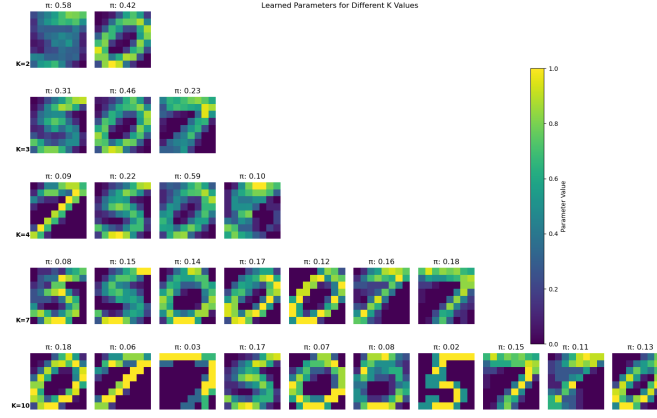


Figure 9: EM Learned Parameters Trial 4

more distinct clusters.

- **Regularization:** Adding regularization to the parameters can help prevent overfitting, especially if the model is too complex for the data.
- **Feature Engineering:** Transforming the data or adding features may help the algorithm better capture the underlying structure.
- **More Iterations:** Allowing for more iterations or implementing a more sophisticated convergence criterion can help the algorithm find better solutions.

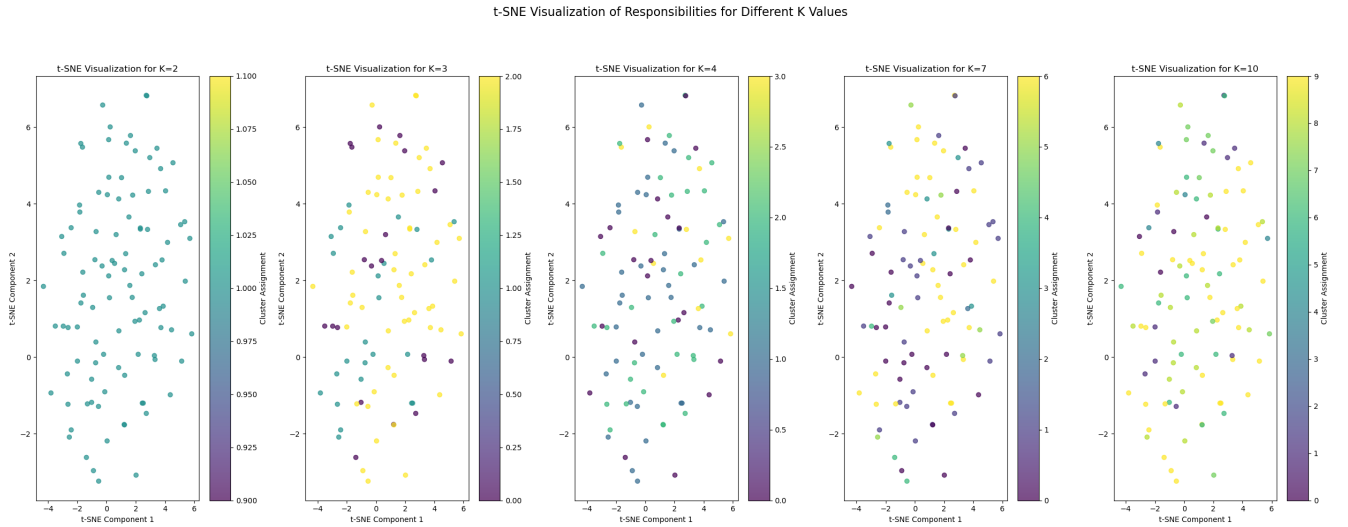


Figure 10: t-SNE Visualization of Representations for Different K values

Moreover, the resulting cluster responsibilities were visualized using t-SNE. Each t-SNE plot in Figure 10 provides a two-dimensional representation of the high-dimensional responsibility distributions, where each point corresponds to a data instance. The colors indicate the cluster assignments derived from the responsibilities, reflecting the degree of membership in each cluster.

The t-SNE visualizations for different values of  $K$  demonstrate the following:

The structure and separability of the clusters depend significantly on the choice of  $K$ . For smaller  $K$  (e.g.,  $K = 2$ ), the data points are grouped into broader clusters, resulting in less granular distinctions between responsibilities. As  $K$  increases (e.g.,  $K = 7, 10$ ), the clusters become more refined, with smaller, more localized groupings. However, overly large  $K$  may result in overfitting, where clusters start representing noise or minor variations in the data.

The t-SNE plots indicate that the clusters are reasonably well-formed for moderate values of  $K$  (e.g.,  $K = 3, 4, 7$ ). For smaller  $K$ , the clusters appear overly broad, potentially failing to capture the finer details of the data structure. For larger  $K$ , while the clusters become more distinct, some may overlap or fragment, suggesting that the increased complexity does not necessarily correspond to better interpretability or meaningful groupings.

The color gradients in the t-SNE plots reflect the responsibilities assigned to each cluster. These responsibilities indicate the likelihood of each data point belonging to a specific cluster. For  $K = 2$ , most points are strongly associated with one cluster or the other. As  $K$  increases, the responsibilities are distributed more finely, reflecting a higher degree of uncertainty or overlapping memberships for some points.

## Python code for (e) (tsne visualization for responsibilities):

```

1 from sklearn.manifold import TSNE
2
3 def visualize_responsibilities_tsne(X, pi_list, P_list, k_values):
4     """
5     Visualize the responsibilities for different K values using t-SNE.
6
7     Parameters:
8     X : numpy.ndarray
9         The input data, shape (N, D) where N is the number of data points and D is the number of
10        features.
11     pi_list : list of numpy.ndarray
12        List of mixing proportions for each K, each shape (K,).
13     P_list : list of numpy.ndarray
14        List of cluster parameters for each K, each shape (K, D).
15     k_values : list of int
16        List of K values corresponding to pi_list and P_list.
17     """
18     num_k = len(k_values)
19     plt.figure(figsize=(25, 10))
20
21     for idx, K in enumerate(k_values):
22         # Compute responsibilities for the current K
23         responsibilities = e_step(X, pi_list[idx], P_list[idx])
24
25         # Apply t-SNE
26         tsne = TSNE(n_components=2, random_state=42)
27         data_tsne = tsne.fit_transform(X)
28
29         # Get cluster assignments based on maximum responsibility
30         cluster_assignments = np.argmax(responsibilities, axis=1)
31
32         plt.subplot(1, num_k, idx + 1)
33         plt.scatter(data_tsne[:, 0], data_tsne[:, 1], c=cluster_assignments, cmap='viridis',
34                    alpha=0.7)
35         plt.title(f't-SNE Visualization for K={K}')
36         plt.xlabel('t-SNE Component 1')
37         plt.ylabel('t-SNE Component 2')
38         plt.colorbar(label='Cluster Assignment')
39
40     plt.suptitle('t-SNE Visualization of Responsibilities for Different K Values', fontsize=16)
41     plt.tight_layout()
42     plt.subplots_adjust(top=0.85)
43     plt.show()
44
45     k_values = [2, 3, 4, 7, 10]
46     pi_list = [np.random.dirichlet(np.ones(K), size=1)[0] for K in k_values]
47     P_list = [np.random.rand(K, D) for K in k_values]
48     visualize_responsibilities_tsne(X, pi_list, P_list, k_values)

```

Listing 5: src/solutions/q1.py

- (f) To express the log-likelihoods obtained in bits and relate it to the length of the naive encoding. First, consider the log-likelihood of the dataset  $\{x^{(n)}\}_{n=1}^N$  given model parameters  $\theta$ . This is defined as:

$$\text{log-likelihood} = \log(P(\{x^{(n)}\}_{n=1}^N|\theta)).$$

Since probabilities  $P(\{x^{(n)}\}_{n=1}^N|\theta)$  are typically values between 0 and 1, their logarithms are negative values (because  $\log(p) < 0$  when  $0 < p < 1$ ).

To represent this log-likelihood as an information cost in bits, we need to convert the log-likelihood from base  $e$  (natural log) to base 2. This conversion is achieved by dividing the log-likelihood by  $\ln(2)$ :

$$\log_2(P(\{x^{(n)}\}_{n=1}^N|\theta)) = \frac{\log(P(\{x^{(n)}\}_{n=1}^N|\theta))}{\ln(2)}.$$

To interpret this result as the total number of bits required to encode the dataset under the model, we use the negative form:

$$-\log_2(P(\{x^{(n)}\}_{n=1}^N|\theta)).$$

The negative sign is necessary because the log-likelihood is generally negative, and using the negative form converts it to a positive quantity. This positive value represents the information cost in bits or the compressed encoding length of the data according to the model.

For the naive encoding length of binary data, where each pixel or feature is represented by a single bit (0 or 1), we calculate the total bit length as:

$$\text{naive encoding length} = N \times D,$$

where  $N$  is the number of data points and  $D$  is the number of binary features per data point. This length represents the uncompressed encoding of the entire dataset.

The compression rate measures how much the model-based encoding reduces the bit length compared to the naive encoding. It is defined as:

$$\text{compression rate} = 1 - \frac{-\log_2(P(\{x^{(n)}\}_{n=1}^N|\theta))}{N \times D}.$$

In this formula:  $-\log_2(P(\{x^{(n)}\}_{n=1}^N|\theta))$  represents the information cost in bits using the model.  $N \times D$  is the naive encoding length.

The negative sign in the log-likelihood expression ensures that we obtain a positive bit length. This bit length can then be compared to the naive encoding length to compute the compression rate. A positive compression rate indicates that the model provides compression relative to naive encoding, and a rate close to 1 suggests a high degree of compression.

K	Final Log-Likelihood in Bits	Compression Rate
2	-4766	0.26
3	-4559	0.29
4	-4209	0.34
7	-3729	0.42
10	-3375	0.47

Table 2: Final Log-Likelihood in Bits and Compression Rate for Different Values of  $K$

From Table 2, as  $K$  increases, both the final log-likelihood in bits and the compression rate improve. The log-likelihood becomes less negative, indicating a better model fit as the number of clusters grows, while the compression rate approaches 0.5 for  $K = 10$ . This trend shows that the model captures more structure in the data with higher  $K$ , resulting in improved compression relative to naive encoding.

Compared to standard compression algorithms like gzip, the compression rates achieved by this model are lower. gzip leverages exact redundancy and repeated patterns in the data, allowing for highly efficient compression, especially on structured or repetitive data. In contrast, our model-based compression relies on capturing underlying probabilistic patterns rather than exact sequence redundancy. While the model captures the statistical structure in the data, it does not reach the compression efficiency of gzip on simple repetitive data, which explains the difference in compression performance.

Python code for (f) (log-likelihood in bits and compression rates):

```
1 import numpy as np
2
3
4 def em_algorithm_compression_rate(X, K, max_iter=20):
5     """
6     Parameters:
7     - X: Data matrix (N x D), where N is the number of data points and D is the dimensionality.
8     - K: Number of clusters/components.
9     - max_iter: Maximum number of iterations for the EM algorithm.
10
11     Returns:
12     - final_log_likelihood_bits: Final log-likelihood value in bits (base 2).
13     - compression_rate: Compression rate compared to the naive encoding.
14     """
15     # Initialize parameters
16     N, D = X.shape
17     pi, P = initialize_parameters(K, D)
18     log_likelihoods = []
19
20     for _ in range(max_iter):
21         # E-step
22         resp = e_step(X, pi, P)
23
24         # M-step
25         pi, P = m_step(X, resp)
26
27         # Compute and store log-likelihood in natural log
28         log_likelihood = compute_log_likelihood(X, pi, P)
29         log_likelihoods.append(log_likelihood)
30
31     # Convert the final log-likelihood to bits
32     final_log_likelihood_bits = log_likelihoods[-1] / np.log(2)
33
34     # Calculate naive encoding length in bits
35     naive_encoding_length = N * D
36
37     # Calculate compression rate
38     compression_rate = 1 - (-final_log_likelihood_bits / naive_encoding_length)
39
40     return final_log_likelihood_bits, compression_rate
41
42 def compute_compression_rates_for_different_K(X, K_values, max_iter=20):
43     """
44     Parameters:
45     - X: Data matrix (N x D).
46     - K_values: List of integers, where each integer is a value of K (number of clusters) to
47       evaluate.
48     - max_iter: Maximum number of iterations for the EM algorithm.
49
50     Returns:
51     - results: Dictionary where keys are K values and values are tuples
52       (final_log_likelihood_bits, compression_rate).
53     """
54     results = {}
55     for K in K_values:
56         final_log_likelihood_bits, compression_rate = em_algorithm_compression_rate(X, K,
57                                           max_iter)
58         results[K] = (final_log_likelihood_bits, compression_rate)
59         print(f"K={K}: Final log-likelihood in bits = {final_log_likelihood_bits},
60               Compression rate = {compression_rate}")
61     return results
62
63 K_values = [2, 3, 4, 7, 10]
64 results = compute_compression_rates_for_different_K(data, K_values, max_iter=20)
65 print(results)
```

Listing 6: src/solutions/q1.py

(g) The total cost of encoding both the data and model parameters can be computed as follows:

$$\text{Total Cost} = -\log_2 \left( P \left( \{\mathbf{x}^{(n)}\}_{n=1}^N \mid \theta \right) \right) + M \cdot (K - 1 + K \cdot D)$$

where:

- $-\log_2 \left( P \left( \{\mathbf{x}^{(n)}\}_{n=1}^N \mid \theta \right) \right)$  represents the final log-likelihood in bits, corresponding to the data encoding cost.
- $M$  is the bit cost of storing a single floating-point value (64 bits for a `float64`).
- $K$  is the number of clusters.
- $D$  is the dimensionality of each data point.
- $(K - 1 + K \cdot D)$  accounts for the number of model parameters, which includes: -
- $K - 1$  mixing proportions (since they sum to 1).
- $K \cdot D$  parameters for the Bernoulli distributions associated with each cluster.

The total compression rate is then defined as the ratio of the total cost to the naive encoding length:

$$\text{Total Compression Rate} = \frac{\text{Total Cost}}{N \cdot D}$$

where:  $N \cdot D$  represents the naive encoding length, assuming each data point can be represented by a single bit per pixel (for binary data).

Thus, the final outputs are the total cost and the total compression rate for each value of  $K$ .

<b>K</b>	<b>Total Cost</b>	<b>Total Compression Rate</b>
2	11523.68	1.8006
3	15700.66	2.4532
4	19808.04	3.0950
7	32117.22	5.0183
10	44521.13	6.9564

Table 3: Total Cost and Total Compression Rate for Different Values of  $K$

The Table 3 shows that the total cost and compression rate increase as  $K$  becomes larger. This indicates that the model complexity and storage requirements grow with the number of clusters, as each additional cluster adds more parameters that must be stored. Compared to a compression algorithm like gzip, which is optimized for minimizing redundancy without explicitly modeling clusters, the EM model with higher  $K$  values is less efficient because it encodes both data and parameters.

As  $K$  increases, the total compression rate exceeds 1, meaning this method requires more bits than naive encoding due to the overhead of storing model parameters. This suggests that while increasing  $K$  provides a finer-grained model, it incurs higher storage costs. This result highlights a trade-off: increasing  $K$  might improve the model's fit but reduces compression efficiency, showing the balance between model complexity and storage efficiency.

## Python code for (g) (Total Cost and Total Compression Rate):

```
1
2 import numpy as np
3 import pandas as pd
4
5 def em_algorithm_total_cost_and_compression_rate_single_run(X, K_values, max_iter=20, M=64):
6     """
7
8     Parameters:
9     - X: Data matrix (N x D), where N is the number of data points and D is the dimensionality.
10    - K_values: List of integers, the different values of K (number of clusters) to evaluate.
11    - max_iter: Maximum number of iterations for the EM algorithm.
12    - M: Cost in bits to store each parameter (default is 64 bits for a float64).
13
14    Returns:
15    - results_df: DataFrame with total compression rates for each K value.
16    """
17    N, D = X.shape
18    naive_encoding_length = N * D
19    results = {"K": [], "Total Cost": [], "Total Compression Rate": []}
20
21    for K in K_values:
22        pi, P = initialize_parameters(K, D)
23        log_likelihoods = []
24
25        # EM Algorithm
26        for _ in range(max_iter):
27            resp = e_step(X, pi, P)
28
29            pi, P = m_step(X, resp)
30
31            log_likelihood = compute_log_likelihood(X, pi, P)
32            log_likelihoods.append(log_likelihood)
33
34        # Final log-likelihood in bits
35        final_log_likelihood_nats = log_likelihoods[-1]
36        final_log_likelihood_bits = final_log_likelihood_nats / np.log(2)
37
38        # Cost of encoding the model parameters
39        num_params = (K - 1) + K * D
40        model_cost_bits = num_params * M
41
42        # Total cost in bits
43        total_cost = -final_log_likelihood_bits + model_cost_bits
44
45        # Total compression rate
46        total_compression_rate = total_cost / naive_encoding_length
47
48        results["K"].append(K)
49        results["Total Cost"].append(total_cost)
50        results["Total Compression Rate"].append(total_compression_rate)
51
52    results_df = pd.DataFrame(results)
53
54    return results_df
55
56 K_values = [2, 3, 4, 7, 10]
57
58 # Run the computation
59 results_df = em_algorithm_total_cost_and_compression_rate_single_run(X, K_values, max_iter=20,
60                                                                    M=64)
61 print(results_df)
```

Listing 7: src/solutions/q1.py

## Question 4

(a) The filtered states exhibit more variability and respond quickly to changes in observations as shown at Figure 11. This is because each filtered estimate is based only on data up to the current time step, making it more sensitive to local fluctuations. The filtered estimates appear noisier, with more frequent and sharper changes in the state values. This reflects the uncertainty in the estimate, as it does not take future information into account. The filtered states are suitable for real-time tracking, as they only require past and current observations, but they may not provide the most accurate state estimate.

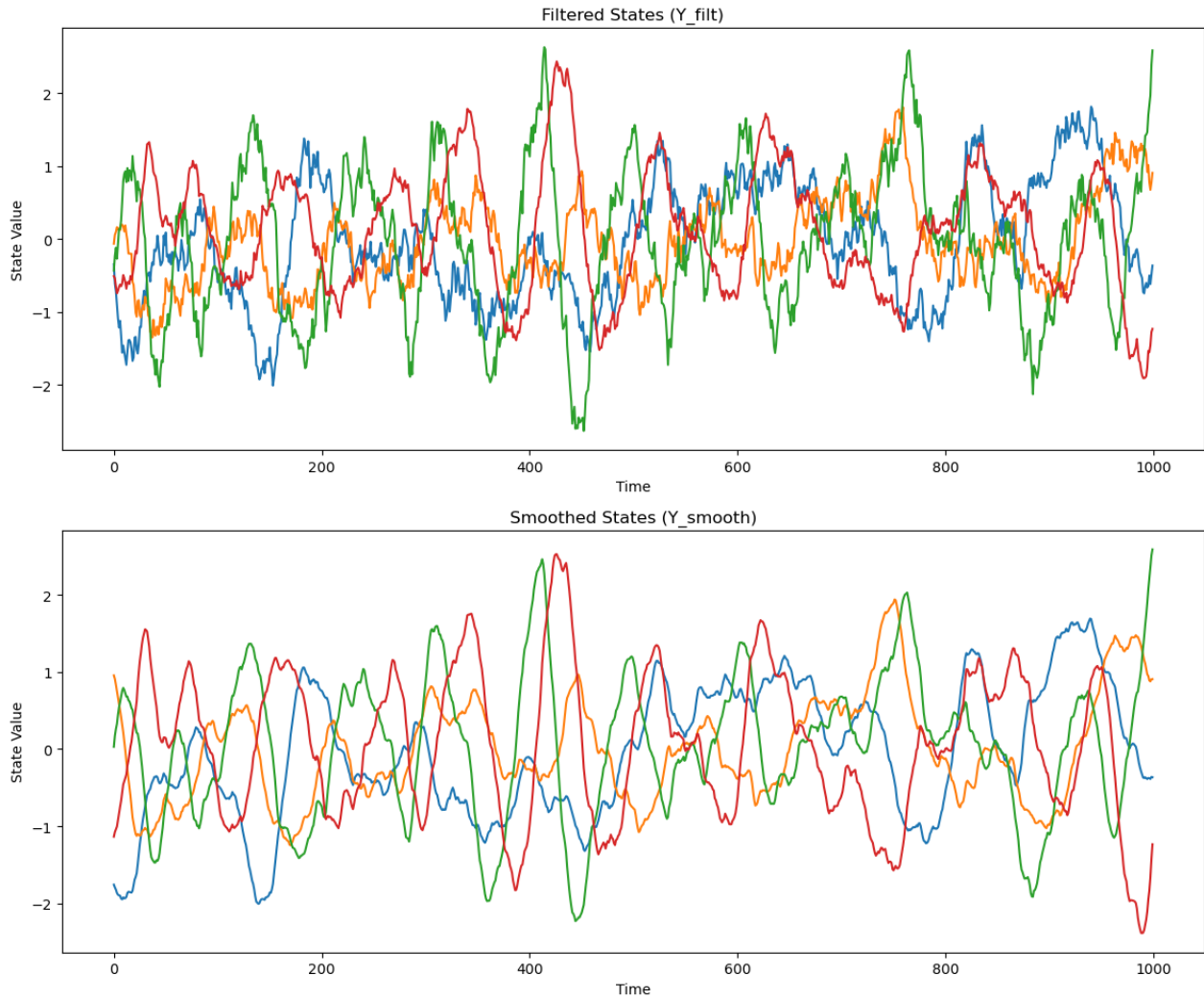


Figure 11: Kalman filtering vs Kalman Smoothing Y

The smoothed states are notably smoother than the filtered states, showing fewer sharp fluctuations as shown at Figure 12. This is because the smoother considers all observations (both past and future) to estimate each state, leading to a more refined trajectory. With the benefit of future data, the smoothed estimates can correct past

fluctuations that might have been due to noise, resulting in a cleaner and more stable pattern over time.

The smoothed states are generally closer to the true underlying states, as they incorporate more information. This makes them more reliable for post-hoc analysis.

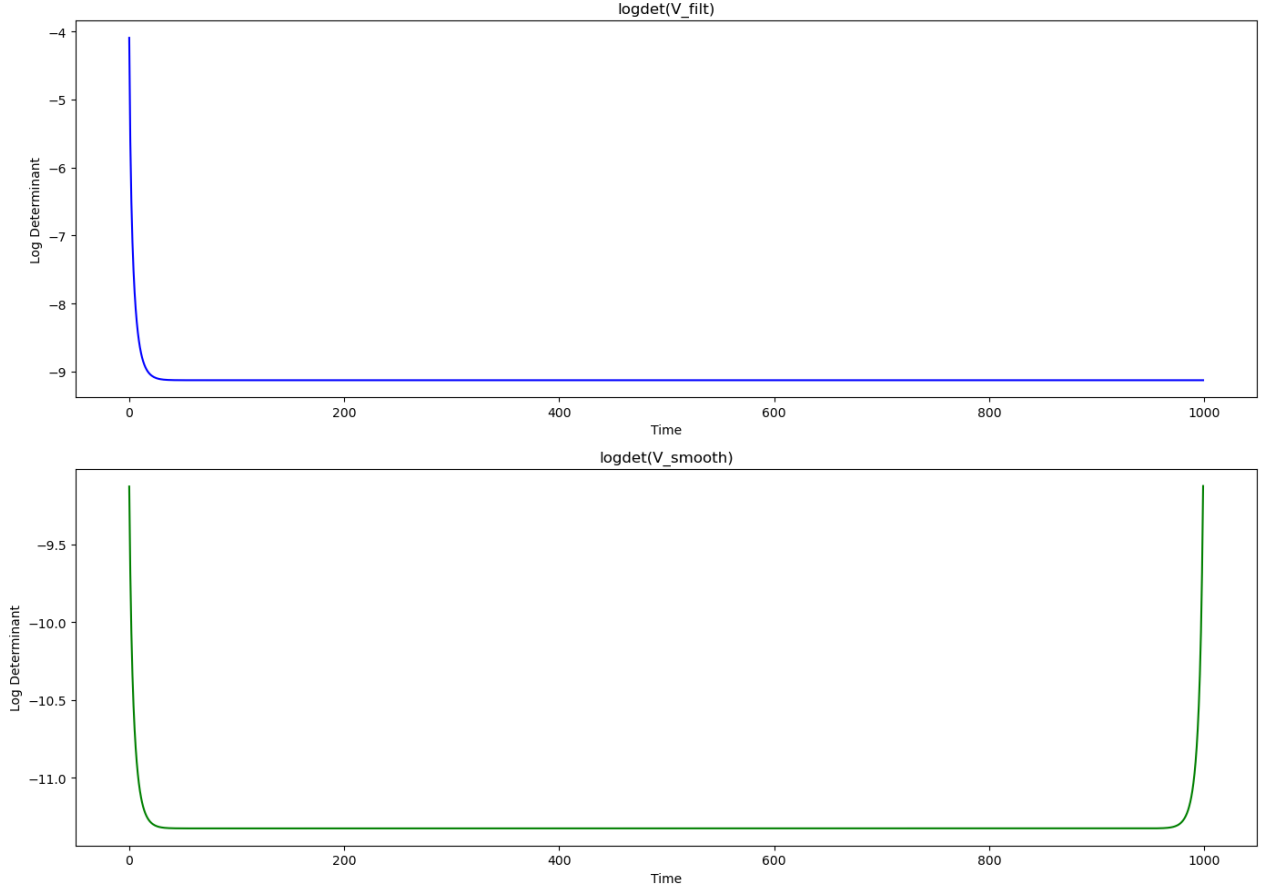


Figure 12: Kalman filtering vs Kalman Smoothing logdet V

The smoothed estimates ( $Y_{\text{smooth}}$ ) are generally more accurate and certain than the filtered estimates ( $Y_{\text{filt}}$ ), as they incorporate information from the entire sequence of observations.

The filtered states ( $Y_{\text{filt}}$ ) are more sensitive to observation noise and local fluctuations, while the smoothed states ( $Y_{\text{smooth}}$ ) average out these fluctuations, providing a better estimate of the underlying state.

The filtered states are suitable for real-time applications where only past and current data is available. In contrast, the smoothed states are more useful for retrospective analysis, where the entire dataset can be used to improve the accuracy of state estimates.



## Python code for (a):

```
1
2 # Define parameters
3 X = np.loadtxt('/Users/richardhuang/Documents/GitHub/Projects/Probabilistic
   Summative/ssm_spins.txt').T
4
5 # Define given parameters
6 A = 0.99 * np.array([
7     [np.cos(2 * np.pi / 180), -np.sin(2 * np.pi / 180), 0, 0],
8     [np.sin(2 * np.pi / 180), np.cos(2 * np.pi / 180), 0, 0],
9     [0, 0, np.cos(2 * np.pi / 90), -np.sin(2 * np.pi / 90)],
10    [0, 0, np.sin(2 * np.pi / 90), np.cos(2 * np.pi / 90)]
11])
12
13 C = np.array([
14     [1, 1, 0, 1],
15     [0, 1, 1, 1],
16     [1, 0, 1, 1],
17     [0, 0, 1, 1],
18     [0.5, 0.5, 0.5, 0.5]
19])
20
21 Q = np.eye(4) - A @ A.T
22 R = np.eye(5)
23 y_init = np.zeros(4)
24 Q_init = np.eye(4)
25
26 def logdet(A):
27     return 2 * np.sum(np.log(np.diag(np.linalg.cholesky(A))))
28
29
30 Y_filt, V_filt, _, L_filt = run_ssm_kalman(X, y_init, Q_init, A, Q, C, R, 'filt')
31 y_smooth, V_smooth, V_joint, L_smooth = run_ssm_kalman(X, y_init, Q_init, A, Q, C, R, 'smooth')
32
33 # Generate the plots
34 plt.figure(figsize=(12, 10))
35
36 plt.subplot(2, 1, 1)
37 plt.plot(y_filt.T)
38 plt.title("Filtered States (Y_filt)")
39 plt.xlabel("Time")
40 plt.ylabel("State Value")
41
42 plt.subplot(2, 1, 2)
43 plt.plot(y_smooth.T)
44 plt.title("Smoothed States (Y_smooth)")
45 plt.xlabel("Time")
46 plt.ylabel("State Value")
47
48 # Calculate log determinants for V_filt and V_smooth
49 logdet_V_filt = [logdet(V_filt[t]) for t in range(len(V_filt))]
50 logdet_V_smooth = [logdet(V_smooth[t]) for t in range(len(V_smooth))]
51
52 plt.figure(figsize=(14, 10))
53 plt.subplot(2, 1, 1)
54 plt.plot(logdet_V_filt, color='blue')
55 plt.title("logdet(V_filt)")
56 plt.xlabel("Time")
57 plt.ylabel("Log Determinant")
58
59 plt.subplot(2, 1, 2)
60 plt.plot(logdet_V_smooth, color='green')
61 plt.title("logdet(V_smooth)")
62 plt.xlabel("Time")
63 plt.ylabel("Log Determinant")
64
65 plt.tight_layout()
66 plt.show()
```

Listing 8: src/solutions/q1.py

(b) To derive the M-step update equations for  $R_{\text{new}}$  and  $Q_{\text{new}}$ .

First, consider  $R_{\text{new}}$ , The observation model is given by:

$$\mathbf{x}_t = C\mathbf{y}_t + \mathbf{v}_t, \quad \mathbf{v}_t \sim \mathcal{N}(0, R)$$

where  $R$  is the observation noise covariance matrix. The goal is to update  $R$  using the current estimates of the latent states  $\mathbf{y}_t$  and the observations  $\mathbf{x}_t$ .

The optimal  $R_{\text{new}}$  can be derived by computing the expected squared difference between  $\mathbf{x}_t$  and  $C\mathbf{y}_t$ , averaged over time. This expectation is:

$$R_{\text{new}} = \frac{1}{T} \sum_{t=1}^T \mathbb{E} [(\mathbf{x}_t - C\mathbf{y}_t)(\mathbf{x}_t - C\mathbf{y}_t)^\top].$$

Expanding this expectation:

$$R_{\text{new}} = \frac{1}{T} \left( \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top - \sum_{t=1}^T \mathbf{x}_t \langle \mathbf{y}_t \rangle^\top C^\top - C \sum_{t=1}^T \langle \mathbf{y}_t \rangle \mathbf{x}_t^\top + C \sum_{t=1}^T \langle \mathbf{y}_t \mathbf{y}_t^\top \rangle C^\top \right).$$

However, by choosing  $C_{\text{new}}$  such that  $C_{\text{new}} = \left( \sum_{t=1}^T \mathbf{x}_t \langle \mathbf{y}_t \rangle^\top \right) \left( \sum_{t=1}^T \langle \mathbf{y}_t \mathbf{y}_t^\top \rangle \right)^{-1}$ , we simplify  $R_{\text{new}}$  to:

$$R_{\text{new}} = \frac{1}{T} \left( \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top - \left( \sum_{t=1}^T \mathbf{x}_t \langle \mathbf{y}_t \rangle^\top \right) C_{\text{new}}^\top \right).$$

For  $Q_{\text{new}}$ , The state transition model is given by:

$$\mathbf{y}_t = A\mathbf{y}_{t-1} + \mathbf{w}_t, \quad \mathbf{w}_t \sim \mathcal{N}(0, Q)$$

where  $Q$  is the process noise covariance matrix. To update  $Q$ , we compute the expected squared difference between  $\mathbf{y}_t$  and  $A\mathbf{y}_{t-1}$ , averaged over time. This expectation is:

$$Q_{\text{new}} = \frac{1}{T-1} \sum_{t=2}^T \mathbb{E} [(\mathbf{y}_t - A\mathbf{y}_{t-1})(\mathbf{y}_t - A\mathbf{y}_{t-1})^\top].$$

Expanding this expectation:

$$Q_{\text{new}} = \frac{1}{T-1} \left( \sum_{t=2}^T \langle \mathbf{y}_t \mathbf{y}_t^\top \rangle - \sum_{t=2}^T \langle \mathbf{y}_t \rangle \langle \mathbf{y}_{t-1}^\top \rangle A^\top - A \sum_{t=2}^T \langle \mathbf{y}_{t-1} \rangle \langle \mathbf{y}_t \rangle^\top + A \sum_{t=2}^T \langle \mathbf{y}_{t-1} \mathbf{y}_{t-1}^\top \rangle A^\top \right).$$

By choosing  $A_{\text{new}}$  such that  $A_{\text{new}} = \left( \sum_{t=2}^T \langle \mathbf{y}_t \rangle \langle \mathbf{y}_{t-1}^\top \rangle \right) \left( \sum_{t=2}^T \langle \mathbf{y}_{t-1} \mathbf{y}_{t-1}^\top \rangle \right)^{-1}$ , we simplify  $Q_{\text{new}}$  to:

$$Q_{\text{new}} = \frac{1}{T-1} \left( \sum_{t=2}^T \langle \mathbf{y}_t \mathbf{y}_t^\top \rangle - \left( \sum_{t=2}^T \langle \mathbf{y}_t \rangle \langle \mathbf{y}_{t-1}^\top \rangle \right) A_{\text{new}}^\top \right).$$

Thus, we obtain the simplified M-step update equations for  $R$  and  $Q$ :

$$R_{\text{new}} = \frac{1}{T} \left( \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top - \left( \sum_{t=1}^T \mathbf{x}_t \langle \mathbf{y}_t \rangle^\top \right) C_{\text{new}}^\top \right)$$

$$Q_{\text{new}} = \frac{1}{T-1} \left( \sum_{t=2}^T \langle \mathbf{y}_t \mathbf{y}_t^\top \rangle - \left( \sum_{t=2}^T \langle \mathbf{y}_t \rangle \langle \mathbf{y}_{t-1}^\top \rangle \right) A_{\text{new}}^\top \right).$$

The two plots below display the progression of the log-likelihood over 100 iterations of the Expectation-Maximization (EM) algorithm under two different initialization conditions. In the first plot, the EM algorithm is initialized with the generating (true) parameters, while in the second plot, it is initialized with 10 different sets of random parameters.

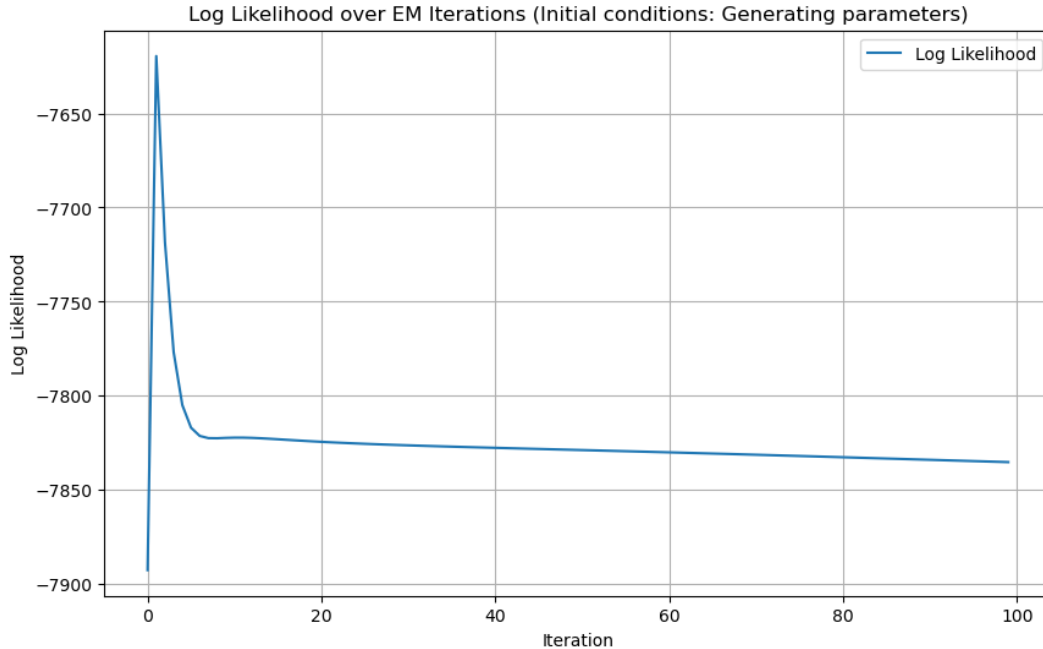


Figure 13: EM Iterations (Initial Conditions: Generating parameters)

In the first Figure 13, we observe a rapid increase in log-likelihood within the initial few iterations, followed by a gradual convergence to a stable value. This behavior suggests that starting with parameters close to the true values allows the EM algorithm to quickly reach a high log-likelihood, as it starts near an optimal point in the parameter space. After the initial rapid increase, the log-likelihood stabilizes, indicating that the algorithm has likely converged and further iterations provide minimal improvement.

The second Figure 14, where EM is initialized with random parameters for 10 different runs, shows a different pattern. Each line represents a separate run with a unique random initialization. Here, the log-likelihoods begin at much lower values, reflecting the random starting points far from the optimum. Over the iterations, each run's log-likelihood gradually increases, but at a slower rate compared to the first plot. Despite the varied starting points, all runs eventually converge to similar log-likelihood values, demonstrating the robustness of the EM algorithm in finding a high-likelihood solution regardless of initialization. However, the convergence process is slower and less direct compared to starting from the true parameters.

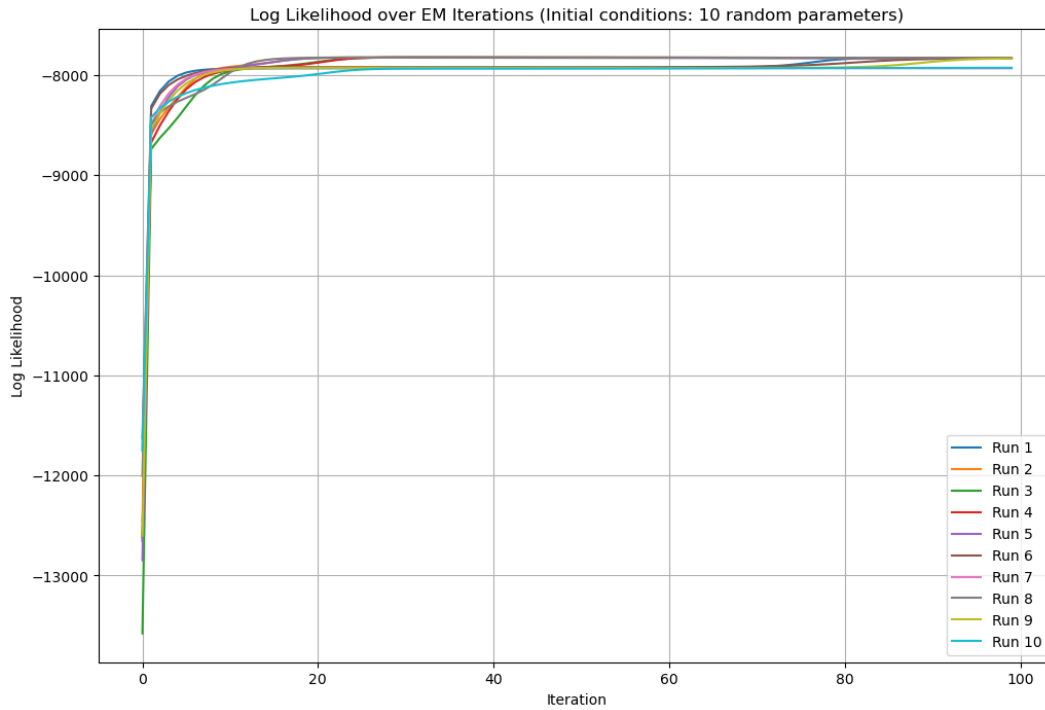


Figure 14: EM Iterations (Initial Conditions: 10 random parameters)

## Python code for (b)

```

1  def cellsum(C):
2      """
3      Sums a list of matrices element-wise.
4      :param C: List of matrices (numpy arrays) to sum
5      :return: Element-wise sum of all matrices in C
6      """
7      return np.sum(np.stack(C, axis=2), axis=2)
8
9
10 def em_algorithm(X, A_init, Q_init, C_init, R_init, y_init, Q_0, num_iterations=100):
11     """
12     EM algorithm for learning the parameters A, Q, C, and R.
13     :param X: Observation data, shape [d, T]
14     :param A_init: Initial A matrix
15     :param Q_init: Initial Q matrix
16     :param C_init: Initial C matrix
17     :param R_init: Initial R matrix
18     :param y_init: Initial state vector
19     :param Q_0: Initial state covariance
20     :param num_iterations: Number of EM iterations
21     :return: Updated A, Q, C, R matrices and log-likelihoods
22     """
23     A, Q, C, R = A_init, Q_init, C_init, R_init
24     T = X.shape[1]
25     log_likelihoods = []
26
27     for i in range(num_iterations):
28         y_hat, V_hat, V_joint, log_likelihood = run_ssm_kalman(X, y_init, Q_0, A, Q, C, R,
29             mode='smooth')
30         log_likelihoods.append(np.sum(log_likelihood))
31
32         # M-step: Update parameters A, Q, C, and R
33
34         # Update C
35         C_new = (cellsum([np.outer(X[:, t], y_hat[:, t]) for t in range(T)]) @
36             np.linalg.inv(cellsum([V_hat[t] + np.outer(y_hat[:, t], y_hat[:, t])

```

```

36         for t in range(T)))))
37
38     # Update R
39     R_new = (1 / T) * (cellsum([np.outer(X[:, t], X[:, t]) for t in range(T)]) -
40                       cellsum([np.outer(X[:, t], y_hat[:, t]) for t in range(T)]) @ C_new.T)
41
42     # Update A
43     A_new = (cellsum([np.outer(y_hat[:, t + 1], y_hat[:, t]) for t in range(T - 1)]) @
44              np.linalg.inv(cellsum([V_hat[t] + np.outer(y_hat[:, t], y_hat[:, t])
45              for t in range(T - 1)])))
46
47     # Update Q
48     Q_new = (1 / (T - 1)) * (cellsum([V_hat[t + 1] + np.outer(y_hat[:, t + 1], y_hat[:, t +
49     1])
50     for t in range(T - 1)]) -
51     cellsum([np.outer(y_hat[:, t + 1], y_hat[:, t]) for t in range(T - 1)]) @ A_new.T)
52     A, Q, C, R = A_new, Q_new, C_new, R_new
53
54     return A, Q, C, R, log_likelihoods
55
56 A_em, Q_em, C_em, R_em, log_likelihoods_para =
57 em_algorithm(X, A, Q, C, R, y_init, Q_0)
58
59 plt.figure(figsize=(10, 6))
60 plt.plot(log_likelihoods_para, label='Log Likelihood')
61 plt.xlabel('Iteration')
62 plt.ylabel('Log Likelihood')
63 plt.title('Log Likelihood over EM Iterations (Initial conditions: Generating parameters)')
64 plt.legend()
65 plt.grid()
66 plt.show()
67
68 # Run EM with 10 random initializations
69 all_log_likelihoods = []
70
71 for i in range(10):
72     A_rand = np.random.randn(*A.shape)
73     Q_rand = np.eye(A.shape[0])
74     C_rand = np.random.randn(*C.shape)
75     R_rand = np.eye(X.shape[0])
76
77     A_ra, Q_ra, C_ra, R_ra, log_likelihoods_rand
78     = em_algorithm(X, A_rand, Q_rand, C_rand, R_rand, y_init, Q_0)
79
80     all_log_likelihoods.append(log_likelihoods_rand)
81
82 plt.figure(figsize=(12, 8))
83
84 for i, log_likelihoods in enumerate(all_log_likelihoods):
85     plt.plot(log_likelihoods, label=f'Run {i+1}')
86
87 plt.xlabel('Iteration')
88 plt.ylabel('Log Likelihood')
89 plt.title('Log Likelihood over EM Iterations (Initial conditions: 10 random parameters)')
90 plt.legend()
91 plt.grid()
92 plt.show()

```

Listing 9: src/solutions/q1.py

(c)

# Question 5

(a) To derive the Maximum Likelihood (ML) estimates for the transition probabilities  $\psi(\alpha, \beta)$  and the stationary distribution  $\phi(\gamma)$ , we proceed as follows:

The transition probability  $\psi(\alpha, \beta)$  represents the probability of observing symbol  $\alpha$  immediately after symbol  $\beta$ .

Estimate  $\psi(\alpha, \beta) = p(s_i = \alpha | s_{i-1} = \beta)$  based on observed data by maximizing the likelihood of observed transitions.

Let:

- $N(\alpha, \beta)$ : the count of observed pairs where  $\beta$  is followed by  $\alpha$ ,
- $N(\beta)$ : the total number of times  $\beta$  appears as a preceding symbol in any pair.

To maximize the likelihood of the observed data, we define the likelihood function for all observed transitions as:

$$L(\psi) = \prod_{\alpha, \beta} \psi(\alpha, \beta)^{N(\alpha, \beta)}$$

where  $\psi(\alpha, \beta)$  is the probability of observing symbol  $\alpha$  given that the previous symbol was  $\beta$ .

Taking the logarithm of the likelihood function, we get the log-likelihood:

$$\log L(\psi) = \sum_{\alpha, \beta} N(\alpha, \beta) \log \psi(\alpha, \beta)$$

Since the probabilities must sum to 1 for each preceding symbol  $\beta$ , we have the constraint:

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

To maximize  $\log L(\psi)$  subject to the above constraint, we introduce a Lagrange multiplier  $\lambda_{\beta}$  for each  $\beta$ :

$$\mathcal{L} = \sum_{\alpha, \beta} N(\alpha, \beta) \log \psi(\alpha, \beta) + \sum_{\beta} \lambda_{\beta} \left( 1 - \sum_{\alpha} \psi(\alpha, \beta) \right)$$

Taking the derivative of  $\mathcal{L}$  with respect to  $\psi(\alpha, \beta)$  and setting it to zero:

$$\frac{\partial \mathcal{L}}{\partial \psi(\alpha, \beta)} = \frac{N(\alpha, \beta)}{\psi(\alpha, \beta)} - \lambda_{\beta} = 0$$

Rearranging, we find:

$$\psi(\alpha, \beta) = \frac{N(\alpha, \beta)}{\lambda_{\beta}}$$

Using the constraint  $\sum_{\alpha} \psi(\alpha, \beta) = 1$ , we find:

$$\sum_{\alpha} \frac{N(\alpha, \beta)}{\lambda_{\beta}} = 1$$

$$\lambda_{\beta} = N(\beta)$$

Thus, the ML estimate for  $\psi(\alpha, \beta)$  is:

$$\psi(\alpha, \beta) = \frac{N(\alpha, \beta)}{N(\beta)}$$

The stationary distribution  $\phi(\gamma)$  represents the long-term probability of observing symbol  $\gamma$  in the text. Let:

- $N(\gamma)$ : the total count of occurrences of symbol  $\gamma$ ,
- $N$ : the total number of symbols in the text.

The probability  $\phi(\gamma)$  can be estimated by the relative frequency of  $\gamma$  in the text:

$$\phi(\gamma) = \frac{N(\gamma)}{N}$$

where:

- $N(\gamma)$  is the count of symbol  $\gamma$  in the entire text.
- $N = \sum_{\gamma} N(\gamma)$ , the total number of symbols in the text.

Table 4: Transition Matrix  $\Psi$  (Sorry.. the numbers are too many to display here)

(b) The latent variables  $\sigma(s)$  represent the mappings of each symbol in the encrypted text to symbols in the decrypted text. Since we assume a uniform prior distribution over all possible permutations, each permutation is equally likely, which means there is no inherent dependency structure enforced by the prior.

However, **the latent variables  $\sigma(s)$  are not independent**. This is because they represent a permutation. In a permutation, choosing one mapping (e.g.,  $\sigma(a) = s$ ) affects the choices for the remaining mappings (e.g.,  $\sigma(b)$  cannot also be  $s$  if it's already mapped to  $a$ ). Therefore, while the prior is uniform, the constraint that each symbol maps to a unique counterpart means that the latent variables  $\sigma(s)$  are not independent.

The joint probability of the encrypted sequence  $e_1 e_2 \cdots e_n$  given the permutation  $\sigma$  can be expressed in terms of the transition probabilities  $\psi$  and the stationary distribution  $\phi$ , which were estimated from English text in part (a).

Define  $d_i = \sigma^{-1}(e_i)$  to be the decrypted symbol corresponding to  $e_i$  under the permutation  $\sigma$ .

The joint probability of the sequence  $e_1 e_2 \cdots e_n$  given  $\sigma$  is:

$$p(e_1 e_2 \cdots e_n | \sigma) = p(d_1 d_2 \cdots d_n) = \phi(d_1) \prod_{i=2}^n \psi(d_i | d_{i-1})$$

where:

- $\phi(d_1)$  is the stationary probability of the first decrypted symbol  $d_1$ ,
- $\psi(d_i | d_{i-1})$  is the transition probability from  $d_{i-1}$  to  $d_i$ .

This joint probability captures the likelihood of observing the decrypted sequence  $d_1 d_2 \cdots d_n$ , which corresponds to the encrypted sequence  $e_1 e_2 \cdots e_n$  under the permutation  $\sigma$ .

(c) Here, the proposal  $\sigma \rightarrow \sigma'$  is defined by randomly choosing two symbols  $s$  and  $s'$  and swapping their mappings in  $\sigma$ . Since both  $s$  and  $s'$  are chosen randomly, each pair of symbols has an equal probability of being selected.

Let  $N$  be the total number of symbols. The probability of selecting any specific pair  $(s, s')$  is:

$$S(\sigma \rightarrow \sigma') = \frac{1}{\binom{N}{2}} = \frac{2}{N(N-1)}$$

because there are  $\binom{N}{2}$  ways to choose two distinct symbols from  $N$  symbols.

This proposal probability  $S(\sigma \rightarrow \sigma')$  does not depend on the specific configurations of  $\sigma$  or  $\sigma'$ ; it is only dependent on the number of symbols  $N$ .

The Metropolis-Hastings acceptance probability  $A(\sigma \rightarrow \sigma')$  is given by:

$$A(\sigma \rightarrow \sigma') = \min \left( 1, \frac{\pi(\sigma') \cdot S(\sigma' \rightarrow \sigma)}{\pi(\sigma) \cdot S(\sigma \rightarrow \sigma')} \right)$$

where:

- $\pi(\sigma)$  is the target probability (posterior) of the current permutation  $\sigma$ ,
- $S(\sigma \rightarrow \sigma')$  is the proposal probability for moving from  $\sigma$  to  $\sigma'$ ,
- $S(\sigma' \rightarrow \sigma)$  is the reverse proposal probability from  $\sigma'$  back to  $\sigma$ .

Since  $S(\sigma \rightarrow \sigma') = S(\sigma' \rightarrow \sigma) = \frac{2}{N(N-1)}$  (symmetric proposal), the acceptance probability simplifies to:

$$A(\sigma \rightarrow \sigma') = \min \left( 1, \frac{\pi(\sigma')}{\pi(\sigma)} \right)$$

Thus, the acceptance probability depends only on the ratio of the posterior probabilities  $\pi(\sigma')$  and  $\pi(\sigma)$ .

- If  $\pi(\sigma') > \pi(\sigma)$ , then  $\sigma'$  has a higher probability than  $\sigma$ , and the move is always accepted.
- If  $\pi(\sigma') < \pi(\sigma)$ , the move is accepted with probability  $\frac{\pi(\sigma')}{\pi(\sigma)}$ .

Therefore,

$S(\sigma \rightarrow \sigma') = \frac{2}{N(N-1)}$ , independent of  $\sigma$  and  $\sigma'$ .

$A(\sigma \rightarrow \sigma') = \min \left( 1, \frac{\pi(\sigma')}{\pi(\sigma)} \right)$ , dependent only on the ratio  $\frac{\pi(\sigma')}{\pi(\sigma)}$ .



(d) To implement the Metropolis-Hastings (MH) sampler for decrypting the provided encrypted text. Below is a description of the process:

### **Initialization**

Defined a `Decrypter` class that takes a decryption dictionary to map encrypted symbols to their decrypted counterparts. The `initialize_parameters_by_frequency` function was used to create an initial mapping based on the frequency of symbols in the encrypted text. This mapping was derived from the frequency of letters in the English language, which helps in starting the decryption process with a reasonable guess.

### **Transition Probability Calculation**

The `calculate_transition_probs` function was implemented to compute the transition probabilities between symbols in the training text. This involved counting the occurrences of each symbol and the pairs of consecutive symbols, applying Laplace smoothing to avoid zero probabilities.

### **Log-Likelihood Calculation**

The `compute_log_likelihood` function was created to evaluate the likelihood of the decrypted text based on the transition probabilities. This function sums the log probabilities of the transitions in the decrypted text.

### **Metropolis-Hastings Algorithm**

The `metropolis_hastings_decrypt` function was implemented to run the MH algorithm. It starts with an initial mapping and iteratively proposes new mappings by swapping two random symbols. For each proposed mapping, the new decrypted text is generated, and its log-likelihood is computed. The new mapping is accepted if it improves the score or based on a probability that considers the temperature of the system (simulated annealing). The temperature is gradually decreased to allow the algorithm to converge over time.

### **Logging Progress**

The algorithm was run for 5000 times, after every 100 iterations, the current decryption of the first 60 symbols is logged.

MH Iteration	Current Decryption
0	fwxm3x3hdwt1 x2wrxxmh 1xydqwl 2eq1x312 gxm3x42ln1 xt2ylxmxgh
100	o2w4243joa102nor2w3012mjdol0ned1241n0g2w42qnli102anm12w12g3
200	e3 s4 4gj3r10 n32 sg01 mjd310nid1 41n0o s4 qnla10 rnm1 s1 og
300	e3 s4 4gh3r0y n32 sgy0 1hd30ynid0 40nyo s4 fnla0y rn10 s0 og
400	e0 s4 4ch0o3y n0d scy3 1h203yni23 43nyt s4 knla3y on13 s3 tc
500	e4 s1 1ch43of n4d scfo rhg4ofnigo 1onft s1 knlaof 3nro so tc
600	o4 s1 1ai4vef n4d safe rig4efnpge 1enft s1 kncuef vnre se ta
700	o3 ml lai3jes n3d mase vig3esnbge lenst ml kncues jnve me ta
800	o3 mg gai3res n3d mase vil3esnbge lenst ml kncues jnve me ta
900	on mg gainres und mase vilnesuble geust mg kuches ruve me ta
1000	on py yaincer und pare vilnerufle yeurt py kusher cuve pe ta
1100	on my yaincer und mare vilneruble yeurt my kusher cuve me ta
1200	an my yoincer und more vilneruble yeurt my kusher cuve me to
1300	an my yoinler und more vilneruble yeurt my cusher luve me to
1400	an cy yoinler und core vilneruble yeurt cy musher luve ce to
1500	an cy yoinker und core vilneruble yeurt cy musher kuve ce to
1600	an cy yoinker und core vilneruble yeurt cy musher kuve ce to
1700	an by yoinker und bore vilneruple yeurt by musher kuve be to
1800	an by yuinker ond bure vilnerople yeort by mosher kove be tu
1900	an by yuinger ond bure vilneromle yeort by posher gove be tu
2000	an by yuinger ond bure vilneromle yeort by fosher gove be tu
2100	an by yuinger ond bure vilneromle yeort by fosher gove be tu
2200	on my yuinger and mure vilnerable years my father gave me su
2300	on my yuinger and mure vilnerable years my father gave me su
2400	on my yuinger and mure vilnerable years my father gave me su
2500	on my yuinger and mure vilnerable years my father gave me su
2600	on my yuinger and mure vilnerable years my father gave me su
2700	on my yuinger and mure vilnerable years my father gave me su
2800	on my yuinger and mire vilnerable years my father gave me si
2900	in my younger and more vulnerable years my father gave me so
3000	in my younger and more vulnerable years my father gave me so
3100	in my younger and more vulnerable years my father gave me so
3200	in my younger and more vulnerable years my father gave me so
3300	in my younger and more vulnerable years my father gave me so
3400	in my younger and more vulnerable years my father gave me so
3500	in my younger and more vulnerable years my father gave me so
3600	in my younger and more vulnerable years my father gave me so
3700	in my younger and more vulnerable years my father gave me so
3800	in my younger and more vulnerable years my father gave me so
3900	in my younger and more vulnerable years my father gave me so
4000	in my younger and more vulnerable years my father gave me so
4100	in my younger and more vulnerable years my father gave me so
4200	in my younger and more vulnerable years my father gave me so
4300	in my younger and more vulnerable years my father gave me so
4400	in my younger and more vulnerable years my father gave me so
4500	in my younger and more vulnerable years my father gave me so
4600	in my younger and more vulnerable years my father gave me so
4700	in my younger and more vulnerable years my father gave me so
4800	in my younger and more vulnerable years my father gave me so
4900	in my younger and more vulnerable years my father gave me so
5000	in my younger and more vulnerable years my father gave me so

Table 5: Decryption Results

## Python code for (d): Report the current decryption

```
1
2 from collections import defaultdict, Counter
3 import random
4 from collections import Counter, defaultdict
5
6 url = 'https://www.gutenberg.org/files/2600/2600-0.txt'
7 response = requests.get(url)
8 text = response.text
9
10 random.seed(42)
11
12 class Decrypter:
13     def __init__(self, decryption_dict):
14         self.decryption_dict = decryption_dict
15
16     def decrypt(self, encrypted_message):
17         return ''.join(self.decryption_dict.get(char, char) for char in encrypted_message)
18
19 def initialize_parameters_by_frequency(encrypted_text):
20     """Initialize mapping based on English letter frequencies."""
21     # English letter frequencies (including space, punctuation, and numbers)
22     english_freq_order = " etaoinshrdlcumwfgypbvkJxqz0123456789.,!?"
23
24     all_symbols = set(encrypted_text)
25
26     mapping = {}
27
28     # First map the symbols that exist in our frequency order
29     available_targets = list(english_freq_order)
30     for symbol in all_symbols:
31         if len(available_targets) > 0:
32             mapping[symbol] = available_targets.pop(0)
33         else:
34             mapping[symbol] = symbol
35
36     return mapping
37
38 def calculate_transition_probs(text):
39     """Calculate transition probabilities with smoothing."""
40     transitions = defaultdict(lambda: defaultdict(lambda: 0.01))
41     char_counts = defaultdict(lambda: 0)
42
43     # Count occurrences
44     for i in range(len(text)-1):
45         curr, next_char = text[i], text[i+1]
46         transitions[curr][next_char] += 1
47         char_counts[curr] += 1
48
49     # Normalize to get probabilities
50     for char in transitions:
51         total = sum(transitions[char].values())
52         for next_char in transitions[char]:
53             transitions[char][next_char] /= total
54
55     return transitions
56
57 def compute_log_likelihood(text, transitions):
58     """Compute log likelihood of text using transition probabilities."""
59     log_prob = 0
60     for i in range(len(text)-1):
61         curr, next_char = text[i], text[i+1]
62         prob = transitions[curr][next_char]
63         log_prob += np.log(max(prob, 1e-10))
64     return log_prob
65
66 def metropolis_hastings_decrypt(encrypted_text, transitions, n_iterations=10000,
67                                report_interval=100):
68     """Run Metropolis-Hastings algorithm with simulated annealing."""
69     current_mapping = initialize_parameters_by_frequency(encrypted_text)
```

```

69     current_decrypter = Decrypter(current_mapping)
70
71     current_decrypt = current_decrypter.decrypt(encrypted_text)
72     current_score = compute_log_likelihood(current_decrypt, transitions)
73
74     best_score = current_score
75     best_mapping = current_mapping.copy()
76     best_decrypt = current_decrypt
77
78     decryption_log = [(0, current_decrypt[:60])]
79
80     temp = 1.0
81     cooling_rate = 0.9999
82
83     symbols = list(set(encrypted_text))
84
85     for i in range(n_iterations):
86         new_mapping = current_mapping.copy()
87         s1, s2 = random.sample(symbols, 2)
88         new_mapping[s1], new_mapping[s2] = new_mapping[s2], new_mapping[s1]
89
90         new_decrypter = Decrypter(new_mapping)
91         new_decrypt = new_decrypter.decrypt(encrypted_text)
92         new_score = compute_log_likelihood(new_decrypt, transitions)
93
94         # Accept or reject new mapping
95         score_diff = new_score - current_score
96         if score_diff > 0 or random.random() < np.exp(score_diff / temp):
97             current_mapping = new_mapping
98             current_score = new_score
99             current_decrypt = new_decrypt
100
101         if current_score > best_score:
102             best_score = current_score
103             best_mapping = current_mapping.copy()
104             best_decrypt = current_decrypt
105
106         if (i + 1) % report_interval == 0:
107             decryption_log.append((i + 1, current_decrypt[:60]))
108
109         temp *= cooling_rate
110
111     return best_mapping, best_decrypt, decryption_log
112
113 def run_decryption(encrypted_message, training_text):
114     """Run the decryption process and display progress."""
115     # Calculate transition probabilities from training text
116     transitions = calculate_transition_probs(training_text)
117
118     # Run Metropolis-Hastings for 5000 iterations
119     _, decrypted_text, log = metropolis_hastings_decrypt(
120         encrypted_message,
121         transitions,
122         n_iterations=5000,
123         report_interval=100)
124     print("\nMH Iteration | Current Decryption")
125     print("-" * 80)
126     for iteration, text in log:
127         print(f"{iteration:>10} | {text}")
128
129     return decrypted_text
130
131 file_path = '/Users/richardhuang/Documents/GitHub/Projects/Probabilistic Summative/message.txt'
132 with open(file_path, 'r') as file:
133     encrypted_message = file.read().strip()
134
135 decrypted_result = run_decryption(
136     encrypted_message=encrypted_message,
137     training_text=text )

```

Listing 10: solutions/q1.py

(e) The presence of zero values in  $\psi(\alpha, \beta)$  can affect the ergodicity of the Markov chain. If these zero values lead to isolated states, the chain becomes reducible and thus non-ergodic, as it would not converge to a unique stationary distribution.

To prove ergodicity, we must show that the chain remains irreducible and aperiodic. If the chain is still ergodic, we can demonstrate that all states can be reached from any other state, and the greatest common divisor of the lengths of all possible return paths is one.

If the chain is found to be non-ergodic, we can restore ergodicity by adjusting transition probabilities to ensure all states are reachable, introducing dummy states to connect isolated states, or reparameterizing  $\alpha$  and  $\beta$  to ensure all  $\psi(\alpha, \beta)$  values are positive.

(f) In analyzing the decoding approach, it becomes clear that relying solely on symbol probabilities, rather than transition probabilities, may not be sufficient for effective decoding. While symbol probabilities provide insights into the likelihood of individual symbols appearing in the text, they fail to account for the context in which these symbols occur. For instance, let  $P(x)$  represent the probability of a symbol  $x$  occurring in the text. This probability alone does not capture the relationships between symbols. Language is inherently sequential, and the meaning of a symbol often depends on the symbols that precede or follow it. Therefore, using only symbol probabilities can lead to ambiguities and incorrect interpretations.

If we were to employ a second-order Markov chain, which considers the probabilities of pairs of consecutive symbols, we would enhance the contextual information captured in the model. This approach allows for the influence of the previous two symbols on the current symbol to be taken into account. Mathematically, this can be represented as:

$$P(x_t|x_{t-1}, x_{t-2}) = \frac{C(x_{t-2}, x_{t-1}, x_t)}{C(x_{t-2}, x_{t-1})}$$

where  $C(a, b)$  denotes the count of occurrences of the symbols  $a$  and  $b$  together. However, several challenges may arise with this method. The state space grows significantly, necessitating more data to accurately estimate the transition probabilities between pairs of symbols, which can lead to sparse data issues, especially in smaller datasets. Additionally, the increased number of states and transitions results in higher computational costs during both training and inference phases, and there is a risk of overfitting the model to the training data, which may not generalize well to unseen text.

Moreover, if the encryption scheme permits two different symbols to be mapped to the same encrypted value, this introduces ambiguity in the decoding process. For example, if symbols  $a$  and  $b$  both map to the same encrypted value  $y$ , we cannot distinguish between them during decoding. This situation can lead to multiple valid interpretations of the decrypted text, making it challenging to determine the correct original symbols. Consequently, the model may struggle to disambiguate between the possible original symbols, resulting in incorrect or nonsensical outputs.

When considering languages like Chinese, which have a vast number of unique symbols (over 10,000), the challenges become even more pronounced. The high dimensionality increases the complexity of the model, making it difficult to estimate probabilities accurately without a large corpus of text. For instance, if we denote the set of symbols as  $S$  and the number of unique symbols as  $|S| > 10000$ , the number of possible pairs in a second-order Markov chain becomes  $|S|^2$ , leading to a combinatorial explosion in the state space. Additionally, Chinese characters can have different meanings based on context, and a simple Markov model may not effectively capture these nuances. The likelihood of encountering certain pairs or sequences of symbols decreases, leading to sparse data issues that can hinder the model's performance.

In summary, while symbol probabilities provide a foundational understanding of the text, they are insufficient for effective decoding without considering the context provided by transitions. A second-order Markov chain can improve context capture but introduces complexity and potential overfitting. Allowing multiple symbols to map to the same encrypted value creates ambiguity, complicating the decoding process. Finally, applying this approach to languages with a large symbol set, such as Chinese, presents significant challenges in terms of data sparsity and contextual understanding.

## Question 6

(a) The sample plots produced by `toysample` with 200 and 500 iterations are shown below in Figure 15 and Figure 16. The code for the four files are also shown below.

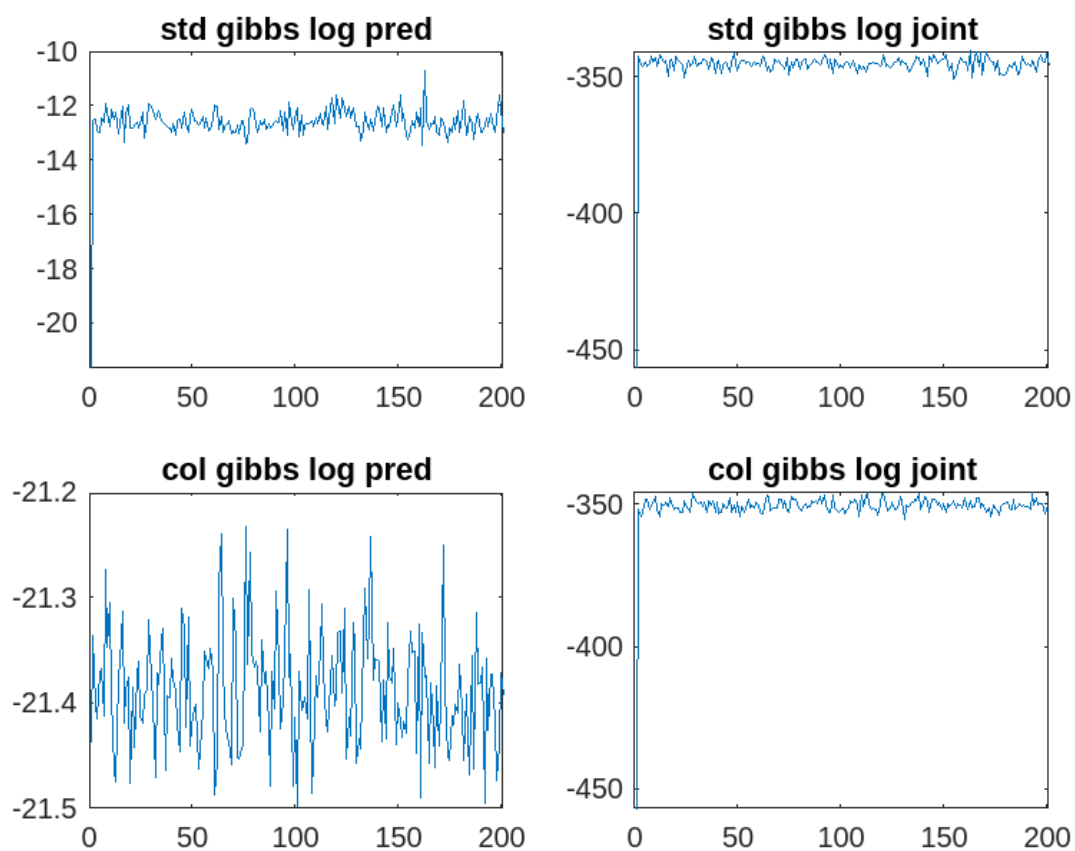


Figure 15: Sample plots produced by `toyexample` try 200 iterations

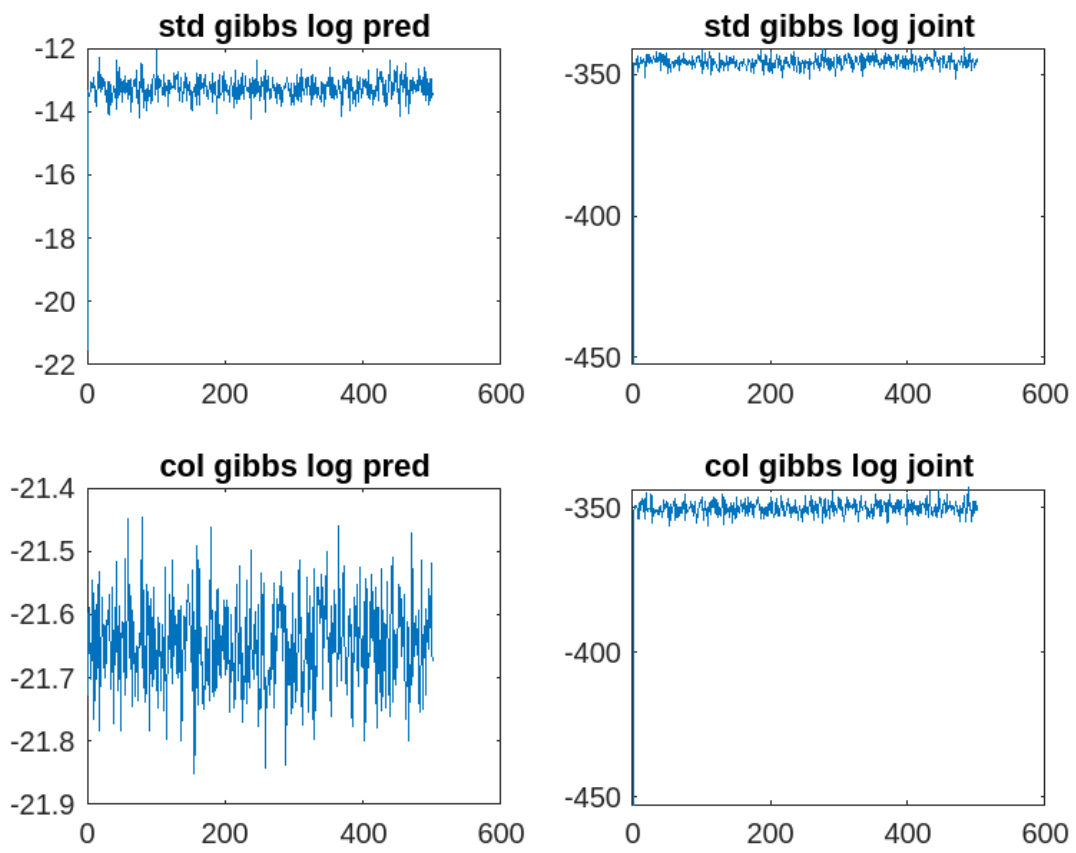


Figure 16: Sample plots produced by toyexample try 500 iterations



## MATLAB code for (a): four files codes

```
1
2 function log_joint = stdgibbs_logjoint(theta, phi, Adk, Bkw, Mk, I, D, K, W, di, wi,
3     ci, citest, Id, Iw, Nd, alpha, beta)
4     % theta: Document-topic distribution (D x K)
5     % phi: Topic-word distribution (K x W)
6     % Adk: Document-topic count matrix (D x K)
7     % Bkw: Topic-word count matrix (K x W)
8     % Mk: Total number of words assigned to each topic (K x 1)
9     % I: Number of (document, word) pairs in training/test sets
10    % D: Number of documents
11    % K: Number of topics
12    % W: Number of unique words
13    % di, wi, ci, citest, Id, Iw, Nd: Various input parameters (details provided in
14        README)
15    % alpha: Dirichlet prior for document-topic distribution
16    % beta: Dirichlet prior for topic-word distribution
17
18    log_joint = 0;
19
20    for d = 1:D
21        theta_d = theta(d, :); % Get theta for document d
22        Adk_d = Adk(d, :) + alpha; % Dirichlet parameters with prior for document d
23        log_joint = log_joint + sum(gammaln(Adk_d)) - gammaln(sum(Adk_d));
24    end
25    log_joint = log_joint - D * K * gammaln(alpha) + D * gammaln(K * alpha);
26
27    % Contribution from phi (topic-word distributions)
28    for k = 1:K
29        phi_k = phi(k, :); % Get phi for topic k
30        Bkw_k = Bkw(k, :) + beta; % Dirichlet parameters with prior for topic k
31        log_joint = log_joint + sum(gammaln(Bkw_k)) - gammaln(sum(Bkw_k));
32    end
33    log_joint = log_joint - K * W * gammaln(beta) + K * gammaln(W * beta);
34 end
```

Listing 11: Code in 'stdgibbs logjoint' file

```

1 function [zi, theta, phi, Adk, Bkw, Mk] = stdgibbs_update(zi, theta, phi, Adk, Bkw, Mk,
2 I, D, K, W, di, wi, ci, citest, Id, Iw, Nd, alpha, beta)
3 % zi: Current topic assignments (cell array with topic for each word in each
4 % document)
5 % theta: Document-topic distribution (D x K)
6 % phi: Topic-word distribution (K x W)
7 % Adk: Document-topic count matrix (D x K)
8 % Bkw: Topic-word count matrix (K x W)
9 % Mk: Total count of words assigned to each topic (K x 1)
10 % I, D, K, W: Number of (doc, word) pairs, documents, topics, and vocabulary size
11 % di, wi: Document and word indices
12 % ci, citest: Training and test counts (if applicable)
13 % Id, Iw, Nd: Document-word pairs and document word counts
14 % alpha, beta: Dirichlet priors for document-topic and topic-word distributions
15
16 for i = 1:I
17     d = di(i);
18     w = wi(i);
19     current_topic = zi{i};
20
21     Adk(d, current_topic) = Adk(d, current_topic) - 1;
22     Bkw(current_topic, w) = Bkw(current_topic, w) - 1;
23     Mk(current_topic) = Mk(current_topic) - 1;
24
25     % Sample theta and phi explicitly using Dirichlet distributions
26     theta(d, :) = dirichrnd(Adk(d, :) + alpha);
27     phi(:, w) = dirichrnd(Bkw(:, w) + beta);
28
29     % Compute conditional probabilities for each topic
30     topic_probs = theta(d, :) .* phi(:, w)';
31
32     % Normalize to get probabilities
33     topic_probs = topic_probs / sum(topic_probs);
34
35     % Sample new topic
36     new_topic = find(mnrnd(1, topic_probs));
37
38     % Update counts with the new topic assignment
39     Adk(d, new_topic) = Adk(d, new_topic) + 1;
40     Bkw(new_topic, w) = Bkw(new_topic, w) + 1;
41     Mk(new_topic) = Mk(new_topic) + 1;
42
43     % Update the topic assignment
44     zi{i} = new_topic;
45 end
46
47 function x = dirichrnd(alpha)
48     x = gamrnd(alpha, 1);
49     x = x / sum(x);
50 end

```

Listing 12: Code in 'stdgibbs update' file

```

1 function [zi, Adk, Bkw, Mk] = colgibbs_update(zi, Adk, Bkw, Mk, I, D, K, W, di, wi, ci,
2         citest, Id, Iw, Nd, alpha, beta)
3     % zi: Current topic assignments (cell array with topic for each word in each
4         document)
5     % Adk: Document-topic count matrix (D x K)
6     % Bkw: Topic-word count matrix (K x W)
7     % Mk: Total count of words assigned to each topic (K x 1)
8     % I, D, K, W: Number of (doc, word) pairs, documents, topics, and vocabulary size
9     % di, wi, ci, citest, Id, Iw, Nd: Various input parameters
10    % alpha, beta: Dirichlet priors for document-topic and topic-word distributions
11
12    % Loop through each (document, word) pair
13    for i = 1:I
14        d = di(i);
15        w = wi(i);
16        current_topic = zi{i}; % Current topic assignment
17
18        Adk(d, current_topic) = Adk(d, current_topic) - 1;
19        Bkw(current_topic, w) = Bkw(current_topic, w) - 1;
20        Mk(current_topic) = Mk(current_topic) - 1;
21
22        % Calculate conditional probabilities for each topic
23        topic_probs = (Adk(d, :) + alpha) .* (Bkw(:, w)' + beta) ./ (Mk + W * beta);
24
25        % Normalize to get probabilities
26        topic_probs = topic_probs / sum(topic_probs);
27
28        % Sample new topic
29        new_topic = find(mnrnd(1, topic_probs));
30
31        % Update counts with the new topic assignment
32        Adk(d, new_topic) = Adk(d, new_topic) + 1;
33        Bkw(new_topic, w) = Bkw(new_topic, w) + 1;
34        Mk(new_topic) = Mk(new_topic) + 1;
35
36        % Update the topic assignment
37        zi{i} = new_topic;
38    end
end

```

Listing 13: Code in 'colgibbs update' file

```

1 function log_joint = colgibbs_logjoint(Adk, Bkw, Mk, I, D, K, W, di, wi, ci, citest,
2   Id, Iw, Nd, alpha, beta)
3   % Adk: Document-topic count matrix (D x K)
4   % Bkw: Topic-word count matrix (K x W)
5   % Mk: Total number of words assigned to each topic (K x 1)
6   % I, D, K, W: Number of (doc, word) pairs, documents, topics, and vocabulary size
7   % di, wi, ci, citest, Id, Iw, Nd: Various input parameters (details provided in
8     README)
9   % alpha, beta: Dirichlet priors for document-topic and topic-word distributions
10
11   log_joint = 0;
12
13   % Contribution from document-topic counts (Adk)
14   for d = 1:D
15       Adk_d = Adk(d, :) + alpha;
16       log_joint = log_joint + sum(gamaln(Adk_d)) - gamaln(sum(Adk_d));
17   end
18   log_joint = log_joint - D * K * gamaln(alpha) + D * gamaln(K * alpha);
19
20   % Contribution from topic-word counts (Bkw)
21   for k = 1:K
22       Bkw_k = Bkw(k, :) + beta;
23       log_joint = log_joint + sum(gamaln(Bkw_k)) - gamaln(sum(Bkw_k));
24   end
25   log_joint = log_joint - K * W * gamaln(beta) + K * gamaln(W * beta);
26 end

```

Listing 14: Code in 'colgibbs logjoint' file

(b) Based on the autocorrelation plots for 200 and 500 iterations in Figure 17 and Figure 18, the plots with 500 iterations show significantly reduced noise in the autocorrelation values compared to those with 200 iterations. Specifically, the autocorrelation values for both the standard and collapsed Gibbs samplers tend to approach zero more consistently within the lag range for 500 iterations, indicating that the chains have decorrelated more effectively. This suggests that increasing the number of iterations helps in reducing the noise in the autocorrelations, thereby providing a more reliable representation of the posterior.

For a representative set of samples from the posterior, it is essential to discard the burn-in period and select samples that are sufficiently uncorrelated. The 500-iteration results indicate that a lag of approximately 10 is sufficient to ensure minimal autocorrelation, based on the rapid drop in the autocorrelation values. Therefore, it is recommended to thin the samples by taking every 10th sample after the burn-in period to construct a representative posterior. This decision is justified as it balances computational efficiency with the need for decorrelation, as evidenced by the reduced noise and consistent trends in the 500-iteration plots.

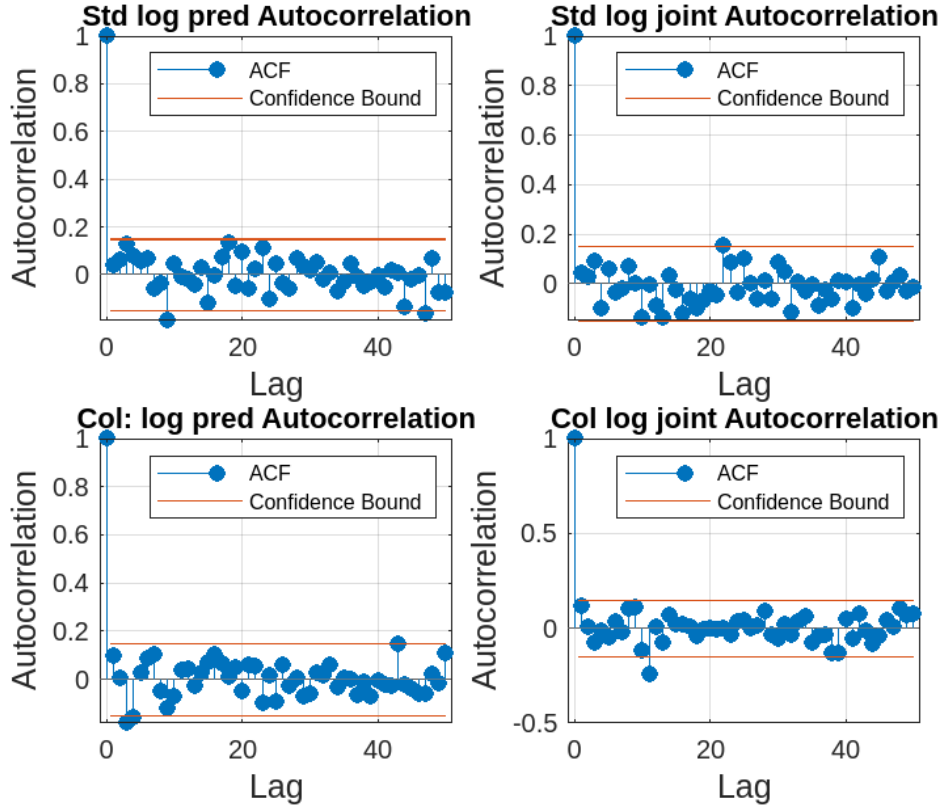


Figure 17: Autocorrelation plot by 200 iterations

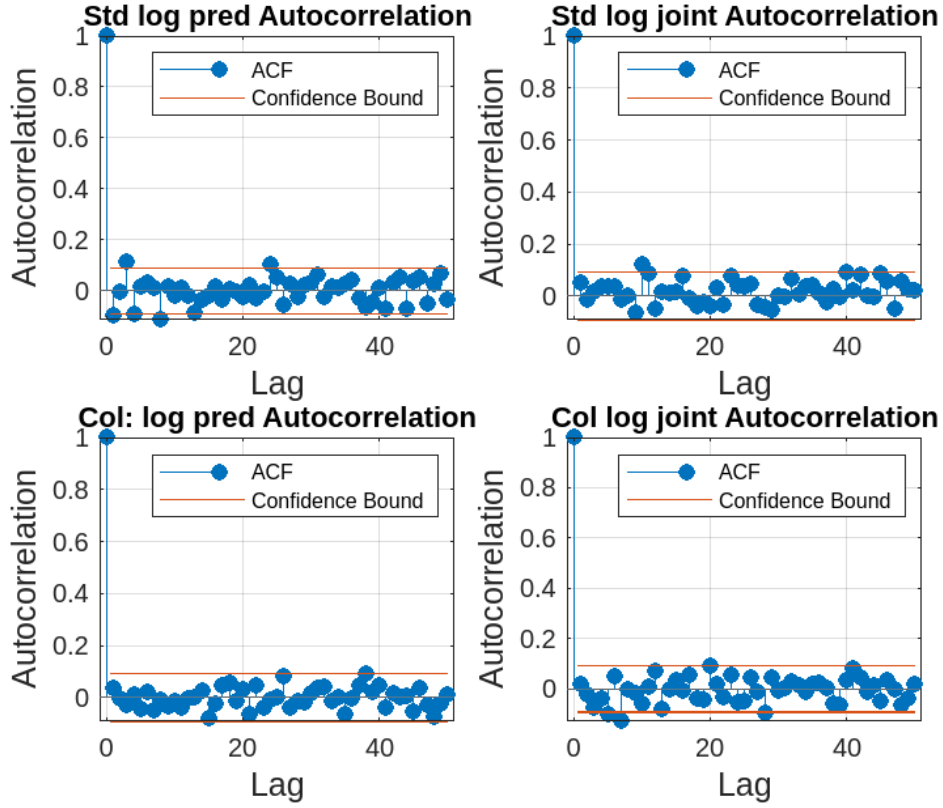


Figure 18: Autocorrelation plot by 500 iterations

(c) Based on the plots, the collapsed Gibbs sampler appears to converge faster compared to the standard Gibbs sampler, particularly in terms of the stability of the log predictive and log joint probabilities. The collapsed Gibbs sampler exhibits less variability in the log joint probabilities and log predictive probabilities after a few iterations, whereas the standard Gibbs sampler shows slightly more fluctuation.

The collapsed Gibbs sampler integrates out parameters during sampling, which reduces the dimensionality of the sampling process and often leads to faster mixing and convergence. This advantage is particularly evident in smaller datasets, such as the toy example used here, where collapsed Gibbs benefits from the reduced variance in the sampling process. The faster convergence is observed in the stability of the collapsed Gibbs plots compared to the standard Gibbs.

(d) In Figure 19 and Figure 20, when varying  $\alpha$ ,  $\beta$ , and  $K$ , the effects on the posterior and predictive performance of the model become evident. Increasing  $K$ , the number of topics, leads to more granular topic modeling. While this can improve the model's capacity to capture finer details in the data, setting  $K$  too high may result in overfitting, as the model tries to assign topics to smaller and less meaningful partitions of the data. For instance, with  $K = 5$ , the model demonstrates stable convergence, and the log predictive probabilities exhibit smoother patterns compared to when  $K$  is lower.

The hyperparameter  $\alpha$ , which governs the Dirichlet prior over topic distributions, plays a crucial role in determining how topics are distributed within documents. Larger values of  $\alpha$  encourage documents to exhibit multiple topics more uniformly, thereby promoting diversity in topic assignments. This is reflected in the smoother predictive log probability plots observed with higher  $\alpha$ , as the regularization effect helps the model generalize better across the dataset. Similarly,  $\beta$ , the hyperparameter for the Dirichlet prior over words within topics, impacts the diversity of words associated with each topic. A higher  $\beta$  ensures that word distributions within topics are more balanced, reducing the prominence of overly specific words. This smoothing effect contributes to more consistent and stable log joint and log predictive scores.

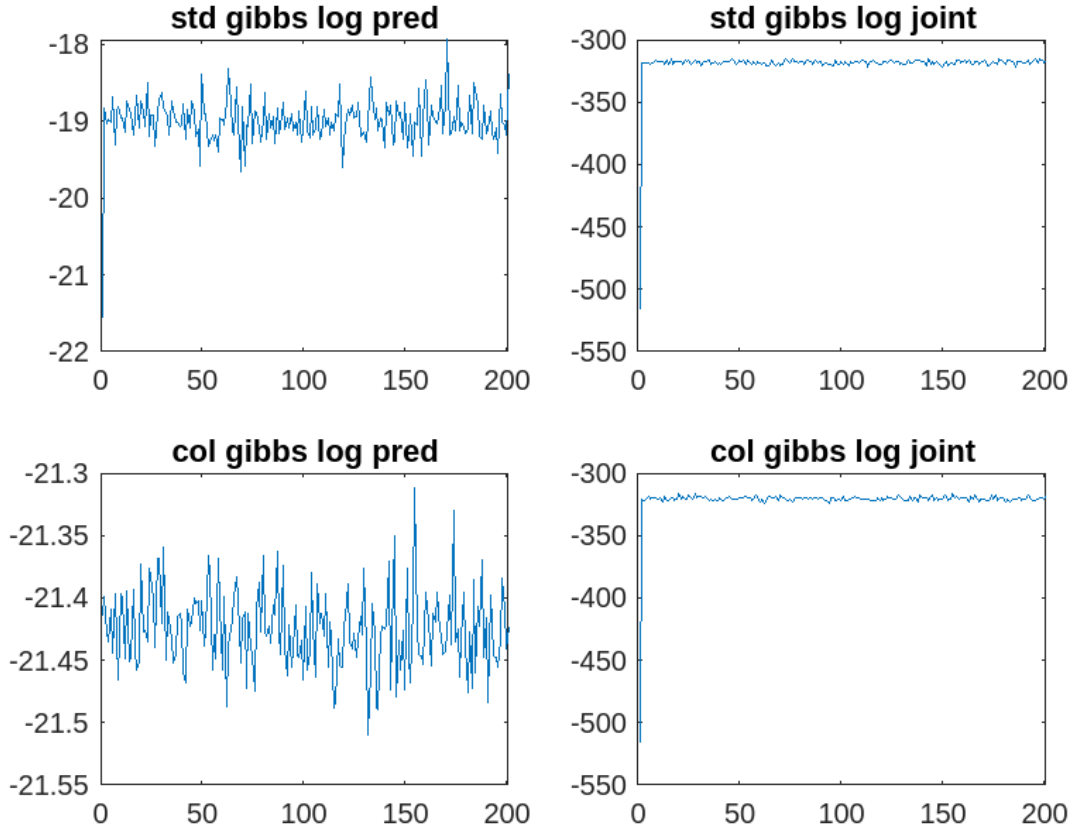


Figure 19:  $K = 5$ ,  $\alpha = 5$ ,  $\beta = 5$

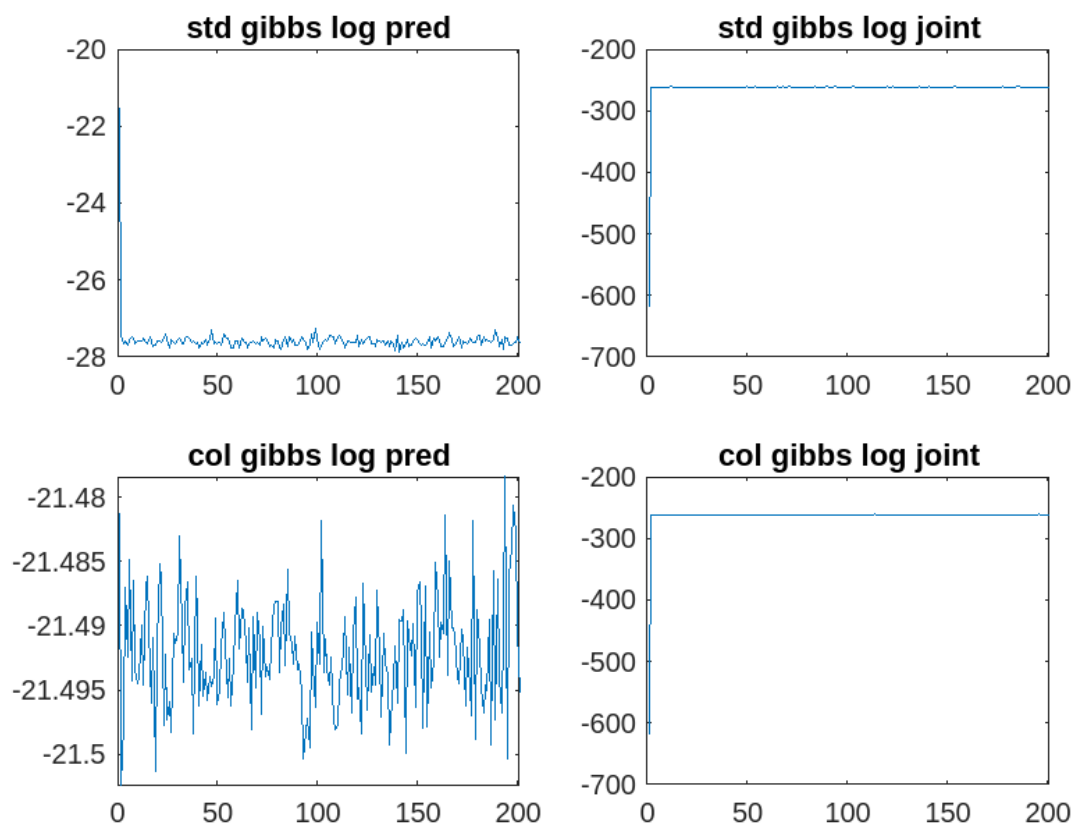


Figure 20:  $K = 10$ ,  $\alpha = 10$ ,  $\beta = 10$



## Question 7

(a) To find the local extrema of the function  $f(x, y) = x + 2y$  subject to the constraint  $y^2 + xy = 1$ , we will use the method of Lagrange multipliers.

We define the Lagrangian function  $\mathcal{L}(x, y, \lambda)$  as follows:

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

where  $\lambda$  is the Lagrange multiplier.

Next, we find the partial derivatives of  $\mathcal{L}$  with respect to  $x$ ,  $y$ , and  $\lambda$ , and set them equal to zero.

Partial derivative with respect to  $x$ :

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

From this, we get:

$$\lambda = \frac{1}{y} \quad (\text{assuming } y \neq 0)$$

Partial derivative with respect to  $y$ :

$$\frac{\partial \mathcal{L}}{\partial y} = 2 - 2\lambda y - \lambda x = 0$$

Substituting  $\lambda = \frac{1}{y}$  into this equation, we get:

$$2 - 2 \cdot \frac{1}{y} \cdot y - \frac{x}{y} = 0$$

Simplifying, this becomes:

$$2 - 2 - \frac{x}{y} = 0 \Rightarrow x = 0$$

Partial derivative with respect to  $\lambda$ :

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

Substituting  $x = 0$ , we obtain:

$$1 - y^2 = 0 \Rightarrow y^2 = 1$$

Thus,  $y = \pm 1$ .

From the equations above, we find the following critical points: When  $y = 1$ ,  $x = 0$ , giving the point  $(0, 1)$ . When  $y = -1$ ,  $x = 0$ , giving the point  $(0, -1)$ .

The local extrema of  $f(x, y) = x + 2y$  subject to the constraint  $y^2 + xy = 1$  occur at the points  $(0, 1)$  and  $(0, -1)$ . These are the points at which the local extrema occur, as required by the problem.

(b) To compute  $\ln(a)$  using Newton's method, we need to define a function  $f(x, a)$  such that the root of this function will give  $x = \ln(a)$ .

Define the Function  $f(x, a)$  We know that if  $x = \ln(a)$ , then  $e^x = a$ . Thus, we can define the function  $f(x, a)$  as:

$$f(x, a) = e^x - a$$

The root of this function (i.e., where  $f(x, a) = 0$ ) will give  $x = \ln(a)$ .

Derive the Update Equation for Newton's Method Newton's method for finding the root of a function  $f(x)$  involves the iterative update:

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

To apply this to our function  $f(x, a) = e^x - a$ , we first need to compute  $f'(x, a)$ .

Compute  $f'(x, a)$ :

$$f'(x, a) = \frac{d}{dx}(e^x - a) = e^x$$

Substitute into the Newton's Method Update Equation: Using  $f(x, a) = e^x - a$  and  $f'(x, a) = e^x$ , the update equation becomes:

$$x_{n+1} = x_n - \frac{e^{x_n} - a}{e^{x_n}}$$

Simplifying this expression, we get:

$$x_{n+1} = x_n - 1 + \frac{a}{e^{x_n}}$$

or equivalently,

$$x_{n+1} = x_n + \frac{a - e^{x_n}}{e^{x_n}}$$

Thus, to compute  $\ln(a)$  using Newton's method:

- i. Define the function  $f(x, a) = e^x - a$ .
- ii. Use the iterative update equation:

$$x_{n+1} = x_n + \frac{a - e^{x_n}}{e^{x_n}}$$

Starting with an initial guess  $x_0$ , we can iterate using this formula until convergence to obtain  $x = \ln(a)$ .

# Question 8

(a) Note that for any nonzero scalar  $\alpha$  and vector  $x \neq 0$ , we have:

$$R_A(\alpha x) = \frac{(\alpha x)^T A(\alpha x)}{(\alpha x)^T (\alpha x)} = \frac{\alpha^2 x^T A x}{\alpha^2 x^T x} = R_A(x).$$

This shows that  $R_A(x)$  is homogeneous of degree 0, meaning it is sufficient to consider vectors  $x$  with  $\|x\| = 1$  because scaling  $x$  by a nonzero constant does not change  $R_A(x)$ . Therefore,

$$\sup_{x \in \mathbb{R}^n} R_A(x) = \sup_{x \in S} R_A(x),$$

where  $S = \{x \in \mathbb{R}^n \mid \|x\| = 1\}$  is the unit sphere.

The unit sphere  $S$  in  $\mathbb{R}^n$  is a compact set. Since  $q_A(x)$  is continuous (as given),  $R_A(x)$  is also continuous on  $S$  because it is the ratio of two continuous functions (and  $\|x\| = 1$  ensures the denominator is nonzero).

By the extreme value theorem, a continuous function on a compact set attains its maximum and minimum. Since  $R_A(x)$  is continuous on the compact set  $S$ , there exists some  $x^* \in S$  such that

$$R_A(x^*) = \max_{x \in S} R_A(x) = \sup_{x \in \mathbb{R}^n} R_A(x).$$

Thus,  $x^*$  maximizes  $R_A(x)$  on  $\mathbb{R}^n$ , and  $R_A(x^*)$  is the largest eigenvalue of  $A$ , with  $x^*$  as a corresponding eigenvector.

(b) To show that  $R_A(x) \leq \lambda_1$  for any vector  $x \in \mathbb{R}^n$ .

Let  $A$  be a symmetric matrix with eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  and corresponding orthonormal eigenvectors  $\xi_1, \xi_2, \dots, \xi_n$ , which form an orthonormal basis (ONB) of  $\mathbb{R}^n$ . Any vector  $x \in \mathbb{R}^n$  can be represented as:

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

The quadratic form  $q_A(x) = x^T A x$  can be expanded as follows:

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

Expanding the transpose, we get:

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

Since  $\xi_i$  are eigenvectors of  $A$  with corresponding eigenvalues  $\lambda_i$ , we have  $A \xi_i = \lambda_i \xi_i$ . Therefore,

$$\xi_i^T A \xi_j = \lambda_i \xi_i^T \xi_j.$$

Because  $\{\xi_i\}$  form an orthonormal basis,  $\xi_i^T \xi_j = 0$  for  $i \neq j$  and  $\xi_i^T \xi_i = 1$ . This means that only terms where  $i = j$  will survive in the summation, simplifying the expression to:

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

The squared norm  $\|x\|^2$  is defined as:

$$\|x\|^2 = x^\top x.$$

Substitute  $x = \sum_{i=1}^n (\xi_i^\top x) \xi_i$ :

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

Expanding the inner product, we get:

$$x^\top x = \sum_{i=1}^n \sum_{j=1}^n (\xi_i^\top x) (\xi_j^\top x) \xi_i^\top \xi_j.$$

Since  $\{\xi_i\}$  form an orthonormal basis,  $\xi_i^\top \xi_j = 0$  for  $i \neq j$  and  $\xi_i^\top \xi_i = 1$ , we get:

$$x^\top x = \sum_{i=1}^n (\xi_i^\top x)^2.$$

Now,  $R_A(x) = \frac{q_A(x)}{\|x\|^2}$ . The norm  $\|x\|^2$  is given by:

$$\|x\|^2 = x^\top x = \sum_{i=1}^n (\xi_i^\top x)^2.$$

Thus:

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

Since  $\lambda_1$  is the largest eigenvalue, we have  $\lambda_i \leq \lambda_1$  for all  $i$ . Therefore:

$$R_A(x) = \frac{\sum_{i=1}^n (\xi_i^\top x)^2 \lambda_i}{\sum_{i=1}^n (\xi_i^\top x)^2} \leq \frac{\sum_{i=1}^n (\xi_i^\top x)^2 \lambda_1}{\sum_{i=1}^n (\xi_i^\top x)^2} = \lambda_1.$$

(c) To show that if  $x \in \mathbb{R}^n$  is not contained in the span of the eigenvectors corresponding to the largest eigenvalue  $\lambda_1$ , then  $R_A(x) < \lambda_1$ .

Since  $A$  is symmetric, it has a complete set of orthonormal eigenvectors  $\{\xi_1, \dots, \xi_k, \xi_{k+1}, \dots, \xi_n\}$ .

Since  $x$  is not contained in the span of  $\{\xi_1, \dots, \xi_k\}$ , it must have a non-zero component in the directions of the eigenvectors  $\{\xi_{k+1}, \dots, \xi_n\}$  associated with eigenvalues less than  $\lambda_1$ . This means that:

$$\sum_{i=k+1}^n (\xi_i^\top x)^2 > 0.$$

In the expression for  $R_A(x)$ , the numerator  $\sum_{i=1}^n \lambda_i (\xi_i^\top x)^2$  includes terms with eigenvalues  $\lambda_{k+1}, \dots, \lambda_n$  that are strictly less than  $\lambda_1$ . Since  $x \notin \text{span}\{\xi_1, \dots, \xi_k\}$ , there exists some components of  $x$  in the directions of eigenvectors  $\{\xi_{k+1}, \dots, \xi_n\}$  associated with eigenvalues smaller than  $\lambda_1$ . This results in the inequality:

$$\sum_{i=1}^n \lambda_i (\xi_i^\top x)^2 < \lambda_1 \sum_{i=1}^n (\xi_i^\top x)^2 = \lambda_1 \|x\|^2.$$

Since the Rayleigh quotient  $R_A(x)$  is defined as

$$R_A(x) = \frac{\sum_{i=1}^n \lambda_i (\xi_i^\top x)^2}{\|x\|^2},$$

and we have shown that

$$\sum_{i=1}^n \lambda_i (\xi_i^\top x)^2 < \lambda_1 \|x\|^2,$$

it follows that:

$$R_A(x) < \lambda_1.$$

Therefore, if  $x \notin \text{span}\{\xi_1, \dots, \xi_k\}$ , then  $R_A(x) < \lambda_1$ , as required.