

# COMP0086 Summative Assignment

Nov 14, 2022

## Question 1

- (a) Our sample space for images is  $\{0, 1\}^D$ , a binary space with  $D$  dimensions, the number of pixels in the image. Thus, picking the exponential family best suited on this sample space is the  $D$ -dimensional multivariate Bernoulli distribution which shares the same sample space. On the other hand, a  $D$ -dimensional multivariate Gaussian has the sample space  $\mathbb{R}^D$ , which does not match the sample space of our data. To match our data sample space, we would have to define additional mapping between our data and model spaces which adds unnecessary complexity. Thus it would be inappropriate to model this dataset of images with a multivariate Gaussian.
- (b) For  $\mathcal{D} := \{x^{(n)}\}_{n=1}^N$  a data set of  $N$  images, the joint likelihood (assuming images are independently and identically distributed) is the product of  $N$   $D$ -dimensional multivariate Bernoulli distributions, one for each image:

$$P(\mathcal{D}|\mathbf{p}) = \prod_{n=1}^N P(x^{(n)}|\mathbf{p})$$

Substituting the  $D$ -dimensional multivariate Bernoulli:

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

Taking the logarithm of this, we get the log likelihood:

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

Note that since the logarithm of the mean is a monotonic increasing function on  $\mathbb{R}_+$ , the maximisers and minimisers of the likelihood do not change.

To solve for the maximum likelihood estimate,  $\hat{p}_d$ , we can take the derivative of  $\mathcal{L}(\mathcal{D}|\mathbf{p})$  with respect to  $p_d$ , the  $d^{th}$  element of  $\mathbf{p}$ :

$$\frac{\partial \mathcal{L}(\mathcal{D}|\mathbf{p})}{\partial p_d} = \sum_{n=1}^N \left( \frac{x_d^{(n)}}{p_d} - \frac{1 - x_d^{(n)}}{1 - p_d} \right)$$

$$\frac{\partial \mathcal{L}(\mathcal{D}|\mathbf{p})}{\partial p_d} = \frac{\sum_{n=1}^N x_d^{(n)}}{p_d} - \frac{\sum_{n=1}^N (1 - x_d^{(n)})}{1 - p_d}$$

and set the derivative to zero and solve for  $\hat{p}_d$ :

$$\begin{aligned} \frac{\sum_{n=1}^N x_d^{(n)}}{\hat{p}_d} - \frac{\sum_{n=1}^N (1 - x_d^{(n)})}{1 - \hat{p}_d} &= 0 \\ \sum_{n=1}^N x_d^{(n)} - \hat{p}_d \sum_{n=1}^N x_d^{(n)} - \hat{p}_d \cdot N + \hat{p}_d \sum_{n=1}^N x_d^{(n)} &= 0 \\ \hat{p}_d &= \frac{1}{N} \sum_{n=1}^N x_d^{(n)} \end{aligned}$$

Moreover, the second derivative with respect to  $p_d$ :

$$\frac{\partial \mathcal{L}(\mathcal{D}|\mathbf{p})}{\partial p_d^2} = \frac{-\sum_{n=1}^N x_d^{(n)}}{p_d^2} + \frac{-\sum_{n=1}^N (1 - x_d^{(n)})}{(1 - p_d)^2}$$

For a maximum, we need to show that the second derivative is negative. Since  $x_d^{(n)} \in \{0, 1\}$ , in the worst case, of  $N = 1$ , the single pixel must either be white ( $\sum_{n=1}^N x_d^{(n)} > 0$ ) or black ( $\sum_{n=1}^N 1 - x_d^{(n)} > 0$ ) so  $\frac{\partial \mathcal{L}(\mathcal{D}|\mathbf{p})}{\partial p_d^2} < 0$  will be guaranteed and  $\hat{p}_d$  is a maximum as required for the maximum likelihood estimate.

Because we assume that each pixel is independent (we are taking the product of  $D$  one dimensional Bernoulli distributions), we can express the maximum likelihood for  $\mathbf{p}$  in matrix form as  $\hat{\mathbf{p}}$ :

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

(c) From Bayes' Theorem:

$$P(\mathbf{p}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathbf{p})P(\mathbf{p})}{P(\mathcal{D})}$$

Taking the logarithm:

$$\mathcal{L}(\mathbf{p}|\mathcal{D}) = \mathcal{L}(\mathcal{D}|\mathbf{p}) + \mathcal{L}(\mathbf{p}) - \mathcal{L}(\mathcal{D})$$

Taking the derivative with respect to  $p_d$ :

$$\frac{\partial \mathcal{L}(\mathbf{p}|\mathcal{D})}{\partial p_d} = \frac{\partial \mathcal{L}(\mathcal{D}|\mathbf{p})}{\partial p_d} + \frac{\partial \mathcal{L}(\mathbf{p})}{\partial p_d}$$

where  $\frac{\partial \mathcal{L}(\mathcal{D})}{\partial p_d} = 0$  because it doesn't depend on  $p_d$ .

We know (b):

$$\frac{\partial \mathcal{L}(\mathcal{D}|\mathbf{p})}{\partial p_d} = \frac{\sum_{n=1}^N x_d^{(n)}}{p_d} - \frac{\sum_{n=1}^N (1 - x_d^{(n)})}{1 - p_d}$$

For the second term  $\frac{\partial \mathcal{L}(\mathbf{p})}{\partial p_d}$ , we start with  $P(\mathbf{p})$ :

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

Assuming independent Beta priors on the parameters  $p_d$ :

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

Taking the logarithm:

$$\mathcal{L}(\mathbf{p}) = \sum_{d=1}^D -\log(B(\alpha, \beta)) + (\alpha - 1) \log p_d + (\beta - 1) \log(1 - p_d)$$

Taking the derivative with respect to  $p_d$ :

$$\frac{\partial \mathcal{L}(\mathbf{p})}{\partial p_d} = \frac{(\alpha - 1)}{p_d} - \frac{(\beta - 1)}{1 - p_d}$$

Since we are only concerned with  $p_d$ , we are only left with a single element of the summation pertaining to  $p_d$ .

Combining to have an expression for the log posterior derivative  $\frac{\partial \mathcal{L}(\mathbf{p}|\mathcal{D})}{\partial p_d}$ :

$$\begin{aligned} \frac{\partial \mathcal{L}(\mathbf{p}|\mathcal{D})}{\partial p_d} &= \frac{\sum_{n=1}^N x_d^{(n)}}{p_d} - \frac{\sum_{n=1}^N (1 - x_d^{(n)})}{1 - p_d} + \frac{(\alpha - 1)}{p_d} - \frac{(\beta - 1)}{1 - p_d} \\ \frac{\partial \mathcal{L}(\mathbf{p}|\mathcal{D})}{\partial p_d} &= \frac{(\alpha - 1) + \sum_{n=1}^N x_d^{(n)}}{p_d} - \frac{(\beta - 1) + \sum_{n=1}^N (1 - x_d^{(n)})}{1 - p_d} \end{aligned}$$

To find the maximum a posteriori (MAP) estimate  $\hat{p}_d$  with  $\frac{\partial \mathcal{L}(\mathbf{p}|\mathcal{D})}{\partial p_d} = 0$ :

$$\begin{aligned} 0 &= \frac{(\alpha - 1) + \sum_{n=1}^N x_d^{(n)}}{\hat{p}_d} - \frac{(\beta - 1) + \sum_{n=1}^N (1 - x_d^{(n)})}{1 - \hat{p}_d} \\ 0 &= (1 - \hat{p}_d)(\alpha - 1) + (1 - \hat{p}_d) \left( \sum_{n=1}^N x_d^{(n)} \right) - \hat{p}_d(\beta - 1) - \hat{p}_d \left( \sum_{n=1}^N (1 - x_d^{(n)}) \right) \\ 0 &= (\alpha - \alpha \hat{p}_d + \hat{p}_d - 1) + \left( \sum_{n=1}^N x_d^{(n)} - \hat{p}_d \sum_{n=1}^N x_d^{(n)} \right) - (\hat{p}_d \beta - \hat{p}_d) - \left( \hat{p}_d \cdot N - \hat{p}_d \sum_{n=1}^N x_d^{(n)} \right) \end{aligned}$$

Cancelling the  $\hat{p}_d \sum_{n=1}^N x_d^{(n)}$  terms:

$$0 = \alpha - \alpha \hat{p}_d + \hat{p}_d - 1 + \sum_{n=1}^N x_d^{(n)} - \hat{p}_d \beta + \hat{p}_d - \hat{p}_d \cdot N$$

$$0 = \hat{p}_d(2 - \alpha - \beta - N) + \alpha - 1 + \sum_{n=1}^N x_d^{(n)}$$

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

To show that this is a maximum, the second derivative is:

$$\frac{\partial^2 \mathcal{L}(\mathbf{p}|\mathcal{D})}{(\partial p_d)^2} = \frac{(1 - \alpha) - \sum_{n=1}^N x_d^{(n)}}{(p_d)^2} + \frac{(1 - \beta) - \sum_{n=1}^N (1 - x_d^{(n)})}{(1 - p_d)^2}$$

For a maximum, we need  $\frac{\partial^2 \mathcal{L}(\mathbf{p}|\mathcal{D})}{(\partial p_d)^2} < 0$  meaning that we need at least one of the strict inequalities  $\alpha < 1 - \sum_{n=1}^N x_d^{(n)}$  or  $\beta < 1 - \sum_{n=1}^N (1 - x_d^{(n)})$  to be satisfied, where the other can be  $\leq$ . The Beta distribution requires  $\alpha > 0$  and  $\beta > 0$  so this requirement will always be satisfied (in the worst case of a single image, either  $x_d^{(1)} = 1$  or  $1 - x_d^{(1)} = 1$ ).

Due to independence of our likelihood and priors for each dimension, we can express the maximum a priori for  $\mathbf{p}$  in matrix form as  $\hat{\mathbf{p}}$ :

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

(d&e) The Python code for MLE and MAP:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4
5 def _compute_maximum_likelihood_estimate(x: np.ndarray) -> np.ndarray:
6     """
7     X: numpy array of shape (N, D)
8     """
9     return np.mean(x, axis=0)
10
11
12 def _compute_maximum_a_priori_estimate(
13     x: np.ndarray, alpha: float, beta: float
14 ) -> np.ndarray:
15     """
16     X: numpy array of shape (N, D)
17     alpha: param of prior distribution
18     beta: param of prior distribution
19     """
20
21     n, _ = x.shape
22     return (alpha - 1 + np.sum(x, axis=0)) / (n + alpha + beta - 2)
23
24
25 def d(x, figure_path, figure_title):
26     maximum_likelihood = _compute_maximum_likelihood_estimate(x)
27     plt.figure()
28     plt.imshow(
29         np.reshape(maximum_likelihood, (8, 8)),
30         interpolation="None",
31     )
32     plt.colorbar()
33     plt.axis("off")
34     plt.title(figure_title)
35     plt.savefig(figure_path)
36
37
38 def e(x, alpha, beta, figure_path, figure_title):
39     maximum_a_priori = _compute_maximum_a_priori_estimate(x, alpha, beta)
40     plt.figure()
41     plt.imshow(
42         np.reshape(maximum_a_priori, (8, 8)),
43         interpolation="None",
44     )
45     plt.colorbar()
46     plt.axis("off")
47     plt.title(figure_title)
48     plt.savefig(f"{figure_path}.png")
49
50     maximum_likelihood = _compute_maximum_likelihood_estimate(x)
51     plt.figure()
52     plt.imshow(
53         np.reshape(maximum_a_priori - maximum_likelihood, (8, 8)),
54         interpolation="None",
55     )
56     plt.colorbar()
57     plt.axis("off")
58     plt.title(f"MAP vs MLE")
59     plt.savefig(f"{figure_path}-mle-vs-map.png")
```

src/solutions/q1.py

Displaying the learned parameters:

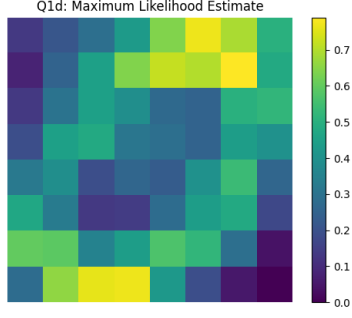


Figure 1: ML parameters

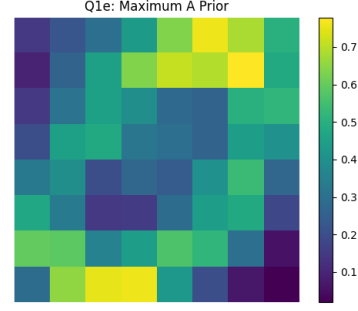


Figure 2: MAP parameters

Comparing the equations:

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

and

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

As the number of data points increases, the MAP estimate approaches  $\frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)}$ , the MLE. This makes sense because as our data set gets bigger, we are less reliant on our prior. Moreover, if a specific pixel in all of the images of our data set are white or all black, the MLE for that pixel will be binary. This may not be representative of our intuitions about image pixels, as there should be some non-zero probability of a pixel being black or white. By introducing an appropriate prior we can ensure that the probability of that pixel will never be exactly zero or one. In our case, with a Beta(3,3) prior on each pixel, our parameter values are biased to be closer to 0.5 and to never be at the extremities 0 and 1. We can see this in Figure 2 where the range of our parameters is smaller than the range of Figure 1. Figure 3 visualises  $\hat{\mathbf{p}}^{MAP} - \hat{\mathbf{p}}^{MLE}$  and we can see that for likelihoods greater than 0.5 in the MLE, the MAP has a lower value and for likelihoods less than 0.5, the MAP has a higher value, confirming our intuitions.

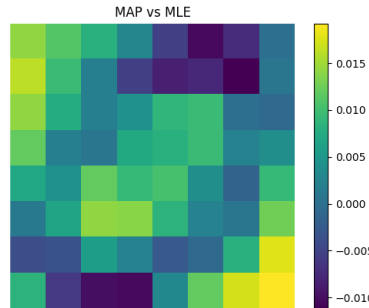


Figure 3:  $\hat{\mathbf{p}}^{MAP} - \hat{\mathbf{p}}^{MLE}$

Priors can also help ensure numerical stability during calculations. The logarithm of zero is negative infinity, so having if the MLE is zero it can be problematic for log-likelihoods calculations whereas MAP can ensure non-zero probabilities. We can see in Figure 1 and Figure 2 that the MLE has parameter values of zero whereas the MAP does not. Interestingly, when  $\alpha = 1$  and  $\beta = 1$ ,  $\hat{\mathbf{p}}^{MLE} = \hat{\mathbf{p}}^{MAP}$ . Intuitively this is when the prior is a uniform distribution and so there aren't any biases on the location of  $\mathbf{p}$  and we recover the MLE.

On the other hand, a mis-specified prior can be problematic, as the estimated parameters might be skewed by the prior and not properly represent the underlying data generating process, this can result in parameter estimates that are worse than using the MLE.

## Question 2

- (a) When all  $D$  components are generated from a Bernoulli distribution with  $p_d = 0.5$ , we have the likelihood function for model  $M_1$ :

$$P(\mathbf{x}^{(n)}|\mathbf{p}^{(1)} = [0.5, 0.5, \dots, 0.5]^T, M_1) = \prod_{d=1}^D (0.5)^{x_d^{(n)}} (0.5)^{1-x_d^{(n)}}$$

- (b) When all  $D$  components are generated from Bernoulli distributions with unknown, but identical,  $p_d$ , we have the likelihood function for model  $M_2$ :

$$P(\mathbf{x}^{(n)}|\mathbf{p}^{(2)} = [p_d, p_d, \dots, p_d]^T, M_2) = \prod_{d'=1}^D p_d^{x_{d'}^{(n)}} (1 - p_d)^{1-x_{d'}^{(n)}}$$

- (c) When each component is Bernoulli distributed with separate, unknown  $p_d$ , we have the likelihood function for model  $M_3$ :

$$P(\mathbf{x}^{(n)}|\mathbf{p}^{(3)} = [p_1, p_2, \dots, p_D]^T, M_3) = \prod_{d=1}^D p_d^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}}$$

For each model  $M_i$ , we can marginalise out  $\mathbf{p}^{(i)}$  to get  $P(\{\mathbf{x}^{(n)}\}_{n=1}^N|M_i)$ :

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N|M_i) = \int_0^1 \dots \int_0^1 P(\{\mathbf{x}^{(n)}\}_{n=1}^N|p_d, M_i) P(p_d|M_i) dp_1 \dots dp_D$$

where  $d = 1, \dots, D$  and  $\{\mathbf{x}^{(n)}\}_{n=1}^N$  is our data set.

Given that the prior of any unknown probabilities is uniform, i.e.  $P(p_d|M_i) = 1$ . We can simplify:

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N|M_i) = \int_0^1 \dots \int_0^1 P(\{\mathbf{x}^{(n)}\}_{n=1}^N|p_d, M_i) dp_1 \dots dp_D$$

For  $M_1$ , we have that all pixels have probability 0.5:

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N|M_1) = \int_0^1 \dots \int_0^1 \prod_{j=1}^N \prod_{d=1}^D (0.5)^{x_d^{(n)}} (1 - 0.5)^{1-x_d^{(n)}} d\theta_1 \dots d\theta_D$$

We can remove the integrals and knowing that either  $x_d^{(n)}$  or  $1 - x_d^{(n)}$  will be 1 and the other zero, we can simplify  $(0.5)^{x_d^{(n)}} (1 - 0.5)^{1-x_d^{(n)}}$  to 0.5:

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N|M_1) = \prod_{j=1}^N \prod_{d=1}^D (0.5)$$

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N|M_1) = (0.5)^{N \cdot D}$$



For  $M_2$ , we have that all pixels share some probability  $p_d$  so we only need to integrate over a single variable  $p_d$ :

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N | M_2) = \int_0^1 \prod_{n=1}^N \prod_{d'=1}^D p_d^{x_{d'}^{(n)}} (1 - p_d)^{1-x_{d'}^{(n)}} dp_d$$

Changing the products to sums:

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N | M_2) = \int_0^1 p_d^{\sum_{n=1}^N \sum_{d'=1}^D x_{d'}^{(n)}} (1 - p_d)^{\sum_{n=1}^N \sum_{d'=1}^D 1-x_{d'}^{(n)}} dp_d$$

Rewriting:

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N | M_2) = \int_0^1 p_d^k (1 - p_d)^{N \cdot D - k} dp_d$$

where  $k = \sum_{n=1}^N \sum_{d'=1}^D x_{d'}^{(n)}$ .

This integral is the beta function:

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N | M_2) = \frac{k!(N \cdot D - k)!}{(N \cdot D + 1)!}$$

For  $M_3$ , we need an integral for each  $p_d$ :

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N | M_3) = \int_0^1 \dots \int_0^1 \prod_{n=1}^N \prod_{d=1}^D p_d^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}} dp_1 \dots dp_D$$

We can separate the integrals to only contain the relevant  $p_d$ :

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N | M_3) = \prod_{d=1}^D \left( \int_0^1 \prod_{n=1}^N p_d^{x_d^{(n)}} (1 - p_d)^{1-x_d^{(n)}} dp_d \right)$$

Changing the products to sums:

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N | M_3) = \prod_{d=1}^D \left( \int_0^1 p_d^{\sum_{n=1}^N x_d^{(n)}} (1 - p_d)^{\sum_{n=1}^N 1-x_d^{(n)}} dp_d \right)$$

In this case, we have the product of integrals where each evaluates to a beta function:

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N | M_3) = \prod_{d=1}^D \frac{k_d!(N - k_d)!}{(N + 1)!}$$

where  $k_d = \sum_{n=1}^N x_d^{(n)}$ .

The posterior probability of a model  $M_i$  can be expressed:

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

We only have three models, so in this case the normalisation  $P(\{\mathbf{x}^{(n)}\}_{n=1}^N)$  can be expressed as a sum:

$$P(M_i|\{\mathbf{x}^{(n)}\}_{n=1}^N) = \frac{P(\{\mathbf{x}^{(n)}\}_{n=1}^N|M_i)P(M_i)}{\sum_{i \in \{1,2,3\}} P(\{\mathbf{x}^{(n)}\}_{n=1}^N|M_i)P(M_i)}$$

Given that  $P(M_i) = \frac{1}{3}$  for all  $i \in \{1, 2, 3\}$ :

$$P(M_i|\{\mathbf{x}^{(n)}\}_{n=1}^N) = \frac{P(\{\mathbf{x}^{(n)}\}_{n=1}^N|M_i)}{\sum_{i \in \{1,2,3\}} P(\{\mathbf{x}^{(n)}\}_{n=1}^N|M_i)}$$

Calculating the posterior probabilities of each of the three models having generated the data in binarydigits.txt using python we can show the values in the table below:

$i$	$P(M_i \{\mathbf{x}^{(n)}\}_{n=1}^N)$
1	1E-1924
2	1E-1858
3	$1-(1\text{E-}1924)-(1\text{E-}1858)$

We can see that for models specified to have the same parameter value for all pixels is very unlikely under the given data set. This makes sense because it is specifying models where the image is essentially blank, which is not reflective of the data. Moreover,  $M_1$  specifies a specific value of 0.5 for all the parameters whereas  $M_2$  specifies any value for all the parameters as long as it's the same. So the model  $M_1$  is a subset of the models specified in  $M_2$  and we can see this reflected in our probabilities when  $P(M_2|\{\mathbf{x}^{(n)}\}_{n=1}^N) > P(M_1|\{\mathbf{x}^{(n)}\}_{n=1}^N)$ .

The Python code for calculating the posterior probabilities of the three models:

```
1 import numpy as np
2 import pandas as pd
3 from scipy.special import betaln, logsumexp
4
5
6 def _log_p_d_given_m1(x):
7     n, d = x.shape
8     return n * d * np.log(0.5)
9
10
11 def _log_p_d_given_m2(x):
12     n, d = x.shape
13     k = np.sum(x, axis=0).astype(int)
14     return betaln(np.sum(k) + 1, n * d - np.sum(k) + 1)
15
16
17 def _log_p_d_given_m3(x):
18     n, _ = x.shape
19     k = np.sum(x, axis=0).astype(int)
20     return logsumexp(betaln(k + 1, n - k + 1))
21
22
23 def c(x, table_path):
24     log_p_d_given_m = np.array(
25         [
26             _log_p_d_given_m1(x),
27             _log_p_d_given_m2(x),
28             _log_p_d_given_m3(x),
29         ]
30     )
31     log_p_m_given_d = log_p_d_given_m - logsumexp(log_p_d_given_m)
32     df = pd.DataFrame(
33         data=np.array(
34             [
35                 np.arange(len(log_p_m_given_d)).astype(int) + 1,
36                 [f"1E{int(x/np.log(10))}" for x in log_p_m_given_d[: -1]]
37                 + [
38                     f"1-{'-'.join([f'(1E{int(x/np.log(10))})' for x in log_p_m_given_d[: -1]])}"
39                 ],
40             ]
41         ).T,
42         columns=["Model", "P(M_i|D)"],
43     )
44     df.set_index("Model", inplace=True)
45     df.to_csv(table_path)
```

src/solutions/q2.py

### Question 3

- (a) The likelihood for a model consisting of a mixture of  $K$  multivariate Bernoulli distributions can be expressed as the product across  $N$  data points:

$$P(\{\mathbf{x}^{(n)}\}_{n=1}^N|\theta) = \prod_{i=1}^N P(x_i|\theta)$$

where  $\{\mathbf{x}^{(n)}\}_{n=1}^N$  is our data set with  $\mathbf{x}^{(n)} \in \mathbb{R}^{D \times 1}$  and  $\theta = \{\pi, \mathbf{P}\}$ ,  $\pi = [\pi_1, \dots, \pi_K] \in \mathbb{R}^{K \times 1}$  our mixing proportions ( $0 \leq \pi_k \leq 1$ ;  $\sum_k \pi_k = 1$ ) and  $\mathbf{P} \in \mathbb{R}^{D \times K}$  the  $K$  Bernoulli parameter vectors with elements  $p_{kd}$  denoting the probability that pixel  $d$  takes value 1 under mixture component  $k$ . We assume the images are iid under the model, and that the pixels are independent of each other within each component distribution.

For each  $P(\mathbf{x}^{(n)}|\theta)$ :

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

The log-likelihood  $\mathcal{L}(\mathbf{x}^{(n)}|\theta)$  can be expressed in matrix form:

$$\mathcal{L}(\mathbf{x}^{(n)}|\theta) = \log \sum_{k=1}^K \pi_k \exp \left( \mathbf{x}^{(n)} \log(\mathbf{P}_k) + (1 - \mathbf{x}^{(n)}) \log(1 - \mathbf{P}_k) \right)$$

which can be further vectorised using Python *scipy*'s *logsumexp* operation.

Moreover, the log-likelihood  $\mathcal{L}(\{\mathbf{x}^{(n)}\}_{n=1}^N|\theta)$  can be expressed:

$$\mathcal{L}(\{\mathbf{x}^{(n)}\}_{n=1}^N|\theta) = \sum_{i=1}^N \left( \log \sum_{k=1}^K \pi_k \exp \left( \mathbf{x}^{(n)} \log(\mathbf{P}_k) + (1 - \mathbf{x}^{(n)}) \log(1 - \mathbf{P}_k) \right) \right)$$

- (b) The expression for the responsibility of mixture component  $k$  for data vector  $x^{(n)}$ , i.e.  $r_{nk} \equiv P(s^{(n)} = k|x^{(n)}, \pi, \mathbf{P})$ . This computation provides the E-step for an EM algorithm.

We know that:

$$P(A|B) \propto P(B|A)P(A)$$

Thus,

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

where  $s^{(n)} \in \{1, \dots, K\}$  a discrete latent variable where  $P(s^{(n)} = k|x^{(n)}|\pi) = \pi_k$ . Note that  $P(s^{(n)} = k|x^{(n)}|\pi) = P(s^{(n)} = k|x^{(n)}|\pi, \mathbf{P})$  as  $s^{(n)}$  isn't dependent on  $\mathbf{P}$ .

Let  $P(s^{(n)} = k|x^{(n)}, \pi, \mathbf{P}) \propto P(s^{(n)})$  be the unnormalised responsibility  $\tilde{r}_{nk}$ . Using the mixture for component  $k$  and the likelihood function of component  $k$ :

$$\tilde{r}_{nk} = \pi_k \prod_{d=1}^D (p_{kd})^{x_d^{(n)}} (1 - p_{kd})^{1 - x_d^{(n)}}$$

Normalising across the components:

$$r_{nk} = \frac{\tilde{r}_{nk}}{\sum_{j=1}^K \tilde{r}_{nj}}$$

and  $r_{nk}$ , we have calculated  $P(s^{(n)} = k | x^{(n)}, \pi, \mathbf{P})$  for the E step.

Moreover,

$$\log \tilde{r}_{nk} = \log \pi_k + \sum_{d=1}^D \left( x_d^{(n)} \log(p_{kd}) + (1 - x_d^{(n)}) \log(1 - \exp(\log(p_{kd}))) \right)$$

and

$$\log r_{nk} = \log \tilde{r}_{nk} - \log \sum_{j=1}^K \exp(\log \tilde{r}_{nj})$$

which can be vectorised as  $\log \mathbf{r}_n$  using Python scipy's *logsumexp* operation.

(c) We know that the expectation log joint can be expressed:

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

Let this quantity be  $E$ . Each term of  $E$  can be expressed:

$$q(s^{(n)}) = \mathbf{r}_n$$

and

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

which is the vectorised version of  $\log \tilde{r}_{nk}$  from part (b) so:

$$\log P(x^{(n)}, s^{(n)} | \pi, \mathbf{P}) = \log(\pi) + \log(\mathbf{P})^T x^{(n)} + \log(1 - \mathbf{P})^T (1 - x^{(n)})$$

Combining:

$$E = \sum_n \mathbf{r}_n^T [\log(\pi) + \log(\mathbf{P})^T x^{(n)} + \log(1 - \mathbf{P})^T (1 - x^{(n)})]$$

To maximise with respect to  $\pi$  and  $\mathbf{P}$  for the M step, we want to take the derivative, set to zero, and solve for  $\hat{\pi}$  and  $\hat{\mathbf{P}}$ .

For the  $k^{th}$  element of  $\pi$ :

$$\frac{\partial E}{\partial \pi_k} = \sum_n r_{nk} \frac{1}{\pi_k}$$

The second derivative:

$$\frac{\partial E}{(\partial \pi_k)^2} = \sum_n r_{nk} \frac{-1}{(\pi_k)^2}$$

is always negative because  $r_{nk} \geq 0$ ,  $\sum_n r_{nk} = 1$ ,  $\pi_k \geq 0$ , and  $\sum_n \pi_k = 1$ , ensuring a maximum in the next step.

We can calculate the maximiser:

$$\frac{\partial E}{\partial \pi_k} + \lambda = 0$$

where  $\lambda$  is a Lagrange multiplier ensuring that the mixing proportions sum to unity.

Thus,

$$\hat{\pi}_k = \frac{\sum_n r_{nk}}{N}$$

For the  $dk^{th}$  element of  $\mathbf{P}$ :

$$\frac{\partial E}{\partial \mathbf{P}_{dk}} = \sum_n r_{nk} \frac{\partial}{\partial \mathbf{P}_{dk}} [x_d^{(n)} \log \mathbf{P}_{dk} + (1 - x_d^{(n)}) \log(1 - \mathbf{P}_{dk})]$$

Simplifying:

$$\frac{\partial E}{\partial \mathbf{P}_{dk}} = \sum_n r_{nk} \left( \frac{x_d^{(n)}}{\mathbf{P}_{dk}} - \frac{1 - x_d^{(n)}}{1 - \mathbf{P}_{dk}} \right)$$

Similar to Question 1, we can see that taking the derivative, the term in the brackets will always be less than zero and with  $r_{nk} \geq 0$  and  $\sum_n r_{nk} = 1$ , the second derivative will always be negative. This ensures that we have a maximum in the next step.

Setting the derivative to zero:

$$\frac{\sum_n x_d^{(n)} r_{nk}}{\mathbf{P}_{dk}} - \frac{\sum_n r_{nk} - \sum_n x_d^{(n)} r_{nk}}{1 - \mathbf{P}_{dk}} = 0$$

Solving for  $\hat{\mathbf{P}}_{dk}$ :

$$\hat{\mathbf{P}}_{dk} \sum_n r_{nk} - \hat{\mathbf{P}}_{dk} \sum_n x_d^{(n)} r_{nk} = \sum_n x_d^{(n)} r_{nk} - \hat{\mathbf{P}}_{dk} \sum_n x_d^{(n)} r_{nk}$$

Thus,

$$\hat{\mathbf{P}}_{dk} = \frac{\sum_n x_d^{(n)} r_{nk}}{\sum_n r_{nk}}$$

We have the maximizing parameters for the expected log-joint

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

thus obtaining an iterative update for the parameters  $\pi$  and  $\mathbf{P}$  in the M-step of EM.

For numerical stability, we can compute the maximisation step for the MAP of  $\mathbf{P}, \hat{\mathbf{P}}_{dk}^{MAP}$ :

$$\frac{\partial E'}{\partial \mathbf{P}_{dk}} = 0$$

where

$$E' = \sum_{n=1}^N q(s^{(n)}) \log P(\mathbf{P} | \pi, x^{(n)}, s^{(n)})$$

and

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

Assuming the same independent Beta prior on each pixel of each component:

$$\log P(\mathbf{P}) = \sum_{k=1}^K \sum_{d=1}^D -\log(B(\alpha, \beta)) + (\alpha - 1) \log \mathbf{P}_{dk} + (\beta - 1) \log(1 - \mathbf{P}_{dk})$$

and

$$\frac{\partial \log P(\mathbf{P})}{\partial \mathbf{P}_{dk}} = \frac{(\alpha - 1)}{\mathbf{P}_{dk}} - \frac{(\beta - 1)}{1 - \mathbf{P}_{dk}}$$

‘ Thus, the derivative can be expressed as:

$$\frac{\partial E'}{\partial \mathbf{P}_{dk}} = \sum_n \left( r_{nk} \left( \frac{\partial \log P(x^{(n)}, s^{(n)}|\pi, \mathbf{P})}{\partial \mathbf{P}_{dk}} + \frac{\partial \log P(\mathbf{P})}{\partial \mathbf{P}_{dk}} \right) \right)$$

Substituting the appropriate expressions:

$$\frac{\partial E'}{\partial \mathbf{P}_{dk}} = \sum_n \left( r_{nk} \left( \frac{x_d^{(n)}}{\mathbf{P}_{dk}} - \frac{1 - x_d^{(n)}}{1 - \mathbf{P}_{dk}} + \frac{(\alpha - 1)}{\mathbf{P}_{dk}} - \frac{(\beta - 1)}{1 - \mathbf{P}_{dk}} \right) \right)$$

Simplifying:

$$\frac{\partial E'}{\partial \mathbf{P}_{dk}} = \frac{(\alpha - 1) + \sum_n r_{nk} x_d^{(n)}}{\mathbf{P}_{dk}} - \frac{(\beta - 1) + \sum_n r_{nk} (1 - x_d^{(n)})}{1 - \mathbf{P}_{dk}}$$

This form is very similar to Question 2 (c). By a similar argument, we can see that the second derivative will always be negative (additionally  $r_{nk} \geq 0$  and  $\sum_n r_{nk} = 1$ ), ensuring a maximum in the calculation of the next step.

Setting  $\frac{\partial E'}{\partial \mathbf{P}_{dk}} = 0$  we can calculate  $\hat{\mathbf{P}}_{dk}^{MAP}$ :

$$\hat{\mathbf{P}}_{dk}^{MAP} = \frac{\alpha - 1 + \sum_n r_{nk} x_d^{(n)}}{(N + \alpha + \beta - 2)}$$

(d) Plotting the posterior likelihood as a function of the iteration number:

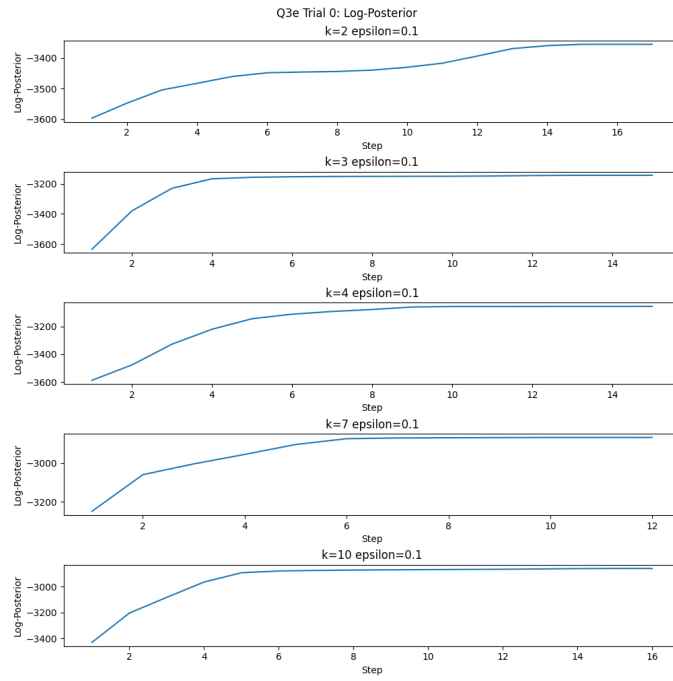


Figure 4: Log Likelihood vs Iteration Number

where *epsilon* is the stopping condition for the posterior posterior converges.



Displaying the parameters found for  $K$  in  $\{2, 3, 4, 7, 10\}$ :

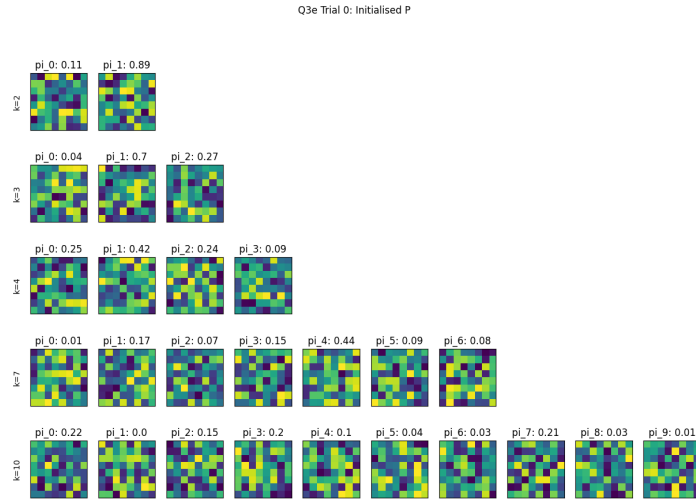


Figure 5: Randomly initialised parameters

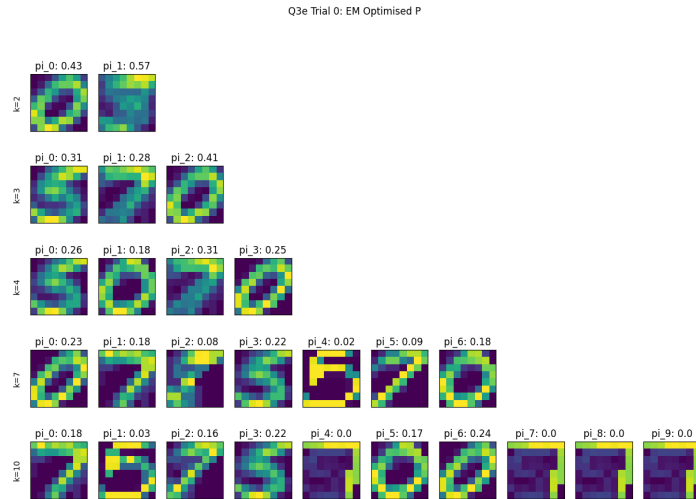


Figure 6: EM optimised parameters

The Python code for the EM algorithm:

```

1 from dataclasses import dataclass
2 from typing import List, Tuple
3
4 import matplotlib.pyplot as plt
5 import numpy as np
6 from scipy.special import logsumexp
7 from sklearn.manifold import TSNE
8
9 from src.constants import DEFAULT_SEED
10
11
12 @dataclass
13 class Theta:
14     """
15     log-pi: the logarithm of the mixing proportions (1, k)
16     log-p-matrix: the logarithm of the probability where the (i,j)th element is the probability that
17                   pixel j takes value 1 under mixture component i (d, k)
18     """
19
20     log_pi: np.ndarray
21     log_p_matrix: np.ndarray
22
23     @property
24     def pi(self):
25         return np.exp(self.log_pi)
26
27     @property
28     def p_matrix(self):
29         d, k = self.log_p_matrix.shape
30         image_dimension = int(np.sqrt(d))
31         return np.exp(self.log_p_matrix).reshape(image_dimension, image_dimension, -1)
32
33     @property
34     def log_one_minus_p_matrix(self) -> np.ndarray:
35         """
36         Compute log(1-P) where P=exp(log-p-matrix)
37         :return: an array of the same shape as log-p-matrix (d, k)
38         """
39         log_of_one = np.zeros(self.log_p_matrix.shape)
40         stacked_sum = np.stack((log_of_one, self.log_p_matrix))
41         weights = np.ones(stacked_sum.shape)
42         weights[1] = -1 # scale p matrix by -1 for subtraction
43         return np.array(logsumexp(stacked_sum, b=weights, axis=0))
44
45     def log_pi_repeated(self, n: int):
46         """
47         Repeats the log_pi vector n times along axis 0
48         :param n: number of repetitions
49         :return: an array of shape (n, k)
50         """
51         return np.repeat(self.log_pi, n, axis=0)
52
53
54 def _init_params(k: int, d: int, seed: int = DEFAULT_SEED) -> Theta:
55     """
56     Random initialisation of theta parameters (log-pi and log-p-matrix)
57     :param k: Number of components
58     :param d: Image dimension (number of pixels in a single image)
59     :param seed: seed initialisation for random methods
60     :return: theta: the parameters of the model
61     """
62     np.random.seed(seed)
63     return Theta(
64         log_pi=np.log(np.random.dirichlet(np.ones(k), size=1)),
65         log_p_matrix=np.log(np.random.uniform(low=0, high=1, size=(d, k))),
66     )
67
68
69 def _compute_log_component_p_x_i_given_theta(x: np.ndarray, theta: Theta) -> np.ndarray:
70     """
71     Compute the unweighted probability of each mixing component for each image
72     :param x: the image data (n, d)
73     :param theta: the parameters of the model
74     :return: an array of the unweighted probabilities (n, k)
75     """
76     return x @ theta.log_p_matrix + (1 - x) @ theta.log_one_minus_p_matrix
77
78
79 def _compute_log_p_x_i_given_theta(x: np.ndarray, theta: Theta) -> np.ndarray:
80     """
81     Computes the log likelihood of each image in the dataset x
82     :param x: the image data (n, d)
83     :param theta: the parameters of the model
84     :return: log-p-x-i-given-theta: a log likelihood array containing the log likelihood of each image (n, 1)
85     """
86     n, _ = x.shape
87     log_component_probabilities = _compute_log_component_p_x_i_given_theta(
88         x, theta
89     ) # (n, k)
90     return np.array(
91         logsumexp(
92             log_component_probabilities
93             + theta.log_pi_repeated(n), # scale each component by component probability

```

```

94         axis=1,
95     )
96 )
97
98
99 def _compute_log_likelihood(x: np.ndarray, theta: Theta) -> float:
100     """
101     Computes the log likelihood of all images in the dataset x
102     :param x: the image data (n, d)
103     :param theta: the parameters of the model
104     :return: log_p_x_given_theta: the log likelihood array across all images
105     """
106     return np.sum(_compute_log_p_x_i_given_theta(x, theta)).item()
107
108
109 def _compute_log_e_step(x: np.ndarray, theta: Theta) -> np.ndarray:
110     """
111     Compute the e step of expectation maximisation
112     :param x: the image data (n, d)
113     :param theta: the parameters of the model
114     :return: an array of the log responsibilities of k mixture components for each image (n, k)
115     """
116     log_r_unnormalised = _compute_log_component_p_x_i_given_theta(x, theta)
117     log_r_normaliser = logsumexp(log_r_unnormalised, axis=1)
118     log_responsibility = log_r_unnormalised - log_r_normaliser[:, np.newaxis]
119     return log_responsibility
120
121
122 def _compute_log_pi_hat(log_responsibility: np.ndarray) -> np.ndarray:
123     """
124     Compute the log of the maximised mixing proportions
125     :param log_responsibility: an array of the log responsibilities of k mixture components for each image
126     (n, k)
127     :return: an array of the maximised log mixing proportions (1, k)
128     """
129     n, _ = log_responsibility.shape
130     return (logsumexp(log_responsibility, axis=0) - np.log(n)).reshape(1, -1)
131
132
133 def _compute_log_p_matrix_hat(
134     x: np.ndarray, log_responsibility: np.ndarray
135 ) -> np.ndarray:
136     """
137     Compute the log of the maximised pixel probabilities
138     :param x: the image data (n, d)
139     :param log_responsibility: an array of the log responsibilities of k mixture components for each image
140     (n, k)
141     :return: an array of the maximised pixel probabilities for each component (d, k)
142     """
143     n, d = x.shape
144     _, k = log_responsibility.shape
145
146     x_repeated = np.repeat(x[:, :, np.newaxis], k, axis=2) # (n, d, k)
147     log_responsibility_repeated = np.repeat(
148         log_responsibility[:, np.newaxis, :], d, axis=1
149     ) # (n, d, k)
150
151     log_p_matrix_unnormalised_likelihood = logsumexp(
152         log_responsibility_repeated, b=x_repeated, axis=0
153     ) # (d, k)
154     log_p_matrix_normaliser_likelihood = np.array(
155         logsumexp(log_responsibility_repeated, axis=0)
156     ) # (d, k)
157
158     alpha = 2
159     beta = 2
160     log_p_matrix_unnormalised_posterior = logsumexp(
161         np.stack(
162             (
163                 (alpha - 1) * np.ones(log_p_matrix_unnormalised_likelihood.shape),
164                 log_p_matrix_unnormalised_likelihood,
165             ),
166             axis=0,
167         ),
168         axis=0,
169     )
170     log_p_matrix_normaliser_posterior = logsumexp(
171         np.stack(
172             (
173                 (alpha + beta - 2) * np.ones(log_p_matrix_normaliser_likelihood.shape),
174                 log_p_matrix_normaliser_likelihood,
175             ),
176             axis=0,
177         ),
178         axis=0,
179     )
180     log_p_matrix_normalised_posterior = log_p_matrix_unnormalised_posterior -
181     log_p_matrix_normaliser_posterior
182     return log_p_matrix_normalised_posterior
183
184
185 def _compute_log_m_step(x: np.ndarray, log_responsibility: np.ndarray) -> Theta:
186     """
187     Compute the m step of expectation maximisation
188     :param x: the image data (n, d)
189     :param log_responsibility: an array of the log responsibilities of k mixture components for each image

```

```

187     (n, k)
188     :return: thetas optimised after maximisation step
189     """
190     return Theta(
191         log_pi=_compute_log_pi_hat(log_responsibility),
192         log_p_matrix=_compute_log_p_matrix_hat(x, log_responsibility),
193     )
194
195 def _run_expectation_maximisation(
196     x: np.ndarray, theta: Theta, max_number_of_steps: int, epsilon: float
197 ) -> Tuple[Theta, np.ndarray, List[float]]:
198     """
199     Run the expectation maximisation algorithm
200     :param x: the image data (n, d)
201     :param theta: initial theta parameters
202     :param max_number_of_steps: the maximum number of steps to run the algorithm
203     :param epsilon: the minimum required change in log likelihood, otherwise the algorithm stops early
204     :return: a tuple containing the optimised thetas, the log responsibilities,
205             and the log likelihood at each step of the algorithm
206     """
207     log_responsibility = None
208     log_likelihooods = []
209     for _ in range(max_number_of_steps):
210         log_responsibility = _compute_log_e_step(x, theta)
211         theta = _compute_log_m_step(x, log_responsibility)
212
213         log_likelihooods.append(_compute_log_likelihood(x, theta))
214
215         # check for early stopping
216         if len(log_likelihooods) > 1:
217             if (log_likelihooods[-1] - log_likelihooods[-2]) < epsilon:
218                 break
219     return theta, log_responsibility, log_likelihooods
220
221
222 def _plot_p_matrix(
223     thetas: List[Theta], ks: List[int], figure_title: str, figure_path: str
224 ):
225     n = len(ks)
226     m = np.max(ks)
227     fig = plt.figure()
228     fig.set_figwidth(15)
229     fig.set_figheight(10)
230     for i, k in enumerate(ks):
231         for j in range(k):
232             ax = plt.subplot(n, m, m * i + j + 1)
233             ax.imshow(
234                 thetas[i].p_matrix[:, :, j],
235                 interpolation="None",
236             )
237             ax.tick_params(
238                 axis="x",
239                 which="both",
240                 bottom=False,
241                 top=False,
242             )
243             ax.tick_params(
244                 axis="y",
245                 which="both",
246                 left=False,
247                 right=False,
248             )
249             ax.xaxis.set_ticklabels([])
250             ax.yaxis.set_ticklabels([])
251             ax.set_title(f"pi-{j}: {np.round(thetas[i].pi[0, j], 2)}")
252             if j == 0:
253                 ax.set_ylabel(f"{k}")
254     fig.suptitle(figure_title)
255     plt.savefig(figure_path)
256
257
258 def _plot_tsne_responsibility_clusters(
259     log_responsibilities: List[np.ndarray], ks: List[int], figure_title: str, figure_path: str
260 ):
261     n = len(ks)
262     fig = plt.figure()
263     fig.set_figwidth(5*n)
264     fig.set_figheight(5)
265     for i, k in enumerate(ks):
266         embedding = TSNE(n_components=2, learning_rate='auto', init='random', perplexity=10).fit_transform(
267             log_responsibilities[i])
268         ax = plt.subplot(1, n, i+1)
269         ax.scatter(embedding[:, 0], embedding[:, 1])
270         ax.set_title(f"{k}")
271     fig.suptitle(figure_title)
272     plt.savefig(figure_path, bbox_inches='tight')
273
274
275 def _plot_log_posteriors(
276     log_posteriors: List[List[float]],
277     ks: List[int],
278     epsilon: float,
279     figure_title: str,
280     figure_path: str,
281 ) -> None:

```

```

282 fig, ax = plt.subplots(len(ks), 1, constrained_layout=True)
283 fig.set_figwidth(10)
284 fig.set_figheight(10)
285 for i, k in enumerate(ks):
286     ax[i].plot(np.arange(1, len(log-posteriors[i]) + 1), log-posteriors[i])
287     ax[i].set_xlabel("Step")
288     ax[i].set_ylabel(f"Log-Posterior")
289     ax[i].set_title(f"{k=} {epsilon=}")
290 plt.suptitle(figure-title)
291
292 plt.savefig(figure-path)
293
294
295 def e(
296     x: np.ndarray,
297     number_of_trials: int,
298     ks: List[int],
299     epsilon: float,
300     max_number_of_steps: int,
301     figure_path: str,
302     figure_title: str,
303 ) -> None:
304     n, d = x.shape
305     seeds = np.random.randint(
306         low=number_of_trials * len(ks), size=(number_of_trials, len(ks))
307     )
308     for i in range(number_of_trials):
309         init_thetas = []
310         em_thetas = []
311         log-posteriors = []
312         log-responsibilities = []
313         for j, k in enumerate(ks):
314             init_theta = _init_params(k, d, seed=seeds[i, j])
315             em_theta, log-responsibility, log-posterior = _run_expectation_maximisation(
316                 x,
317                 theta=init_theta,
318                 epsilon=epsilon,
319                 max_number_of_steps=max_number_of_steps,
320             )
321             init_thetas.append(init_theta)
322             em_thetas.append(em_theta)
323             log-responsibilities.append(log-responsibility)
324             log-posteriors.append(log-posterior)
325
326         _plot_p_matrix(
327             init_thetas,
328             ks,
329             figure_title=f"{figure_title} Trial {i}: Initialised P",
330             figure_path=f"{figure_path}-{i}-initialised-p.png",
331         )
332         _plot_p_matrix(
333             em_thetas,
334             ks,
335             figure_title=f"{figure_title} Trial {i}: EM Optimised P",
336             figure_path=f"{figure_path}-{i}-optimised-p.png",
337         )
338         _plot_tsne_responsibility_clusters(
339             log-responsibilities,
340             ks,
341             figure_title=f"{figure_title} Trial {i}: TSNE Responsibility Visualisation",
342             figure_path=f"{figure_path}-{i}-tsne.png",
343         )
344         _plot_log-posteriors(
345             log-posteriors,
346             ks,
347             epsilon,
348             figure_title=f"{figure_title} Trial {i}: Log-Posterior",
349             figure_path=f"{figure_path}-{i}-log-pos.png",
350         )

```

src/solutions/q3.py

- (e) Running the algorithm a few times starting from randomly chosen initial conditions and visualising the parameters:

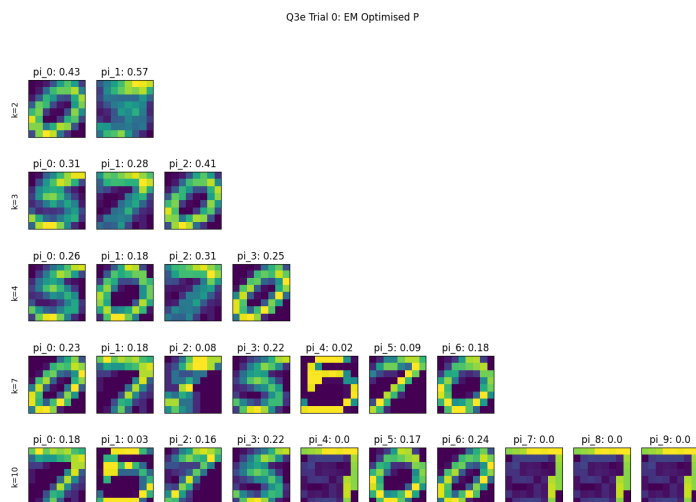


Figure 7: EM optimised parameters: Trial 0

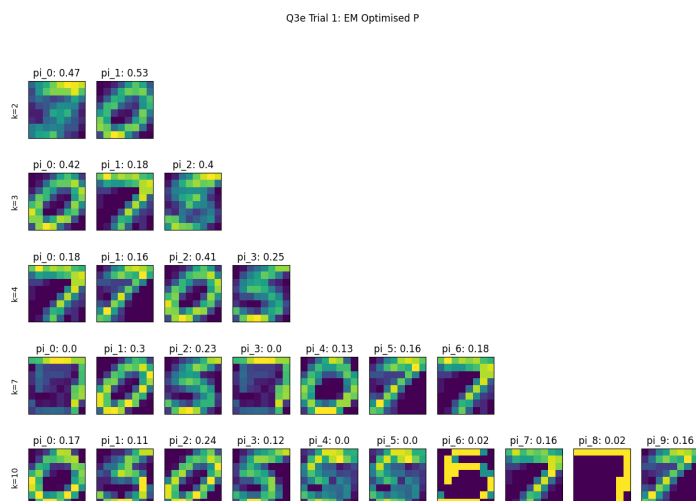


Figure 8: EM optimised parameters: Trial 1

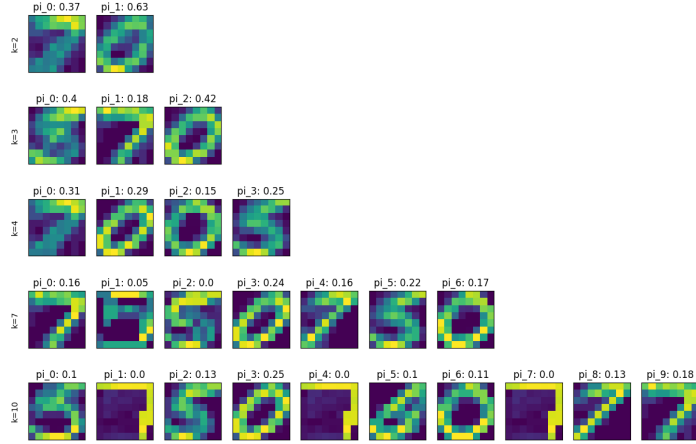


Figure 9: EM optimised parameters: Trial 2

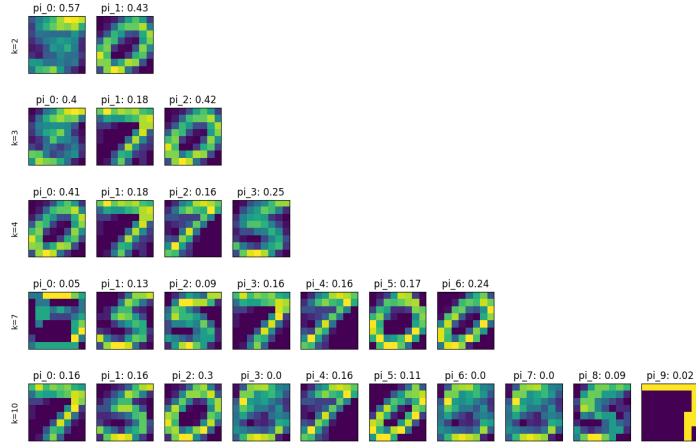


Figure 10: EM optimised parameters: Trial 3

For smaller  $k$ , we can visually see that we obtain very similar solutions (a 7 and a 0 for  $k = 2$ ). However for higher  $K$ , we see that this may not always be the case. For Trial 0 of  $k = 10$ , we have two 0's whereas in Trial 1 we have three 0's. Interestingly, the zeros are all different, representing different variants of the written digit (i.e. a slanted zero, a slightly slanted zero, and a symmetric zero).

Moreover, looking at the responsibilities of each mixture component, we can see that when  $k$  is relatively small they are relatively evenly distributed. However for  $k = 7$  and especially  $k = 10$ , we can see some components have very small or zero probability (i.e.  $\pi_1$  and  $\pi_2$  of trial 2). The corresponding parameter visualisations will essentially never be utilised because no probability is assigned to the component. This can be verified when we perform a TSNE visualisation of the responsibility vector for each of the images. We can see that for large  $k$ , qualitatively the number of clusters no longer matches the  $k$  value, indicating that some clusters are redundant. For example for  $k = 7$  and  $k = 10$  we can only qualitatively see four or five clusters with TSNE.

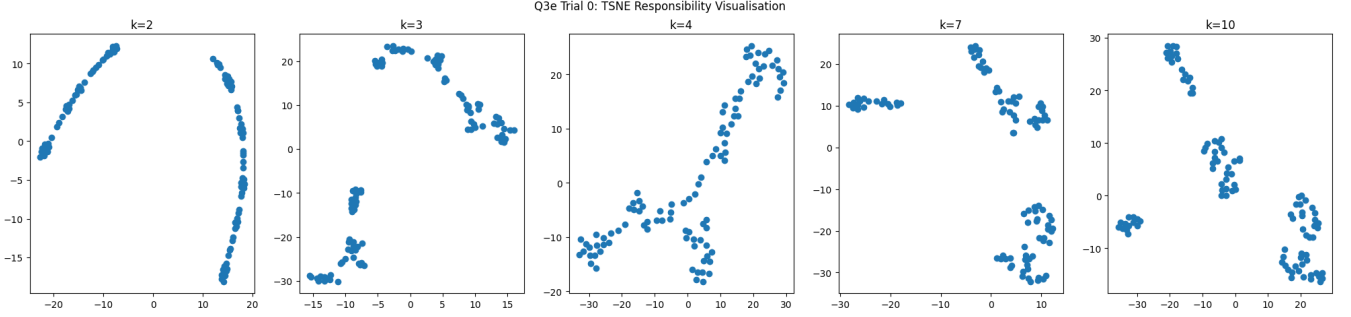


Figure 11: TSNE Visualisation of Image responsibilities: Trial 0

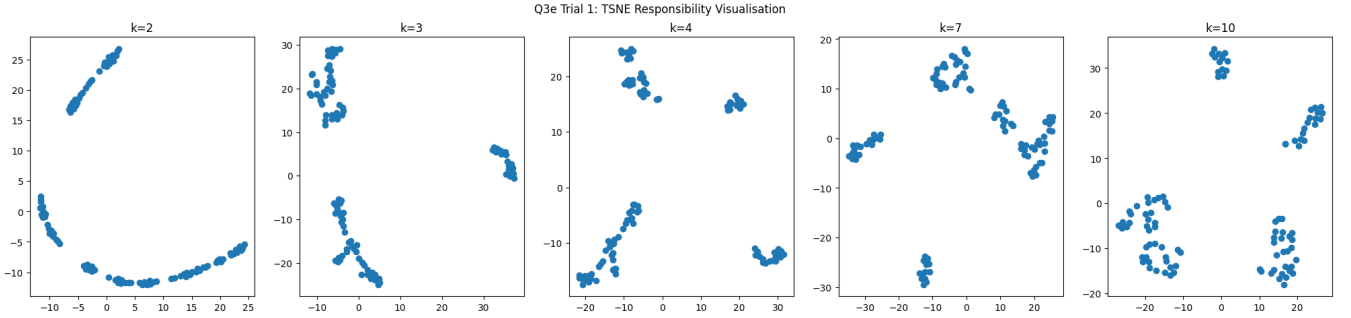


Figure 12: TSNE Visualisation of Image responsibilities: Trial 1

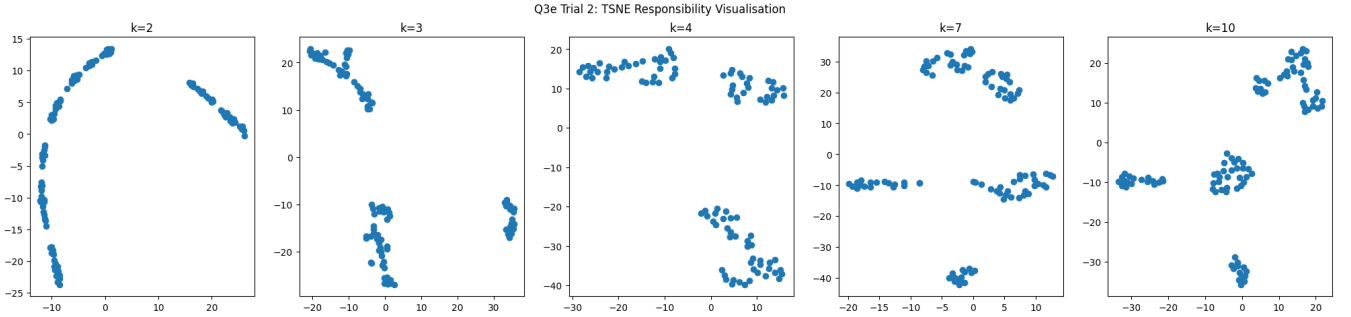


Figure 13: TSNE Visualisation of Image responsibilities: Trial 2



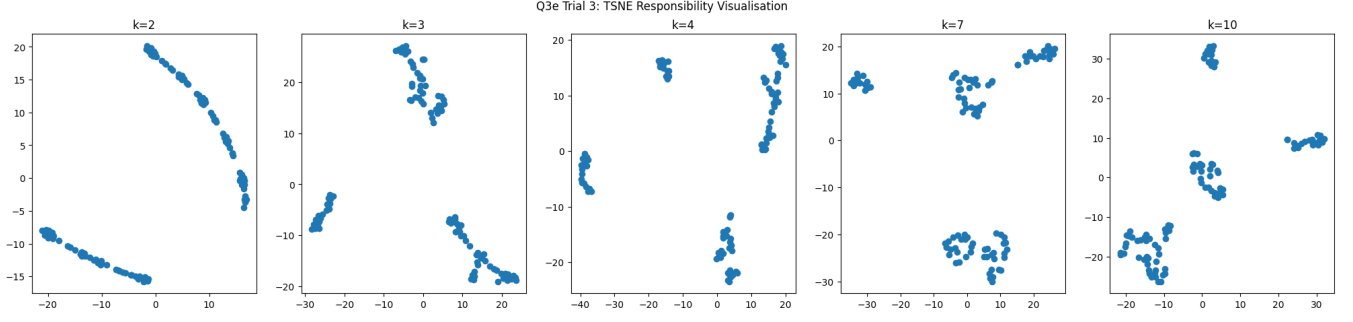


Figure 14: TSNE Visualisation of Image responsibilities: Trial 3

Improvements to the model could include searching for an optimal  $k$  by maximising the log posterior with regularisation on the magnitude of  $k$  to balancing maximising log posterior with minimising model complexity. Additionally, adding a prior on the responsibility components can be helpful to ensure non-zero mixing components unlike the components visualised here. This can help promote more meaningful clusters as  $k$  increases.

## Question 5

- (a) Let  $p(\mathbf{s}_i, \mathbf{s}_{i-1})$  be the probability of the pair of symbols  $\mathbf{s}_i$  and  $\mathbf{s}_{i-1}$  occurring together in the text ( $\mathbf{s}_{i-1}$  followed by  $\mathbf{s}_i$  where order matters). We can model  $p(\mathbf{s}_i, \mathbf{s}_{i-1})$  as a multinomial distribution with  $N = 1$  and  $D = 53^2$ :

$$p(\mathbf{s}_i, \mathbf{s}_{i-1}) = \frac{1}{s^1! s^2! \dots s^{53^2}} \prod_{j \in \{1, \dots, 53\}, k \in \{1, \dots, 53\}} p_{s_i, s_{i-1}}^{t^{s_i, s_{i-1}}}$$

where  $t^{s_i, s_{i-1}}$  is an indicator of transition  $s_{i-1}$  to  $s_i$ . For convenience we will denote  $t^{(\alpha, \beta)}$  as the transition from A multinomial distribution is appropriate because there can only be only one of 53 symbols chosen as  $s$  (i.e. a 53 dimensional dice). Thus,  $p(s = s_i)$  represents the probability of the symbol  $s_i$  in the text.

We can convert  $p(\mathbf{s})$  into exponential family form:

$$p(\mathbf{s}) = \frac{1}{s^1! s^2! \dots s^D} \exp(\mathbf{s}^T \log(\mathbf{p}))$$

where  $\mathbf{p}$  is the vector of  $p_i$ 's. Thus the sufficient statistic is  $T(\mathbf{s}) = \mathbf{s}^T$ .

The ML estimate

- (b) The latent variables  $\sigma(s)$  for different symbols  $s$  are not independent. This is because by choosing an encoding for one symbol  $e = \sigma(s)$ , the encoding for a second symbol  $\sigma(s')$  cannot be  $e$ . We have 53 symbols but only 52 degrees of freedom, because once we have defined the encoding for 52 symbols, the encoding for the 53<sup>rd</sup> symbol cannot be chosen. Thus, there exists a dependence between the symbols for a given  $\sigma$ .

The joint probability of the encrypted text  $e_1 e_2 \dots e_n$  given  $\sigma$ :

$$P(e_1, e_2, \dots, e_n | \sigma) = \psi(\gamma = \sigma^{-1}(e_1)) \prod_{i=2}^n \psi(\alpha = \sigma^{-1}(e_i), \beta = \sigma^{-1}(e_{i-1}))$$

because  $\sigma$  is the encoding function, mapping to a symbol  $s$  into the encoded text as  $e$ , we require  $\sigma^{-1}$  the decoding function mapping the encoded symbol  $e$  back to  $s$ .

- (c) The proposal probability  $S(\sigma \rightarrow \sigma')$  depends on the permutations of  $\sigma$  and  $\sigma'$  because we only choose a proposal  $\sigma'$  that differs at *two* spots:

$$\begin{aligned} \sigma'(s^i) &= \sigma(s^j) \\ \sigma'(s^j) &= \sigma(s^i) \end{aligned}$$

for any two symbols  $s^i$  and  $s^j$  of the 53 possible symbols ( $s^i \neq s^j$ ).

Therefore, if the above doesn't hold for  $\sigma'$ ,  $S(\sigma \rightarrow \sigma') = 0$ , because with our method of choosing proposals, it is not possible to choose  $\sigma'$ . At  $\sigma$  there are  $\binom{53}{2}$  possible proposal  $\sigma'$ 's with the above property. Because we are assuming a uniform prior distribution over  $\sigma$ 's, the transition probability of a  $\sigma'$  that satisfies the above property is  $S(\sigma \rightarrow \sigma') = \frac{1}{\binom{53}{2}}$ .

The MH acceptance probability is given as:

$$A(\sigma \rightarrow \sigma'|\mathcal{D}) = \min\{1, \frac{S(\sigma' \rightarrow \sigma)P(\sigma'|\mathcal{D})}{S(\sigma \rightarrow \sigma')P(\sigma|\mathcal{D})}\}$$

where  $S(\sigma \rightarrow \sigma')$  is the conditional transition probability of  $\sigma'$  given  $\sigma$  and  $\mathcal{D}$  is our encrypted text  $e_1, e_2, \dots, e_n$ .

$S(\sigma \rightarrow \sigma') = S(\sigma' \rightarrow \sigma)$  for all  $\sigma$  and  $\sigma'$  that differ only at two spots because the probability in this case will always be  $\frac{1}{\binom{53}{2}}$ , we can simplify:

$$A(\sigma \rightarrow \sigma'|\mathcal{D}) = \min\{1, \frac{P(\sigma'|\mathcal{D})}{P(\sigma|\mathcal{D})}\}$$

From Bayes' Theorem:

$$P(\sigma|\mathcal{D}) = \frac{P(\mathcal{D}|\sigma)P(\sigma)}{\sum_{\sigma'} P(\mathcal{D}|\sigma')P(\sigma')}$$

We are assuming a uniform prior for  $\sigma$ , so  $P(\sigma)$  is a constant and we can simplify further:

$$A(\sigma \rightarrow \sigma'|\mathcal{D}) = \min\{1, \frac{P(\mathcal{D}|\sigma')}{P(\mathcal{D}|\sigma)}\}$$

This is the acceptance probability for a given proposal  $\sigma'$ . The expression for  $P(\mathcal{D}|\sigma)$  is  $P(e_1, e_2, \dots, e_n|\sigma)$  described in the previous part.

(d) The Python code for the MH sampler:

```
src/solutions/q5.py
```

Implement the MH sampler, and run it on the provided encrypted text. Report the current decryption of the first 60 symbols after every 100 iterations. Your Markov chain should converge to give you a fairly sensible message. (Hint: it may help to initialize your chain intelligently and to try multiple times; in any case, please describe what you did). [30 marks]

TODO

- (e) Note that some  $\Psi(\alpha, \beta)$  values may be zero. Does this affect the ergodicity of the chain? If the chain remains ergodic, give a proof; if not, explain and describe how you can restore ergodicity. [5 marks]

TODO

- (f) Analyse this approach to decoding. For instance, would symbol probabilities alone (rather than transitions) be sufficient? If we used a second order Markov chain for English text, what problems might we encounter? Will it work if the encryption scheme allows two symbols to be mapped to the same encrypted value? Would it work for Chinese with  $> 10000$  symbols? [13 marks]

TODO

## Question 7

- (a) To find the local extrema of the function  $f(x, y) = x + 2y$  subject to the constraint  $y^2 + xy = 1$ , first we define  $g(x, y)$ :

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

where  $g(x, y) = 0$  is an equivalent representation of the given constraint.

We can therefore construct the optimisation problem:

$$\min_{\mathbf{x}} f(\mathbf{x})$$

such that  $g(\mathbf{x}) = 0$  and  $\mathbf{x} := [x, y]^T$ .

We can calculate  $\nabla f(\mathbf{x})$ :

$$\begin{aligned}\nabla f(\mathbf{x}) &= \left[ \frac{\partial}{\partial x}(x + 2y), \frac{\partial}{\partial y}(x + 2y) \right]^T \\ \nabla f(\mathbf{x}) &= [1, 2]^T\end{aligned}$$

and calculating  $\nabla g(\mathbf{x})$ :

$$\begin{aligned}\nabla g(\mathbf{x}) &= \left[ \frac{\partial}{\partial x}(y^2 + xy - 1), \frac{\partial}{\partial y}(y^2 + xy - 1) \right]^T \\ \nabla g(\mathbf{x}) &= [y, 2y + x]^T\end{aligned}$$

Solving the constraint optimisation problem with Lagrange multipliers, we set up the equations:

$$\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = \mathbf{0}$$

and

$$g(\mathbf{x}) = 0$$

Giving us the three equations:

$$\begin{aligned}1 + \lambda y &= 0 \\ 2 + \lambda(2y + x) &= 0 \\ y^2 + xy - 1 &= 0\end{aligned}$$

Substituting  $\lambda y = -1$  from the first equation into the second equation:

$$2 + 2(-1) + x = 0$$

We see that  $x = 0$ . Solving for  $y$  in our third equation with  $x = 0$ :

$$y^2 - 1 = 0$$

We see that  $y = \pm 1$  and from the first equation  $\lambda \mp 1$ .

The local extrema are  $(x = 0, y = 1)$  when our  $\lambda = -1$  and  $(x = 0, y = -1)$  when our  $\lambda = 1$ .

(b)

(i) Given that  $g(a) = \ln(a)$ , we want to transform this to the form  $f(x, a) = 0$ :

$$x = \ln(a)$$

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

Thus,

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

(ii) We know that for Newton's method's

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

where  $f(x_n) = f(x_n, a) = \exp(x_n) - a$

We can calculate:

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

Assuming we can evaluate  $\exp(x)$ , our update equation:

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

Simplifying:

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

# Appendix: main.py

```
1 import os
2
3 import numpy as np
4
5 from src.constants import BINARY_DIGITS_FILE_PATH, OUTPUTS_FOLDER
6 from src.solutions import q1, q2, q3
7
8 if __name__ == "__main__":
9
10     if not os.path.exists(OUTPUTS_FOLDER):
11         os.makedirs(OUTPUTS_FOLDER)
12
13     x = np.loadtxt(BINARY_DIGITS_FILE_PATH)
14     # Question 1
15     Q1_OUTPUT_FOLDER = os.path.join(OUTPUTS_FOLDER, "q1")
16     if not os.path.exists(Q1_OUTPUT_FOLDER):
17         os.makedirs(Q1_OUTPUT_FOLDER)
18
19     q1.d(
20         x,
21         figure_path=os.path.join(Q1_OUTPUT_FOLDER, "q1d.png"),
22         figure_title="Q1d: Maximum Likelihood Estimate",
23     )
24     q1.e(
25         x,
26         alpha=3,
27         beta=3,
28         figure_path=os.path.join(Q1_OUTPUT_FOLDER, "q1e"),
29         figure_title="Q1e: Maximum A Prior",
30     )
31
32     # Question 2
33     Q2_OUTPUT_FOLDER = os.path.join(OUTPUTS_FOLDER, "q2")
34     if not os.path.exists(Q2_OUTPUT_FOLDER):
35         os.makedirs(Q2_OUTPUT_FOLDER)
36     q2.c(x, table_path=os.path.join(Q2_OUTPUT_FOLDER, "q2c.csv"))
37
38     # Question 3
39     Q3_OUTPUT_FOLDER = os.path.join(OUTPUTS_FOLDER, "q3")
40     if not os.path.exists(Q3_OUTPUT_FOLDER):
41         os.makedirs(Q3_OUTPUT_FOLDER)
42     q3.e(
43         x,
44         number_of_trials=4,
45         ks=[2, 3, 4, 7, 10],
46         epsilon=1e-1,
47         max_number_of_steps=int(1e2),
48         figure_path=os.path.join(Q3_OUTPUT_FOLDER, "q3e"),
49         figure_title="Q3e",
50     )
```

main.py