

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 :

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

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

Simplifying:

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

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 $\frac{\partial \mathcal{L}(\mathbf{p}|\mathcal{D})}{\partial p_d} = 0$:

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

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

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:

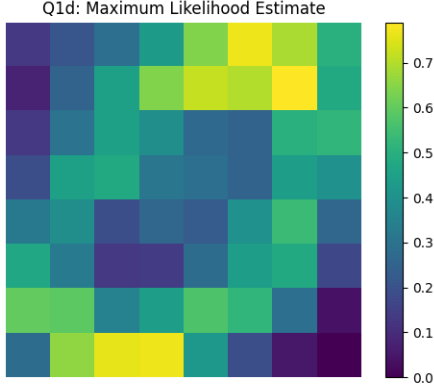


Figure 1: ML parameters

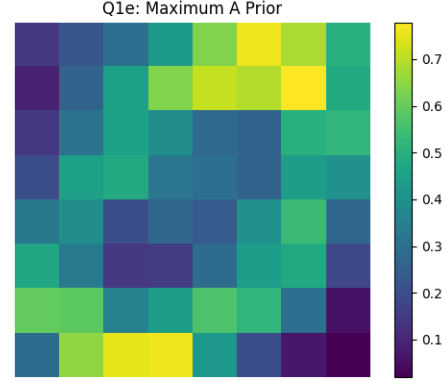


Figure 2: MAP parameters

Comparing the equations:

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

and

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

We see that when the number of data points increases, the MAP estimate approaches $\frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)}$, the MLE. This result makes sense because as our data set gets bigger, are less reliant on our prior. Moreover, if in our data set of images, a specific pixel in all of these images are white or all black, the MLE for that pixel will be binary. This may not be representative of our intuitions about images, where there should be some non-zero probability of a pixel being black or white. Thus, introducing an appropriate prior will ensure that the probability of that pixel will never be zero or one. Moreover, a prior can also help ensure numerical stability during calculations. The logarithm of zero is negative infinity, so having probability zero can be problematic for log-likelihoods calculations. A prior can ensure non-zero probabilities. In our case, by using a small Beta distribution prior on each pixel, our parameter values are biased be close to 0.5 and to never be at the extremities 0 and 1. We can see this in our MAP parameters figure where the range of our parameters is smaller than the range for the ML parameters. Also, the MLE has 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 misspecified prior can be problematic, as the estimated parameters might be skewed by the prior and not properly represent the underlying data generating process and this can be worse than using the MLE.

We can also visualise $\hat{\mathbf{p}}^{MAP} - \hat{\mathbf{p}}^{MLE}$ to 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.

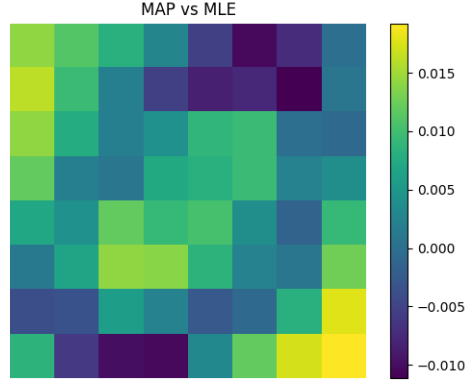


Figure 3: $\hat{\mathbf{p}}^{MAP} - \hat{\mathbf{p}}^{MLE}$

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 :

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

- (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 :

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

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

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

For each model M_i , we can marginalise out P 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:

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

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:

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

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 ave three models, so in this case the normalisation $P(D)$ 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 | $P(M_i \mathcal{D})$ |
|-----|------------------------|
| 1 | 1E-1924 |
| 2 | 1E-1858 |
| 3 | 1-(1E-1924)-(1E-1858) |

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 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 .

Filling in the appropriate quantities:

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

where $q(s^{(n)}) = r_n^T$ and $\log P(x^{(n)}, s^{(n)}) = \log(\pi) + \log(\mathbf{P})^T x^{(n)} + \log(1 - \mathbf{P})^T (1 - x^{(n)})$. Noting that $P(x^{(n)}, s^{(n)} | \pi, \mathbf{P}) = P(x^{(n)} | s^{(n)}) \pi, \mathbf{P}) P(s^{(n)} | \pi, \mathbf{P}) = \tilde{r}_n$

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

For the k^{th} element of π :

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

and

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

where λ is a Lagrange multiplier ensuring that the mixing proportions sum to unity.

Thus,

$$\hat{\pi}_m = \frac{1}{N} \sum_n r_{nk}$$

For the dk^{th} element of P:

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

Simplifying:

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

Setting the derivative to zero:

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

Solving for P_{dk} :

$$P_{dk} \sum_n r_{nk} - P_{dk} \sum_n x_d^{(n)} r_{nk} = \sum_n x_d^{(n)} r_{nk} - P_{dk} \sum_n x_d^{(n)} r_{nk}$$

Thus,

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

We have the maximizing parameters for the expected log-joint

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

thus obtaining an iterative update for the parameters π and P in the M-step of EM.

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

```

```

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

```

```

184 """
185 Run the expectation maximisation algorithm
186 :param x: the image data (n, d)
187 :param theta: initial theta parameters
188 :param max_number_of_steps: the maximum number of steps to run the algorithm
189 :param epsilon: the minimum required change in log likelihood, otherwise the algorithm stops early
190 :return: a tuple containing the optimised thetas and the log likelihood at each step of the algorithm
191 """
192 log_likelihoods = []
193 for _ in range(max_number_of_steps):
194     log_r = _compute_log_e_step(x, theta)
195     theta = _compute_log_m_step(x, log_r)
196
197     log_likelihoods.append(_compute_log_likelihood(x, theta))
198
199     # check for early stopping
200     if len(log_likelihoods) > 1:
201         if (log_likelihoods[-1] - log_likelihoods[-2]) < epsilon:
202             break
203 return theta, log_likelihoods
204
205
206 def _plot_p_matrix(
207     thetas: List[Theta], ks: List[int], figure_title: str, figure_path: str
208 ):
209     n = len(ks)
210     m = np.max(ks)
211     fig = plt.figure()
212     for i, k in enumerate(ks):
213         for j in range(k):
214             ax = plt.subplot(n, m, m * i + j + 1)
215             ax.imshow(
216                 thetas[i].p_matrix[:, :, j],
217                 interpolation="None",
218             )
219             ax.tick_params(
220                 axis="x",
221                 which="both",
222                 bottom=False,
223                 top=False,
224             )
225             ax.tick_params(
226                 axis="y",
227                 which="both",
228                 left=False,
229                 right=False,
230             )
231             ax.xaxis.set_ticklabels([])
232             ax.yaxis.set_ticklabels([])
233             if j == 0:
234                 ax.set_ylabel(f"{k}")
235     fig.suptitle(figure_title)
236     plt.savefig(figure_path)
237
238
239 def _plot_log_likelihoods(
240     log_likelihoods: List[List[float]],
241     ks: List[int],
242     figure_title: str,
243     figure_path: str,
244 ) -> None:
245     fig, ax = plt.subplots(len(ks), 1, constrained_layout=True)
246     fig.set_figwidth(10)
247     fig.set_figheight(15)
248     for i, k in enumerate(ks):
249         ax[i].plot(np.arange(1, len(log_likelihoods[i]) + 1), log_likelihoods[i])
250         ax[i].set_xlabel("Step")
251         ax[i].set_ylabel("Log-Likelihood")
252         ax[i].set_title(f"{k}")
253     plt.suptitle(figure_title)
254
255     plt.savefig(figure_path)
256
257
258 def e(
259     x: np.ndarray,
260     number_of_trials: int,
261     ks: List[int],
262     epsilon: float,
263     max_number_of_steps: int,
264     figure_path: str,
265     figure_title: str,
266 ) -> None:
267     n, d = x.shape
268     seeds = np.random.randint(
269         low=number_of_trials * len(ks), size=(number_of_trials, len(ks))
270     )
271     for i in range(number_of_trials):
272         init_thetas = []
273         em_thetas = []
274         log_likelihoods = []
275         for j, k in enumerate(ks):
276             init_theta = _init_params(k, d, seed=seeds[i, j])
277             em_theta, log_likelihood = _run_expectation_maximisation(
278                 x,
279                 theta=init_theta,

```



```

280         epsilon=epsilon ,
281         max_number_of_steps=max_number_of_steps ,
282     )
283     init_thetas.append(init_theta)
284     em_thetas.append(em_theta)
285     log_likelihoods.append(log_likelihood)
286
287     _plot_p_matrix(
288         init_thetas ,
289         ks,
290         figure_title=f"{figure_title} Trial {i}: Initialised P",
291         figure_path=f"{figure_path}-{i}-initialised-p.png",
292     )
293     _plot_p_matrix(
294         em_thetas ,
295         ks,
296         figure_title=f"{figure_title} Trial {i}: EM Optimised P",
297         figure_path=f"{figure_path}-{i}-optimised-p.png",
298     )
299     _plot_log_likelihoods(
300         log_likelihoods ,
301         ks,
302         figure_title=f"{figure_title} Trial {i}: Negative Log-Likelihood",
303         figure_path=f"{figure_path}-{i}-log-like.png",
304     )

```

src/solutions/q3.py

Run your algorithm on the data set for values of K in 2, 3, 4, 7, 10. Plot the log likelihood as a function of the iteration number, and display the parameters found. [30 marks]

- (e) Run the algorithm a few times starting from randomly chosen initial conditions. Do you obtain the same solutions (up to permutation)? Does this depend on K ? Show the learned probability vectors as images.

Comment on how well the algorithm works, whether it finds good clusters (look at the cluster means and responsibilities and try to interpret them), and how you might improve the model. [10 marks]

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 53rd symbol cannot be chosen. Thus, there exists a dependence between the symbols for a given σ .

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

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

because σ is the encoding function, mapping to a symbol s into the encoded text as e , we require σ^{-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 σ and σ' because we only choose a proposal σ' 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 σ' , $S(\sigma \rightarrow \sigma') = 0$, because with our method of choosing proposals, it is not possible to choose σ' . At σ there are $\binom{53}{2}$ possible proposal σ' 's with the above property. Because we are assuming a uniform prior distribution over σ 's, the transition probability of a σ' 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 σ' given σ and \mathcal{D} is our encrypted text e_1, e_2, \dots, e_n .

$S(\sigma \rightarrow \sigma') = S(\sigma' \rightarrow \sigma)$ for all σ and σ' 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 σ , 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 σ' . 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(x, figure_path=os.path.join(Q1_OUTPUT_FOLDER, "q1d.png"), figure_title="Q1d: Maximum Likelihood Estimate")
20     q1.e(x, alpha=3, beta=3, figure_path=os.path.join(Q1_OUTPUT_FOLDER, "q1e"), figure_title="Q1e: Maximum A Prior")
21
22     # Question 2
23     Q2_OUTPUT_FOLDER = os.path.join(OUTPUTS_FOLDER, "q2")
24     if not os.path.exists(Q2_OUTPUT_FOLDER):
25         os.makedirs(Q2_OUTPUT_FOLDER)
26     q2.c(x, table_path=os.path.join(Q2_OUTPUT_FOLDER, "q2c.csv"))
27
28     # Question 3
29     Q3_OUTPUT_FOLDER = os.path.join(OUTPUTS_FOLDER, "q3")
30     if not os.path.exists(Q3_OUTPUT_FOLDER):
31         os.makedirs(Q3_OUTPUT_FOLDER)
32     q3.e(
33         x,
34         number_of_trials=3,
35         ks=[2, 3, 4, 7, 10],
36         epsilon=1e-1,
37         max_number_of_steps=int(1e2),
38         figure_path=os.path.join(Q3_OUTPUT_FOLDER, "q3e"),
39         figure_title="Q3e",
40     )
```

main.py