

EM for Binary Data.

Consider the data set of binary (black and white) images used in the previous assignment. Each image is arranged into a vector of pixels by concatenating the columns of pixels in the image. The data set has N images $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$ and each image has D pixels, where D is (number of rows \times number of columns) in the image. For example, image $\mathbf{x}^{(n)}$ is a vector $(x_1^{(n)}, \dots, x_D^{(n)})$ where $x_d^{(n)} \in \{0, 1\}$ for all $n \in \{1, \dots, N\}$ and $d \in \{1, \dots, D\}$.

- (a) Write down the likelihood for a model consisting of a mixture of K multivariate Bernoulli distributions. Use the parameters π_1, \dots, π_K to denote the mixing proportions ($0 \leq \pi_k \leq 1; \sum_k \pi_k = 1$) and arrange the K Bernoulli parameter vectors into a matrix P with elements p_{kd} denoting the probability that pixel d takes value 1 under mixture component k . Assume the images are iid under the model, and that the pixels are independent of each other within each component distribution.

Just as we can for a mixture of Gaussians, we can formulate this mixture as a latent variable model, introducing a discrete hidden variable $s^{(n)} \in \{1, \dots, K\}$ where $P(s^{(n)} = k | \boldsymbol{\pi}) = \pi_k$.

- (b) Write down the expression for the responsibility of mixture component k for data vector $\mathbf{x}^{(n)}$, i.e. $r_{nk} \equiv P(s^{(n)} = k | \mathbf{x}^{(n)}, \boldsymbol{\pi}, \mathsf{P})$. This computation provides the E-step for an EM algorithm.
- (c) Find the maximizing parameters for the expected log-joint

$$\underset{\boldsymbol{\pi}, \mathsf{P}}{\operatorname{argmax}} \left\langle \sum_n \log P(\mathbf{x}^{(n)}, s^{(n)} | \boldsymbol{\pi}, \mathsf{P}) \right\rangle_{q(\{s^{(n)}\})}$$

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

- (d) Implement the EM algorithm for a mixture of K multivariate Bernoullis.

Your code should take as input the number K , a matrix X containing the data set, and a maximum number of iterations to run. The algorithm should run for that number of iterations or until the log likelihood converges (does not increase by more than a very small amount).

Hand in clearly commented code and a description of how the code works.

Run your algorithm on the data set for values of K in $\{2, 3, 4, 7, 10\}$. Report the log likelihoods obtained and display the parameters found.

[Hints: Although the implementation may seem simple enough given the equations, there are many numerical pitfalls (that are common in much of probabilistic learning). A few suggestions:

- Likelihoods can be very small; it is often better to work with log-likelihoods.

- You may still encounter numerical issues computing responsibilities. Consider scaling the numerator and denominator of the equation by a suitable constant while still in the log domain.
 - It may also help to introduce (weak) priors on P and π and use EM to find the MAP estimates, rather than ML.
 - If working in MATLAB, consider “vectorizing” your code, avoiding `for` loops where possible. This often [though not always] makes the code more readable and lets it run faster.
 - It is always useful to plot the log-likelihood after every iteration. If it ever decreases (when implementing EM for ML) you have a bug! The same is true for the log-joint when implementing MAP.
- (e) Express the log-likelihoods obtained in bits and relate these numbers to the length of the naive encoding of these binary data. How does your number compare to gzip (or another compression algorithm)? Why the difference?
- (f) 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 ?
 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.
- (g) Consider the *total* cost of encoding both the model parameters and the data given the model. How does this total cost compare to gzip (or similar)? How does it depend on K ? What might this tell you?

1. (a) $\mathcal{D} = \{\underline{x}^{(1)}, \dots, \underline{x}^{(n)}\}$, $\underline{x} = \left(\begin{array}{c} \underline{x}^{(1)} \\ \vdots \\ \underline{x}^{(n)} \end{array}\right) \in \mathbb{R}^{D \times n}$ $d: \# \text{ of pixel}$
 $\underline{P} = \begin{pmatrix} P_{11} & \cdots & P_{1D} \\ \vdots & \ddots & \vdots \\ P_{K1} & \cdots & P_{KD} \end{pmatrix} \in \mathbb{R}^{K \times D}$, $\underline{p}_k = (P_{k1} \dots P_{KD}) \in \mathbb{R}^{1 \times D}$ $n: \# \text{ of image}$.
 $\underline{\pi} = (\pi_1, \dots, \pi_K) \in \mathbb{R}^{1 \times K}$

$$P(\underline{x}^{(n)} | \underline{p}_k) = \prod_{d=1}^D p_{kd}^{\underline{x}^{(n)}_d} \cdot (1 - p_{kd})^{(1 - \underline{x}^{(n)}_d)}$$

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

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

(b) $r_{nk} = P(S^{(n)}=k | \underline{x}^{(n)}, \underline{\pi}, \underline{P})$ (definition)

$$= \frac{P(S^{(n)}=k | \underline{\pi}, \underline{P}) \cdot P(\underline{x}^{(n)} | S^{(n)}=k, \underline{\pi}, \underline{P})}{\sum_{j=1}^K P(S^{(n)}=j | \underline{\pi}, \underline{P}) \cdot P(\underline{x}^{(n)} | S^{(n)}=j, \underline{\pi}, \underline{P})}$$

$$= \frac{\pi_k \cdot P(\underline{x}^{(n)} | S^{(n)}=k, \underline{\pi}, \underline{P})}{\sum_{j=1}^K \pi_j \cdot P(\underline{x}^{(n)} | S^{(n)}=j, \underline{\pi}, \underline{P})} = \frac{\pi_k \cdot \prod_{d=1}^D p_{kd}^{\underline{x}^{(n)}_d} \cdot (1 - p_{kd})^{(1 - \underline{x}^{(n)}_d)}}{\sum_{j=1}^K \pi_j \prod_{d=1}^D p_{jd}^{\underline{x}^{(n)}_d} \cdot (1 - p_{jd})^{(1 - \underline{x}^{(n)}_j)}}$$

(c) $E = \left\langle \sum_{n=1}^N \log P(\underline{x}^{(n)}, S^{(n)} | \underline{\pi}, \underline{P}) \right\rangle_{f(\{S^{(n)}\})}$. π_k from definition

$$= \sum_{k=1}^K \sum_{n=1}^N f(S^{(n)}=k) \log \left\{ P(\underline{x}^{(n)} | S^{(n)}=k, \underline{\pi}, \underline{P}) \cdot P(S^{(n)}=k | \underline{\pi}, \underline{P}) \right\}$$

$$= \sum_{k=1}^K \sum_{n=1}^N r_{nk} \log \left\{ \prod_{d=1}^D p_{kd}^{\underline{x}^{(n)}_d} \cdot (1 - p_{kd})^{(1 - \underline{x}^{(n)}_d)} \cdot \pi_k \right\}$$

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

* $\frac{\partial E}{\partial p_k} = \sum_{n=1}^N r_{nk} \left\{ \frac{\underline{x}^{(n)}}{p_k} - \frac{(1 - \underline{x}^{(n)})}{1 - p_k} \right\} = 0$.

$$\Leftrightarrow \sum_n r_{nk} \underline{x}^{(n)} - p_k \sum_n r_{nk} \underline{x}^{(n)} = p_k \sum_n r_{nk} - p_k \sum_n r_{nk} \underline{x}^{(n)}$$

$$\Leftrightarrow \underline{p}_k = \frac{\sum_n r_{nk} \underline{x}^{(n)}}{\sum_n r_{nk}} = \frac{1}{N_k} \sum_n r_{nk} \cdot \underline{x}^{(n)}$$

* $\frac{\partial}{\partial \pi_k} \left\{ \left\langle \sum_{n=1}^N \log P(\underline{x}^{(n)}, S^{(n)} | \underline{\pi}, \underline{P}) \right\rangle_{f(\{S^{(n)}\})} + \lambda \left(\sum_{j=1}^K \pi_j - 1 \right) \right\} = 0$

$$\Leftrightarrow \lambda = -N \quad \text{and} \quad \pi_k = \frac{1}{N} \sum_n r_{nk} = \frac{N_k}{N}$$

λ: Lagrangian multiplier for constrained optimisation

(d) Body, run this in MATLAB

```
1 clear; clc; clf; load binarydigits.txt -ascii;
2 Y=binarydigits;
3 [N, D] = size(Y); % D = 64 (number of pixels), N = 100 (number of images)
4 X = Y'; % dataset X, with size(X) = [D N]
5 K = 10; % chosen K value
6 % Random initialise pi and P, where pi is the vector of mixing coefficient;
7 % P is the matrix of Bernoulli parameters
8 pi = rand(1,K);
9 pi = pi/sum(pi); % normalise and update pi, so that sum(pi) = 1
10 P = rand(K,D);
11 % Iteration
12 num_steps = 20;
13 log_likelihood_store = zeros(num_steps,1);
14 for i = 1:num_steps
15     R = calculateResponsibilities(pi, P, X, K); % E step
16     P = calculateParameter(R,X,K); % M step
17     pi = calculateMixingCoefficient(R,X,K); % M step
18     log_likelihood_store(i) = evaluateLogLikelihood(R,P,pi,X,K); % M step
19 end
20 log_likelihood_store % print to check whether it converges
21 plot(log_likelihood_store)
```

```
1 function R = calculateResponsibilities(pi, P, X, K)
2 % Calculate responsibility matrix R
3 % Input:
4 % pi: mixing coefficient vector
5 % P: parameter matrix
6 % X: data
7 % K: number of mixtures
8 [D, N] = size(X);
9 R = zeros(N,K);
10 product_store = zeros(K,1);
11 for n = 1:N
12     for k = 1:K
13         product_store(k,1) = prod(P(k,:)' .^ X(:,n) .* (1-P(k,:))' .^(1-X(:,n)));
14     end
15     denominator = pi * product_store;
16     for k = 1:K
17         numerator = pi(k) * product_store(k,1);
18         R(n,k) = numerator / denominator;
19     end
20 end
```

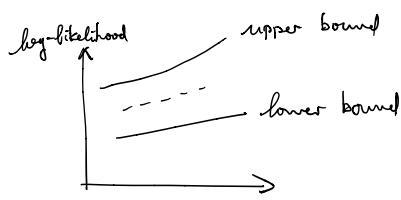
```
1 function P = calculateParameter(R,X,K)
2 % Calculate parameter matrix P
3 % Input:
4 % R: responsibility matrix
5 % X: data
6 % K: number of mixtures
7 [D N] = size(X);
8 P = zeros(K,D);
9 for k = 1:K
10    P(k,:) = (R(:,k)'*X'./sum(R(:,k)) )';
11 end
12
```

```
1 function pi = calculateMixingCoefficient(R,X,K)
2 % Calculate Mixing coefficient vector pi
3 % Input:
4 % R: responsibility matrix
5 % X: data
6 % K: number of mixtures
7 [D, N] = size(X);
8 pi = zeros(1,K);
9 for k = 1:K
10    pi(k) = sum(R(:,k))/N;
11 end
12
```

```
1 function log_likelihood = evaluateLogLikelihood(R,P,pi,X,K)
2 % Calculate the updated log likelihood
3 % Input:
4 % R: responsibility matrix
5 % P: parameter matrix
6 % pi: mixing coefficient
7 % X: data
8 % K: number of mixtures
9 [D, N] = size(X);
10 log_likelihood = 0;
11 product_store = zeros(K,1);
12 for n = 1:N
13    for k = 1:K
14       product_store(k,1) = prod(P(k,:)'.^X(:,n).*(1-P(k,:))'.^(1-X(:,n)));
15    end
16    log_likelihood = log_likelihood + log(pi * product_store);
17 end
```

(e) log-likelihood in bits can be expressed by : - $\frac{\text{log likelihood}}{\log 2}$. We lose some information here. While gzip compression does not loss any information in compression, it searches for repeated patterns from the data.

(f) We obtain different log-likelihood results and different parameters with different initial conditions with same K . But the resulting log-likelihood are very close. For example, for $K=4$, the log-likelihood converges around -2800 and -3000 . If K increases, the log-likelihood increases, and both upper bound and lower bound increases as well. But difference between upper bound and lower bound increases.



We notice that if K is small, parameter converges to very different values, i.e. it is unstable with small number of clusters. When K increases, EM converges more stable.

Hence the suggestion is to increase K , so other more "features" in the training data could be "explained". (Also likelihood increases).

Also we do multiple initialisation to find a "better" local maxima, hopefully the global maxima.

(g). the total cost of encoding in our algorithm is less than that of gzip.

if we increase the value of K , the total cost increases, and it become close to the cost in gzip if K continues to increase.

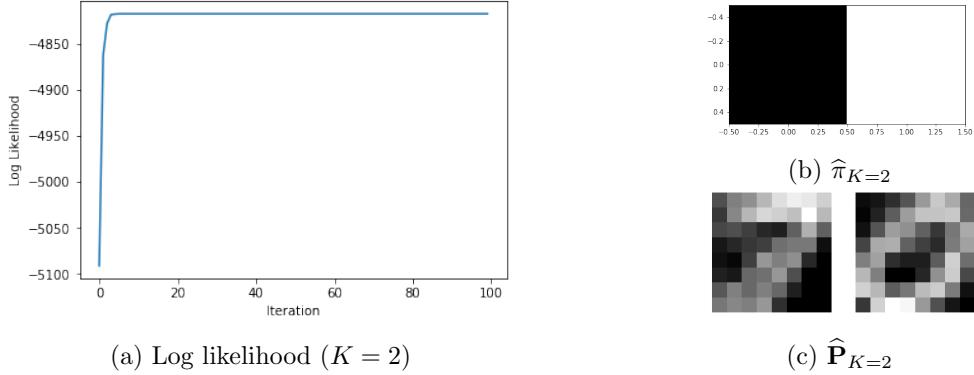


Figure 7: Log Likelihood and estimated parameters when $K = 2$

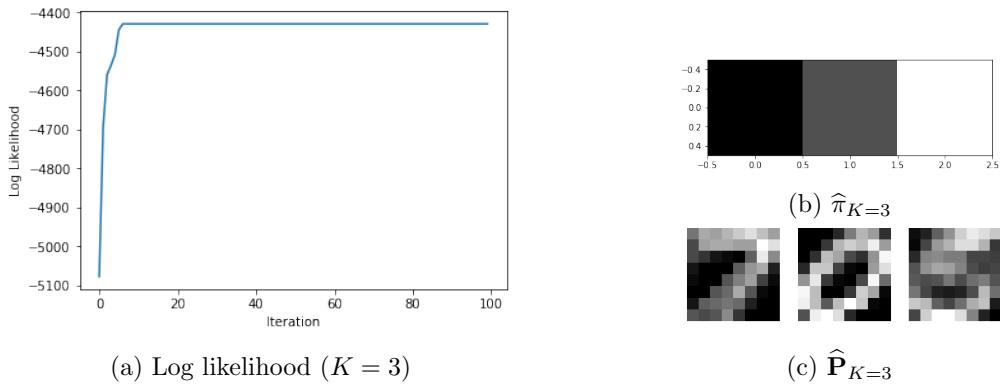


Figure 8: Log Likelihood and estimated parameters when $K = 3$

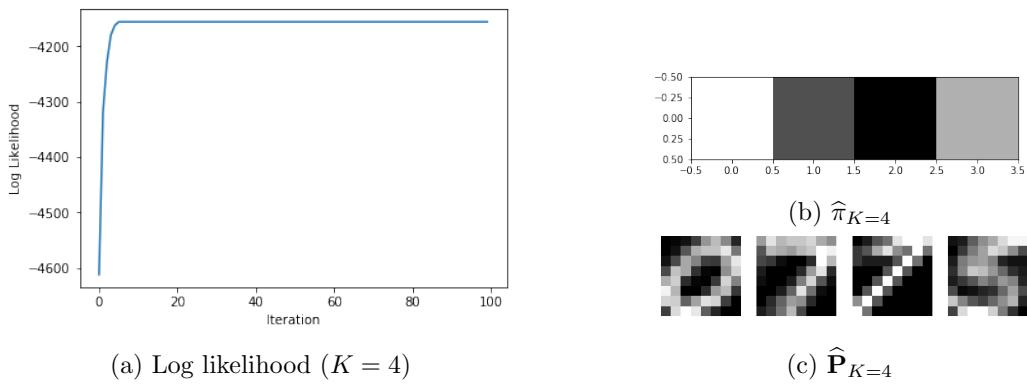


Figure 9: Log Likelihood and estimated parameters when $K = 4$

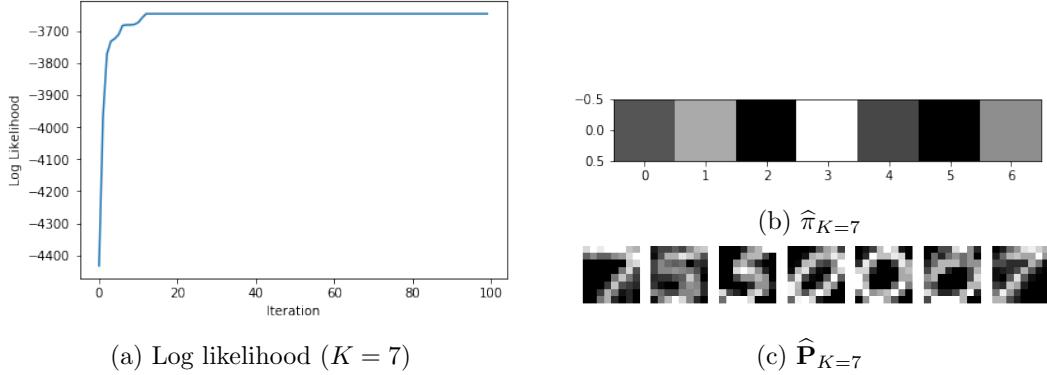


Figure 10: Log Likelihood and estimated parameters when $K = 7$

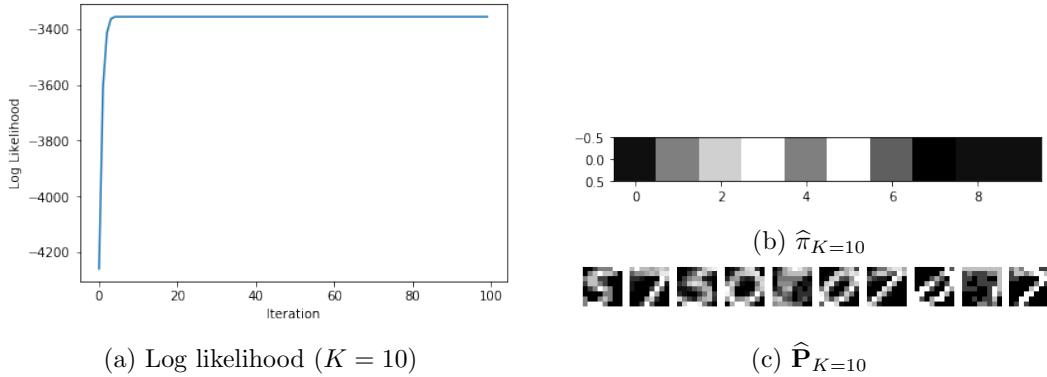


Figure 11: Log Likelihood and estimated parameters when $K = 10$

(e)

Here, the length of the naive encoding of these binary data is:

$$8 \times 8 \times 100 = 6400$$

By applying gzip, the 800 bytes large `binarydigits.bin` file can be compressed to 690 bytes (5520 bits). And the log-likelihoods obtained in bits can be calculated as:

$$-\log_2 \text{likelihood}$$

The comparison of these numbers is summarized in Table 2:

According to Table 2, the log-likelihoods in bits for $K = 2, 3, 4, 7, 10$ are all smaller than the length of naive encoding and are also smaller than the gzip size. This fact shows that the results of Mixtures of Multivariate Bernoulli modeling here **perform data compression**.

Since gzip algorithm **performs lossless data compression**, which means that we can reconstruct the original data sets accurately from the compressed one. We can somehow view 5520 bits obtained

K	Log Likelihood in bits	Naive Encoding	Gzip
2	4817		
3	4429		
4	4155		
7	3645	6400	5520
10	3355		

Table 2: Comparison of log-likelihoods in bits, naive encoding, and gzip compression

here as the minimal compressed information we need to carry out the lossless reconstruction of our `binarydigits` data.

However, the log-likelihoods in bits are all **smaller than this number**. Therefore, it can be implied that the Mixtures of Multivariate Bernoulli Model perform a kind of **lossy data compression**, through which we can not reconstruct the original data from the compressed one.

	Mixtures of Multivariate Bernoulli					gzip
	$K = 2$	$K = 3$	$K = 4$	$K = 7$	$K = 10$	
Compression Rate	75.27 %	69.20%	64.92%	56.95%	52.42%	86.25%

Table 3: Compression Rate of Mixtures of Multivariate Bernoulli Model and gzip

By further examining the compression rates shown in Table 3, we know that **the more mixture component K we set for the Mixtures of Multivariate Bernoulli Model, the higher compression we gain, and the more information we loss** for reconstruction of the original data.

(f)

Here, we run the algorithm for 10 times starting from randomly generated initial conditions. By displaying the estimated \mathbf{P} , it can be seen that the estimated clusters indicate different (hand-written) numbers.

The solutions $\{\log P(\mathcal{D}|\hat{\boldsymbol{\pi}}, \hat{\mathbf{P}}), \hat{\boldsymbol{\pi}}, \hat{\mathbf{P}}\}$ obtained under different initial conditions are all different. However, if we consider the $\hat{\mathbf{p}}_k$ (row vector of \mathbf{P}) to be the same when they indicate the same hand written digits, then we actually gain some same solutions of $\hat{\mathbf{P}}$ upon permutation within the 10 runs. Also, if we round off $\hat{\boldsymbol{\pi}}$ to 2 digit decimals, we gain some same solutions of $\hat{\boldsymbol{\pi}}$ upon permutation within the 10 runs as well.

And we are more likely to obtain the same solution of $\hat{\mathbf{P}}$ and $\hat{\boldsymbol{\pi}}$ increases when K is small.

Results of the solutions when $K = 2, 3, 10$ are displayed respectively in Table 4, Table 5 and Figure 12.

Table 4: Cluster Mean and Log Likelihoods obtained in 10 run ($K = 2$)

Run	Cluster Mean		Run	Cluster Mean	
	digit 7	digit 0		digit 5	digit 0
1	0.34	0.66	2	0.58	0.42
3	0.29	0.71	4	0.58	0.42
7	0.29	0.71	5	0.59	0.41
8	0.32	0.68	6	0.58	0.42
			9	0.6	0.4
			10	0.59	0.41

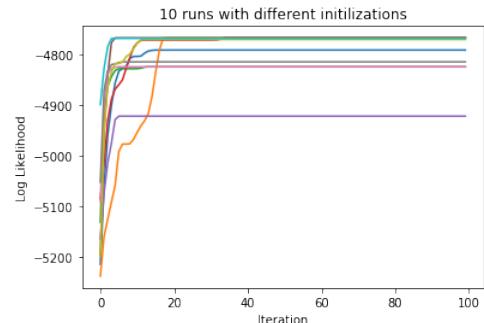
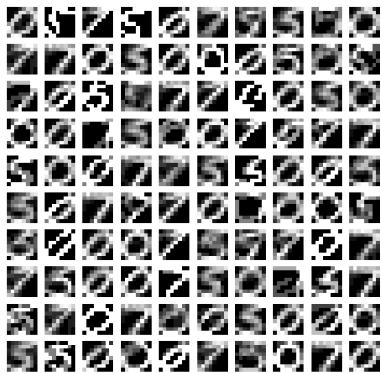
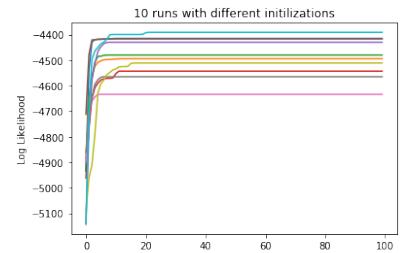
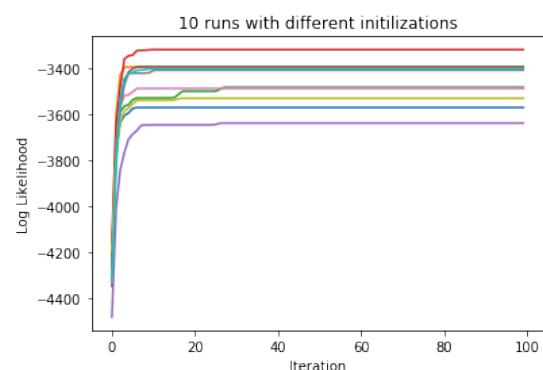


Table 5: Cluster Mean and Log Likelihoods obtained in 10 runs ($K = 3$)

Run	Cluster Mean			Run	Cluster Mean		
	digit 7	digit 5	digit 0		digit 7	digit 0	digit 0
1	0.15	0.4	0.45	2	0.49	0.29	0.22
3	0.28	0.23	0.49	8	0.35	0.45	0.22
4	0.32	0.41	0.27				
5	0.32	0.28	0.4				
6	0.37	0.33	0.3				
9	0.39	0.17	0.44				
10	0.25	0.35	0.4				
				7	0.28	0.11	0.61



(a) $\hat{\mathbf{P}}$ in 10 runs



(b) Log likelihoods in 10 runs

Figure 12: $\hat{\mathbf{P}}$ and Log Likelihoods in 10 runs ($K = 10$) under different initialization

We could deduce from the results that the algorithm performs more stable when K is small (most stable when $K = 2$), in the sense that the algorithm finds better clusters and the proportion of the data set is more stable under small K . For instance, when $K = 2$, the algorithm finds clusters of either **digit 7 and 0 pair** or **digit 5 and 0 pair**, and the proportion is (almost) 3 : 7 and 6 : 4

respectively.

However, when K gets larger, the algorithm could not find good clusters. As is shown in Table 5, the algorithm finds clusters of **digit 7, 5 and 0 pair** or **digit 7, 7 and 0 pair** or **digit 7, 0 and 0 pair**, where there is already duplication in digits (7 and 0 appear twice as different clusters). And the cluster mean also becomes very unstable. For even larger K such as $K = 10$, it can be seen from Figure 6 that there are even more duplication of digit clusters (too many digit 7, 5 and 0 !) in each run.

Therefore, we could conclude that the algorithm tends to fall into different local maximum under different initialization, and the larger the K , the more likely solutions will be different. Moreover, the algorithm works well and finds better clusters on relatively small K .

One way to improve the model is to introduce Dirichlet priors over the K multivariate Bernoulli model, including the prior information about what kind of digits we can observe from the data (say, only digit 0, 7 and 5 are observed).

(g)[BONUS]

The total cost of encoding both the model parameters and the data can be computed as

$$\text{total cost} = -\log_2 \text{likelihood} + K \times D(\text{size of } \mathbf{P}) + K(\text{size of } \pi)$$

and the results are summarized in Table 6:

K	total cost	gzip
2	4947	
3	4624	
4	4415	5520
7	4100	
10	4005	

Table 6: total cost of encoding when $K = 2, 3, 4, 7, 10$ versus gzip

The total cost of encoding when $K = 2, 3, 4, 7, 10$ are still smaller than gzip. However, as the encoding cost of parameters increase when K increase, the total cost does not decrease that much when K gets larger. This might tell us that there exists a trade-off relation between compression of data and complexity of model. When we use complex model (with more parameters), the lossy compression rate might be high but we also have to pay high costs for large parameters size.