

Assignment13: Illustration of the Expectation-Maximization Algorithm using MNIST Dataset

yifan

2025-05-14

Introduction

In class, we learned the general expectation-maximization (E-M) algorithm:¹

Given a joint distribution $p(X, Z|\theta)$ over observed variables X and latent variables Z , governed by parameters θ , the goal is to maximize the likelihood function $p(X|\theta)$ with respect to θ .

1. Choose an initial setting for the parameters θ^{old} .
2. **E step:** Evaluate $p(Z|X, \theta^{\text{old}})$.
3. **M step:** Evaluate θ^{new} given by

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$$

where

$$Q(\theta, \theta^{\text{old}}) = \sum_Z p(Z|X, \theta^{\text{old}}) \ln p(X, Z|\theta).$$

4. Check for convergence of either the log likelihood or the parameter values. If the convergence criterion is not satisfied, then let

$$\theta^{\text{old}} \leftarrow \theta^{\text{new}}$$

and return to step 2.

However, to illustrate this algorithm, we would probably need a more specific example of how this algorithm can be applied. Let's consider a finite mixture of Bernoulli distributions, as mentioned in the textbook:²

$$p(\mathbf{x}|\mu, \pi) = \sum_{k=1}^K \pi_k p(\mathbf{x}|\mu_k) \quad (1)$$

Here, $\mathbf{x} = (x_1, \dots, x_D)^T \in R^D$, $\mu = (\mu_1, \dots, \mu_k) \in R^{D \times k}$, $\pi = (\pi_1, \dots, \pi_k) \in R^k$ and

$$p(\mathbf{x}|\mu_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{(1-x_i)} \quad (2)$$

Scenario Description

How could we understand the above model more intuitively? In textbook, a scenario in machine learning is considered. We consider an MNIST dataset where each pixel's grayscale value has been binarized to 0 or 1. In this case, since the images in the MNIST dataset have a resolution of 28×28 pixels, here $D = 784$. N represents the number of images, and K represents the number of image class (without selection, $K = 10$). Our goal is to estimate the unknown parameters using maximum-likelihood.

R Implementation

First, since MNIST data is in grey scale $[0, 255]$, we have to convert it to a binary scale using a threshold of 128:

```

1 # Load MNIST data
2 mnist <- getOMLDataSet(data.id = 554)
3 mnist_data <- mnist$data
4 mnist_data <- mnist_data[, -ncol(mnist_data)] # Remove target column
5
6 # Convert to binary scale [0,1] with threshold 128
7 mnist3 <- as.data.frame(ifelse(mnist_data/128 >= 1, 1, 0))

```

Implementation of the function to plot the image:

```

1 # Function to plot MNIST data
2 show <- function(image) {

```

```

3 image_matrix <- matrix(as.numeric(image), nrow = 28, byrow = TRUE)
4 image_matrix <- t(apply(image_matrix, 2, rev)) # Rotate for proper
  orientation
5
6 ggplot(data = reshape2::melt(image_matrix), aes(x = Var1, y = Var2, fill
  = value)) +
7   geom_tile() +
8   scale_fill_gradient(low = "white", high = "black") +
9   theme_void() +
10  theme(legend.position = "none") +
11  coord_equal()
12 }

```

Implementation of the function to calculate the Bernoulli probability (2):

```

1 # Bernoulli probability calculation
2 bernoulli <- function(data, means) {
3   N <- nrow(data)
4   K <- nrow(means)
5   prob <- matrix(0, nrow = N, ncol = K)
6
7   for (i in 1:N) {
8     for (k in 1:K) {
9       prob[i, k] <- prod((means[k, ] ^ data[i, ]) * ((1 - means[k, ]) ^ (1
        - data[i, ])))
10    }
11  }
12
13  return(prob)
14 }

```

Implementation of the function to calculate the responsibilities (E step)

$$\gamma(z_{nk}) = \frac{\pi_k p(\mathbf{x}_n | \mu_k)}{\sum_{j=1}^k \pi_j p(\mathbf{x}_n | \mu_j)} \quad (3)$$

```

1 # Responsibility calculation
2 respBernoulli <- function(data, weights, means) {
3   # Step 1: Calculate p(x|means)
4   prob <- bernoulli(data, means)
5

```

```

6   # Step 2: Calculate numerator
7   prob <- prob * matrix(weights, nrow = nrow(data), ncol =
      length(weights), byrow = TRUE)
8
9   # Step 3: Calculate denominator
10  row_sums <- rowSums(prob)
11
12  # Step 4: Calculate responsibilities
13  if (any(row_sums == 0)) {
14    warning("Division by zero occurred in responsibility calculations!")
15    return(prob)
16  }
17
18  prob <- prob / row_sums
19  return(prob)
20 }

```

Implementation of the function used to re-estimate the parameters using the current responsibilities (M step)

$$\begin{cases} \mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(Z_{nk}) \mathbf{x}_n \\ \pi_k = \frac{N_k}{N} \end{cases} \quad (4)$$

where $N_k = \sum_{n=1}^N \gamma(Z_{nk})$:

```

1  # M-step for Bernoulli mixture
2  bernoulliMStep <- function(data, resp) {
3    N <- nrow(data)
4    D <- ncol(data)
5    K <- ncol(resp)
6
7    Nk <- colSums(resp)
8    mus <- matrix(0, nrow = K, ncol = D)
9
10   for (k in 1:K) {
11     mus[k, ] <- colSums(resp[, k] * data)
12     if (Nk[k] == 0) {
13       warning("Division by zero occurred in Mixture of Bernoulli Dist
          M-Step!")
14       next
15     }

```

```

16     mus[k, ] <- mus[k, ] / Nk[k]
17   }
18
19   return(list(weights = Nk / N, means = mus))
20 }

```

Implementation of log-likelihood calculation for determination of convergence:

$$E_Z[\ln p(\mathbf{X}, \mathbf{Z} | \mu, \pi)] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left\{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\} \quad (5)$$

```

1  # Log-likelihood calculation
2  llBernoulli <- function(data, weights, means) {
3    N <- nrow(data)
4    K <- nrow(means)
5
6    resp <- respBernoulli(data, weights, means)
7
8    ll <- 0
9    for (i in 1:N) {
10     sumK <- 0
11     for (k in 1:K) {
12       temp1 <- (means[k, ] ^ data[i, ]) * ((1 - means[k, ]) ^ (1 - data[i,
13         ]))
14       temp1 <- log(pmax(temp1, 1e-50)) # Avoid log(0)
15       sumK <- sumK + resp[i, k] * (log(weights[k]) + sum(temp1))
16     }
17     ll <- ll + sumK
18   }
19   return(list(ll = ll, resp = resp))
20 }

```

R function to run the EM algorithm:

```

1  # EM algorithm for mixture of Bernoulli distributions
2  mixOfBernoulliEM <- function(data, init_weights, init_means, maxiters =
3    1000, relgap = 1e-4, verbose = FALSE) {
4    N <- nrow(data)
5    D <- ncol(data)

```

```

5   K <- nrow(init_means)
6
7   # Initialize
8   weights <- init_weights
9   means <- init_means
10  ll_result <- llBernoulli(data, weights, means)
11  ll_old <- ll_result$ll
12  resp <- ll_result$resp
13
14  for (i in 1:maxiters) {
15    if (verbose && (i %% 5 == 0)) {
16      cat(sprintf("iteration %d:\n", i))
17      cat("  ", weights, "\n")
18      cat(sprintf(" %.6f\n", ll_old))
19    }
20
21    # M Step
22    mstep_result <- bernoulliMStep(data, resp)
23    weights <- mstep_result$weights
24    means <- mstep_result$means
25
26    # Convergence check
27    ll_result <- llBernoulli(data, weights, means)
28    ll <- ll_result$ll
29    resp <- ll_result$resp
30
31    if (abs(ll - ll_old) < relgap) {
32      cat(sprintf("Relative gap: %.8f at iterations %d\n", ll - ll_old, i))
33      break
34    } else {
35      ll_old <- ll
36    }
37  }
38
39  return(list(weights = weights, means = means))
40 }

```

R code to pick specific digits from the MNIST dataset:

```

1 # Function to pick specific digits
2 pickData <- function(digits, N) {

```

```

3  set.seed(30)
4  shuffled_idx <- sample(nrow(mnist3))
5  sData <- mnist3[shuffled_idx, ]
6  sTarget <- mnist$data[shuffled_idx, ncol(mnist$data)]
7
8  selected <- sData[as.integer(sTarget) %in% digits, ]
9  selected <- selected[sample(nrow(selected), N), ]
10 return(selected)
11 }

```

Implementation of the experiments function:

```

1  # Experiment function
2  experiments <- function(digits, K, N, iters = 100, verbose = FALSE) {
3    expData <- pickData(digits, N)
4    D <- ncol(expData)
5
6    initWts <- runif(K, 0.25, 0.75)
7    initWts <- initWts / sum(initWts)
8
9    initMeans <- as.matrix(expData[sample(nrow(expData), K, replace =
10      FALSE), ])
11    initMeans <- initMeans * (0.6 + 0.8 * matrix(runif(K * D), nrow = K,
12      ncol = D))
13    initMeans <- pmin(pmax(initMeans, 0.05), 0.95)
14
15    return(mixOfBernoulliEM(expData, initWts, initMeans, maxiters = iters,
16      relgap = 1e-15, verbose = verbose))
17  }

```

Note that in contrast to the python code which initialized the mean values using uniform distribution, I added randomness to the initialization of the mean values. In this way, the algorithm can be more stable, especially when the number of samples is small.

Clustering Performance Evaluation for Different Choice of K

Of course, we could easily run the code to produce the result illustrated in Figure 9.10 in the textbook. However, it is hard to judge the performance of our algorithm based on that figure.

Therefore, we could take a step further, and sort our data into training set and testing set. However, our task is somewhat different from traditional machine learning: our output can be seen as a result of **clustering**, rather than **classification**. In other words, it would be difficult to calculate the accuracy of our model on the testing set, since we only cluster the images in the testing set, without labeling them.

However, there are other ways to evaluate our clustering performance. Say that K types of digits are included. Since in this section, we assume we know the ground truth K , the mixture model would output K clusters.

We evaluate the performance of our mixture model by the metric **adjusted rand index**(ARI).³ It is a widely used metric for the evaluation of clustering performance:

$$ARI = \frac{RI - E(RI)}{\max(RI) - E(RI)} \quad (6)$$

which can be understood intuitively as normalizing RI into the range of $[-1, 1]$. Here,

$$RI = \frac{TP + TN}{TP + TN + FP + FN} \quad (7)$$

Where we consider the pairs in the testing set and categorize them into:

- True Positives (TP): These are pairs of elements that are correctly grouped together in both the predicted and actual clustering.
- True Negatives (TN): These are pairs of elements that are correctly separated in both the predicted and actual clustering.
- False Positives (FP): These are pairs of elements that are grouped together in the predicted clustering but were not grouped together in the actual clustering.
- False Negatives (FN): These are pairs of elements that are grouped together in the actual clustering but were not grouped together in the predicted clustering.

A positive ARI value indicates that the performance of the model is better than random grouping.

Now, we implement this in R.

First, we need to revise the pickData function to return the true label of the testing set:

```
1 pickData <- function(digits, N) {
2   # Shuffle data and labels while maintaining correspondence
3   set.seed(30)
4   shuffled_indices <- sample(nrow(mnist3))
5   sData <- mnist3[shuffled_indices, ]
6   sTarget <- mnist.target[shuffled_indices]
7
8   # Filter data for specified digits and collect corresponding labels
9   filtered_data <- list()
10  filtered_labels <- list()
11  for (i in 1:nrow(sData)) {
12    if (as.integer(sTarget[i]) %in% digits) {
13      filtered_data <- append(filtered_data, list(as.numeric(sData[i, ])))
14      filtered_labels <- append(filtered_labels, as.integer(sTarget[i]))
15    }
16  }
17
18  # Convert to matrices
19  filtered_data <- do.call(rbind, filtered_data)
20  filtered_labels <- unlist(filtered_labels)
21
22  # Randomly select N samples while maintaining data-label correspondence
23  set.seed(30)
24  selected_indices <- sample(nrow(filtered_data), N)
25
26  selected_data <- filtered_data[selected_indices, ]
27  selected_labels <- filtered_labels[selected_indices]
28
29  return(list(data = selected_data, labels = selected_labels))
30 }
```

Then we implement a function to divide data into training set and testing set:

```
1 dataPartition <- function(digits, N1, N2) {
2   # Get combined data and labels
3   result <- pickData(digits, N1 + N2)
4   data <- result$data
```

```

5  labels <- result$labels
6
7  # Split into training and testing sets
8  training_set <- data[1:N1, ]
9  testing_set <- data[(N1+1):(N1+N2), ]
10 testing_set_labels <- labels[(N1+1):(N1+N2)]
11
12 return(list(
13   training_set = training_set,
14   testing_set = testing_set,
15   testing_set_labels = testing_set_labels
16 ))
17 }

```

Implementation for generated the predicted cluster for each data sample in the testing set (simply picking out the largest value in the posterior distribution) and the implementation of ARI:

```

1  library(gtools) # For comb() function
2
3  convert_cluster <- function(posterior) {
4    # Convert posterior probabilities to cluster assignments
5    # Equivalent to np.argmax(posterior, axis=1)
6    apply(posterior, 1, which.max) - 1 # Subtract 1 for 0-based indexing
7  }
8
9  adjusted_rand_index <- function(true_labels, pred_labels) {
10   # Calculate the adjusted Rand index for clustering evaluation
11
12   # Create contingency table
13   contingency <- table(factor(pred_labels, levels = 0:max(pred_labels)),
14                        factor(true_labels, levels = 0:max(true_labels)))
15
16   # Convert to matrix
17   contingency <- as.matrix(contingency)
18
19   # Calculate sums
20   sum_comb_c <- sum(comb(contingency, 2))
21   sum_comb_a <- sum(comb(rowSums(contingency), 2))
22   sum_comb_b <- sum(comb(colSums(contingency), 2))
23

```

```

24  n <- length(true_labels)
25  expected_ri <- sum_comb_a * sum_comb_b / comb(n, 2)
26  max_ri <- 0.5 * (sum_comb_a + sum_comb_b)
27
28  ari <- (sum_comb_c - expected_ri) / (max_ri - expected_ri)
29  return(ari)
30 }

```

Finally, we note that the EM algorithm is in fact not very stable, since we initialize the initial value of π_k and μ_k randomly. Sometimes, the algorithm would fall into locally optimal solution. Therefore, we choose to run the algorithm B times and calculate the mean and standard error of $ARI_i, i = 1, \dots, B$.

```

1  # Run EM algorithm B times and calculate Adjusted Rand Index for each run
2  run_EM <- function(digits, K, N1, N2, B) {
3    results <- numeric(B) # Initialize vector to store results
4
5    for (i in 1:B) {
6      # Partition data into training and testing sets
7      partitioned_data <- dataPartition(digits, N1, N2)
8      training_set <- partitioned_data$training_set
9      testing_set <- partitioned_data$testing_set
10     testing_set_labels <- partitioned_data$testing_set_labels
11
12     # Run EM experiments and get final means
13     finMeans <- experiments(K, training_set = training_set, verbose = TRUE)
14
15     # Calculate cluster assignments for test data
16     cluster_label <- convert_cluster(respBernoulli(testing_set,
17       finMeans[[1]], finMeans[[2]]))
18
19     # Compute Adjusted Rand Index and store result
20     results[i] <- adjusted_rand_index(testing_set_labels, cluster_label)
21   }
22   return(results)
23 }
24
25 # Calculate mean and standard deviation while handling missing/zero values
26 mean_and_std <- function(vector) {
27   # Convert input to numeric vector and handle missing values

```

```

28  vec <- as.numeric(vector)
29
30  # Filter out NA/NaN and zero values
31  filtered_vec <- vec[!is.na(vec) & vec != 0]
32
33  # Return NA if no valid values remain
34  if (length(filtered_vec) == 0) {
35    return(list(mean = NA, std = NA))
36  }
37
38  # Calculate mean and sample standard deviation (n-1 denominator)
39  mean_val <- mean(filtered_vec)
40  std_val <- sd(filtered_vec)
41
42  return(list(mean = mean_val, std = std_val))
43 }

```

In practice, to test the performance of our mixture model on the task of k -clustering, we simply choose the digits $[0, \dots, k - 1]$. We make sure that for each digit category, there are approximately 1000 training images and 1000 testing images, i.e. $N_1 = N_2 = 1000k$. The number of experiments conducted for each clustering task is set to $B = 10$, and the experiment is conducted on a L4 GPU on colab.

The results are shown in Figure 1:

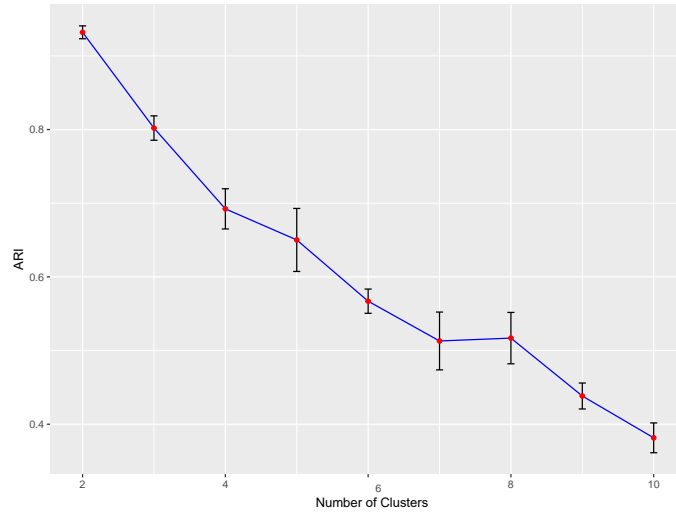


Figure 1: Model Performance

Note that the data is represented as mean and 95% confidence intervals. Here,

$$\begin{cases} ARI_{mean_i} = \sum_{j=1}^B ARI_{ij}/10 \\ ARI_{sd_i} = \sum_{j=1}^B (ARI_{ij} - ARI_{mean_i})^2/9 \\ ARI_{CI_i} = [ARI_{mean_i} - t_{\alpha/2,df} ARI_{sd_i}/\sqrt{10}, ARI_{mean_i} + t_{\alpha/2,df} ARI_{sd_i}/\sqrt{10}], \end{cases} \quad (8)$$

Where $\alpha = 0.05$, $df = 9$, $t_{\alpha/2,df} \approx 2.2622$.

From the figure above, it is obvious that the clustering performance decreases as the number of clusters increases, though the number of training images remains at approximately 1000. However, there isn't an obvious pattern for the width of 95% confidence intervals, subject to the value of K . This may be due to the small number of experiments we conducted on each choice of K ($B = 10$).

Criterion for the Selection of K

Now in this section, we assume that the ground truth K is unknown, and we try to come up with a criterion to estimate K .

First, since we could obtain the posterior likelihood:

$$\ln p(\mathbf{X}|\mu, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k p(\mathbf{X}_n|\mu_k) \right\} \quad (9)$$

```

1 loglikelihood <- function(data, weights, means) {
2   # Calculate probability of each sample under each component: p(x_n|_k)
3   prob <- bernoulli(data, means) # Shape (N, K)
4
5   # Weighted sum: sum_{k=1}^K _k p(x_n|_k)
6   weighted_prob <- prob * matrix(weights, nrow = nrow(data), ncol =
7     length(weights), byrow = TRUE)
8
9   # Compute log-likelihood: sum_n ln{sum_k _k p(x_n|_k)}
10  # Using log(sum) instead of logSumExp for simplicity (add protection if
11    needed)
12  log_likelihood <- sum(log(rowSums(weighted_prob)))
13  return(log_likelihood)
14 }
```

One may want to choose K using the Akaike information criterion (AIC):⁴

$$\hat{k}_{AIC} = \operatorname{argmin}_k \{2p_k - 2\ln(\hat{L})\} \quad (10)$$

Or the Bayesian information criterion (BIC):⁵

$$\hat{k}_{BIC} = \operatorname{argmin}_k \{p_k \ln(n) - 2\ln(\hat{L})\} \quad (11)$$

Here, $\ln(\hat{L})$ represents the maximum value of the log-likelihood function, $p_k = kD + (k-1)$ represents the number of parameters included in the mixture model when $K = k$. Note that although there are k parameters included in $\pi = (\pi_1, \dots, \pi_k)$, since $\sum_{i=1}^k \pi_i = 1$, the number of free valuable is only $k - 1$.

These two criteria both view $-2\ln(\hat{L})$ as a cost function, with additional penalty terms $2p_k, \ln(n)p_k$ to penalize the complexity of the model, respectively. However, after trying both criteria, I found that the cost function dominates overwhelmingly in both, and the penalty term hardly has any effect, leading them to consistently overestimate the value of K .

Therefore, we would probably need a different selection criterion. When the ground truth K is fixed, $\ln(\hat{L})$ is a monotonically increasing function, and in practice, I have found that it is generally a concave function. Therefore, we could apply the **elbow criterion**. Actually, this criterion is widely applied to the selection of clustering number. The idea behind the elbow criterion is to identify the inflection point where the log-likelihood transitions from rapid growth to slow growth.

However, this may require manually identifying the inflection point. Is there an algorithm to automatically select the value K ? Our method here is to draw a straight line linking the first point and the last point of the curve. Then, we find out the point that has the maximum distance to that line (Figure 2):

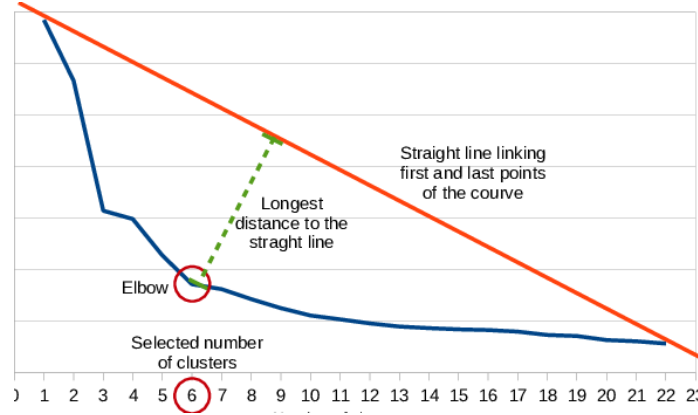


Figure 2: Illustration of the Elbow Criterion

Note that in contrast to Figure 2, our curve is monotonically increasing and concave, instead of monotonically decreasing and convex.

Implementation of the elbow criterion using R:

```
1 library(ggplot2)
2 library(kneed) # For KneeLocator equivalent
3
4 estimate_K_by_elbow <- function(digits, N, max_K = 10, iters = 50) {
5   # Partition data into training and testing sets
6   data_partition <- dataPartition(digits, N, N)
7   training_set <- data_partition$training_set
8   testing_set <- data_partition$testing_set
9   testing_labels <- data_partition$testing_labels
10
11   K_values <- 1:max_K
12   log_likelihoods <- numeric(length(K_values))
13
14   for (i in seq_along(K_values)) {
15     K <- K_values[i]
16     tryCatch({
17       # Train model
18       model <- experiments(K, training_set, iters = iters, verbose = FALSE)
19       weights <- model$weights
20       means <- model$means
21
22       # Calculate log-likelihood (with stable implementation)
23       if (K == 1) {
24         # For K=1, use direct Bernoulli log-likelihood
25         ll <- sum(training_set * log(means + 1e-100) +
26                 sum((1 - training_set) * log(1 - means + 1e-100)))
27       } else {
28         # For K>1, use mixture log-likelihood
29         ll <- loglikelihood(training_set, weights, means)
30       }
31
32       log_likelihoods[i] <- ll
33       cat(sprintf("K=%d: log-likelihood=%.1f\n", K, ll))
34
35     }, error = function(e) {
36       cat(sprintf("K=%d failed: %s\n", K, conditionMessage(e)))
37       log_likelihoods[i] <- -Inf # Mark failure
```

```

38   })
39 }
40
41 # Automatic elbow detection
42 kneedle <- KneeLocator(
43   x = K_values,
44   y = log_likeliheids,
45   curve = "concave", # Log-likelihood curve is typically concave
46   direction = "increasing"
47 )
48 best_K <- kneedle$elbow
49
50 # Visualization
51 elbow_plot <- ggplot(data.frame(K = K_values, LL = log_likeliheids),
52   aes(x = K, y = LL)) +
53   geom_line(color = "magenta", linewidth = 1) +
54   geom_point(color = "magenta", size = 3) +
55   labs(x = "Number of clusters (K)", y = "Log-Likelihood",
56     title = "Elbow Method for K Selection") +
57   theme_minimal() +
58   theme(panel.grid.major = element_line(color = "gray", alpha = 0.3))
59
60 if (!is.na(best_K)) {
61   elbow_plot <- elbow_plot +
62     geom_vline(xintercept = best_K, color = "red", linetype = "dashed",
63       alpha = 0.5) +
64     annotate("text", x = best_K + 0.1, y = median(log_likeliheids, na.rm
65       = TRUE),
66       label = paste("Elbow K =", best_K), color = "red")
67 }
68
69 print(elbow_plot)
70
71 return(list(best_K = best_K, log_likeliheids = log_likeliheids))
72 }

```

Now, we test the performance of the elbow criterion. Set the ground truth of cluster numbers to $K = 4$, and run the R code:

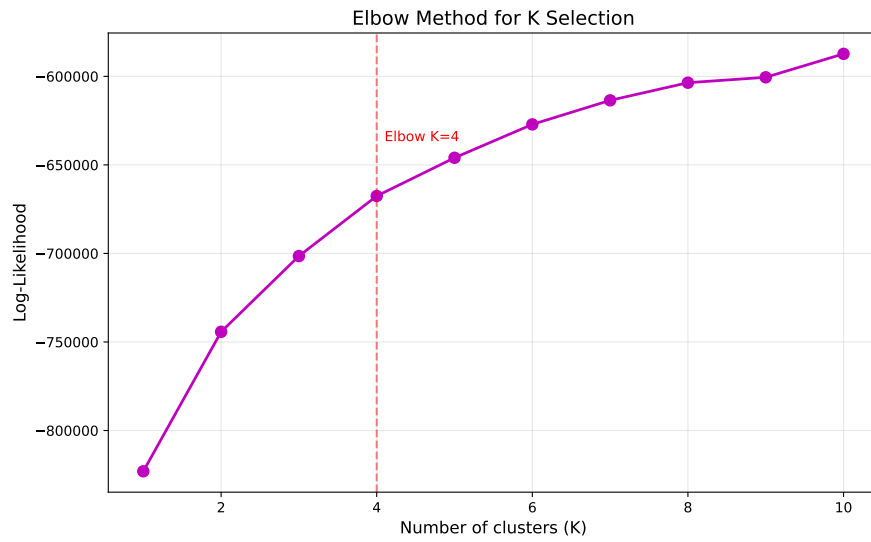


Figure 3: Elbow Criterion when $K = 4$

Set the ground truth of cluster numbers to $K = 5$, and run the R code:

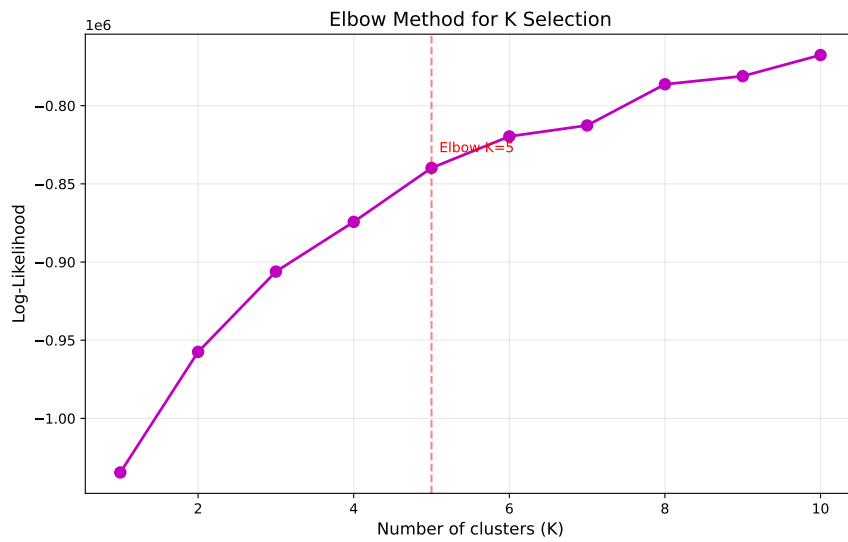


Figure 4: Elbow Criterion when $K = 5$

We could see that the elbow criterion successfully identifies the ground truth of K for both cases, demonstrating ideal performance.

Discussion

In this assignment, we applied the E-M algorithm to the MNIST dataset, and tested its clustering performance for different values of K . We could see that as K increases, the performance of our mixture model decreases rapidly. In fact, MNIST is a classical machine-learning dataset for classification, and current neural-networks achieves an accuracy of over 99% on testing sets.⁶ It is obvious that neural-networks would outperform our traditional mixture model, which isn't surprising.

Moreover, we developed a criterion for estimation of the clustering number K using the elbow criterion. Due to time constraint, we only ran the algorithm two times, and both time we obtained the correct result. However, repetition is needed to further validify our method.

References

- [1] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society: series B (methodological)*, 39(1):1–22, 1977.
- [2] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer, New York, 2006.
- [3] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of classification*, 2:193–218, 1985.
- [4] Hirotugu Akaike. Information theory and an extension of the maximum likelihood principle. In *Selected papers of hirotugu akaike*, pages 199–213. Springer, 1998.
- [5] Gideon Schwarz. Estimating the dimension of a model. *The annals of statistics*, pages 461–464, 1978.
- [6] Feiyang Chen, Nan Chen, Hanyang Mao, and Hanlin Hu. Assessing four neural networks on handwritten digit recognition dataset (mnist). *arXiv preprint arXiv:1811.08278*, 2018.