

Exponential Mixture Models applied to Simulation Study (Part II)

Rebecca C. Steorts

This document contains the solution to Part II of the simulation study, where there is a function below called `exponentialMixture` that implements our derivations from Part I. Furthermore, you will see that the estimates are quite similar to the true values. The function also monitors the log-loglikelihood, and we provide a plot for this over the number of EM iterations.

```
# Load necessary libraries
library(ggplot2)
exponentialMixture <- function(data, K, max_iter = 1000, tol = 1e-5) {
  n <- length(data)

  # Initialize mixing proportions (pi) and rate parameters (lambda)
  pi <- rep(1/K, K) # Mixing proportions
  lambda <- runif(K, 0.1, 1) # Rate parameters

  log_likelihoods <- numeric(max_iter) # To store log-likelihood values

  for (iter in 1:max_iter) {
    # E-step: Compute responsibilities or rather the gamma values
    gamma <- matrix(NA, nrow = n, ncol = K)
    for (k in 1:K) {
      gamma[, k] <- pi[k] * dexp(data, rate = lambda[k])
    }
    row_sums <- rowSums(gamma)
    gamma <- gamma / row_sums # Normalize probabilities

    # M-step: Update mixing proportions and rate parameters
    pi_old <- pi
    lambda_old <- lambda

    pi <- colMeans(gamma) # Update mixing proportions
    for (k in 1:K) {
      lambda[k] <- sum(gamma[, k]) / sum(gamma[, k] * data) # Update rate parameters
    }

    # Calculate log-likelihood
    log_likelihoods[iter] <- sum(log(row_sums))

    # Check for convergence
    if (max(abs(pi - pi_old)) < tol && max(abs(lambda - lambda_old)) < tol) {
      log_likelihoods <- log_likelihoods[1:iter] # Trim to the number of iterations
      cat("Convergence reached at iteration", iter,
          "with log-likelihood:", log_likelihoods[iter], "\n")
      break
    }
  }
}
```

```

}

return(list(pi = pi, lambda = lambda, log_likelihood = log_likelihooods))
}

```

Below, we simulate some data, so that we can sanity check our code and our results before applying it to a real data set. The simulated data set is set up such that we hopefully will not run into convergence issues.

```

set.seed(1234)
# sample
n_samples <- 2000
true_rates <- c(0.5, 2.0)
true_proportions <- c(0.5, 0.5)

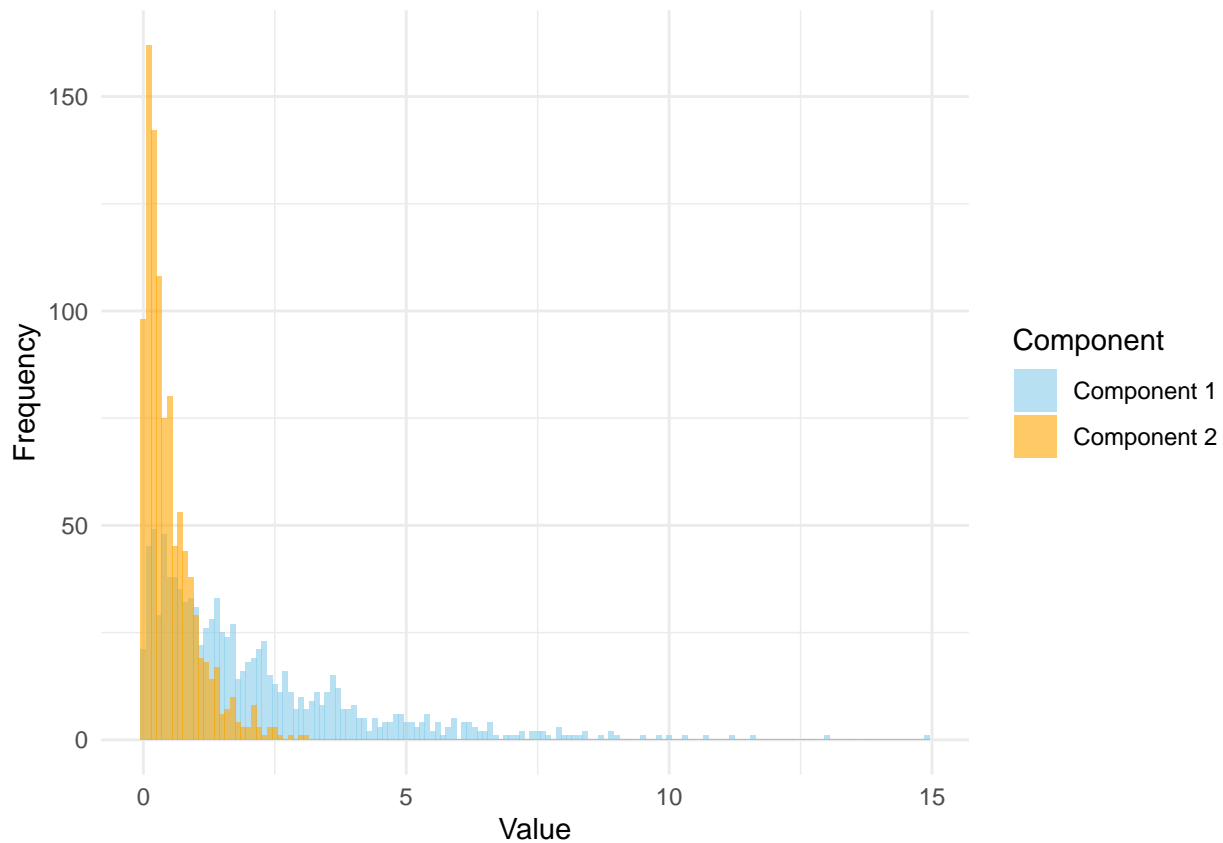
# Sample group assignments (1 or 2) based on mixing proportions
component_labels <- sample(1:2, n_samples, replace = TRUE, prob = true_proportions)

# Stage 2: Generate data from the corresponding exponential
# distribution based on the group assignment
data <- numeric(n_samples) # Initialize the data vector
# Generate data for group 1 (rate = 0.5)
data[component_labels == 1] <- rexp(sum(component_labels == 1), rate = true_rates[1])
# Generate data for group 2 (rate = 2.0)
data[component_labels == 2] <- rexp(sum(component_labels == 2), rate = true_rates[2])

# Create a data frame for plotting
plot_data <- data.frame(
  Value = data, # Data points
  Component = factor(component_labels, labels = c("Component 1", "Component 2")) # Component labels
)

# Plot the data using ggplot2
ggplot(plot_data, aes(x = Value, fill = Component)) +
  geom_histogram(binwidth = 0.1, position = "identity", alpha = 0.6) +
  scale_fill_manual(values = c("skyblue", "orange")) + # Set colors for the components
  labs(
    x = "Value",
    y = "Frequency",
    fill = "Component"
  ) +
  theme_minimal() +
  theme(plot.title = element_blank()) # Remove the title

```



```
set.seed(1234)
library(knitr)
# Fit the mixture model
result <- exponentialMixture(data, K = 2)

## Convergence reached at iteration 229 with log-likelihood: -2390.612
# Print results
print(paste("Estimated Mixing Proportions:", round(result$pi, 4)))

## [1] "Estimated Mixing Proportions: 0.5135"
## [2] "Estimated Mixing Proportions: 0.4865"
print(paste("Estimated Rate Parameters:", round(result$lambda, 4)))

## [1] "Estimated Rate Parameters: 0.4998" "Estimated Rate Parameters: 1.9721"
# Create a data frame to display the results in a table
results_table <- data.frame(
  "Component" = 1:2,
  "Estimated Mixing Proportion" = round(result$pi, 4),
  "Estimated Rate Parameter (lambda)" = round(result$lambda, 4)
)

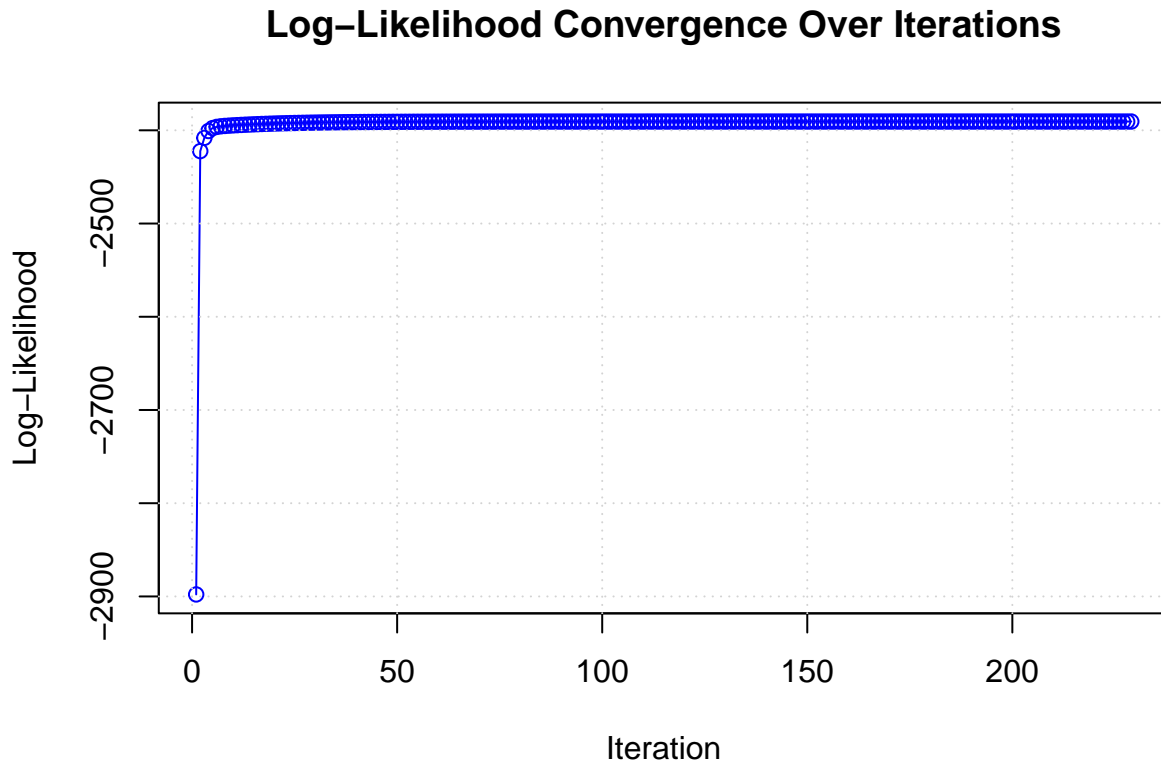
# move results to a table
kable(results_table, caption = "Estimated Parameters for the Mixture Model", format = "pipe")
```

Table 1: Estimated Parameters for the Mixture Model

Component	Estimated.Mixing.Proportion	Estimated.Rate.Parameter..lambda.
1	0.5135	0.4998
2	0.4865	1.9721

The mixture model recovers the true parameters pretty well in this situation.

```
log_likelihoods <- result$log_likelihood
plot(log_likelihoods, type = "o", col = "blue", xlab = "Iteration",
      ylab = "Log-Likelihood", main = "Log-Likelihood Convergence Over Iterations")
grid()
```



Hopefully, you have learned that coding up a function can be hard if you have not done this before. This of course is good practice for understanding the functions that you are utilizing in R and having a better understanding of machine learning methods in general and how they work if this is something that you might have an interest in moving forward. What would next steps look like? We would want to test our function out on many simulation studies and then we could reapply this example on the Snoq. Falls case study to see if it has better performance than a Gaussian mixture model. (Return to that module and see if you can tie together these concepts regarding mixture models in general).