

Project 2 EM Algorithm

Problem 1

For each observation D_i we can compute the responsibility of each component $P(C_A|D_i)$ as follows:

$$P(C_A|D_i) = \frac{P(C_A)P(D_i|\theta_A)}{P(C_A)P(D_i|\theta_A) + P(C_B)P(D_i|\theta_B)}$$

$$P(C_B|D_i) = \frac{P(C_B)P(D_i|\theta_B)}{P(C_A)P(D_i|\theta_A) + P(C_B)P(D_i|\theta_B)}$$

We know that $\sum_k P(C_k|D_i) = 1$ and therefore that $P(C_A|D_i) = 1 - P(C_B|D_i)$ for observation i . We also know that $P(D_i|\theta) = \binom{10}{k}\theta^k(1-\theta)^{10-k}$, where k is the number of heads in observation D_i . Given that $P(C_A) = 0.6$ and $P(C_B) = 0.4$, and that θ_A and θ_B are 0.7 and 0.4 respectively; for the first observation D_1 we have that $k = 4$ and therefore:

$$P(C_A|D_1) = \frac{0.6 * 0.7^4 * 0.3^6}{0.6 * 0.7^4 * 0.3^6 + 0.4 * 0.4^4 * 0.6^6} \approx 0.18$$

```
a = 0.6 * 0.7^4 * 0.3^6
b = 0.4 * 0.4^4 * 0.6^6
a/(a+b)
```

```
## [1] 0.1802056
```

$$P(C_B|D_1) = 1 - P(C_A|D_1) \approx 0.82$$

```
b/(a+b)
```

```
## [1] 0.8197944
```

For D_2 we have that $k = 8$ we have:

$$P(C_A|D_1) = \frac{0.6 * 0.7^8 * 0.3^2}{0.6 * 0.7^8 * 0.3^2 + 0.4 * 0.4^8 * 0.6^2} \approx 0.97$$

```
a = 0.6 * 0.7^8 * 0.3^2
b = 0.4 * 0.4^8 * 0.6^2
a/(a+b)
```

```
## [1] 0.9705765
```

$$P(C_B|D_1) = 1 - P(C_A|D_1) \approx 0.03$$

```
b/(a+b)
```

```
## [1] 0.02942349
```

Lastly, the update of the mixing weights $P(C_A)$ and $P(C_B)$ is given by the following expression:

$$P(C_i) = \frac{P(C_i|D_1) + P(C_i|D_2)}{2}$$

Therefore:

$$P(C_A) = \frac{0.97 + 0.18}{2} = 0.575$$

```
(0.97+0.18)/2
```

```
## [1] 0.575
```

$$P(C_B) = \frac{0.82 + 0.103}{2} = 0.425$$

```
(0.82+0.103)/2
```

```
## [1] 0.4615
```

Problem 2

In this problem, you will implement the EM algorithm for the coin toss problem in R.

Below we provide you with a skeleton of the algorithm. You can either fill this skeleton with the required functions or write your own version of the EM algorithm. If you choose to do the latter, please also present your results using Rmarkdown in a clear fashion.

```
set.seed(2022)
```

(a) Load data

We first read the data stored in the file “coinflip.csv”.

```
# read the data into D
D <- vroom("./coinflip.csv", show_col_types = FALSE)
# check the dimension of D
all(dim(D) == c(200, 100))
```

```
## [1] TRUE
```

(b) Initialize parameters

Next, we will need to initialize the mixture weights and the probabilities of obtaining heads. You can choose your own values as long as they make sense.

```
# Number of coins
k <- 2
# Mixture weights (a vector of length k)
lambda <- c(0.5, 0.5)
# Probabilities of obtaining heads (a vector of length k)
theta <- runif(k)
```

(c) The EM algorithm

Now we try to implement the EM algorithm. Please write your code in the indicated blocks.

```
##' This function implements the EM algorithm for the coin toss problem
##' @param D Data matrix of dimensions 100-by-N, where N is the number of observations
##' @param k Number of coins
##' @param lambda Vector of mixture weights
##' @param theta Vector of probabilities of obtaining heads
##' @param tolerance A threshold used to check convergence
coin_EM <- function(D, k, lambda, theta, tolerance = 1e-2) {

  # expected complete-data (hidden) log-likelihood
  ll_hid <- -Inf
  # observed log-likelihood
  ll_obs <- -Inf
  # difference between two iterations
  diff <- Inf
  # number of observations
  N <- nrow(D)
  # responsibilities
  gamma <- matrix(0, nrow = k, ncol = N)

  # run the E-step and M-step until convergence
  while (diff > tolerance) {

    # store old likelihood
    ll_obs_old <- ll_obs

    ##### E-step #####

    ### YOUR CODE STARTS ###

    # Compute the responsibilities
    P <- rbind(theta[1]^rowSums(D) * (1 - theta[1])^(dim(D)[2] - rowSums(D)),
               theta[2]^rowSums(D) * (1 - theta[2])^(dim(D)[2] - rowSums(D)))
    a <- lambda[1] * P[1, ]
    b <- lambda[2] * P[2, ]
    gamma[1, ] <- a/(a + b)
    gamma[2, ] <- b/(a + b)

    # Update expected complete-data (hidden) log-likelihood
    ll_hid = sum(gamma * log(lambda*P))

    # Update observed log-likelihood
    if (prod(a) + prod(b) == 0) {
      ll_obs = 0
    } else {
      ll_obs = log(prod(a) + prod(b))
    }

    # Recompute difference between two iterations
    diff = ll_obs - ll_obs_old

    ### YOUR CODE ENDS ###
  }
}
```

```

##### M-step #####

### YOUR CODE STARTS ###

# Recompute priors (mixture weights)
lambda = rowSums(gamma) / N

# Recompute probability of heads for each coin
num_h=rowSums(D)
count_h=gamma%*%num_h
count_t=(1-gamma)%*%(dim(D)[2]-num_h)

theta=t(count_h/(count_h+count_t))

### YOUR CODE ENDS ###

}

return(list(ll_hid = ll_hid, ll_obs = ll_obs, lambda = lambda, theta = theta, gamma = gamma))
}

```

Run the EM algorithm:

```
res <- coin_EM(D, k, lambda, theta)
```

(d) Results

Probability of heads:

```
## YOUR CODE ##
res$theta
```

```
##           [,1]      [,2]
## [1,] 0.7286612 0.267633
```

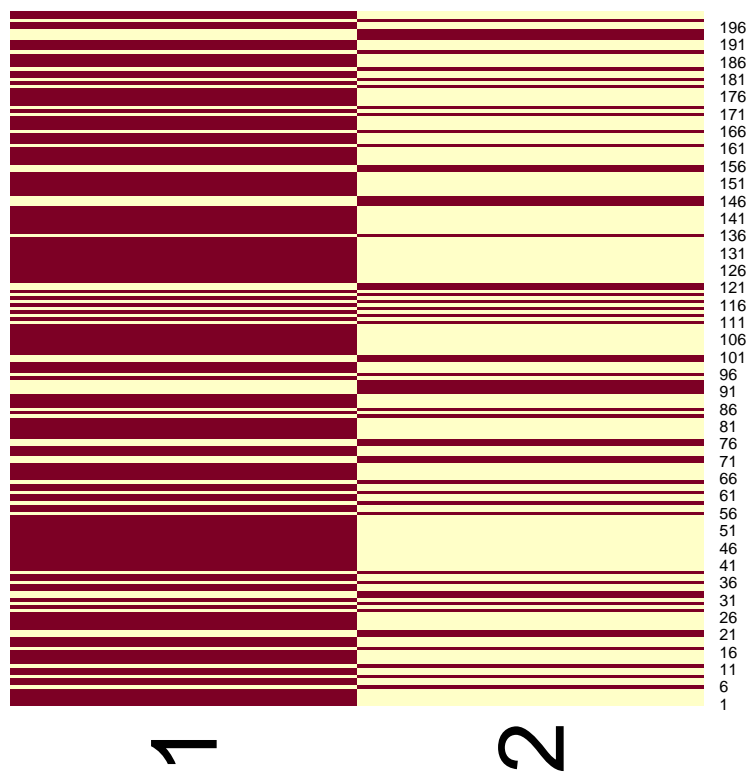
Mixture weights:

```
## YOUR CODE ##
res$lambda
```

```
## [1] 0.7307371 0.2692629
```

Heatmap of responsibilities:

```
## YOUR CODE ##
heatmap(t(res$gamma),Rowv=NA,Colv=NA)
```



How many observations belong to each coin?

```
## YOUR CODE ##
```

```
# Coin 1
```

```
obs_1 = sum(res$gamma[1,]>res$gamma[2,])
```

```
obs_1
```

```
## [1] 146
```

```
# Coin 2
```

```
dim(D)[1]-obs_1
```

```
## [1] 54
```