

machine learning(732A99) lab1 Block 2

Anubhav Dikshit(anudi287)

4 December 2018

Contents

Assignment 1	1
Loading The Libraries	1
1. Your task is to evaluate the performance of Adaboost classification trees and random forests on the spam data. Specifically, provide a plot showing the error rates when the number of trees considered are 10,20,...,100. To estimate the error rates, use 2/3 of the data for training and 1/3 as hold-out test data.	2
Loading Input files	2
Splitting into Train and Test with 66% and 33% ratio.	2
Training the Model	2
2 Your task is to implement the EM algorithm for mixtures of multivariate Bernoulli distributions. Please use the template in the next page to solve the assignment. Then, use your implementation to show what happens when your mixture models has too few and too many components, i.e. set $K = 2,3,4$ and compare results. Please provide a short explanation as well.	4
Function for EM Algorithm	4
$K = 2$	6
$K = 3$	19
$K = 4$	66
Appendix	85

Assignment 1

Loading The Libraries

1. Your task is to evaluate the performance of Adaboost classification trees and random forests on the spam data. Specifically, provide a plot showing the error rates when the number of trees considered are 10, 20, ..., 100. To estimate the error rates, use 2/3 of the data for training and 1/3 as hold-out test data.

Loading Input files

```
spam_data <- read.csv(file = "spambase.data", header = FALSE)
colnames(spam_data)[58] <- "Spam"
spam_data$Spam <- factor(spam_data$Spam, levels = c(0,1), labels = c("0", "1"))
```

Splitting into Train and Test with 66% and 33% ratio.

```
set.seed(12345)
n = NROW(spam_data)
id = sample(1:n, floor(n*(2/3)))
train = spam_data[id,]
test = spam_data[-id,]
```

Training the Model

Adaboost with varying depth

```
final_result <- NULL
for(i in seq(from = 10, to = 100, by = 10)){

  ada_model <- mboost::blackboost(Spam~.,
                                data = train,
                                family = AdaExp(),
                                control=boost_control(mstop=i))

  forest_model <- randomForest(Spam~., data = train, ntree = i)

  prediction_function <- function(model, data){
    predicted <- predict(model, newdata = data, type = c("class"))
    predict_correct <- ifelse(data$Spam == predicted, 1, 0)
    score <- sum(predict_correct)/NROW(data)
    return(score)
  }

  train_ada_model_predict <- predict(ada_model, newdata = train, type = c("class"))
  test_ada_model_predict <- predict(ada_model, newdata = test, type = c("class"))
  train_forest_model_predict <- predict(forest_model, newdata = train, type = c("class"))
  test_forest_model_predict <- predict(forest_model, newdata = test, type = c("class"))
}
```

```

test_predict_correct <- ifelse(test$Spam == test_forest_model_predict, 1, 0)
train_predict_correct <- ifelse(train$Spam == train_forest_model_predict, 1, 0)

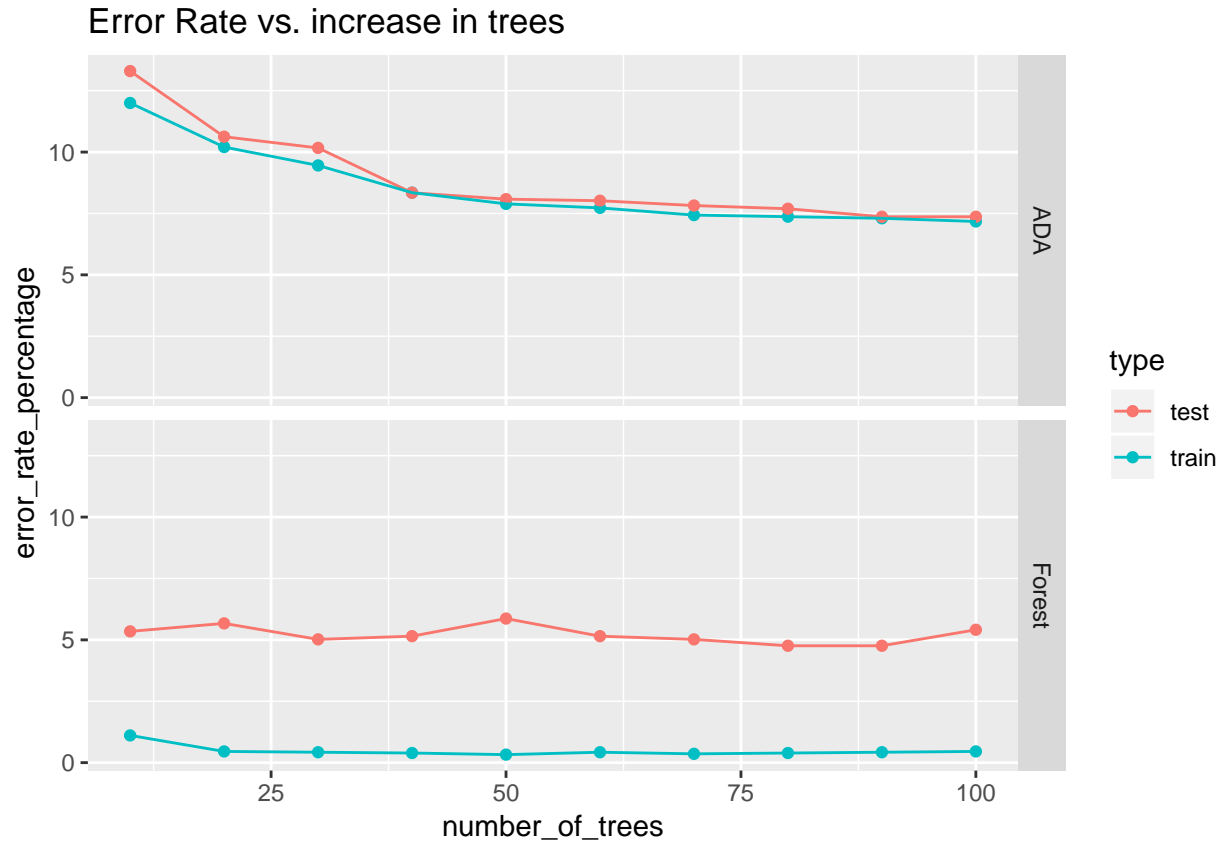
train_ada_score <- prediction_function(ada_model, train)
test_ada_score <- prediction_function(ada_model, test)
train_forest_score <- prediction_function(forest_model, train)
test_forest_score <- prediction_function(forest_model, test)

iteration_result <- data.frame(number_of_trees = i,
                              accuracy = c(train_ada_score,
                                             test_ada_score,
                                             train_forest_score,
                                             test_forest_score),
                              type = c("train", "test", "train", "test"),
                              model = c("ADA", "ADA", "Forest", "Forest"))

final_result <- rbind(iteration_result, final_result)
}

final_result$error_rate_percentage <- 100*(1 - final_result$accuracy)
ggplot(data = final_result, aes(x = number_of_trees,
                                y = error_rate_percentage,
                                group = type, color = type)) +
  geom_point() +
  geom_line() +
  ggtitle("Error Rate vs. increase in trees") + facet_grid(rows = vars(model))

```



Analysis:

From the plots we can clearly see that ADA boosted methods uses more trees (~50) to reduce the test error, while randomforest achieves saturation in short number of trees (~10). We also see that random forest achieves less error than ADA tree for both tree and test cases.

2 Your task is to implement the EM algorithm for mixtures of multivariate Bernoulli distributions. Please use the template in the next page to solve the assignment. Then, use your implementation to show what happens when your mixture models has too few and too many components, i.e. set $K = 2, 3, 4$ and compare results. Please provide a short explanation as well.

Function for EM Algorithm

```
myem <- function(K){
  set.seed(1234567890)

  max_it <- 100 # max number of EM iterations
  min_change <- 0.1 # min change in log likelihood between two consecutive EM iterations
  N=1000 # number of training points
  D=10 # number of dimensions
```

```

x <- matrix(nrow=N, ncol=D) # training data
true_pi <- vector(length = K) # true mixing coefficients
true_mu <- matrix(nrow=K, ncol=D) # true conditional distributions
true_pi=c(rep(1/3, K))

if(K == 2){
  true_mu[1,]=c(0.5,0.6,0.4,0.7,0.3,0.8,0.2,0.9,0.1,1)
  true_mu[2,]=c(0.5,0.4,0.6,0.3,0.7,0.2,0.8,0.1,0.9,0)
}else if(K == 3){
  true_mu[1,]=c(0.5,0.6,0.4,0.7,0.3,0.8,0.2,0.9,0.1,1)
  true_mu[2,]=c(0.5,0.4,0.6,0.3,0.7,0.2,0.8,0.1,0.9,0)
  true_mu[3,]=c(0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5)
}else {
  true_mu[1,]=c(0.5,0.6,0.4,0.7,0.3,0.8,0.2,0.9,0.1,1)
  true_mu[2,]=c(0.5,0.4,0.6,0.3,0.7,0.2,0.8,0.1,0.9,0)
  true_mu[3,]=c(0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5)
  true_mu[4,]=c(0.5,0.4,0.5,0.7,0.5,0.5,0.2,0.5,0.5,0.5)}

# Producing the training data
for(n in 1:N) {
  k <- sample(1:K,1,prob=true_pi)
  for(d in 1:D) {
    x[n,d] <- rbinom(1,1,true_mu[k,d])
  }
}

z <- matrix(nrow=N, ncol=K) # fractional component assignments
pi <- vector(length = K) # mixing coefficients
mu <- matrix(nrow=K, ncol=D) # conditional distributions
llik <- vector(length = max_it) # log likelihood of the EM iterations
# Random initialization of the paramters
pi <- runif(K,0.49,0.51)
pi <- pi / sum(pi)

for(k in 1:K) {
  mu[k,] <- runif(D,0.49,0.51)
}

for(it in 1:max_it) {

if(K == 2){
  plot(mu[1,], type="o", col="blue", ylim=c(0,1))
  points(mu[2,], type="o", col="red")
}else if(K == 3){
  plot(mu[1,], type="o", col="blue", ylim=c(0,1))
  points(mu[2,], type="o", col="red")
  points(mu[3,], type="o", col="green")
}else{
  plot(mu[1,], type="o", col="blue", ylim=c(0,1))
  points(mu[2,], type="o", col="red")

```

```

points(mu[3,], type="o", col="green")
points(mu[4,], type="o", col="yellow")}

Sys.sleep(0.5)
# E-step: Computation of the fractional component assignments

for(k in 1:K)
prod <- exp(x %*% log(t(mu))) * exp((1-x) %*% t(1-mu))

num = matrix(rep(pi,N), ncol = K, byrow = TRUE) * prod
dem = rowSums(num)
poster = num/dem

#Log likelihood computation.
llik[it] = sum(log(dem))
# Your code here
cat("iteration: ", it, "log likelihood: ", llik[it], "\n")
flush.console()
# Stop if the log likelihood has not changed significantly
if( it != 1){
if(abs(llik[it] - llik[it-1]) < min_change){break}
}
#M-step: ML parameter estimation from the data and fractional component assignments
# Your code here
num_pi = colSums(poster)
pi = num_pi/N
mu = (t(poster) %*% x)/num_pi

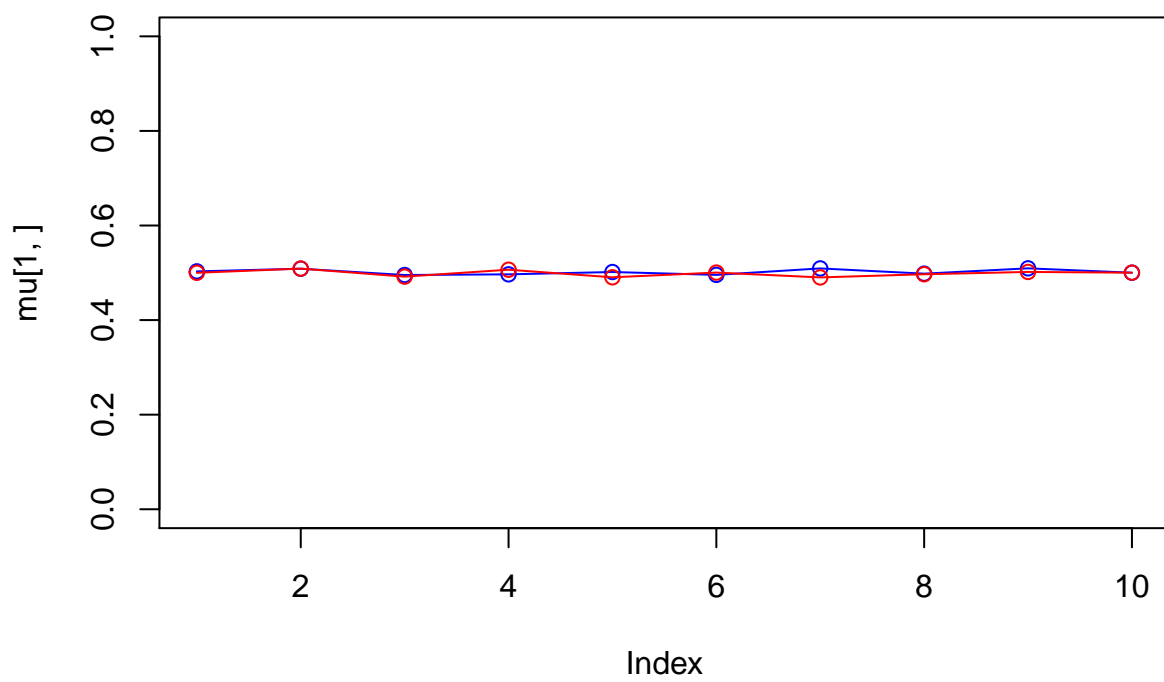
}

a <- pi
b <- mu
c <- plot(llik[1:it], type="o")
result <- list(c(a,b,c))
return(result)
}

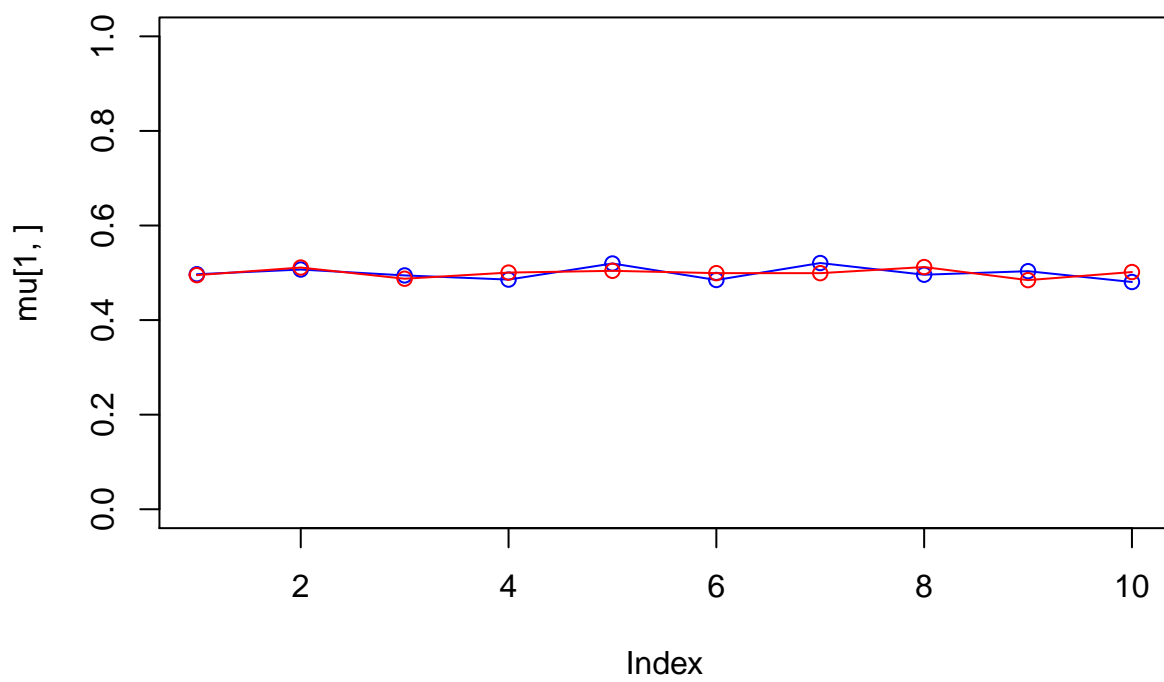
```

$K = 2$

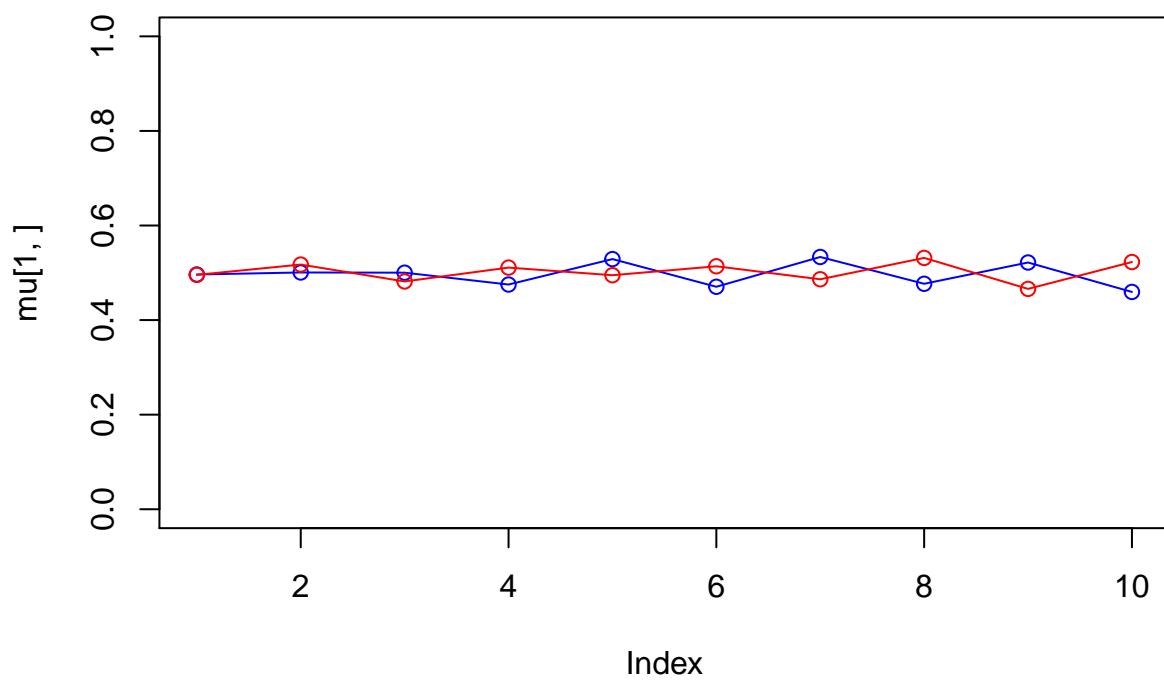
```
myem(K=2)
```



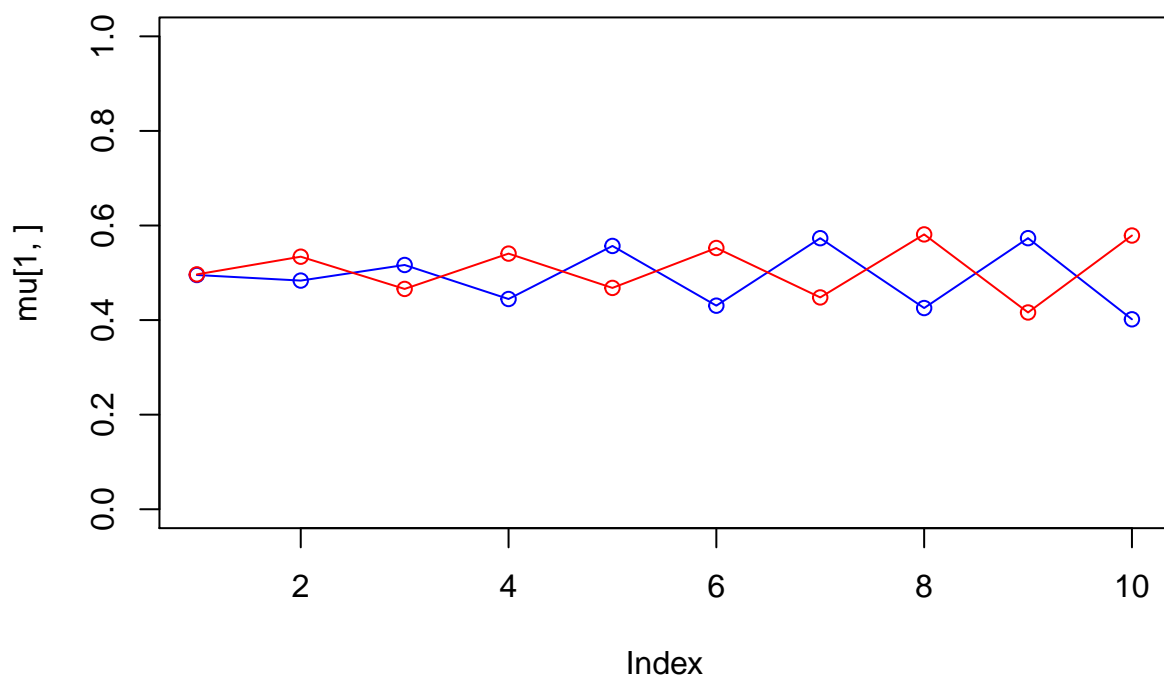
iteration: 1 log likelihood: -954.7133



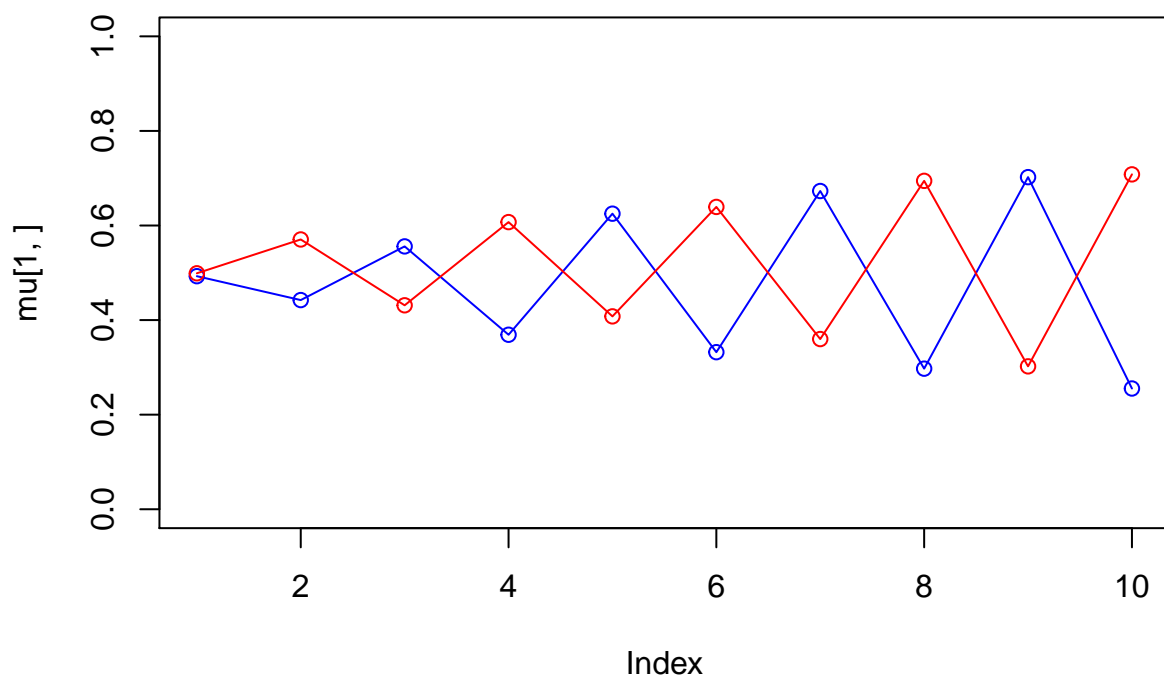
iteration: 2 log likelihood: -957.1002



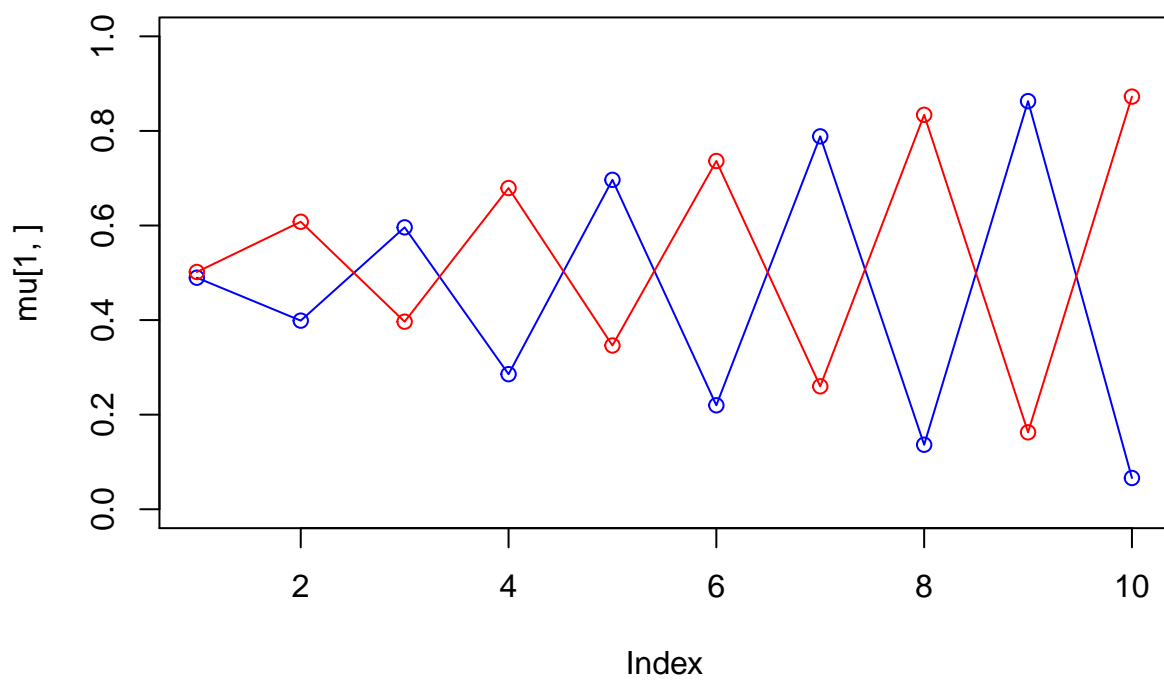
iteration: 3 log likelihood: -944.9229



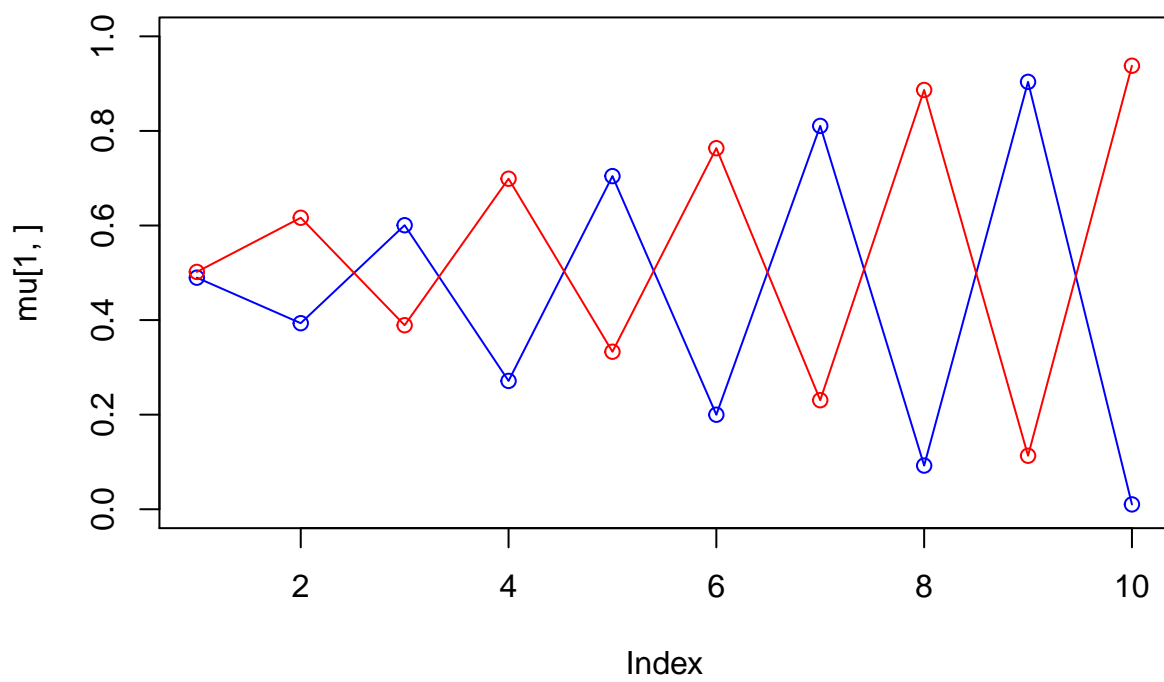
iteration: 4 log likelihood: -857.3443



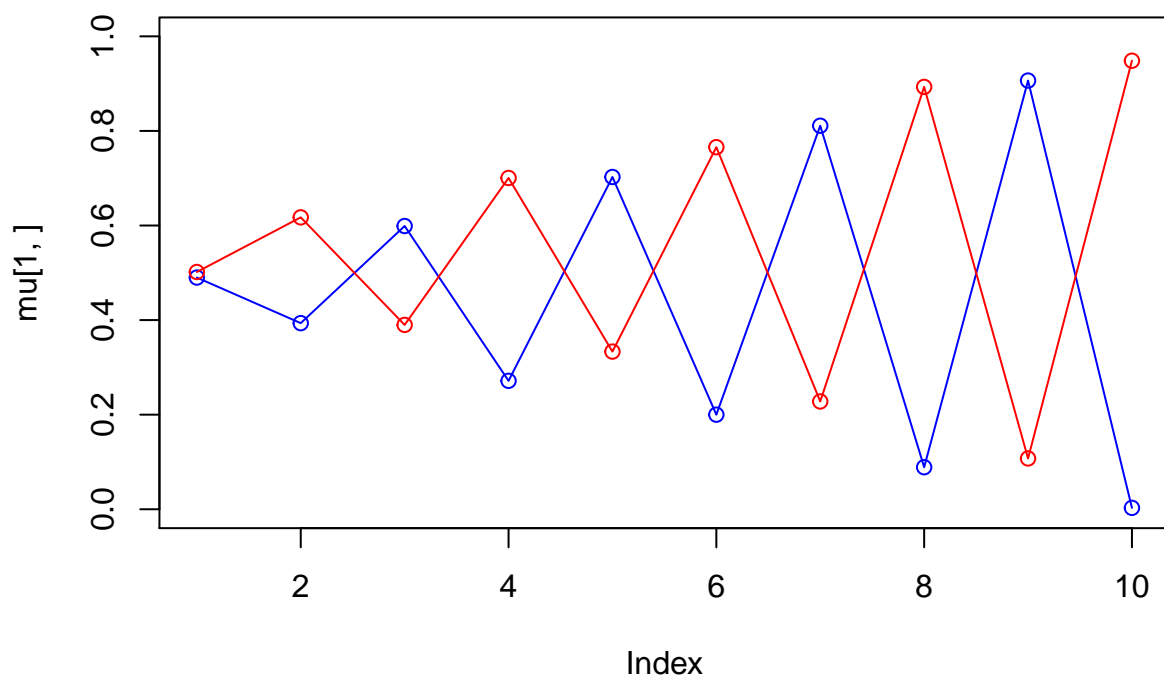
iteration: 5 log likelihood: -464.0063



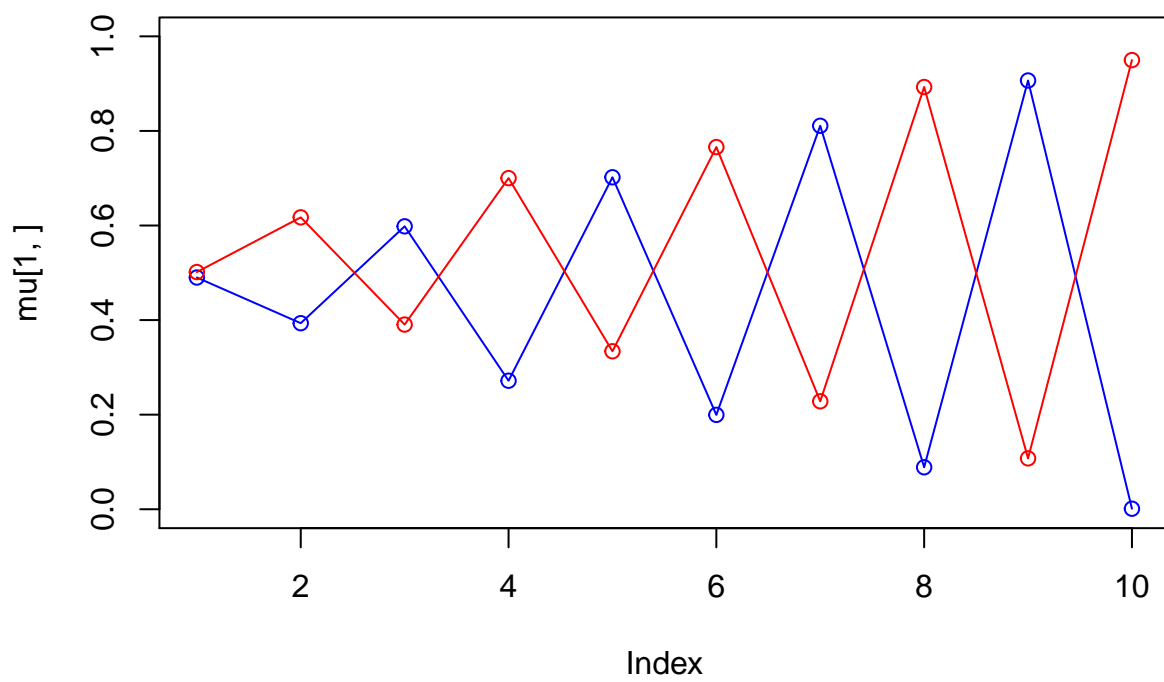
iteration: 6 log likelihood: 50.2616



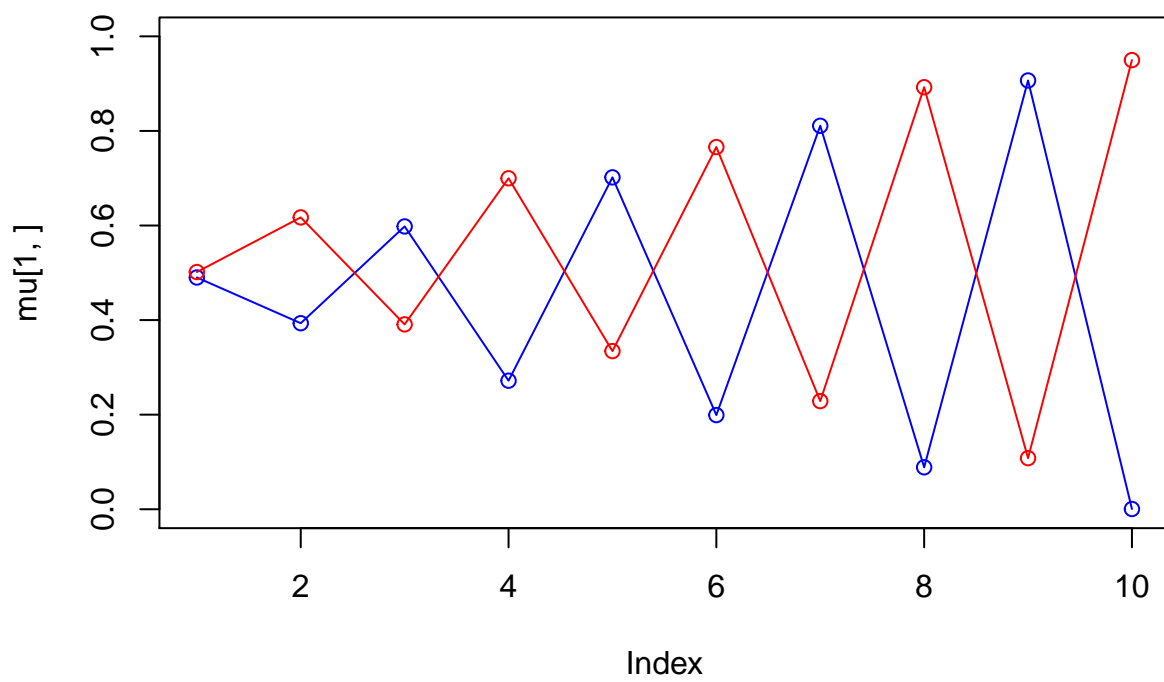
iteration: 7 log likelihood: 177.0235



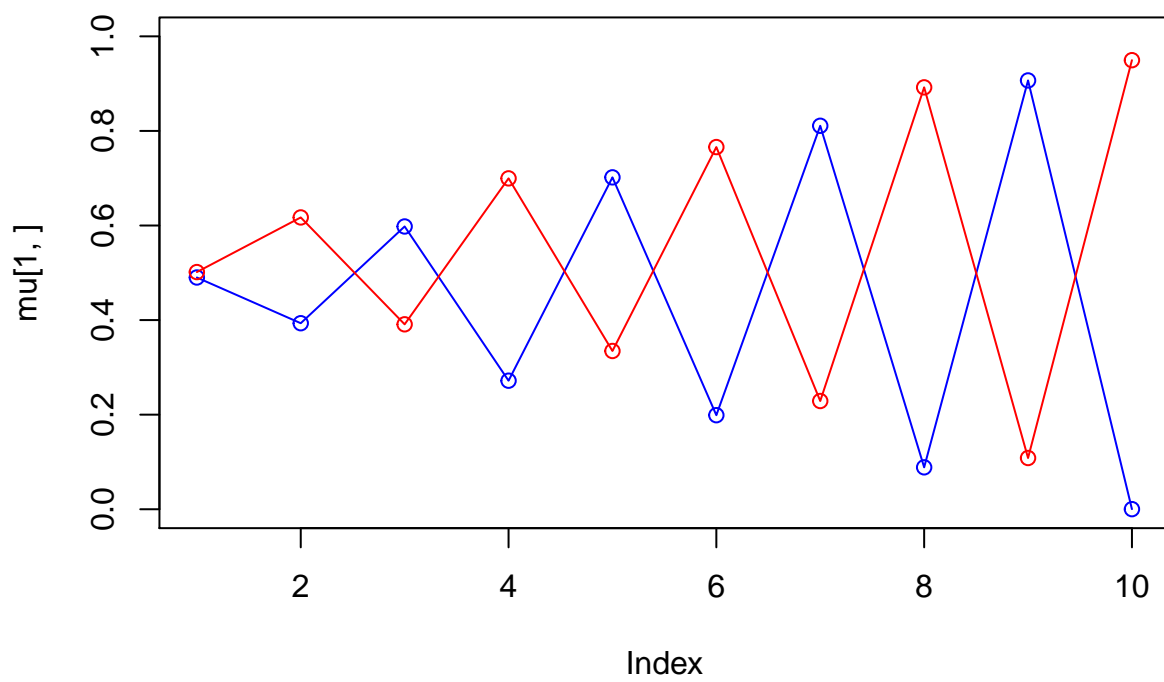
iteration: 8 log likelihood: 189.1059



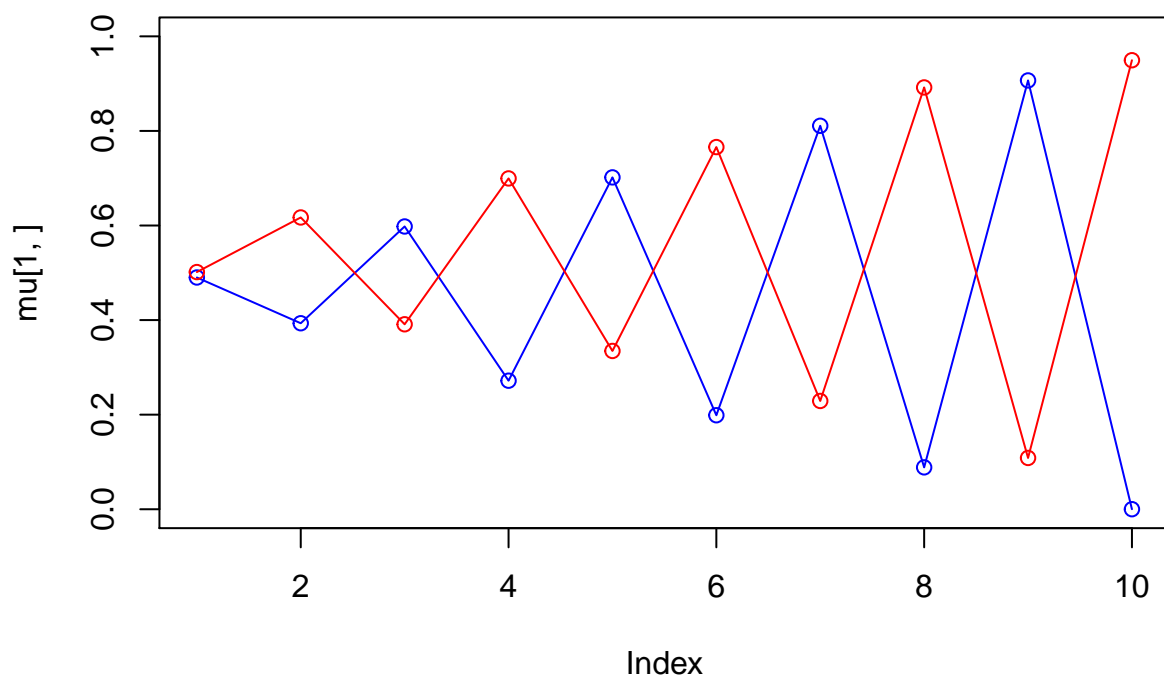
iteration: 9 log likelihood: 190.0362



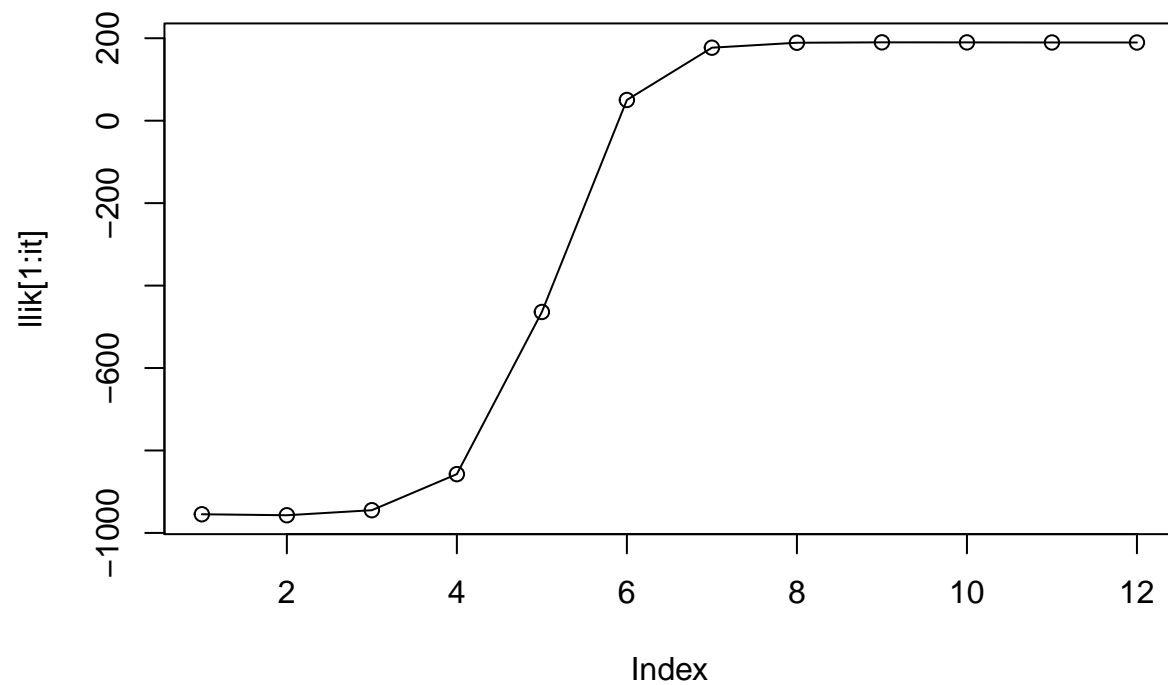
iteration: 10 log likelihood: 189.9033



iteration: 11 log likelihood: 189.7476



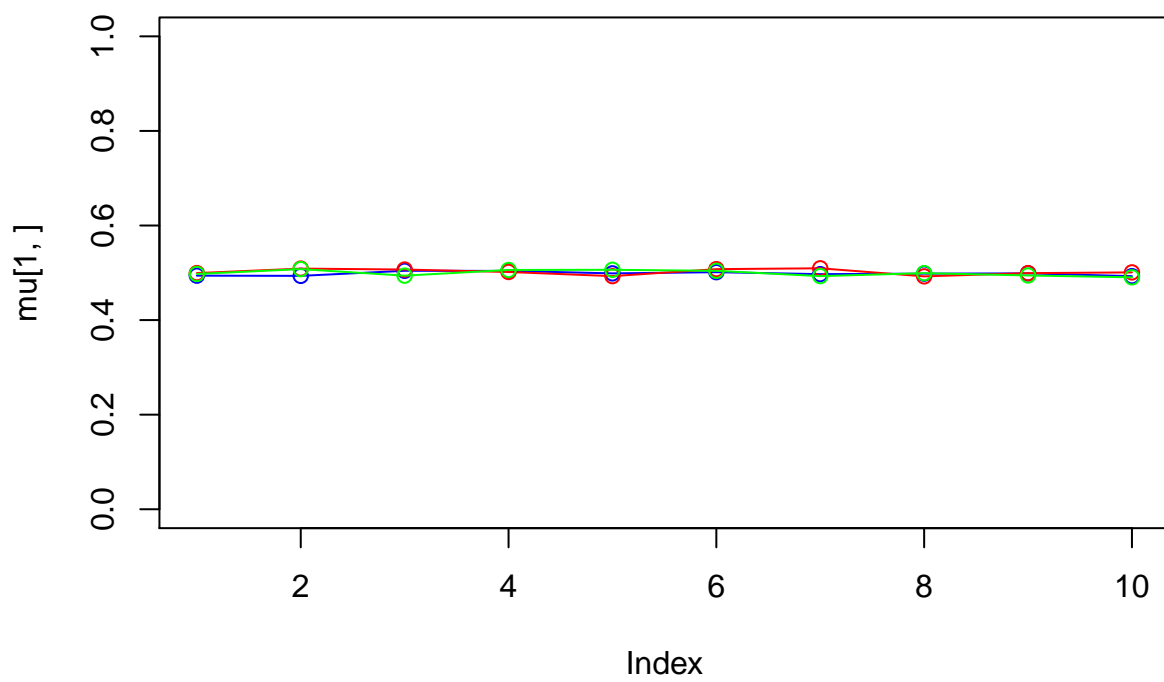
iteration: 12 log likelihood: 189.6572



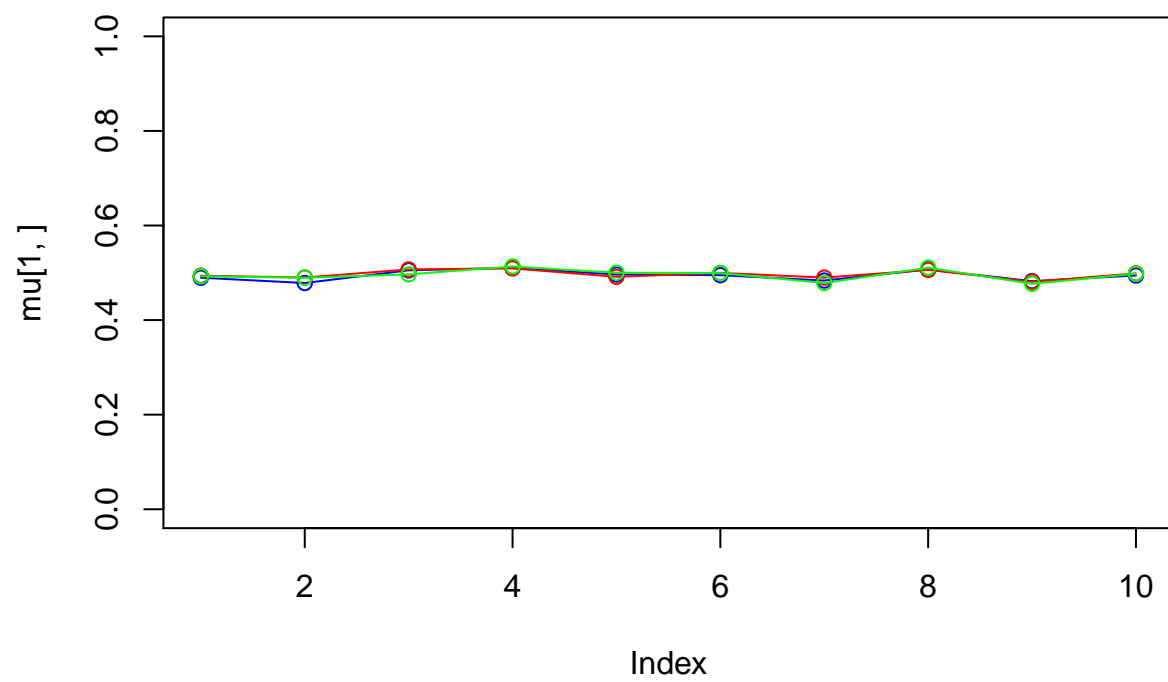
```
## [[1]]
## [1] 0.48293133107 0.51706866893 0.49007966733 0.50152946699 0.39332820172
## [6] 0.61703504228 0.59792428535 0.39113495199 0.27182539774 0.69957245640
## [11] 0.70178609664 0.33474375189 0.19877624614 0.76586485831 0.81088543925
## [16] 0.22897928765 0.08857255800 0.89200054923 0.90675481733 0.10849562636
## [21] 0.00008952955 0.94950012035
```

K = 3

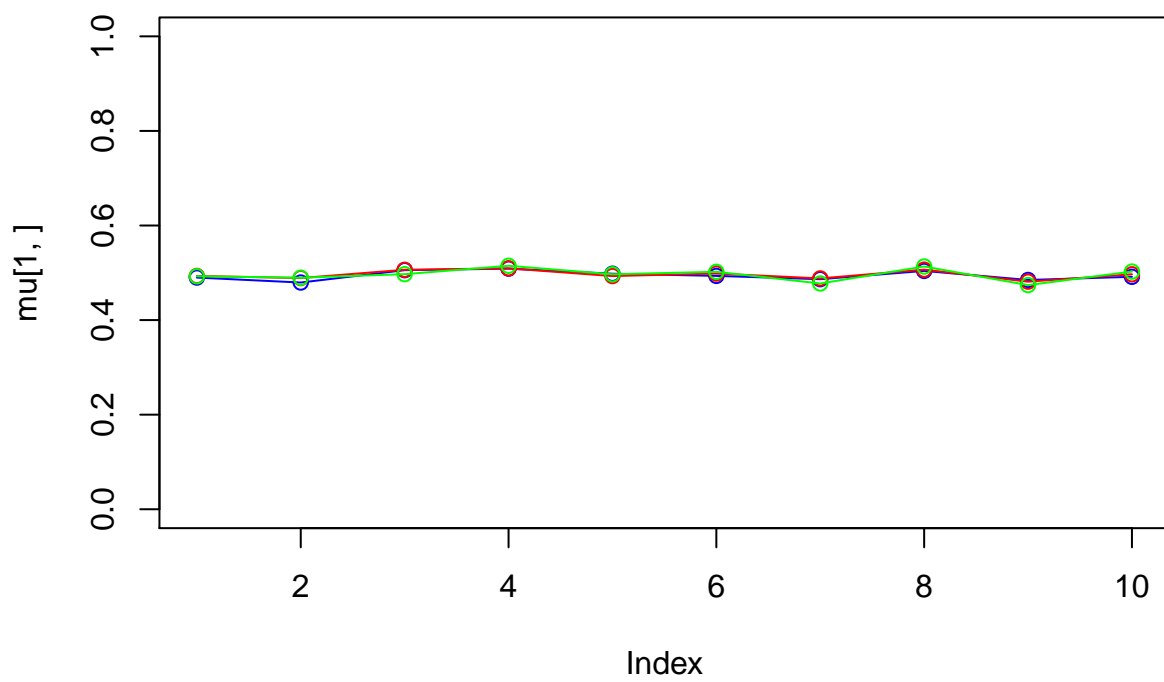
```
myem(K=3)
```



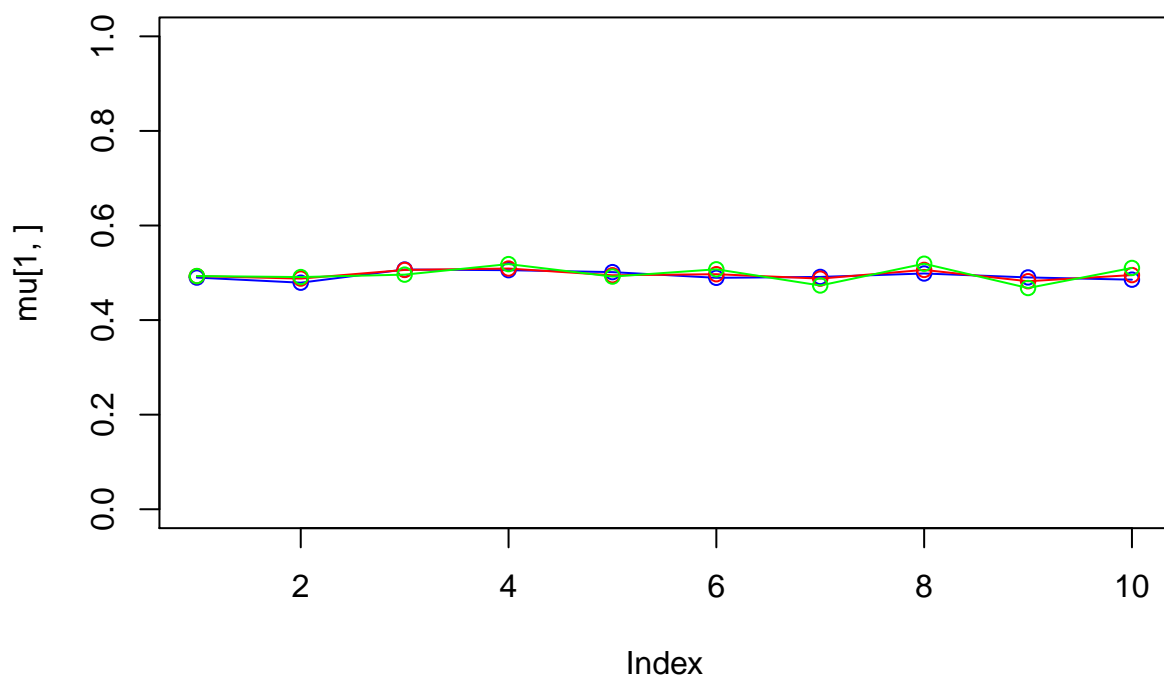
iteration: 1 log likelihood: -912.7567



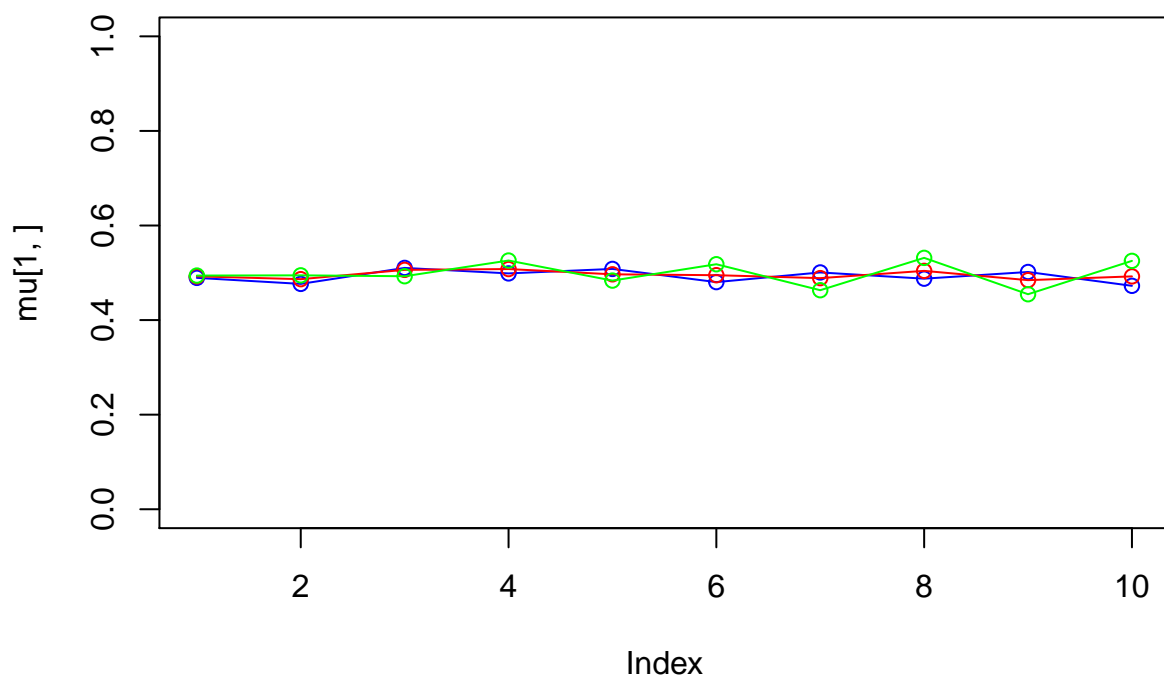
iteration: 2 log likelihood: -932.1921



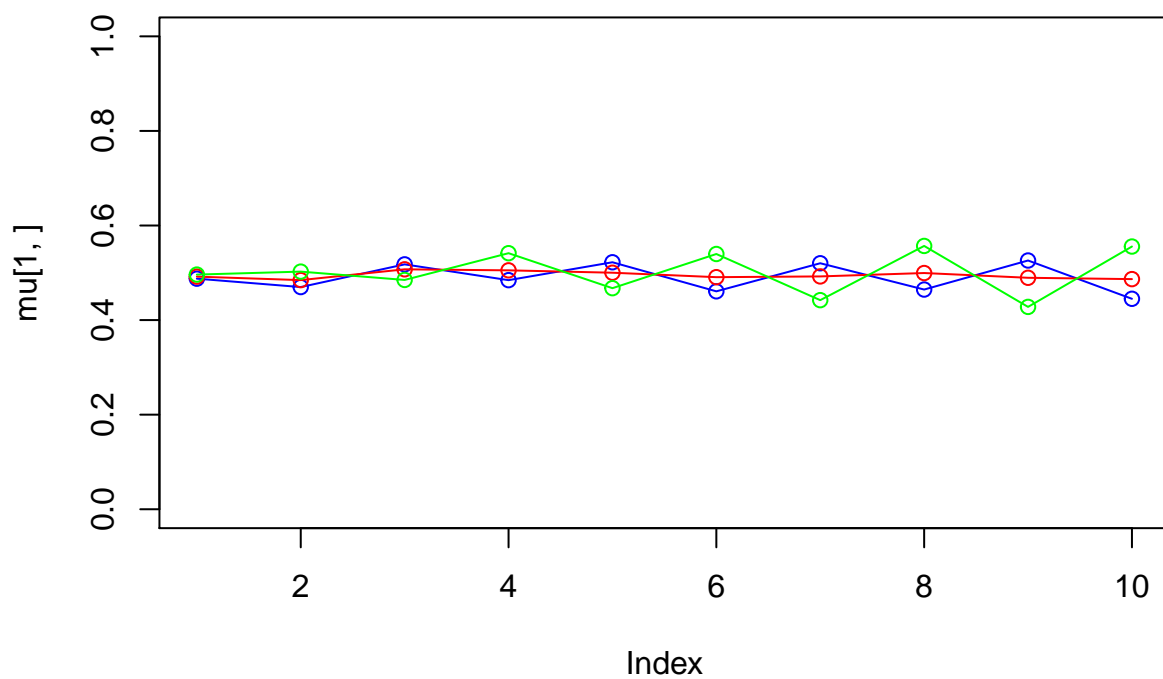
iteration: 3 log likelihood: -932.0234



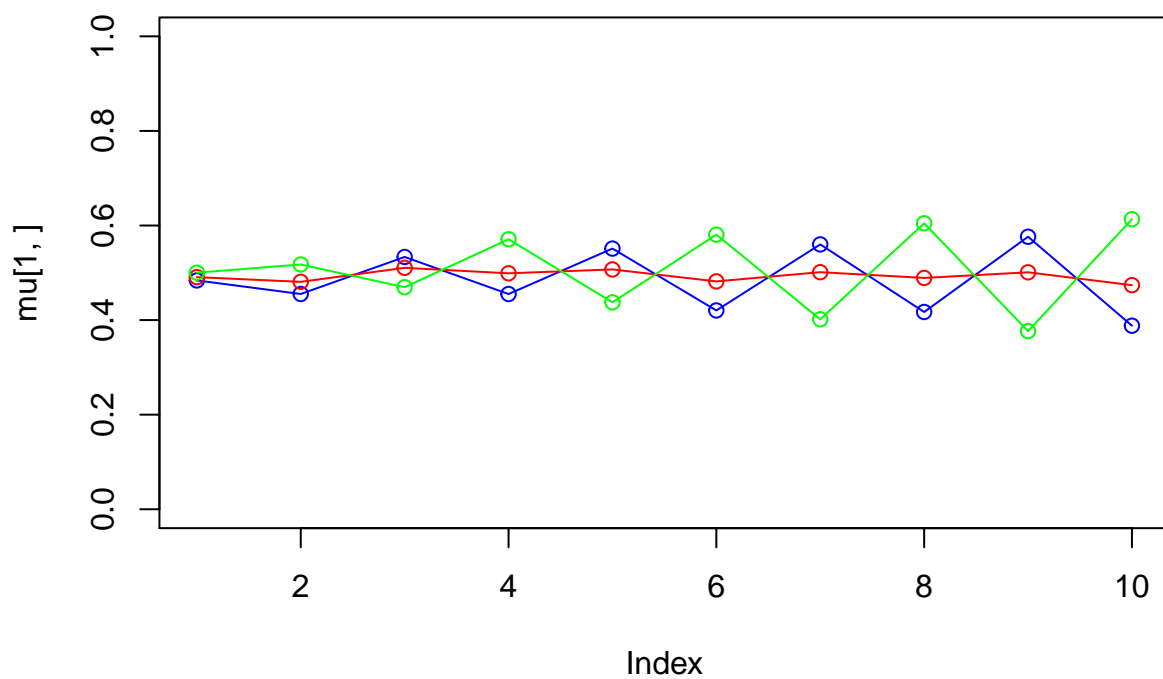
iteration: 4 log likelihood: -931.2587



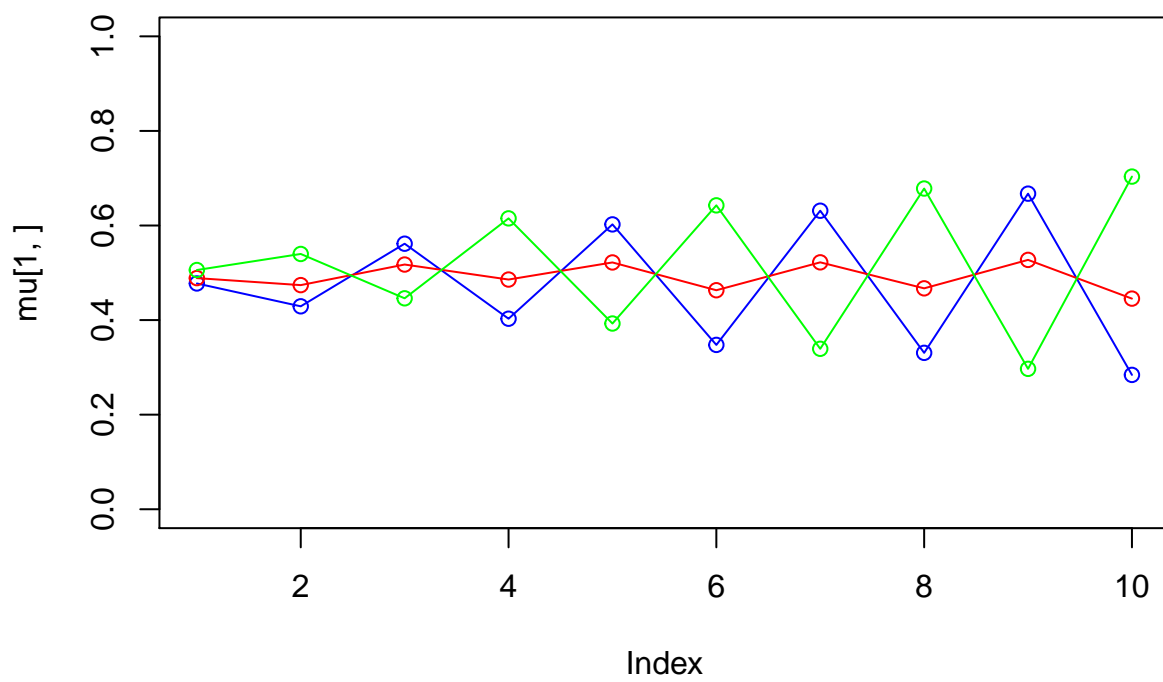
iteration: 5 log likelihood: -927.8881



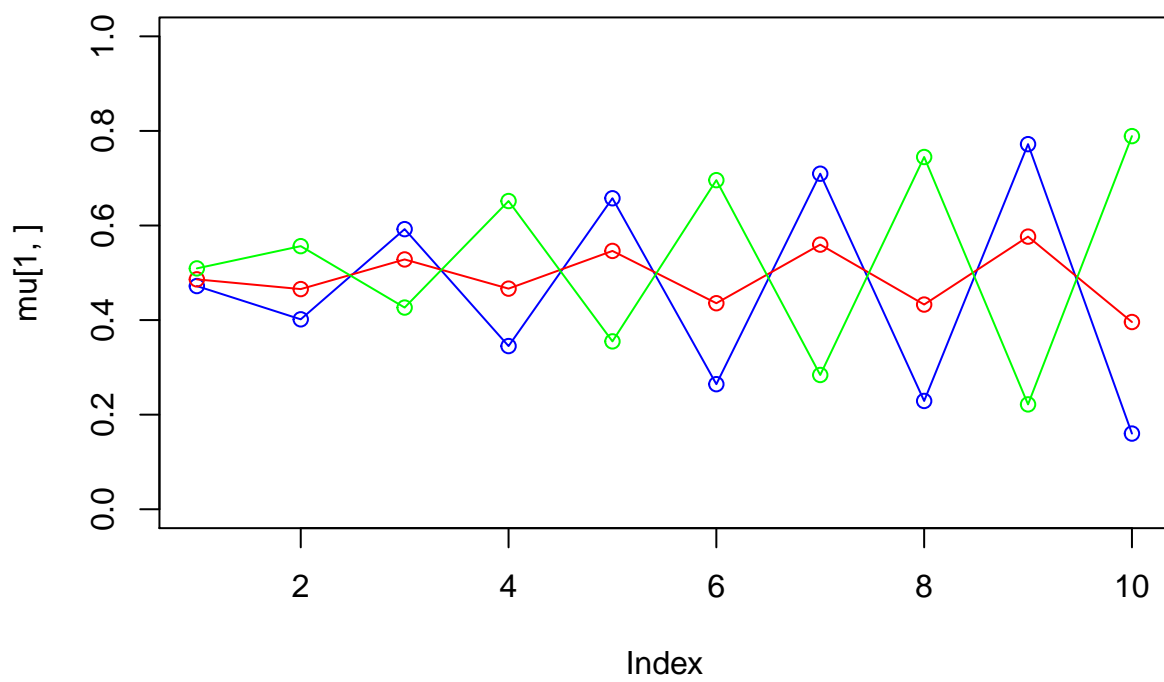
iteration: 6 log likelihood: -913.454



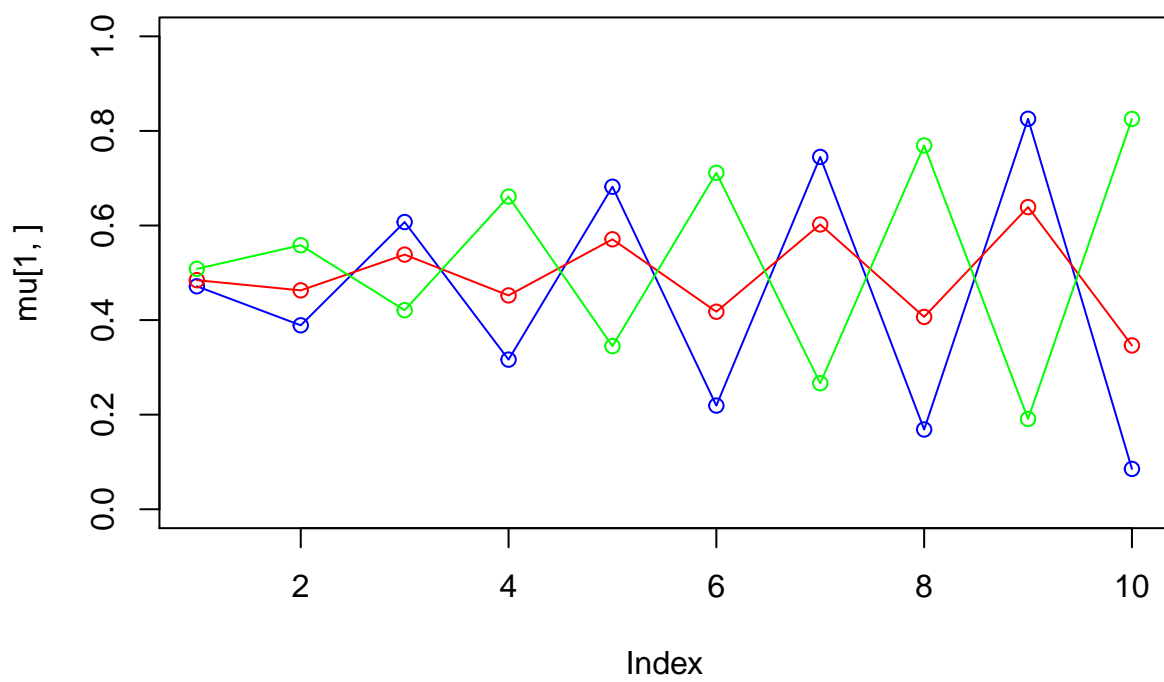
iteration: 7 log likelihood: -858.0583



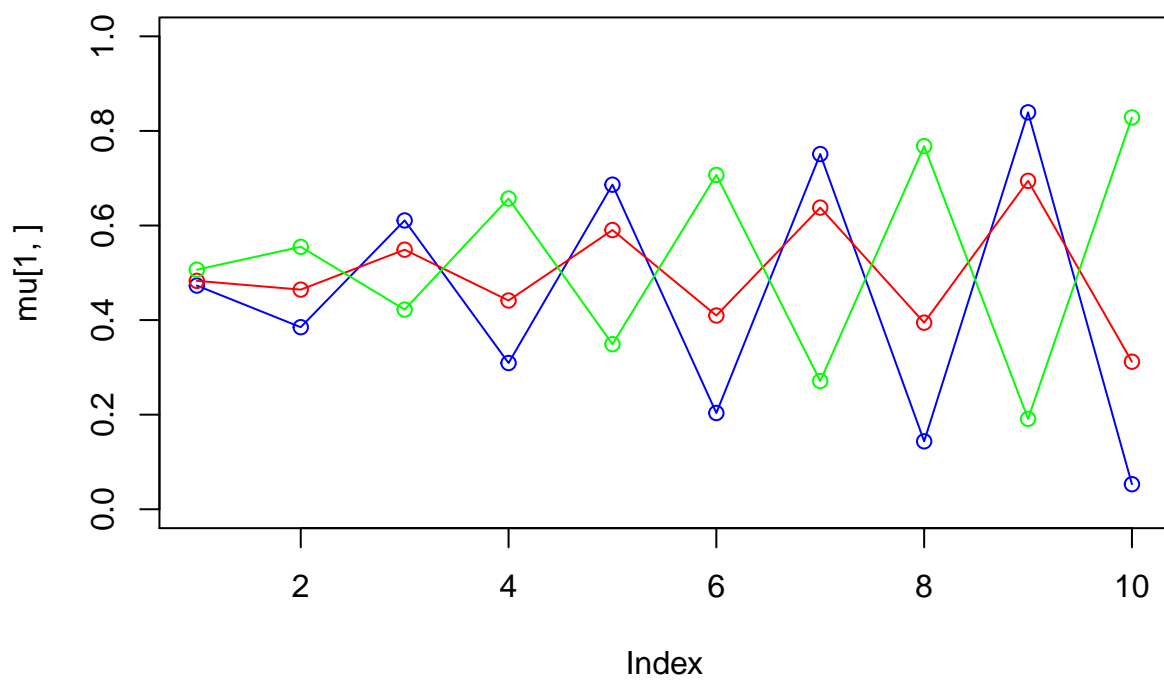
iteration: 8 log likelihood: -709.6665



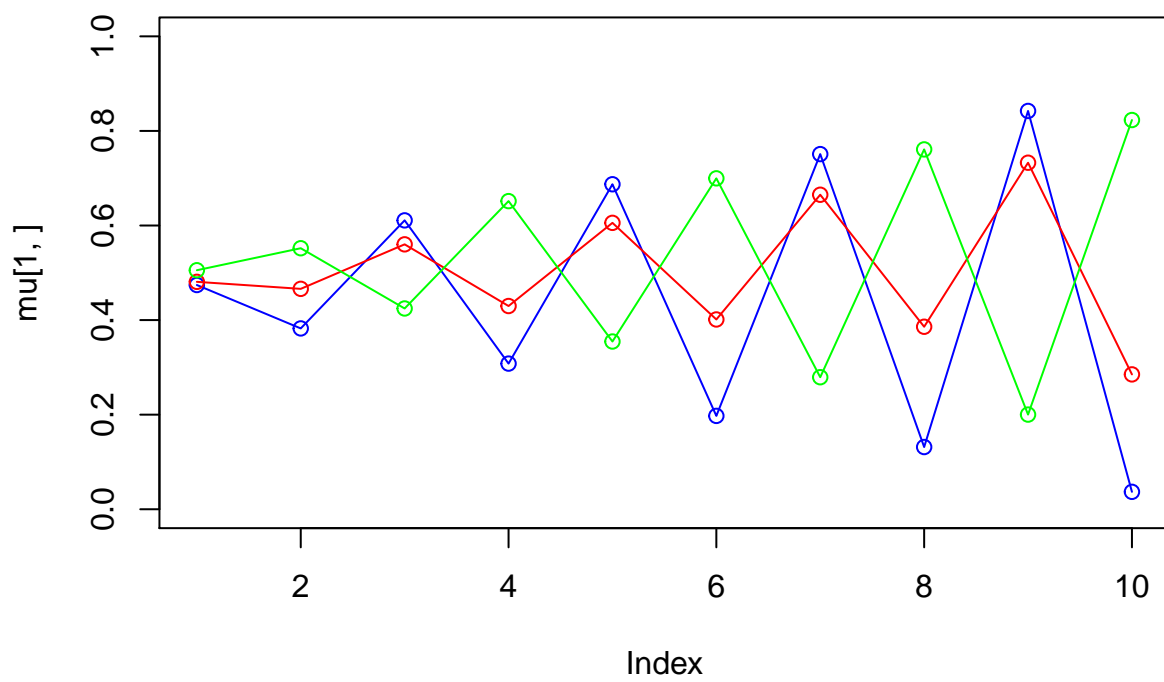
iteration: 9 log likelihood: -524.1097



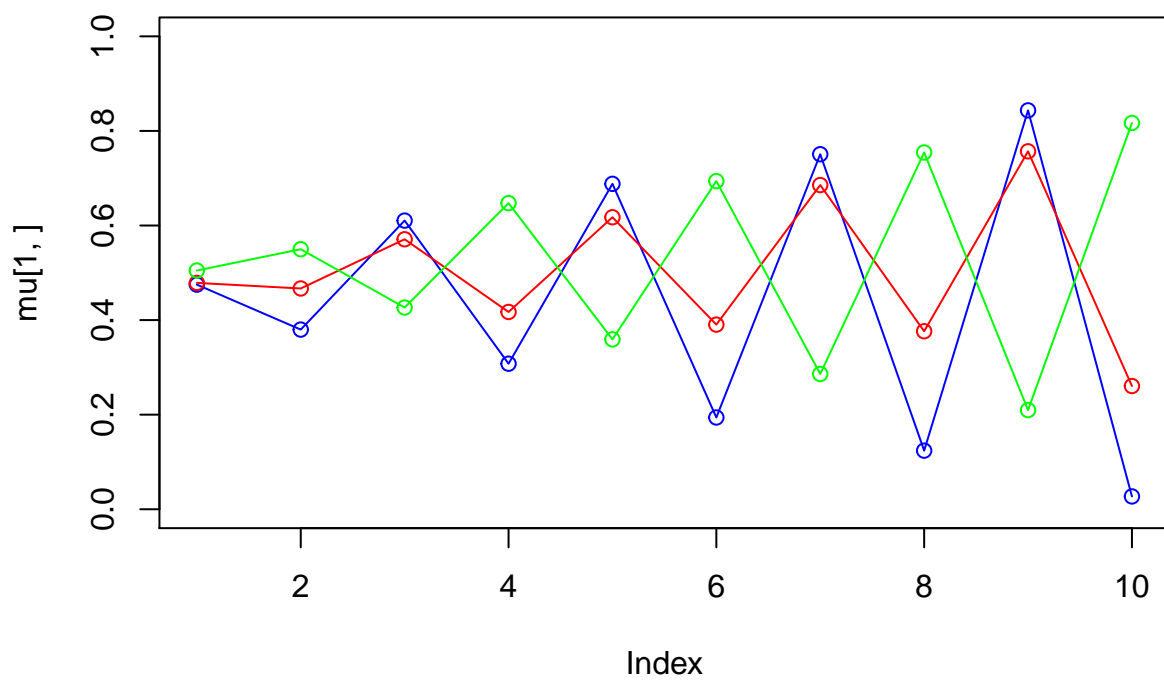
iteration: 10 log likelihood: -433.1614



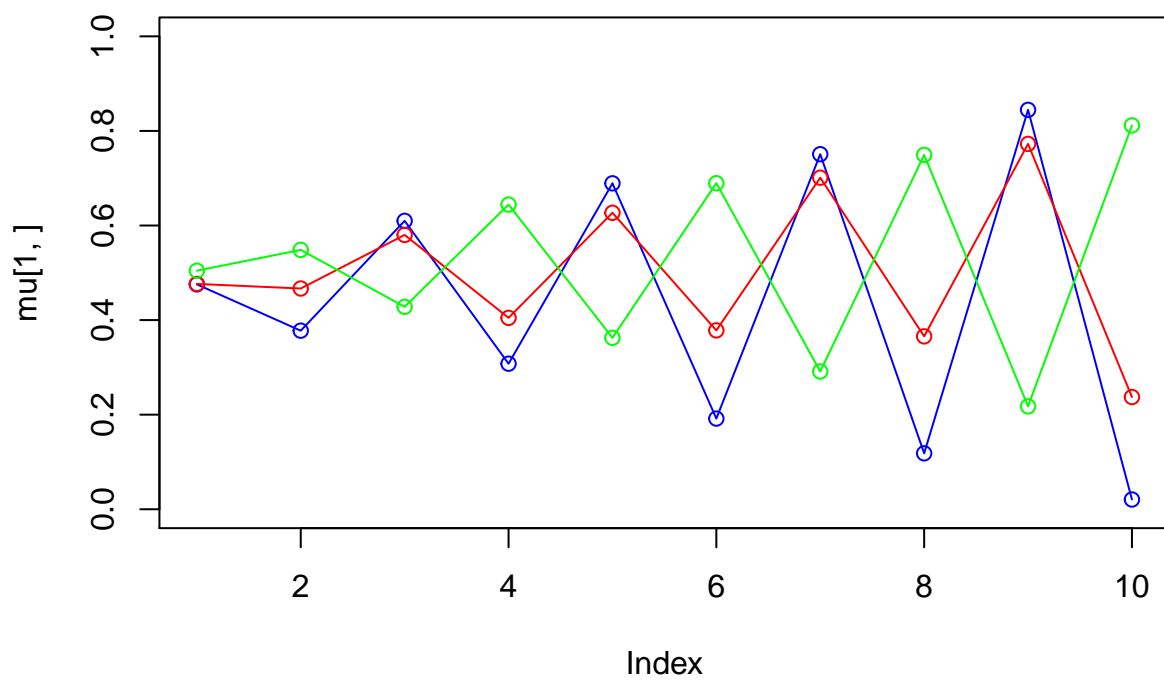
iteration: 11 log likelihood: -409.3331



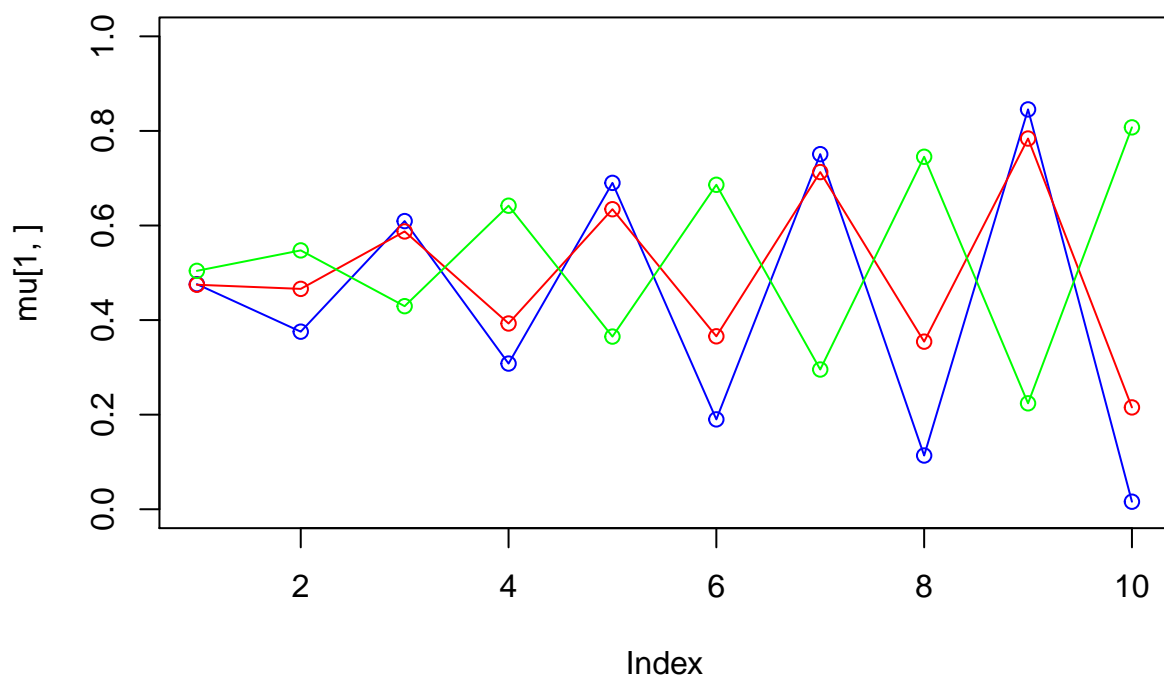
iteration: 12 log likelihood: -405.2132



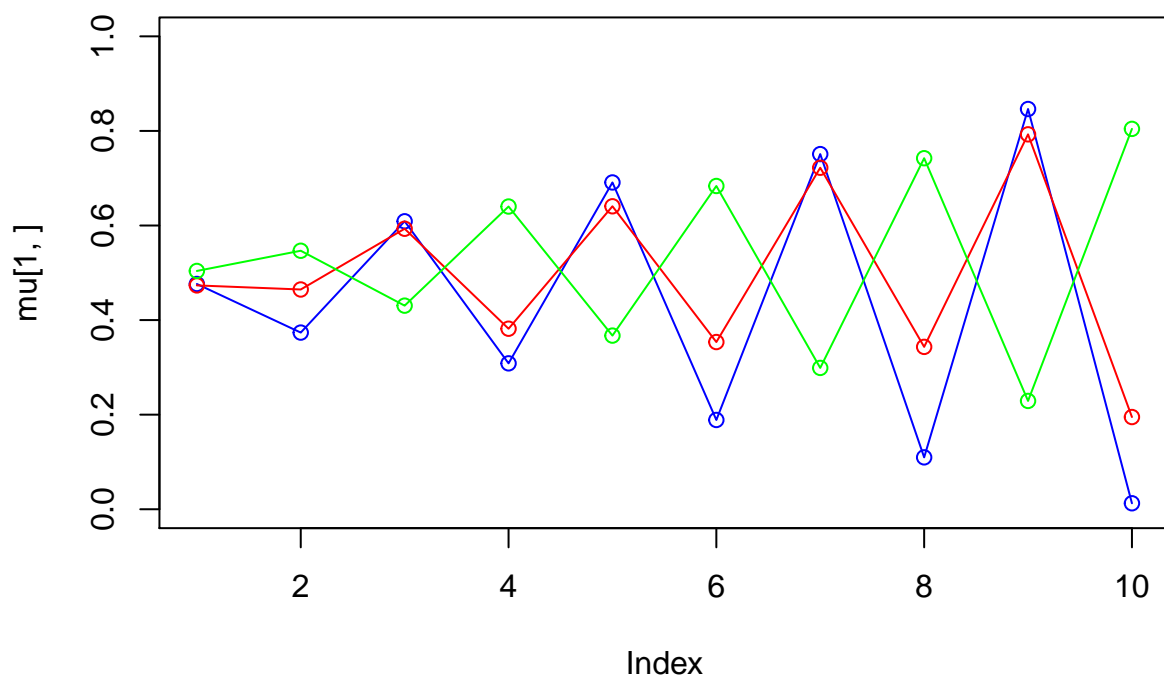
iteration: 13 log likelihood: -405.7233



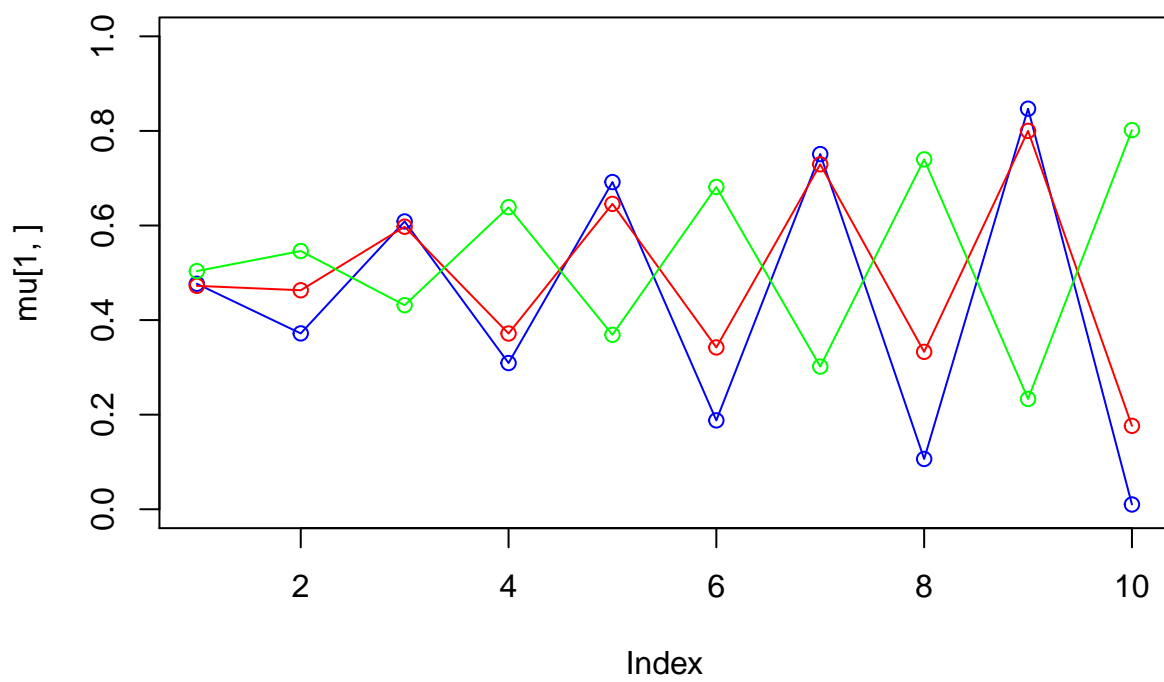
iteration: 14 log likelihood: -407.1621



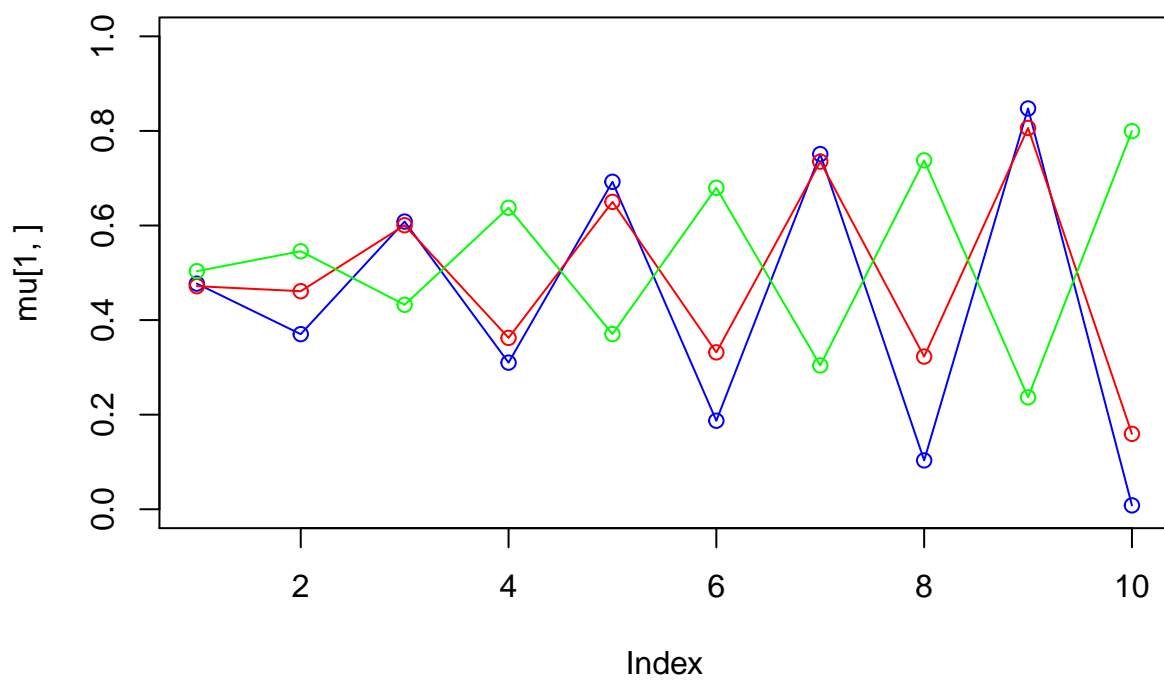
iteration: 15 log likelihood: -408.6475



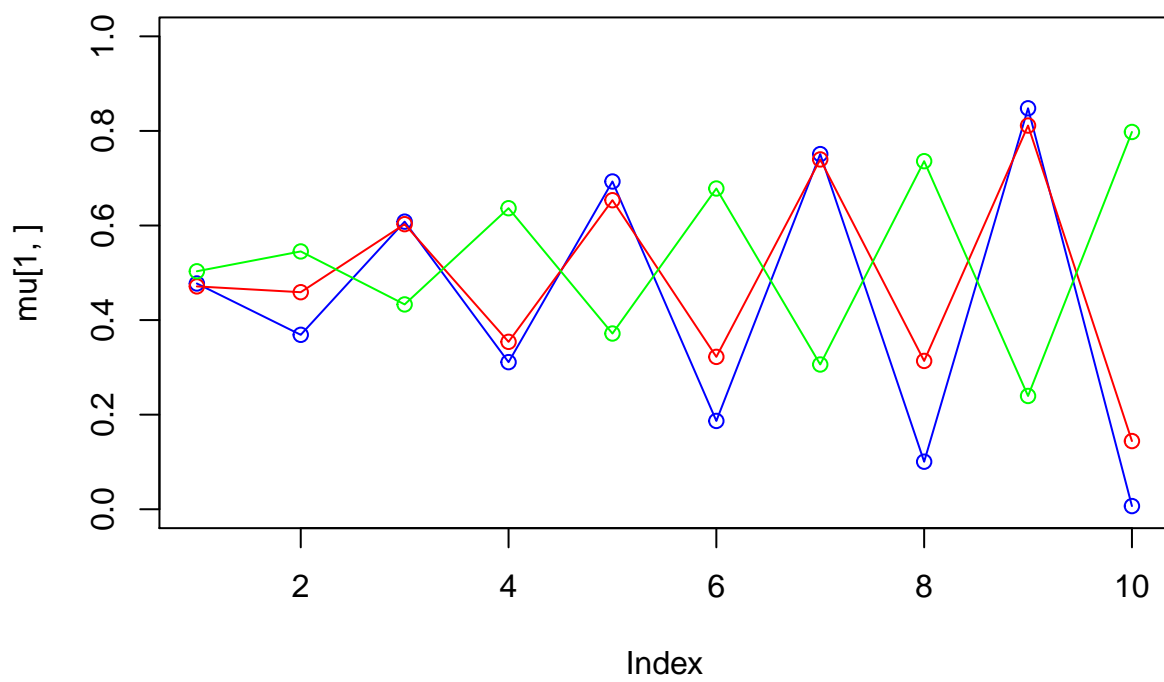
iteration: 16 log likelihood: -409.9879



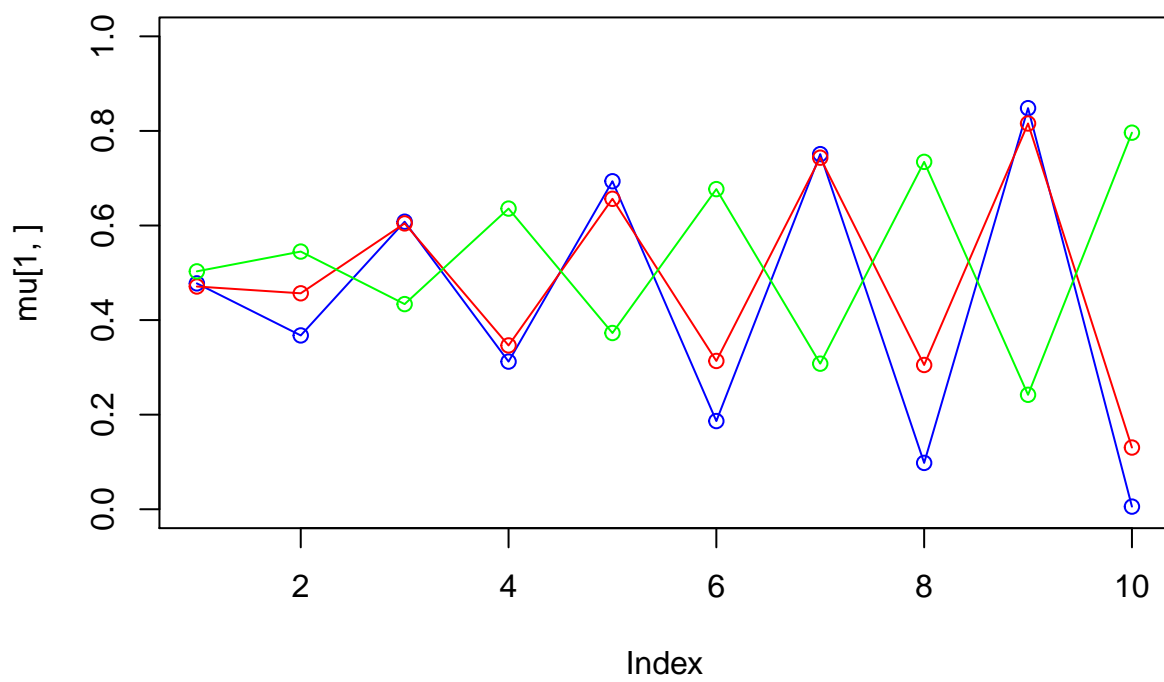
iteration: 17 log likelihood: -411.1645



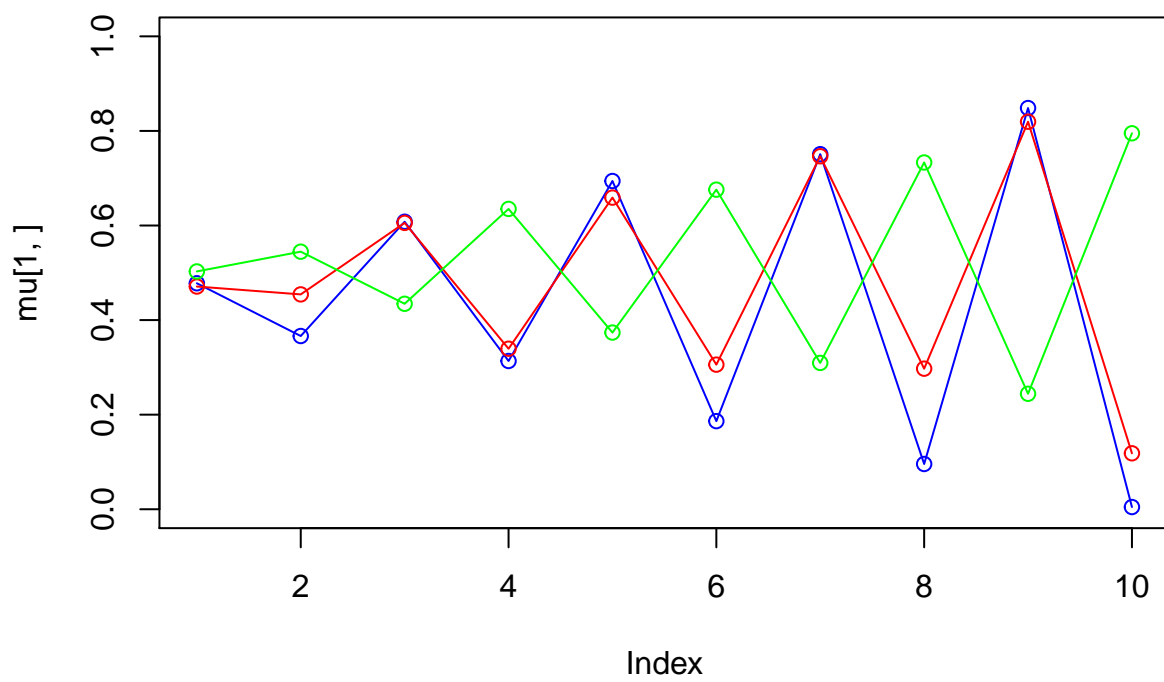
iteration: 18 log likelihood: -412.1979



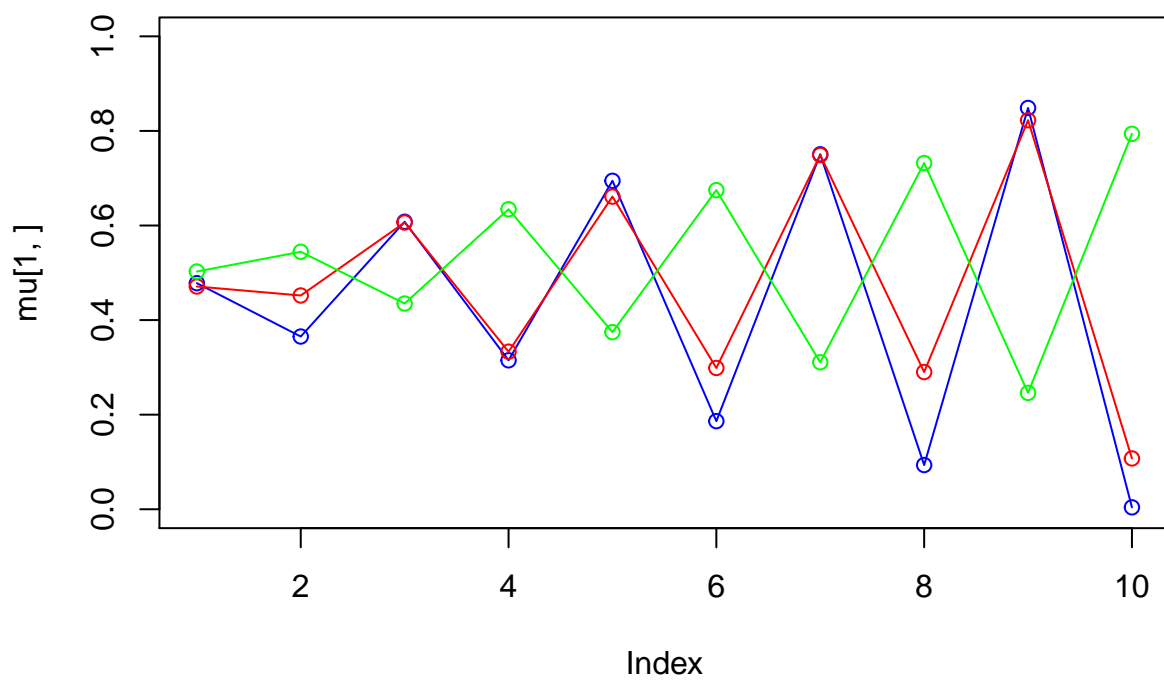
iteration: 19 log likelihood: -413.1139



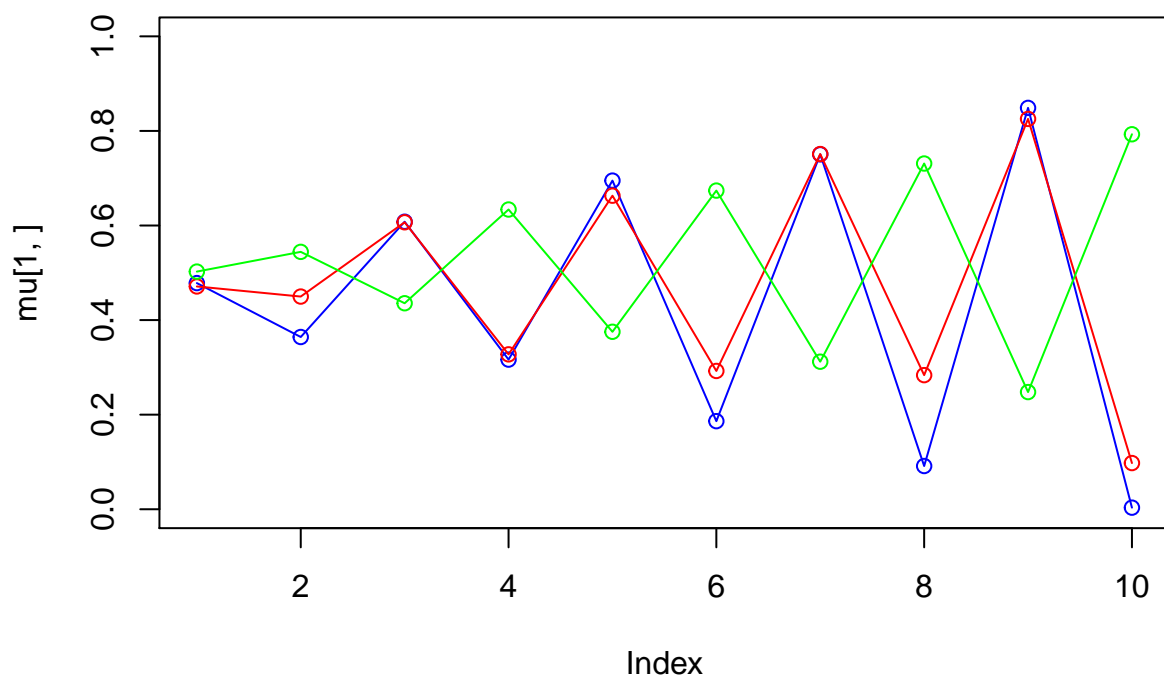
iteration: 20 log likelihood: -413.934



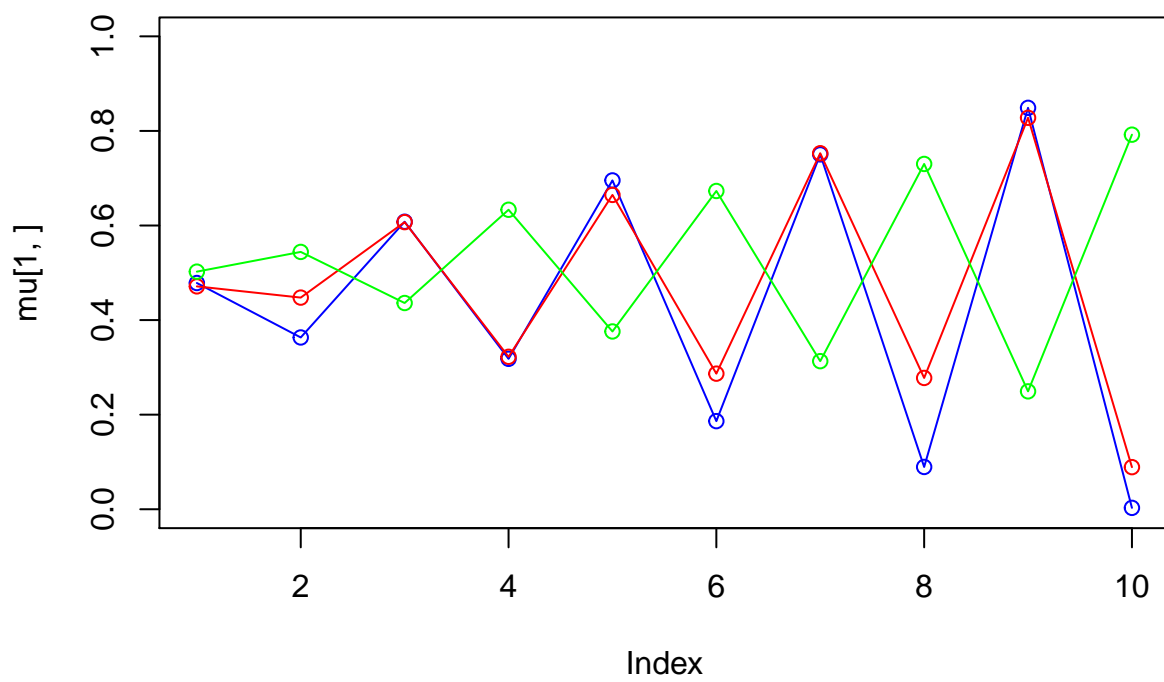
iteration: 21 log likelihood: -414.675



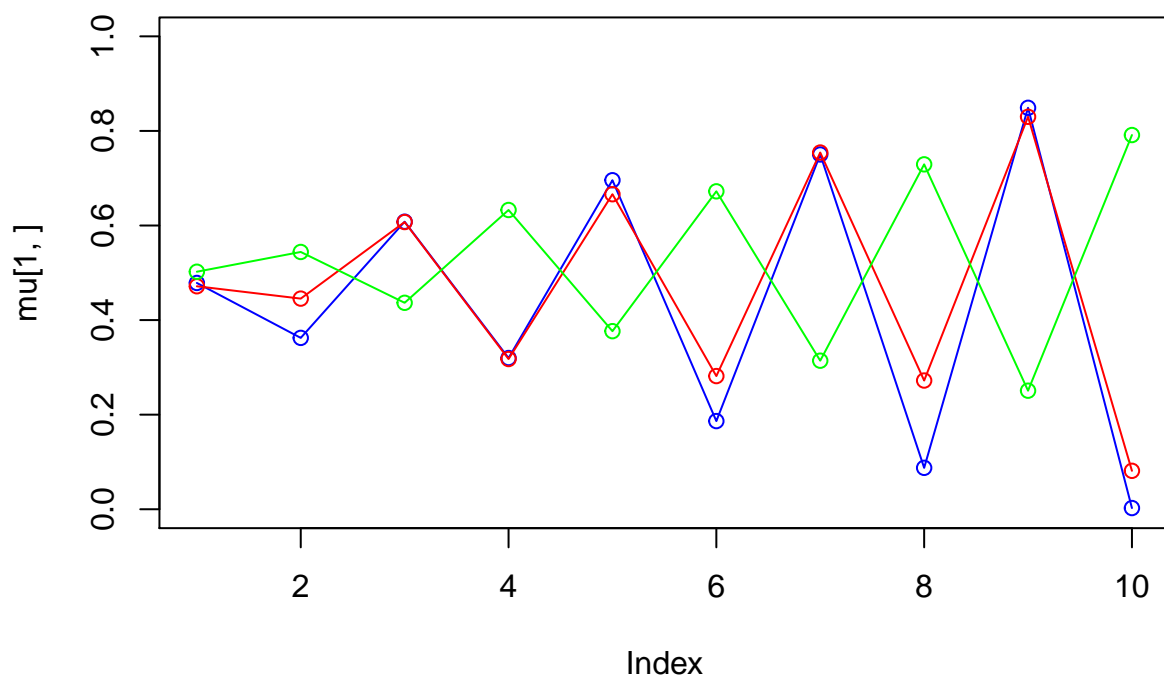
iteration: 22 log likelihood: -415.3492



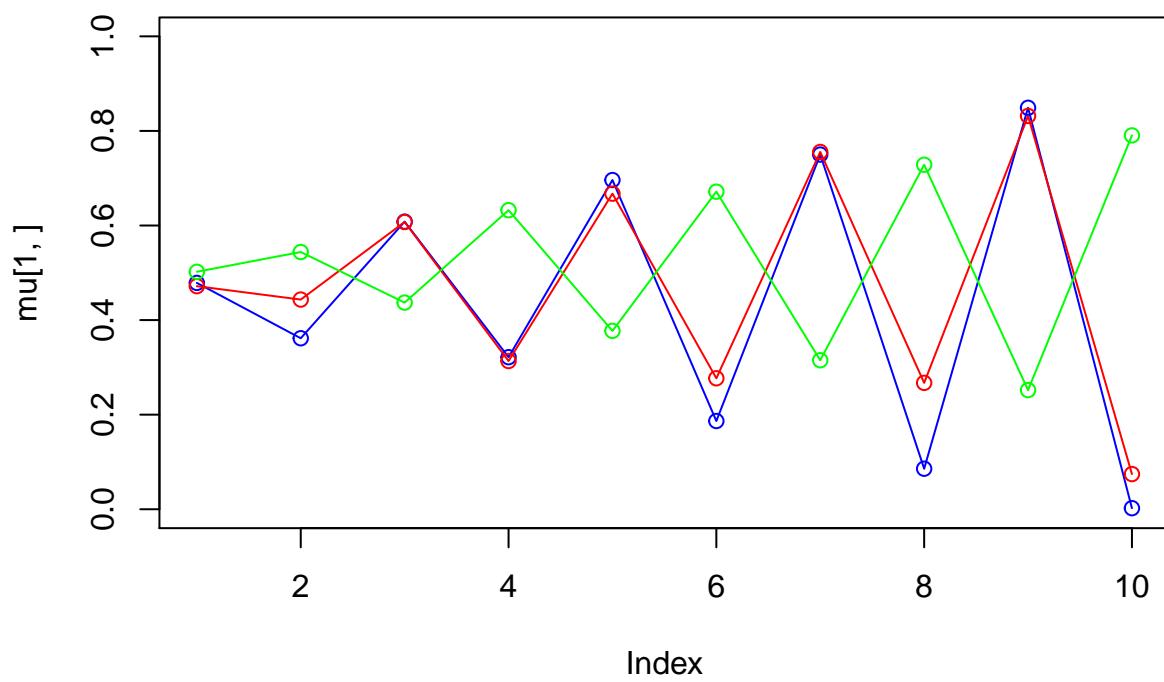
iteration: 23 log likelihood: -415.9659



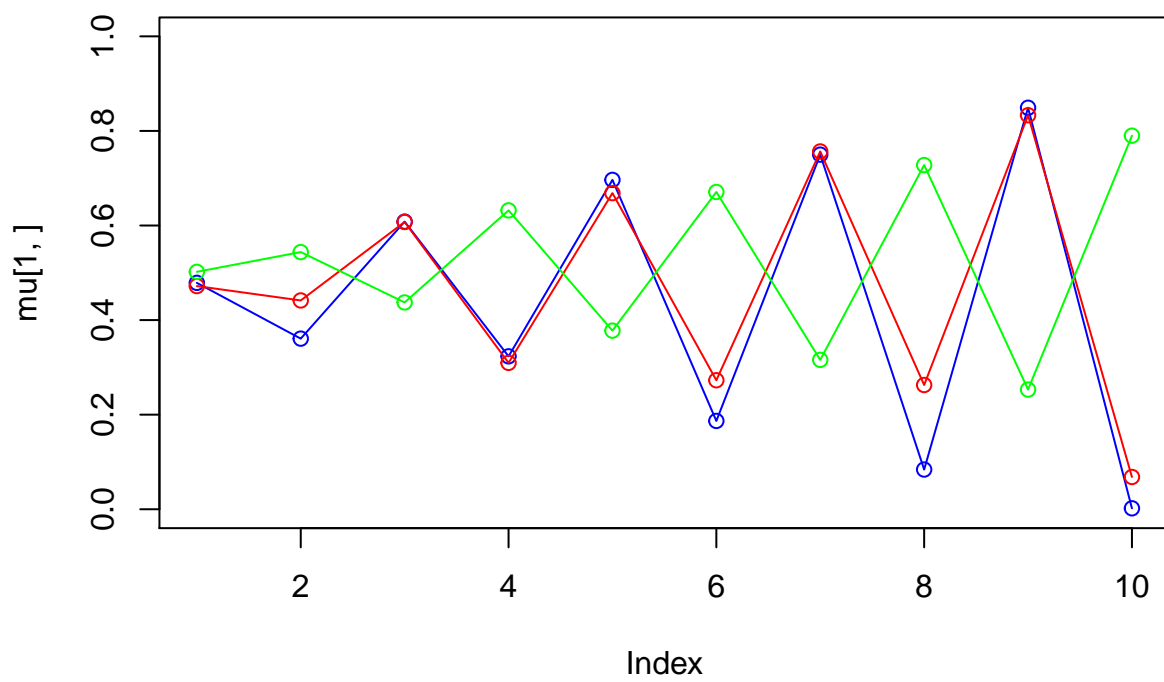
iteration: 24 log likelihood: -416.532



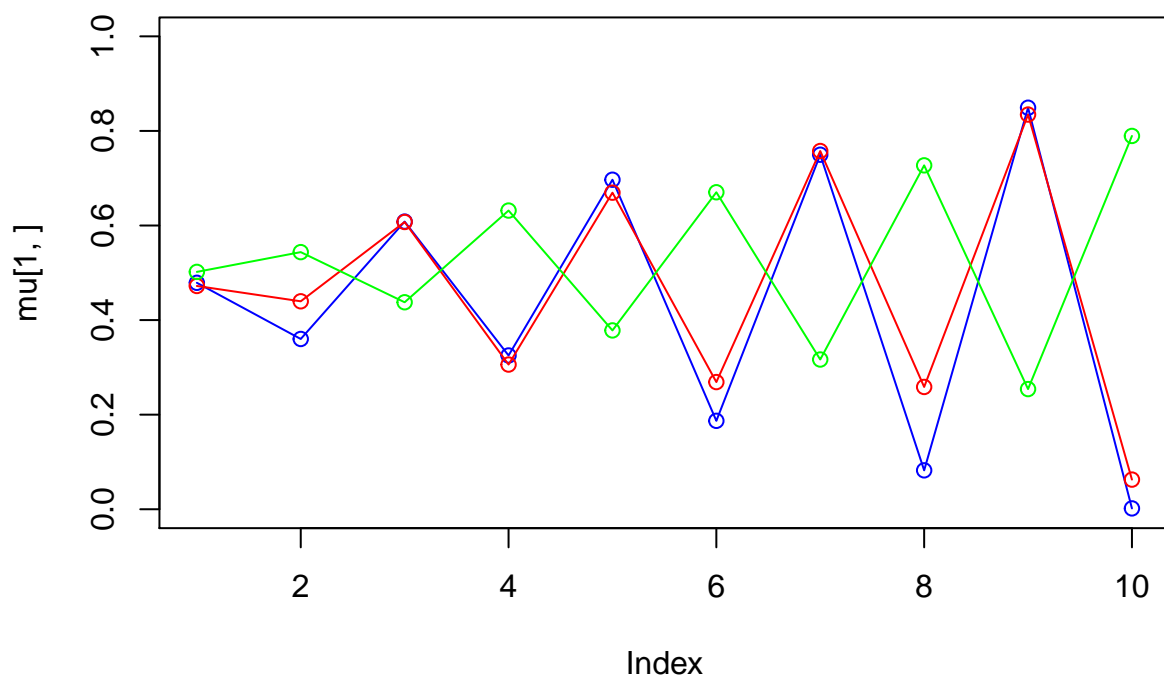
iteration: 25 log likelihood: -417.0528



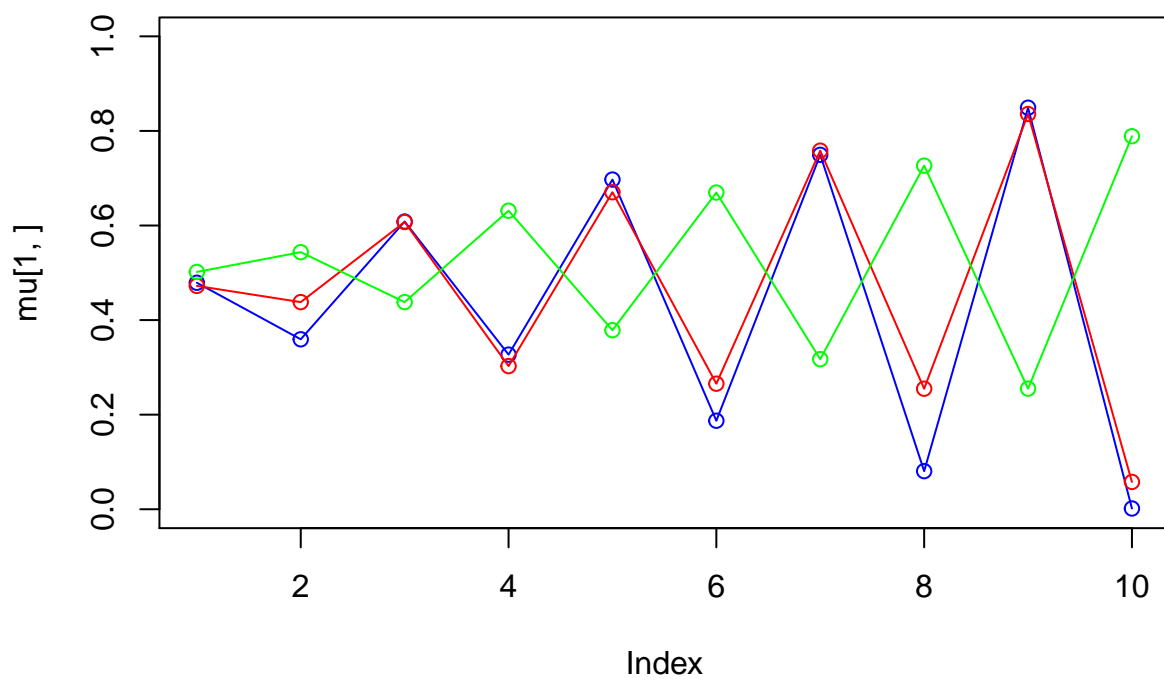
iteration: 26 log likelihood: -417.5328



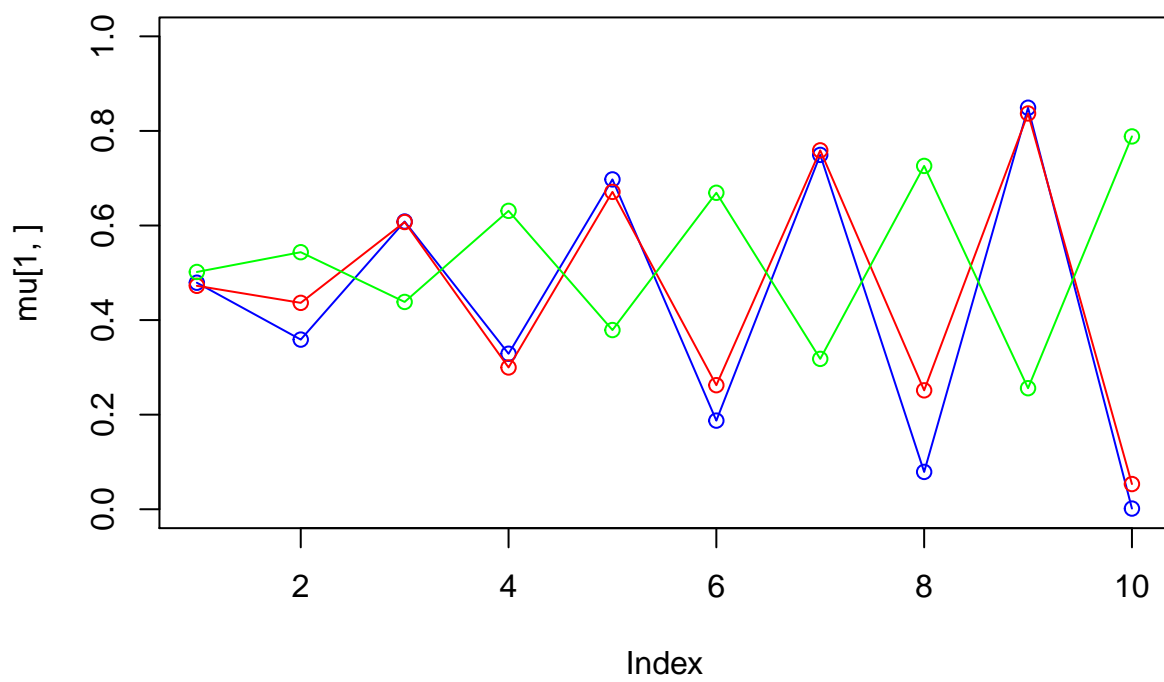
iteration: 27 log likelihood: -417.9753



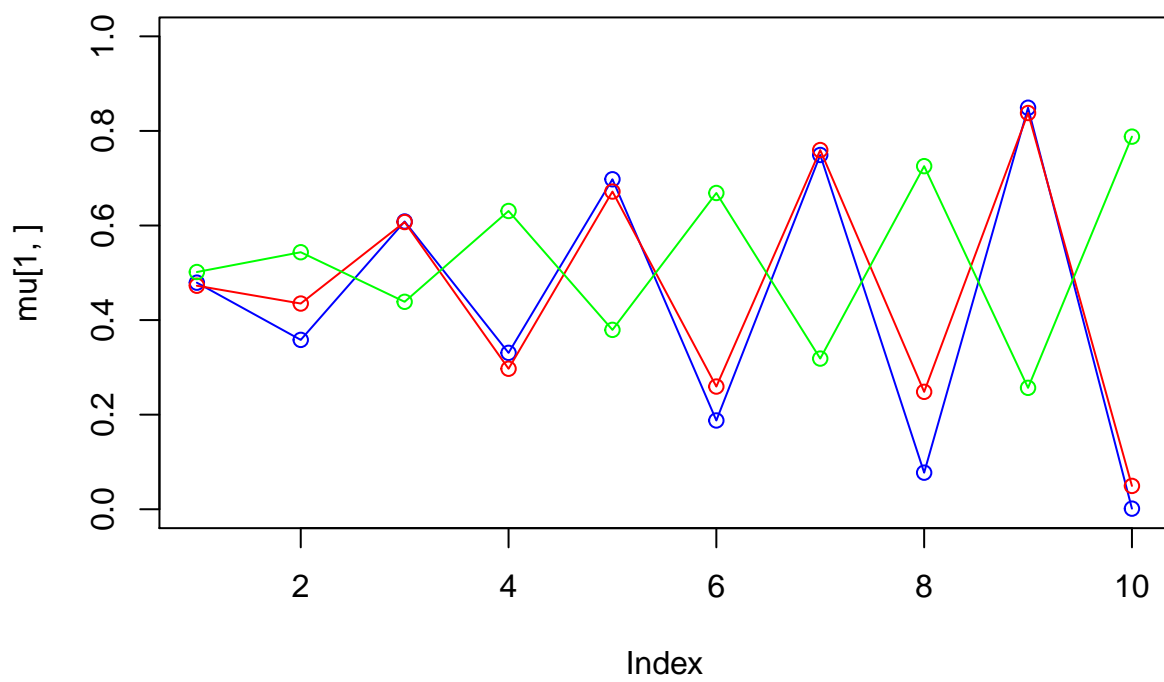
iteration: 28 log likelihood: -418.3836



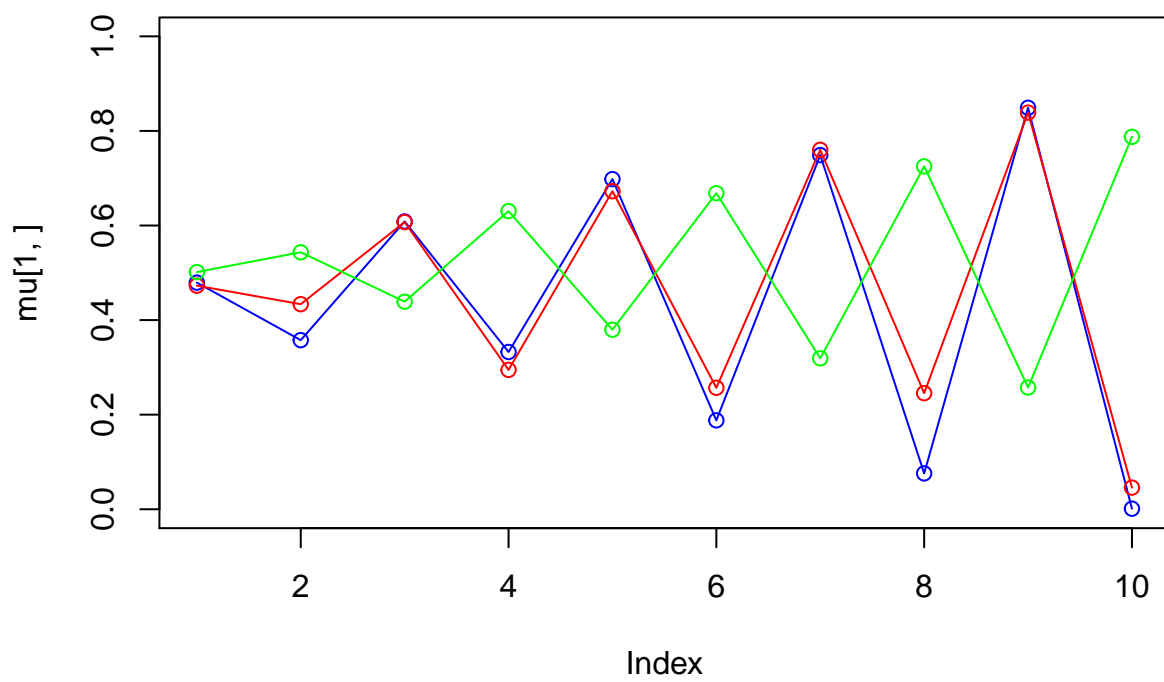
iteration: 29 log likelihood: -418.7601



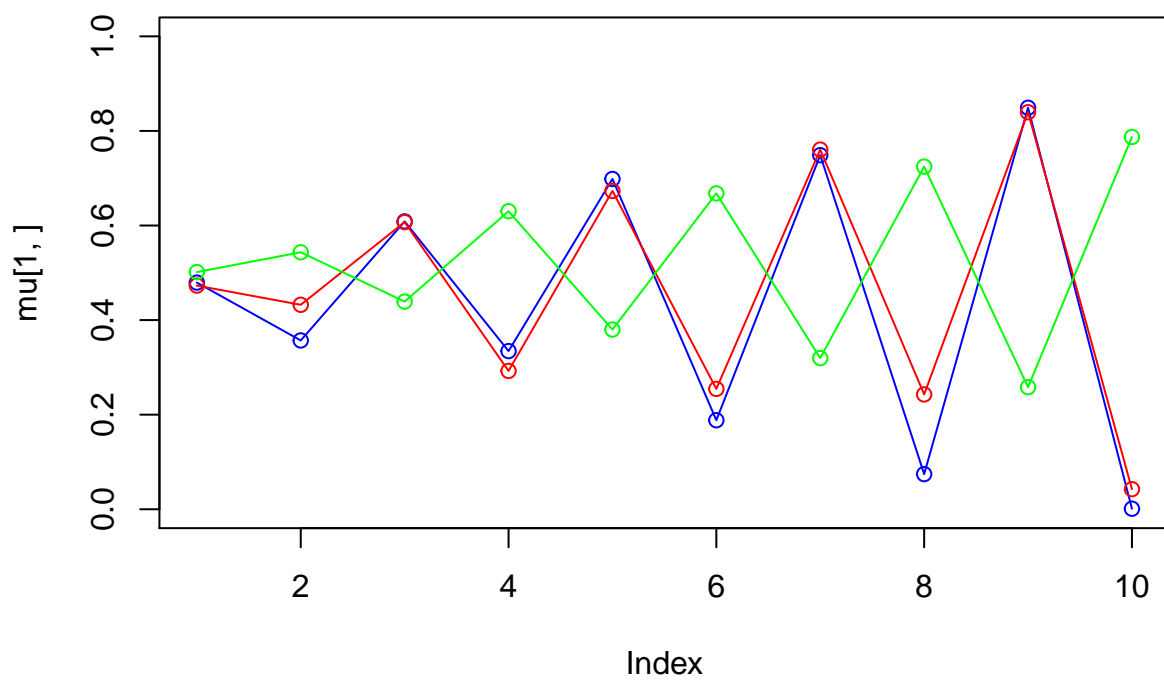
iteration: 30 log likelihood: -419.1074



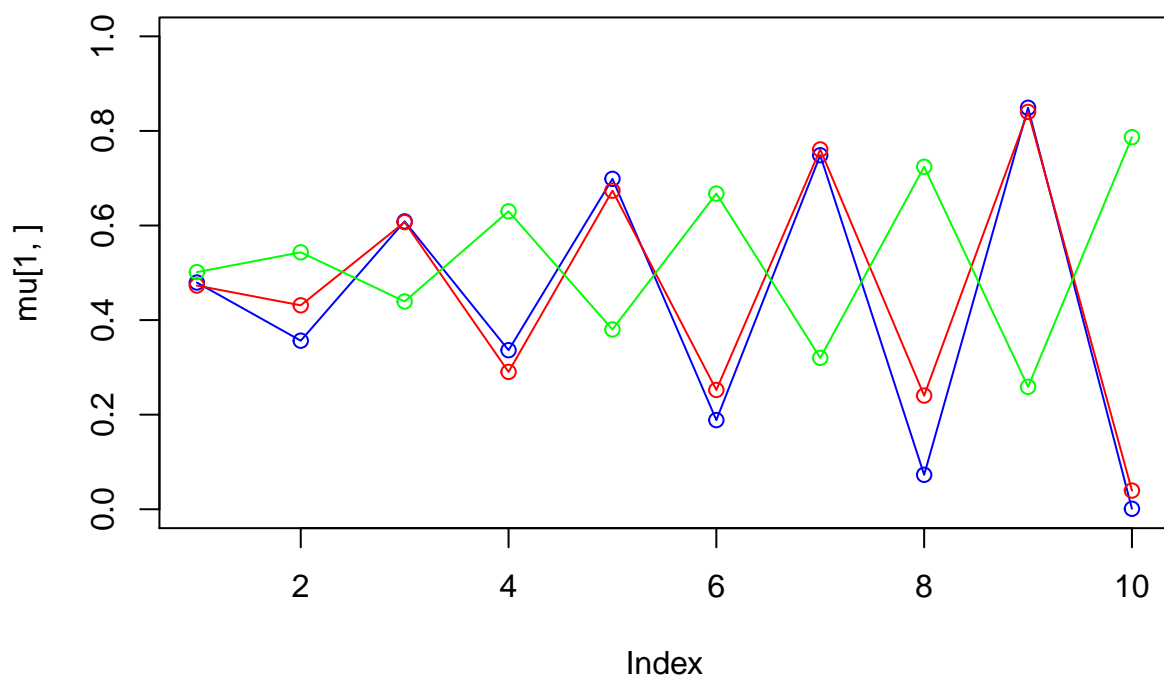
iteration: 31 log likelihood: -419.4277



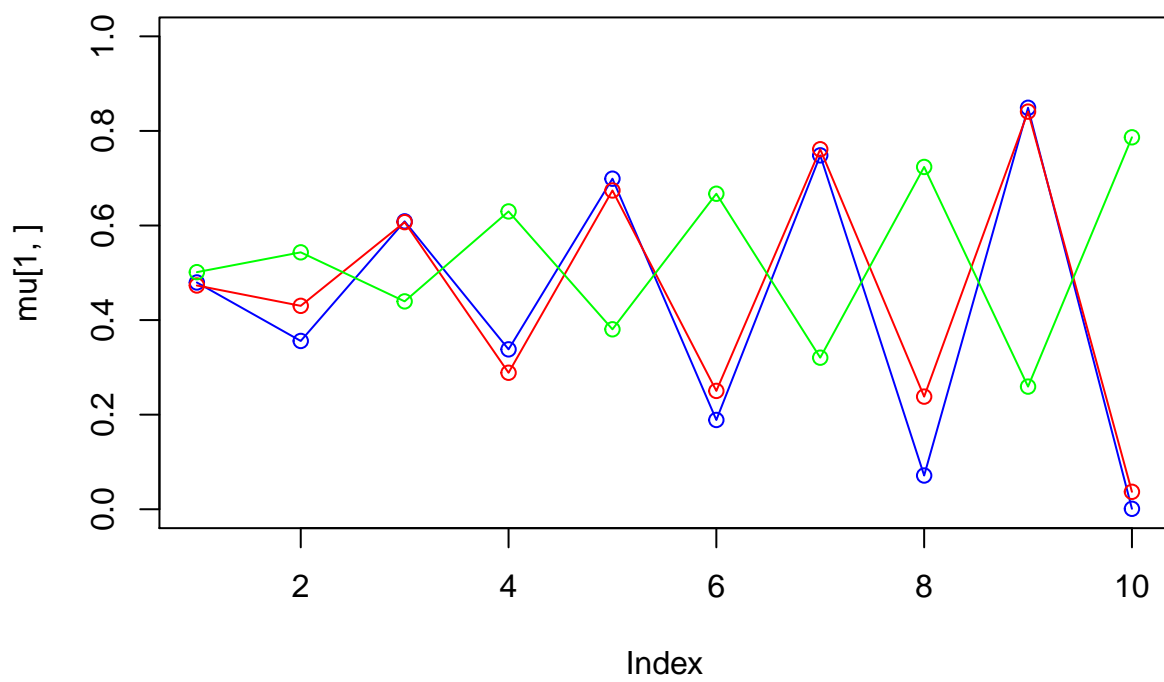
iteration: 32 log likelihood: -419.7229



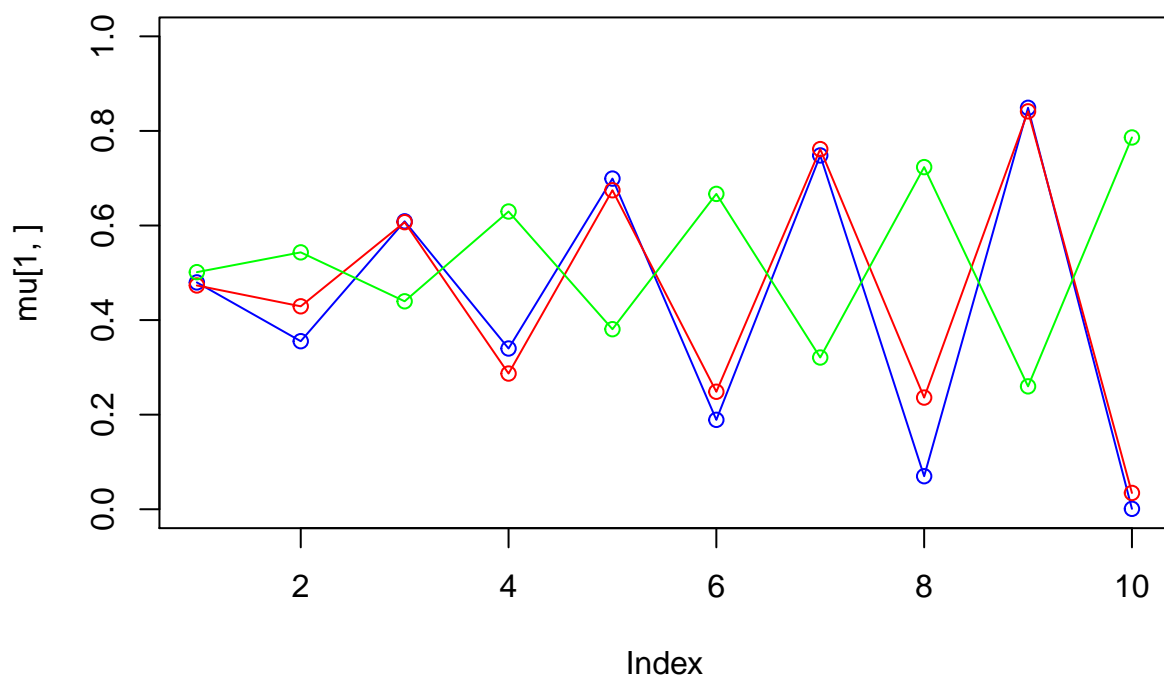
iteration: 33 log likelihood: -419.995



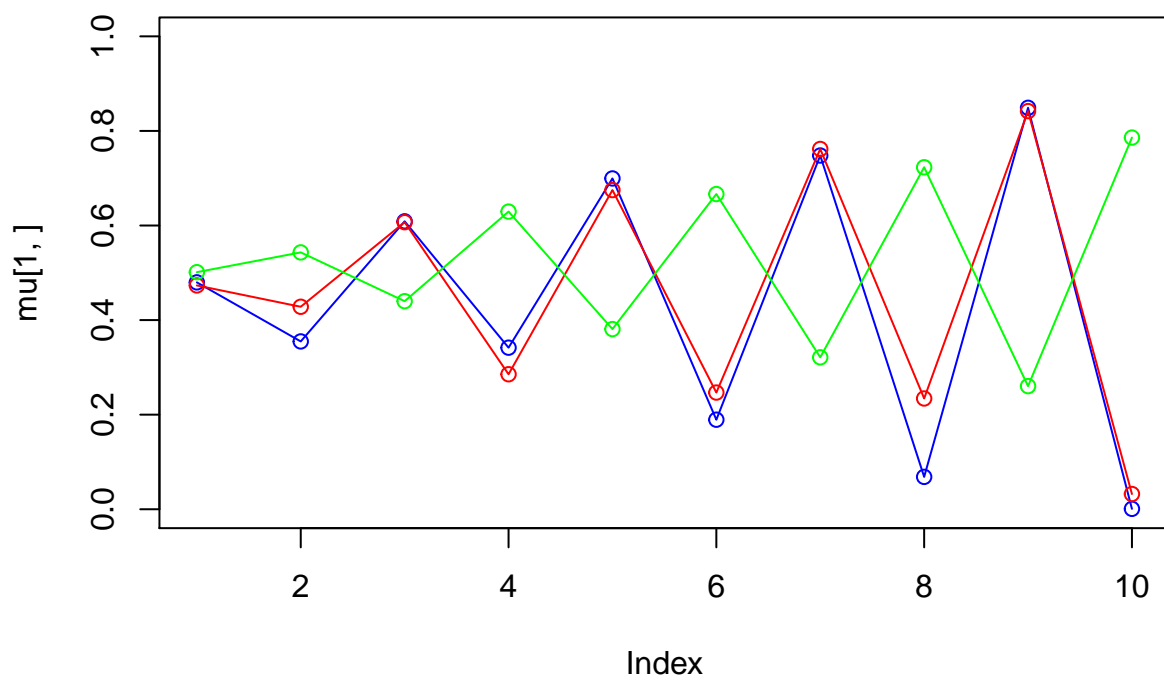
iteration: 34 log likelihood: -420.2457



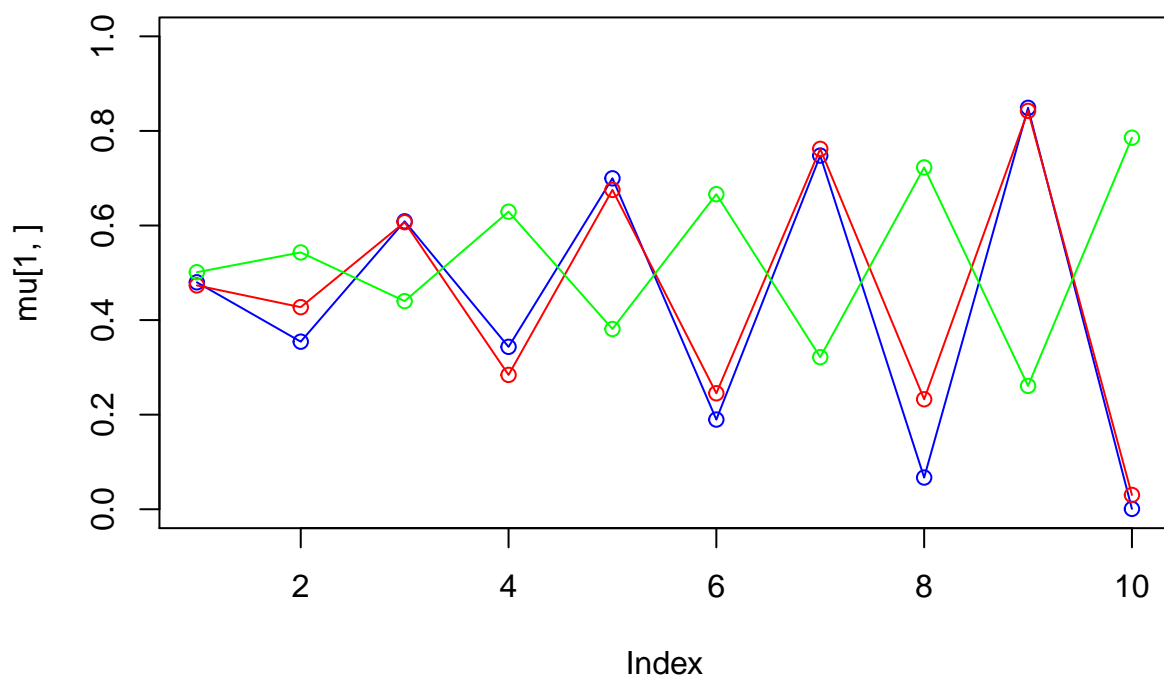
iteration: 35 log likelihood: -420.4767



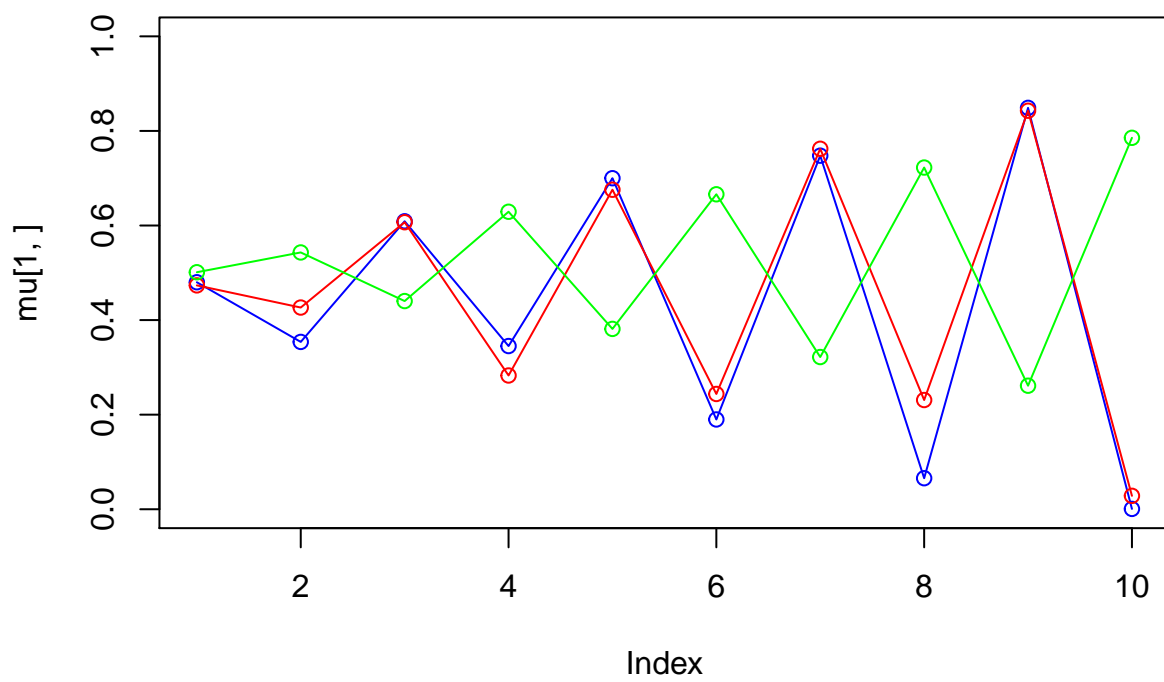
iteration: 36 log likelihood: -420.6895



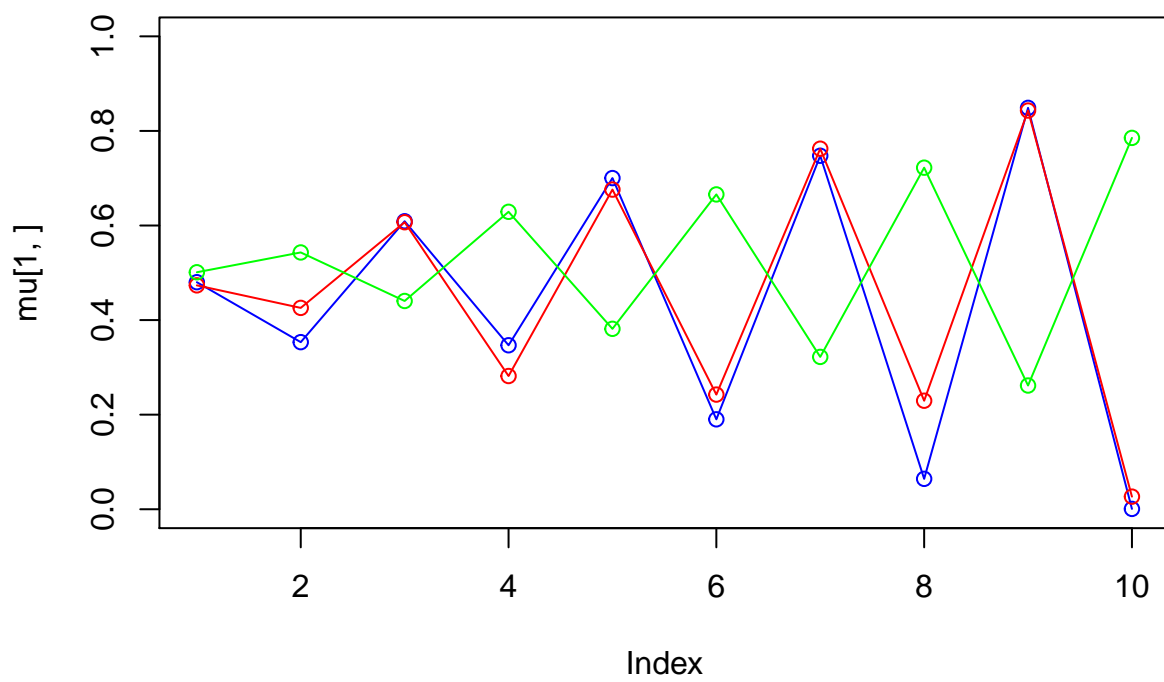
iteration: 37 log likelihood: -420.8856



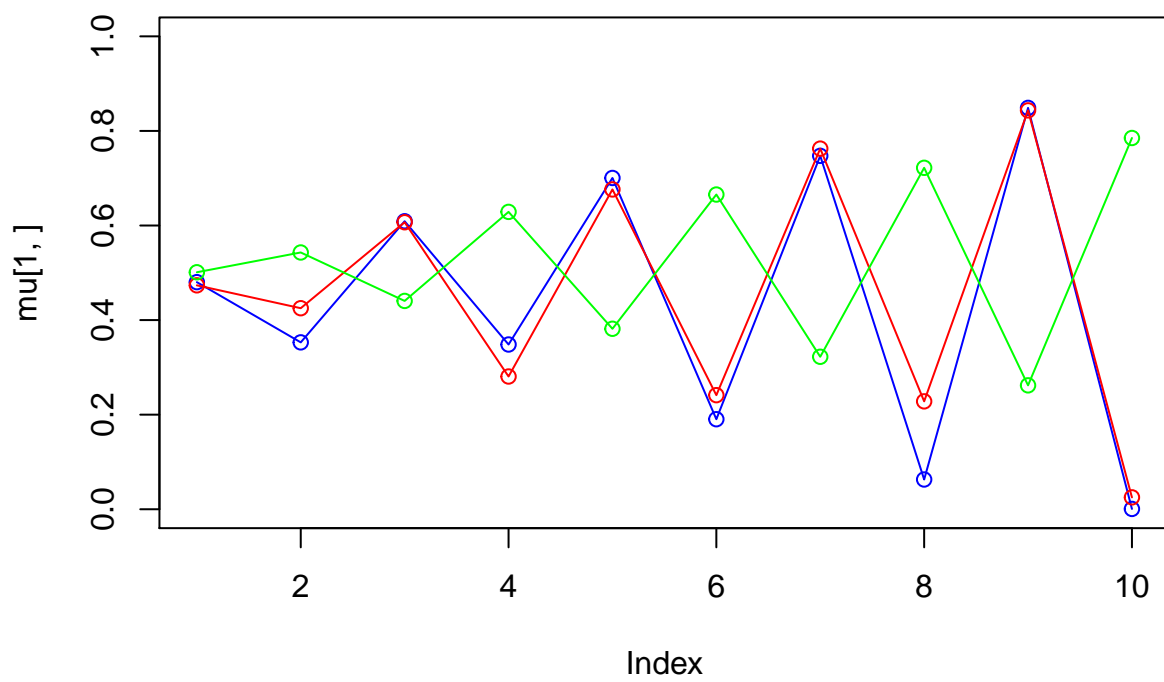
iteration: 38 log likelihood: -421.0663



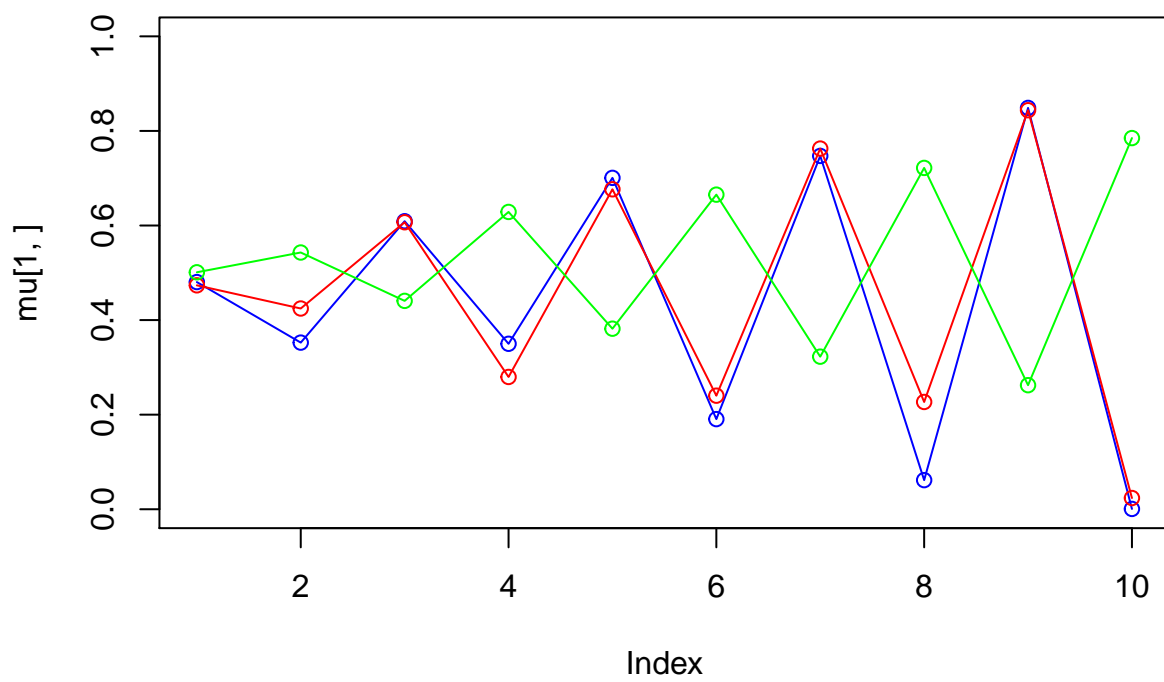
iteration: 39 log likelihood: -421.2329



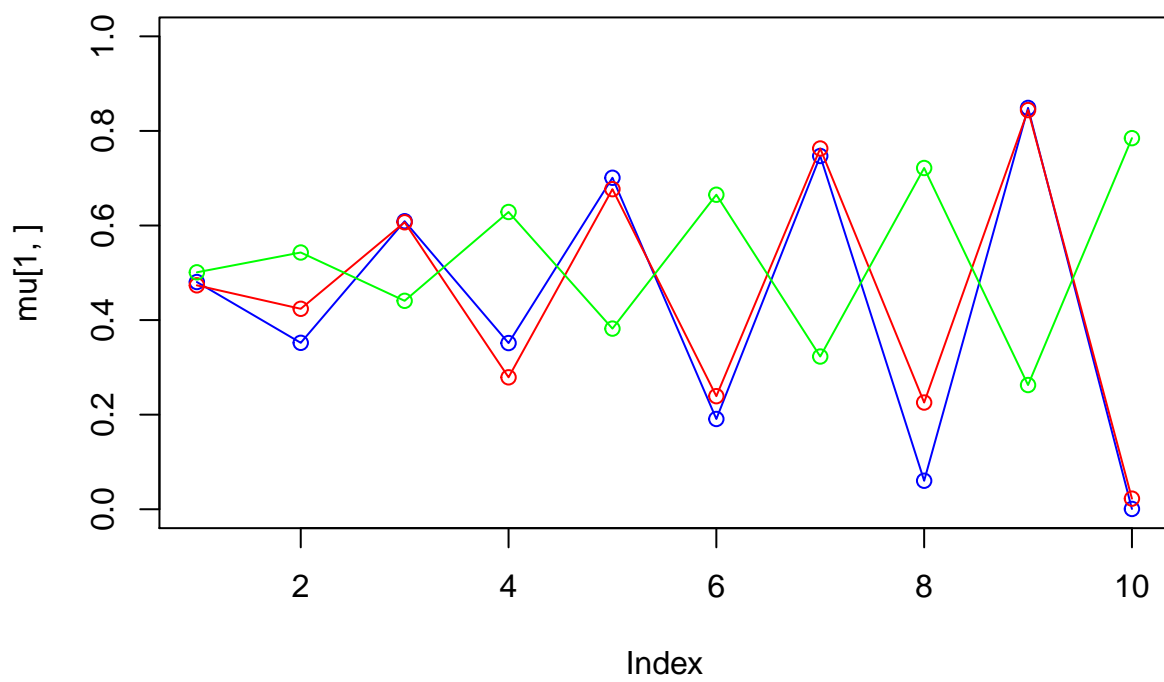
iteration: 40 log likelihood: -421.3865



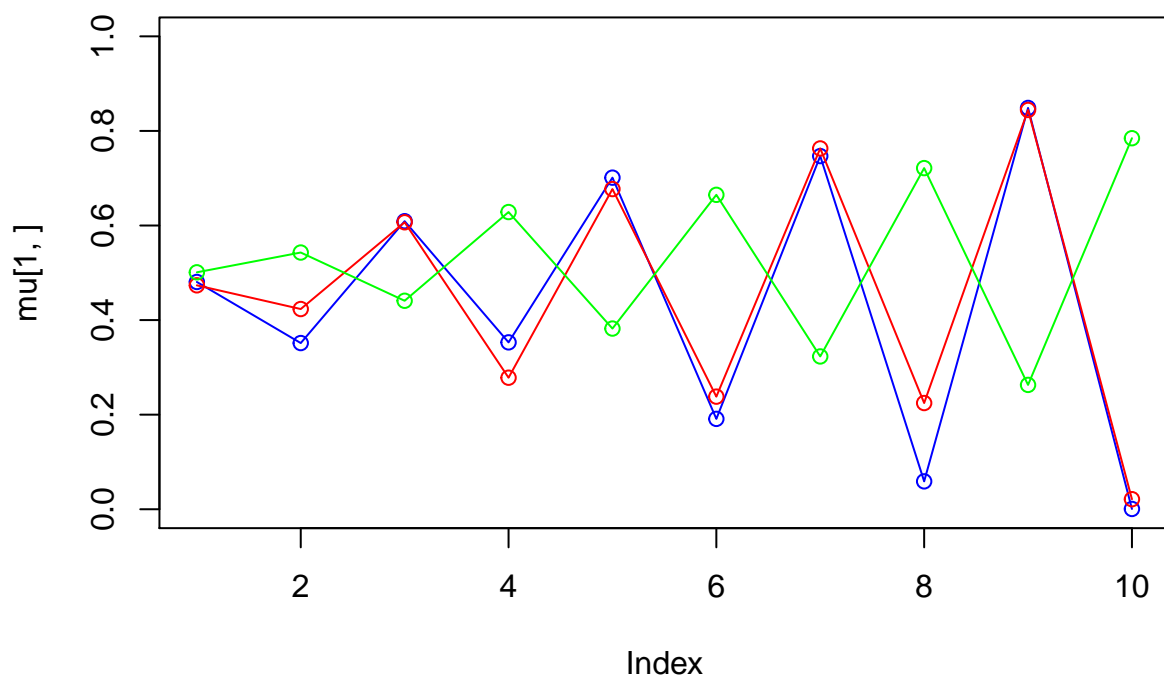
iteration: 41 log likelihood: -421.5282



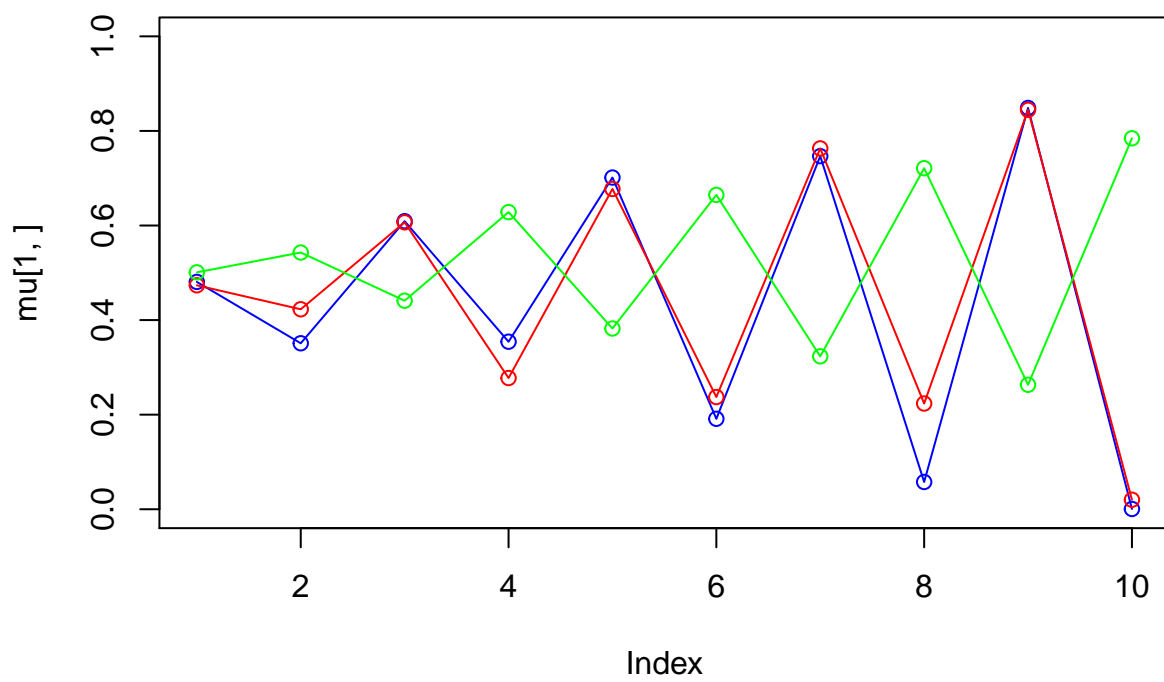
iteration: 42 log likelihood: -421.659



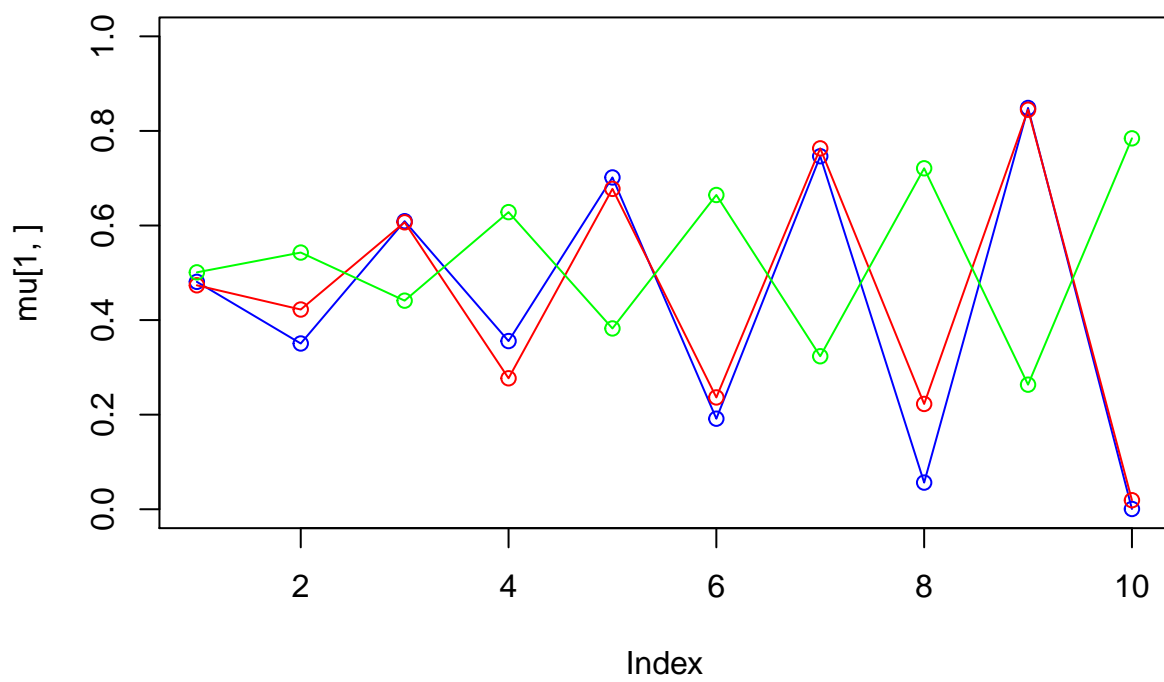
iteration: 43 log likelihood: -421.7797



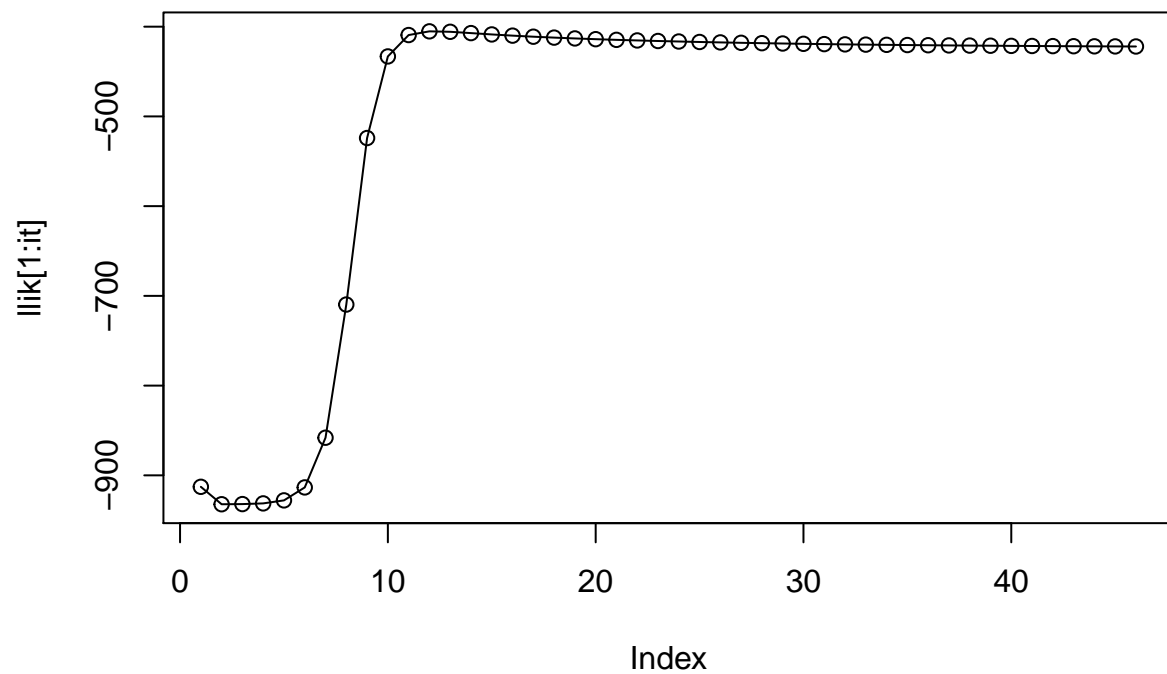
iteration: 44 log likelihood: -421.8913



iteration: 45 log likelihood: -421.9945



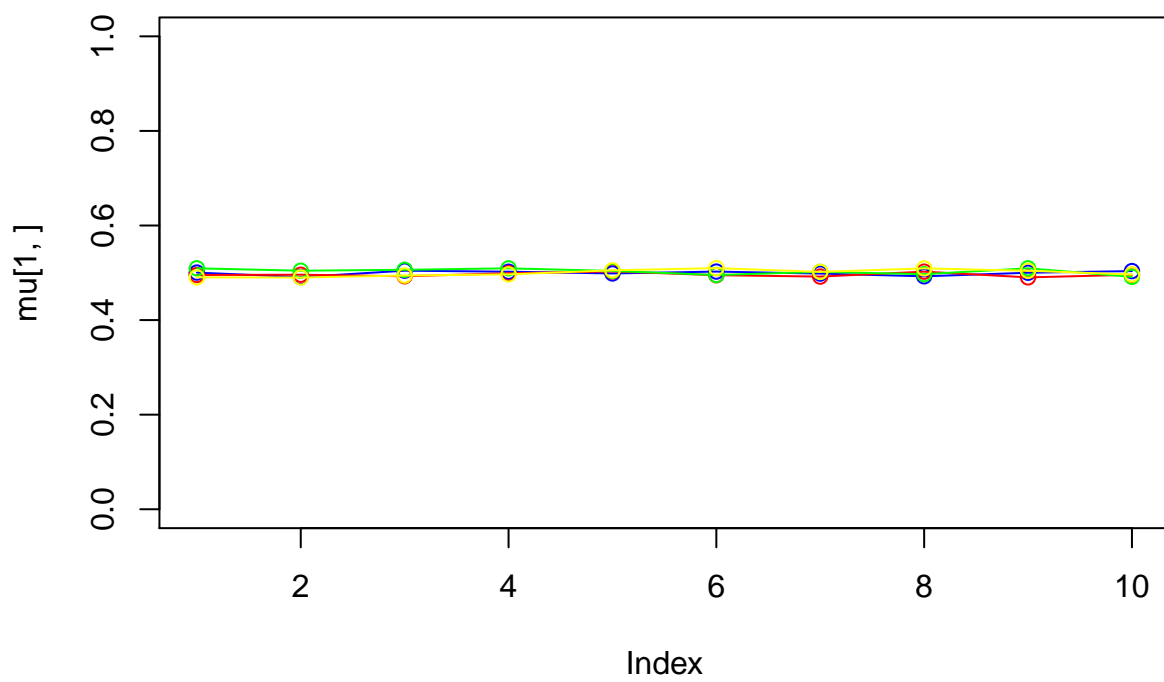
iteration: 46 log likelihood: -422.09



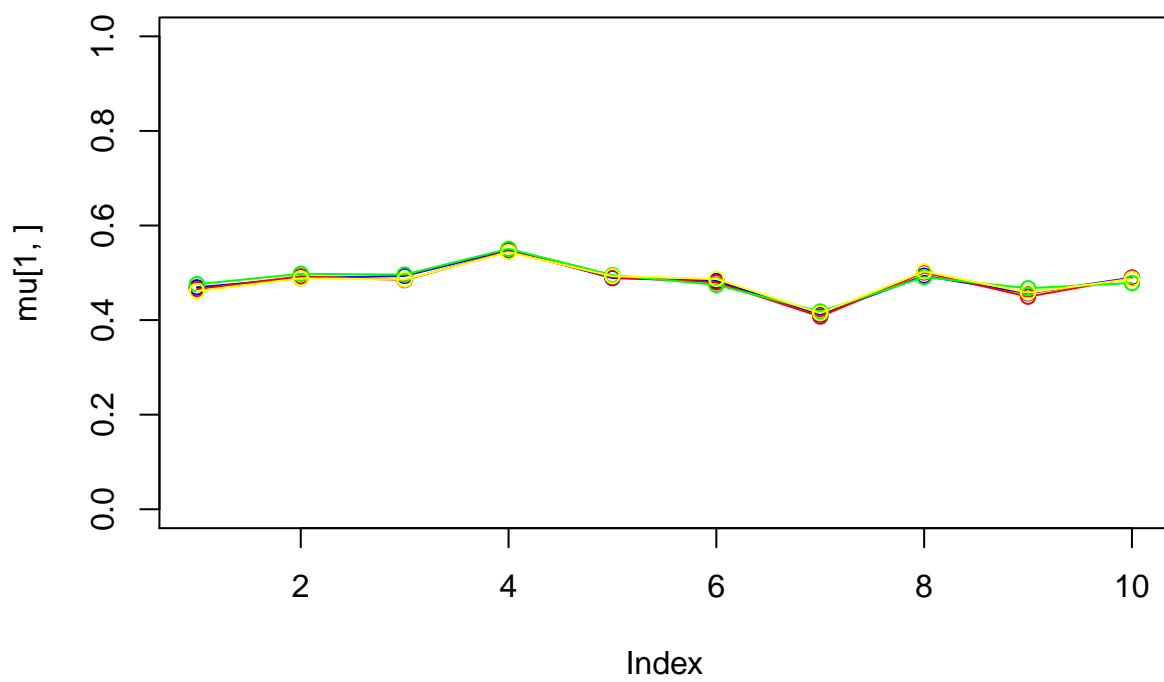
```
## [[1]]
## [1] 0.1679716522 0.2034249451 0.6286034026 0.4808697447 0.4735292886
## [6] 0.5009515437 0.3505033441 0.4223594793 0.5428015802 0.6091318393
## [11] 0.6067582335 0.4410624521 0.3556548071 0.2768901990 0.6282716562
## [16] 0.7016525159 0.6775124264 0.3823067747 0.1914725060 0.2364291675
## [21] 0.6645564670 0.7465112147 0.7631736115 0.3235087926 0.0563854887
## [26] 0.2226418271 0.7210236740 0.8485479479 0.8448194739 0.2634580777
## [31] 0.0005534402 0.0190935069 0.7843147843
```

K = 4

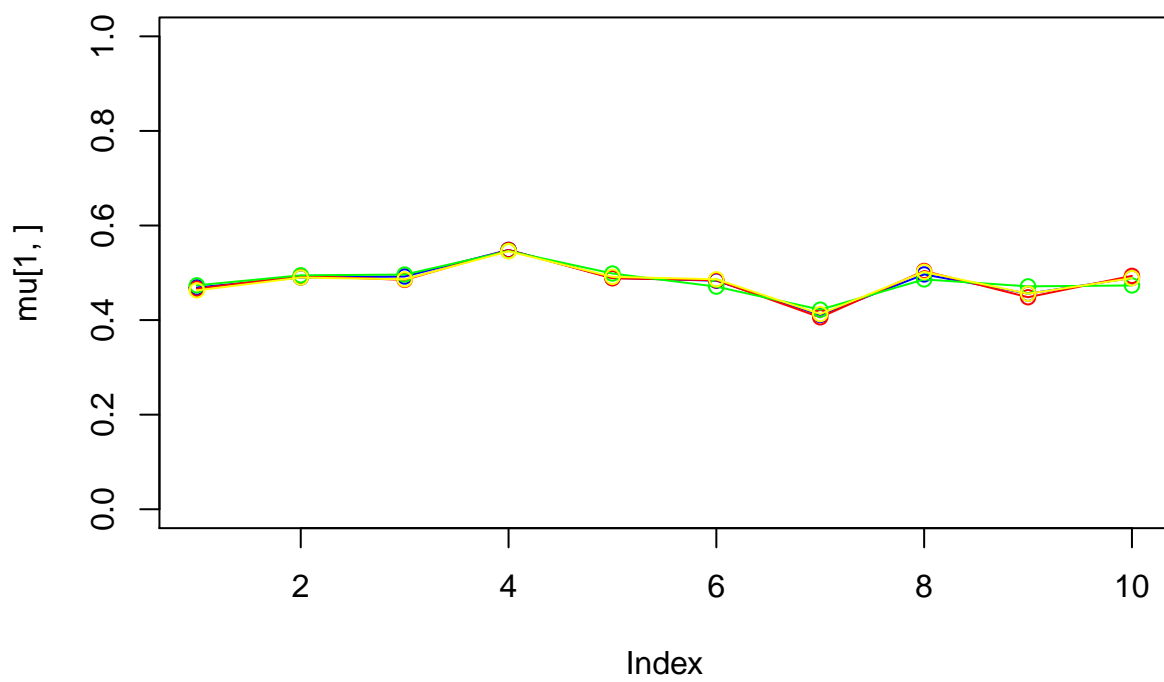
```
myem(K=4)
```



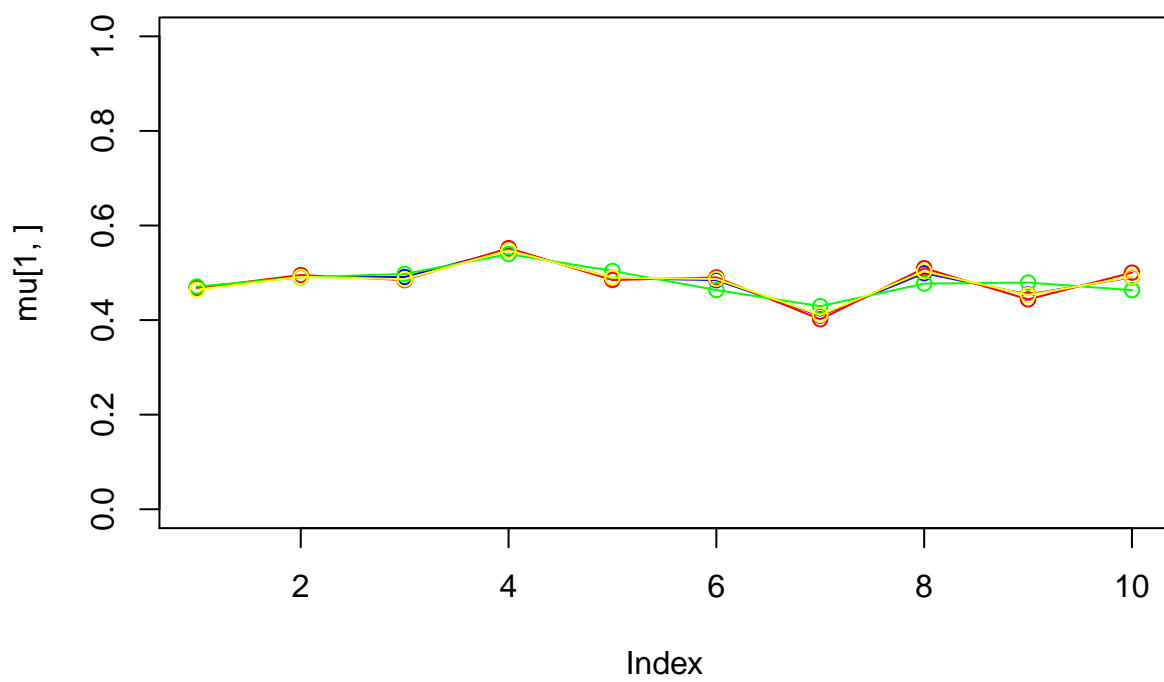
iteration: 1 log likelihood: -757.0451



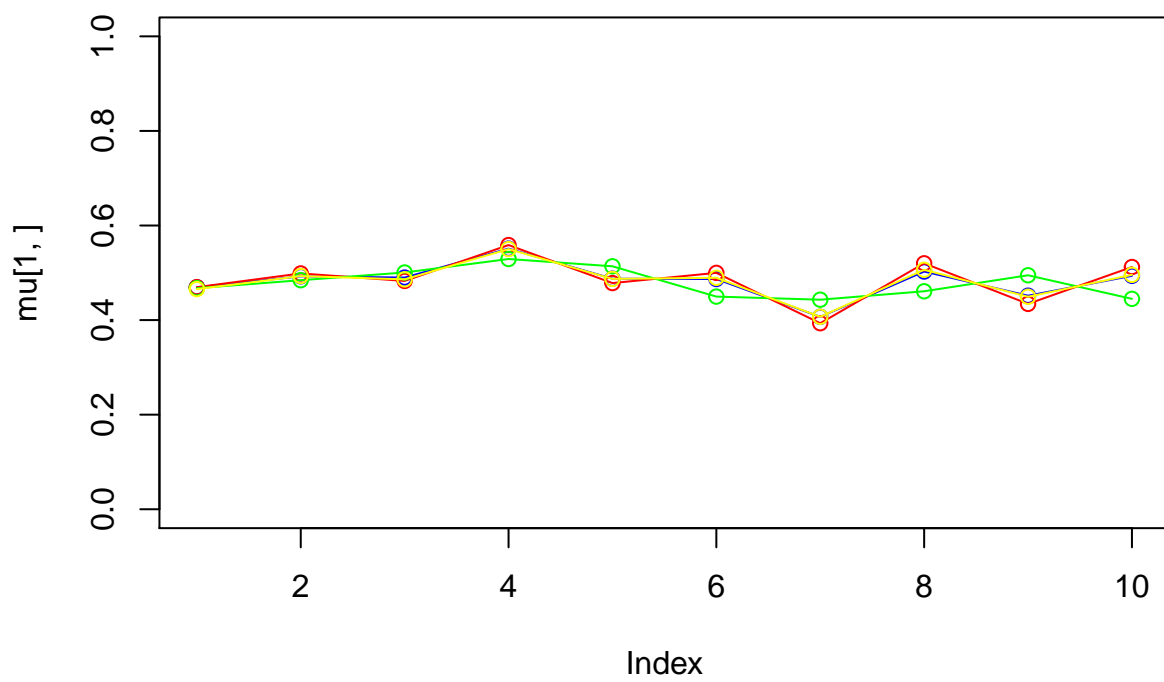
iteration: 2 log likelihood: -816.3628



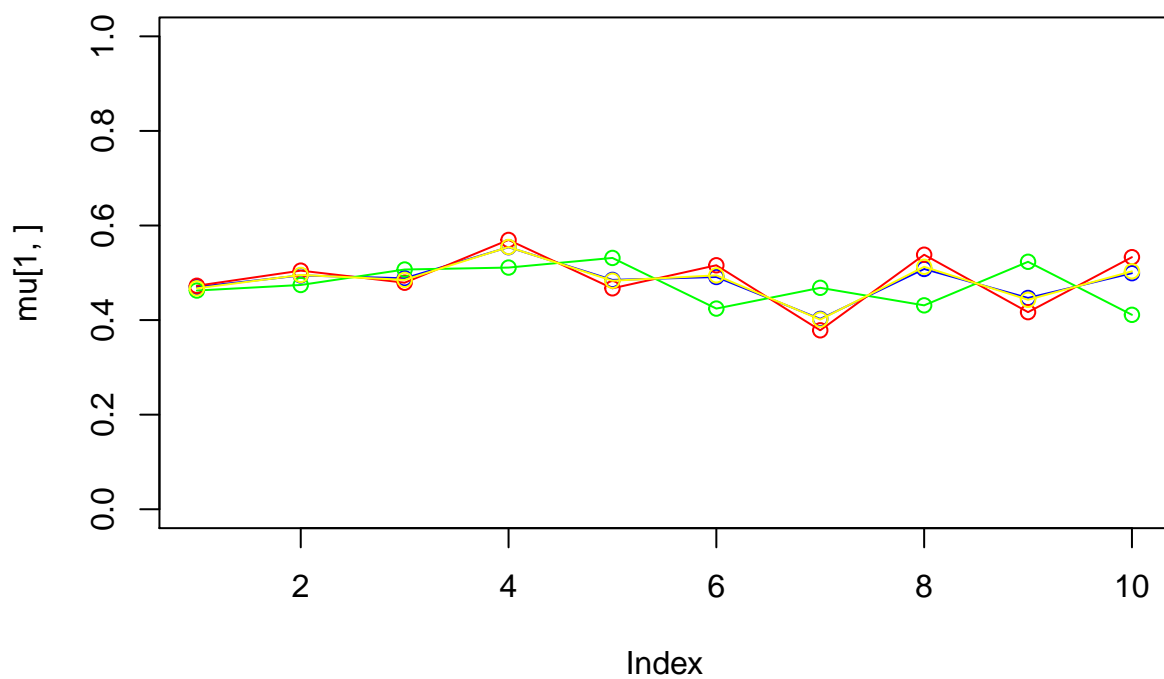
iteration: 3 log likelihood: -816.0773



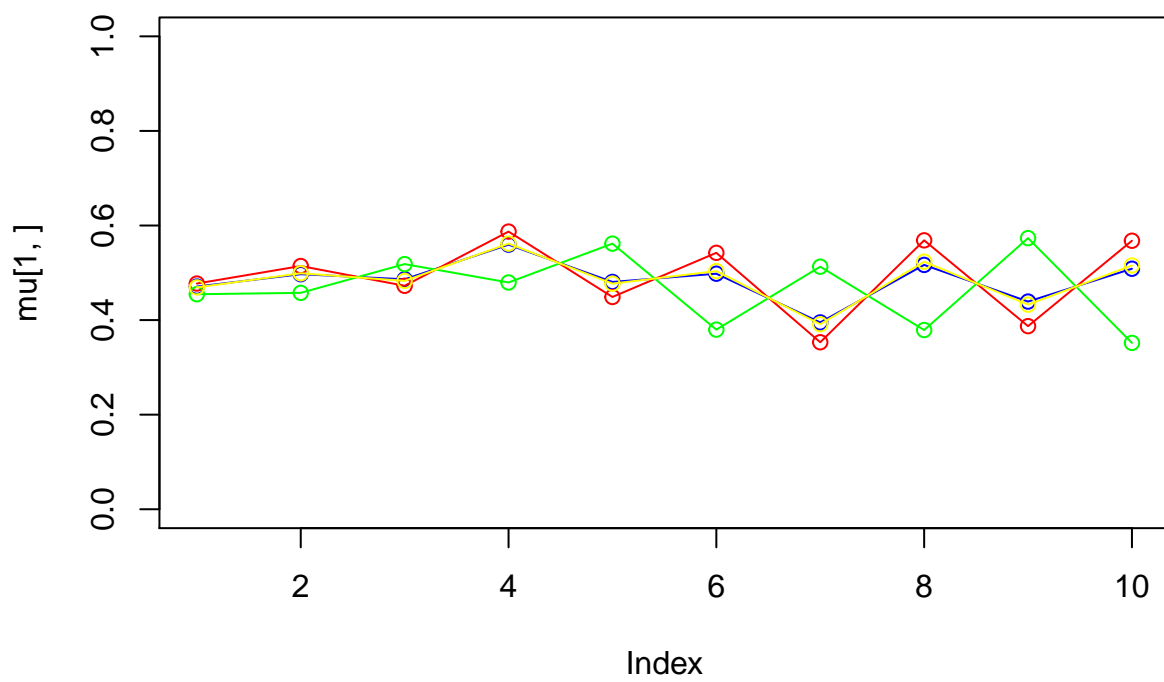
iteration: 4 log likelihood: -815.0484



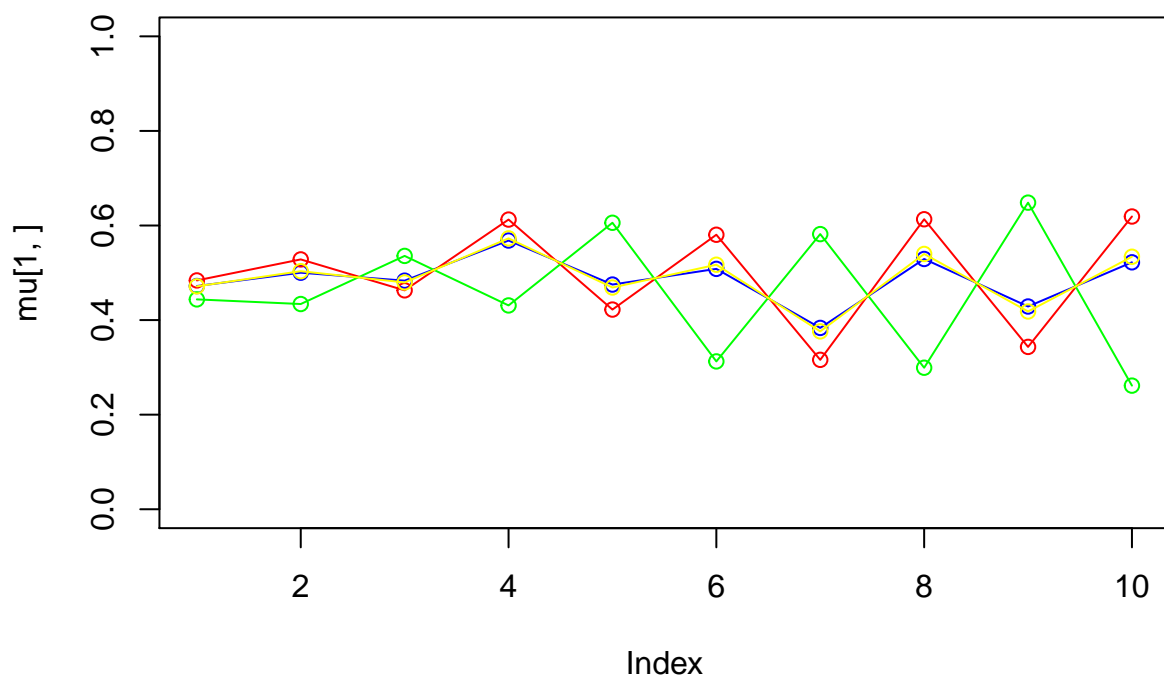
iteration: 5 log likelihood: -811.5926



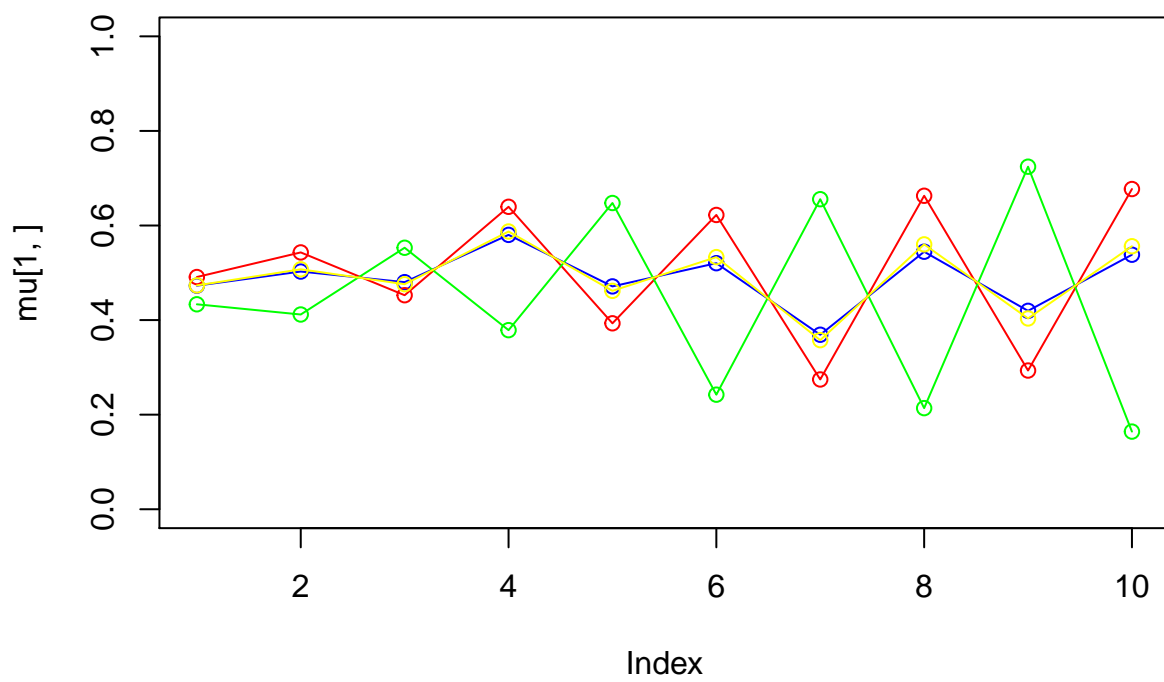
iteration: 6 log likelihood: -800.5387



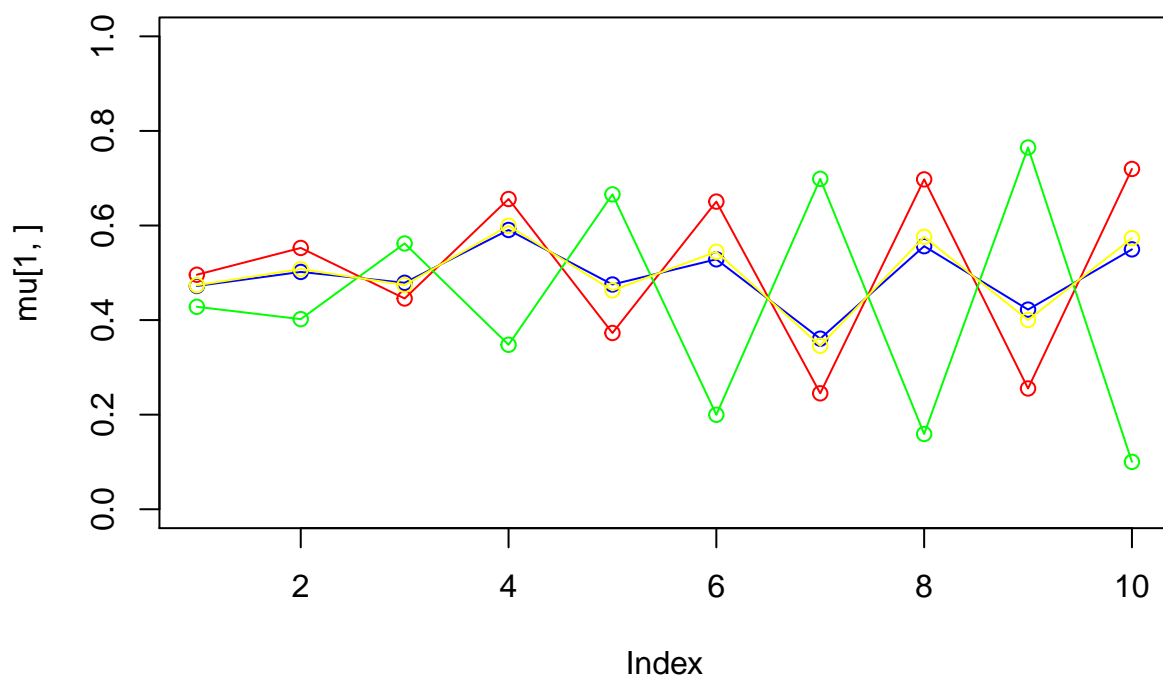
iteration: 7 log likelihood: -768.6586



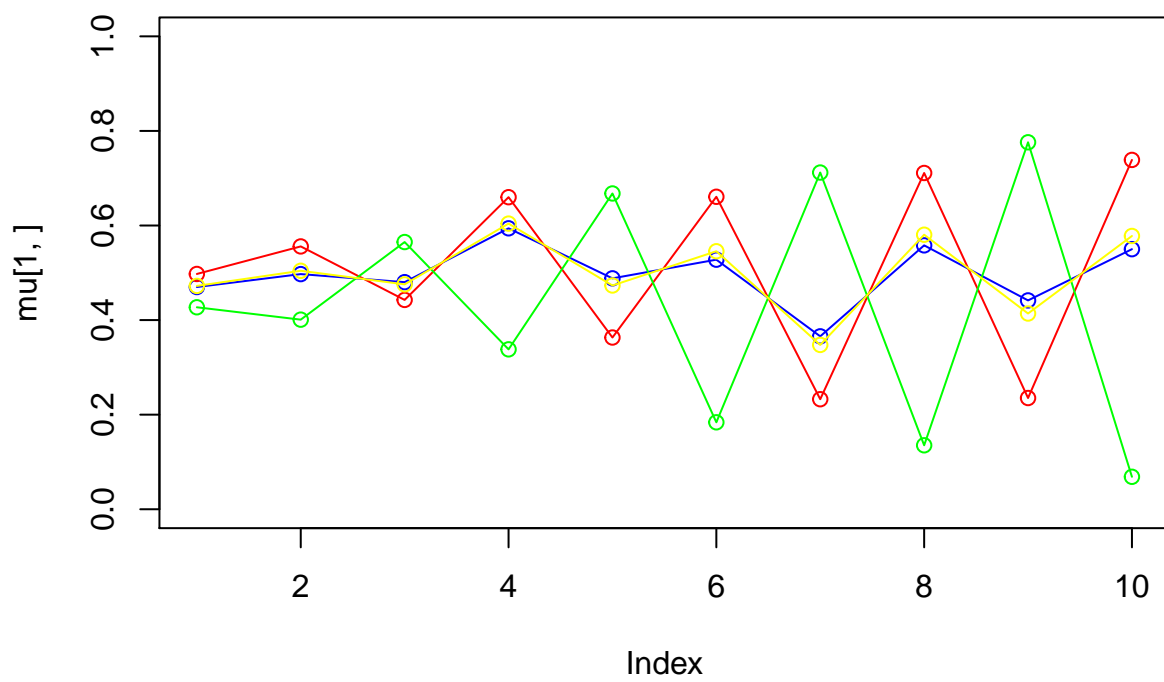
iteration: 8 log likelihood: -697.4363



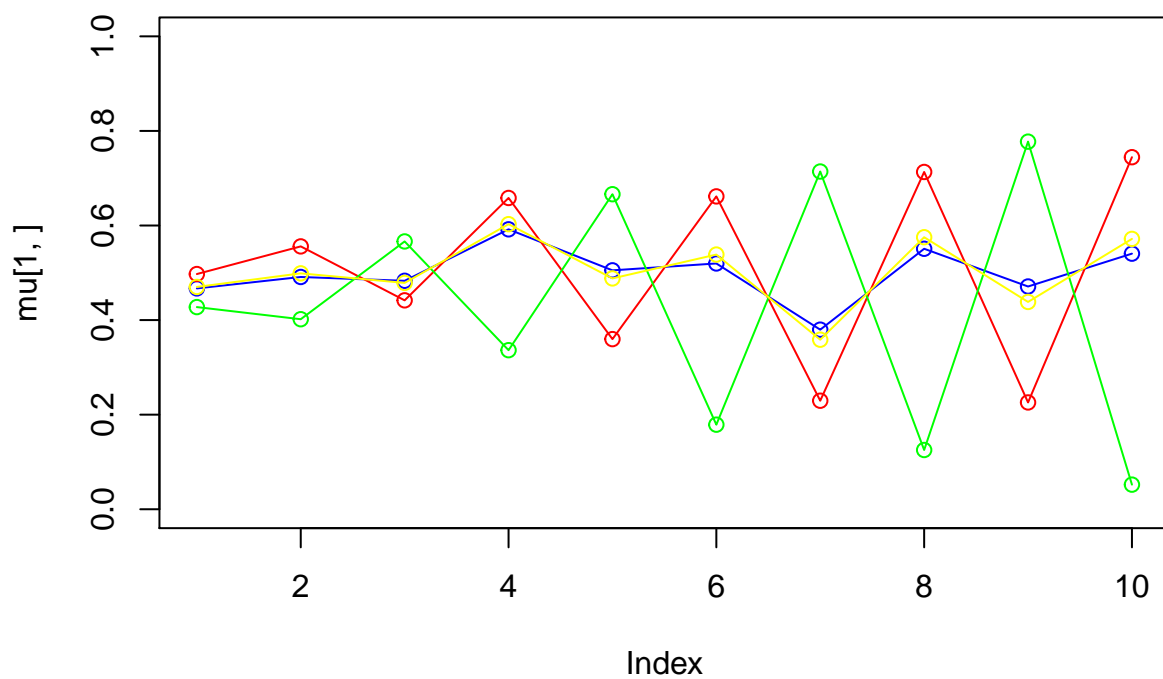
iteration: 9 log likelihood: -601.1825



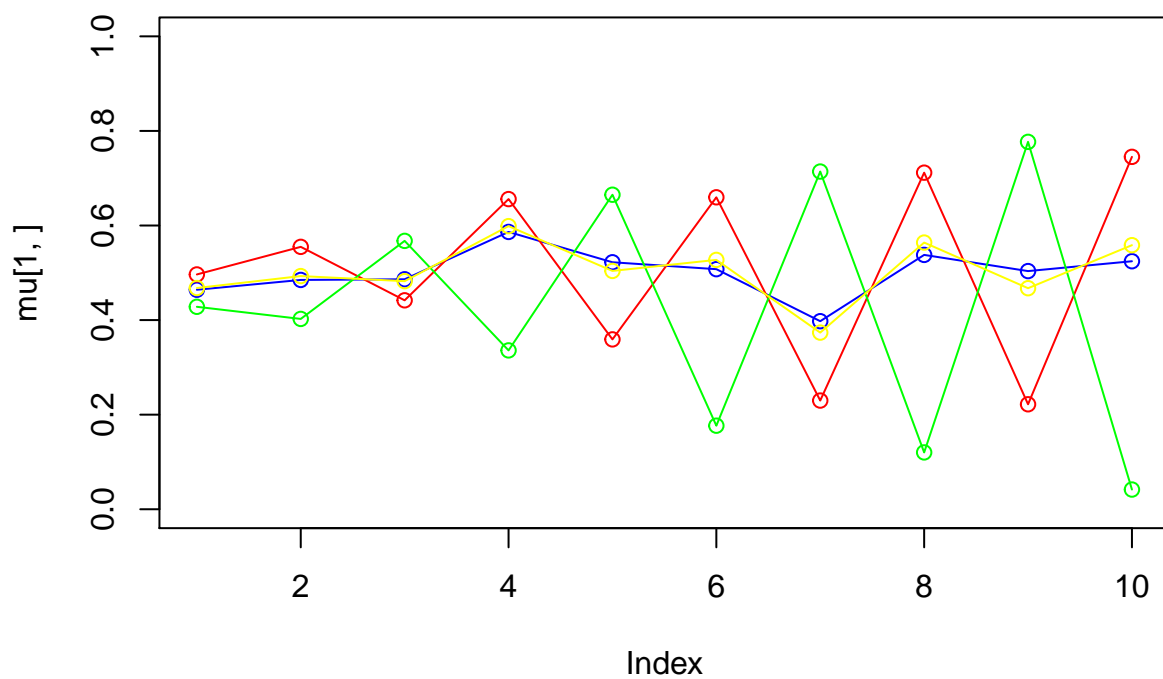
iteration: 10 log likelihood: -533.3314



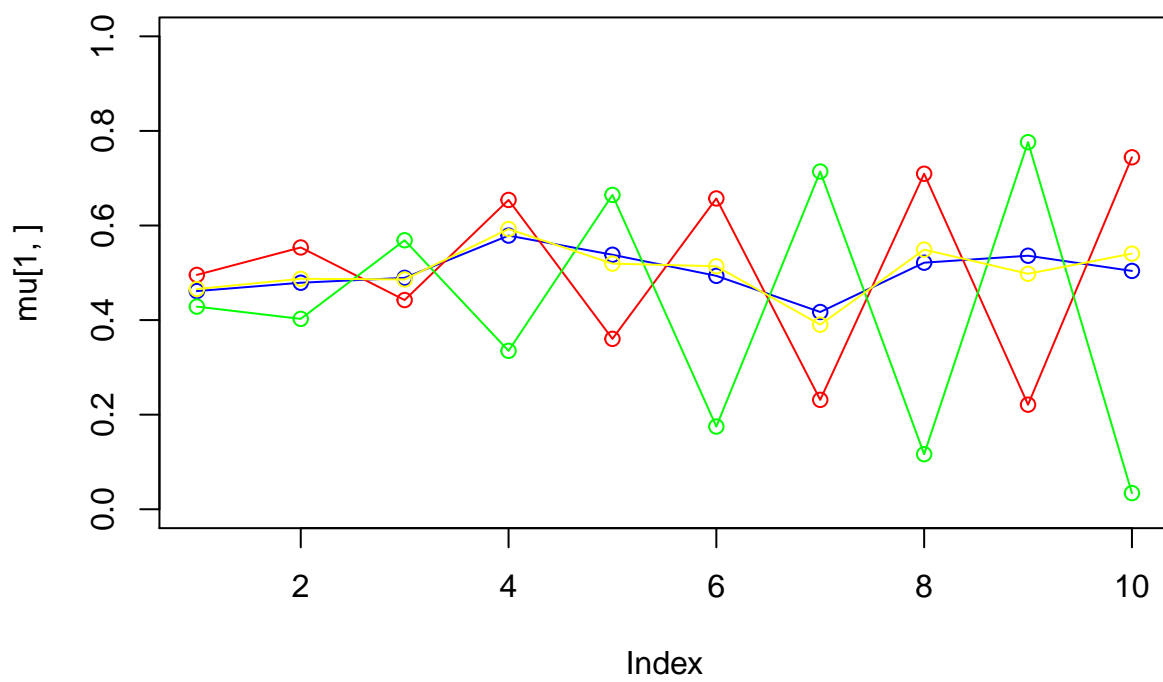
iteration: 11 log likelihood: -503.7221



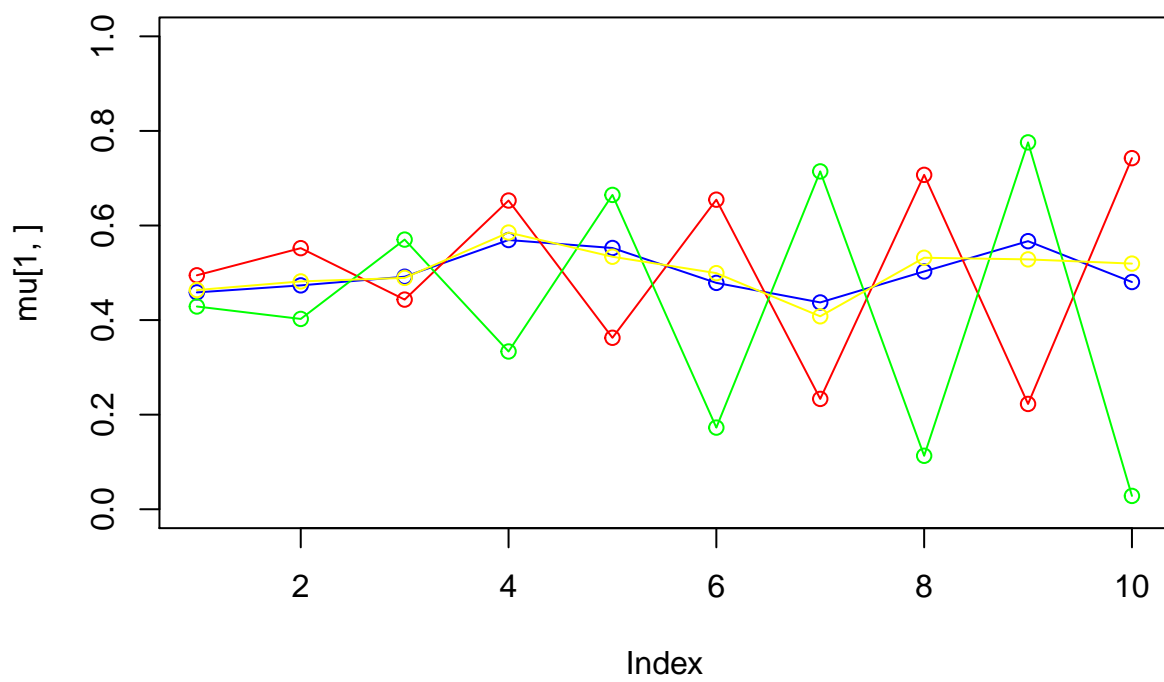
iteration: 12 log likelihood: -492.4649



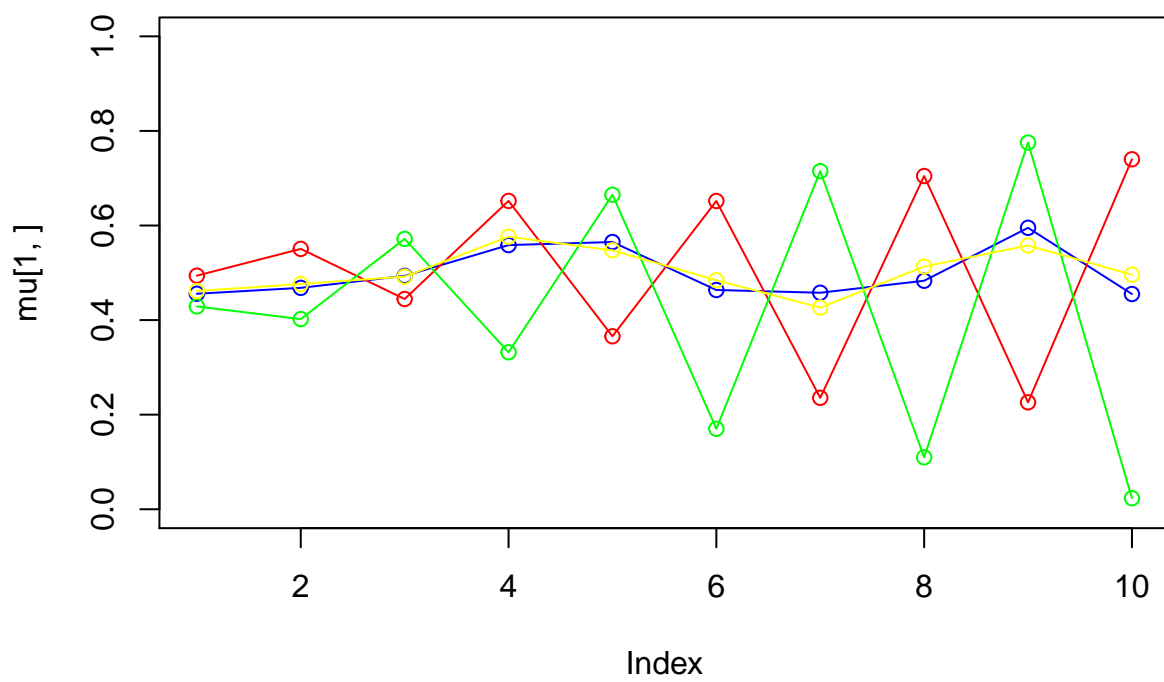
iteration: 13 log likelihood: -487.6123



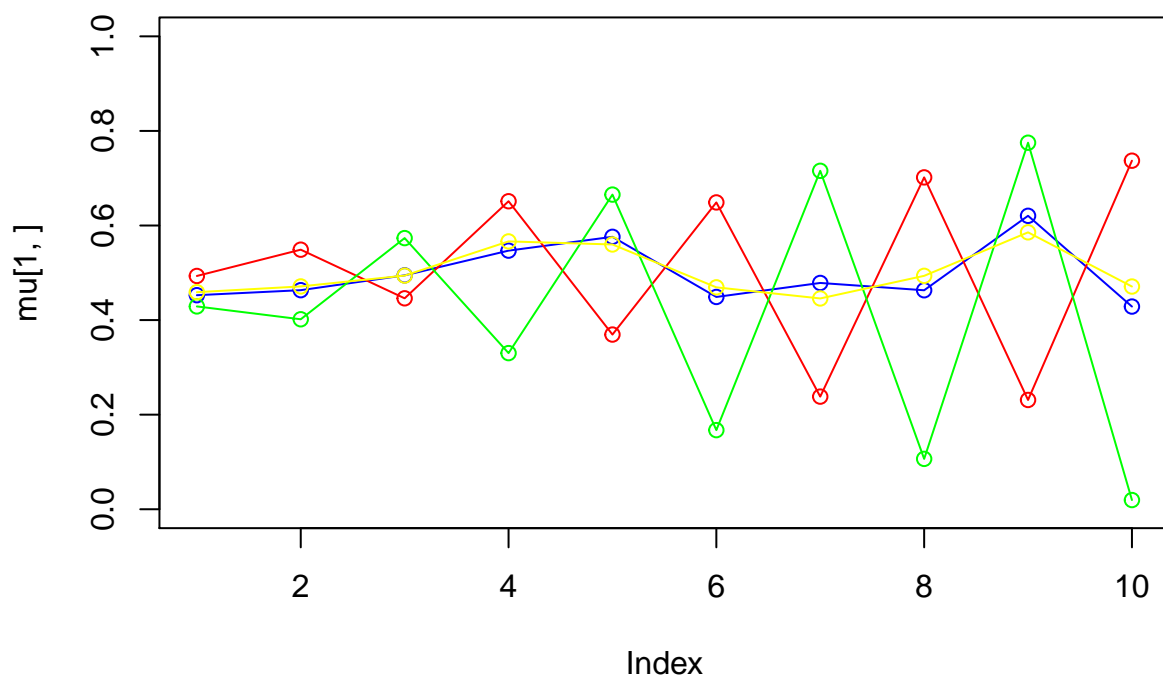
iteration: 14 log likelihood: -484.9568



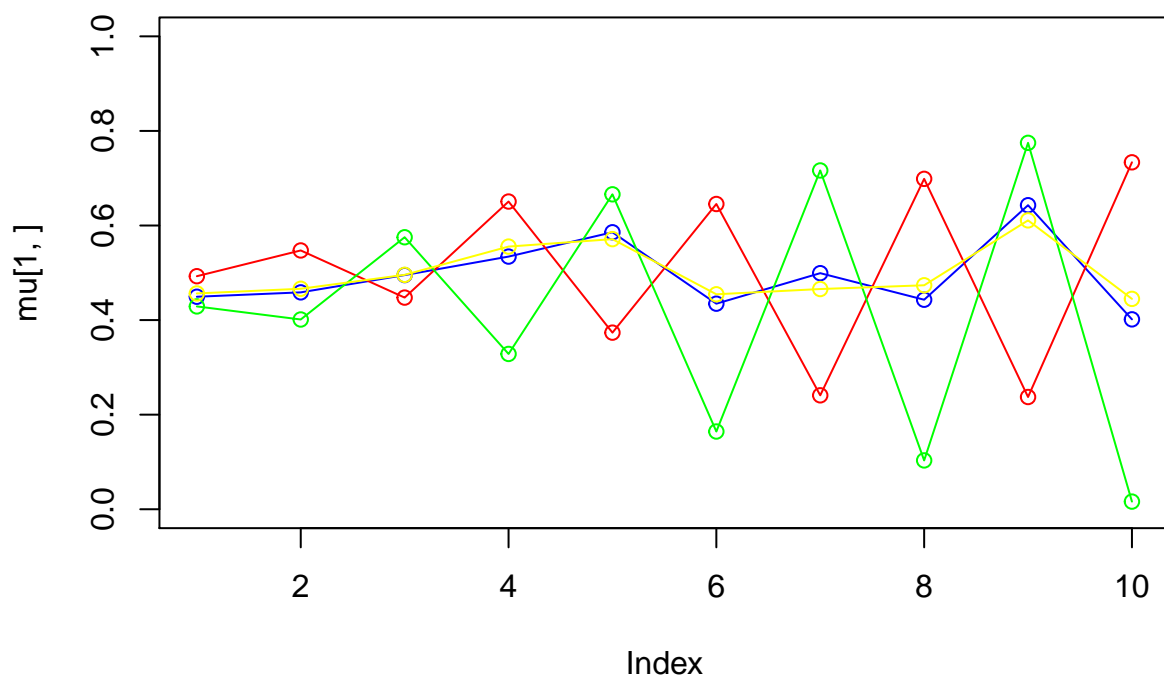
iteration: 15 log likelihood: -483.2528



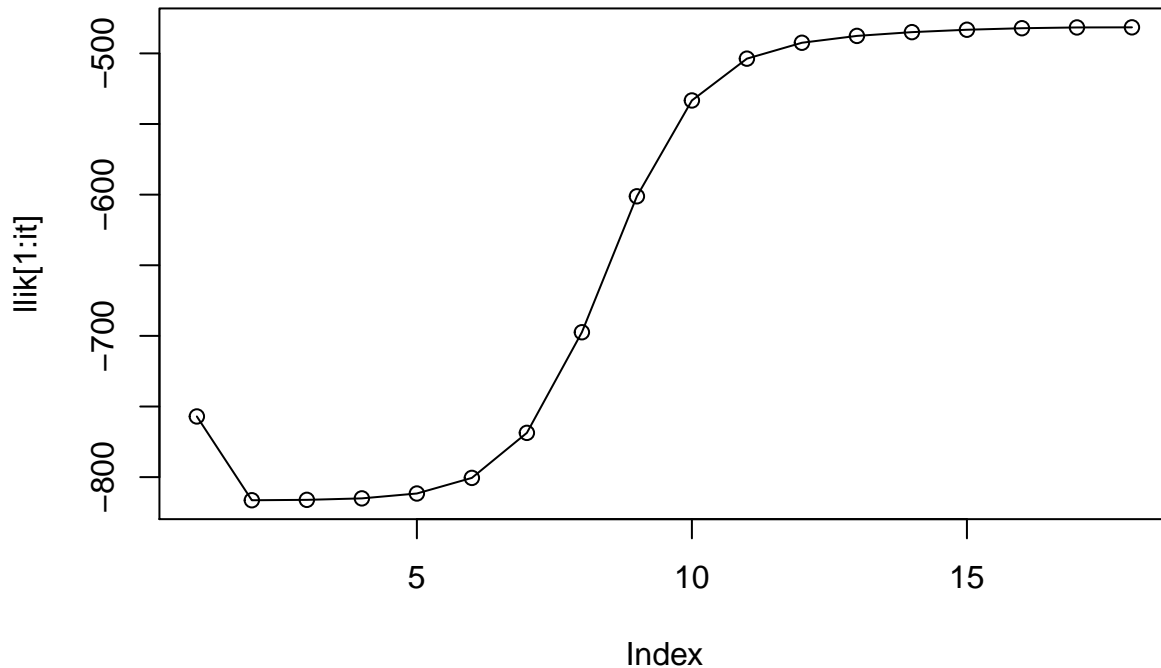
iteration: 16 log likelihood: -482.1709



iteration: 17 log likelihood: -481.6236



iteration: 18 log likelihood: -481.563



```
## [[1]]
## [1] 0.10548260 0.52735680 0.24858075 0.11857984 0.44936107 0.49277582
## [7] 0.42891247 0.45633523 0.45861698 0.54724067 0.40127012 0.46622403
## [13] 0.49516607 0.44756557 0.57514103 0.49563953 0.53427986 0.65068227
## [19] 0.32841859 0.55542754 0.58571495 0.37345536 0.66596065 0.57115980
## [25] 0.43467454 0.64546553 0.16437493 0.45453183 0.49919120 0.24090431
## [31] 0.71641466 0.46563220 0.44304425 0.69866151 0.10324431 0.47358884
## [37] 0.64287442 0.23733737 0.77477105 0.61083980 0.40154795 0.73361331
## [43] 0.01616773 0.44483717
```

Appendix

```
knitr::opts_chunk$set(echo = TRUE)
if (!require("pacman")) install.packages("pacman")
pacman::p_load(mboost, randomForest, dplyr, ggplot2)

options(scipen = 999)

spam_data <- read.csv(file = "spambase.data", header = FALSE)
colnames(spam_data)[58] <- "Spam"
spam_data$Spam <- factor(spam_data$Spam, levels = c(0,1), labels = c("0", "1"))
set.seed(12345)
n = NROW(spam_data)
id = sample(1:n, floor(n*(2/3)))
```

```

train = spam_data[id,]
test = spam_data[-id,]

final_result <- NULL
for(i in seq(from = 10, to = 100, by = 10)){

  ada_model <- mboost::blackboost(Spam~.,
                                data = train,
                                family = AdaExp(),
                                control=boost_control(mstop=i))

  forest_model <- randomForest(Spam~., data = train, ntree = i)

  prediction_function <- function(model, data){
    predicted <- predict(model, newdata = data, type = c("class"))
    predict_correct <- ifelse(data$Spam == predicted, 1, 0)
    score <- sum(predict_correct)/NROW(data)
    return(score)
  }

  train_ada_model_predict <- predict(ada_model, newdata = train, type = c("class"))
  test_ada_model_predict <- predict(ada_model, newdata = test, type = c("class"))
  train_forest_model_predict <- predict(forest_model, newdata = train, type = c("class"))
  test_forest_model_predict <- predict(forest_model, newdata = test, type = c("class"))

  test_predict_correct <- ifelse(test$Spam == test_forest_model_predict, 1, 0)
  train_predict_correct <- ifelse(train$Spam == train_forest_model_predict, 1, 0)

  train_ada_score <- prediction_function(ada_model, train)
  test_ada_score <- prediction_function(ada_model, test)
  train_forest_score <- prediction_function(forest_model, train)
  test_forest_score <- prediction_function(forest_model, test)

  iteration_result <- data.frame(number_of_trees = i,
                                accuracy = c(train_ada_score,
                                              test_ada_score,
                                              train_forest_score,
                                              test_forest_score),
                                type = c("train", "test", "train", "test"),
                                model = c("ADA", "ADA", "Forest", "Forest"))

  final_result <- rbind(iteration_result, final_result)
}

final_result$error_rate_percentage <- 100*(1 - final_result$accuracy)
ggplot(data = final_result, aes(x = number_of_trees,
                                y = error_rate_percentage,
                                group = type, color = type)) +
  geom_point() +

```

```

geom_line() +
ggtitle("Error Rate vs. increase in trees") + facet_grid(rows = vars(model))

myem <- function(K){
  set.seed(1234567890)

  max_it <- 100 # max number of EM iterations
  min_change <- 0.1 # min change in log likelihood between two consecutive EM iterations
  N=1000 # number of training points
  D=10 # number of dimensions
  x <- matrix(nrow=N, ncol=D) # training data
  true_pi <- vector(length = K) # true mixing coefficients
  true_mu <- matrix(nrow=K, ncol=D) # true conditional distributions
  true_pi=c(rep(1/3, K))

  if(K == 2){
    true_mu[1,]=c(0.5,0.6,0.4,0.7,0.3,0.8,0.2,0.9,0.1,1)
    true_mu[2,]=c(0.5,0.4,0.6,0.3,0.7,0.2,0.8,0.1,0.9,0)
  }else if(K == 3){
    true_mu[1,]=c(0.5,0.6,0.4,0.7,0.3,0.8,0.2,0.9,0.1,1)
    true_mu[2,]=c(0.5,0.4,0.6,0.3,0.7,0.2,0.8,0.1,0.9,0)
    true_mu[3,]=c(0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5)
  }else {
    true_mu[1,]=c(0.5,0.6,0.4,0.7,0.3,0.8,0.2,0.9,0.1,1)
    true_mu[2,]=c(0.5,0.4,0.6,0.3,0.7,0.2,0.8,0.1,0.9,0)
    true_mu[3,]=c(0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5,0.5)
    true_mu[4,]=c(0.5,0.4,0.5,0.7,0.5,0.5,0.2,0.5,0.5,0.5)}

  # Producing the training data
  for(n in 1:N) {
    k <- sample(1:K,1,prob=true_pi)
    for(d in 1:D) {
      x[n,d] <- rbinom(1,1,true_mu[k,d])
    }
  }

  z <- matrix(nrow=N, ncol=K) # fractional component assignments
  pi <- vector(length = K) # mixing coefficients
  mu <- matrix(nrow=K, ncol=D) # conditional distributions
  llik <- vector(length = max_it) # log likelihood of the EM iterations
  # Random initialization of the paramters
  pi <- runif(K,0.49,0.51)
  pi <- pi / sum(pi)

  for(k in 1:K) {
    mu[k,] <- runif(D,0.49,0.51)
  }

  for(it in 1:max_it) {

```

```

if(K == 2){
  plot(mu[1,], type="o", col="blue", ylim=c(0,1))
  points(mu[2,], type="o", col="red")
}else if(K == 3){
  plot(mu[1,], type="o", col="blue", ylim=c(0,1))
  points(mu[2,], type="o", col="red")
  points(mu[3,], type="o", col="green")
}else{
  plot(mu[1,], type="o", col="blue", ylim=c(0,1))
  points(mu[2,], type="o", col="red")
  points(mu[3,], type="o", col="green")
  points(mu[4,], type="o", col="yellow")}

Sys.sleep(0.5)
# E-step: Computation of the fractional component assignments

for(k in 1:K)
prod <- exp(x %*% log(t(mu))) * exp((1-x) %*% t(1-mu))

num = matrix(rep(pi,N), ncol = K, byrow = TRUE) * prod
dem = rowSums(num)
poster = num/dem

#Log likelihood computation.
llik[it] = sum(log(dem))
# Your code here
cat("iteration: ", it, "log likelihood: ", llik[it], "\n")
flush.console()
# Stop if the log likelihood has not changed significantly
if( it != 1){
  if(abs(llik[it] - llik[it-1]) < min_change){break}
}

#M-step: ML parameter estimation from the data and fractional component assignments
# Your code here
num_pi = colSums(poster)
pi = num_pi/N
mu = (t(poster) %*% x)/num_pi

}

a <- pi
b <- mu
c <- plot(llik[1:it], type="o")
result <- list(c(a,b,c))
return(result)
}

myem(K=2)
myem(K=3)
myem(K=4)

```