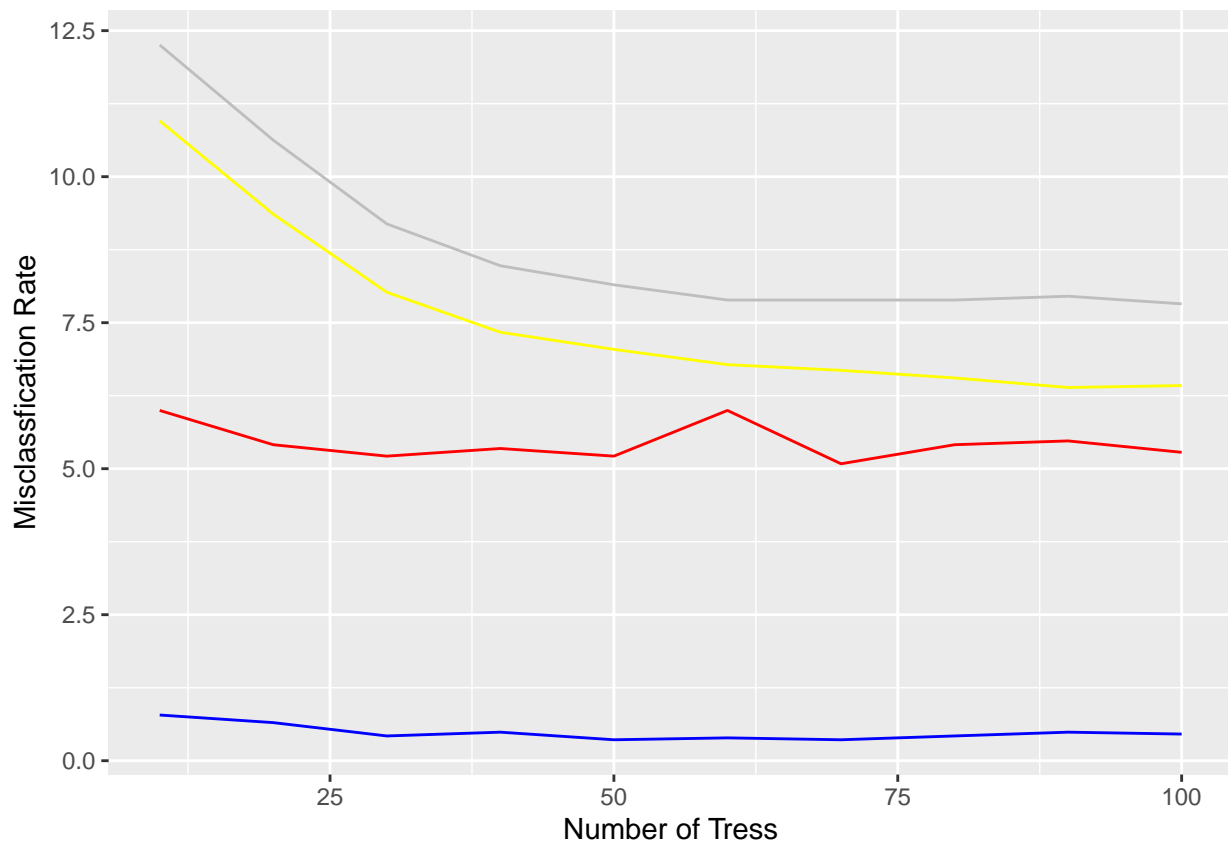# Block 2 Lab 1 Report

*Zuxiang Li(zuxli371)*

*12/3/2019*

## 1. ENSEMBLE METHODS

The file spambase.csv contains information about the frequency of various words, char- acters, etc. for a total of 4601 e-mails. Furthermore, these e-mails have been classified as spams (spam = 1) or regular e-mails (spam = 0). 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. To learn Adaboost classification trees, use the function blackboost() of the R package mboost. Specify the loss function corresponding to Adaboost with the parameter family. To learn random forests, use the function randomForest of the R package randomForest.
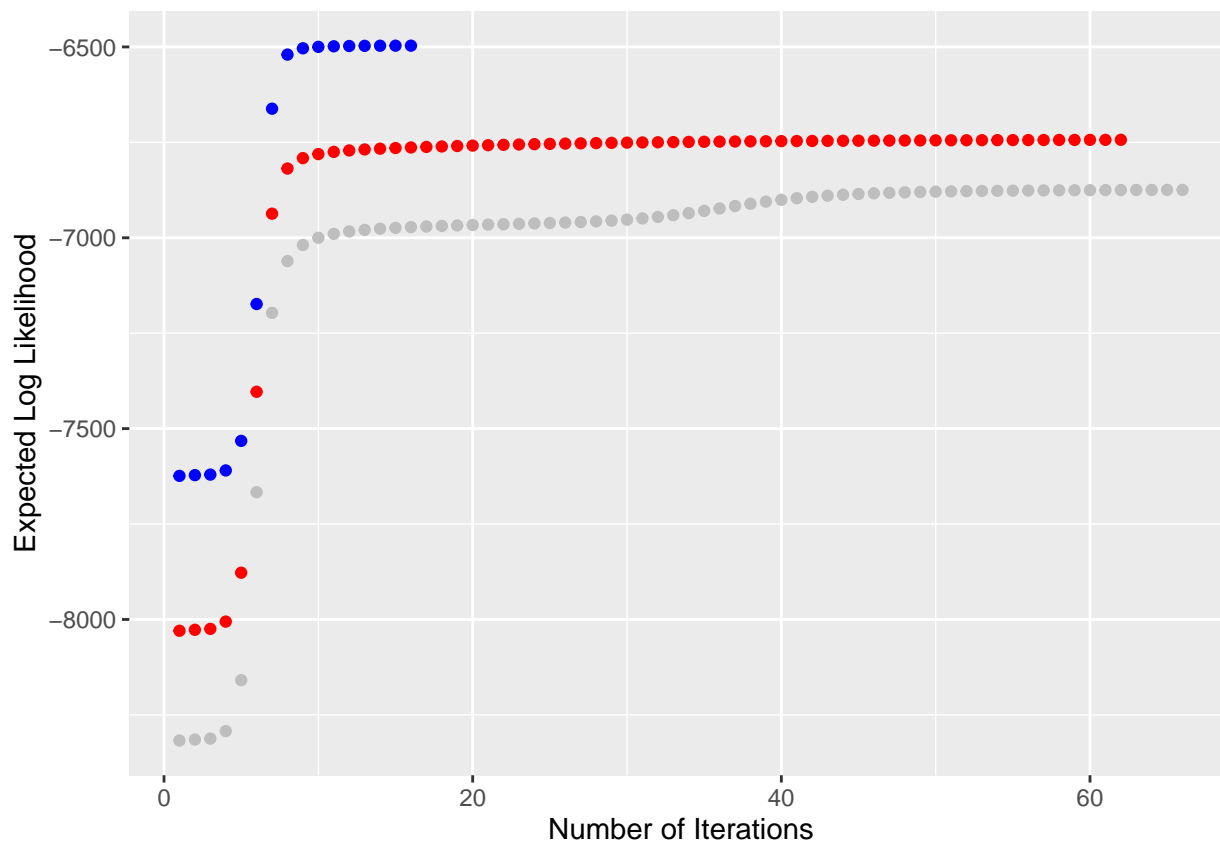


The plot above showed the Misclassfication rate using different models and number of trees. Grey and yellow represent Adaboost model for train and test data. Red and blue lines are Randomforest model for train and test data.

The misclassfication values of Adaboost with incresaing number of trees are constantly decreasing for both train and test data. In every iteration a classfier is produced and given different weight. Then use the classifer generated in previous iteration to produced.It returns a weighted average of the classfiers. So the rate keeps decrease.

As for Random forest, number of tress seems have no effect on the misclassfication rate for both train and test data. Since Random forest used bagging, individual classifiers are then combined by taking a simple majority vote of their decisions.

## 2. MIXTURE MODELS

Your task is to implement the EM algorithm for mixtures of multivariate Benoulli distributions. Please use the template in the next page to solve the assignment. Then, use your implementa- tion 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.



In this case, we generate data using true_mu and true_pi. For K=2, the blue dots, which means too few components, the blue points in the plot, only 16 iterations were made, this model is underfitted. For K=4, the grey dots, too many components, grey dots in the plot, 66 iterations were made, then this model is overfitted. Only when K=3, the red dots in plot, we got similar mu and pi compare to the true values.

# Appendix

## 1. ENSEMBLE METHODS

```r
library(mboost)
library(randomForest)
# Import data
sp <- read.csv2("material/spambase.csv")
sp$Spam <- as.factor(sp$Spam)
n <- dim(sp)[1]
set.seed(1234567890)
id <- sample(1:n, floor(n*2/3))
train <- sp[id,]
test <- sp[-id,]

# Generate confusion matrix and calculate misclassfication rate
cft<-function(Pred,Actual){
  cft<-table(Pred,Actual)
  # True positive
  tp <- cft[2, 2]
  # True negetive
  tn <- cft[1, 1]
  # False positive
  fp <- cft[2, 1]
  # False negetive
  fn <- cft[1, 2]

  misclassfication <- (1-(tp + tn)/(tp + tn + fp + fn))*100
  return(misclassfication)
}


ensemble_methods<-function(method,data){
  err_rate<-c()
  # select methods mboost
  if(substitute(method)=="mboost"){
    # Loop for 10 times each time the number of trees increase 10
    for(i in seq(10,100,10)){
      # Generate mboost model using train data,specify loss function as Binomial()
      bb<-blackboost(Spam~.,data=train,control = boost_control(mstop = i),family = AdaExp())
      # Get the predict value based on model we generated
      pred<-predict.mboost(bb,data,type="class")
      # Get the missclassfication rate for current number of trees
      res<-cft(Pred=pred,Actual=data$Spam)
      err_rate<-c(err_rate,res)
    }
  }
  # select methods randomforest
  else if(substitute(method)=="randomforest"){
    # Loop for 10 times each time the number of trees increase 10
    for(i in seq(10,100,10)){
      # Generate randomforest model
```

```r
    rf<-randomForest(Spam~.,data=train,ntree=i)
    # Get predict value
    pred<-predict(rf,data)
    # Get the missclassfication rate for current number of trees
    res<-cft(Pred=pred,Actual=data$Spam)
    err_rate<-c(err_rate,res)
  }
 }
    return(err_rate)
}
err_rate1<-ensemble_methods("mboost",train)
err_rate2<-ensemble_methods("mboost",test)
err_rate3<-ensemble_methods("randomforest",train)
err_rate4<-ensemble_methods("randomforest",test)


df<-data.frame(step=seq(10,100,10),err1=err_rate1,err2=err_rate2,err3=err_rate3,err4=err_rate4)
ggplot()+xlab("")+geom_line(data = df, aes(x = step, y = err3), color = "blue") +
  geom_line(data = df, aes(x = step, y = err4), color = "red") +geom_line(data = df, aes(x = step, y = 
  geom_line(data = df, aes(x = step, y = err2), color = "grey")+xlab("Number of Tress")+ylab("Misclassf
```

## 2. MIXTURE MODELS

```r
RNGversion('3.5.1')
```

```
## Warning in RNGkind("Mersenne-Twister", "Inversion", "Rounding"): non-uniform
## 'Rounding' sampler used
```

```r
EM_algorithm<-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 = 3) # true mixing coefficients
  true_mu <- matrix(nrow=3, ncol=D) # true conditional distributions
  true_pi=c(1/3, 1/3, 1/3)
  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)
  plot(true_mu[1,], type="o", col="blue", ylim=c(0,1))
  points(true_mu[2,], type="o", col="red")
  points(true_mu[3,], type="o", col="green")
  # Producing the training data
  for(n in 1:N) {
    k <- sample(1:3,1,prob=true_pi)
    for(d in 1:D) {
      x[n,d] <- rbinom(1,1,true_mu[k,d])
    }
  }
  #K=3 # number of guessed components
```

```r
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) {
  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
  # Your code here
  for(n in 1:N){
    denominator<-c()
    numerator<-c()
    for(k in 1:K){
      new_numerator<-pi[k]*prod(mu[k,]^x[n,]*(1-mu[k,])^(1-x[n,]))
      numerator<-c(numerator,new_numerator)
      denominator<-sum(numerator)
    }
    z[n,]<-numerator/denominator
  }
  #Log likelihood computation.
  sum=0
  for(n in 1:N){
    for(k in 1:K){
      sum=sum+z[n,k]*(log(pi[k])+sum(x[n,]*log(mu[k,])+(1-x[n,])*log(1-mu[k,])))
    }
  }
  llik[it]<-sum
  # Your code here
  #cat("iteration: ", it, "log likelihood: ", llik[it], "\n")
  flush.console()
  # Stop if the lok likelihood has not changed significantly
  if(it>1){
    if(abs(llik[it]-llik[it-1])<min_change)
      break
  }
  # Your code here
  #M-step: ML parameter estimation from the data and fractional component assignments
  # Your code here
  pi=apply(z,2,mean)
  for(k in 1:K){
    sum=0
    for (n in 1:N) {
      sum = sum + x[n,] * z[n,k]
```

```
      }
      mu[k,] = sum / sum(z[,k])
    }
  }
  pi
  mu
  plot(llik[1:it], type="o")
  return(llik[1:it])
  #return(list(pi=pi,mu=mu,res=,llik=llik[1:it])
}

#em1<-EM_algorithm(K=2)
em2<-EM_algorithm(K=3)
```
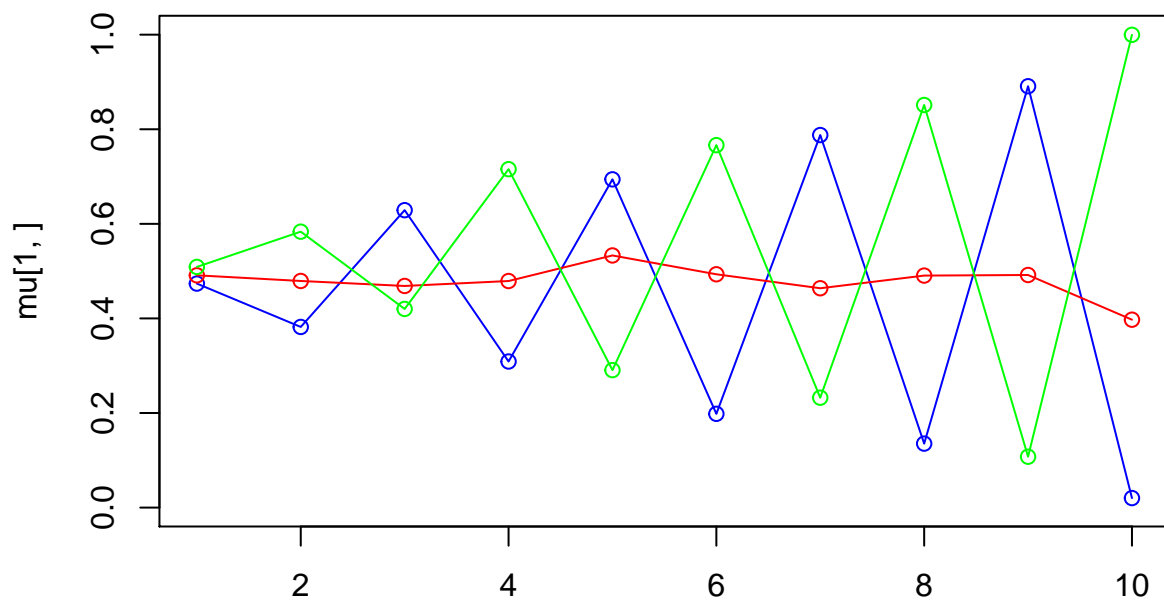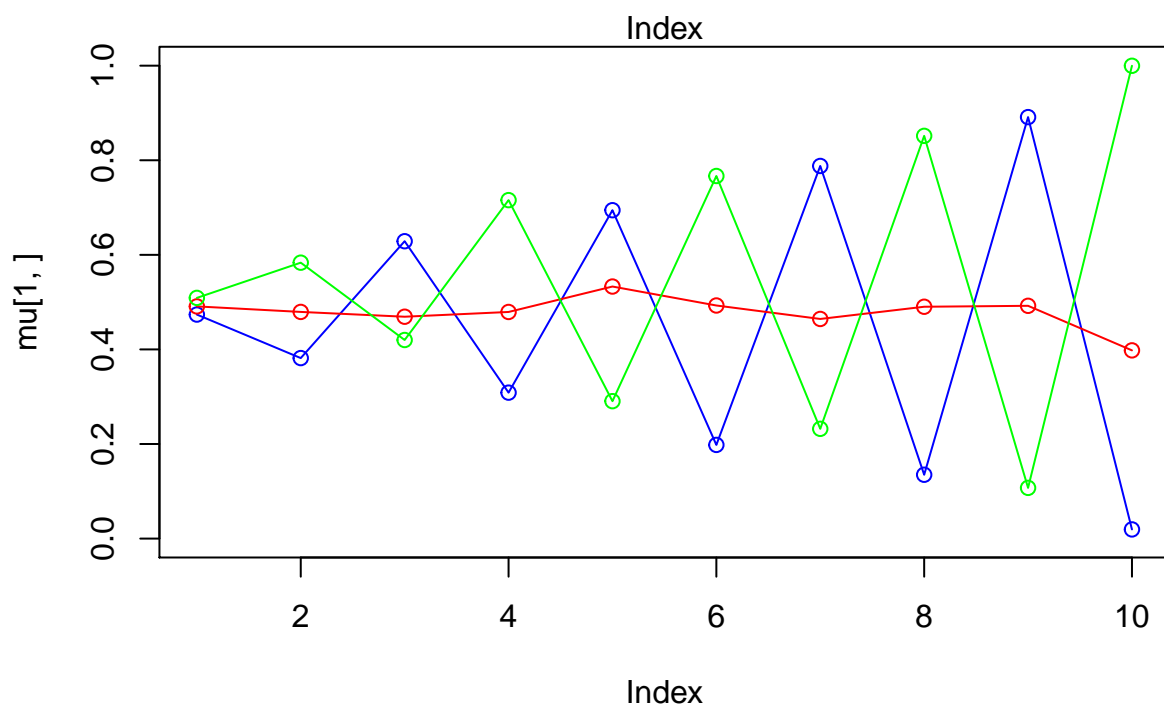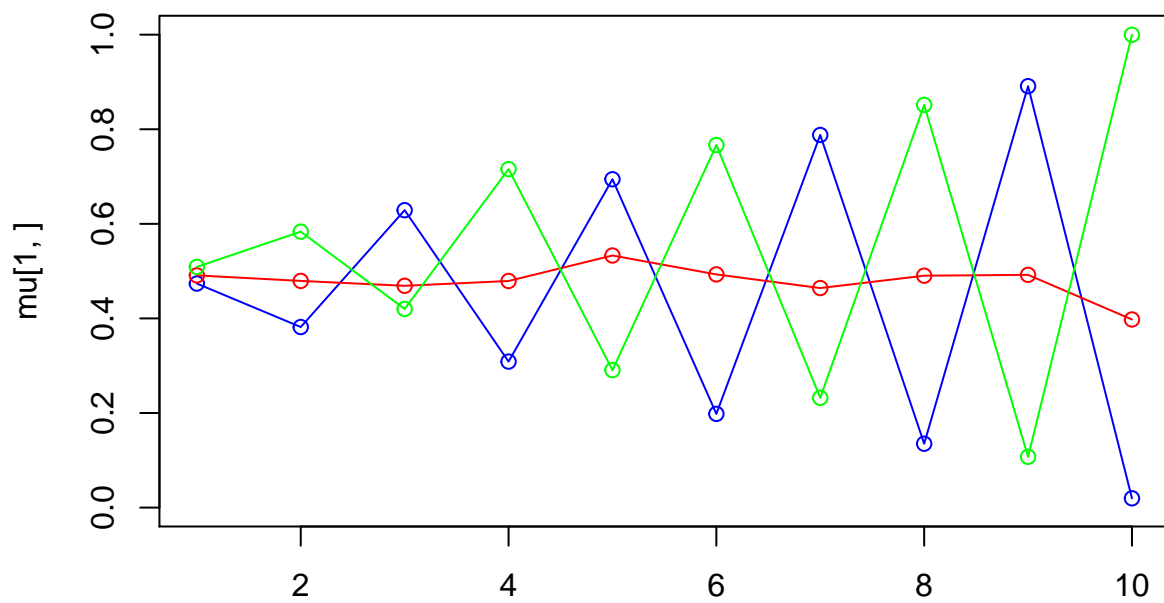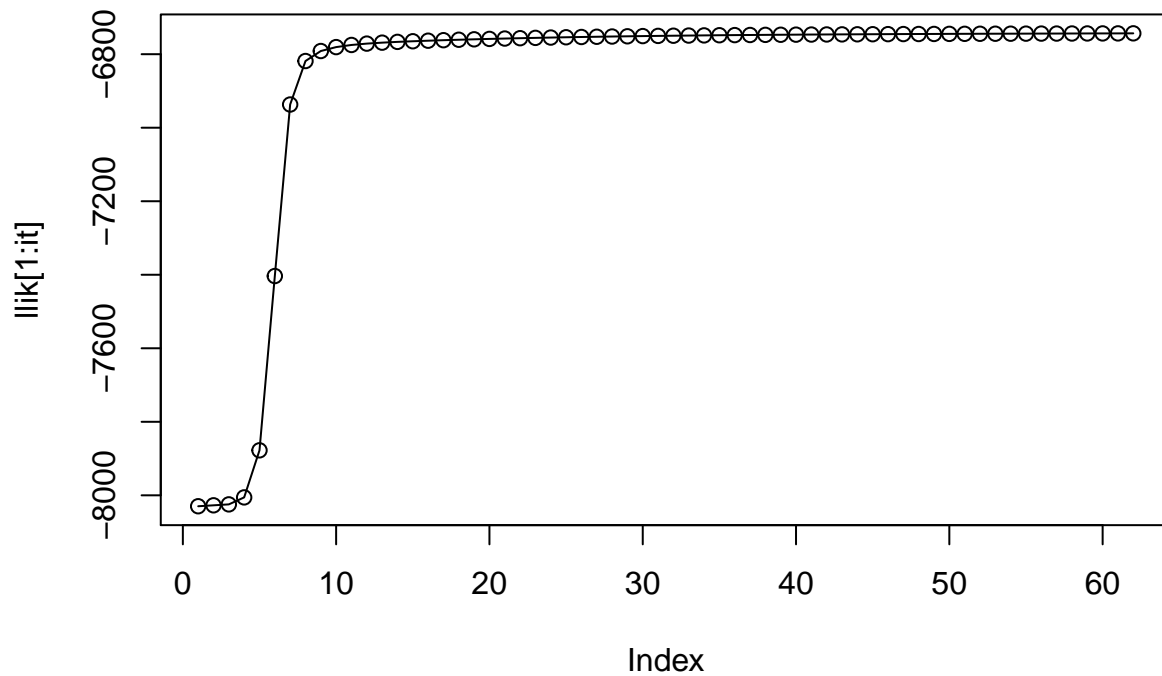
```
#em3<-EM_algorithm(K=4)

# ggplot()+geom_point(aes(x=c(1:length(em1)),y=em1),color="blue")+
#   geom_point(aes(x=c(1:length(em2)),y=em2),color="red")+
#   geom_point(aes(x=c(1:length(em3)),y=em3),color="grey")+
#   xlab("Number of Iterations")+ylab("Expected Log Likelihood")
```