

# main1

Group4

4/9/2018

## Step 0: Load the packages

```
#if (!require("")) install.packages("")  
  
#default the wd to the fold this rmd file exists
```

## Step 1: Load and process the data

```
MS_train    <- read.csv("../data/data_sample/MS_sample/data_train.csv")  
MS_test     <- read.csv("../data/data_sample/MS_sample/data_test.csv")  
movie_train <- read.csv("../data/data_sample/eachmovie_sample/data_train.csv")  
movie_test  <- read.csv("../data/data_sample/eachmovie_sample/data_test.csv")
```

## Step 2 : Transformation

Convert the original dataset to a matrix which rows represents users and columns represents items For dataset 1 (MS), we assign 0 to those items which users never visited. For dataset 2 (Movie), we assign NA to those items which users never rated.

```
source("../lib/MemoryBased.R")  
  
MS_train <- Transform_ms(MS_train)  
MS_test  <- Transform_ms(MS_test)  
# save(MS_train, file = "../output/MS_train.RData")  
# save(MS_test, file = "../output/MS_test.RData")  
  
movie_train <- Transform_m(movie_train)  
movie_test  <- Transform_m(movie_test)  
# save(movie_train, file = "../output/movie_train.RData")  
# save(movie_test, file = "../output/movie_test.RData")
```

## Memory-based Algorithm

## Step 3 : Similarity Weight

Pearson Correlation & Mean-square-difference & SimRank

```

load("../output/MS_train.RData")
load("../output/MS_test.RData")
load("../output/movie_train.RData")
load("../output/movie_test.RData")

#For dataset 1 (MS)

#Pearson Correlation
#pearson_corr_MS <- pearson_corr(MS_train)
#save(pearson_corr_MS, file = "pearson_corr_MS.RData")

#Mean-square-difference
#MSD_w_1 <- MSD_Weight(MS_train)
#save(MSD_w_1, file = "../output/MSD_Data1.RData")

#SimRank
#SR_MS <- simrank(MS_train)
#save(SR-MS, file = "../output/simrank_MS_train")

#For dataset 2 (Movie)

#Pearson Correlation
#pearson_corr_movie <- pearson_corr(movie_train)
#save(pearson_corr_movie, file = "pearson_corr_movie.RData")

#Mean-square-difference
#MSD_w_2 <- MSD_Weight(movie_train)
#save(MSD_w_2, file = "../output/MSD_Data2.RData")

```

# No Variance Weighting

## Step 4: Selecting Neighbors

```

# Implementation on Dataset 1
#BNN_1 <- Select_BNN(MSD_w_1, 20)
#save(BNN_1, file = "../output/BNN_Data1.RData")

#BNN_SR <- Select_BNN(SR_MS, 20)
#save(BNN_SR, file = "../output/BNN_SR_Data1.RData")

# Implementation on Dataset 2
#BNN_2 <- Select_BNN(MSD_w_2, 20)
#save(BNN_2, file = "../output/BNN_Data2.RData")

```

# Step5 : Valuation

```
# Implementation on Dataset 1: ranked scoring

## mean-squared-difference + best-n neighbors
# ZScore_mat_1 <- ZScore_Mat(MSD_w_1, BNN_1, MS_train, MS_test)
# ms_msd_bnn_rs <- Rank_Score(ZScore_mat_1, MS_test)
#save(ZScore_mat_1, file = "../output/ZScore_Data1.RData")

# Implementation on Dataset 2: MAE

## mean-squared-difference + best-n neighbors
# ZScore_mat_2 <- ZScore_Mat(MSD_w_2, BNN_2, movie_train, movie_test)
# movie_msd_bnn_mae <- MAE(ZScore_mat_2, movie_test)

#save(ZScore_mat_2, file = "../output/ZScore_Data2.RData")
```

## Model-based Algorithm

### Step 3: Cluster Model

```
load("../output/movie_train.RData")
load("../output/movie_test.RData")
train <- movie_train
test <- movie_test

N <- nrow(train)
M <- ncol(train)

user <- rownames(train)
movie <- colnames(train)

### cluster model
em_fun <- function(data, C, thres){
  #Input: train_data, number of classes, threshold to determine convergence
  #Output: parameters for cluster models:
  #   mu: probability of belonging to class c, vector
  #   gamma: probability of scores for a movie given the class, 3 dimentions

  #=====
  # Step 1 - initialization
  #=====
  set.seed(2)
```

```

mu <- runif(C)
mu <- mu/sum(mu)
gamma <- array(NA,c(M,C,6)) #each matrix represents a class
#the i,j-th element means the probability of rating jth movie with score i in the c
lass
for(m in 1:M){
  for(c in 1:C){
    gamma[m,c,] <- runif(6)
    gamma[m,c,] <- gamma[m,c,]/sum(gamma[m,c,])
  }
}

v <- array(0, c(M,N,7))
for(k in 1:6){
  v[, ,k] <- ifelse(t(data)==(k-1), 1, 0)
  v[, ,k] <- ifelse(is.na(v[, ,k]), 0, v[, ,k])
  v[, ,7] <- v[, ,7] + v[, ,k]
}

mu_new <- mu
gamma_new <- gamma

## Iterations based on the stop criterion
thres1 <- 1000
thres2 <- 1000
thres1_new <- 0
thres2_new <- 0
count <- 0

while((thres1>thres|thres2>thres)&(abs(thres1-thres1_new)>thres|abs(thres2-thres2_n
ew)>thres))
{
  count <- count + 1
  print(paste0("iteration = ", count))

  thres1_new <- thres1
  thres2_new <- thres2

  mu <- mu_new
  gamma <- gamma_new

  #=====
  # Step 2 - Expectation
  #=====
  #expectation pi with rows meaning classes and columns meaning users
  phi <- matrix(0, C, N)

  for(k in 1:6){
    phi <- phi + t(log(gamma[, ,k]))%*%v[, ,k]
  }
}

```

```

phi <- phi-rep(colMeans(phi),each=C)

for(c in 1:C){
  phi[c,] <- mu[c]*exp(phi)[c,]
}
phi <- ifelse(phi == rep(colSums(phi),each=C), 1, phi/rep(colSums(phi), each=C))

#####
# Step 3 - Maximization
#####
mu_new <- rowSums(phi)/N #update mu vector

for(k in 1:6){
  gamma_new[,k] <- v[,k]%*%t(phi)/v[,7]%*%t(phi) #update gamma
}

gamma_new[gamma_new == 0] <- 10^(-100)
if(sum(is.na(gamma_new)) != 0){
  is_zero <- which(is.na(gamma_new))
  gamma_new[is_zero] <- rep(1/6, length(is_zero))
}

## Check convergence
thres1 <- mean(abs(mu_new - mu)) #mean absolute difference of mu
thres2 <- 0
for(c in 1:C){
  thres2 <- max(thres2,norm(as.matrix(gamma_new[,c,] - gamma[,c,]), "O"))
}
print(paste0("threshold1 = ", thres1, " threshold2 = ", thres2))
}
return(list(mu = mu, gamma = gamma))
}

#predict score estimate function
cm_predict <- function(train_df, test_df, par){
  set.seed(2)
  mu <- par$mu
  gamma <- par$gamma
  C <- length(mu)

  w <- array(0, c(M,N,7))
  for(k in 1:6){
    w[,k] <- ifelse(t(train_df)==k, 1, 0)
    w[,k] <- ifelse(is.na(w[,k]),0,w[,k])
  }
  w[,7] <- ifelse(!is.na(t(test_df)), 1, 0)

  ##using Naive Bayes formula
  prob <- array(0,c(N,M,7))
  prob_mu <- matrix(mu, N, C, byrow = TRUE)

```

```

phi <- matrix(0, C, N)
for(k in 1:6){
  phi <- phi + t(log(gamma[, ,k]))%*%w[, ,k]
}

phi <- exp(phi)

den <- matrix(diag(prob_mu%*%phi), N, M, byrow=FALSE)  #denominator in equation (2)
) of cluster model notes

for(k in 1:6){
  print(paste0("k = ", k))

  num <- (t(phi)*prob_mu)%*%t(gamma[, ,k]) #numerator in equation (2) of cluster model notes
  prob[, ,k] <- ifelse(num==den & num == 0, runif(1)/6, num/den)
  prob[, ,7] <- prob[, ,7] + k*prob[, ,k]
}
return(prob[, ,7]*t(w[, ,7]))
}

### 5-fold cross validation to find best class number C
set.seed(2)
K <- 5
n <- ncol(train)
m <- nrow(train)
n.fold <- floor(n/K)
m.fold <- floor(m/K)
s <- sample(rep(1:K, c(rep(n.fold, K-1), n-(K-1)*n.fold)))
s1 <- sample(rep(1:K, c(rep(m.fold, K-1), m-(K-1)*m.fold)))

c_list <- c(2,3,6,12)
validation_error <- matrix(NA, K, length(c_list))

train_data <- data.frame(matrix(NA, N, M))
colnames(train_data) <- movie
rownames(train_data) <- user

test_data <- data.frame(matrix(NA, N, M))
colnames(test_data) <- movie
rownames(test_data) <- user

#cv 5 folds
#calculate cv error
# for(i in 1:K){
#   train_data[s1 != i, ] <- train[s1 != i, ]
#   train_data[s1 == i, s != i] <- train[s1==i, s != i]

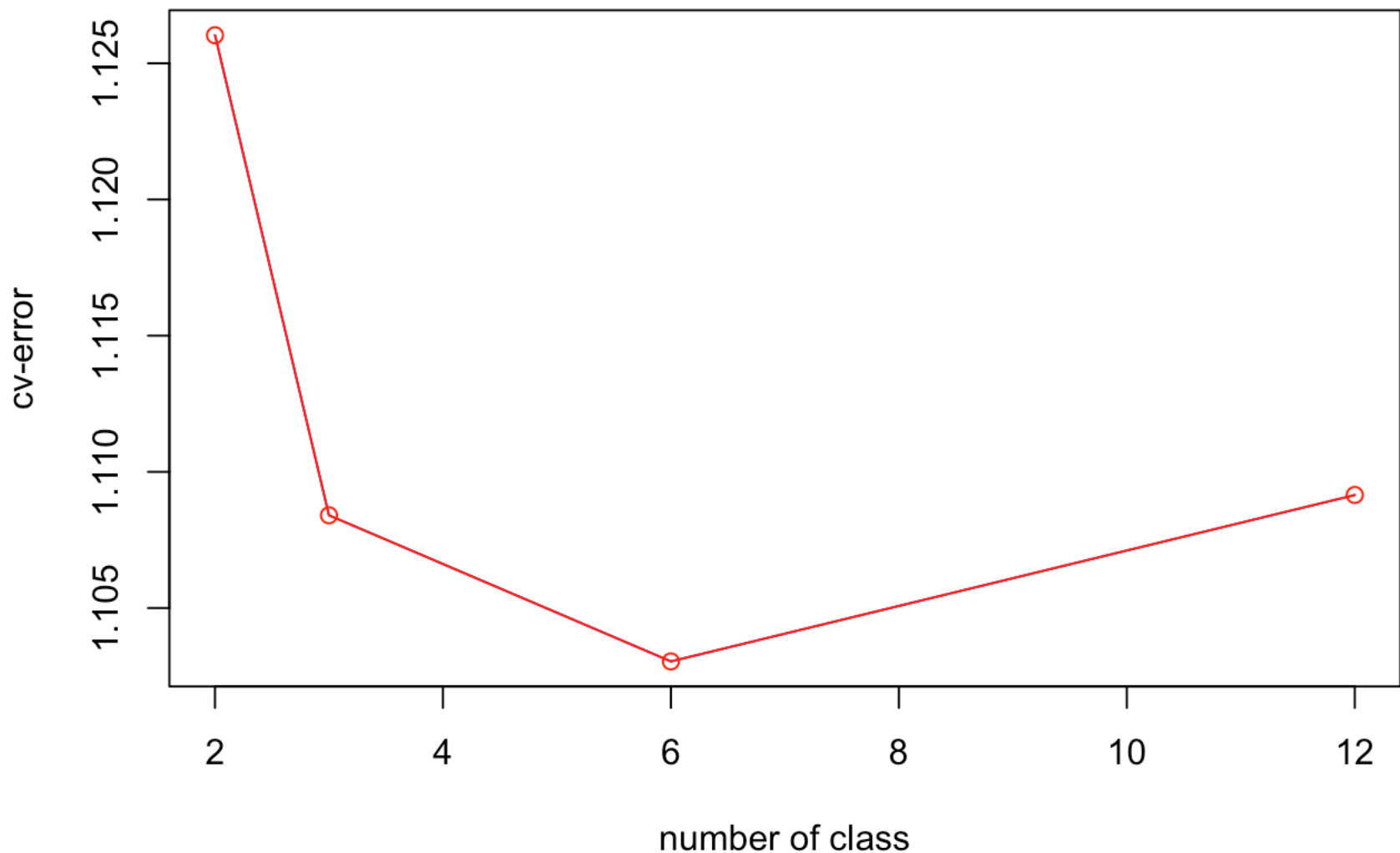
```

```

# test_data[s1 == i,s == i] <- train[s1 == i ,s == i]
# #write.csv(train_data,paste0("../output/cluster_model_subtrain.csv"))
# #write.csv(test_data,paste0("../output/cluster_model_validation.csv"))
#
# estimate_data <- test_data
#
# for(c in 1:length(c_list)){
#   cm_par <- em_fun(data = train_data, C = c_list[c], thres = 0.05)
#   estimate_data <- cm_predict(train_df = train_data, test_df = test_data, par = c
m_par)
#   validation_error[i,c] <- sum(abs(estimate_data-test_data),na.rm = T)/sum(!is.na
(estimate_data-test_data))
# }
# }

#save(validation_error, file=paste0("../output/validation_err.RData"))
load("../output/validation_err.RData")
cv_error<-colMeans(validation_error)
plot(c_list,cv_error,xlab="number of class",ylab="cv-error",col="blue",type="l")
points(c_list,cv_error,col="red",type="o")

```



```
class = c_list[which.min(cv_error)]  
print(paste("Best class number is", class))
```

```
## [1] "Best class number is 6"
```

```
class <- 6  
  
#best_par <- em_fun(data = train, C = class, thres = 0.01)  
#save(best_par, file = "../output/best_par.RData")  
  
load("../output/best_par.RData")  
  
###estimate scores  
  
#estimate <- cm_predict(train, test, best_par)  
  
#write.csv(estimate, paste0("../output/cluster_model_estimate.csv"))  
  
estimate <- read.csv("../output/cluster_model_estimate.csv")  
estimate <- estimate[,-1]  
  
# MAE of EM algorithm  
coltest <- colnames(test)  
colnames(estimate) <- movie  
  
coltest <- which(is.element(movie, coltest))  
estimate <- estimate[,coltest]  
  
MAE <- function(pred, true){  
  mae <- sum(abs(pred-test),na.rm = T)/sum(!is.na(abs(pred-test)))  
  return(mae)  
}  
  
error_em<- MAE(estimate,test)  
error_em
```

```
## [1] 2.803069
```