

4-13

XU, Xin

## 目录

<b>1</b>	<b>Data Pre-processing</b>	<b>2</b>
1.1	ISBN . . . . .	2
1.2	Romve Some Entries . . . . .	3
1.3	User Groups . . . . .	4
<b>2</b>	<b>book average</b>	<b>4</b>
<b>3</b>	<b>PMF</b>	<b>6</b>
3.1	Implementation . . . . .	7
3.2	Simulation . . . . .	9
3.3	Real data set . . . . .	12
<b>4</b>	<b>Logistic PMF</b>	<b>12</b>
4.1	Simulation . . . . .	13
4.2	Real data set . . . . .	13
<b>5</b>	<b>Constrained PMF</b>	<b>17</b>
5.1	Simulation . . . . .	18
5.2	Real Data . . . . .	19

```
library(dplyr)
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':  
##  
##   filter, lag  
  
## The following objects are masked from 'package:base':  
##  
##   intersect, setdiff, setequal, union  
  
library(Rcpp)
```

## 1 Data Pre-processing

```
Ratings <- read.csv("archive/Ratings.csv")  
Ratings_by_users <- group_by(Ratings, User.ID)  
user.Rating <- Ratings_by_users %>% summarise(  
  num = length(Book.Rating),  
  avg.rating = mean(Book.Rating)  
)  
num.Rating <- group_by(user.Rating, num) %>% summarise(  
  users = length(avg.rating)  
)
```

### 1.1 ISBN

```
N <- max(Ratings$User.ID)  
book_idx <- unique(Ratings$ISBN)  
M <- length(book_idx)  
book_idx <- 1:M  
names(book_idx) <- unique(Ratings$ISBN)  
train <- read.csv("archive/train.csv")  
train$ISBN <- book_idx[train$ISBN]  
test <- read.csv("archive/test.csv")
```

```
test$ISBN <- book_idx[test$ISBN]
```

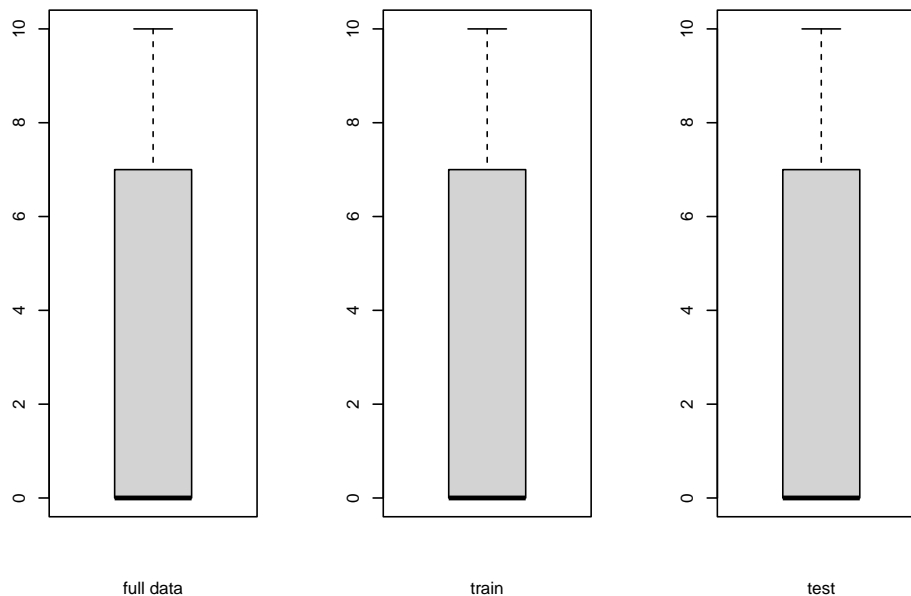
## 1.2 Remove Some Entries

```
train.book.avg <- group_by(train, ISBN) %>% summarise(  
  avg.rating = mean(Book.Rating)  
)  
test <- test[which(test$ISBN %in% train.book.avg$ISBN),]
```

```
## There are 278854 users and 340556 books.
```

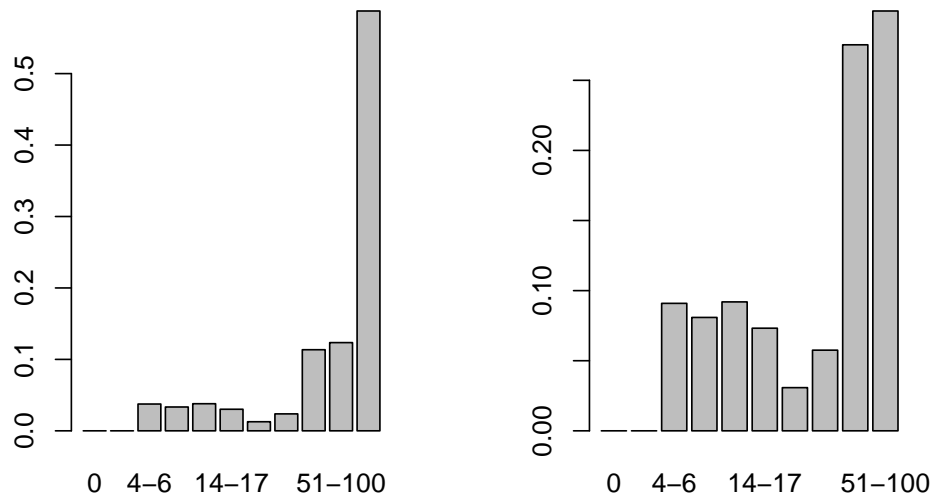
```
## There are 847349 examples in the training set and 240117 examples in the test set.
```

```
par(mfrow = c(1, 3))  
boxplot(Ratings$Book.Rating, xlab = "full data")  
boxplot(train$Book.Rating, xlab = "train")  
boxplot(test$Book.Rating, xlab = "test")
```



### 1.3 User Groups

```
par(mfrow = c(1, 2))
barplot(grp_num/sum(grp_num))
barplot(grp_num[-11]/sum(grp_num[-11]))
```



## 2 book average

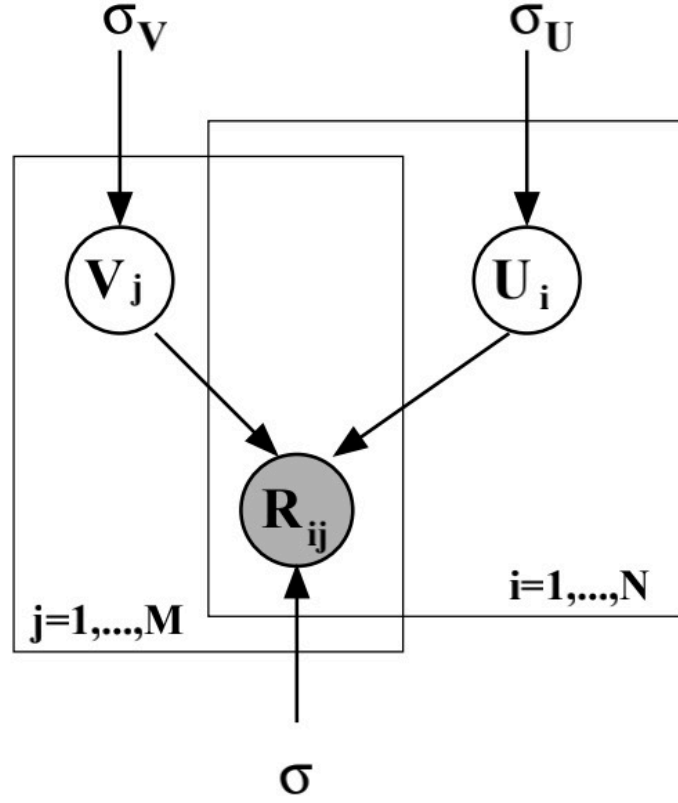
```
train$Book.Rating <- train$Book.Rating/10
test$Book.Rating <- test$Book.Rating/10
train.book.avg <- group_by(train, ISBN) %>% summarise(
  avg.rating = mean(Book.Rating)
)
book_avg <- train.book.avg$avg.rating
names(book_avg) <- train.book.avg$ISBN
book_idx1 <- names(book_avg)
Rcpp::sourceCpp("cpp/pred_bookavg.cpp")
bookavg_rslt <- pred_bookavg(test, book_avg, book_idx1)
save(bookavg_rslt, file = "bookavg_rslt.Rda")
```

```
## MAE: 0.3266 RMSE: 0.4104
```

```
head(cbind(test[, -1], bookavg_rslt$prediction), 10)
```

```
##      User.ID ISBN Book.Rating bookavg_rslt$prediction
## 1         8 8974           0             0.4
## 2         8 8975           0             0.2
## 3         8 8979           0             0.0
## 4         8 8982           5             0.0
## 5         8 8984           0             0.5
## 6        14 453           0             0.1
## 7        17 8994           0             0.3
## 8        17 8998           0             0.3
## 11       53 9017           3             0.5
## 13       67 9028           0             0.1
```

## 3 PMF



$$p(R|U, V, \sigma^2) = \prod_{i=1}^N \prod_{j=1}^M \left[ \mathcal{N}(R_{ij} | U_i^T V_j, \sigma^2) \right]^{I_{ij}},$$

$$p(U | \sigma_U^2) = \prod_{i=1}^N \mathcal{N}(U_i | 0, \sigma_U^2 \mathbf{I}), \quad p(V | \sigma_V^2) = \prod_{j=1}^M \mathcal{N}(V_j | 0, \sigma_V^2 \mathbf{I}).$$

$$E = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^M I_{ij} (R_{ij} - U_i^T V_j)^2 + \frac{\lambda_U}{2} \sum_{i=1}^N \|U_i\|_{Fro}^2 + \frac{\lambda_V}{2} \sum_{j=1}^M \|V_j\|_{Fro}^2,$$

### 3.1 Implementation

```

4 // [[Rcpp::export]]
5 List grad_pmf(DataFrame df,
6               NumericMatrix Ut,
7               NumericMatrix Vt,
8               double u_lam,
9               double v_lam) {
10   NumericVector user = as<NumericVector>(df["User.ID"]);
11   NumericVector book = as<NumericVector>(df["ISBN"]);
12   NumericVector rate = as<NumericVector>(df["Book.Rating"]);
13   NumericVector u, v, du, dv;
14   NumericMatrix grad_Ut = clone(Ut), grad_Vt = clone(Vt);
15   grad_Ut = u_lam * grad_Ut;
16   grad_Vt = v_lam * grad_Vt;
17   int i, j, r;
18   double aux;
19   for(int k = 0; k < df.rows(); k++){
20     i = user[k];
21     j = book[k];
22     r = rate[k];
23     u = Ut.row(i-1);
24     v = Vt.row(j-1);
25     //du = grad_Ut.row(i-1);
26     //dv = grad_Vt.row(j-1);
27     aux = 0;
28     for(int l = 0; l < u.size(); l++){
29       aux += u[l]*v[l];
30     }
31     aux -= r;
32     grad_Ut.row(i-1) = aux*v + grad_Ut.row(i-1);
33     grad_Vt.row(j-1) = aux*u + grad_Vt.row(j-1);
34   }
35   return List::create(
36     Named("Ut") = grad_Ut,
37     Named("Vt") = grad_Vt
38   );
39 }

```

```

4 double loss_pmf(DataFrame df,
5                 NumericMatrix Ut,
6                 NumericMatrix Vt,
7                 double u_lam,
8                 double v_lam) {
9     NumericVector user = as<NumericVector>(df["User.ID"]);
10    NumericVector book = as<NumericVector>(df["ISBN"]);
11    NumericVector rate = as<NumericVector>(df["Book.Rating"]);
12    NumericVector u, v;
13    int i, j, r;
14    double sum = 0, sum_u = 0, sum_v = 0;
15    double aux;
16    for(int k = 0; k < df.rows(); k++){
17        i = user[k];
18        j = book[k];
19        r = rate[k];
20        u = Ut.row(i-1);
21        v = Vt.row(j-1);
22        aux = 0;
23        for(int l = 0; l < u.size(); l++){
24            aux += u[l]*v[l];
25        }
26        sum += (r-aux)*(r-aux);
27        //std::printf("%f\n", sum);
28    }
29    for(int k = 0; k < Ut.size(); k++){
30        sum_u += Ut[k]*Ut[k];
31        sum_v += Vt[k]*Vt[k];
32    }
33    sum_u *= u_lam;
34    sum_v *= v_lam;
35    sum += (sum_u + sum_v);
36    sum /= 2;
37    return sum;
38 }

```



```

4 List pred_pmf(DataFrame df,
5               NumericMatrix Ut,
6               NumericMatrix Vt) {
7   NumericVector user = as<NumericVector>(df["User.ID"]);
8   NumericVector book = as<NumericVector>(df["ISBN"]);
9   NumericVector rate = as<NumericVector>(df["Book.Rating"]);
10  NumericVector u, v;
11  IntegerVector pred(rate.size());
12  int i, j, r;
13  double sum = 0, sum1 = 0;
14  double aux;
15  for(int k = 0; k < df.rows(); k++){
16    i = user[k];
17    j = book[k];
18    r = rate[k];
19    u = Ut.row(i-1);
20    v = Vt.row(j-1);
21    aux = 0;
22    for(int l = 0; l < u.size(); l++){
23      aux += u[l]*v[l];
24    }
25    pred[k] = floor(aux);
26    if(pred[k] > 10){
27      pred[k] = 10;
28    }
29    if(pred[k] < 0){
30      pred[k] = 0;
31    }
32    sum += (r-aux)*(r-aux);
33    if(r-aux > 0){
34      sum1 += (r-aux);
35    }else{
36      sum1 += (aux-r);
37    }
38  }
39  return List::create(
40    Named("prediction") = pred,

```

### 3.2 Simulation

```

set.seed(123)
A <- matrix(rnorm(1e4*4, sd = sqrt(1)), ncol = 1e4)
B <- matrix(rnorm(1e2*4, sd = sqrt(1)), ncol = 1e2)
E <- matrix(rnorm(1e4*1e2, sd = sqrt(1e-2)), ncol = 1e2)

```

```

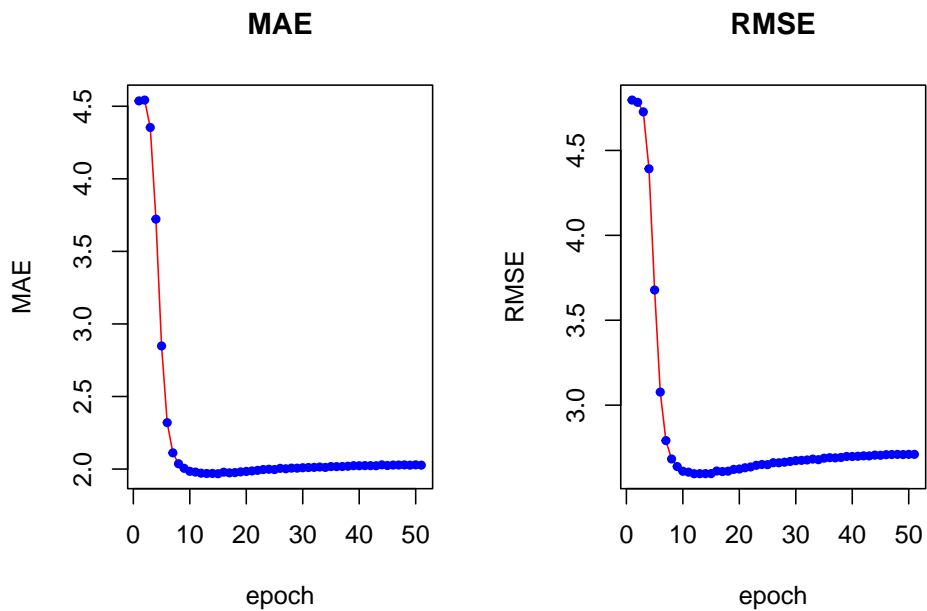
aux <- t(A) %*% B
aux <- 1/(1+exp(-aux))
obs <- aux + E
obs <- ifelse(obs < 0, 0, obs)
obs <- ifelse(obs > 1, 1, obs)
obs <- floor(10*obs)/10
obs <- 10*obs
W <- matrix(sample(0:2, 1e4*1e2, prob = c(0.9, 0.07, 0.03), replace = TRUE), ncol = 1e2)
strain <- data.frame(User.ID = c(), ISBN = c(), Book.Rating = c())
stest <- data.frame(User.ID = c(), ISBN = c(), Book.Rating = c())
for(i in 1:1e4){
  for(j in 1:1e2){
    if(W[i, j] == 1){#train
      strain <- rbind(strain, c(i, j, obs[i, j]))
    }
    if(W[i, j] == 2){#test
      stest <- rbind(stest, c(i, j, obs[i, j]))
    }
  }
}
colnames(strain) <- c("User.ID", "ISBN", "Book.Rating")
colnames(stest) <- c("User.ID", "ISBN", "Book.Rating")
strain.book.avg <- group_by(strain, ISBN) %>% summarise(
  avg.rating = mean(Book.Rating)
)
stest <- stest[which(stest$ISBN %in% strain.book.avg$ISBN),]
write.csv(strain, file = "archive/strain.csv")
write.csv(stest, file = "archive/stest.csv")

strain <- read.csv("archive/strain.csv")
stest <- read.csv("archive/stest.csv")
strain$Book.Rating <- strain$Book.Rating/10
stest$Book.Rating <- stest$Book.Rating/10

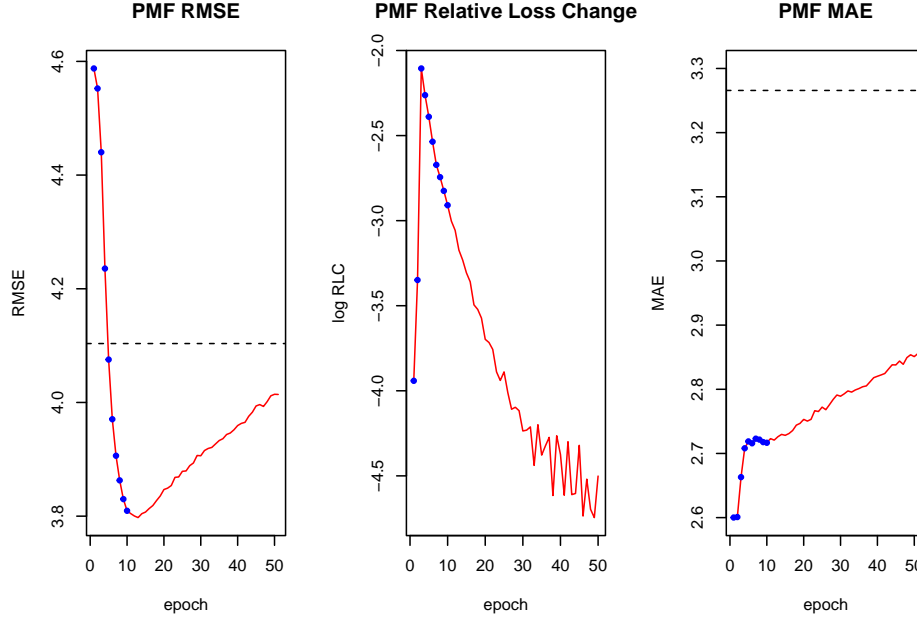
```

```
strain.book.avg <- group_by(strain, ISBN) %>% summarise(  
  avg.rating = mean(Book.Rating)  
)  
sbook_avg <- strain.book.avg$avg.rating  
names(sbook_avg) <- strain.book.avg$ISBN  
sbook_idx1 <- names(sbook_avg)  
Rcpp::sourceCpp("cpp/pred_bookavg.cpp")  
sbookavg_rslt <- pred_bookavg(stest, sbook_avg, sbook_idx1)
```

```
## MAE: 0.2638 RMSE: 0.3098
```



### 3.3 Real data set



##

## PMF MAE: 2.726 RMSE: 3.7974

##

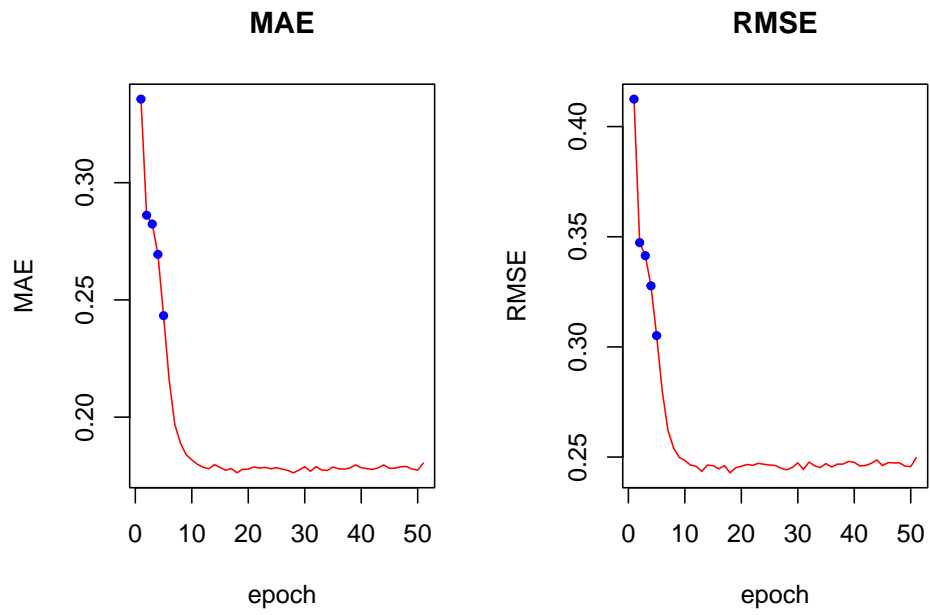
## Book Average MAE: 3.2657 RMSE: 4.1037

## 4 Logistic PMF

$$p(R|U, V, \sigma^2) = \prod_{i=1}^N \prod_{j=1}^M \left[ \mathcal{N}(R_{ij} | g(U_i^T V_j), \sigma^2) \right]^{I_{ij}}.$$

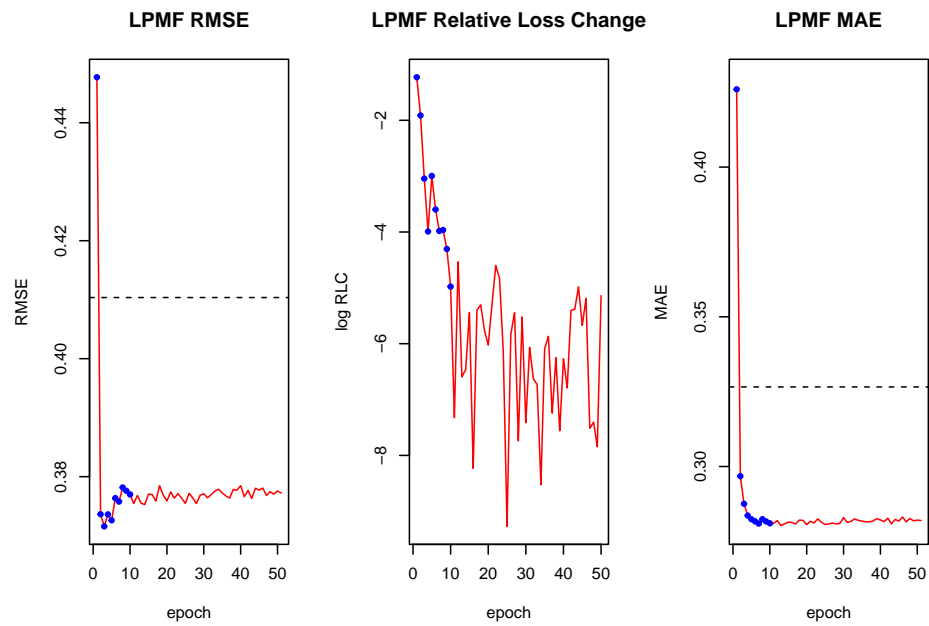
We map the ratings  $1, \dots, K$  to the interval  $[0, 1]$  using the function  $t(x) = (x - 1)/(K - 1)$ , so that the range of valid rating values matches the range of predictions our model makes. Minimizing

### 4.1 Simulation



### 4.2 Real data set

## 4.2.1 Momentum



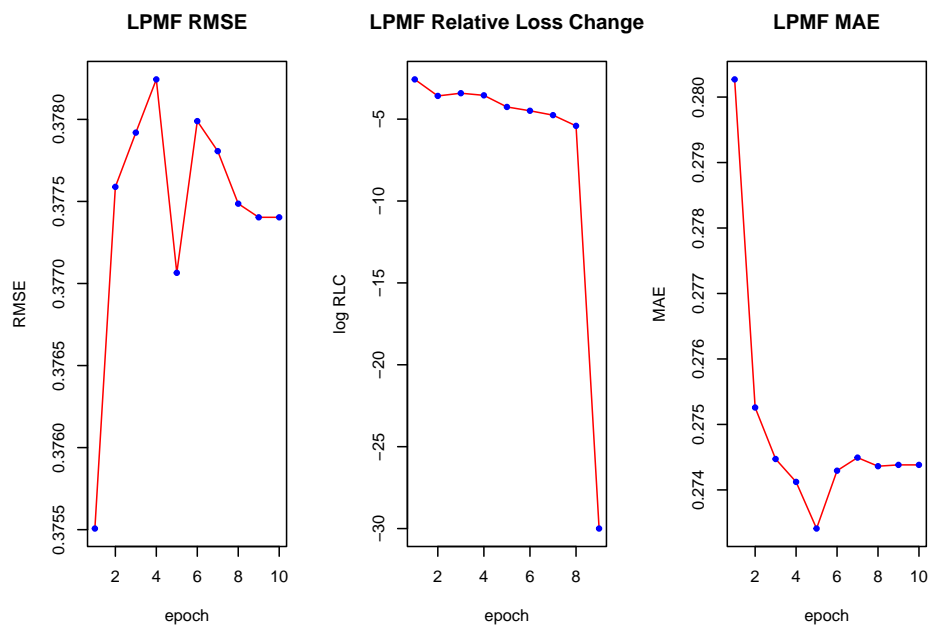
##

## LPMF Momentum      MAE: 0.2803      RMSE: 0.3755

##

## Book Average      MAE: 0.3266      RMSE: 0.4104

## 4.2.2 SGD



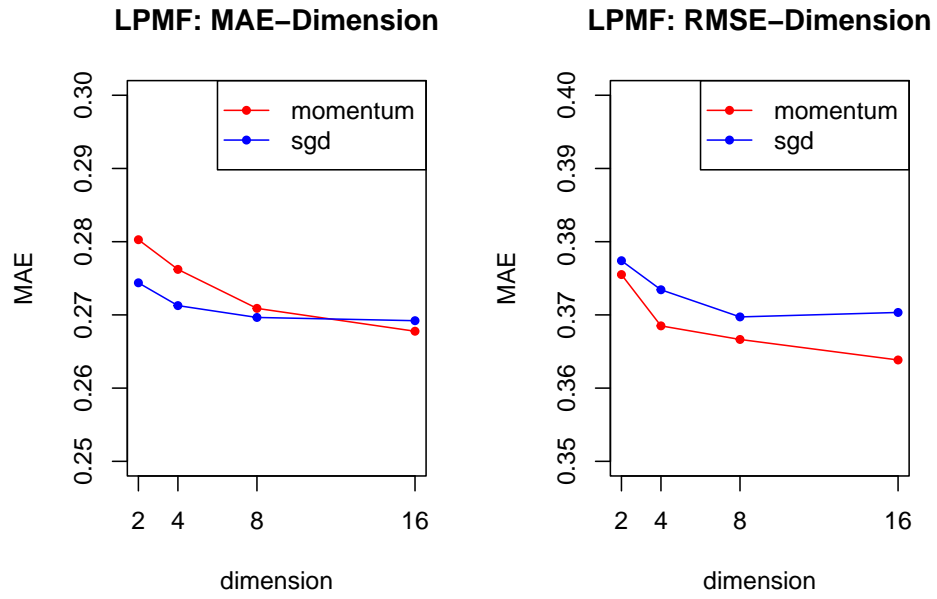
##

## LPMF SGD      MAE: 0.2744      RMSE: 0.3774

##

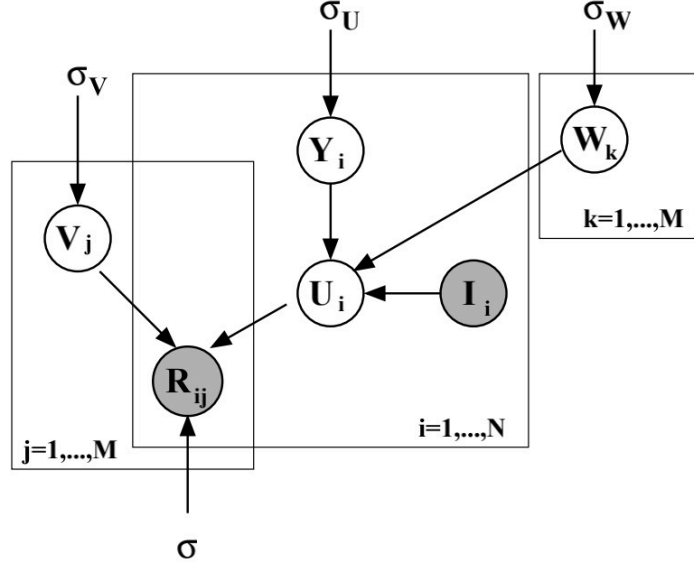
## Book Average      MAE: 0.3266      RMSE: 0.4104

## 4.2.3 Dimension of the Feature Vectors





## 5 Constrained PMF



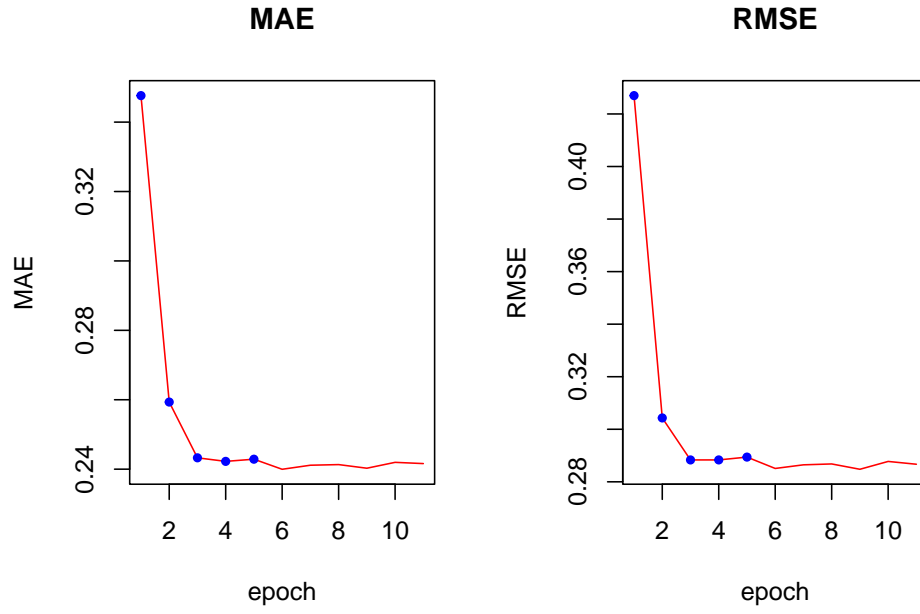
$$p(W|\sigma_W) = \prod_{k=1}^M \mathcal{N}(W_k|0, \sigma_W^2 \mathbf{I}).$$

$$U_i = Y_i + \frac{\sum_{k=1}^M I_{ik} W_k}{\sum_{k=1}^M I_{ik}}.$$

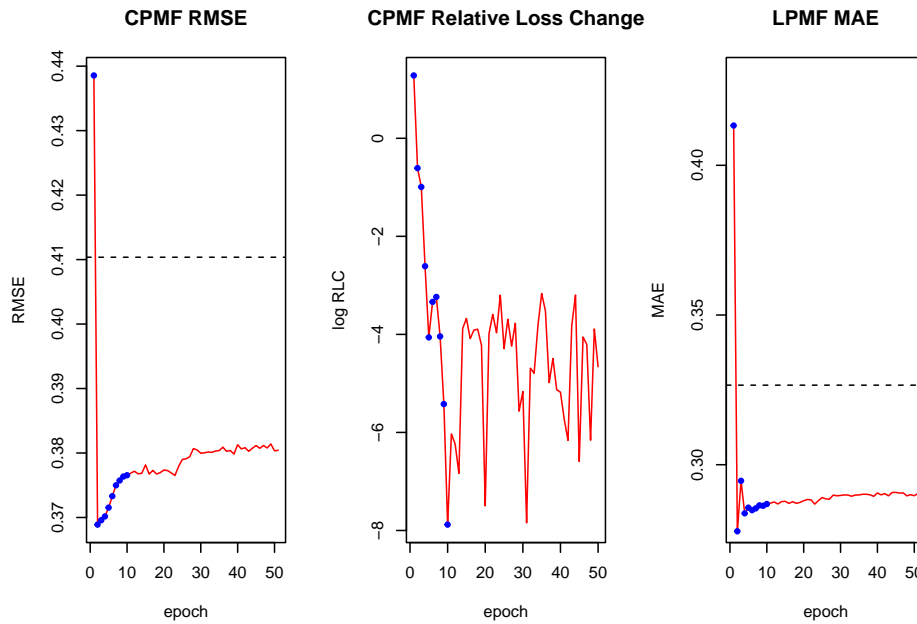
$$p(R|Y, V, W, \sigma^2) = \prod_{i=1}^N \prod_{j=1}^M \left[ \mathcal{N}(R_{ij} | g([Y_i + \frac{\sum_{k=1}^M I_{ik} W_k}{\sum_{k=1}^M I_{ik}}]^T V_j), \sigma^2) \right]^{I_{ij}}.$$

$$\begin{aligned}
E = & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^M I_{ij} \left( R_{ij} - g \left( \left[ Y_i + \frac{\sum_{k=1}^M I_{ik} W_k}{\sum_{k=1}^M I_{ik}} \right]^T V_j \right) \right)^2 \\
& + \frac{\lambda_Y}{2} \sum_{i=1}^N \| Y_i \|_{Fro}^2 + \frac{\lambda_V}{2} \sum_{j=1}^M \| V_j \|_{Fro}^2 + \frac{\lambda_W}{2} \sum_{k=1}^M \| W_k \|_{Fro}^2,
\end{aligned}$$

### 5.1 Simulation



## 5.2 Real Data



```
## CPMF      MAE: 0.2778      RMSE: 0.3689

##
## LPMF      MAE: 0.2762      RMSE: 0.3685

##
## Book Average      MAE: 0.3266      RMSE: 0.4104

##
## > dim = 2
##
## > Y <- matrix(rnorm(N * dim), nrow = dim)
##
## > W <- matrix(rnorm(M * dim), nrow = dim)
##
## > aaa <- user_cpmf(read, t(Y), t(W))

##
## > dim = 2
```

```
##  
## > Y <- matrix(rnorm(N * dim), nrow = dim)  
##  
## > V <- matrix(rnorm(M * dim), nrow = dim)  
##  
## > W <- matrix(rnorm(M * dim), nrow = dim)  
##  
## > Ut <- user_cpmf(read, t(Y), t(W))  
##  
## > pred_cpmf(test, Ut, t(V))$mae  
## [1] 2.744441  
  
##  
## > dim = 2  
##  
## > U <- matrix(rnorm(N * dim), nrow = dim)  
##  
## > V <- matrix(rnorm(M * dim), nrow = dim)  
##  
## > pred_lpmf(test, t(U), t(V))$mae  
## [1] 2.745198
```

