

Brain data and Music data application.

Chanwoo Lee

Dec 15, 2019

1 Algorithm modification for Brain data.

When I applied our algorithm in the brain data, there are some problems I had to cope with. The first problem was memory issue. My codes kept getting killed by the statistic server saying that it is out of the memory. Main reason for this error is that our algorithm generates really large matrices whose dimension increase to 628864×2000 (the size is almost 10 GB) to update the core tensor. This matrix is made to get a gradient value for each iteration and has the following form.

$$W = A_3 \otimes A_2 \otimes A_1.$$

However, this matrix is needed to find a gradient and Hessian of the loss function. I tried to find a way not using this large matrix to calculate a gradient but there is no way to get gradient without W because gradient of the core tensor has the following elements.

$$\sum_{k \in [K]} \sum_{i \in [d_1 d_2 d_3]} W_i q_{ik}. \quad (1)$$

where W_i is the i -th row of the matrix W and q_{ik} is a constant. Instead, I modified the gradient algorithm not saving whole W .

The first trial is to calculate each row of W and add each above term without calculating whole W and saving it. For this method, I used ‘for’ loop and made new algorithm to do parallel computing. However, this approach becomes so slow that I cannot implement this algorithm for the brain data size: the server spent 3 days to get an output but failed.

The second approach is from the fact that it doesn’t take too much space to save matrices

W_1, W_2, W_3 where

$$W_1 = A_3 \otimes A_2, \quad W_2 = A_3 \otimes A_1, \quad W_3 = A_2 \otimes A_1$$

. To get (1), I saved the matrix W_3 and do parallel computing based on the following formula.

$$\sum_{k \in [K]} \sum_{j \in d_3} ([A_3]_j \otimes W_3)^T q_{jk}. \quad (2)$$

where $[A_3]_j$ is j -th row of A_3 and q_{jk} is a constant vector. Although this approach is still slower than saving the whole matrix W , running times is improved a lot from the first method and I put this algorithm in the server now. Also, I changed to use ‘BFGS’ method instead of using Newton’s method not to use Hessian matrix. The following is the modified algorithm to get cost function and gradient.

```

1 ##### cost function with arbitrary k #####
2 # previous cost function
3 h1 = function(A_1,W1,ttnsr,omega){
4   k = length(omega)
5   thet = W1%%c(A_1)
6   p = matrix(nrow = length(thet),ncol = k)
7   for (i in 1:k) {
8     p[,i] = as.numeric(logistic(thet + omega[i]))
9   }
10  p = cbind(p,rep(1,length(thet)))-cbind(rep(0,length(thet)),p)
11  l = lapply(1:(k+1),function(i) -log(p[which(c(ttnsr)==i),i]))
12  return(sum(unlist(l)))
13 }
14
15 # current cost function not saving W
16 hc = function(A_1,A_2,A_3,C,ttnsr,omega){
17   k = length(omega)
18   thet = c(ttl(C,list(A_1,A_2,A_3),ms=1:3)@data)
19   p = matrix(nrow = length(thet),ncol = k)
20   for (i in 1:k) {
21     p[,i] = as.numeric(logistic(thet + omega[i]))
22   }

```

```

23 p = cbind(p,rep(1,length(thet)))-cbind(rep(0,length(thet)),p)
24 l = lapply(1:(k+1),function(i) -log(p[which(c(ttnsr)==i),i]))
25 return(sum(unlist(l)))
26 }
27
28
29 ##### gradient with arbitrary k #####
30 # previous gradient function
31 g1 = function(A_1,W,ttnsr,omega){
32   k = length(omega)
33   thet =W%*%c(A_1)
34   p = matrix(nrow = length(thet),ncol = k)
35   for (i in 1:k) {
36     p[,i] = as.numeric(logistic(thet + omega[i]))
37   }
38   q = matrix(nrow = length(thet),ncol = k+1)
39   q[,1] <- p[,1]-1
40   for (i in 2:k) {
41     q[,i] <- (p[,i]*(1-p[,i])-p[,i-1]*(1-p[,i-1]))/(p[,i-1]-p[,i])
42   }
43   q[,k+1] <- p[,k]
44   l <- Reduce("+",lapply(1:(k+1),function(i) apply(rbind(W[which(c(ttnsr)
45     ==i),]),)*q[which(c(ttnsr)==i),i],2,sum)))
46   return(l)
47 }
48
49 # current gradient function not saving W
50 gc = function(A_1,A_2,A_3,C,ttnsr,omega){
51   k = length(omega)
52   thet = c(ttl(C,list(A_1,A_2,A_3),ms=1:3)@data)
53   p = matrix(nrow = length(thet),ncol = k)
54   for (i in 1:k) {
55     p[,i] = as.numeric(logistic(thet + omega[i]))
56   }
57   q = matrix(nrow = length(thet),ncol = k+1)
58   q[,1] <- p[,1]-1
59   for (i in 2:k) {
60     q[,i] <- (p[,i]*(1-p[,i])-p[,i-1]*(1-p[,i-1]))/(p[,i-1]-p[,i])

```

```

61  q[,k+1] <- p[,k]
62  W = kronecker(A_2,A_1)
63  d = c(nrow(A_1),nrow(A_2),nrow(A_3))
64  cl <- makeCluster(20)
65  registerDoParallel(cl)
66  l <- foreach(j = 1:d[3],.combine = "+") %dopar% {
67    Reduce("+",lapply(1:(k+1),function(i) apply(rbind(kronecker(A_3[j,,
drop = F],W)[which(c(ttnsr)[(d[1]*d[2]*(j-1)+1):(d[1]*d[2]*(j-1)+d[1]*d
[2])]==i),])*)
68                                                                    q[which(c(ttnsr)[(d[1]*d
[2]*(j-1)+1):(d[1]*d[2]*(j-1)+d[1]*d[2])]==i)+d[1]*d[2]*(j-1),i],2,sum)
))
69  }
70  stopCluster(cl)
71  return(l)
72 }

```

2 Current result for the brain data.

For small rank, saving whole W matrix and calculations based on it require almost 20 GB memory which is doable in the server. So I got the result from the rank (5,5,5) to the rank (5,5,20) where the first rank is fixed as 5. Also, I obtained the result from the rank (6,6,4) to (21,21,4) where the second rank is fixed as 4. The result is in the following figures. Based on this figure, my guess for the best rank according to BIC would be around (23,23,6) because in the second plot we can expect that decreasing rate becomes almost zero around the first rank being 21. Also, from the first plot, we can expect that the best second rank would not be large.

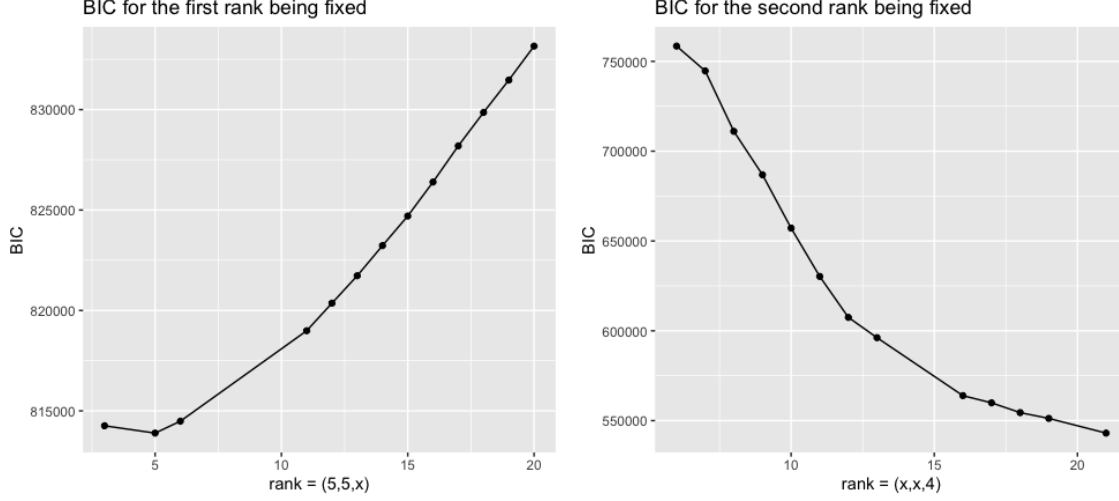


Figure 1: BIC result for one rank being fixed.

3 Questions and discussions for the Brain data application

1. Expected rank size issue: Based on the outputs I get from the server, my guess of the rank is around (23,23,6). However, if you consider the dimension of the brain data (68,68,136), I do not know how to interpret this rank size based on the meaning of each mode.
2. Symmetry issue: From the data output with the rank (21,21,4), I realized A_1 and A_2 are not necessarily the same. The main reason for this is that A_1 and A_2 are not identifiable which means I can manipulate A_1 and A_2 without changing $\Theta = \mathcal{C} \times_1 A_1 \times_2 A_2 \times_3 A_3$. However, if I check $\Theta_{\cdot i \cdot} - \Theta_{i \cdot \cdot}$ where $\Theta_{\cdot i \cdot}$ is a slice with the second mode being fixed as i , the maximum entry for all $i \in [68]$ is 0.91. Considering that the maximum entry of Θ is 145, We can say that Θ is almost symmetric. The proportion between the two is 0.006239625
3. Clustering issue: Based on the following formula, if a row of $(A_1)_{i \cdot \cdot}$ and a row of $(A_1)_{i_2 \cdot \cdot}$ is close enough, $\theta_{i_1 j k}$ and $\theta_{i_2 j k}$ have close value. However, it is not necessary condition and since the range of C is much larger than the range of A_1 , we need to consider the

value of C_{i_1jk} and C_{i_2jk} . So I am thinking about do clustering with weight C_{jk}

$$\theta_{ijk} = \sum_{p=1}^{68} \sum_{q=1}^{68} \sum_{r=1}^{136} C_{pqr} (A_1)_{ip} (A_2)_{jq} (A_3)_{kr}$$

4 Missing data handling

4.1 Algorithmic problem

If I do not set threshold α when implementing the algorithm with the music survey data, The error pops up saying that the parameter for estimation goes infinity. If you consider simple case estimation with missing data, you can figure out why this error happens. Let us consider the following model.

$$P(y_{ijk} = 1) = \text{logit}(a_i b_j c_k) \text{ where } i, j, k \in \{1, 2\} \quad (3)$$

Suppose we have data such that y_{1jk} have only 2 data available and the values are 1 while y_{2jk} have all data available. Since a_1 is only component for y_{1jk} and do not affect any values in y_{2jk} , the MLE for a_1 should go to infinity. This situation happens in the Music data. Therefore, we have to set the threshold α to make the algorithm run without errors.

4.2 Method of the estimation

My method to infer the missing data us to input missing data with EM. This approach consists of mainly two steps.

1. In expectation step, we compute expectation over missing values using our algorithm with threshold α from initial music data.
2. In maximization step, we update our data set replacing missing values with the values we estimated in expectation step. Based on this new data set, we obtain optimal parameters that minimize the loss function with the new data set.
3. Repeat expectation steps and maximization steps until converge.

The algorithm for the misssing data is listed below.

```

1 source("functions2.R")
2 load("InCar_Music.RData")
3
4 fit_missing = function(tensor,ttnsr,C,A_1,A_2,A_3,omega=TRUE,alpha){
5   alphbound <- alpha+10^-4
6   result = list()
7   error<- 3
8   iter = 0
9   cost=NULL
10  omg = omega
11  while ((error > 10^-4)&(iter<50) ) {
12    iter = iter +1
13
14    #update omega
15    prevtheta <- ttl(C,list(A_1,A_2,A_3),ms=1:3)@data
16    if(sum(omg)==TRUE) {
17      omega <- polr(as.factor(c(ttnsr[ttnsr>0]))~offset(-c(prevtheta[ttnsr
18    >0])))$zeta
19    }
20
21    prev <- likelihood(ttnsr[tensor>0],prevtheta[tensor>0],omega)
22
23
24    #update A_1
25    W =kronecker(A_3,A_2)%*%t(k_unfold(C,1)@data)
26    A_1 <- comb(A_1,W,ttnsr,1,omega,alphbound)
27    #orthogonalize A_1
28    qr_res=qr(A_1)
29    A_1=qr.Q(qr_res)
30    C=ttm(C,qr.R(qr_res),1)
31    new <- likelihood(ttnsr[tensor>0],theta[tensor>0],omega)
32    if(max(abs(ttl(C,list(A_1,A_2,A_3),ms=1:3)@data))>=alpha) break
33
34    # update A_2
35    W <- kronecker(A_3,A_1)%*%t(k_unfold(C,2)@data)
36    A_2 <- comb(A_2,W,ttnsr,2,omega,alphbound)
37    #orthogonalize A_2
38    qr_res=qr(A_2)

```

```

39   A_2=qr.Q(qr_res)
40   C=ttm(C,qr.R(qr_res),2)
41   new <- likelihood(ttnsr[tensor>0],theta[tensor>0],omega)
42   if(max(abs(ttl(C,list(A_1,A_2,A_3),ms=1:3)@data))>=alph) break
43
44
45   # update A_3
46   W <- kronecker(A_2,A_1)%*%t(k_unfold(C,3)@data)
47   A_3 <- comb(A_3,W,ttnsr,3,omega,alphbound)
48   #orthogonalize A_3
49   qr_res=qr(A_3)
50   A_3=qr.Q(qr_res)
51   C=ttm(C,qr.R(qr_res),3)
52   if(max(abs(ttl(C,list(A_1,A_2,A_3),ms=1:3)@data))>=alph) break
53
54   # update C
55   W <- kronecker(kronecker(A_3,A_2),A_1)
56   C <- corecomb(C,W,c(ttnsr),omega)
57   new <- likelihood(ttnsr[tensor>0],theta[tensor>0],omega)
58
59   theta <- ttl(C,list(A_1,A_2,A_3),ms=1:3)@data
60   new <- likelihood(ttnsr[tensor>0],theta[tensor>0],omega)
61   cost = c(cost,new)
62   error <- abs((new-prev)/prev)
63   if(max(abs(ttl(C,list(A_1,A_2,A_3),ms=1:3)@data))>=alph) break
64 }
65 result$C <- C; result$A_1 <- A_1; result$A_2 <- A_2; result$A_3 <- A_3
66 result$iteration <- iter
67 result$cost = cost; result$omega=omega; result$lastcost = new
68 return(result)
69 }
70
71
72
73
74 fit_ordinal_missing = function(tensor,C,A_1,A_2,A_3,omega = TRUE,alph =
    20){
75
76   ttnsr = tensor

```



```

77 output = list()
78 error<- 3
79 iter = 0
80 cost=NULL
81 # Maximization step
82 result = fit_missing(tensor,ttnsr,C,A_1,A_2,A_3,omega=TRUE,alph =20)
83 pcost = result$lastcost
84 theta = ttl(result$C,list(result$A_1,result$A_2,result$A_3),ms=1:3)@data
85 # Expectation step
86 ttnsr[which(tensor== -1)] = realization(theta,result$omega)@data[which(
    tensor== -1)]
87 while ((error > 10^-4)&(iter<50)) {
88     # Maximization step
89     result = fit_missing(tensor,ttnsr,C,A_1,A_2,A_3,omega=TRUE,alph =20)
90     ncost = result$lastcost
91     theta = ttl(result$C,list(result$A_1,result$A_2,result$A_3),ms=1:3)
    @data
92     # Expectation step
93
94     ttnsr[which(tensor== -1)] = realization(theta,result$omega)@data[which(
    tensor== -1)]
95     error = abs((ncost-pcost)/ncost)
96     cost = c(cost,ncost)
97     pcost = ncost
98     iter = iter+1
99
100 }
101 output$C = result$C; output$A_1 = result$A_1; output$A_2 = result$A_2;
    output$A_3 = result$A_3
102 output$iteration = iter; output$error= error; output$costvar = cost
103 return(output)
104 }
105
106
107
108
109 #Initialization
110 A_1 = randortho(d[1])[,1:r[1]]
111 A_2 = randortho(d[2])[,1:r[2]]

```

```
112 A_3 = randortho(d[3])[,1:r[3]]
113 C = rand_tensor(modes = r)
114 result = fit_ordinal_missing(tensor, C, A_1, A_2, A_3, omega=TRUE, alph =20)
```