

Brain data and Music data application 3

Chanwoo Lee

Dec 25, 2019

1 Missing data handling

In the new algorithm, we define $Q(\Theta, \Theta_0) = \frac{1}{4} \|(\Theta - (\Theta_0 - 2\text{Gradient}(\Theta_0)))\|_F^2$ and updated Θ as a $\hat{\Theta}$ such that

$$\hat{\Theta} = \arg \min_{\Theta: \text{rank}(\Theta)=r} \|(\Theta - (\Theta_0 - 2\text{Gradient}(\Theta_0)))\|_F^2.$$

When Θ is complete, then we can use Tucker decomposition. In particular,

$$\hat{\Theta} = \text{Tucker}(\Theta_0 - 2\text{Gradient}(\Theta_0)).$$

However, when Θ is not complete, we have to find minimizer for

$$\sum_{(i,j,k) \in \Omega} \frac{1}{4} (\theta_{ijk} - ((\theta_0)_{ijk} - 2\text{Gradient}((\theta_0)_{ijk})))^2.$$

where Ω is a index set of available data. In this case we can rephrase the above term by

$$\text{defining } w_{ijk} = \begin{cases} 1, & \text{when } (i, j, k) \in \Omega \\ 0, & \text{when } (i, j, k) \in \Omega^c \end{cases}.$$

Then the minimizer Θ can be obtained by

$$\hat{\Theta} = \arg \min_{\Theta: \text{rank}(\Theta)=r} \frac{1}{2} \|(W * (\Theta - (\Theta_0 - 2\text{Gradient}(\Theta_0))))\|_F^2.$$

We can find the minimizer numerically using gradients. Let us denote $\Theta = \mathcal{C} \times_1 A_1 \times_2 A_2 \times_3 A_3$ and $f_W(\Theta) \stackrel{\text{def}}{=} \frac{1}{2} \|(W * (\Theta - (\Theta_0 - 2\text{Gradient}(\Theta_0))))\|_F^2 = \frac{1}{2} \|(W * (\Theta - M))\|_F^2$. Then gradients for each factor are as follows.

$$\nabla_{A_1} f_W = [W * (\mathcal{C} \times_1 A_1 \times_2 A_2 \times_3 A_3 - M)]_{(1)} [\mathcal{C} \times_2 A_2 \times_3 A_3]_{(1)}^T. \quad (1)$$

$$\nabla_{\mathcal{C}} f_W = W * (\mathcal{C} \times_1 A_1 \times_2 A_2 \times_3 A_3 - M) \times_1 A_1^T \times_2 A_2^T \times_3 A_3^T. \quad (2)$$

Based on those gradients, we can minimize f_W with block optimization, where in f_W is optimized with respect to one set of parameters with other parameters being fixed. However, in this case, the computation cost is expensive considering that updating Θ here is a just one iteration out of whole iterations. In particular, our previous algorithm uses block optimization as a main iteration. However, this sub iteration algorithm can be used as a Tucker decomposition for tensors with missing values.

2 Continuous Tucker algorithm with the missing values

I constructed continuous Tucker algorithm for the missing values. This algorithm is a gradient based optimization using the gradients (1) and (2) with M replaced by D which is the data tensor. It took about 30 mins for this algorithm to converge in the application for the brain data with the rank (23,23,8). The algorithm code is as follows.

```

1 library(MASS)
2 library(rTensor)
3 library(pracma)
4
5
6 # initial point
7 A_1 = randorth(d[1])[,1:r[1]]
8 A_2 = randorth(d[2])[,1:r[2]]
9 A_3 = randorth(d[3])[,1:r[3]]
10 C = rand_tensor(modes = r)
11 prevtheta <- ttl(C,list(A_1,A_2,A_3),ms=1:3)@data
12 theta=prevtheta-2*gradient_tensor(A_1,A_2,A_3,C,ttnsr,omega)
13
14 # constant term for the missing values
15 W = array(0, dim(tensor))
16 W[tensor>0] = 1
17
18 costf = function(C,A_1,A_2,A_3,theta,W){
19 return(sum((W*(ttl(C,list(A_1,A_2,A_3),ms = 1:3)@data-theta))^2)/2)
20 }
21
22 cgm = function(C,A_1,A_2,A_3,theta,W,i){
23 val = as.tensor(W*(ttl(C,list(A_1,A_2,A_3),ms = 1:3)@data-theta))

```

```

24 if(i == 1){
25   g = k_unfold(val,1)@data%*%t(k_unfold(ttl(C,list(A_2,A_3),ms = c(2,3)),1)
      @data)
26 }else if(i == 2){
27   g = k_unfold(val,2)@data%*%t(k_unfold(ttl(C,list(A_1,A_3),ms = c(1,3)),2)
      @data)
28 }else if(i == 3){
29   g = k_unfold(val,3)@data%*%t(k_unfold(ttl(C,list(A_1,A_2),ms = c(1,2)),3)
      @data)
30 }else{
31   g = ttl(val,list(t(A_1),t(A_2),t(A_3)),ms = 1:3)
32 }
33 return(g)
34 }
35
36 cgma = function(A,G,theta,W,i){
37   g = (k_unfold(as.tensor(W),i)@data*(A%*%G-k_unfold(as.tensor(theta),i)
      @data))%*%t(G)
38   return(g)
39 }
40
41
42 tucker_missing = function(A_1,A_2,A_3,C,theta,W,alph = TRUE){
43   alphbound <- alph+10^-4
44   result = list()
45   error<- 3
46   iter = 0
47   cost=NULL
48   if (alph == TRUE) {
49     while ((error > 10^-4)&(iter<50) ) {
50       iter = iter +1
51       (prev = costf(C,A_1,A_2,A_3,theta,W))
52
53       # update A_1
54       # tic()
55       # l <- lapply(1:nrow(A_1),
56       #             function(i){optim(A_1[i,],
57       #                               function(x) costf(C,matrix(x,ncol =
58       # length(x)),A_2,A_3,theta[i,,,drop = F],W[i,,,drop= F]),

```

```

58 #                                     function(x) cgm(C,matrix(x,ncol = length
    (x)),A_2,A_3,theta[i,,,drop = F],W[i,,,drop = F],1),
59 #                                     method = "BFGS")$par})
60 # #A_1 = matrix(unlist(l),nrow = nrow(A_1),byrow = T) #It is slower
61 # toc()
62 #
63 # tic()
64 # G = k_unfold(C,1)@data%*%t(kronecker(A_3,A_2))
65 # f = function(x) return(costf(C,matrix(x,nrow = dim(A_1)[1]),A_2,A_3,
    theta,W))
66 # g = function(x) return(cgma(matrix(x,nrow = dim(A_1)[1]),G,theta,W,1))
67 # A_1 <- matrix(optim(c(A_1),f,g,method = "BFGS" )$par,nrow = dim(A_1)
    [1])
68
69 # l <- lapply(1:nrow(A_1),
70 #             function(i){optim(A_1[i,],
71 #                                 function(x) costf(C,matrix(x,ncol =
    length(x)),A_2,A_3,theta[i,,,drop = F],W[i,,,drop= F]),
72 #                                 function(x) cgma(x,G,theta[i,,,drop = F
    ],W[i,,,drop = F],1),
73 #                                 method = "BFGS")$par}) # It is slower
74 # toc()
75 #
76 f = function(x) return(costf(C,matrix(x,nrow = dim(A_1)[1]),A_2,A_3,
    theta,W))
77 g = function(x) return(cgm(C,matrix(x,nrow = dim(A_1)[1]),A_2,A_3,theta,
    W,1))
78 A_1 <- matrix(optim(c(A_1),f,g,method = "BFGS" )$par,nrow = dim(A_1)[1])
79 #orthogonalize A_1
80 qr_res=qr(A_1)
81 A_1=qr.Q(qr_res)
82 C=ttm(C,qr.R(qr_res),1)
83
84 # update A_2
85
86 f = function(x) return(costf(C,A_1,matrix(x,nrow = dim(A_2)[1]),A_3,
    theta,W))
87 g = function(x) return(cgm(C,A_1,matrix(x,nrow = dim(A_2)[1]),A_3,theta,
    W,2))

```

```

88   A_2 <- matrix(optim(c(A_2),f,g,method = "BFGS")$par,nrow = dim(A_2)[1])
89
90   #orthogonalize A_2
91   qr_res=qr(A_2)
92   A_2=qr.Q(qr_res)
93   C=ttm(C,qr.R(qr_res),2)
94
95   # update A_3
96   f = function(x) return(costf(C,A_1,A_2,matrix(x,nrow = dim(A_3)[1]),
    theta,W))
97   g = function(x) return(cgm(C,A_1,A_2,matrix(x,nrow = dim(A_3)[1]),theta,
    W,3))
98   A_3 <- matrix(optim(c(A_3),f,g,method = "BFGS")$par,nrow = dim(A_3)[1])
99   #orthogonalize A_3
100  qr_res=qr(A_3)
101  A_3=qr.Q(qr_res)
102  C=ttm(C,qr.R(qr_res),3)
103
104  # update C
105  f = function(x) return(costf(new("Tensor",C@num_modes,C@modes,x),A_1,A_
    2,A_3,theta,W))
106  g = function(x) return(c(cgm(new("Tensor",C@num_modes,C@modes,x),A_1,A_
    2,A_3,theta,W,4)@data))
107  C <- new("Tensor",C@num_modes,C@modes,data =optim(c(C@data),f,g,method =
    "BFGS")$par)
108  (new = costf(C,A_1,A_2,A_3,theta,W))
109  cost = c(cost,new)
110  (error <- abs((new-prev)/prev))
111 }
112 }else{
113 while ((error > 10^-4)&(iter<50) ) {
114   iter = iter +1
115   (prev = costf(C,A_1,A_2,A_3,theta,W))
116   #update A_1
117   G = k_unfold(C,1)@data%*%t(kronecker(A_3,A_2))
118   l <- lapply(1:nrow(A_1),
119               function(i){constrOptim(A_1[i,],
120                                       function(x) costf(C,matrix(x,ncol =
    length(x)),A_2,A_3,theta[i,,,drop = F],W[i,,,drop= F]),

```

```

121         function(x) cgm(x,G,theta[i,,,drop
= F],W[i,,,drop = F],1),
122         ui = as.matrix(rbind(t(G),-t(G))),ci
= rep(-alph,2*nrow(t(G))),
123         method = "BFGS")$par})
124 A_1 = matrix(unlist(l),nrow = nrow(A_1),byrow = T)
125
126 #orthogonalize A_1
127 qr_res=qr(A_1)
128 A_1=qr.Q(qr_res)
129 C=ttm(C,qr.R(qr_res),1)
130 if(max(round(abs(ttl(C,list(A_1,A_2,A_3),ms=1:3)@data)),digits = 3)>=
alph) break
131
132 # update A_2
133 G = k_unfold(C,2)@data%*%t(kronecker(A_3,A_1))
134 l <- lapply(1:nrow(A_2),
135             function(i){constrOptim(A_2[i,],
136                                     function(x) costf(C,A_1,matrix(x,
ncol =length(x)),A_3,theta[,i,,drop = F],W[,i,,drop= F]),
137                                     function(x) cgm(x,G,theta[,i,,drop
= F],W[,i,,drop = F],2),
138                                     ui = as.matrix(rbind(t(G),-t(G))),ci
= rep(-alph,2*nrow(t(G))),
139                                     method = "BFGS")$par})
140 A_2 = matrix(unlist(l),nrow = nrow(A_2),byrow = T)
141 #orthogonalize A_2
142 qr_res=qr(A_2)
143 A_2=qr.Q(qr_res)
144 C=ttm(C,qr.R(qr_res),2)
145 if(max(round(abs(ttl(C,list(A_1,A_2,A_3),ms=1:3)@data)),digits = 3)>=
alph) break
146
147 # update A_3
148 G = k_unfold(C,3)@data%*%t(kronecker(A_2,A_1))
149 l <- lapply(1:nrow(A_3),
150             function(i){constrOptim(A_3[i,],
151                                     function(x) costf(C,A_1,A_2,matrix(x
,ncol =length(x)),theta[, ,i,drop = F],W[, ,i,drop= F]),

```

```

152         function(x) cgm(x,G,theta[, ,i,drop
= F],W[, ,i,drop = F],3),
153         ui = as.matrix(rbind(t(G),-t(G))),ci
= rep(-alpha,2*nrow(t(G))),
154         method = "BFGS")$par})
155 A_3 = matrix(unlist(l),nrow = nrow(A_3),byrow = T)
156 #orthogonalize A_3
157 qr_res=qr(A_3)
158 A_3=qr.Q(qr_res)
159 C=ttm(C,qr.R(qr_res),3)
160 if(max(round(abs(ttl(C,list(A_1,A_2,A_3),ms=1:3)@data)),digits = 3)>=
alpha) break
161
162 # update C
163 f = function(x) return(costf(new("Tensor",C@num_modes,C@modes,x),A_1,A_
2,A_3,theta,W))
164 g = function(x) return(c(cgm(new("Tensor",C@num_modes,C@modes,x),A_1,A_
2,A_3,theta,W,4)@data))
165 C <- new("Tensor",C@num_modes,C@modes,data =optim(c(C@data),f,g,method =
"BFGS")$par)
166 (new = costf(C,A_1,A_2,A_3,theta,W))
167 cost = c(cost,new)
168 (error <- abs((new-prev)/prev))
169 }
170 }
171 result$C <- C; result$A_1 <- A_1; result$A_2 <- A_2; result$A_3 <- A_3
172 result$iteration <- iter
173 result$cost = cost
174 return(result)
175 }

```

3 Current algorithm modification

Our main problem for the algorithm was the memory and time cost for getting gradients for the core tensor. The reason for the expensive cost is that the algorithm uses kronecker product to get the gradients. Instead, the modified algorithm uses Tensor n-mode product

to get a gradient which is a tensor structure. Previous gradient has the following formula.

$$\text{Gradient} = q(A_3 \otimes A_2 \otimes A_1) \text{ where } q \in \mathbb{R}^{d_1 d_2 d_3}.$$

This gradient can have different formula if I made q as a tensor. The modified gradient has the following formula.

$$\text{Gradient} = \text{Tensor}(q) \times_1 A_1^T \times_2 A_2^T \times_3 A_3^T \text{ where } \text{Tensor}(q) \in \mathbb{R}^{d_1 \times d_2 \times d_3}.$$

From this formula, our algorithm does not need to compute kronecker product and takes less memory and time in the end.

The following is the modified gradient algorithm.

```

1 gc = function(A_1,A_2,A_3,C,ttnsr,omega){
2 k = length(omega)
3 thet = c(ttl(C,list(A_1,A_2,A_3),ms=1:3)@data)
4 p = matrix(nrow = length(thet),ncol = k)
5 for (i in 1:k) {
6 p[,i] = as.numeric(logistic(thet + omega[i]))
7 }
8 q = matrix(nrow = length(thet),ncol = k+1)
9 q[,1] <- p[,1]-1
10 for (i in 2:k) {
11 #q[,i] <- (p[,i]*(1-p[,i])-p[,i-1]*(1-p[,i-1]))/(p[,i-1]-p[,i])
12 q[,i] <- p[,i]+p[,i-1]-1
13 }
14 q[,k+1] <- p[,k]
15 g = ttnsr
16 for(i in 1:(k+1)){
17 g[which(ttnsr==i)] = q[which(ttnsr==i),i]
18 }
19
20 g = ttl(as.tensor(g),list(t(A_1),t(A_2),t(A_3)),ms = 1:3)
21
22 # d = c(nrow(A_1),nrow(A_2),nrow(A_3))
23 # cl <- makeCluster(20)
24 # registerDoParallel(cl)
25 # l <- foreach(j = 1:d[3],.combine = "+") %dopar% {

```



```

26 #   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),])*)
27 #
      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)
      ))
28 # }
29 # stopCluster(cl)
30 return(g)
31 }

```

4 BIC result from the new algorithm and modified algorithm

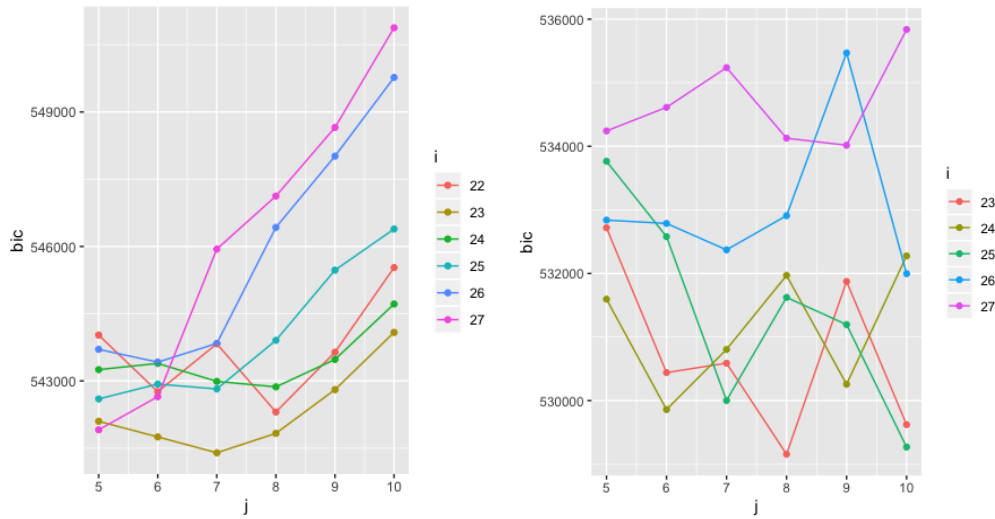


Figure 1: the left figure is the BIC from the new algorithm and the right figure is the BIC from the modified figure. Both figure show the rank around (23,23,8) is good.

5 Cross Validation result

I did 5-fold cross validation to compare performances of the two methods: Continuous Tucker method and Ordinal GLM Tucker Method. The following is the result for Continuous Tucker model and ordinal GLM Tucker model with the rank (25,25,7).

MSE	MAE	Error_rate
0.1593824	0.1591757	0.1590723
0.1575139	0.1574344	0.1573947
0.1587304	0.1586668	0.1586350
0.1576093	0.1575298	0.1574901
0.1592472	0.1591041	0.1590325

MSE	MAE	Error_rate
0.1504456	0.1503343	0.1502787
0.1495392	0.1493961	0.1493246
0.1491338	0.1490383	0.1489906
0.1507557	0.1506444	0.1505888
0.1510102	0.1509306	0.1508909

Table 1: Continuous Tucker decomposition Table 2: Ordinal GLM Tucker decomposition

6 Clustering

I used loading based clustering for the output with the rank (24,24,8). Also, I used ‘tucker’ function to make the output normalized before doing clustering. The largest difference between entries of before and after ‘tucker’ decomposition is 1.477929e-12, which means θ values are almost the same. I choose the number of cluster groups 8 because it does not have the tendency that most of points are included in the same one group as you check in the following figure. I used ‘kmeans’ function with **nstart** = 20 option to reproduce the same result.

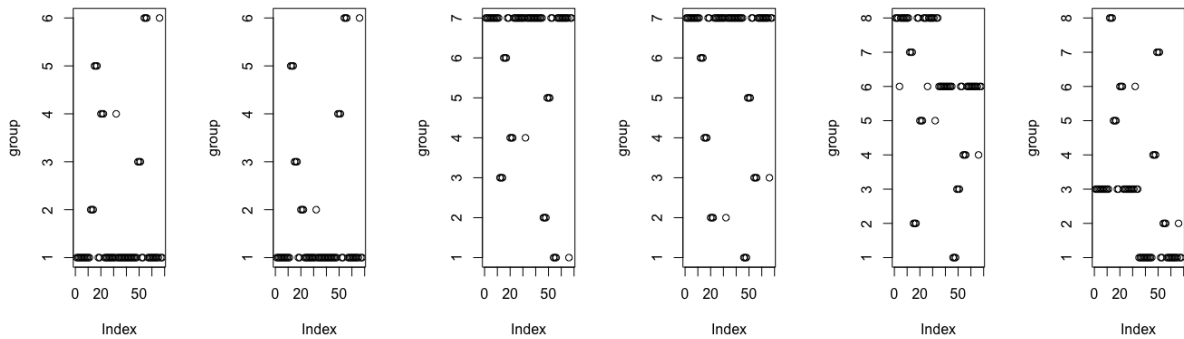


Figure 2: The number of the groups is from 6 to 8 in order. You can see that cluster with mode 1 and with mode 2 have exactly the same clusters.

7 Music data application

There are some changes I made in the algorithm for the music data.

1. I made algorithm that change the order of updating factor matrices at every iteration. The main reason for this is that it usually hits the boundary α in one factor update.
2. The algorithm estimates the missing value as a mean value of realization in expectation step.

However, it sometimes gives me error saying that initial value hits the boundary. This is because after each iteration, output value becomes really close to the boundary and ‘constrOptim’ function is quite sensitive to this initial value. So I am working on fixing the bug now. One thing I want to discuss is that continuous Tucker method works algorithmically without any modification. One reason for this is that the parameters of continuous Tucker method cannot be infinity while the parameters of ordinal GLM method can be infinity because it only becomes probability 1 in realization. Therefore, if we use EM algorithm for the missing data in ordinal GLM method, it is not fair to compare with continuous Tucker method because we do one extra step more. My suggestion for this problem is that we set α , make all factor matrices and core tensor updated at least one time and estimate the missing value as the mean value from the given parameters.