

# New version of the theorem for the convergence and The new algorithm for ordinal tensors

Chanwoo Lee

Nov 19, 2019

## 1 Estimation accuracy for more general randomized SVD

We generalize from random normal test matrices to arbitrary test matrices. We can guarantee the convergence of estimators under the certain conditions in the next theorem.

**Theorem 1.** *Let  $\mathcal{A} = \mathcal{C} \times_1 M^{(1)} \times_2 \cdots \times_N M^{(N)}$  be a signal tensor, where  $\mathcal{C} \in \mathbb{R}^{r_1 \times \cdots \times r_N}$  is a core tensor and  $M^{(i)}$  is an orthonormal matrix in  $\mathbb{R}^{d_i \times r_i}$  for all  $i \in [N]$ . Let  $\mathcal{D} = \mathcal{A} + \mathcal{E}$  be a noisy tensor with a noise tensor  $\mathcal{E}$  with i.i.d. entries from  $N(0, \sigma^2)$ . Suppose,  $(\hat{\mathcal{C}}, \hat{M}^{(1)}, \dots, \hat{M}^{(N)})$  is obtained from the randomized algorithms with test matrices  $\Omega^{(i)}$  for  $i \in [N]$ . If  $s_{\min}(C_{(i)}) \sqrt{\mathbb{1}_{r_i}^T (\Omega^{(i)})^T P^{(i)} (P^{(i)})^T \Omega^{(i)} \mathbb{1}_{r_i}} \gg \sigma \sqrt{r_i \max(d_i, \frac{\prod_{j \neq i} d_j}{d_i})} \|\Omega^{(i)}\|_\sigma$  as  $d_1, \dots, d_N \rightarrow \infty$ , where  $s_{\min}(C_{(i)})$  is the smallest singular value of  $C_{(i)}$ ,  $P^{(i)} = [(M^{(N)} \otimes \cdots \otimes M^{(i+1)} \otimes M^{(i-1)} \otimes \cdots \otimes M^{(1)}) V_{r_i}^{(i)}]$  and  $V_{r_i}^{(i)}$  is the matrix of the  $r$  largest left singular vectors of  $C_{(i)}$ . Then, the following holds true.*

$$\cos \theta(M^{(i)}, \hat{M}^{(i)}) \rightarrow 1 \text{ in probability for } i \in [N].$$

$$\|\mathcal{A} - \hat{\mathcal{A}}\|_F \rightarrow 0 \text{ in probability.}$$

where  $\hat{\mathcal{A}} = \hat{\mathcal{C}} \times_1 \hat{M}^{(1)} \times_2 \cdots \times_N \hat{M}^{(N)}$ .

*Proof.* It suffices to show when  $i = 1$ . Notice,

$$\begin{aligned} A_{(1)} &= M^{(1)} (\mathcal{C} \times_2 M^{(2)} \times_3 \cdots \times_N M^{(N)})_{(1)} \\ &= M^{(1)} C_{(1)} (M^{(N)} \otimes \cdots \otimes M^{(2)})^T. \end{aligned}$$

Define  $B = (M^{(N)} \otimes \dots \otimes M^{(2)})$ . the randomized algorithms generates a test matrix  $\Omega^{(1)}$  and captures the image space of unfolded matrix  $A_{(1)}$ . Having this procedure in mind, we obtain,

$$\begin{aligned} A_{(1)}\Omega^{(1)} &= M^{(1)}C_{(1)}(M^{(N)} \otimes \dots \otimes M^{(2)})^T\Omega^{(1)} \\ &= M^{(1)}C_{(1)}B^T\Omega^{(1)}. \end{aligned}$$

However, since the input is  $\mathcal{D} = \mathcal{A} + \mathcal{E}$ , we have the image space of  $A_{(1)} + E_{(1)}$  instead of  $A_{(1)}$ . Therefore, the estimator  $\hat{M}^{(1)}$  is obtained from the following equality.

$$\begin{aligned} (A_{(1)} + E_{(1)})\Omega^{(1)} &= M^{(1)}C_{(1)}B^T\Omega^{(1)} + E_{(1)}\Omega^{(1)} \\ &= \hat{M}^{(1)}R \quad (\text{QR decomposition}). \end{aligned}$$

From the relationship that  $\text{span}(A_{(1)}\Omega^{(1)}) \subset \text{span}(M^{(1)})$  and  $\text{span}(A_{(1)}\Omega^{(1)} + E_{(1)}\Omega^{(1)}) = \text{span}(\hat{M}^{(1)})$ , we have the following.

$$\begin{aligned} \cos \theta(M^{(1)}, \hat{M}^{(1)}) &= \max_{u \in \text{span}(M^{(1)}), v \in \text{span}(\hat{M}^{(1)})} \cos(u, v) \\ &\geq \max_{u \in \text{span}(A_{(1)}\Omega^{(1)}), v \in \text{span}((A_{(1)} + E_{(1)})\Omega^{(1)})} \cos(u, v) \\ &= \max_{x \in R^{r_1}, y \in R^{r_1}, \|x\|_2 = \|y\|_2 = 1} \cos(A_{(1)}\Omega^{(1)}x, (A_{(1)} + E_{(1)})\Omega^{(1)}y). \end{aligned} \tag{1}$$

The first argument in the theorem holds true by (1) if

$$\max_{x \in R^{r_1}, y \in R^{r_1}, \|x\|_2 = \|y\|_2 = 1} \cos(A_{(1)}\Omega^{(1)}x, (A_{(1)} + E_{(1)})\Omega^{(1)}y) \rightarrow 1. \tag{2}$$

Also (2) holds true, if

$$\cos(A_{(1)}\Omega^{(1)}x, (A_{(1)} + E_{(1)})\Omega^{(1)}y) \rightarrow \infty \text{ for some } x, y \text{ such that } \|x\| = \|y\| = 1. \tag{3}$$

So the main proof of this theorem is to show (3). We prove (3) by the following inequality.

$$\cos(A_{(1)}\Omega^{(1)}x, (A_{(1)} + E_{(1)})\Omega^{(1)}y) \geq \frac{\|A_{(1)}\Omega^{(1)}x\|_2}{\|E_{(1)}\Omega^{(1)}y\|_2} \geq \frac{s_{\min}(C_{(1)})\sqrt{\mathbb{1}_{r_1}^T(\Omega^{(1)})^T P^{(1)}(P^{(1)})^T \Omega^{(1)} \mathbb{1}_{r_1}}}{\sqrt{r_1}\|E\|_F\|\Omega^{(1)}\|_\sigma}. \tag{4}$$

for some  $x$  and  $y$ .

To get the numerator part in (4),

$$\begin{aligned}
\|A_{(1)}\Omega^{(1)}x\|_2 &= \|M^{(1)}C_{(1)}B^T\Omega^{(1)}x\|_2 \\
&\stackrel{(i)}{=} \|C_{(1)}B^T\Omega^{(1)}x\|_2 \\
&\stackrel{(ii)}{=} \|U^{(1)}\Sigma^{(1)}(V^{(1)})^TB^T\Omega^{(1)}x\|_2 \\
&= \|\Sigma^{(1)}(V_{r_1}^{(1)})^TB^T\Omega^{(1)}x\|_2 \\
&\geq s_{\min}(C_{(1)})\|(V_{r_1}^{(1)})^TB^T\Omega^{(1)}x\|_2 \\
&= s_{\min}(C_{(1)})\|(P^{(1)})^T\Omega^{(1)}x\|_2 \\
&\stackrel{(iii)}{=} s_{\min}(C_{(1)})\frac{1}{\sqrt{r_1}}\|(P^{(1)})^T\Omega^{(1)}\mathbb{1}_{r_1}\|_2 \\
&= \frac{s_{\min}(C_{(1)})\sqrt{\mathbb{1}_{r_1}^T(\Omega^{(1)})^TP^{(1)}(P^{(1)})^T\Omega^{(1)}\mathbb{1}_{r_1}}}{\sqrt{r_1}}.
\end{aligned}$$

(i) is from the orthonormality of  $M^{(1)}$ . Singular value decomposition of  $C_{(1)}$  is used in (ii). Notice, In (iii), we put  $x = \mathbb{1}_{r_1}/\sqrt{r_1}$ . Therefore, we get the numerator part. For the denominator of (4), we know

$$\|E\|_F \asymp (2 + o(1))\sigma\sqrt{\max(d_1, d_2 \cdots d_N)}.$$

Also, notice that

$$\|\Omega y\|_2 \leq \|\Omega\|_\sigma.$$

Therefore, we get (4). The last argument that  $\|\mathcal{A} - \hat{\mathcal{A}}\| \rightarrow 0$  is derived directly from Theorem 2 and Theorem 3 in the 7th meeting note.  $\square$

When all entries of the test matrices are i.i.d. from  $N(0, 1)$ , we have the following corollary.

**Corollary 1.** *Let  $\mathcal{A} = \mathcal{C} \times_1 M^{(1)} \times_2 M^{(2)} \times_3 M^{(3)}$  be a signal tensor, where  $\mathcal{C} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$  is a core tensor and  $M^{(1)}, M^{(2)}, M^{(3)}$  are orthonormal matrices in  $\mathbb{R}^{d_1 \times r_1}, \mathbb{R}^{d_2 \times r_2}, \mathbb{R}^{d_3 \times r_3}$  respectively. Suppose we use standard normal random matrices as test matrices in Theorem 1. If  $s_{\min}(C_{(i)}) \gg \sigma\sqrt{\max(d_i, \frac{d_1 d_2 d_3}{d_i}) \frac{d_1 d_2 d_3}{d_i r_i}}$  as  $d_1, d_2, d_3 \rightarrow \infty$ , where  $s_{\min}(C_{(i)})$  is the smallest singular value of  $C_{(i)}$ . Then, the following holds true.*

$$\cos \theta(M^{(i)}, \hat{M}^{(i)}) \rightarrow 1 \text{ in probability for } i \in [3].$$

$$\|\mathcal{A} - \hat{\mathcal{A}}\|_F \rightarrow 0 \text{ in probability.}$$

*Proof.* It is enough to show that the condition in Theorem 1 implies the condition in the Corollary 1. We will show this argument when  $i = 1$  with fixed  $\Omega^{(1)}$  replaced by a standard normal random matrix. First, let us check

$$\sqrt{\mathbb{1}_{r_1}^T (\Omega^{(1)})^T P^{(1)} (P^{(1)})^T \Omega^{(1)} \mathbb{1}_{r_1}}. \quad (5)$$

where  $\Omega^{(1)}$  is a standard normal random matrix. Notice,

$$(P^{(1)})^T \Omega^{(1)} \mathbb{1}_{r_1} \sim (P^{(1)})^T N_{r_1}(0, r_1 I_{r_1}) \stackrel{d}{=} N_{r_1}(0, r_1 (P^{(1)})^T P^{(1)}) = N_{r_1}(0, r_1 I_{r_1}).$$

Because  $(P^{(1)})^T P^{(1)} = I_{r_1}$  by the definition of  $P^{(1)}$ . We get the following distribution for (5).

$$\sqrt{\mathbb{1}_{r_1}^T (\Omega^{(1)})^T P^{(1)} (P^{(1)})^T \Omega^{(1)} \mathbb{1}_{r_1}} \stackrel{d}{=} \sqrt{N_{r_1}(0, r_1 I_{r_1})^T N_{r_1}(0, r_1 I_{r_1})} \stackrel{d}{=} \sqrt{r_1 \chi_{r_1}^2}. \quad (6)$$

Secondly, we have

$$\begin{aligned} \|\Omega^{(1)}\|_\sigma &\geq \|\Omega^{(1)} y\|_2 \quad \text{where } y \in \mathbb{R}^{r_1} \text{ such that } \|y\|_2 = 1 \\ &\stackrel{d}{=} \sqrt{N_{d_2 d_3}(0, I_{d_2 d_3})^T N_{d_2 d_3}(0, I_{d_2 d_3})} \\ &\stackrel{d}{=} \sqrt{\chi_{d_2 d_3}^2} \\ &\asymp (1 + o(1)) \sqrt{d_2 d_3}. \end{aligned}$$

Therefore, the condition in Theorem 1 can be rewritten as,

$$s_{\min}(C_{(1)}) \sqrt{\chi_{r_1}^2} \gg \sigma \sqrt{\max(d_1, d_2 d_3) d_2 d_3}. \quad (7)$$

By the following inequality, we have the condition of this corollary from (7). (7) implies the left side of the following equation converges to 1.

$$\begin{aligned} P(s_{\min}(C_{(1)}) \sqrt{\chi_{r_1}^2} > \sigma \sqrt{\max(d_1, d_2 d_3) d_2 d_3}) &= P(\chi_{r_1}^2 \geq \frac{\sigma^2 d_2 d_3 \max(d_1, d_2 d_3)}{s_{\min}(C_{(1)})^2}) \\ &\stackrel{(i)}{\leq} 1 - \left( \lambda e^{1-\lambda} \right)^{\frac{r_1}{2}}. \end{aligned}$$

□

In (i), we defined  $\lambda \stackrel{def}{=} \frac{\sigma^2 d_2 d_3 \max(d_1, d_2 d_3)}{r_1 s_{min}(C_{(1)})^2}$  and used Chernoff bounds,

$$P(\chi_r^2 \geq t) \leq 1 - \left(\frac{t}{r} e^{1-\frac{t}{r}}\right)^{\frac{r}{2}} \text{ for any } t \geq 0.$$

Therefore,  $\lambda$  should converge to 0 when (7) holds. Now we have the condition of the corollary true.

## 2 Extended angle simulation for an arbitrary rank

This simulation investigates the accuracy of estimators in terms of angles and MSE for an arbitrary rank. We consider an order-3 dimension  $(20, 20, 20)$  signal tensor  $X$ . We assume  $X$  has Tucker decomposition as  $X = \mathcal{C} \otimes_1 B_1 \otimes_2 B_2 \otimes_3 B_3$ , where  $B_i \in \mathbb{R}^{20 \times 3}$  for all  $i$ . and  $\mathcal{C} \in \mathbb{R}^{3 \times 3 \times 3}$  a core tensor. All entries of  $\mathcal{C}$  are i.i.d. drawn from  $N(0, 1)$ .  $B_1, B_2, B_3$  are randomly drawn from orthonormal matrices. We vary the noise level  $\sigma \in \{0.01, 0.02, \dots, 0.49, 0.5\}$ . We use target rank 3 and estimate the signal matrices according to each algorithms. We compare the principal angles between the true signal matrices and estimators. Figure 1 shows that Method 3 outperforms the other methods in all respects.

## 3 Improved ordinal tensors algorithm

First, I constructed stochastic gradient descent(SGD) algorithm for updating the core tensor. However, there were some problems to implement SGD method. First, I used various batch sizes from  $B = 100$  to  $B = 1000$ . But this algorithm had more than 100 the number of iterations in all batch sizes. Secondly, I picked the tolerance size as  $10^{-4}$ . this algorithm reached to this tolerance size fast but had the smaller likelihood value that the value with the true parameter. This means it did not get to the optimal point but had small improvement on each update. In addition, I found that it also has variation issues, which is inevitable. There were some cases that this algorithm worked quite well with moderate iteration numbers and time. However, it sometimes performed poorly with large iteration numbers and bad outputs. Instead of using stochastic gradient method, I constructed the algorithm which calculate hessian in each update for the core tensor. This hessian function reduced iteration time dramatically for updating the core tensor. Also, it converged with the less the number

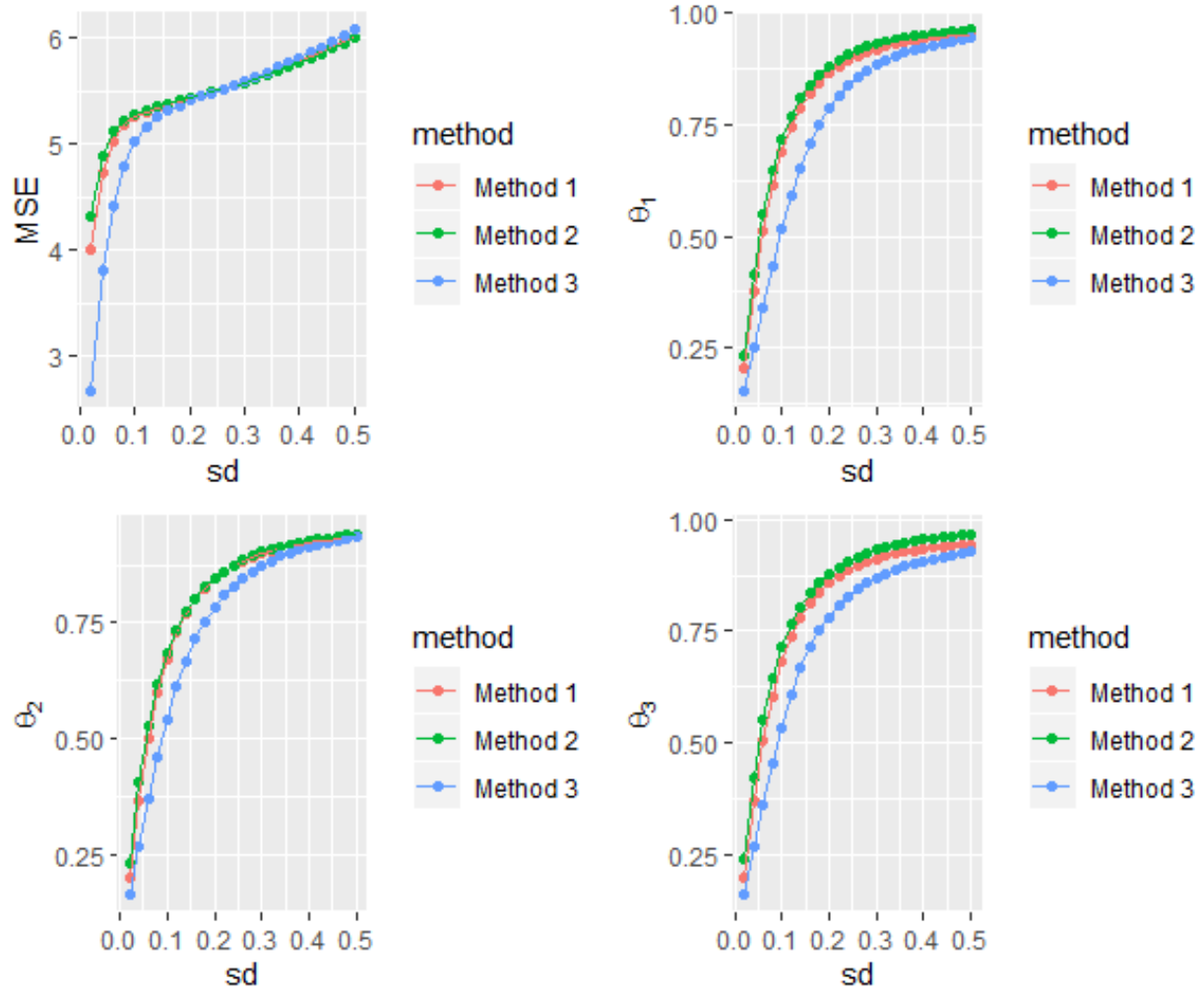


Figure 1: Each  $y$  axis means the principal angles for each mode. We do not have consistent simulation results: Method 3 has the best performance in MSE and the principal angles.

of iterations than the previous algorithm. All output had the greater likelihood value than the value with the true parameters.

### 3.1 Simulation for ordinal tensors

I performed the simulations having the same setting as in section 3.2 in the 7-th note. I summarized the output as follows.

1. When  $d = 20$  and  $r = 3$  with  $\max(\Theta_{\text{True}}) = 5.78633$  and  $\omega = (-0.2, 0.2)$ .

We have  $L(\Theta_{\text{True}}) = 6459.568$  and  $L(\Theta_0) = 7414.672$ .

	(With $\omega$ information)	(Without $\omega$ information)
$L(\hat{\Theta})$	6373.461	6373.191
Computation time	51 sec	52 sec

When we implement algorithm without  $\omega$  information, we have an estimate  $\hat{\omega} = (-1.8011, 0.2513)$

2. When  $d = 30$  and  $r = 3$  with  $\max(\Theta_{\text{True}}) = 6.8348$  and  $\omega = (-0.2, 0.2)$ .

We have  $L(\Theta_{\text{True}}) = 21895.92$  and  $L(\Theta_0) = 24917.37$ .

	(With $\omega$ information)	(Without $\omega$ information)
$L(\hat{\Theta})$	21761.86	21761.92
Computation time	432 sec	385.56 sec

When we implement algorithm without  $\omega$  information, we have an estimate  $\hat{\omega} = (-0.221403, 0.1868)$

The following is the scatter plot between true parameters and estimators.

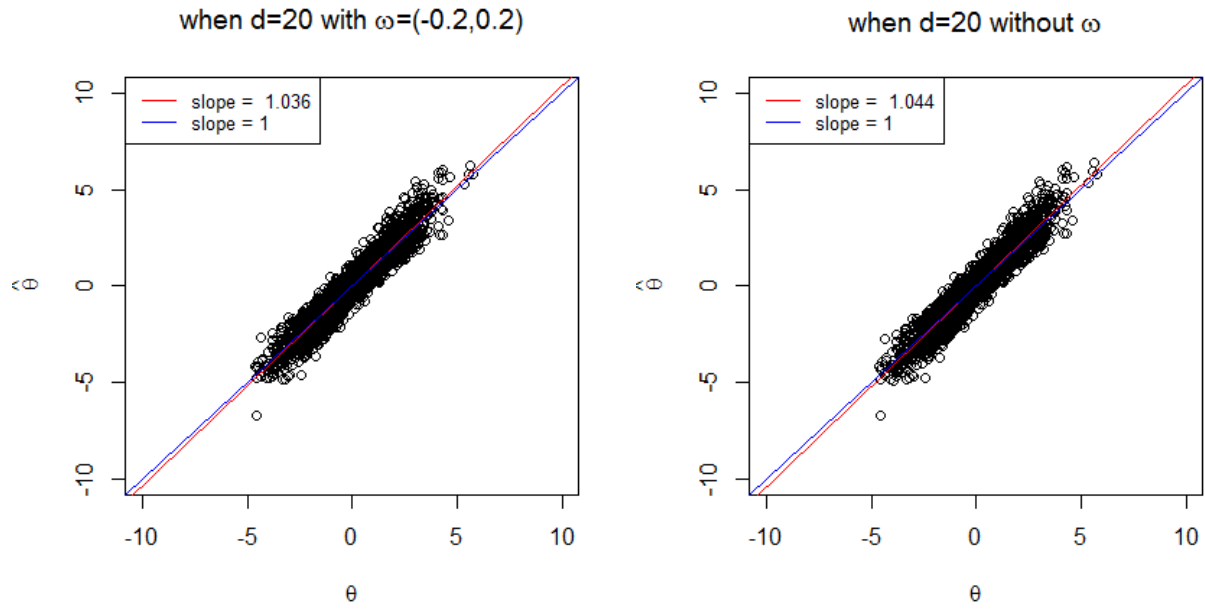


Figure 2: When  $d = 20$ . Red lines are slopes of ordinary least square estimators. Blue lines are line of  $y = x$ .

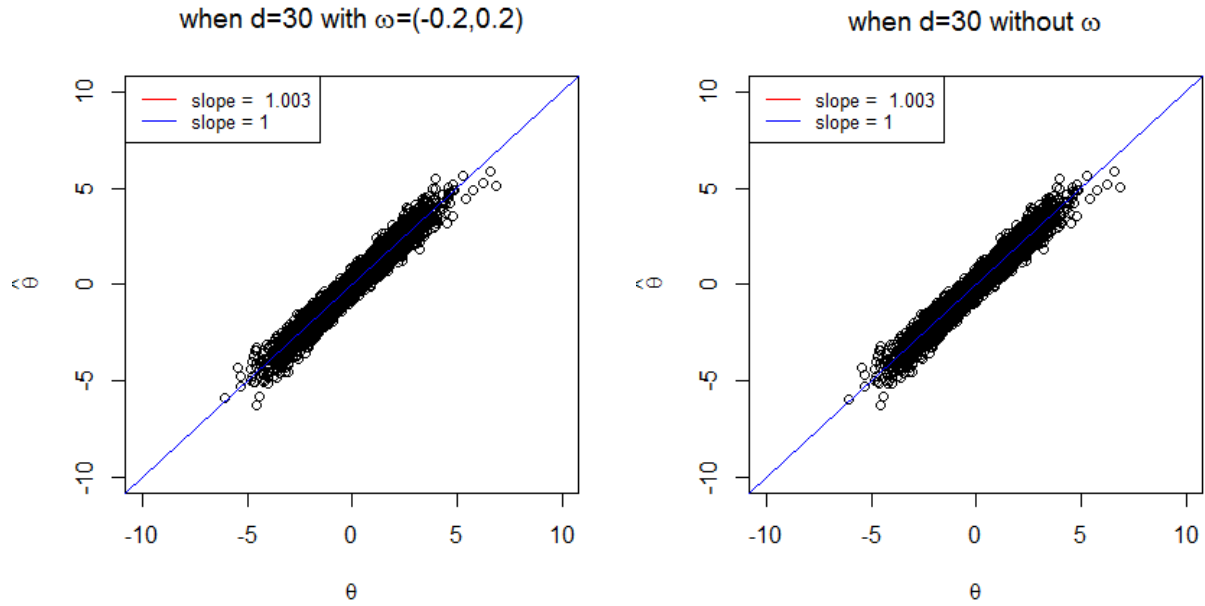


Figure 3: When  $d = 30$ . Red lines are slopes of ordinary least square estimators. Blue lines are line of  $y = x$ .

We can check that our estimators has a tendency to overestimate the true parameters.



The reason we had the opposite tendency in the 7th note is that I made mistake to calculate the slope of each scatter plot. I put  $\Theta_{\text{True}}$  into response variables and  $\hat{\Theta}$  into explanatory variables, which is the opposite direction from what I wanted to calculate. I checked we have the same overestimation tendency from the previous algorithm too if I calculate the slope in the right way. The overall comparisons between the previous algorithm and the current algorithm are shown in the following table. The table shows that the number of the iterations and the slope between the true parameters and the estimators are almost the same. However they have huge difference in the computation time. To find out the reason, I checked the core tensor updating inner function, which is the only different part between the two. For the previous algorithm, we use “optim” function to update the core tensors while we use “nlminb” function in the current algorithm. I checked that “optim” function needs 16.39 sec and 21 iteration numbers to update the core tensor. On the other hand, “nlminb” function takes 7.34 sec and 4 iteration numbers under the same condition. We can conclude that this computation time difference causes difference between the two main algorithms. However, I checked that “optim” and “nlminb” have the really close outputs. This means that the two main algorithm have the similar output for each update for the core tensor. Therefore, we can explain why the two main algorithms do not have big distinction for the output values and the iteration number in the main loop.

	$d = 20$			$d = 30$		
	Iteration #	Computation time	Slope	Iteration #	Computation time	Slope
Algorithm 1	7	121.79 sec	1.053	6	1470.78 sec	0.99698
Algorithm 2	7	45.5 sec	1.053	6	431 sec	0.99712

Table 1: Algorithm 1 is the previous algorithm and Algorithm 2 is the current algorithm using Hessian. We set the tolerance as  $10^{-4}$ .

## 4 Rank selection with BIC

In practice, we have no information about true rank of the model. Estimating an appropriate rank is an important part to deal with. We use the Bayesian information criterion(BIC),

and choose the rank that has the minimum BIC. Our estimated rank can be written as

$$\hat{R} = \arg \min_{R \in \mathbb{R}_+} BIC(R) = \arg \min_{R \in \mathbb{R}_+} (-2\mathcal{L}_{\mathcal{Y}}(\hat{\Theta}(R)) + (\prod_k r_k + \sum_k (d_k r_k - r_k^2)) \log(\prod_k d_k)).$$

where  $\hat{\Theta}(R)$  is the estimated tensor  $\Theta$  under the rank size  $R$ . The front term of log is the effective number of parameters and we proved this formula in the 7th note. We simulated the tensor rank selection by BIC. We considered the tensor dimension  $d = 20$  case and varied the rank  $R \in \{3, 4\}$ . The following table shows the BIC values for each rank according to true rank  $R$ . BICs with the close ranks to true rank have the small values as we expected.

	$R = 2$	$R = 3$	$R = 4$	$R = 5$	$R = 6$
$R_{\text{True}} = 3$	14483.72	14559.58	15067.75	15688.72	16749.21
$R_{\text{True}} = 4$	14267.40	13985.33	13818.74	14448.10	15219.82

Table 2: BIC values in ordinal tensor decomposition according to true rank and various ranks.

## 5 Algorithms

### 5.1 Extended angle simulation

```

1 library(rTensor)
2 library(pracma)
3 B_1 = matrix(rnorm(20*3), nrow = 20)
4 B_2 = matrix(rnorm(20*3), nrow = 20)
5 B_3 = matrix(rnorm(20*3), nrow = 20)
6 C = as.tensor(array(rnorm(3^3), dim = c(3,3,3)))
7 X = ttm(ttm(ttm(C, B_1, 1), B_2, 2), B_3, 3)
8 sd = 0.02*1:25
9 result = data.frame(matrix(0, nrow = 75, ncol = 6))
10 names(result) <- c("sd", "angle1", "angle2", "angle3", "method", "MSE")
11
12
13 for (i in 1:25) {
14   s=sd[i]
```

```

15 result[i,1] = s
16 result[i+25,1] = s
17 result[i+50,1] = s
18 for (j in 1:200) {
19   set.seed(j)
20   e = as.tensor(array(rnorm(8000,mean = 0,sd = s),dim = c(20,20,20)))
21   D = X+e
22   est1 = tensor_svd(D,3,3,3,0)
23   est2 = tensor_svd3(D,3,3,3,0)
24   est3 = tensor_svd4(D,3,3,3,0)
25   result[i,2] <- result[i,2]+subspace(est1$U[[1]],B_1)
26   result[i,3] <- result[i,3]+subspace(est1$U[[2]],B_2)
27   result[i,4] <- result[i,4]+subspace(est1$U[[3]],B_3)
28   result[i,6] <- result[i,6]+tensor_resid(X,est1)
29   result[i+25,2] <- result[i+25,2]+subspace(est2$U[[1]],B_1)
30   result[i+25,3] <- result[i+25,3]+subspace(est2$U[[2]],B_2)
31   result[i+25,4] <- result[i+25,4]+subspace(est2$U[[3]],B_3)
32   result[i+25,6] <- result[i+25,6]+tensor_resid(X,est2)
33   result[i+50,2] <- result[i+50,2]+subspace(est3$U[[1]],B_1)
34   result[i+50,3] <- result[i+50,3]+subspace(est3$U[[2]],B_2)
35   result[i+50,4] <- result[i+50,4]+subspace(est3$U[[3]],B_3)
36   result[i+50,6] <- result[i+50,6]+tensor_resid(X,est3)
37 }
38 result[i,5] = "Method 1"
39 result[i+25,5] = "Method 2"
40 result[i+50,5] = 'Method 3'
41 }
42 result[,2:4] <- result[,2:4]/200
43 result[,6] <- result[,6]/200
44
45 library(gridExtra)
46 library(ggplot2)
47 g1 <- ggplot(data = result, aes(x=sd,y = MSE ,color = method))+
48   geom_point(aes(x=sd, y = MSE))+geom_line(aes(x=sd, y = MSE))
49 g2 <- ggplot(data = result, aes(x=sd,y = abs(angle1) ,color = method))+
50   geom_point(aes(x=sd, y = abs(angle1)))+geom_line(aes(x=sd, y = abs(
    angle1)))+ylab(expression(theta[1]))
51 g3 <- ggplot(data = result, aes(x=sd,y = abs(angle2) ,color = method))+

```

```

52 geom_point(aes(x=sd, y = abs(angle2)))+geom_line(aes(x=sd, y = abs(
    angle2)))+ylab(expression(theta[2]))
53 g4 <- ggplot(data = result, aes(x=sd,y = abs(angle3) ,color = method))+
54 geom_point(aes(x=sd, y = abs(angle3)))+geom_line(aes(x=sd, y = abs(
    angle3)))+ylab(expression(theta[3]))
55 grid.arrange(g1,g2,g3,g4)

```

## 5.2 New ordinal tensor algorithms

```

1 library(MASS)
2 library(rTensor)
3 library(pracma)
4 library(ggplot2)
5 library(ggthemes)
6 library(gridExtra)
7
8 realization = function(tnsr,alpha){
9   thet <- k_unfold(tnsr,1)@data
10  theta1 <- thet + alpha[1]
11  theta2 <- thet + alpha[2]
12  result <- k_unfold(tnsr,1)@data
13  p1 <- logistic(theta1)
14  p2 <- logistic(theta2)-logistic(theta1)
15  p3 <- matrix(1,nrow = nrow(thet),ncol = ncol(thet))-logistic(theta2)
16  for (i in 1:nrow(thet)) {
17    for(j in 1:ncol(thet)){
18      result[i,j] <- sample(c(1,2,3),1,prob= c(p1[i,j],p2[i,j],p3[i,j]))
19    }
20  }
21  return(k_fold(result,1,modes = tnsr@modes))
22 }
23
24 #Hessian function
25 Hessi = function(A_1,W4,ttnsr,omega){
26   thet =W4*%c(A_1)
27   p1 = logistic(thet + omega[1])
28   p2 = logistic(thet + omega[2])
29   q1 = p1*(1-p1)

```

```

30 q2 = p2*(1-p2)+p1*(1-p1)
31 q3 = p2*(1-p2)
32 H = t(W4[which(c(ttnsr)==1),])%%diag(q1[which(c(ttnsr)==1)])%%W4[which(
    (c(ttnsr)==1),]+
33 t(W4[which(c(ttnsr)==2),])%%diag(q2[which(c(ttnsr)==2)])%%W4[which(c
    (ttnsr)==2),]+
34 t(W4[which(c(ttnsr)==3),])%%diag(q3[which(c(ttnsr)==3)])%%W4[which(c
    (ttnsr)==3),]
35 return(H)
36 }
37
38 #Function
39 h1 = function(A_1,W1,ttnsr,omega){
40   thet =W1%%c(A_1)
41   p1 = logistic(thet + omega[1])
42   p2 = logistic(thet + omega[2])
43   p = cbind(p1,p2-p1,1-p2)
44   return(-sum(log(c(p[which(c(ttnsr)==1),1],p[which(c(ttnsr)==2),2],p[
       which(c(ttnsr)==3),3]))))
45 }
46
47 #Gradient
48 g1 = function(A_1,W1,ttnsr,omega){
49   thet =W1%%c(A_1)
50   p1 = logistic(thet + omega[1])
51   p2 = logistic(thet + omega[2])
52   q1 <- p1-1
53   q2 <- (p2*(1-p2)-p1*(1-p1))/(p1-p2)
54   q3 <- p2
55   gd = apply(diag(q1[which(c(ttnsr)==1)])%%W1[which(c(ttnsr)==1),],2,sum)
      +
56   apply(diag(q2[which(c(ttnsr)==2)])%%W1[which(c(ttnsr)==2),],2,sum)+
57   apply(diag(q3[which(c(ttnsr)==3)])%%W1[which(c(ttnsr)==3),],2,sum)
58   return(gd)
59 }
60
61
62 comb = function(A,W,ttnsr,k,omega,alph=TRUE){
63   nA = A

```

```

64  tnsr1 <- k_unfold(as.tensor(ttnsr),k)@data
65  if (alph==TRUE) {
66    l <- lapply(1:nrow(A),function(i){optim(A[i,],
67                                              function(x) h1(x,W,tnsr1[i,],
68                                              omega),
69                                              function(x) g1(x,W,tnsr1[i,],
70                                              omega),
71                                              method = "BFGS")$par}})
72    nA <- matrix(unlist(l),nrow = nrow(A),byrow = T)
73  }else{
74    l <- lapply(1:nrow(A),function(i){constrOptim(A[i,],
75                                                    function(x) h1(x,W,tnsr1
76                                                    [i,],omega),function(x) g1(x,W,tnsr1[i,],omega),
77                                                    ui = rbind(W,-W),ci =
78                                                    rep(-alph,2*nrow(W)),method = "BFGS")$par}})
79    nA <- matrix(unlist(l),nrow = nrow(A),byrow = T)
80  }
81  return(nA)
82
83 optim(Cvec,h,g,method = "BFGS")
84 nlminb(Cvec,h,g,H)
85
86 corecomb = function(C,W,ttnsr,omega,alph=TRUE){
87   Cvec <- c(C@data)
88   h <- function(x) h1(x,W,ttnsr,omega)
89   g <- function(x) g1(x,W,ttnsr,omega)
90   H <- function(x) Hessi(x,W,ttnsr,omega)
91   d <- nlminb(Cvec,h,g,H)
92   C <- new("Tensor",C@num_modes,C@modes,data =d$par)
93
94   return(C)
95 }
96
97 #previous core tensor updating algorithm.
98 prevcorecomb = function(C,W,ttnsr,omega,alph=TRUE){
99   Cvec <- c(C@data)
100   h <- function(x) h1(x,W,ttnsr,omega)
101   g <- function(x) g1(x,W,ttnsr,omega)

```

```

99   H <- function(x) Hessi(x,W,ttnsr,omega)
100
101
102   if (alph==TRUE) {
103     d <- nlminb(Cvec,h,g,H)
104     C <- new("Tensor",C@num_modes,C@modes,data =d$par)
105   }else{
106     d <- constrOptim(Cvec,h,g,ui = rbind(W,-W),ci = rep(-alph,2*nrow(W)),
107       method = "BFGS")
108     C <- new("Tensor",C@num_modes,C@modes,data =d$par)
109   }
110   return(C)
111 }
112
113
114
115 #####
116 fit_ordinal = function(ttnsr,C,A_1,A_2,A_3,omega,alph = TRUE){
117   alphbound <- alph+10^-4
118   result = list()
119   error<- 3
120   iter = 0
121   d1 <- nrow(A_1); d2 <- nrow(A_2); d3 <- nrow(A_3)
122   r1 <- ncol(A_1); r2 <- ncol(A_2); r3 <- ncol(A_3)
123   if (alph == TRUE) {
124     while ((error > 10^-4)&(iter<50) ) {
125       iter = iter +1
126
127       #update A_1
128       prevtheta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
129       prev <- likelihood(ttnsr,prevtheta,omega)
130       W1 =kronecker(A_3,A_2)%*%t(k_unfold(C,1)@data)
131       A_1 <- comb(A_1,W1,ttnsr,1,omega)
132
133
134       # update A_2
135       W2 <- kronecker(A_3,A_1)%*%t(k_unfold(C,2)@data)
136       A_2 <- comb(A_2,W2,ttnsr,2,omega)

```

```

137
138   # update A_3
139   W3 <- kronecker(A_2,A_1)%*%t(k_unfold(C,3)@data)
140   A_3 <- comb(A_3,W3,ttnsr,3,omega)
141
142   # update C
143   W4 <- kronecker(kronecker(A_3,A_2),A_1)
144   C <- corecomb(C,W4,c(ttnsr),omega)
145   theta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
146   new <- likelihood(ttnsr,theta,omega)
147   (error <- abs((new-prev)/prev))
148 }
149 }else{
150   while ((error > 10^-4)&(iter<50) ) {
151     iter = iter +1
152
153     #update A_1
154     prevtheta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
155     prev <- likelihood(ttnsr,prevtheta,omega)
156     W1 =kronecker(A_3,A_2)%*%t(k_unfold(C,1)@data)
157     A_1 <- comb(A_1,W1,ttnsr,1,omega,alphbound)
158     if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
159
160
161     # update A_2
162     W2 <- kronecker(A_3,A_1)%*%t(k_unfold(C,2)@data)
163     A_2 <- comb(A_2,W2,ttnsr,2,omega,alphbound)
164     if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
165
166     # update A_3
167     W3 <- kronecker(A_2,A_1)%*%t(k_unfold(C,3)@data)
168     A_3 <- comb(A_3,W3,ttnsr,3,omega,alphbound)
169     if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
170
171     # update C
172     W4 <- kronecker(kronecker(A_3,A_2),A_1)
173     C <- corecomb(C,W4,c(ttnsr),omega)
174     theta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
175     new <- likelihood(ttnsr,theta,omega)

```



```

176     error <- abs((new-prev)/prev)
177     if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
178   }
179 }
180
181 result$C <- C; result$A_1 <- A_1; result$A_2 <- A_2; result$A_3 <- A_3
182 result$iteration <- iter
183 return(result)
184 }
185
186
187 fit_ordinal2 = function(ttnsr,C,A_1,A_2,A_3,omega=TRUE,alph = TRUE){
188   omega <- sort(rnorm(2))
189   alphbound <- alph+10^-4
190   result = list()
191   error<- 3
192   iter = 0
193   d1 <- nrow(A_1); d2 <- nrow(A_2); d3 <- nrow(A_3)
194   r1 <- ncol(A_1); r2 <- ncol(A_2); r3 <- ncol(A_3)
195   if (alph == TRUE) {
196     while ((error > 10^-4)&(iter<50) ) {
197       iter = iter +1
198
199       #update A_1
200       prevtheta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
201       prev <- likelihood(ttnsr,prevtheta,omega)
202       W1 =kronecker(A_3,A_2)%*%t(k_unfold(C,1)@data)
203       A_1 <- comb(A_1,W1,ttnsr,1,omega)
204
205
206       # update A_2
207       W2 <- kronecker(A_3,A_1)%*%t(k_unfold(C,2)@data)
208       A_2 <- comb(A_2,W2,ttnsr,2,omega)
209
210       # update A_3
211       W3 <- kronecker(A_2,A_1)%*%t(k_unfold(C,3)@data)
212       A_3 <- comb(A_3,W3,ttnsr,3,omega)
213
214       # update C

```

```

215     W4 <- kronecker(kronecker(A_3,A_2),A_1)
216     C <- corecomb(C,W4,c(ttnsr),omega)
217
218     #update omega
219     theta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
220     m <- polr(as.factor(c(ttnsr))~offset(-c(theta@data)))
221     omega <- m$zeta
222
223
224
225     theta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
226     new <- likelihood(ttnsr,theta,omega)
227     error <- abs((new-prev)/prev)
228   }
229 }else{
230   while ((error > 10^-4)&(iter<50) ) {
231     iter = iter +1
232
233     #update A_1
234     prevtheta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
235     prev <- likelihood(ttnsr,prevtheta,omega)
236     W1 =kronecker(A_3,A_2)%*%t(k_unfold(C,1)@data)
237     A_1 <- comb(A_1,W1,ttnsr,1,omega,alphbound)
238     if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
239
240
241     # update A_2
242     W2 <- kronecker(A_3,A_1)%*%t(k_unfold(C,2)@data)
243     A_2 <- comb(A_2,W2,ttnsr,2,omega,alphbound)
244     if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
245
246     # update A_3
247     W3 <- kronecker(A_2,A_1)%*%t(k_unfold(C,3)@data)
248     A_3 <- comb(A_3,W3,ttnsr,3,omega,alphbound)
249     if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
250
251     # update C
252     W4 <- kronecker(kronecker(A_3,A_2),A_1)
253     C <- corecomb(C,W4,c(ttnsr),omega)

```

```

254     if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
255
256     #update omega
257     theta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
258     m <- polr(as.factor(c(ttnsr))~offset(-c(theta@data)))
259     omega <- m$zeta
260
261
262     theta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
263     new <- likelihood(ttnsr,theta,omega)
264     error <- abs((new-prev)/prev)
265   }
266 }
267
268 result$C <- C; result$A_1 <- A_1; result$A_2 <- A_2; result$A_3 <- A_3
269 result$iteration <- iter; result$omega <- omega
270 return(result)
271 }

```

## 6 Stochastic gradient method.

```

1 fit_ordinal = function(ttnsr,C,A_1,A_2,A_3,omega,alph = TRUE){
2   alphbound <- alph+10^-4
3   Batchsize <- 1000
4   ransample <- sample(1:length(c(ttnsr)),Batchsize)
5   result = list()
6   error<- 3
7   iter = 0
8   d1 <- nrow(A_1); d2 <- nrow(A_2); d3 <- nrow(A_3)
9   r1 <- ncol(A_1); r2 <- ncol(A_2); r3 <- ncol(A_3)
10  if (alph == TRUE) {
11    while ((error > 10^-4)&(iter<50) ) {
12      iter = iter +1
13
14      #update A_1
15      prevtheta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
16      prev <- likelihood(ttnsr,prevtheta,omega)

```

```

17     W1 =kronecker(A_3,A_2)%*%t(k_unfold(C,1)@data)
18     A_1 <- comb(A_1,W1,ttnsr,1,omega)
19
20
21     # update A_2
22     W2 <- kronecker(A_3,A_1)%*%t(k_unfold(C,2)@data)
23     A_2 <- comb(A_2,W2,ttnsr,2,omega)
24
25     # update A_3
26     W3 <- kronecker(A_2,A_1)%*%t(k_unfold(C,3)@data)
27     A_3 <- comb(A_3,W3,ttnsr,3,omega)
28
29     # update C
30     W4 <- kronecker(kronecker(A_3,A_2),A_1)
31     C <- corecomb(C,W4[ransample,],c(ttnsr)[ransample],omega)
32     theta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
33     new <- likelihood(ttnsr,theta,omega)
34     (error <- abs((new-prev)/prev))
35 }
36 }else{
37     while ((error > 10^-4)&(iter<50) ) {
38         iter = iter +1
39
40         #update A_1
41         prevtheta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
42         prev <- likelihood(ttnsr,prevtheta,omega)
43         W1 =kronecker(A_3,A_2)%*%t(k_unfold(C,1)@data)
44         A_1 <- comb(A_1,W1,ttnsr,1,omega,alphbound)
45         if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
46
47
48         # update A_2
49         W2 <- kronecker(A_3,A_1)%*%t(k_unfold(C,2)@data)
50         A_2 <- comb(A_2,W2,ttnsr,2,omega,alphbound)
51         if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
52
53         # update A_3
54         W3 <- kronecker(A_2,A_1)%*%t(k_unfold(C,3)@data)
55         A_3 <- comb(A_3,W3,ttnsr,3,omega,alphbound)

```

```

56     if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
57
58     # update C
59     W4 <- kronecker(kronecker(A_3,A_2),A_1)
60     C <- corecomb(C,W4[ransample,],c(ttnsr)[ransample],omega)
61     theta <- ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)
62     new <- likelihood(ttnsr,theta,omega)
63     error <- abs((new-prev)/prev)
64     if(max(abs(ttm(ttm(ttm(C,A_1,1),A_2,2),A_3,3)@data))>=alph) break
65   }
66 }
67
68 result$C <- C; result$A_1 <- A_1; result$A_2 <- A_2; result$A_3 <- A_3
69 result$iteration <- iter
70 return(result)
71 }

```