

Rank estimation

Chanwoo Lee, August 31, 2020

1 Simulation 1

I reduce the size of feature brain connection matrices to 18 by 18 such that all nodes from left and right side match i.e. 9 nodes from each side of brain. First, I made ground truth coefficient \mathbf{B} with rank = 3, 5, 8, 10 and assign y_i with the following rule

$$y_i = \text{sign}(\langle \mathbf{B}, \mathbf{X}_i \rangle), \quad i = 1, \dots, n. \quad (1)$$

Second, I perform 5-folded cross validation with the same test set and training set with different rank and cost combinations such that $(\text{rank}, \text{cost}) \in \{1, \dots, 18\} \times \{2, 4, 6, 8, 10\}$. Last, I plot mean accuracy rate from 5 test sets according to rank and cost.

The following figure plot the simulation results. I observed the tendency that as rank increases, cross validation results become similar regardless of cost value. Unfortunately, this simulation shows that we fail to estimate close rank to real rank of ground truth matrix \mathbf{B} .

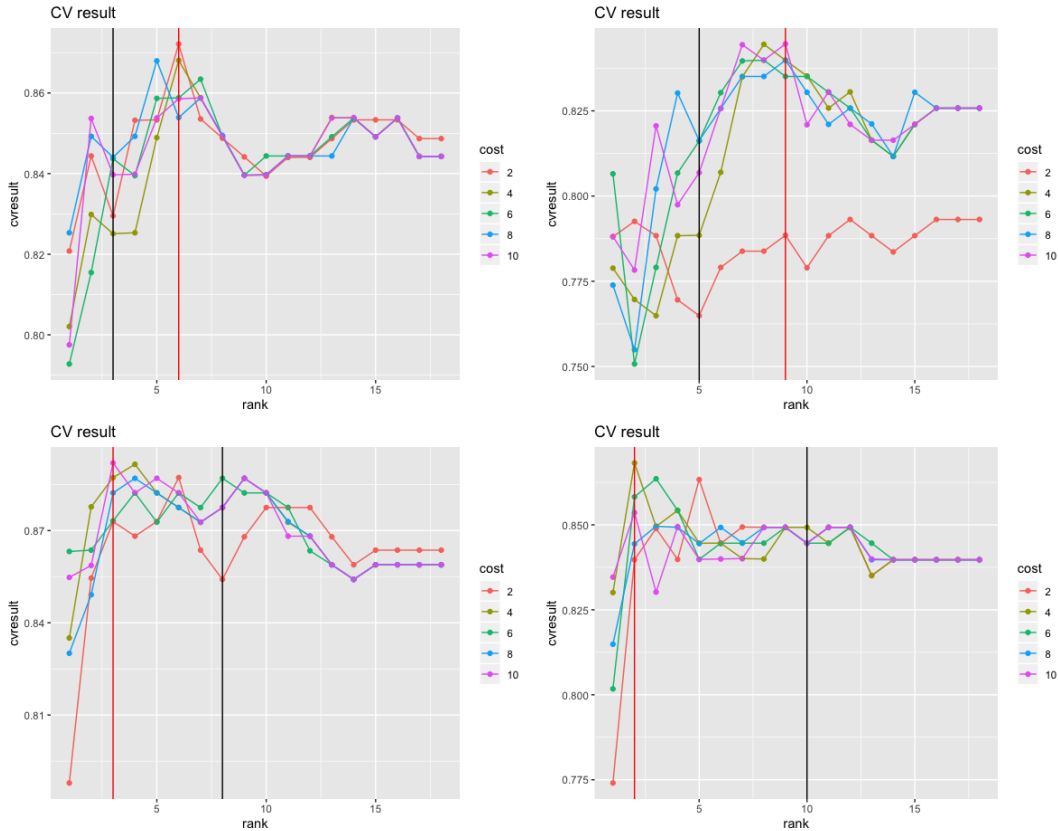


Figure 1: The figures plot mean accuracy rates according to rank and cost values. Black horizontal lines show true rank and red one show the rank which maximizes accuracy rate with certain cost values.

2 Simulation 2

Since Simulation 1 is not working well, I used centered feature matrices. We define new centered feature matrices as

$$\mathbf{X}'_i = \mathbf{X}_i - \bar{\mathbf{X}}, \quad \text{where } \bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i.$$

By similar way, I choose ground truth coefficient \mathbf{B} with rank = 3, 5, 8. and assign y_i with the rule, $y_i = \text{sign}(\langle \mathbf{B}, \mathbf{X}'_i \rangle) = \text{sign}(\langle \mathbf{B}, \mathbf{X}_i - \bar{\mathbf{X}} \rangle)$ for $i = 1, \dots, n$. From this data set, I perform 5-folded cross validation by the same way in Simulation 1.

The following figure plot the simulation results. When rank is small (3,5), estimated rank and true rank differ only by 1. However, as rank increase, rank estimation is still bad considering true rank. The tendency of insensitivity to cost value still happens as in Simulation 1 when the rank is large.

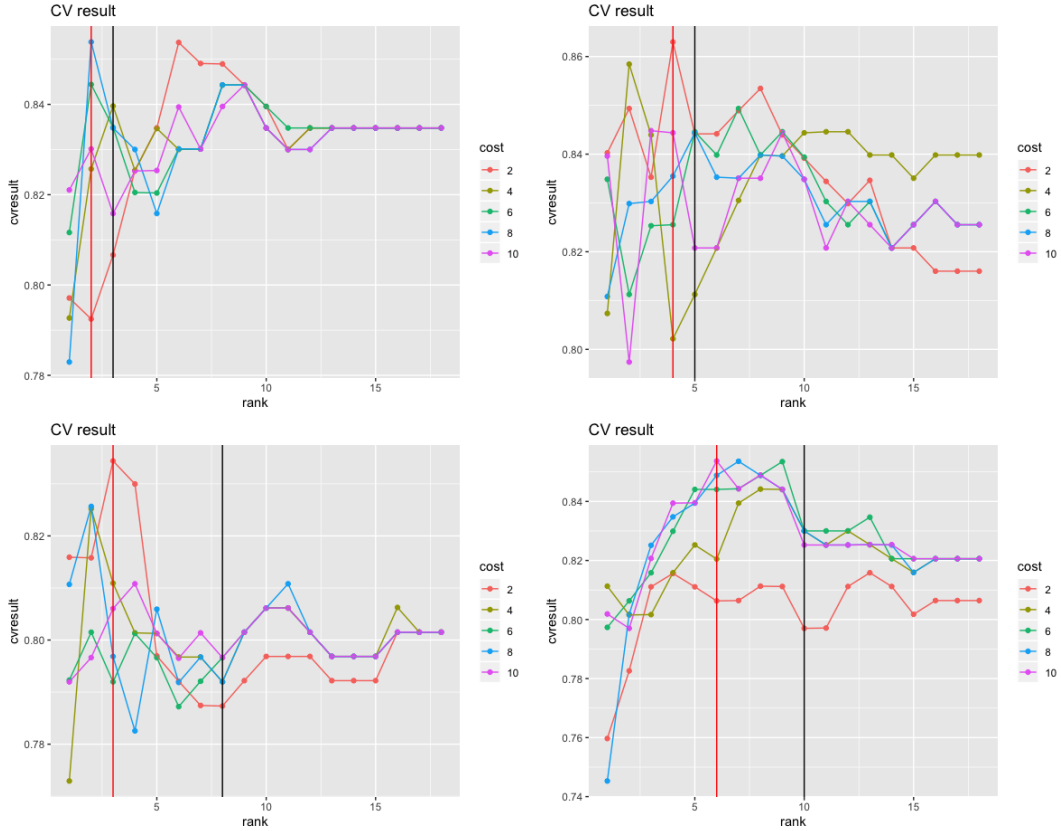


Figure 2: The figures plot mean accuracy rates according to rank and cost values. Black horizontal lines show true rank and red one show the rank which maximizes accuracy rate with certain cost values.

One positive perspective is that simulation 2 performs good when cost = 6. The following table shows the estimated rank when cost = 6 according to true rank.

	Rank 3	Rank 5	Rank 8	Rank 10
Cost = 6	2	7	10	9

Table 1: The estimated rank from 5 folded CV when cost = 6 according to true rank.

From this perspective, another way to find a good tuning parameter is to have simulation on real sized dataset with ground truth coefficient \mathbf{B} . From the simulation result, we set a tuning parameter that has the best performance for estimating true rank. To be specific, I randomly generated the coefficient matrix \mathbf{B} with rank $r = 3, 5, 8, 10$. By the similar way in Simulation 1, I assigned y_i according to (1). Lastly, I perform 5-folded CV with the same test set and training set with different rank and cost combination such that $(\text{rank}, \text{cost}) \in \{1, \dots, 20\} \times \{2, 4, \dots, 18, 20\}$. My current thought is to choose, my plan is to apply 5-folded CV to estimate the rank which minimizes error rate on test dataset.

3 Simulation 3

Feature matrices have the size 10 by 10. By similar way in Simulation 1, I plotted three different plots for each true rank $\in \{3, 5\}$. First plot A is 0-1 loss according to different combinations of rank and cost. Plot B shows Frobenius norm error of $\|\mathbf{B}_{\text{true}} - \hat{\mathbf{B}}\|_F$. The last one plots the hinge loss at a convergent point given rank and cost.

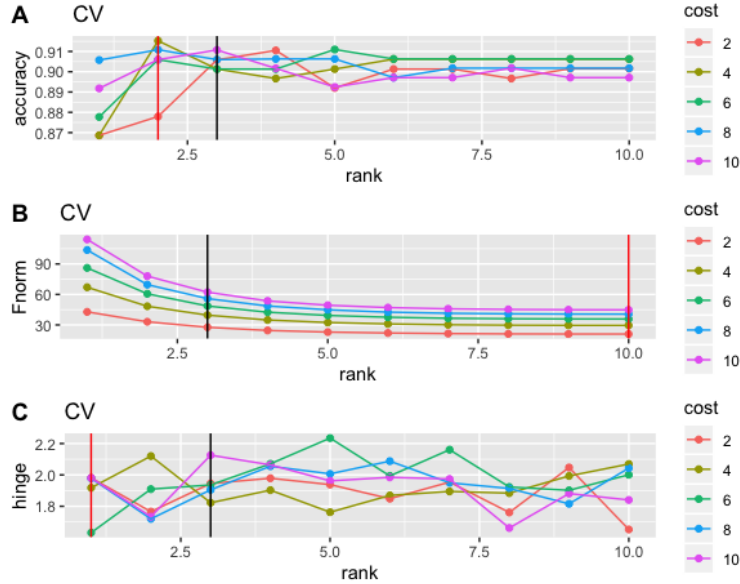


Figure 3: When true rank is 3

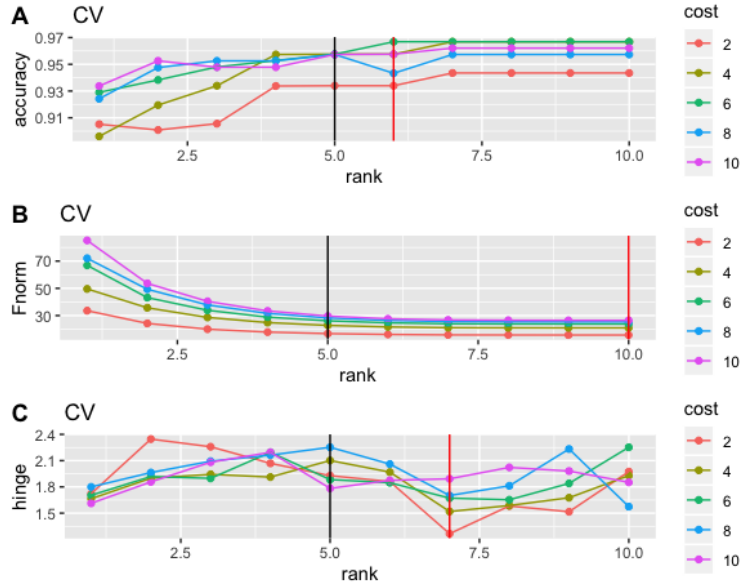


Figure 4: When true rank is 5

4 Simulation codes

```

1 # #####3
2 load(file = "bbnet68_spatial_orientation.RData")
3
4 X = list()
5 ind = c(seq(1,34,8),seq(35,68,8))
6 # X is 10 by 10
7 for(n in 1:212){
8   X[[n]] = A[ind,ind,n]
9 }
10
11 mX = matrix(0,nrow(X[[1]]),ncol(X[[1]]))
12 for(n in 1:212){
13   mX = mX+X[[n]]
14 }
15 mX = mX/212
16 # standardize
17 for(n in 1:212){
18   X[[n]] = X[[n]] - mX
19 }
20
21
22 r = 8
23
24 # Ground truth B
25 rU = matrix(runif(10*r,-1,1),nrow = 10)
26 rV = matrix(runif(10*r,-1,1),nrow = 10)
27
28 B = rU*%t(rV)
29
30 # classification
31 y = list()

```

```

32 for (i in 1:212) {
33   y[[i]] = sign(sum(B*X[[i]]))
34 }
35 y = unlist(y); length(which(y==1));length(which(y==-1)) #95/117 rank = 5, 113/99
   rank = 10, 116/96 rank =8
36 #99/113 rank = 3 #113/99 rank = 15
37
38 # standardize rank=5 111/101, rank = 3 105/107, rank = 8 99/113, rank 10 102/110
39 set.seed(1)
40 l1 = split(sample(which(y==1),length(which(y==1))),as.factor(1:5))
41 l2 = split(sample(which(y==-1),length(which(y==-1))),as.factor(1:5))
42 cindex = list()
43 for (k in 1:5) {
44   cindex[[k]] = c(l1[[k]],l2[[k]])
45 }
46
47
48 indset =matrix(nrow =50,ncol=2)
49 s = 0
50 for(p in 1:10){
51   for(q in 1:5){
52     s = s+1
53     indset[s,] = c(p,2*q)
54   }
55 }
56 }
57
58 cvresult = matrix(nrow = 50, ncol = 5)
59 colnames(cvresult) = c("accuracy","Fnorm","hinge","rank","cost")
60 for(BATCH in 1:50){
61   r = indset[BATCH,][1]
62   cost = indset[BATCH,][2]
63   cvresult[BATCH,4:5] = c(r,cost)
64   obj = 0; fnorm = 0; accuracy = 0
65   for(k in 1:5){
66     test_index = cindex[[k]]
67     train_index = setdiff(1:212,test_index)
68     train_X = X[train_index]; train_y = y[train_index]
69     con=SMMK_con(train_X,train_y,r,kernel_row="linear",kernel_col="linear",cost,
70     rep = 1, p = .5)
71     y_predict = unlist(lapply(X[test_index],con$predict))
72
73     #accuracy
74     accuracy = accuracy + length(which(y_predict-y[test_index]==0))/length(y_
75     predict)
76
77     #objective
78     obj = obj + con$obj[con$iter]
79
80     # Frobenius norm difference
81     Pr = con$P_row; Pc = con$P_col; alpha = con$alpha
82     W_row = Pr%%t(Pr);W_col= Pc%%t(Pc)
83     B1 =0; B2 = 0;
84     for(i in 1:length(train_y)){
85       B1= B1+ alpha[i]*train_y[i]*W_row%%X[[i]]
86       B2 = B2+ t(alpha[i]*train_y[i]*W_col%%t(X[[i]]))
87     }
88     approx = svd(B1+B2)
89     aB = approx$u[,1,drop = F]%%diag(approx$d[1],1)%%t(approx$v[,1,drop = F])

```

```

88
89
90     fnorm = fnorm + sqrt(sum((aB-B)^2))
91
92     print(paste("cv",r,"-",cost,"-",k," is done",sep = ""))
93
94
95 }
96 # average 5 test results
97 cvresult[BATCH,1] = accuracy/5
98 cvresult[BATCH,2] = obj/5
99 cvresult[BATCH,3] = fnorm/5
100
101 }
102
103
104
105 cvr = as.data.frame(cvresult)
106 cvr[,5] = as.factor(cvr[,5])
107 library(ggplot2)
108 g1 = ggplot(cvr,aes(x = rank,y = accuracy, color = cost))+geom_point()+
109     geom_line()+labs(title = "CV")+
110     geom_vline(xintercept = 8)+geom_vline(xintercept = 6,color = "red")
111 g2 = ggplot(cvr,aes(x = rank,y = Fnorm, color = cost))+geom_point()+
112     geom_line()+labs(title = "CV")+
113     geom_vline(xintercept = 8)+geom_vline(xintercept = 10,color = "red")
114 g3 = ggplot(cvr,aes(x = rank,y = hinge, color = cost))+geom_point()+
115     geom_line()+labs(title = "CV")+
116     geom_vline(xintercept = 8)+geom_vline(xintercept = 1,color = "red")
117
118 cvresult[which.max(cvresult[,1]),4:5] #Based on accuracy => rank 2/5, 8/3, 6/8
119 cvresult[which.min(cvresult[,2]),4:5] #Based on fnorm => rank 10/5, 10/3, 10/2
120 cvresult[which.min(cvresult[,3]),4:5] #Based on hinge loss => rank 1/5, 1/3, 1/8
121
122 ggarrange( g1,g2,g3,labels = c("A", "B","C"),ncol = 1, nrow = 3)

```