

Manual for program

Hui Zhou

Division of Biostatistics, School of Public Health, University of Minnesota

email: zhoux292@umn.edu

Nov 23, 2009

R_common_diag_cov.R and *C_common_diag_cov.cpp* tries to fit penalized model based likelihood to the data assuming common diagonal covariance matrix across different clusters.

- Since the core algorithm is carried out through calling the C++ program from R, the first step is to compile the C++ code. Type at the command prompt in Linux or Unix terminal:

```
R CMD SHLIB C_common_diag_cov.cpp
```

which would generate a *C_common_diag_cov.so* file. Now, read the function into R by typing the following in the R prompt:

```
source("R_common_diag_cov.R")
```

- Next the function is run using the command:

```
common_diag_cov(Y,Y_tune,n,n_tune,k,k0,TRUE_INDEX,BIC,num_cluster,lambda)
```

<i>Y</i>	Training data of size $n \times k$, with each row represents one observation.
<i>Y_tune</i>	Tuning data of size $n_tune \times k$
<i>n</i>	Number of observations for training data
<i>n_tune</i>	Number of observations for tuning data
<i>k</i>	Dimension for both training and tuning data
<i>k0</i>	First k_0 variables are informative; if this quantity is unknown, then z_1 and z_2 in the output should be ignored
<i>TRUE_INDEX</i>	The index showing the true group membership of each observation; if this quantity is unknown, then <i>RI</i> and <i>aRI</i> in the output should be ignored
<i>BIC</i>	Boolean variable, if TRUE then select number of cluster, penalty parameters based on BIC, otherwise, use tuning data to select them. Default is FALSE
<i>num_cluster</i>	A vector containing number of clusters to be considered. Default is from 1 to 10
<i>lambda</i>	A vector containing penalty coefficients of the mean vector. Default is from 0 to 20, with step size 1
<i>MAX_iter</i>	Maximum iteration allowed the EM algorithm. Default is 100
<i>threshold</i>	Threshold for convergence. Default is 0.0001

where each argument means:

Therefore when knowing which variables are informative, the user should move those to the left most k_0 columns of *Y* and *Y_tune*.

The function will return several R objects, which can be assigned to a variable. For example, with all default options, to save the results in the variable *out*, type the command:

```
out = common_diag_cov(Y,Y_tune,n,n_tune,k,k0,TRUE_INDEX,BIC,
num_cluster,lambda,MAX_iter,threshold)
```

To see the results, use the “\$” operator (VariableName\$ObjectName). *common_diag_cov()* returns the following objects:

- Example:

<i>num_cluster</i>	optimal number of clusters
<i>lambda</i>	optimal penalty parameter on the mean structure
<i>group_member</i>	posterior probability for each observation belonging to each cluster
<i>z_1</i>	number of estimated non-informative variables among truly informative variables
<i>z_2</i>	number of estimated non-informative variables among truly non-informative variables
<i>RI</i>	rand index between the estimated group membership and TRUE_INDEX
<i>aRI</i>	adjusted rand index between the estimated group membership and TRUE_INDEX

In this example, we first generate a dataset consist of two clusters, with the same diagonal covariance structure but with different mean vector in the first 21 variables, while the remaining 279 variables are noninformative, i.e simulated from standard Normal distribution. The details of how to generate the data is shown below. After data generation, call the function and save the results:

```
out = common_diag_cov(Y,Y_tune,n,n_tune,k,k0,TRUE_INDEX,BIC=T,
  num_cluster=seq(1:5),lambda=seq(0:20),MAX_iter=100,threshold=1e-4)
```

Then to see the output:

```
> out$num_cluster
[1] 2
> out$lambda
[1] 11
> out$z_1
[1] 0
> out$z_2
[1] 277
> out$RI
[1] 1
> out$aRI
```



```

#cluster 1      n-n0 |      Y11      Y12
#cluster 2      n0   |      Y21      Y22
Y11<-matrix(rnorm(k0*(n-n0),u1,sd1),nrow=n-n0,ncol=k0) # simulated
Y12<-matrix(rnorm((k-k0)*(n-n0),u1,sd1),nrow=n-n0,ncol=k-k0) # simulated
Y21<-matrix(rnorm(n0*k0,u2,sd2), nrow=n0,ncol=k0)+du
Y22<-matrix(rnorm( n0*(k-k0),u2,sd2), nrow=n0,ncol=k-k0)
Y<-rbind(cbind(Y11,Y12),cbind(Y21,Y22))
#
Y11<-matrix(rnorm(k0*(n_tune-n0_tune),u1,sd1),nrow=n_tune-n0_tune,ncol=k0)
Y12<-matrix(rnorm((k-k0)*(n_tune-n0_tune),u1,sd1),nrow=n_tune-n0_tune,ncol=k-k0)
Y21<-matrix(rnorm(n0_tune*k0,u2,sd2), nrow=n0_tune,ncol=k0)+du_tune
Y22<-matrix(rnorm(n0_tune*(k-k0),u2,sd2), nrow=n0_tune,ncol=k-k0)
Y_tune<-rbind(cbind(Y11,Y12),cbind(Y21,Y22))
#
mu_Y<-apply(Y,2,mean)
Y<-t(t(Y)-mu_Y)
Y_tune<-t(t(Y_tune)-mu_Y)
sd_Y<-apply(Y,2,sd)
Y<-t(t(Y)/sd_Y)
Y_tune<-t(t(Y_tune)/sd_Y)
#####
out <- common_diag_cov(Y,Y_tune,n,n_tune,k,k0,TRUE_INDEX,BIC=T,
num_cluster=seq(1:5),lambda=seq(0:20),MAX_iter=100,threshold=1e-4)
out$num_cluster
out$lambda
out$z_1
out$z_2

```

out\$RI

out\$aRI

out\$group_member