

# Using MBatch Assessments: SupervisedClustering\_Pairs\_Structures

*Tod Casasent*

*2018-06-21*

## 1 Introduction

These instructions are aimed at people familiar with R and familiar with TCGA/GDC platforms and data types. They are intended to introduce the reader to producing the given assessment. These instructions will only rarely, if ever, touch on the appropriateness of the assessment algorithm or interpretation of output. See `MBatch_01_InstallLinux` for instructions on downloading test data.

## 2 Algorithm

`SupervisedClustering_Pairs_Structures` is a function used to perform batch effects assessments using the supervised clustering algorithm for each pair of batch types provided.

## 3 Output

The primary output method for MBatch is to view results in the Batch Effects Website, described elsewhere. The PNG files are rough versions of the website output.

Graphical output is a heatmap of the correlation values, topped by a covariate bar with the batch information, and at the top dendrograms for the clustering. The columns are batch values for a single batch type. The rows are sample ids.

## 4 Usage

`SupervisedClustering_Pairs_Structures(theData, theTitle, theOutputPath, theDoHeatmapFlag, theListOfBatchPairs, theBatchTypeAndValuePairsToRemove=list(), theBatchTypeAndValuePairsToKeep=list() )`

## 5 Arguments

### 5.1 theData

An instance of `BEA_DATA`.

`BEA_DATA` objects can be created by calls of the form `new("BEA_DATA", theData, theBatches, theCovariates)`. If you have no covariate data, use an empty `data.frame` created with `data.frame()`

`mData`: Object of class "matrix" A matrix where the `colnames` are sample ids and the `rownames` are gene equivalents. All names should be strings, not factors.

`mBatches`: Object of class "data.frame" A `data.frame` where the column "names" are batch types. The first batch "type" is "Sample". All names and values should be strings, not factors or numeric.

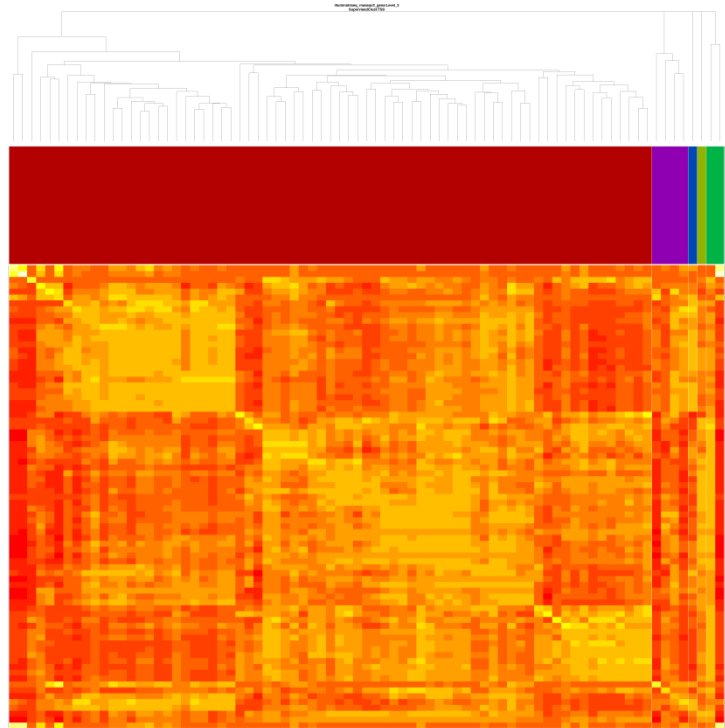


Figure 1: Supervised Clustering Example

**mCovariates:** Object of class “data.frame” A data.frame where the column “names” are covariate types. The first covariate “type” is “Sample”. All names and values should be strings, not factors or numeric.

## 5.2 theTitle

A string title to use in PNG files.

## 5.3 theOutputPath

String giving directory in which to place output PNG files.

## 5.4 theDoHeatmapFlag

A flag indicating whether or not to create HC heatmap, where TRUE means to create heatmap.

## 5.5 theListOfBatchPairs

A vector of strings, where pairs of strings give batch types to use for pairs assessment.

## 5.6 theBatchTypeAndValuePairsToRemove

A list of vectors containing the batch type (or \* for all types) and the value to remove. list() indicates none while NULL will cause an error.

## 5.7 theBatchTypeAndValuePairsToKeep

A list of vectors containing the batch type (or \* for all types) and a vector of the the value(s) to keep. list() indicates none while NULL will cause an error.

## 6 Example Call

The following code is adapted from the tests/SupervisedClustering\_Pairs\_Structures file. Data used is from the testing data as per the MBatch\_01\_InstallLinux document. In the future, we plan to make the output from MBatch more user friendly, but currently, this produces the following output at the command line.

This output can generally be skipped as very long and generally obscure. After the output is an explanation of files and directories created.

```
{
  library(MBatch)

  # set the paths
  theGeneFile="/bea_testing/MATRIX_DATA/matrix_data-Tumor.tsv"
  theBatchFile="/bea_testing/MATRIX_DATA/batches-Tumor.tsv"
  theOutputDir="/bea_testing/output/SupervisedClustering_Pairs_Structures"
  theRandomSeed=314

  # make sure the output dir exists and is empty
  unlink(theOutputDir, recursive=TRUE)
  dir.create(theOutputDir, showWarnings=FALSE, recursive=TRUE)

  # load the data and reduce the amount of data to reduce run time
  myData <- mbatchLoadFiles(theGeneFile, theBatchFile)
  myData@mData <- mbatchTrimData(myData@mData, 100000)

  # here, we take most defaults
  SupervisedClustering_Pairs_Structures(theData=myData,
    theTitle="Test Data Title",
    theOutputPath=theOutputDir,
    theDoHeatmapFlag=TRUE,
    theListOfBatchPairs=c("PlateId", "TSS", "BatchId", "TSS"),
    theBatchTypeAndValuePairsToRemove=list(),
    theBatchTypeAndValuePairsToKeep=list() )
}
```

```
## 2018 06 21 10:31:02.964 DEBUG megazone23 Changing LC_COLLATE to C for duration of run
## 2018 06 21 10:31:02.964 INFO megazone23 \ / \ / \ / \ / \ / \ / \ / \ / \ /
## 2018 06 21 10:31:02.965 INFO megazone23 Starting mbatchLoadFiles
## 2018 06 21 10:31:02.965 INFO megazone23 MBatch Version: 2017-09-19-1530
## 2018 06 21 10:31:02.965 INFO megazone23 read batch file= /bea_testing/MATRIX_DATA/batches-Tumor.tsv
## 2018 06 21 10:31:02.967 INFO megazone23 read gene file= /bea_testing/MATRIX_DATA/matrix_data-Tumor.t
## 2018 06 21 10:31:11.840 INFO megazone23 filter samples in batches using gene samples
## 2018 06 21 10:31:11.841 INFO megazone23 sort batches by gene file samples
## 2018 06 21 10:31:11.933 INFO megazone23 Finishing mbatchLoadFiles
## 2018 06 21 10:31:11.934 INFO megazone23 ~~~~~
## 2018 06 21 10:31:11.934 DEBUG megazone23 Changing LC_COLLATE to C for duration of run
## 2018 06 21 10:31:11.934 INFO megazone23 \ / \ / \ / \ / \ / \ / \ / \ / \ /
## 2018 06 21 10:31:11.935 INFO megazone23 mbatchTrimData Starting
```

```

## 2018 06 21 10:31:11.935 INFO megazone23 MBatch Version: 2017-09-19-1530
## 2018 06 21 10:31:20.301 INFO megazone23 mbatchTrimData Finishing
## 2018 06 21 10:31:20.302 INFO megazone23 ~~~~~
## 2018 06 21 10:31:20.304 DEBUG megazone23 Changing LC_COLLATE to C for duration of run
## 2018 06 21 10:31:20.304 INFO megazone23 \ / \ / \ / \ / \ / \ / \ / \ / \ /
## 2018 06 21 10:31:20.304 INFO megazone23 mbatchFilterData Starting
## 2018 06 21 10:31:20.305 INFO megazone23 MBatch Version: 2017-09-19-1530
## 2018 06 21 10:31:20.305 DEBUG megazone23 rows pre filter 1250
## 2018 06 21 10:31:20.599 DEBUG megazone23 rows post filter 1250
## 2018 06 21 10:31:20.600 DEBUG megazone23 mbatchFilterData Prefilter, gene data had 1250 while post
## 2018 06 21 10:31:20.600 DEBUG megazone23 mbatchFilterData Prefilter, batch data had 80 while post
## 2018 06 21 10:31:20.601 INFO megazone23 mbatchFilterData Finishing
## 2018 06 21 10:31:20.601 INFO megazone23 ~~~~~
## 2018 06 21 10:31:20.602 INFO megazone23 createBatchEffectsOutput_SupervisedClustering_pairs - batchT
## 2018 06 21 10:31:20.602 INFO megazone23 createBatchEffectsOutput_SupervisedClustering_pairs - batchT
## 2018 06 21 10:31:20.602 DEBUG megazone23 checkCreateDir: /bea_testing/output/SupervisedClustering_P
## 2018 06 21 10:31:20.603 INFO megazone23 makeBiasClust - starting
## 2018 06 21 10:31:20.744 INFO megazone23 makeBiasClust - quantile dat dim = 1250,80
## 2018 06 21 10:31:20.746 INFO megazone23 makeBiasClust - quantile U.data is.data.frame = FALSE
## 2018 06 21 10:31:20.746 INFO megazone23 makeBiasClust - quantile U.data is.array = TRUE
## 2018 06 21 10:31:20.747 INFO megazone23 makeBiasClust - quantile U.data is.list = FALSE
## 2018 06 21 10:31:20.747 INFO megazone23 makeBiasClust - quantile U.data nrow = 312
## 2018 06 21 10:31:20.747 INFO megazone23 makeBiasClust - quantile U.data ncol = 80
## 2018 06 21 10:31:20.747 INFO megazone23 makeBiasClust - quantile U.data length = 24960
## 2018 06 21 10:31:20.748 INFO megazone23 makeBiasClust - quantile U.data dim = 312,80
## 2018 06 21 10:31:20.748 INFO megazone23 makeBiasClust - quantile U.data is.null = FALSE
## 2018 06 21 10:31:20.748 INFO megazone23 makeBiasClust - data frame
## 2018 06 21 10:31:20.748 INFO megazone23 makeBiasClust - U.dend1 <- bias.clust
## 2018 06 21 10:31:20.752 INFO megazone23 makeBiasClust new.dis size - 80-80
## 2018 06 21 10:31:20.755 INFO megazone23 makeBiasClust orig - 80-80
## 2018 06 21 10:31:20.755 INFO megazone23 makeBiasClust is.na - 80-80
## 2018 06 21 10:31:20.755 INFO megazone23 makeBiasClust is.infinite - 80-80

## 2018 06 21 10:31:23.028 DEBUG megazone23 mbatchStandardLegend - Calling .jinit /home/linux/R/x86_64
## 2018 06 21 10:31:23.297 DEBUG megazone23 mbatchStandardLegend - .jinit complete
## 2018 06 21 10:31:23.298 DEBUG megazone23 mbatchStandardLegend - theTitle PlateId
## 2018 06 21 10:31:23.298 DEBUG megazone23 mbatchStandardLegend - theVersion MBatch 1.4.16
## 2018 06 21 10:31:23.299 DEBUG megazone23 mbatchStandardLegend - theFilePath /bea_testing/output
## 2018 06 21 10:31:23.299 DEBUG megazone23 mbatchStandardLegend - theLegendNames A29J (80)
## 2018 06 21 10:31:23.299 DEBUG megazone23 mbatchStandardLegend - theLegendNames 1
## 2018 06 21 10:31:23.300 DEBUG megazone23 mbatchStandardLegend - theLegendColors 1
## 2018 06 21 10:31:23.300 DEBUG megazone23 mbatchStandardLegend - theLegendSymbols 0
## 2018 06 21 10:31:23.300 DEBUG megazone23 mbatchStandardLegend - myColors #b30000
## 2018 06 21 10:31:23.301 DEBUG megazone23 mbatchStandardLegend before java
## 2018 06 21 10:31:23.331 DEBUG megazone23 mbatchStandardLegend after java
## 2018 06 21 10:31:23.334 DEBUG megazone23 mbatchStandardLegend - Calling .jinit /home/linux/R/x86_64
## 2018 06 21 10:31:23.616 DEBUG megazone23 mbatchStandardLegend - .jinit complete
## 2018 06 21 10:31:23.616 DEBUG megazone23 mbatchStandardLegend - theTitle TSS
## 2018 06 21 10:31:23.616 DEBUG megazone23 mbatchStandardLegend - theVersion MBatch 1.4.16
## 2018 06 21 10:31:23.617 DEBUG megazone23 mbatchStandardLegend - theFilePath /bea_testing/output
## 2018 06 21 10:31:23.617 DEBUG megazone23 mbatchStandardLegend - theLegendNames OR - University of M
## (72), OU - Roswell Park (1), P6 - Translational Genomics
## Research Institute (2), PA - University of Minnesota
## (1), PK - University Health

```

```

## Network (4)
## 2018 06 21 10:31:23.618 DEBUG megazone23 mbatchStandardLegend - theLegendNames 5
## 2018 06 21 10:31:23.618 DEBUG megazone23 mbatchStandardLegend - theLegendColors 5
## 2018 06 21 10:31:23.618 DEBUG megazone23 mbatchStandardLegend - theLegendSymbols 0
## 2018 06 21 10:31:23.619 DEBUG megazone23 mbatchStandardLegend - myColors #b30000,#8fb300,#00b347,#00
## 2018 06 21 10:31:23.619 DEBUG megazone23 mbatchStandardLegend before java
## 2018 06 21 10:31:23.684 DEBUG megazone23 mbatchStandardLegend after java
## 2018 06 21 10:31:23.685 INFO megazone23 createBatchEffectsOutput_SupervisedClustering_pairs - batchT
## 2018 06 21 10:31:23.686 INFO megazone23 createBatchEffectsOutput_SupervisedClustering_pairs - batchT
## 2018 06 21 10:31:23.688 DEBUG megazone23 checkCreateDir: /bea_testing/output/SupervisedClustering_P
## 2018 06 21 10:31:23.689 INFO megazone23 makeBiasClust - starting
## 2018 06 21 10:31:23.868 INFO megazone23 makeBiasClust - quantile dat dim = 1250,80
## 2018 06 21 10:31:23.870 INFO megazone23 makeBiasClust - quantile U.data is.data.frame = FALSE
## 2018 06 21 10:31:23.871 INFO megazone23 makeBiasClust - quantile U.data is.array = TRUE
## 2018 06 21 10:31:23.871 INFO megazone23 makeBiasClust - quantile U.data is.list = FALSE
## 2018 06 21 10:31:23.871 INFO megazone23 makeBiasClust - quantile U.data nrow = 312
## 2018 06 21 10:31:23.872 INFO megazone23 makeBiasClust - quantile U.data ncol = 80
## 2018 06 21 10:31:23.872 INFO megazone23 makeBiasClust - quantile U.data length = 24960
## 2018 06 21 10:31:23.872 INFO megazone23 makeBiasClust - quantile U.data dim = 312,80
## 2018 06 21 10:31:23.872 INFO megazone23 makeBiasClust - quantile U.data is.null = FALSE
## 2018 06 21 10:31:23.873 INFO megazone23 makeBiasClust - data frame
## 2018 06 21 10:31:23.873 INFO megazone23 makeBiasClust - U.dend1 <- bias.clust
## 2018 06 21 10:31:23.877 INFO megazone23 makeBiasClust new.dis size - 80-80
## 2018 06 21 10:31:23.878 INFO megazone23 makeBiasClust orig - 80-80
## 2018 06 21 10:31:23.879 INFO megazone23 makeBiasClust is.na - 80-80
## 2018 06 21 10:31:23.879 INFO megazone23 makeBiasClust is.infinite - 80-80

## 2018 06 21 10:31:26.153 DEBUG megazone23 mbatchStandardLegend - Calling .jinit /home/linux/R/x86_64
## 2018 06 21 10:31:26.431 DEBUG megazone23 mbatchStandardLegend - .jinit complete
## 2018 06 21 10:31:26.431 DEBUG megazone23 mbatchStandardLegend - theTitle BatchId
## 2018 06 21 10:31:26.432 DEBUG megazone23 mbatchStandardLegend - theVersion MBatch 1.4.16
## 2018 06 21 10:31:26.432 DEBUG megazone23 mbatchStandardLegend - theFilenamePath /bea_testing/output
## 2018 06 21 10:31:26.432 DEBUG megazone23 mbatchStandardLegend - theLegendNames 00304 (80)
## 2018 06 21 10:31:26.433 DEBUG megazone23 mbatchStandardLegend - theLegendNames 1
## 2018 06 21 10:31:26.433 DEBUG megazone23 mbatchStandardLegend - theLegendColors 1
## 2018 06 21 10:31:26.434 DEBUG megazone23 mbatchStandardLegend - theLegendSymbols 0
## 2018 06 21 10:31:26.434 DEBUG megazone23 mbatchStandardLegend - myColors #b30000
## 2018 06 21 10:31:26.434 DEBUG megazone23 mbatchStandardLegend before java
## 2018 06 21 10:31:26.449 DEBUG megazone23 mbatchStandardLegend after java
## 2018 06 21 10:31:26.454 DEBUG megazone23 mbatchStandardLegend - Calling .jinit /home/linux/R/x86_64
## 2018 06 21 10:31:26.744 DEBUG megazone23 mbatchStandardLegend - .jinit complete
## 2018 06 21 10:31:26.744 DEBUG megazone23 mbatchStandardLegend - theTitle TSS
## 2018 06 21 10:31:26.745 DEBUG megazone23 mbatchStandardLegend - theVersion MBatch 1.4.16
## 2018 06 21 10:31:26.745 DEBUG megazone23 mbatchStandardLegend - theFilenamePath /bea_testing/output
## 2018 06 21 10:31:26.745 DEBUG megazone23 mbatchStandardLegend - theLegendNames OR - University of M
## (72), OU - Roswell Park (1), P6 - Translational Genomics
## Research Institute (2), PA - University of Minnesota
## (1), PK - University Health
## Network (4)
## 2018 06 21 10:31:26.746 DEBUG megazone23 mbatchStandardLegend - theLegendNames 5
## 2018 06 21 10:31:26.746 DEBUG megazone23 mbatchStandardLegend - theLegendColors 5
## 2018 06 21 10:31:26.746 DEBUG megazone23 mbatchStandardLegend - theLegendSymbols 0
## 2018 06 21 10:31:26.747 DEBUG megazone23 mbatchStandardLegend - myColors #b30000,#8fb300,#00b347,#00
## 2018 06 21 10:31:26.747 DEBUG megazone23 mbatchStandardLegend before java

```

```
## 2018 06 21 10:31:26.805 DEBUG megazone23 mbatchStandardLegend after java
```

## 7 Example File Output

The above code creates the following subdirectories and files. The subdirectories correspond to the Batch Type Pairs on which assessments were requested.

```
/bea_testing/output/SupervisedClustering_Pairs_Structures$ ls -l
total 8
drwxr-xr-x 2 linux linux 4096 Jun 14 12:56 BatchId-TSS
drwxr-xr-x 2 linux linux 4096 Jun 14 12:56 PlateId-TSS
```

Looking at the “BatchId-TSS” subdirectory, it contains the following diagram and legend files. This algorithm does not currently generate data usable with dynamic displays.

```
/bea_testing/output/SupervisedClustering_Pairs_Structures/BatchId-TSS$ ls -l
total 276
-rw-r--r-- 1 linux linux 261168 Jun 19 09:58 SupervisedClust_Diagram.png
-rw-r--r-- 1 linux linux 2701 Jun 19 09:58 SupervisedClust_Legend-BatchId.png
-rw-r--r-- 1 linux linux 12899 Jun 19 09:58 SupervisedClust_Legend-TSS.png
```

Here is the diagram generated from this code.

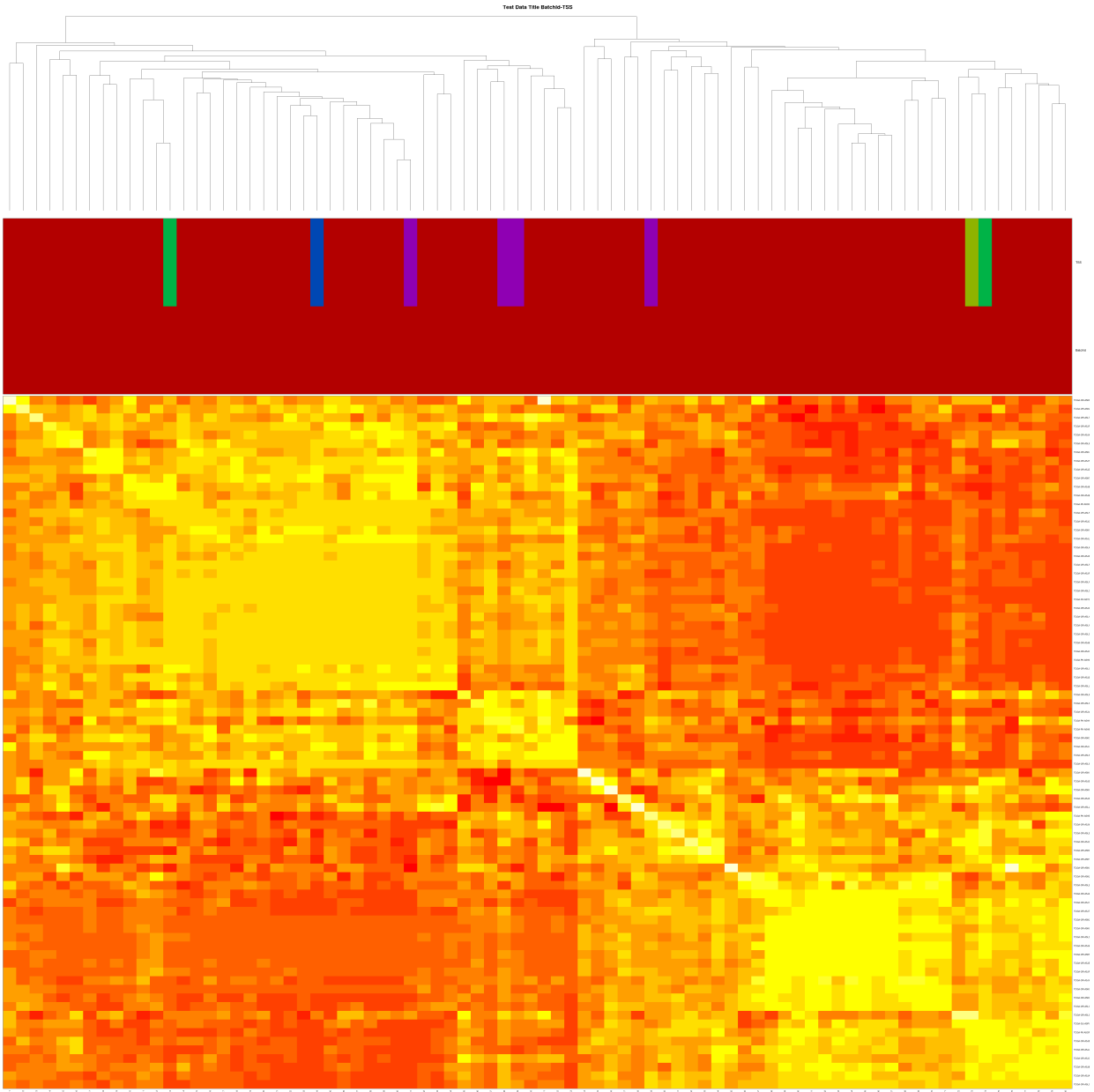


Figure 2: SupervisedClustering\_Pairs\_Structures Output