

CSFA 0.1 - Vignette

Ewoud De Troyer

1 Introduction

One of the many challenges in today’s omics data is the goal of connecting those compounds/molecules/samples together which have similar properties by gene expression. Techniques like this allow the discovery of new molecule properties by connecting their signatures with those derived from already well-known ones.

Papers such as Lamb et al. (2006) and Zhang and Gant (2008) both already took up the challenge of dealing with this problem. In Lamb et al. (2006), a reference collection of gene expression profiles from human cells treated with bioactive small molecules was created in order to design a systematic approach to discover these functional connections. While their approach achieved a good degree of succes, it was unable to measure statistical significance. This is where the paper from Zhang and Gant (2008) continued for example. Their paper offers a more principled statistical procedure to test connections between the compounds which allows the valuation of statistical significance.

The CSFA package accompanies the paper/report by Shkedy, Z. and De Troyer. E. (ADD REAL REF), which proposes the usage of *factor analysis* methods (Principal Component Analysis (=PCA), (Sparse) Multiple Factor Analysis (MFA) (Abdi et al., 2013) and FABIA (Factor Analysis for Bicluster Acquisition) (Hochreiter et al., 2010)) to derive the connectivity between compounds. Using these methods, not only do you obtain information about the connectivity between the compounds, you also get information about which genes are responsible for guiding this connectivity.

Further instead of computing a pairwise correlation/connection score between the compounds, now the entire available data is being used to look for dominant structures on both dimenstions. This is very similar to try to discover biclusters in the data. Consequently, it is not necessary anymore to decide upon a cut-off for up- and downregulated genes since you will be using all the genes to do the factor analysis.

It should also be noted that the setting in which the *factor analysis* is applied, is slightly different from the one in the Connectivity Map (Lamb et al., 2006). In the Connectivity Map there is a large data set of references profiles to which the query signatures are compared. In this setting, the meaning of ‘reference’ and ‘query’ will be switched around. You start with a small set of references, namely a small set of samples of which they are similar. These are compared with a larger set of queries in order to try to discover samples or compounds similar to the reference set.

Further, since the methods will be applied on a matrix which consists out of both the reference and the query profiles, the number of genes for these signatures will have to be the same.

$$\begin{array}{c} \begin{array}{cc} \text{Reference Samples} & \text{Query Samples} \\ \left[\begin{array}{cc} \mathbf{X}_1 & \mathbf{X}_2 \end{array} \right] & \begin{array}{c} \text{\textit{g genes}} \end{array} \end{array} \\ \text{\textit{n samples}} \end{array} \quad (1)$$

Finally in order to easily compare these methods with the Zhang and Gant Score, CSFA also includes an implementation of this algorithm together with the ability to compare the scores with the FA scores.

2 Data

In order to showcase the functionality of CSFA, some simulated microarray data will be used. The data contains 1000 genes and 341 compounds of which 6 will be used as reference signatures. The remaining query signatures consist out of 5 strongly positive connected compounds, 20 weakly positive connected compounds, 10 strongly negative connected compounds and 300 compounds which are not connected at all.

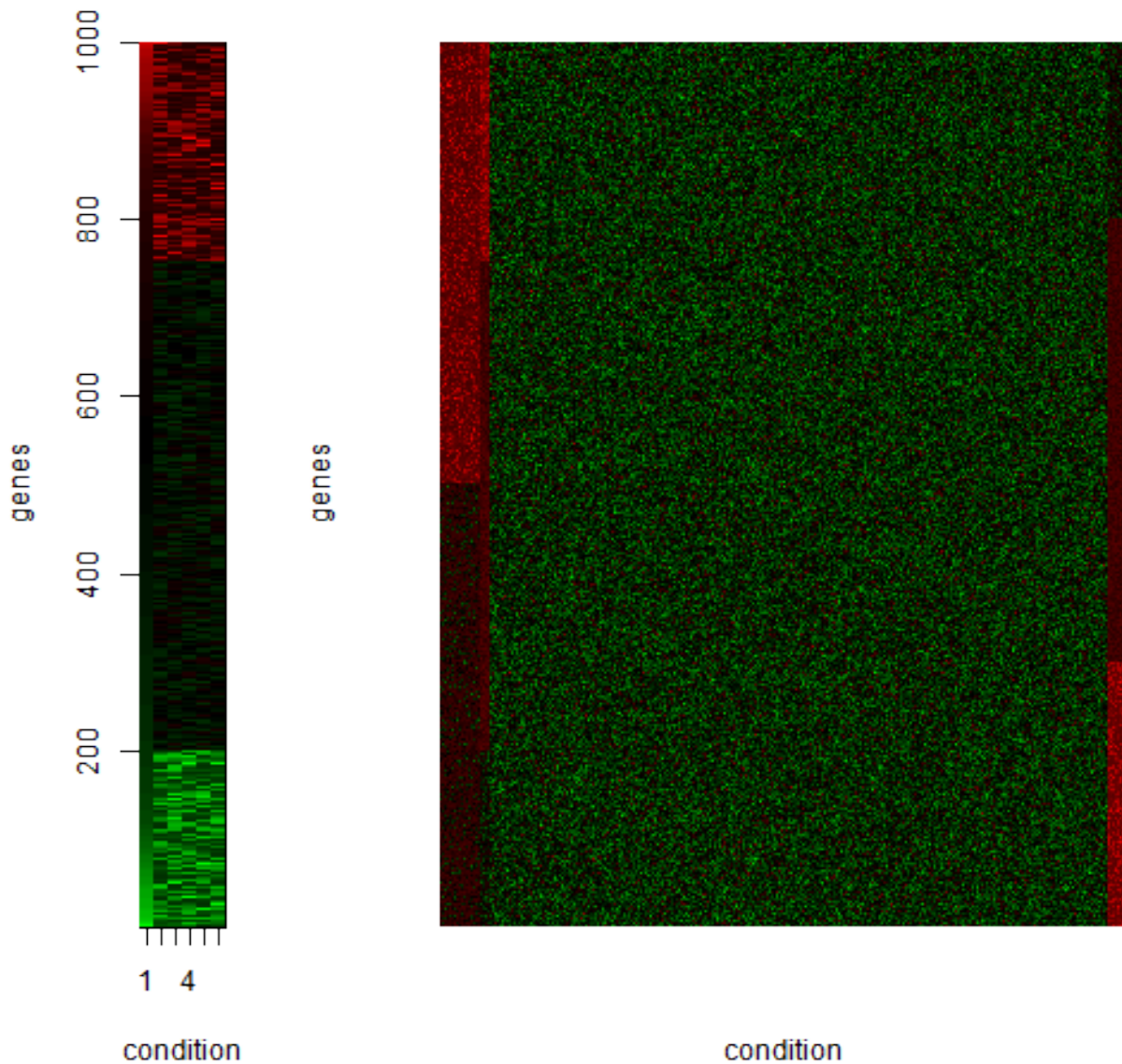


Figure 1: Heatmap of Reference and Query Matrix

3 Example CS Analysis

Start by first loading both the CSFA library and the example data available in the package. The simulated data is split up in the reference and query matrix.

```
library(CSFA)
data("dataSIM",package="CSFA")

refMat <- dataSIM[,c(1:6)]
querMat <- dataSIM[,-c(1:6)]
```

Next, the Connectivity Scores from Zhang and Gant, MFA and FABIA will be computed with the package. The last two methods will also provide scores for the genes involved in the structure.

More details about the connectivity and gene scores as well as the decision making of which component to look at can be found in Shkedy, Z. and De Troyer, E. (ADD REF).

3.1 Zhang and Gant

The Zhang and Gant scores are computed with the default parameters. This means all the genes will be used (no cut-off) and the query signature will be considered as an ordered signature. Also no permutation will be applied by default.

Note that for the vignette, which is a sweave document, we use the "sweave" `plot.type`. Normally you would be using either "device" or "pdf".

```
out_ZG <- CSanalysis(refMat,querMat,"CSzhang",plot.type="sweave")

##      posname  posscore negname   negscore
## 1    cSP-7  0.8159043   cSN-6 -0.74898237
## 2    cSP-10 0.8137249   cSN-8 -0.74288930
## 3    cSP-9  0.8093167   cSN-3 -0.73490411
## 4    cSP-8  0.8081590   cSN-2 -0.73414170
## 5    cSP-6  0.8081177  cSN-10 -0.73401815
## 6    cWP-17 0.5702188   cSN-7 -0.72983443
## 7    cWP-19 0.5694001   cSN-1 -0.72942985
## 8    cWP-12 0.5672755   cSN-9 -0.72934793
## 9    cWP-11 0.5659705   cSN-4 -0.72832312
## 10   cWP-7  0.5658364   cSN-5 -0.72709705
## 11   cWP-4  0.5649765    c-25 -0.03391145
## 12   cWP-9  0.5648644   c-264 -0.03185064
## 13   cWP-3  0.5621351    c-28 -0.03167404
## 14   cWP-1  0.5590061    c-16 -0.03005021
## 15   cWP-8  0.5581202   c-220 -0.02962786
## 16   cWP-5  0.5578107    c-32 -0.02801159
## 17   cWP-18 0.5485421    c-94 -0.02700210
## 18   cWP-15 0.5481818    c-27 -0.02650816
## 19   cWP-16 0.5473567   c-228 -0.02583543
## 20   cWP-2  0.5473340   c-169 -0.02576697
```

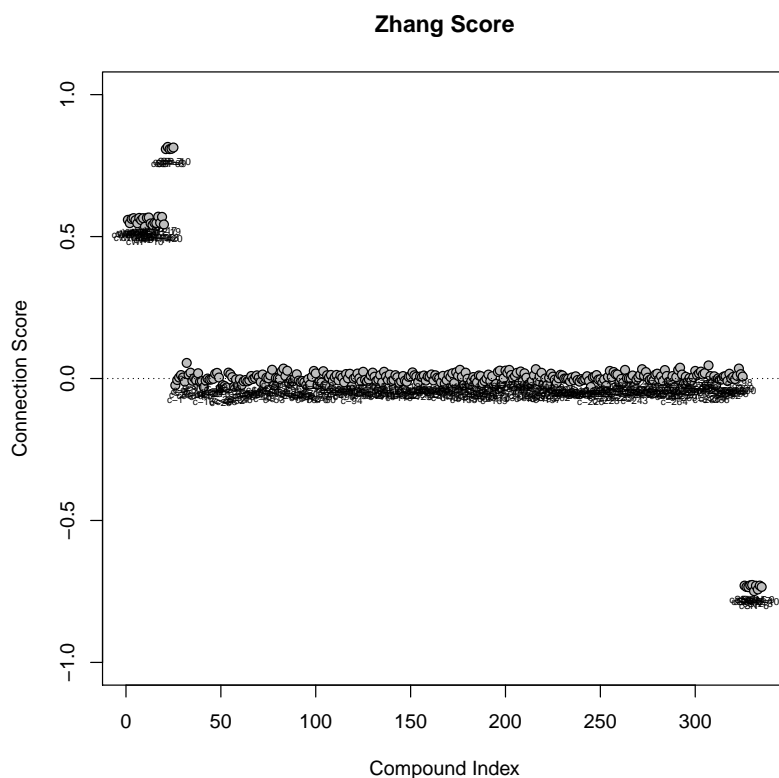


Figure 2: CSanalysis Graphs for CSzhang

While all the connectivity scores can be found in the `out_ZG` object, the function already prints by default the top 20 positive and negative connectivity scores. Figure 2 clearly shows the positive (weak and strong) and negative connected compounds.

3.2 MFA

The next CS analysis which is applied is the one using Multiple Factor Analysis (MFA) by setting the `type` to `"CSmfa"`. Three of the available plots were chosen, namely the Reference Loadings, Gene Scores, Compound Loadings (Connectivity Scores) and Compound Profiles (`which=c(2,4,5)`).

Note that in the R-code we already preselected which component to investigate with `factor.plot`. Further we also already decided which columns of the query matrix we would like to draw in the compound profiles graph with `column.interest`. Indices 1, 2 and 3 coincide with 3 weakly positive connected compounds.

However, if you are not sure beforehand what you want to investigate, you can also decide upon these parameters on the fly interactively. To do this simply set these parameters to `NULL` or leave them out.

To determine `factor.plot`, you will be able to click on the factor you want to observe in the *"Loadings for Ref..."* plot. After all, this graph will be your main guideline on which factor is capturing the structure of your reference set of signatures. As shown in Figure 3 below, this is clearly the first factor.

Next, the `column.interest` parameter can be chosen in the *"Compound Loadings"* plot. You can left-click on multiple compounds you wish to draw in the compound profiles graph (and right-click to stop the selection procedure).

```
out_MFA <- CSanalysis(refMat,querMat,"CSmfa",plot.type="sweave",which=c(2,4,5),
                      factor.plot=1,column.interest=c(1,2,3))

## Echoufier Rv Correlation:
##           Reference      Query      MFA
## Reference 1.0000000 0.4947578 0.8001384
## Query      0.4947578 1.0000000 0.9171329
## MFA        0.8001384 0.9171329 1.0000000
```

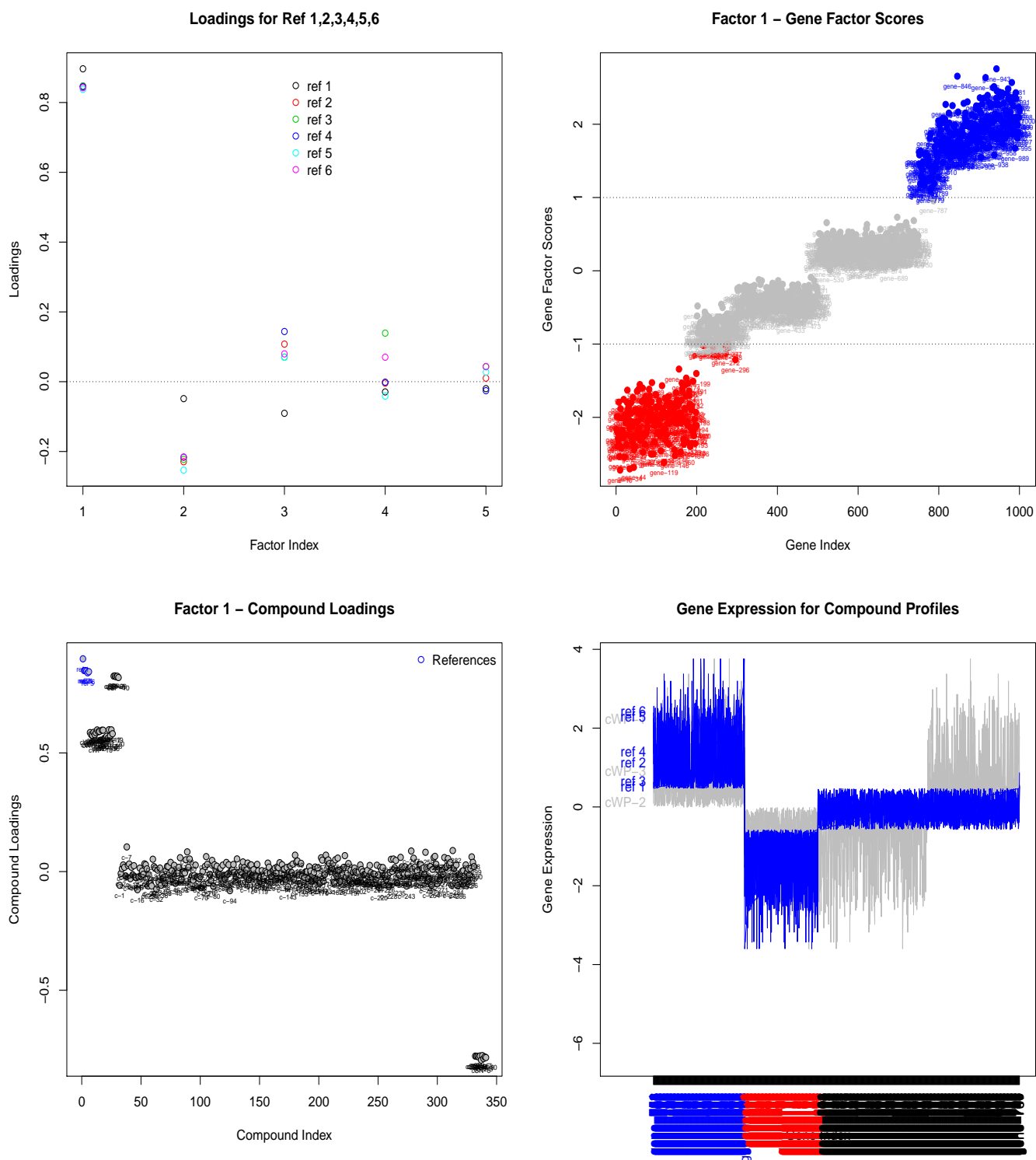


Figure 3: CSanalysis Graphs for CSfma

Just like in the Zhang and Gant plot, we again see that the simulated positive and negative connected compounds are appearing in the Compound Loadings plot. However now we also get a plot showing the scores of the genes involved in the structure of the first factor in the MFA analysis. Finally also note that the `CSanalysis` function automatically prints the Echoufier Rv Correlation matrix between the Reference, Query and MFA matrix.

3.3 FABIA

The last analysis is done with FABIA, Factor Analysis for Bicluster Acquisition (`type="CSfabia"`). We will only select 2 plots this time, namely the reference loadings and compound loadings (`which=c(2,5)`). However in contrary with the MFA analysis, we can select 2 components for this analysis. Based on the reference loadings we decide to select bicluster 1 and 2 (`BC.plot=c(1,2)`).

This time we also do some manual coloring of the columns to highlight some strongly connected compounds with `color.columns`. We start by making a vector of length 341 (column dimension of example data) and fill it with the color black. Next we fill in the color blue for the 6 reference compounds and red for 3 of the strongly positive connected compounds. We also change the legend according to this coloring.

Note that we have also set a seed just before the FABIA analysis in order to have a reproducible result.

```
color.columns <- rep("black",dim(dataSIM)[2])
color.columns[1:6] <- "blue"
color.columns[c(29,30,31)] <- "red"

set.seed(8956)
out_FABIA <- CSanalysis(refMat,querMat,"CSfabia",plot.type="sweave",which=c(2,5),
                        color.columns=color.columns,
                        legend.names=c("References","SP Connected"),
                        legend.cols=c("blue","red"), BC.plot=c(1,2),
                        gene.thresP=2,gene.thresN=-2)
```

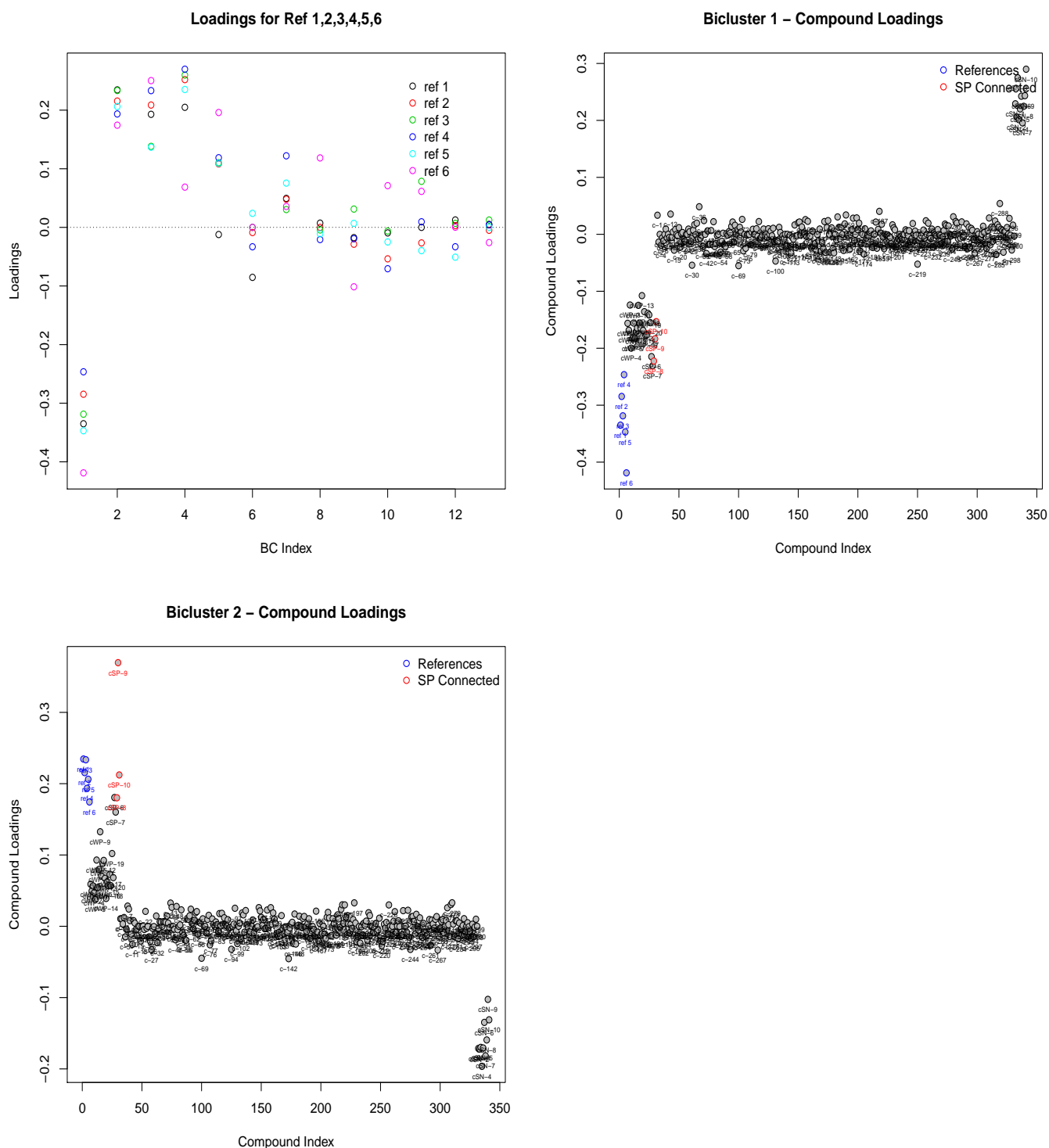


Figure 4: CSanalysis Graphs for CSfabia

The results in Figure 4 are very comparable with the Zhang and MFA graphs.

4 Example CS permutation

The CSFA package also contains a function called `CSpermute`. With this function it is possible to compute p-values through permutation for the MFA and Zhang & Gant results. These results will be added to the `CS` slot of both the MFA and ZG results. More information is also entered in the `permutation.object` slot.

First, let us apply the permutation on the MFA and ZG result without plotting any plots just yet by putting

which to `c()`. The number of permutation was chosen to only be 100 in this case. Further, the p-values are also adjusted for multiplicity by setting a value for `method.adjust` different than "none".

(Note: it is possible for MFA results to compute the p-values in a different factor than the one chosen in the CS analysis. This is accomplished through the `mfa.factor` parameter.)

```
out_MFA <- CSpermute(refMat,querMat,CSresult=out_MFA,B=100,method.adjust="BH",
                     which=c(),verbose=FALSE)
out_ZG <- CSpermute(refMat,querMat,CSresult=out_ZG,B=100,method.adjust="BH",
                     which=c(),verbose=FALSE)
```

```
head(out_MFA@CS$CS.query)
```

```
##          Factor1 pvalues pvalues.adjusted
## cWP-1 0.5858096      0              0
## cWP-2 0.5771857      0              0
## cWP-3 0.5739775      0              0
## cWP-4 0.5867910      0              0
## cWP-5 0.5770335      0              0
## cWP-6 0.5852554      0              0
```

```
head(out_ZG@CS$CS.query)
```

```
##          ZhangScore pvalues pvalues.adjusted
## cWP-1 0.5590061      0              0
## cWP-2 0.5473340      0              0
## cWP-3 0.5621351      0              0
## cWP-4 0.5649765      0              0
## cWP-5 0.5578107      0              0
## cWP-6 0.5468158      0              0
```

Next, we can actually re-use the updated `out_MFA` and `out_ZG` in `CSpermute`. As long as the number of permutations (`B`) is not changed, the permutation will not need to be computed all over again. This means you can plot the available graphs (which: 1, volcano plot ; 2, connectivity score compound distribution histogram under null hypothesis with p-value) as many times as needed. The parameter `cmpd.hist` decides which compounds should be used for the second type of plot. If this parameter is not given (`NULL`), you can interactively choose them on the volcano plot by left-clicking on them (and right-click to stop). In the code below, we plot both type of graphs for the MFA result with a pre-determined `cmpd.hist`.

```
out_MFA <- CSpermute(refMat,querMat,out_MFA,B=100,method.adjust="BH",
                     which=c(1,2),cmpd.hist=c(29,105),plot.type="sweave")
```

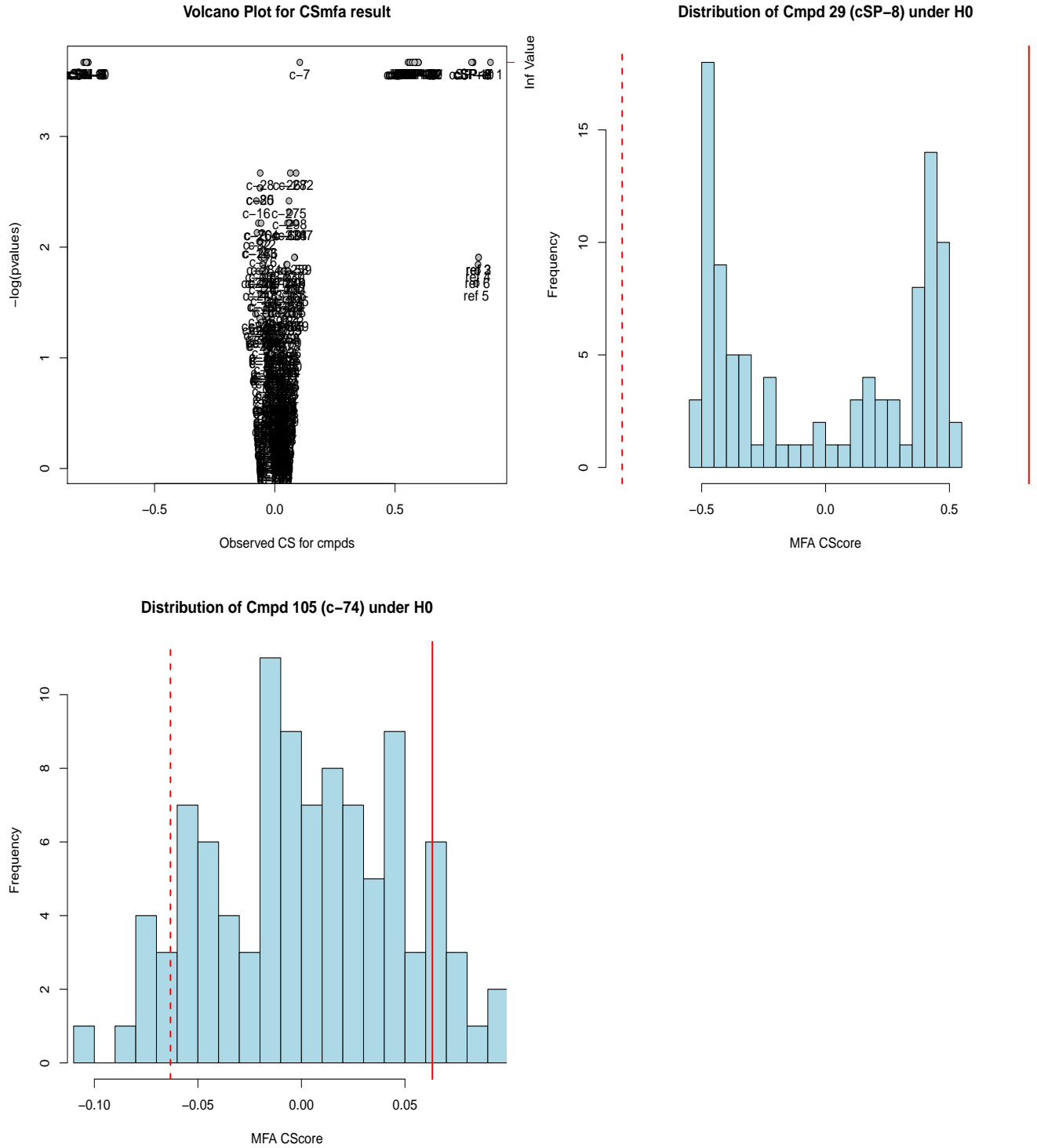



Figure 5: CSpermute graphs for MFA result

5 Example Compare CS Results

Finally, CSFA also provides a way to quickly compare the 2 results on the same data.

In the R-code below, we first compare the Zhang and Gant results with the MFA result. With `component2.plot=1` we choose the first component for the second result which is the first factor for the MFA result in this example. Since the Zhang and Gant analysis only provides connectivity scores, only 1 comparison graph will be created. The second example in the code compares the MFA with the FABIA results. For both results we choose the first component which corresponds with the first factor and first bicluster. Further, we also set some positive and negative gene thresholds for both of the results. In this example we keep them the same for both the MFA

and FABIA results namely 2 for the upper threshold and -2 for the lower one. This time since both results also contain gene scores, 2 graphs will be created. Further because we set thresholds for the genes, the gene score comparison plot will be coloring according to these thresholds.

```
corr_ZG_MFA <- CScompare(out_ZG,out_MFA,component2.plot=1,plot.type="sweave")

corr_MFA_FABIA <- CScompare(out_MFA,out_FABIA,component1.plot=1,component2.plot=1,
                             gene.thresP=c(2,2),gene.thresN=c(-2,-2),plot.type="sweave")

corr_ZG_MFA$correlation

## CS Correlation
##      0.9931169

corr_MFA_FABIA$correlation

## CS Correlation GS Correlation
##      -0.9573119      -0.7752930
```

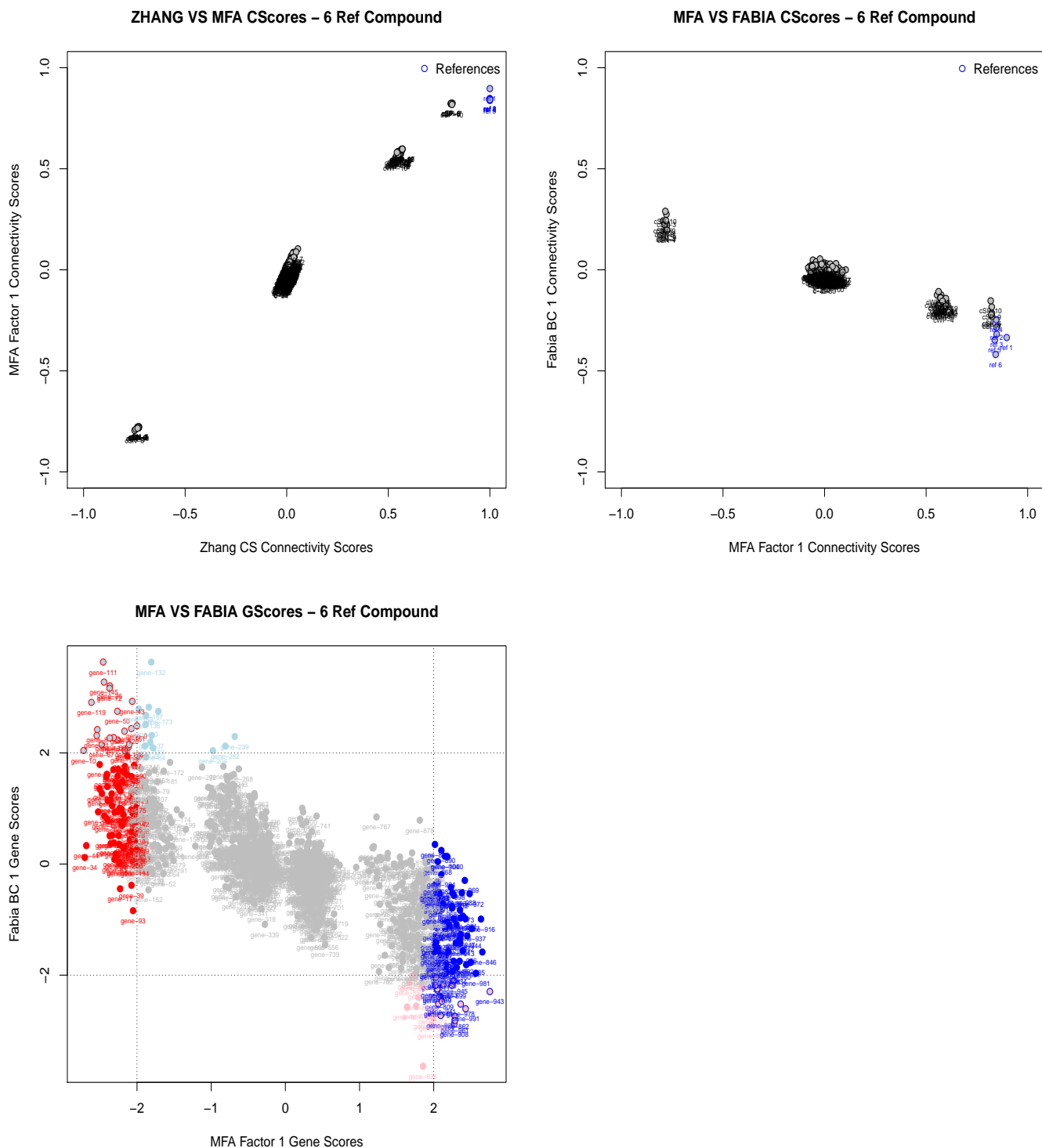


Figure 6: Compare CSresults

Note that apart from the scatter plots in Figure 6, the function also returns the pearson correlation between the scores.

Further because both the MFA and ZG contain p-values and adjusted p-values, the returned object also contains a small comparison between the number of significant p-values. The significance threshold can be changed with the `threshold.pvalues` parameter and is defaulted to 0.05.

```
corr_ZG_MFA$compare.pvalues
##          Result1.Sign Result1.NotSign
## Result2.Sign          37            0
```

```
## Result2.NotSign      47      257
      corr_ZG_MFA$compare.pvalues.adjusted
##      Result1.Sign Result1.NotSign
## Result2.Sign      37      0
## Result2.NotSign    11     293
```

References

- Abdi, H., Williams, L. J., and Valentin, D. (2013), “Multiple factor analysis: principal component analysis for multitable and multiblock data sets,” *WIREs Comput Stat*, 1–31.
- Hochreiter, S., Bodenhoger, U., Heusel, M., Mayr, A., Mitterecker, A., Kasim, A., Khamiakova, T., Sanden, S., Lin, D., Talloen, W., Bijmens, Göhlmann, H., Shkedy, Z., and Clevert, D.-A. (2010), “FABIA: Factor Analysis for Bicluster acquisition,” *Bioinformatics*, 26, 1520–1527.
- Lamb, J., Crawford, E. D., Peck, D., Modell, J. D., Blat, I. C., Wrobel, M. J., Lerner, J., Brunet, J.-P., Subramanian, A., Ross, K. N., Reich, M., Hieronymus, H., Wei, G., Armstrong, S. A., Haggarty, S. J., Clemons, P. A., Wei, R., Carr, A., Lander, E. S., and Golub, T. R. (2006), “The Connectivity Map: Using Gene-Expression Signatures to Connect Small Molecules, Genes, and Disease,” *Science*, 313, 1929–1934.
- Zhang, S.-D. and Gant, T. W. (2008), “A simple and robust method for connecting small-molecule drugs using gene-expression signatures,” *BMC Bioinformatics*, 9, 10.