

Tutorial on studying module preservation:

IV. Preservation of female mouse liver modules in male data

Peter Langfelder and Steve Horvath

October 25, 2010

Contents

1	Overview	1
1.a	Setting up the R session	1
2	Data input and rudimentary network analysis	2
2.a	Loading expression data	2
2.b	Calculating topological overlap and clustering	3
3	Module preservation	5
3.a	Calculation of module preservation	5
3.b	Analysis and display of module preservation results	5
4	Alternative methods: clusterRepro	10
4.a	Calculation of preservation statistics	10
4.b	Analysis of results	11

1 Overview

This tutorial presents an introductory analysis of module preservation between two independent data sets. The data we use here are expression measurements in livers of a mouse F2 cross described in [1]. We use the data set of female mice as the reference set, and the male mice as the test set. Because male and female mice are clearly very similar, one would expect that most modules will be preserved. In this sense, module preservation serves as validation for the modules found in the female data. However, metabolism is known to be a sexually dimorphic trait, and it is possible that some of the modules found in the female samples are female-specific and will not be found in the male data. Thus, module preservation can also serve as a tool for differential network analysis.

Here we first take a look at the data to get a feeling for the module structure, and to check visually how much similarity there is between the male and female networks. We then use the function `modulePreservation` to quantify module preservation using multiple preservation statistics.

We encourage readers unfamiliar with any of the functions used in this tutorial to open an R session and type

```
help(functionName)
```

(replace `functionName` with the actual name of the function) to get a detailed description of what the functions does, what the input arguments mean, and what is the output.

1.a Setting up the R session

After starting R we execute a few commands to set the working directory and load the requisite packages:

```

# Display the current working directory
getwd();
# If necessary, change the path below to the directory where the data files are stored.
# "." means current directory. On Windows use a forward slash / instead of the usual \.
workingDir = ".";
setwd(workingDir);
# Load the package
library(WGCNA);
# The following setting is important, do not omit.
options(stringsAsFactors = FALSE);

```

2 Data input and rudimentary network analysis

In this section we illustrate data loading and basic network construction. We assume the user has downloaded the data for this tutorial and saved them in a local directory. At the end of this section we save the results necessary for the subsequent analysis, so this section may be skipped in subsequent runs.

2.a Loading expression data

The expression data is contained in two files `cnew_liver_bxh_f2female_8000mvgenes_p3600_UNIQUE_tommodules.xls.bz2` and `LiverMaleFromLiverFemale3600.csv.bz2` that come with this tutorial.

```

#Read in the female liver data set
# The following 3421 probe set were arrived at using the following steps
#1) reduce to the 8000 most varying, 2) 3600 most connected, 3) focus on unique genes
file = bzfile("cnew_liver_bxh_f2female_8000mvgenes_p3600_UNIQUE_tommodules.xls.bz2");
dat0=read.table(file, header=TRUE)
names(dat0)

```

In addition to expression data, the data files contain extra information about the surveyed genes we do not need. One can inspect larger data frames such as `femData` by invoking R data editor via `fix(femData)`. The expression data set contains 135 samples. Note that each row corresponds to a gene and column to a sample or auxiliary information. We now remove the auxiliary data and transpose the expression data for further analysis.

```

# this contains information on the genes
datSummary=dat0[,c(1:8,144:150)]
# the following data frame contains
# the gene expression data: columns are genes, rows are arrays (samples)
datExprFemale <- t(dat0[,9:143])
no.samples <- dim(datExprFemale)[[1]]
dim(datExprFemale)
# Set the columns names to probe names
colnames(datExprFemale) = datSummary$substanceBXH
# This module assignment was obtained by Ghazalpour et al
colorsFemale = dat0$module

```

We now similarly read in the male data and put the data in a suitable form:

```

file = bzfile("LiverMaleFromLiverFemale3600.csv.bz2");
data = read.csv(file, header = TRUE);
data1 = t(data[, substring(colnames(data), 1, 3)=="F2_"]);
colnames(data1) = data$substanceBXH

```

We now make sure that the genes (column names) in the female and male expression data agree:

```
female2male = match(colnames(datExprFemale), colnames(data1));
table(is.finite(female2male))
datExprMale = data1[, female2male];
all.equal(colnames(datExprFemale), colnames(datExprMale))
```

The last command should result in a TRUE result.

2.b Calculating topological overlap and clustering

Although the network analysis has been performed by Ghazalpour *et al* [1], we repeat here the basic network construction to get a feel for what the network and modules look like. We calculate the pairwise probe dissimilarity based on the Topological Overlap Matrix [3, 4] and cluster genes using average-linkage hierarchical clustering:

```
dissTOMfemale = 1-TOMsimilarityFromExpr(datExprFemale, power = 6, TOMType = "unsigned");
treeFemale = flashClust(as.dist(dissTOMfemale), method = "a");
```

The module colors in `colorsFemale` correspond to branches of this dendrogram:

```
sizeGrWindow(12,9)
plotDendroAndColors(treeFemale, colorsFemale, "Modules", main = "Female gene dendrogram and module colors",
                    dendroLabels = FALSE, addGuide = TRUE);
```

We now repeat the TOM calculation and clustering for the male data:

```
dissTOMmale = 1-TOMsimilarityFromExpr(datExprMale, power = 6, TOMType = "unsigned");
treeMale = flashClust(as.dist(dissTOMmale), method = "a");
```

Lastly, we plot the two clustering trees (dendrograms) together with the female module colors:

```
sizeGrWindow(11, 7)
# Set up appropriate screen sectioning
layout(matrix(c(1:4), 4, 1), heights = rep(c(0.8, 0.2), 2));
# Plot the female dendrogram
plotDendroAndColors(treeFemale, colorsFemale,
                    "Female modules", main = "Female gene dendrogram and module colors",
                    dendroLabels = FALSE, setLayout = FALSE, marAll = c(2,10,3,0), cex.colorLabels = 1.4,
                    cex = 1.4, cex.main = 2, cex.lab = 1.4, cex.axis = 1.4,
                    addGuide = TRUE);
# Plot the male dendrogram with female module colors
plotDendroAndColors(treeMale, colorsFemale,
                    "Female modules", main = "Male gene dendrogram and female module colors",
                    dendroLabels = FALSE, setLayout = FALSE, marAll = c(2,10,3,0), cex.colorLabels = 1.4,
                    cex = 1.4, cex.main = 2, cex.lab = 1.4, cex.axis = 1.4,
                    addGuide = TRUE);
```

The resulting plot is shown in Figure 1. The female dendrogram serves a check that our network construction is truly the same as the one by Ghazalpour *et al* [1]. Since the module colors agree with the branches exactly, we have the same network that Ghazalpour *et al* used. On the other hand, the male dendrogram and the colors do not agree perfectly, but there is still a large amount of overlap (blocks of solid or almost solid colors under the male branches). Thus, we expect that most female modules will show significant preservation in the male data. Before we end, we save the data in R-internal format. This may be useful for subsequent re-runs of this tutorial.

```
save(datExprFemale, datExprMale, colorsFemale,
     file = "BxHLiver-FemaleVsMale-expressionDataAndModuleColors.RData");
```

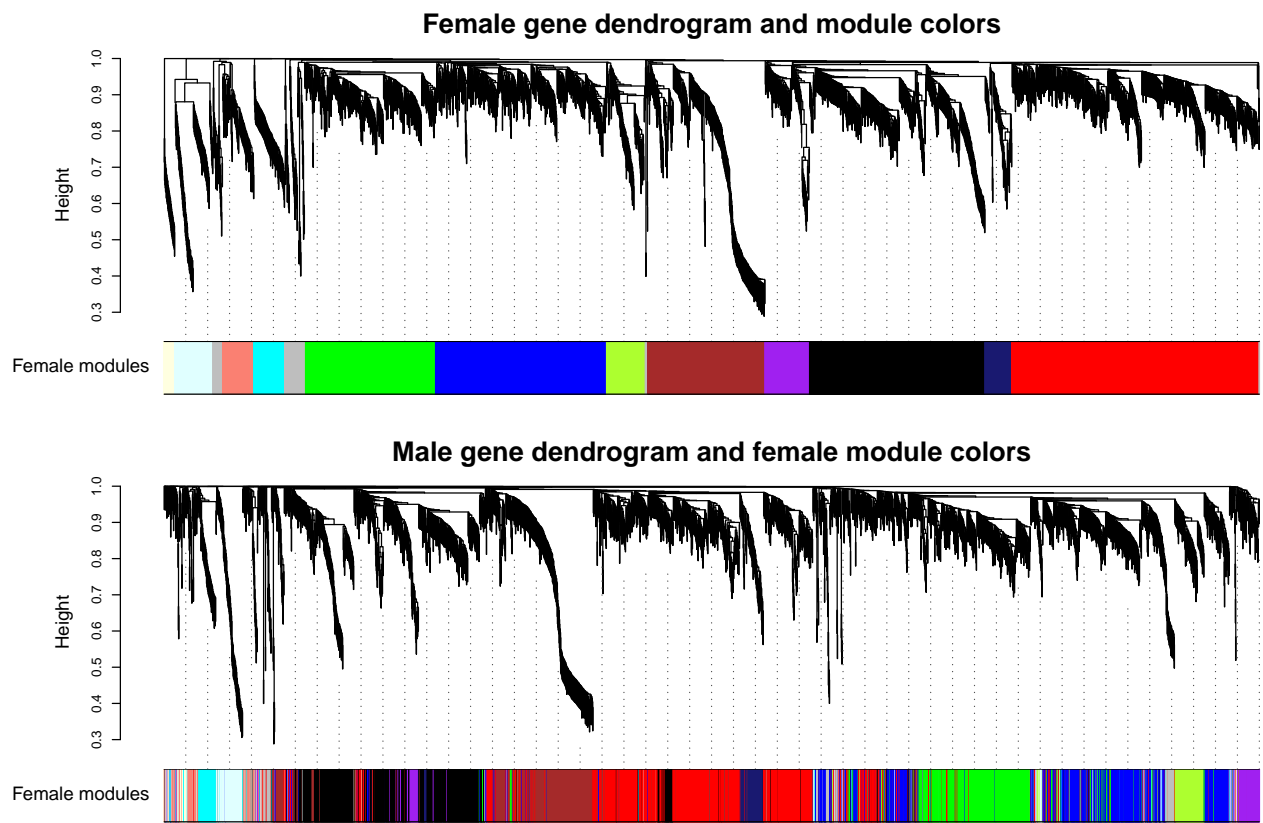


Figure 1: Gene clustering dendrograms in female and male samples, together with the module assignment (colors) obtained from female data by Ghazalpour *et al* [1].

3 Module preservation

In this section we use the function `modulePreservation` to study module preservation of female modules in the male data. The function call is very simple and most of this section will be devoted to the analysis of the results. If this is not the first run of the code in this tutorial, we may load the data saved in the previous section:

```
lnames = load("BxHLiver-FemaleVsMale-expressionDataAndModuleColors.RData");
# Verify the names of loaded variables
lnames
```

3.a Calculation of module preservation

We now set up the multi-set expression data and corresponding module colors:

```
setLabels = c("Female", "Male");
multiExpr = list(Female = list(data = datExpr), Male = list(data = datExprMale));
multiColor = list(Female = colors, Male = maleColors);
nSets = 2
```

Roll the drums, here comes the calculation of module preservation. This calculation takes about 1.5 hours on a fast desktop PC. If this code has been already executed once and the result is saved, it can be loaded from disk, thus saving time.

```
system.time( {
  mp = modulePreservation(multiExpr, multiColor,
    referenceNetworks = c(1:2),
    nPermutations = 200,
    randomSeed = 1,
    verbose = 3)
} );
# Save the results
save(mp, file = "BxHLiverFemaleOnly-modulePreservation.RData");
```

Alternatively, if the data has already been calculated before, load them from disk:

```
load(file = "BxHLiverFemaleOnly-modulePreservation.RData");
```

3.b Analysis and display of module preservation results

We now analyze the data. Isolate the observed statistics and their Z scores:

```
ref = 1
test = 2
statsObs = cbind(mp$quality$observed[[ref]][[test]][, -1], mp$preservation$observed[[ref]][[test]][, -1])
statsZ = cbind(mp$quality$Z[[ref]][[test]][, -1], mp$preservation$Z[[ref]][[test]][, -1]);
```

We look at the main output: the preservation *medianRank* and $Z_{summary}$ statistics.

```
# Compare preservation to quality:
print( cbind(statsObs[, c("medianRank.pres", "medianRank.qual")],
  signif(statsZ[, c("Zsummary.pres", "Zsummary.qual")], 2)) )
```

The numbers are nice, but (to paraphrase a saying) a picture is worth a thousand numbers. We plot the preservation *medianRank* and $Z_{summary}$ for the female modules as a function of module size. The plotting code may seem a bit involved, but it is worth going through.

```

# Module labels and module sizes are also contained in the results
modColors = rownames(mp$preservation$observed[[ref]][[test]])
moduleSizes = mp$preservation$Z[[ref]][[test]][, 1];
# leave grey and gold modules out
plotMods = !(modColors %in% c("grey", "gold"));
# Text labels for points
text = modColors[plotMods];
# Auxiliary convenience variable
plotData = cbind(mp$preservation$observed[[ref]][[test]][, 2], mp$preservation$Z[[ref]][[test]][, 2])
# Main titles for the plot
mains = c("Preservation Median rank", "Preservation Zsummary");
# Start the plot
sizeGrWindow(10, 5);
#pdf(fi="Plots/BxHLiverFemaleOnly-modulePreservation-Zsummary-medianRank.pdf", wi=10, h=5)
par(mfrow = c(1,2))
par(mar = c(4.5,4.5,2.5,1))
for (p in 1:2)
{
  min = min(plotData[, p], na.rm = TRUE);
  max = max(plotData[, p], na.rm = TRUE);
  # Adjust plotting ranges appropriately
  if (p==2)
  {
    if (min > -max/10) min = -max/10
    ylim = c(min - 0.1 * (max-min), max + 0.1 * (max-min))
  } else
    ylim = c(max + 0.1 * (max-min), min - 0.1 * (max-min))
  plot(moduleSizes[plotMods], plotData[plotMods, p], col = 1, bg = modColors[plotMods], pch = 21,
        main = mains[p],
        cex = 2.4,
        ylab = mains[p], xlab = "Module size", log = "x",
        ylim = ylim,
        xlim = c(10, 2000), cex.lab = 1.2, cex.axis = 1.2, cex.main = 1.4)
  labelPoints(moduleSizes[plotMods], plotData[plotMods, p], text, cex = 1, offs = 0.08);
  # For Zsummary, add threshold lines
  if (p==2)
  {
    abline(h=0)
    abline(h=2, col = "blue", lty = 2)
    abline(h=10, col = "darkgreen", lty = 2)
  }
}
# If plotting into a file, close it
dev.off();

```

The resulting plot is shown in Figure 2. We find that most modules show strong evidence of preservation. The exceptions are the lightyellow and salmon modules whose $Z_{summary}$ statistics are 2.0 (borderline between no preservation and very weak preservation) and 7.6 (moderate preservation), respectively. The lightyellow and salmon modules also have to lowest median rank statistic, meaning their observed preservation statistics tend to be lowest among all modules. In Supplementary Text S5 published alongside the main module preservation paper we analyze the weak to non-existent preservation of these two modules and show that their presence in the female data is mainly due to an outlier sample. Thus, in this example module preservation statistics flag modules that are driven by technical artefacts.

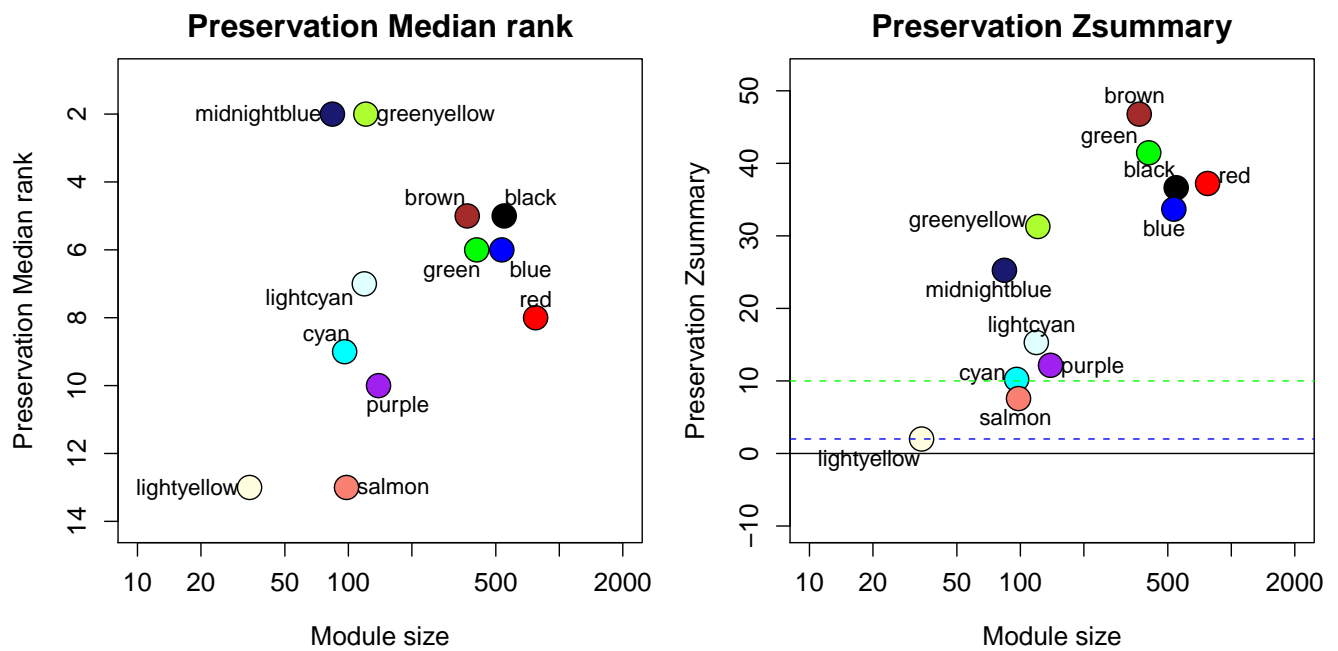


Figure 2: The $medianRank$ and $Z_{summary}$ statistics of module preservation of female modules in male modules (y -axis) vs. module size (x -axis).

We now plot the density and connectivity statistics all in one plot. We include the module quality measures for comparison:

```
# Re-initialize module color labels and sizes
modColors = rownames(statsZ)
moduleSizes = mp$quality$Z[[ref]][[test]][, 1];
# Exclude improper modules
plotMods = !(modColors %in% c("grey", "gold"));
# Create numeric labels for each module
labs = match(modColors[plotMods], standardColors(50));
# Start the plot: open a suitably sized graphical window and set sectioning and margins. Alternatively,
# plot into a pdf file.
sizeGrWindow(10, 9);
#pdf(file="Plots/PreservationZStatistics.pdf", w=10, h=9)
par(mfrow = c(4,4))
par(mar = c(3,3,2,1))
par(mgp = c(1.6, 0.4, 0));
for (s in 1:ncol(statsZ))
{
  min = min(statsZ[plotMods, s], na.rm = TRUE);
  max = max(statsZ[plotMods, s], na.rm = TRUE);
  if (min > -max/12) min = -max/12
  plot(moduleSizes[plotMods], statsZ[plotMods, s], col = 1, bg = modColors[plotMods], pch = 21,
       main = colnames(statsZ)[s],
       cex = 2.2,
       ylab = colnames(statsZ)[s], xlab = "Module size", log = "x",
       ylim = c(min - 0.1 * (max-min), max + 0.1 * (max-min)),
       xlim = c(30, 800),
       cex.lab = 1.2, cex.axis = 1.2)
  labelPoints(moduleSizes[plotMods], statsZ[plotMods, s], labs, cex = 1, offs = 0.06);
  #text(moduleSizes[-1], statsZ[-c(1:2), s], labels = letter[-c(1:2)], col = "black"); #modColors[-2]);
  abline(h=0)
  abline(h=2, col = "blue", lty = 2)
  abline(h=10, col = "darkgreen", lty = 2)
}
# If plotting into a file, close it, otherwise it is unreadable.
dev.off();
```

The result is shown in Figure 3. The translation table between module colors and the numeric labels can be obtained for example as follows:

```
data.frame(color = modColors[plotMods], label = labs)
```

The result is

```
> data.frame(color = modColors[plotMods], label = labs)
  color label
1  black    7
2  blue    2
3  brown    3
4  cyan   14
5  green    5
6 greenyellow 11
7  lightcyan 16
8  lightyellow 19
9  midnightblue 15
10 purple    10
11  red      6
12  salmon   13
```

Of note is the fact the the quality (for example, $Z_{\text{summary.qual}}$) of the lightyellow and salmon modules is quite high, definitely on par with other modules, while their preservation is much lower. The all quality statistics point to bradly the same conclusion, but they measure different aspects of quality and hence are all somewhat different from each other. Hence the usefulness of the Z_{summary} statistic. The same can be said about the preservation statistics - they are broadly similar but in details they differ, and summary is very useful to form an overall conclusion. As a programming/graphics side note, some of the plots illustrate graphically the limitations of the function `labelPoints`. There are no miracles, not even in WGCNA.

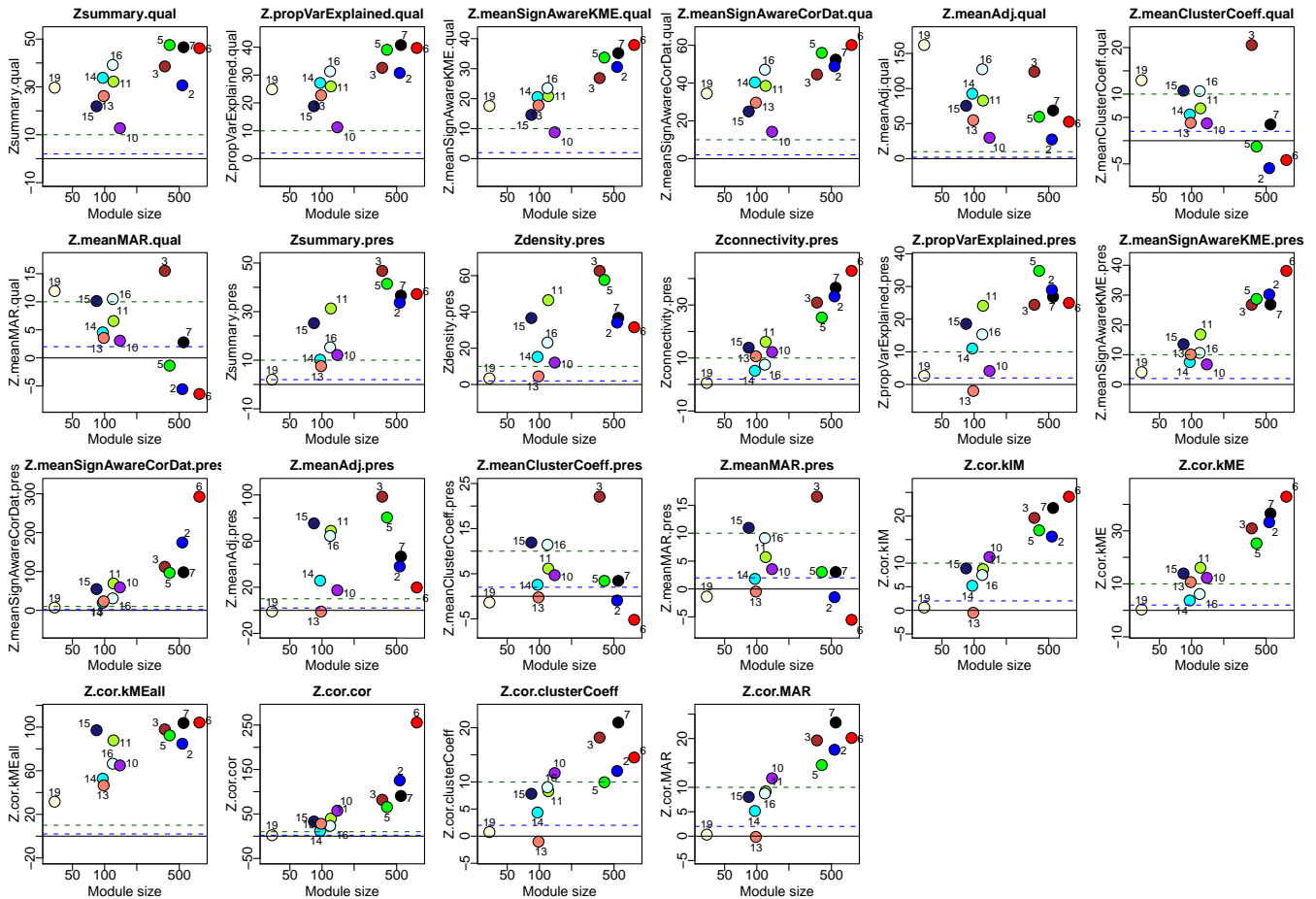


Figure 3: Statistics of module quality (first row; names end with suffix `.qual`) and module preservation (second and third row) of female modules. Modules are labeled by their color and by numeric labels (see text).

4 Alternative methods: clusterRepro

4.a Calculation of preservation statistics

In this section we provide the code for reproducing the results of `clusterRepro` [2] in this application. Running this code requires the `clusterRepro` package, available from CRAN. In the first block we apply the function `clusterRepro` and save the results.

```
nSets = 2
doClusterRepro = TRUE;
if (doClusterRepro)
{
  cr = list();
  library(clusterRepro);
  set.seed(20);
  for (ref in 1:nSets)
  {
    eigengenes = multiSetMEs(multiExpr, universalColors = multiColor[[ref]]);
    cr[[ref]] = list();
    for (set in 1:nSets)
    {
      printFlush(paste("Working on reference set", ref, " and test set", set));
      rownames(eigengenes[[set]]$data) = rownames(multiExpr[[set]]$data);
      print(system.time({ cr[[ref]][[set]] = clusterRepro(Centroids = as.matrix(eigengenes[[set]]$data),
                                                         New.data = as.matrix(multiExpr[[set]]$data),
                                                         Number.of.permutations = 2000)}))
    }
  }
}
# Save the results so they can be re-used if necessary.
save(cr, file = "BxHLiverFemaleOnly-cr.RData");
```

Alternatively, if the calculation has been performed already, simply load the results and save several hours:

```
load(file = "BxHLiverFemaleOnly-cr.RData");
```

4.b Analysis of results

Here we plot the clusterRepro results.

```
# Find the maximum and minimum
max = 0; min = 1;
for (ref in c(1:nSets))
{
  p = 1;
  for (test in 1:nSets)
  {
    stats = cr[[ref]][[test]]$Actual.IGP;
    max = max(max, stats, na.rm = TRUE);
    min = min(min, stats, na.rm = TRUE);
  }
}
max = max + 0.05*(max-min);
min = min - 0.03*(max-min);
refNames = setLabels
# Plot the results on a uniform scale
sizeGrWindow(10,9);
#pdf(file="Plots/BxHLiverFemaleOnly-clusterRepro-IGP.pdf", w=10, h=9)
par(mfrow = c(2,2))
par(mar = c(3.5, 3.5, 3, 0.4))
par(mgp = c(2.0, 0.5, 0));
for (ref in 1:nSets)
{
  p = 1;
  for (test in 1:nSets)
  {
    stats = cr[[ref]][[test]]$Actual.IGP;
    moduleSizes = table(multiColor[[ref]])
    labelsX = names(moduleSizes);
    #modNumbers = match(labelsX, standardColors(20))
    plotMods = !(labelsX %in% c("grey", "gold"));
    labelsX = labelsX[plotMods]

    xmin = min(moduleSizes[plotMods]);
    xmax = max(moduleSizes[plotMods]);
    xlim = c(xmin * (xmin/xmax)^.20, xmax * (xmax/xmin)^0.15);
    plot(moduleSizes[plotMods], stats, bg = labelsX, pch = 21,
         main = spaste(LETTERS[ref], p, ". Modules: ", refNames[ref], "\n Test data: ", setLabels[test]),
         cex = 2, cex.axis = 1.2, cex.lab = 1.2,
         ylab = "Actual IGP", xlab = "Module size", log = "x", ylim = c(min, max), xlim = xlim)
    #abline(h = stats[colors=="orange",s], col = "grey30", lty = 2)
    abline(h=0)
    labelPoints(moduleSizes[plotMods], stats, labels = labelsX, offs = 0.070,
               jiggle = 0, cex = 1)
    p = p+ 1;
  }
}
# If plotting into a file, close it.
dev.off()
```

References

- [1] A. Ghazalpour, S. Doss, B. Zhang, C. Plaisier, S. Wang, E.E. Schadt, A. Thomas, T.A. Drake, A.J. Lusis, and S. Horvath. Integrating genetics and network analysis to characterize genes related to mouse weight. *PloS*

Genetics, 2(2):8, 2006.

- [2] Amy V. Kapp and Robert Tibshirani. Are clusters found in one dataset present in another dataset? *Biostat*, 8(1):9–31, 2007.
- [3] E Ravasz, A L Somera, D A Mongru, Z N Oltvai, and A L Barabasi. Hierarchical organization of modularity in metabolic networks. *Science*, 297(5586):1551–5, 2002.
- [4] B Zhang and S. Horvath. General framework for weighted gene coexpression analysis. *Statistical Applications in Genetics and Molecular Biology.*, 4(17), 2005.