

Tutorial for the WGCNA package in R: Consensus network, GO analysis, cell type enrichment, module network plots

Elisabeth Rebboah

December 2018

Part 1: Data Input, cleaning, and pre-processing

Set working directory and load libraries necessary for WGCNA and other steps. `enableWGCNAThreads` allows parallel execution with multiple cores.

```
1 setwd("~/WGCNA")
2 getwd()
3 library(WGCNA)
4 library(flashClust)
5 library(gplots)
6 library(cluster)
7 library(igraph); #for part 8
8 library(RColorBrewer); #for part 8
9 options(stringsAsFactors=FALSE)
10 #Only when not on server, comment out:
11 enableWGCNAThreads()
```

Expression matrix `datExpr` needs to have genes in columns and samples in rows, i.e. final `datExpr` might be something like 20,000 genes by 200 patients. Also must have Ensembl annotation for list of genes. Expression values are from microarray or RNA-seq, already pre-processed and normalized. Target matrices contain samples in rows and traits in columns, such as diagnosis, age, sex, etc.

For a consensus network, the genes between datasets must overlap. For example, for three datasets, `datExpr.Ref.DS`, `datExpr.Ref.AD_1`, and `datExpr.Ref.AD_2` :

```
1 gnS=intersect(colnames(datExpr.Ref.DS),intersect(colnames(datExpr.Ref.AD_1),colnames(datExpr
  .Ref.AD_2)))
2 datExpr.Ref.DS=datExpr.Ref.DS[,match(gnS,colnames(datExpr.Ref.DS))]
3 datExpr.Ref.AD_1 =datExpr.Ref.AD_1 [,match(gnS,colnames(datExpr.Ref.AD_1))]
4 datExpr.Ref.AD_2 =datExpr.Ref.AD_2 [,match(gnS,colnames(datExpr.Ref.AD_2))]
```

From two or more expression matrices, combine into a `multiExpr` dataframe, a vector of lists, in each list there must be a component named "data" whose content is the individual expression matrices.

```
1 #Make multiExpr list of 3 data frames
2 nSets=3;
3 setLabels=c("Downs syndrome samples", "Alzheimer's samples 1", "Alzheimer's samples 2")
4 #Form multi-set expression data
5 multiExpr = vector(mode="list", length=nSets)
6
7 #Remove rownames and column names from datExpr dfs, re-assign in multiExpr
8 genes = colnames(datExpr.Ref.DS);
9 samples = rownames(datExpr.Ref.DS);
10 rownames(datExpr.Ref.DS)=NULL;
11 colnames(datExpr.Ref.DS)=NULL;
12 multiExpr[[1]] = list(data = datExpr.Ref.DS);
13 colnames(multiExpr[[1]]$data) = genes;
14 rownames(multiExpr[[1]]$data) = samples;
15
16 genes = colnames(datExpr.Ref.AD_1);
```

```

17 samples = rownames(datExpr.Ref.AD_1);
18 rownames(datExpr.Ref.AD_1)=NULL;
19 colnames(datExpr.Ref.AD_1)=NULL;
20 multiExpr[[2]] = list(data = datExpr.Ref.AD_1);
21 names(multiExpr[[2]]$data) = genes;
22 rownames(multiExpr[[2]]$data) = samples;
23
24 genes = colnames(datExpr.Ref.AD_2);
25 samples = rownames(datExpr.Ref.AD_2);
26 rownames(datExpr.Ref.AD_2)=NULL;
27 colnames(datExpr.Ref.AD_2)=NULL;
28 multiExpr[[3]] = list(data = datExpr.Ref.AD_2);
29 colnames(multiExpr[[3]]$data) = genes;
30 rownames(multiExpr[[3]]$data) = samples;
31
32 exprSize = checkSets(multiExpr)
33 exprSize #Check data is in correct format; checkSets function checks whether given sets have
           the correct format and retrieves dimensions.
34
35 ##$nGenes
36 #[1] 11810
37 ##$nSamples
38 #[1] 110 104 465

```

Also make list of lists for all traits in multiMeta object, and check that all genes and samples don't have too many missing values (gsg\$allOK should return TRUE if all genes and samples are okay).

```

1 #Make multiMeta of traits for all sets - list of lists
2 multiMeta=list(aging = list(data=targets.Ref.DS), AD1 = list(data=targets.Ref.AD_1), AD2 =
   list(data=targets.Ref.AD_2))
3
4 #Check that all genes and samples have sufficiently low numbers of missing values
5 gsg = goodSamplesGenesMS(multiExpr, verbose = 3);
6 gsg$allOK

```

Save input data in rda file.

```

1 #Save everything
2 nGenes=exprSize$nGenes;
3 nSamples=exprSize$nSamples;
4 save(multiExpr, multiMeta, nGenes, nSamples, setLabels,
5       targets.Ref.AD_1, targets.Ref.AD_2, targets.Ref.DS,
6       datExpr.Ref.AD_1, datExpr.Ref.AD_2, datExpr.Ref.DS, file = "Consensus_dataInput_112418.
   rda");

```

Part 2: Choose soft thresholding power

Next step is to set a power (β) for a scale-free topology network, which has to be done by the user based on plots of scale-free topology fit and network connectivity. It is best to choose a power that crosses some threshold in the scale-free topology fit plot (such as > 0.8) but does not lose too much connectivity.

Input a set of soft thresholding powers to test and call pickSoftThreshold on data for the type of network you want to build with the type of correlation function to be used in adjacency calculation; in this case it is signed with a bicor correlation function.

```

1 # Choose a set of soft-thresholding powers

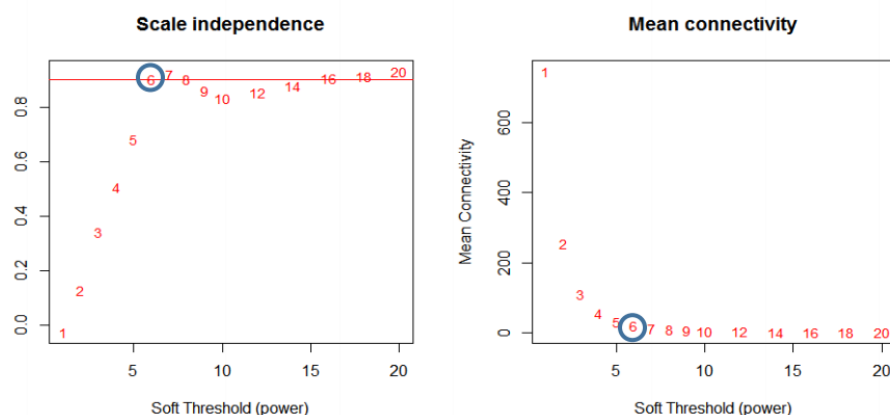
```

```

2 powers = c(seq(1,10,by=1), seq(12,30, by=2));
3 # Initialize a list to hold the results of scale-free analysis
4 powerTables = vector(mode = "list", length = nSets);
5 # Call the network topology analysis function for each set in turn
6 for (set in 1:nSets)
7   powerTables[[set]] = list(data = pickSoftThreshold(multiExpr[[set]]$data, powerVector=
8     powers,
                                verbose = 5, networkType="signed",
                                corFnc="bicor")[[2]]);

```

Example plots, where 6 is chosen since it crosses an arbitrary threshold (red line) but doesn't lose all the connectivity.



Part 3: Network Construction

From plots generated in part 2, set the soft thresholding power. Next, run `blockwiseConsensusModules` to generate network across datasets. This function has many parameters than can be adjusted. Again, in this case `networkType` is signed. Consult [function documentation](#) for more information on parameters.

```

1 #Set soft thresholding power to number based on plots
2 softpower=16;
3 # Auto network
4 load(file = "Consensus_dataInput_112418.rda");
5 net=blockwiseConsensusModules(multiExpr, blocks = NULL,
6                               maxBlockSize = 30000, ## This should be set to a smaller size
6                               if the user has limited RAM
7                               randomSeed = 12345,
8                               corType = "pearson", ## no use for bicor
9                               power = softpower,
10                              consensusQuantile = 0.2,
11                              networkType = "signed",
12                              TOMType = "unsigned",
13                              TOMDenom = "min",
14                              scaleTOMs = TRUE, scaleQuantile = 0,
15                              sampleForScaling = TRUE, sampleForScalingFactor = 1000,
16                              useDiskCache = TRUE, chunkSize = NULL,
17                              deepSplit = 2,
18                              detectCutHeight = 0.99999999, minModuleSize = 20,
19                              mergeCutHeight = 0.2,
20                              saveConsensusTOMs = TRUE,
21                              consensusTOMFilePattern = "ConsensusTOM-block.%.b.rda")
22 save(list=ls(),file="consensus_112418.rda")

```

Get module eigengenes (MEs) from “net” dataframe. To cluster the network into a dendrogram, use hclust on the dissimilarity matrix, 1-TOM.

```
1 consMEs = net$multiMEs;
2 moduleColors = net$colors;
3 table(moduleColors)
4
5 load("ConsensusTOM-block.1.rda") # consensus TOM
6 consTree= hclust(1-consTomDS,method="average");
7 save(list=ls(),file="consensus_112418.rda")
```

To identify modules, cutreeHybrid detects clusters in the dendrogram based on a set of parameters, some of which are looped through so that the user can determine which set of modules makes the most sense. As a personal preference, I also plotted how each gene correlates to all traits to see any obvious trends between modules and traits. Trait data has to be converted to numbers, then each gene is correlated to the trait. For example, diagnosis is 0 or 1, control or diseased, and genes could be positively or negatively correlated to having the disease.

Examples of converting traits to binaries, followed by calculating biweight midcorrelation between each gene across all samples and each trait value as a number across all samples, so between two numeric vectors. Bicor is median-based rather than mean-based, therefore it is more robust to outliers than Pearson correlation.

```
1 Diagnosis=factor(as.character(targets.Ref.AD_2$Disease))
2 Diagnosis<-relevel(Diagnosis,'Control')
3 Diagnosis=as.numeric(Diagnosis)
4 Age=as.numeric(targets.Ref.AD_2$Age)
5 GEO=as.numeric(factor(targets.Ref.AD_2$geo_accession))
6 Gender=as.numeric(factor(targets.Ref.AD_2$Sex))
7
8 geneSigsAD2=matrix(NA,nrow=4,ncol=ncol(datExpr.Ref.AD_2))
9 for(i in 1:ncol(geneSigsAD2)) {
10   exprvec=as.numeric(datExpr.Ref.AD_2[,i])
11   ager=bicor(Age,exprvec,use="pairwise.complete.obs")
12   sexr=bicor(exprvec, Gender,use="pairwise.complete.obs")
13   conditionr=bicor(exprvec, Diagnosis,use="pairwise.complete.obs")
14   geor=bicor(GEO,exprvec,use="pairwise.complete.obs")
15   geneSigsAD2[,i]=c(ager, sexr, conditionr, geor)
16   cat('Done for gene...',i,'\n')
17 }
```

Finding different sets of modules based on different values for parameters mms, dthresh, and ds using cutreeHybrid, which are then plotted:

```
1 # Calculate modules for each set of parameters
2 load("consensus_112418.rda")
3 mColorh <- mLabelh <- colorLabels <- NULL
4 for (minModSize in c(40,100,160)) {
5   for (dthresh in c(0.1,0.2,0.25)) {
6     for (ds in c(2,4)) {
7       print("Trying parameters:")
8       print(c(minModSize,dthresh,ds))
9       tree = cutreeHybrid(dendro = consTree, pamStage=FALSE,
10                          minClusterSize = minModSize, cutHeight = 0.99999999,
11                          deepSplit = ds, distM = as.matrix(1-consTomDS))
12
13       merged <- mergeCloseModules(exprData = multiExpr,colors = tree$labels,
14                                  cutHeight = dthresh)
15       mColorh <- cbind(mColorh,labels2colors(merged$colors))
16       mLabelh <- c(mLabelh,paste("DS=",ds," mms=\n",minModSize," dcor=",dthresh))
17     }
18   }
```

```
19 }
```

Once mms, ds, and dthresh are chosen based on best fit of clusters, run cutreeHybrid again and merge chosen modules so that genes are associated with a given module. These genes are summarized by a “module eigengene” that captures the most variance in the module expression data, can be thought of as the first principal component of the module.

```
1 mms=160
2 ds=4
3 dthresh=0.1
4
5 tree = cutreeHybrid(dendro = consTree, pamStage=FALSE,
6                     minClusterSize = mms, cutHeight = 0.99999999,
7                     deepSplit = ds, distM = as.matrix(1-consTomDS))
8
9 merged <- mergeCloseModules(exprData = multiExpr, colors = tree$labels,
10                             cutHeight = dthresh)
11
12 # Eigengenes of the new merged modules:
13 MES_DS = merged$newMES[[1]]$data;
14 MES_AD1 = merged$newMES[[2]]$data;
15 MES_AD2 = merged$newMES[[3]]$data;
```

Next, connectivity of each gene across datasets to the module eigengene is determined by the consensusKME function, then consensus kMEs are determined.

```
1 #Calculate kMEs and p value, Z score of kMEs for each gene (connectivity of each gene to a
   module eigengene) across datasets
2 consKME1=consensusKME(multiExpr=multiExpr, moduleColor.cons,
3                       multiEigengenes = NULL,
4                       consensusQuantile = 0.2,
5                       signed = TRUE)
6 #Consensus kMEs:
7 consensus.KMEs=consKME1[, regexpr('consensus.kME', names(consKME1))>0]
```

Part 4: Get annotated gene lists and their associated kMEs

Convert Ensemble gene IDs to gene symbols, output Excel file of list of genes and their associated kMEs. This can be used as an input file for further analyses.

```
1 load("Consensus_MES_112418.rda")
2 ensembl=read.csv("/home/vivek/FTD_Seeley/Analysis_Nov2017/ENSG85_Human.csv.gz")
3 consensus.KMEs$Ensembl.Gene.ID=paste(rownames(consensus.KMEs))
4
5 merged=merge(consensus.KMEs, ensembl, by.x="Ensembl.Gene.ID", by.y="Ensembl.Gene.ID", all.x=T)
6 ind=match(consensus.KMEs$Ensembl.Gene.ID, merged$Ensembl.Gene.ID)
7 merged1=merged[ind,]
8 consensus.KMEs.annot=merged1
9
10 geneInfo.cons=as.data.frame(cbind(consensus.KMEs.annot$Ensembl.Gene.ID, consensus.KMEs.annot$
   Associated.Gene.Name,
11                                 moduleColor.cons, consensus.KMEs))
12 geneInfo.cons=geneInfo.cons[, -ncol(geneInfo.cons)] # check if last column is Ensembl gene id
13
14 colnames(geneInfo.cons)[1]= "Ensembl.Gene.ID"
15 colnames(geneInfo.cons)[2]= "GeneSymbol"
16 colnames(geneInfo.cons)[3]= "Initially.Assigned.Module.Color"
17
18 write.csv(geneInfo.cons, 'geneInfo.cons.DSAD_112418.csv')
19 save(list=ls(), file='geneInfo.cons.112418.rda')
```

Part 5: Module-trait relationships

Use correlations to make a heatmap of each module eigengene across traits of interest, display significant p-values and correlations in heatmap boxes. These values can be saved to Excel sheets later.

```
1 pdf(paste('NetworkPlot_AD1_consensus_112418.pdf'),width=16,height=30)
2 par( mar = c(8, 12, 3, 3) );
3 labeledHeatmap( Matrix = moduleTraitCor_AD1,
4                 xLabels = colnames(factors1_AD1),
5                 yLabels = rownames(moduleTraitPvalue_AD1),
6                 ySymbols = rownames(moduleTraitPvalue_AD1),
7                 colorLabels = FALSE,
8                 colors = blueWhiteRed(50),
9                 textMatrix = textMatrix1,
10                setStdMargins = FALSE,
11                cex.text = 1.5,
12                zlim = c(-1, 1),
13                cex.lab.x = 1.2,
14                main = paste("Module-trait relationships")
15 );
```

This code plots eigengene network in the form of a dendrogram and heatmap. This shows how related the eigengenes are to each other. Each row and column corresponds to an eigengene (labeled by consensus module color). Red indicates high adjacency (positive correlation) and white no adjacency (no correlation).

```
1 #Plot eigengene heatmap
2 par(cex = 1.0)
3 plotEigengeneNetworks(MEs_DS, "Eigengene Network", marHeatmap = c(3,4,2,2), marDendro = c
  (0,4,1,2),cex.adjacency = 0.3,plotDendrograms = TRUE, xLabelsAngle = 90,heatmapColors=
  blueWhiteRed(100)[51:100])
```

Also make boxplots and scatterplots of trait enrichment for each module. For the binary factor traits (male/female, diseased/control) boxplots indicate enrichment. For the numeric traits such as age, a scatterplot shows if there is enrichment for a specific age and plots a trendline.

```
1 #Plot boxplots, scatterplots
2 toplot=t(MEs_AD2)
3 cols=substring(colnames(MEs_AD2),3,20)
4 par(mfrow=c(4,4))
5 par(mar=c(5,6,4,2))
6 for (i in 1:nrow(toplot)) {
7   boxplot(toplot[i,]~factor(as.vector(as.factor(targets.Ref.AD_2$Disease))),c('Control','AD')
8   ),col=cols[i],ylab="ME",main=rownames(toplot)[i],xlab=NULL,las=2)
9   verboseScatterplot(x=as.numeric(targets.Ref.AD_2$Age),y=toplot[i,],xlab="Age",ylab="ME",
10  abline=TRUE,cex.axis=1,cex.lab=1,cex=1,col=cols[i],pch=19)
11   boxplot(toplot[i,]~factor(targets.Ref.AD_2$Sex),col=cols[i],ylab="ME",main=rownames(toplot
12  ) [i],xlab=NULL,las=2)
13 }
14 dev.off()
```

Part 6: GO Analysis

Finally, find GO enrichment terms for each module. This is the most interesting part of the analysis since you can start to make biological connections to the modules.

First, create directories for inputs and outputs.

```
1 dir.create("./geneInfo")
2 dir.create("./geneInfo/background/")
3 dir.create("./geneInfo/input/")
4 dir.create("./geneInfo/output/")
```

After adding some reference files into these directors, analysis is run as a nohupped bash shell script that uses inputs and outputs from these new directories.

```
1 # Run GO elite as nohupped shell script:
2 codedir <- "/home/vivek/bin/GO-Elite_v.1.2.5-Py"
3 pathname <- "/home/erebboah/WGCNA/geneInfo"
4 nperm=10000
5 system(paste("nohup python ",codedir,"/GO-Elite.py --species Hs --mod Ensembl --permutations
6     ",
7     nperm," --method \"z-score\" --zscore 1.96 --pval 0.01 --num 5 --input ",
8     pathname,
9     "/input --denom ",pathname,"/background --output ",pathname,"/output &",sep="")
10 )
```

In this way, code can be run in the background for several hours. It should only take a few hours or so for 16 modules (depending on available CPU). Plot the results:

```
1 pdf("GOElite_plot_Modules_112718.pdf",height=8,width=12)
2 for(i in 1:length(uniquemodcolors)){
3   thismod= uniquemodcolors[i]
4   tmp=read.csv(file=paste(pathname,"/ORA_pruned/",thismod,"_Module-GO_z-score_elite.txt",sep
5     =""),sep="\t")
6   tmp=subset(tmp,Ontology.Type!='cellular_component')
7   tmp=tmp[,c(2,9)] ## Select GO-terms and Z-score
8   tmp=tmp[order(tmp$Z.Score,decreasing=T),] #
9   if (nrow(tmp)<10){
10     tmp1=tmp ## Take top 10 Z-score
11     tmp1 = tmp1[order(tmp1$Z.Score),] ##Re-arrange by increasing Z-score
12     par(mar=c(5,40,5,2))
13     barplot(tmp1$Z.Score,horiz=T,col="blue",names.arg= tmp1$Ontology.Name,cex.names=1.2,las
14       =1,main=paste("Gene Ontology Plot of",thismod,"Module"),xlab="Z-Score")
15     abline(v=2,col="red")
16   } else {
17     tmp1=tmp[c(1:10),] ## Take top 10 Z-score
18     tmp1 = tmp1[order(tmp1$Z.Score),] ##Re-arrange by increasing Z-score
19     par(mar=c(5,40,5,2))
20     barplot(tmp1$Z.Score,horiz=T,col="blue",names.arg= tmp1$Ontology.Name,cex.names=1.2,las
21       =1,main=paste("Gene Ontology Plot of",thismod,"Module"),xlab="Z-Score")
22     abline(v=2,col="red")
23   }
24   cat('Done ...',thismod,'\n')
25 }
26 dev.off()
```

Part 7: Cell type enrichment

Make heatmap plus p-values and correlation scores across modules for these cell types: neurons, microglia, myelinating oligodendrocytes, astrocytes, and endothelial cells. Also useful for making biological insights into modules.

```

1 pdf("CellTypeEnrich_WGCNAMods_110618.pdf", width=6,height=10)
2 labeledHeatmap(Matrix=dispMat,
3                 yLabels=rownames(dispMat),
4                 yColorLabels=TRUE,
5                 xLabels= colnames(dispMat),
6                 colors=blueWhiteRed(40),
7                 textMatrix = textMatrix1,
8                 cex.lab.x=1.0,
9                 zlim=c(-0.1,3),
10                main="Cell-type enrichment Heatmap")
11 dev.off()

```

Part 8: TOM network plot

Will make node and edge plots with hub genes in the center surrounded by all other genes in each module. Also somewhat useful for biological insights, hub genes can serve as sanity checks that the WGCNA actually finds relevant genes to the diseases being studied.

```

1 pdf("ModuleNetworks.pdf",height=9,width=10);
2 for (mod in uniquemodcolors) {
3   numgenesingraph = 100;
4   numconnections2keep = 1500;
5   cat('module:',mod,'\n');
6   geneInfo.cons=geneInfo.cons[geneInfo.cons$GeneSymbol!="NA",]
7   colind = which(colnames(geneInfo.cons)==paste("consensus.kME",mod, sep=""));
8   rowind = which(geneInfo.cons[,3]==mod);
9   cat(' ',length(rowind),'probes in module\n');
10  submatrix = geneInfo.cons[rowind,];
11  orderind = order(submatrix[,colind],decreasing=TRUE);
12  if (length(rowind) < numgenesingraph) {
13    numgenesingraph = length(rowind);
14    numconnections2keep = numgenesingraph * (numgenesingraph - 1);
15  }
16  cat('Making network graphs, using top',numgenesingraph,'probes and',numconnections2keep,'
17    connections of TOM\n');
18  submatrix = submatrix[orderind[1:numgenesingraph],];
19  #Identify the columns in the TOM that correspond to these hub probes
20  matchind = match(submatrix$Ensembl.Gene.ID,colnames(multiExpr));
21  reducedTOM = TOM.matrix[matchind,matchind];
22
23  orderind = order(reducedTOM,decreasing=TRUE);
24  connections2keep = orderind[1:numconnections2keep];
25  reducedTOM = matrix(0,nrow(reducedTOM),ncol(reducedTOM));
26  reducedTOM[connections2keep] = 1;
27
28  g0 <- graph.adjacency(as.matrix(reducedTOM[1:10,1:10]),mode="undirected",weighted=TRUE,
29    diag=FALSE)
30  layoutMata <- layout.circle(g0)
31
32  g0 <- graph.adjacency(as.matrix(reducedTOM[11:50,11:50]),mode="undirected",weighted=TRUE,
33    diag=FALSE)
34  layoutMatb <- layout.circle(g0)
35
36  g0 <- graph.adjacency(as.matrix(reducedTOM[51:ncol(reducedTOM),51:ncol(reducedTOM)]),mode=
37    "undirected",weighted=TRUE,diag=FALSE)
38  layoutMatc <- layout.circle(g0)
39
40  g1 <- graph.adjacency(as.matrix(reducedTOM),mode="undirected",weighted=TRUE,diag=FALSE)
41  layoutMat <- rbind(layoutMata*0.25,layoutMatb*0.8, layoutMatc)
42
43  plot(g1,edge.color="grey",vertex.color=mod,vertex.label=as.character(submatrix$GeneSymbol)
44    ,vertex.label.cex=0.7,vertex.label.dist=0.45,vertex.label.degree=-pi/4,vertex.label.

```



```
39     color="black", layout= layoutMat, vertex.size=submatrix[,colind]^2*8, main=paste(mod, "  
40     module"))  
41 dev.off();
```