Stat4DS / Homework 02 Luca Scofano / Davide Zingaro 19 dicembre 2019 Exercise: Connect your brain In this task, our goal is to identify and evaluate the dependencies that exist between the various cortical regions. Importing libraries... library(igraph) library(ggraph) library(boot) library(SIN) Quick analysis Let's start with a quick analysis of the dataset mts: load("C:\\Users\\Luca\\Desktop\\WIP\\Homework 2 SDS\\.RData") data <- mts # Load Data hist(data) # Visualize Data Histogram of data 1500 500 200 400 800 1000 600 data # Variables n.features <- ncol(data)</pre> n <- nrow(data) # number of observations</pre> data.cor <- cor(data) # matrix of correlations of the sample (Pearson's)</pre> feature.names <- colnames(data)</pre> Initialize edge matrix # Initialize edge matrix filled with NA edge <- matrix(NA, n.features, n.features)</pre> colnames(edge) <- feature.names</pre> rownames(edge) <- feature.names</pre> Non Parametric Bootstrap procedure and ecdf With this data, let's consider the Pearson correlation, and implement the bootstrap procedure described according to "simultaneous bootstrapped ci's for a generic association measure  $\rho$ ". set.seed(123) library(boot) B <- 1000 delta.beta <- function(data, indices) {</pre> d <- data[indices,] # allows boot to select sample</pre> pearson <- cor(d) # creates a correlation matrix for each boostrap</pre> data.cor <- cor(data) # matrix of correlations of the sample</pre> delta <- sqrt(nrow(data))\*max(abs(pearson-data.cor))</pre> return(delta) bootdelta <- boot(data = data, statistic = delta.beta, R = B)</pre> Now we have the list of delta.b and we can compute the ecdf delta.b <- bootdelta\$t f.hat <- ecdf(delta.b)</pre> # Plot the ecdf plot(f.hat, ylab="Fn(x)", verticals = FALSE, lwd = 3)ecdf(delta.b) 1.0 0.8 9.0 0.4 0.2 0.0 7 X Confidence Intervals We compute he confidence intervals to create the EDGE matrix. We place 1 in the matrix (the edge exists) only if:  $[-\epsilon, +\epsilon] \cap C_{n,\alpha} = \varnothing$ # let's start with alpha = 0.95 m <- choose(81, 2) alpha1 <- 0.05 t.alpha = quantile(f.hat, probs = (1-alpha1), type =1) # inverse of the ECDF function ci.lower = data.cor-t.alpha/sqrt(n) # lower bound of the CI ci.upper = data.cor+t.alpha/sqrt(n) # upper bound of the CI # and alpha = 0.99alpha2 <- 0.01 t.alpha2 = quantile(f.hat, probs = (1-alpha2), type =1) # inverse of the ECDF function ci.lower2 = data.cor-t.alpha2/sqrt(n) # lower bound of the CI ci.upper2 = data.cor+t.alpha2/sqrt(n) # upper bound of the CI Since we have the correlaction matrix and the confidence intervals relative to alpha = 0.95, we can plot our association graph. Edge matrix 1. EPSILON = 0 eps <- 0 # we initially set epsilon as 0 for (i in 1:n.features){ for(j in 1:n.features){ # If 0 is not included in the CI it means there is an edge between the two nodes if ((ci.lower[i,j] < eps & eps < ci.upper[i, j]) == FALSE) edge[i,j] = 1 # we set one in the adjacency matrix else edge[i,j] = 0 # if 0 is included in CI we set 0 in the adjacency matrix if (i == j) edge[i,j] = 0 # self edge is excluded graph.eps1 = graph.adjacency(edge, mode = "undirected") # Density density.eps1 <- edge\_density(graph.eps1, loops=FALSE)</pre> density.eps1 ## [1] 0.862963 As we can see the graph is very dense, since the maximum value is 1 it means that not all the nodes are connected to each other although it's very close to this result. # Connected components components(graph.eps1) # all nodes are connected ## \$membership ## 1R 2R 3R 4R 5R 6R 7R 8R 9R 10R 11R 13R 17R 18R 19R 20R 21R 22R 23R 24R ## 25R 28R 29R 30R 31R 32R 34R 35R 36R 37R 38R 39R 40R 41R 42R 43R 44R 46R 47R 1L ## 2L 3L 4L 5L 6L 7L 8L 9L 10L 11L 13L 17L 18L 19L 20L 21L 22L 23L 24L 25L ## 27L 28L 29L 30L 31L 32L 33L 34L 35L 36L 37L 38L 39L 40L 41L 42L 43L 44L 45L 46L ## 47L 1 ## \$csize ## [1] 81 ## \$no ## [1] 1 All the nodes are connected. # Number of edges sum(edge) ## [1] 5592 2. EPSILON = 0.5 eps2 <- 0.5 # we initially set epsilon as 0 edge2 <- matrix(NA, n.features, n.features)</pre> for (i in 1:n.features){ for(j in 1:n.features){ # If 0 is not included in the CI it means there is an edge between the two nodes if (eps2 < ci.lower[i, j] | -eps2 > ci.upper[i, j]) edge2[i,j] = 1 **else** edge2[i,j] = 0if (i == j) edge2[i,j] = 0 # self edge is excluded } graph.eps2 = graph.adjacency(edge2, mode = "undirected") # Components components(graph.eps2) ## \$membership ## [76] 1 1 1 1 1 1 ## \$csize ## [1] 71 1 2 1 1 1 1 1 1 1 ## \$no ## [1] 10 # Number of edges sum(edge2) # number of edges ## [1] 2958 3. EPSILON = 0.65 eps3 <- 0.65 # we initially set epsilon as 0 edge3 <- matrix(NA, n.features, n.features)</pre> for (i in 1:n.features){ for(j in 1:n.features){ # If 0 is not included in the CI it means there is an edge between the two nodes if (eps3 < ci.lower[i, j] | -eps3 > ci.upper[i, j]) edge3[i,j] = 1 **else** edge3[i,j] = 0if (i == j) edge3[i,j] = 0 # self edge is excluded graph.eps3 = graph.adjacency(edge3, mode = "undirected") # Components components(graph.eps3) ## \$membership ## [1] 1 1 1 2 3 1 4 1 1 5 1 6 7 8 1 1 1 1 1 1 1 9 10 11 12 ## [26] 1 1 1 13 14 1 15 16 1 1 17 1 1 18 1 1 1 1 1 1 1 1 1 1 20 ## [51] 21 1 1 1 22 1 1 1 1 23 24 1 25 26 27 1 1 20 28 1 1 1 29 30 31 ## [76] 1 1 1 1 1 1 ## \$csize ## [26] 1 1 1 1 1 1 ## \$no ## [1] 31 # Number of edges sum(edge3) # number of edges ## [1] 998 **PLOTS** # Change colours V(graph.eps1)\$color[1:39] <- "red"</pre> ## Warning in vattrs[[name]][index] <- value: number of items to replace is not a</pre> ## multiple of replacement length V(graph.eps1)\$color[40:81] <- "green" # Change colours V(graph.eps2)\$color[1:39] <- "red"</pre> ## Warning in vattrs[[name]][index] <- value: number of items to replace is not a</pre> ## multiple of replacement length V(graph.eps2)\$color[40:81] <- "green" # Plot only nodes with more than two edges clean <- function(graph){</pre> Isolated = which(degree(graph)<2)</pre> G2 = delete.vertices(graph, Isolated) return (G2) # Epsilon = 0 plot.igraph(clean(graph.eps1), ylab = expression(epsilon == 0), vertex.color=V(graph.eps1)\$colorx, vertex.label = legend("topright", c("right side","left side"), fill=c("red", "green")) right side left side # Epsilon = 0.5 plot.igraph(clean(graph.eps2), ylab = expression(epsilon == 0.5), vertex.color=V(graph.eps2)\$color, vertex.labellegend("topright", c("right side","left side"), fill=c("red", "green")) right side left side # Epsilon = 0.65 plot.igraph(clean(graph.eps3), ylab = expression(epsilon == 0.65), vertex.color=V(graph.eps2)\$color, vertex.labellegend("topright", c("right side","left side"), fill=c("red", "green")) right side left side

Graphically it's very cleat that the more epsilon grows while:  $\alpha=0.5$  the less edges we have (since the threshold for the correlation grows). But there is another interesting result to analyze, most of the nodes that

recap <- matrix(c(sum(edge), sum(edge2), sum(edge3)), nrow = 1, ncol = 3)</pre> colnames(recap) <- c("Eps = 0", "Eps = 0.5", "Eps = 0.65")rownames(recap) <- c("Number of nodes")</pre> recap Eps = 0 Eps = 0.5 Eps = 0.65## Number of nodes 5592 2958 Relationship between epsilon and edges # we repeat the same function for different values of epsilon n.edges <- function(epsilon, n.features, 1, u){</pre> e <- matrix(NA, n.features, n.features)</pre> for (i in 1:n.features){ for(j in 1:n.features){ # If 0 is not included in the CI it means there is an edge between the two nodes if (epsilon < l[i, j] | -epsilon > u[i, j]) e[i,j] = 1**else** e[i,j] = 0 if (i == j) e[i,j] = 0 # self edge is excluded

remain are the ones corrisponding to the **right** part of the brain.

Summary

return (sum(e))

# alpha = 0.05

for (eps in a){

# alpha = 0.01

for (eps in a){

n.edges.vec <- list()</pre>

n.edges.vec2 <- list()</pre>

out <- sinUG(data.cor, n)</pre>

# Change colours

## \$csize

## [26] 1

## \$no ## [1] 26

V(G.SIN)\$color[1:39] <- "red"</pre>

## multiple of replacement length

V(G.SIN)\$color[40:81] <- "green"

out <- round(out,3) # estimated partial correlation matrix</pre>

plotUGpvalues(out) # take a look at the p-values, they are almost all at 1

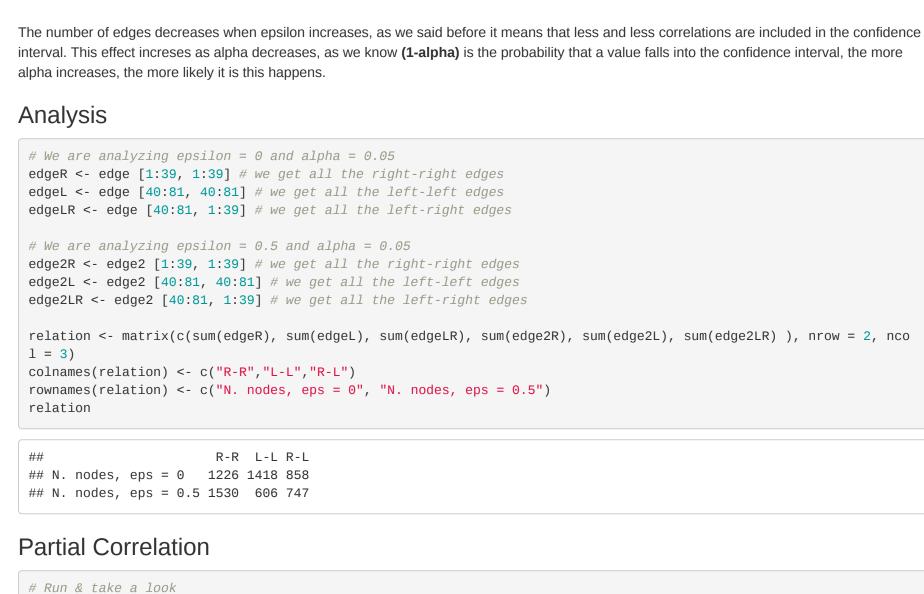
a <- seq(0,1, 0.01) # epsilon from 0 to 1

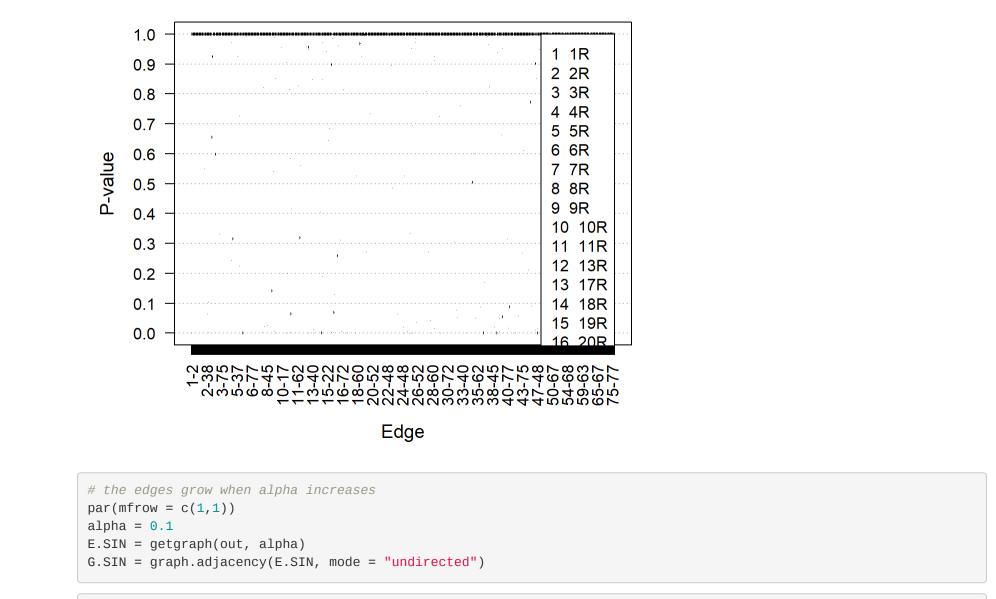
n.edges.vec <- append(n.edges.vec, n.edges1)</pre>

n.edges1 <- n.edges(eps, n.features, ci.lower, ci.upper)</pre>

n.edges2 <- n.edges(eps, n.features, ci.lower2, ci.upper2)</pre>

n.edges.vec2 <- append(n.edges.vec2, n.edges2)</pre> Plot par(mfrow=c(1,1))plot(a, lwd = 4, main="Different alphas", n.edges.vec, type = 'l',xlab = 'epsilon', ylab = 'n. edges', col = 'bl lines(a, n.edges.vec2, col="red", lwd = 4) legend("topright", c("alpha = 0.05", "alpha = 0.01"),fill=c("blue", "red")) **Different alphas**  $\blacksquare$  alpha = 0.05 alpha = 0.01 4000 3000 2000 1000 0 0.2 0.6 0.0 0.4 8.0 1.0 epsilon





plot.igraph(clean(G.SIN), ylab = expression(alpha == 0.1), vertex.color=V(G.SIN)\$color, vertex.label = NA)

## Warning in vattrs[[name]][index] <- value: number of items to replace is not a</pre>