

FACULTAD DE CIENCIAS



Universidad Zaragoza

APPENDIX: APPLICATION OF MENDELIAN RANDOMIZATION IN THE INFERENCE OF GENE REGULATORY NETWORKS

ANEXO: APLICACIÓN DE LA ALEATORIZACIÓN MENDELIANA EN INFERENCIA DE
REDES DE REGULACIÓN GENÉTICA

Trabajo Fin de Grado

Autor:

Claudia Llop Moreno

Tutor:

Dr. Joaquín Sanz Remón

Grado de Física.
2021-2022

Índice

1. Pipeline ingredients (i): data types under analysis	2
1.1. Measuring gene expression through RNA-seq	2
1.2. Genotyping chips	3
1.3. Genomic feature data and samples meta data	4
2. Pipeline ingredients (ii): Analytic modules	4
2.1. Phenotypes correlations: Co-expression networks analysis.	4
2.2. Control variables: eQTL-mapping	5
3. Trimmed M Means algorithm (TMM)	6
4. Benjamini-Hochberg against Storey-Tibshirani as the method for FDR estimation	7
5. Implementation of balanced-SNPs algorithm	9
6. Link direction inference under different genetic variable selection criteria	9
7. Codes	11
7.1. Prepare input tables	11
7.2. WCGNA	15
7.3. Identify eQTLs	19
7.4. Mendelian Randomization for Top-SNPs method	28
7.5. Mendelian Randomization for Balanced-SNPs method	36
7.6. Implementation of empirical null models to quantify directional bias in link prediction for Top-SNPs method	44
7.7. Implementation of empirical null models to quantify directional bias in link prediction for Balanced-SNPs method	46

1. Pipeline ingredients (i): data types under analysis

1.1. Measuring gene expression through RNA-seq

The main goal of this TFG is the analysis of the co-variation patterns of pairs of genes in a given experimental cohort, and the implementation of a statistical method to infer causal directions between the genes conforming the different pairs in a genome-wide gene co-expression network. Gene Co-expression Networks (GCN) are concerned with describing the correlation patterns of gene expression across sets of samples in which the expression of a large number of genes has been measured simultaneously.

Nowadays, the technologies of choice to characterize genome-wide expression levels in arbitrarily large cohorts of experimental samples are the so-called Next Generation sequencing technologies, or NGS, whose output is the sequencing of all the messenger RNA present in a given sample (RNA-seq). Using NGS technologies, it is nowadays possible to identify which of the genes encoded in our DNA are expressed in a given tissue, and to what extent, unlocking the interrogation of tens of thousands of genes in thousands of samples at once (1).

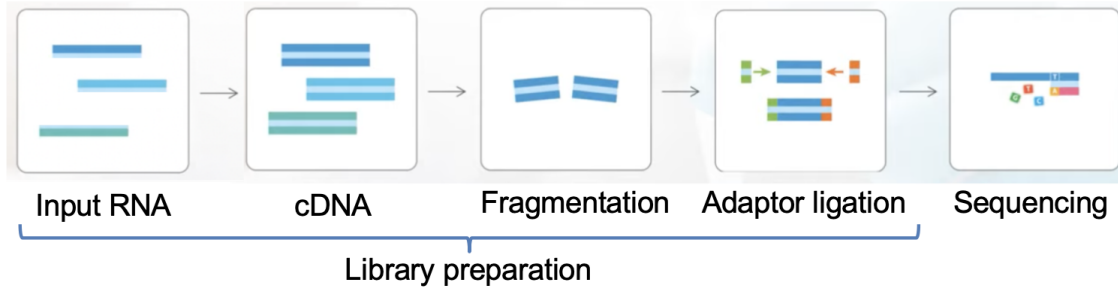


Figura 1: RNA-sequencing process.

RNA sequencing follows 5 main steps. From a sample of single-stranded RNA molecules we obtain cDNA (complementary DNA). The cDNA strand is synthesized to have double-stranded cDNA molecules whose nucleotide sequence is complementary to the original RNA. It is then chopped, adapters are placed at the ends of each fragment to facilitate sequencing and the resulting fragments, which constitute the so called sequencing library, are finally sequenced. In the following, we will look further into the process described.

The first 4 steps are part of what is called library preparation. From a given biological sample, cells are broken and their total RNA is extracted. This RNA is of many types: rRNA (ribosomal), tRNA (transfer RNA; small RNAs that carry amino acids to ribosomes for translation), mRNA (messenger) and other types. Since RNA-seq is typically bound to the sequencing of messenger RNA alone, we need a method to separate them so that we are left with only the mRNA that we are looking to sequence. This is done by an enrichment process. In our case we start with eukaryotic cells, whose mRNA molecules have a poly-A tail, which is leveraged using a Thymine-type bait to precipitate the mRNA we need and isolate.

Next, the sample, already strongly enriched in mRNA content, is incubated with reverse transcriptase to synthesize a cDNA strand. When it is completed, the RNA strand is hydrolyzed and a second (complementary) cDNA strand is generated. Then cDNA molecules are fragmented, generating short reads of homogeneous size. At this point, adapters are added to the fragments. These are small oligonucleotides that have different functions, including ensuring the hybridization of the fragments to complementary sequences that have been seeded in the surface of the

sequencing flowcell. Others are primer sequences that serve to start a complementary synthesis reaction. As sequencing denatures the strands (separates them), primers are needed to synthesize the complementary strands.

Once the library is prepared we move on to sequencing. In sequencing by synthesis -which is the method used in Illumina sequencing, the technology used in this study-, we start from a flow cell, that is, a surface seeded with DNA sequences complementary to the adapters placed at the ends of the fragments. The adapters adhere to these sequences. The number of strands in the sample needs to be increased for the signal to have sufficient amplitude to be detected. To do this, the amplification process begins. The parallel strand that remains attached to the floor of the flow cell is sequenced. This is folded and hybridizes with the nearby complementary adapter and generates the chain again as can be seen in Figure 2a. Then, the corresponding sequences are denatured which translates into two strands with the same sequence covalently attached to the flow cell. Finally, this process is repeated until enough copies of the original fragment are synthesized. This process is named bridge amplification.

For sequencing we pull nucleotides (Adenine, Guanine, Cytosine, Thymine) which are labeled with different colored fluorescent labels. They also have a radical that prevents further nucleotides from binding to the sequencing chain, therefore, when the corresponding nucleotide is anchored, it prevents the binding of the next nucleotide in the chain. The flow cell is then stimulated with a laser in such a way that the nucleotide returns the color of its fluorescent label and is recognized, and sequenced (Figure 2b). Then, the radical stopping further reactions is removed, and the process is iterated to continue reading nucleotides. In this way the sequence of nucleotides in each cluster of the flow cell, is registered.

Once all the fragments have been sequenced, a final step prior to statistical analysis is the mapping to a reference genome. This step briefly consists of the identification of the most likely genomic location, or locations, each fragment maps to, counting the number of fragments mapping into each gene or transcript. The result of that process is the matrix of expression data that we use as input for our analyses.

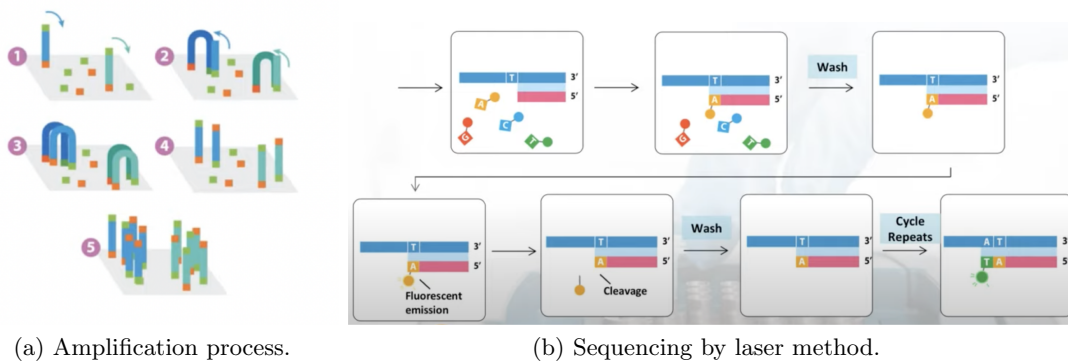


Figura 2: Parts of the RNA-sequencing process.

1.2. Genotyping chips

A single nucleotide polymorphism (SNP) is a type of single nucleotide variant (that is: a genomic position whose genotype is variable in a population), whose less frequent allele, is present in more than 1 % of the population. SNP chips are DNA microarrays that test genetic variation at thousands, or millions of specific locations across the genome(2), which makes them an excellent tool for studying common genetic variation, which can be used to assess ancestry as well as predisposition to many complex multifactorial diseases. (3)

A SNP chip briefly consists of an array of dwells containing a matrix of beads, on top of which, pairs of specific DNA sequences are seeded. These sequences are different in each dwell, and correspond to the complementary sequences immediately preceding the SNP you want to detect in each dwell (Figure 3).

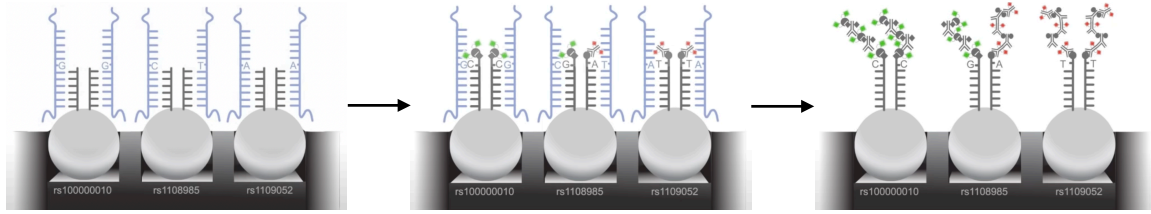


Figure 3: Genotyping chips for SNP sequencing process.

After loading the chip with a solution containing fragments of DNA of a given individual, the fragments that will hybridize to each bead will have the SNP of interest as the first free nucleotide that is not hybridized. Once you have hybridized the molecules, you pour a soup of labelled nucleotides of the 4 types onto the chip and each one is added to the strands that were originally linked to the SNP. Guanines and Cytosines are labeled with one color marker and Adenines and Thymines with another color. The original molecule hybridizing to the chip is washed and the fluorescent signal is amplified by the addition of further molecules to each bead. Finally, the distinction between nucleotides of each color allows the identification of the genotype of each SNP, either homozygous or heterozygous, what is commonly referred to as SNP-calling.

1.3. Genomic feature data and samples meta data

In addition to raw expression tables and SNP genotypes, in order to complete our analytical pipeline, we need to know certain attributes of the genomic features (genes and SNPs), as well as certain key information of our samples (sample metadata).

The only genomic features needed for completing these study are three: chromosome and genomic locations for genes and SNPs, (in order to distinguish between putatively cis- and -trans tests from the gene-SNO distances, and to exclude co-expression pairs involving nearby genes), and gene biotypes as well, since we will focus only on protein coding genes in our study.

Concerning sample meta.data, our samples come from an homogeneous panel of blood derived macrophages from male donors with balanced biological determinants, such as age and sex (all male). The only relevant metadata for our study in the sequencing batch each sample corresponds to among n=9 different batch levels.

2. Pipeline ingredients (ii): Analytic modules

2.1. Phenotypes correlations: Co-expression networks analysis.

A gene co-expression network (GCN) is an undirected graph, where each node corresponds to a gene, and a pair of nodes is connected with an edge if there is a significant co-expression relationship between them. (4) Having gene expression profiles of a number of genes for several samples or experimental conditions, a gene co-expression network can be constructed by looking

for pairs of genes which show a similar expression pattern across samples, since the transcript levels of two co-expressed genes often show significant co-variation patterns across samples, even when these are biological replicates. Gene co-expression networks are usually constructed using datasets generated by high-throughput gene expression profiling technologies such as RNA-Seq. (5). One of the most relevant approaches to the analysis of co-expression networks is based on modularity analysis, which, broadly speaking, aim at identifying groups of proteins whose interaction patterns in the network are stronger and more frequent than with the remaining of the system. Their characterization and functional characterization through ontology enrichment analyses constitutes a valuable tool to understand gene variation patterns in a given cohort.

In spite of its wide utility, co-expression networks are, per se, barely informative of causality relations across genes. To that end, more sophisticated inference approaches are needed to complete the transition from the (undirected) co-expression networks, -capturing correlations- to the (directed) regulatory networks (GRNs) -capturing causality-. In a GRN, a directed edge connects two genes, representing a biochemical process such as a reaction, transformation, interaction, activation or inhibition, or a set of them, that translates into a causal effect bounding the activity of the link sender to that of the link receiver. Compared to a GRN, a GCN does not attempt to infer the causality relationships between genes and in a GCN the edges represent only a correlation or dependency relationship among genes.

In this work, through Mendelian Randomization we seek to infer causality in a gene co-expression network (GCN).

2.2. Control variables: eQTL-mapping

It is known that SNPs located in regulatory regions, e.g. transcription factor (TF) binding sites, are often eQTLs (Expression quantitative trait loci), as they modulate gene expression. (6) An eQTL is a locus that explains a fraction of the genetic variance of a gene expression phenotype. Standard eQTL analysis involves a direct association test between markers of genetic variation with gene expression levels typically measured in tens or hundreds of individuals. This association analysis can be performed proximally (*cis*) or distally (*trans*) to the gene. (7)

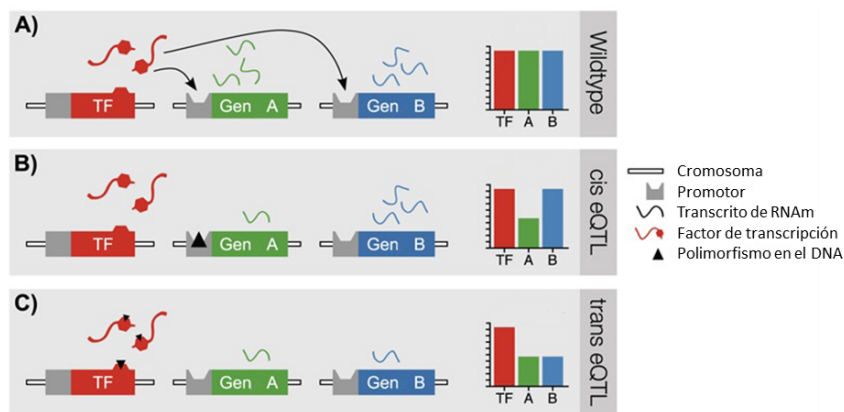


Figura 4: Illustration of an cis- and trans-eQTL concept (8).

Starting point at Figure 4 shows a transcription factor that regulates two copies, A and B, of the same gene: one inherited from the mother and the other from the father. In the first case we see what happens when you have a SNP (triangle) attached at the place where a TF has to bind nearby the target gene. The top row represents a case where regulation of both copies works

equally. In the bottom row, the presence of a cis-genetic variant nearby the gene inherited from the mother, translates into a difference in the expression levels between both copies (cis-eQTL). This situation represents what we know as cis-eQTL: it is a genetic variant that affects the expression level of a gene only in the same molecule in which the variant is found and not in another.

Finally, in the last row, we see a case when a variant is found in the coding part of the regulator, which, specially if this changes an amino acid of the TF which, in turn, modifies its affinity for the binding site of the target gene, resulting in less expression in both cases. Here the variant works in trans because it lies on one molecule but generates effects on other molecules. In this case the expression of both genes would be the same, if interrogated in a single individual, and we would only see differences when comparing with the levels across individuals. Typically cis-eQTLs are related as genetic variants close to a gene and trans as distant. Although not correct by definition, in practice this is true in the vast majority of cases.

3. Trimmed M Means algorithm (TMM)

The weighted and trimmed mean of M values, or TMM(9), is an algorithm implemented in the *edgeR*(10) software package that use scale factors per sample to normalize gene expression. These factors allow to transform the raw reads data to a relative measurement such as read counts per million reads sequenced, enabling comparisons of gene expression across samples of different complexity. To divide the observed number of counts ($Y_{g,k}$) of gene g in library k summarized from the raw reads by the total number of reads in the k -th library, N_k , is not enough to produce reliable comparisons across samples. Instead, we can define $\mu_{g,k}$ as the true and unknown expression level of the gene of interest g in sample k ; then, if L_k is the length of gene g , we can model the expected value of as follows:

$$E(Y_{gk}) = N_k \mu_{g,k} L_g / S_k \quad (1)$$

Where

$$S_k = \sum_{g'} \mu_{g',k} L_{g'} \quad (2)$$

Although the total RNA production $S_{k'}$ cannot be directly estimated, it is possible to estimate the relative production of two samples, $S_k/S_{k'} = \rho(k, k')$. Clearing the true expression level of gene g in library k from 1 as below:

$$\mu_{g,k} \approx \frac{Y_{gk} \cdot S_k}{N_k L_g} \quad (3)$$

the ratio of expression of that gene between two libraries k and k' remains:

$$\frac{\mu_{g,k}}{\mu_{g,k'}} = \frac{Y_{gk}}{N_k} / \frac{Y_{gk'}}{N'_k} \rho_g(k, k') \quad (4)$$

Now we use the approximation $\mu_{g,k} = \mu_{g,k'}$ assuming that all genes in our dataset have the same true expression between samples k and k' , i.e. the genes are not differentially expressed, and we can rewrite an independent estimation for the true ratio:

$$\frac{Y_{gk'}}{N'_k} / \frac{Y_{gk}}{N_k} = \rho_g(k, k') \quad (5)$$

This approximation is not quite correct, as it is not valid for all genes, for some may be strongly differentially expressed between samples in many cases. To overcome this limitation,

the extreme values of the distribution of $\rho_g(k, k')$ values are trimmed (30 % of them, by default), and the rest are averaged. The statistic that will be trimmed is:

$$\log(\rho_g(k, k')) = \log\left(\frac{Y_{gk'}}{N'_k}\right) - \log\left(\frac{Y_{gk}}{N_k}\right) \quad (6)$$

In this way, the most extreme estimates across all genes that are sufficiently expressed are excluded. The trimmed values will be averaged, and the weights of each observation will be extracted from the expected statistical uncertainty, expressed as the expected variance, associated to each of the observations:

$$\text{Var}(\log(\rho_g(k, k'))) = \left(\frac{1}{E(\rho_g(k, k'))}\right)^2 \cdot \text{Var}(\rho_g(k, k')) = \frac{N_{k'} - Y_{gk'}}{N_{k'} \cdot Y_{gk'}} + \frac{N_k - Y_{gk}}{N_k \cdot Y_{gk}} \quad (7)$$

From this equation, the weights will be obtained as the inverse values of those variances:

$$w_{g,k'} = \text{Var}(\log(\rho_g(k, k')))^{-1} \quad (8)$$

Now, once the weights, and the set G^* of genes that survive the trimming are defined, we defined the global trimmed-M-means, which is the final estimator of :

$$TMM_{k'} = \frac{\sum_{g \in G^*} w_{g,k'} \rho_g(k, k')}{\sum_{g \in G^*} w_{g,k'}} \quad (9)$$

Finally, the coefficients of all samples k' are calculated with respect to the common reference sample k using these $TMM_{k'}$ values as relative normalization coefficients (in the sense that $TMM_k = 1$). Once multiplied by the library depths (number of fragments per sample) they constitute suitable denominators for normalization of gene expression for comparisons across samples in each library. These normalization factors, in our case, are used to build \log_2 counts per million (*cpm*) expression estimates as:

$$\log_2(\text{cpm})_{gk} = \log_2(1E6 \cdot \frac{Y_{gk} + 0,5}{N_k TMM_k + 1}) \quad (10)$$

Where the " + 0,5" and " + 1" are introduced to avoid computational singularities when processing big datasets. In the context of our pipeline; coefficients are estimated using the function *calcnormfactors*, in the bioinformatic package *edgeR*(10), while the $\log_2(\text{cpm})$ transformation of eq. 10 is done using *voom*, in the *limma* (11) package.

4. Benjamini-Hochberg against Storey-Tibshirani as the method for FDR estimation

Now comes the first alteration we make to the work proposed by Regina Santesteban(12). In her case, she uses the Benjamini-Hochberg (BH)(13) correction on the vector of p-values. We will use the Storey-Tibshirani (ST)(14) for multiple testing correction method for greater precision. What are the differences between these two methods?

The BH method [5a] is equally useful for analyzing uniform null distributions and empiric nulls coming from permutations. However, this method is too conservative. The approach implicitly assumes that the true fraction of null hypothesis in your data (π_0) is 1. This is $f(p) = \pi_0 f_0$. On the other hand, the ST method [5b] drops the assumption of $\pi_0 = 1$ and move towards a two component model: $f(p) = \pi_0 f_0(p) + (1 - \pi_0) f_A(p)$.

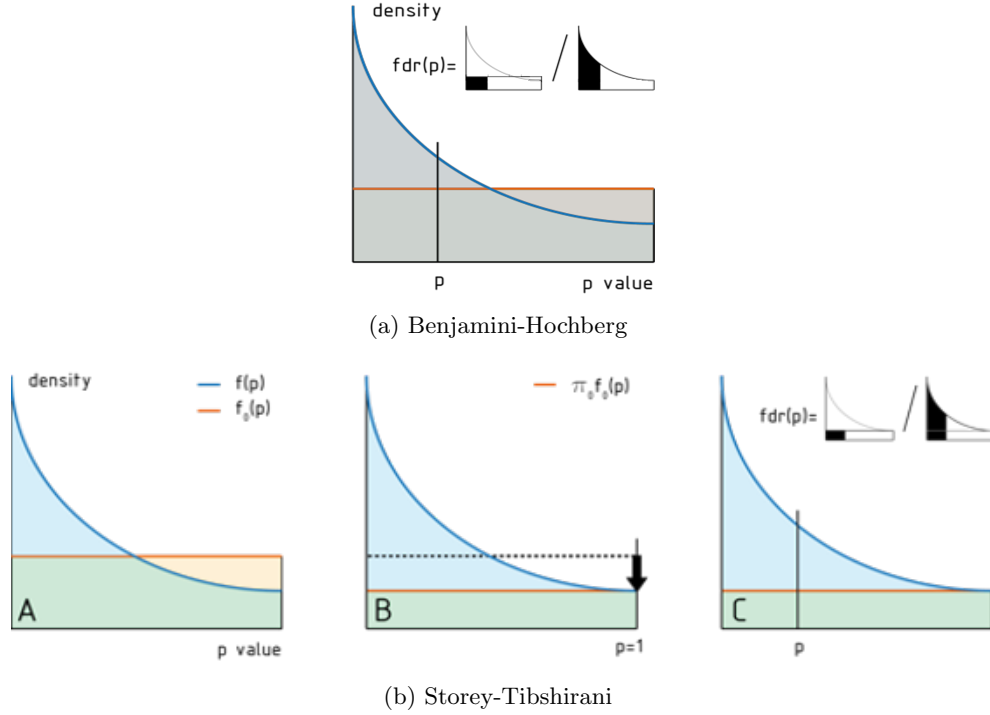


Figure 5: Differences between the two methods: BH & ST

Both methods offer their own definition of the FDR. For BH: $fdr(p) = F_0(p)/F(p)$ and for ST: $fdr(p) = \pi_0(p)F_0(p)/F(p)$.

The ST method provides an often considerably larger amount of statistical power with respect to BH. Contrary to BH, it also allows a straightforward definition for the False Non Discovery Rate: $FnDR(p) = 1 - \pi_0/\hat{\pi}_0(p)$ where $\hat{\pi}_0(p) = \frac{1-F(p)}{1-F_0(p)}$. On the downside, ST is only implemented against a uniform null distribution: a problem in some cases, when flat null distributions are not a good idea.

But what exactly do FDR and FnDR mean? If you look at the table below, the test correctly deemed significant are given by a . The ones correctly deemed not-significant are given by d . Type I errors (not-significant test called significant) in this case are given by b and type II (significant test called not-significant) by c . The value of the π_0 used before in this terms would be $\pi_0 = b + d$. Finally, the FDR is defined by the fraction of type I errors one makes after calling significant everything under a give p-value threshold, i.e. $fdr = b/(a + b)$. The FnDR is the fraction of type II errors one makes after calling non-significant everything above a given p-value threshold: $fndr = c/(d + c)$.

	True alternative hypothesis H_A	True null hypothesis H_0
Called significant	a	b
Called not-significant	c	d

Tabla 1: Significance elements through true null and alternative hypothesis.

False non discovery rates should be used for defining condition-specific tests: e.g. a gene is condition specific if we can confidently assure we can accept it as an alternative hypothesis under condition A (i.e. it is under FDR threshold at A), and, at the same time, we can confidently reject it at condition B (It is under FnDR threshold at B). This will be used later when the

definition of directionality is imposed.

5. Implementation of balanced-SNPs algorithm

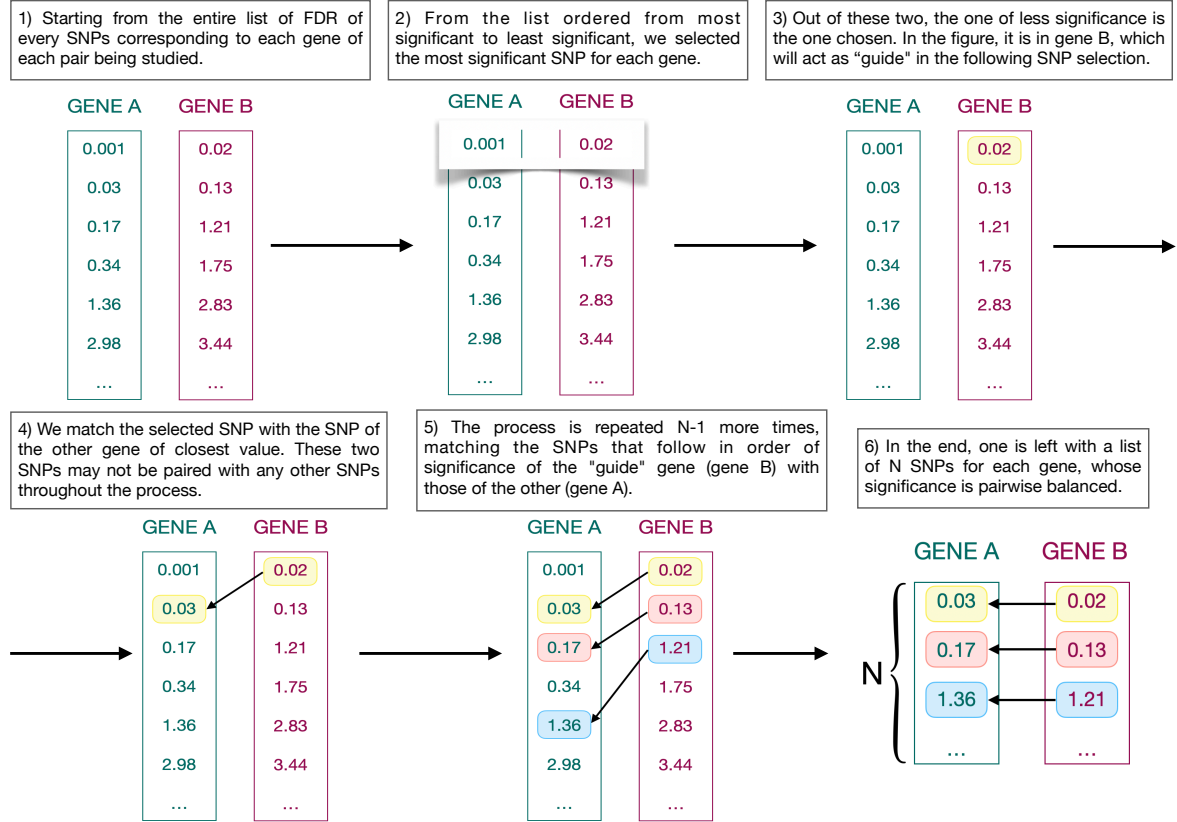


Figura 6: Process of balancing SNPs FDR for each pair of genes.

6. Link direction inference under different genetic variable selection criteria

The results presented in the report are reproduced by changing the forward/backward linkage definition strategy. In this case, instead of comparing the two most significant SNPs of each pair, we took the logarithmic sum of the significance of the N SNPs of each gene. These results are similar in magnitude to those presented in the main paper, noting again that the bias given by the balanced-SNPs algorithm is much smaller.

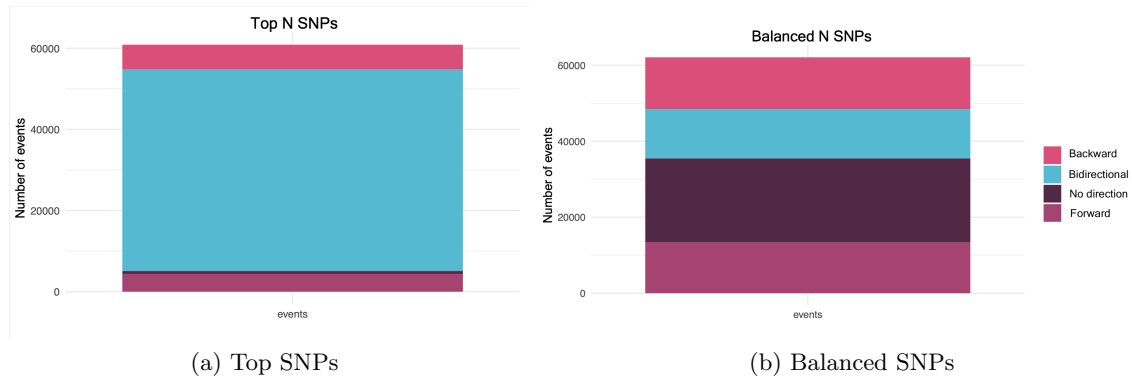


Figure 7: Fraction of links of defined direction results using the mean of each N SNP to discern hypothesis null or alternative.

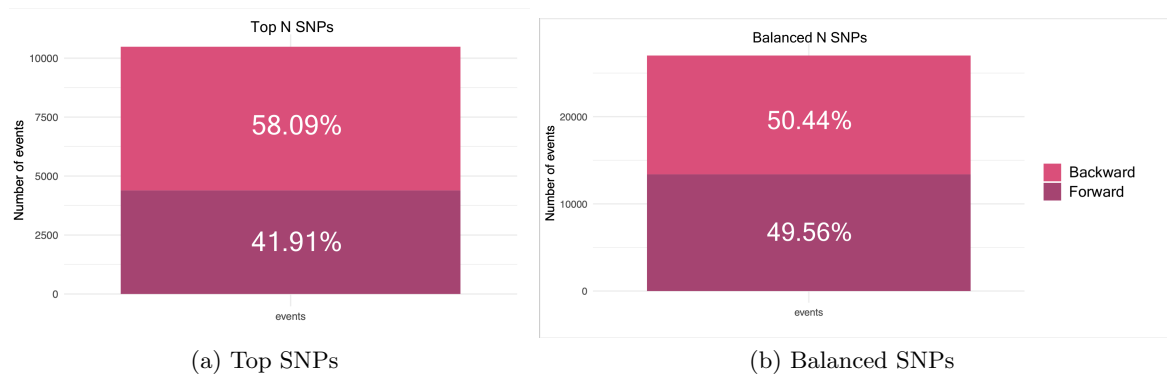


Figure 8: Fraction of links of defined forward/backward using the mean of each N SNP to discern hypothesis null or alternative.

7. Codes

7.1. Prepare input tables

```
1 #####
2 #### 1. Load dependencies #####
3 #####
4
5 first_run=0
6 {
7     ## bioMart is needed to get genes biotypes (identify protein coding genes)
8     ## SNPRelate is needed to calculate genotype Principal components.
9     library(edgeR)
10    library(limma)
11    library(biomaRt)
12    library(SNPRelate)
13
14    ul=function(tab,n=5)
15    {
16        rs=min(nrow(tab),n) #Coge valor mínimo entre 5 y el número de filas que hay
17        cs=min(ncol(tab),n) #Coge valor mínimo entre 5 y el número de columnas que hay
18        return(tab[1:rs,1:cs]) #Devuelve la matriz desde 1 hasta el numero minimo anterior para
19                                #columnas y filas
20    }
21    ur=function(tab,n=5)
22    {
23        rs=min(nrow(tab),n)
24        cs=min(ncol(tab),n)
25        return(tab[1:rs,(ncol(tab)-cs):ncol(tab)]) #¿Para qué sirve? ¿Por qué 5?
26    }
27    dcols=function(tab){data.frame(colnames(tab))} #Guarda los nombres de las columnas
28    drows=function(tab){data.frame(rownames(tab))} #Guarda los nombres de las filas
29 }

30 #####
31 #### 1. Load reads and metadata, subset relevant samples for our needs & make sample orders
32 concurrent ###
33 #####
34 {
35     ## Load the set of individuals that passed initial QC steps.
36     item=load("inputs/raw/individuals_sets.Rdata")
37     individuals=individuals_set$no_dup
38     ##94 guys.
39
40
41     ## Load the metadata; and subset only the samples corresponding to those guys, at time=2h and
42     condition="NI".
43     ## As you see we will be looking to only a subset of this rather big dataset.
44     metadata = read.table("inputs/raw/samples_data.txt")
45     metadata=metadata[which(metadata$TimePoint=="02h" & metadata$Treatment=="NI"),]
46     rownames(metadata)=paste0(metadata$Individual,"_",metadata$Treatment,"_",metadata$TimePoint)
47     ## In this condition/time, only 90 out of the 94 individuals were represented.
48
49     ## Load the raw expression matrix; then subset the individuals;
50     reads = read.table("inputs/raw/reads_Regina.txt",header=TRUE,check.names=FALSE)
51     reads=reads[,which(colnames(reads) %in% rownames(metadata))]
52
53     dim(reads)
54     dim(metadata)
55
56     reads=reads[,order(colnames(reads))]
57     metadata=metadata[order(rownames(metadata)),]
58
59     reads=reads[,order(colnames(reads))]
60     metadata=metadata[order(rownames(metadata)),]
```



```

63 #####
64 #### 2. Load genes dictionary; filter protein coding, autosomal ####
65 ##### genes & translate unfriendly gene names to HUGOs #####
66 #####
67
68 {
69   genes = read.table("inputs/raw/genes_data.txt", header = TRUE, stringsAsFactors = FALSE)
70   if(first_run==0)
71   {
72     human <- useMart("ENSEMBL_MART_ENSEMBL", dataset="hsapiens_gene_ensembl");
73     genesV2 = getLDS(attributes = c("ensembl_gene_id"), filters = "ensembl_gene_id", values =
74       as.character(genes$gene_id), mart = human, attributesL = c("gene_biotype"), martL = human,
75       uniqueRows=T)
76
77     from_biomart=list(mart=human,LDS=genesV2)
78     system("mkdir -p inputs/processed")
79     save(from_biomart, file = "inputs/processed/biomart_data.RData")
80   }else{
81     load("inputs/processed/biomart_data.RData")
82     human=from_biomart$mart
83     genesV2=from_biomart$LDS
84   }
85
86   protein_coding_genes=genesV2$Gene.stable.ID[which(genesV2$Gene.type=="protein_coding")]
87
88   genes=genes[which(genes$gene_id %in% protein_coding_genes),]
89   reads=reads[which(rownames(reads) %in% protein_coding_genes),]
90
91   genes=genes[order(genes$gene_id),]
92   reads=reads[order(rownames(reads)),]
93   length(which(rownames(reads)!=genes$gene_id))
94
95   length(unique(genes$gene_name))
96   ## There is some non-unique gene_names; this often happens, let us define a label ID that will
97   ## be the gene name for the first gene_name appearance, and, for the repeated ones, the ENSEMBL
98   ## IDs.
99   genes$label=genes$gene_name
100   genes$label[which(duplicated(genes$gene_name))]=genes$gene_id[which(duplicated(genes$gene_name))]
101   rownames(reads)=genes$label
102   rownames(genes)=genes$label
103
104   ## And order alphabetically both so it looks nice.
105   reads=reads[order(rownames(reads)),]
106   genes=genes[order(genes$label),]
107   length(which(rownames(reads)!=genes$label))
108
109   ## In this dictionary file I also have the positions of the genes, which will be needed to do
110   ## the EQTL mapping step.
111
112   ## Filter out non autosomal genes, and modify the chromosome notations:
113   colnames(genes)[3]="chr"
114   genes$chr=as.character(genes$chr)
115   genes=genes[which(genes$chr %in% as.character(c(1:22))),]
116   reads=reads[which(rownames(reads) %in% genes$label),]
117   length(which(rownames(reads)!=genes$label))
118   genes$chr=paste0("chr",genes$chr)
119   ## Filter only the columns needed by matrixEQTL, in the order they are needed too.
120   genes=genes[,c(7,3,4,5)]
121 }

```

```

120 #####
121 ### 3. Remove lowly expressed genes (and subset the genes list file accordingly) ###
122 #####
123
124 {
125     {
126         design=model.matrix(~1,data=metadata)
127         dge <- DGEList(counts=reads)
128         dge <- calcNormFactors(dge)
129         v <- voom(dge,design,plot=FALSE)
130         exp=v$E
131
132         voom_x=apply(exp,1,mean)
133         length(which(voom_x>1))
134         # 10586
135     }
136     reads=reads[which(voom_x>1),]
137     genes=genes[which(voom_x>1),]
138
139     length(which(rownames(reads)!=genes$label))
140 }

142 #####
143 ### 4. Load genotype data and genotype positions files. Subset genotype data. ###
144 #####
145 {
146     {
147         snpspos = read.table(paste0("inputs/raw/SNP_positions.txt"),header = TRUE, stringsAsFactors =
            FALSE)
148         rownames(snpspos)=snpspos$snp
149         unique(snpspos$chr)
150         ## Just to see that the file only contains snps in autosomes.
151         gtypes = read.table(paste0("inputs/raw/SNP_genotypes_data.txt"),header = TRUE, stringsAsFactors
            = FALSE)
152
153         gtypes=gtypes[,which(colnames(gtypes) %in% metadata$Individual)]
154         gtypes=gtypes[,order(colnames(gtypes))]
155         length(which(colnames(gtypes)!=metadata$Individual))
156         ## We will have to change the order of the samples:
157         gtypes=gtypes[,colnames(gtypes)]
158         reads=reads[,order(as.character(metadata$Individual))]
159         metadata=metadata[order(as.character(metadata$Individual)),]
160
161         length(which(rownames(metadata)!=colnames(reads)))
162         length(which(as.character(metadata$Individual)!=colnames(gtypes)))
163         ## Everything is ok, so let us simplify sample notation, given that we will not analyze samples
            from other times or treatments here.
164         ## We could filter out SNPs out from the cis-windows of the genes that we want to analyze, but
            it is easier doing that within the matrixEQTL run.
165         ## Now, filter for minimum allele frequency:
166         mafs=apply(gtypes,1,function(x){length(which(x>0))/length(x)})
167         gtypes=gtypes[which(mafs>0.05),]
168         snpspos=snpspos[which(mafs>0.05),]
169         length(which(rownames(gtypes)!=rownames(snpspos)))
170     }

```

```

174 #####
175 #### 5. Get the genotype principal components and configure the useful metadata file ####
176 #####
177 {
178     samples=colnames(gtypes) #Individuos I_xx
179     gtypes_pca=data.frame(snp_id=rownames(gtypes),gtypes) #Se crea una fila con el primer elemento
180     snp_id y los siguientes I_xx
181     system("mkdir -p inputs/aux") #Crea un directorio
182     snpgdsCreateGeno("inputs/aux/GDS_genotypes.gds", genmat = as.matrix(gtypes_pca[, samples]),
183     sample.id = unique(samples), snp.id = gtypes_pca$snp_id, snpfirstdim=TRUE)
184     #Create a SNP genotype dataset from a matrix
185     snpgdsSummary("inputs/aux/GDS_genotypes.gds")
186     genofile <- snpgdsOpen("inputs/aux/GDS_genotypes.gds")
187     pca <- snpgdsPCA(genofile) #Principal Component Analysis (PCA) on SNP genotype data
188     tab <- data.frame(sample.id = pca$sample.id,
189     PC1 = pca$eigenvect[,1], # the first eigenvector
190     PC2 = pca$eigenvect[,2], PC3 = pca$eigenvect[,3], PC4 = pca$eigenvect[,4], PC5 =
191     pca$eigenvect[,5],
192     stringsAsFactors = FALSE)
193     metadata=cbind(metadata,tab)
194     pca$varprop[1]*90
195     # 1.568026 first PC does only explain 1.7% of total variance, 57% more variance than a random
196     direction: there is no population structure.
197     length(which(metadata$Individual!=tab$sample.id))
198     metadata=cbind(Individual=tab$sample.id,Batch=metadata$Batch,tab[,c(2:6)])
199     rownames(metadata)=tab$sample.id
200     {
201         colnames(reads)=metadata$Individual
202         colnames(gtypes)=metadata$Individual
203         rownames(metadata)=metadata$Individual
204     }
205 }

210 #####
211 #### 6. Check the congruence & Write the procees input tables. ####
212 #####
213 {
214     dim(reads)
215     #[1] 10586 90
216     dim(metadata)
217     #[1] 90 7
218     dim(gtypes)
219     #[1] 6159725 90
220     dim(genes)
221     #[1] 10586 4
222     dim(snpupos)
223     #[1] 6159725 3
224
225     ## Samples congruence
226     length(which(colnames(reads)!=rownames(metadata)))
227     length(which(colnames(reads)!=colnames(gtypes)))
228     ## Genes congruence
229     length(which(rownames(reads)!=rownames(genes)))
230     ## SNPs congruence.
231     length(which(rownames(gtypes)!=rownames(snpupos)))
232
233     system("mkdir -p inputs/processed")
234     write.table(reads,"inputs/processed/reads_ok.txt")
235     write.table(metadata,"inputs/processed/metadata_ok.txt")
236     write.table(genes,"inputs/processed/genes_ok.txt")
237     write.table(snpupos,"inputs/processed/snpupos_ok.txt")
238     write.table(gtypes,"inputs/processed/gtypes_ok.txt")
239 }

```

7.2. WCGNA

```
1 #####
2 #### 1. Load dependencies #####
3 #####
4 {
5     ## gdata is needed for the function of upper & lower triangles.
6     ## cowplot aligns ggplot panels nicely.
7     ## ggrepel has algorithm to add labels in data in ggplot plots a way that they repel each other.
8
9     library(edgeR)
10    library(limma)
11    library(ggplot2)
12    library(cowplot)
13    library(gdata)
14    library(ggrepel)
15
16    ul=function(tab,n=5)
17    {
18        rs=min(nrow(tab),n)
19        cs=min(ncol(tab),n)
20        return(tab[1:rs,1:cs])
21    }
22    ur=function(tab,n=5)
23    {
24        rs=min(nrow(tab),n)
25        cs=min(ncol(tab),n)
26        return(tab[1:rs,(ncol(tab)-cs):ncol(tab)])
27    }
28    dcols=function(tab){data.frame(colnames(tab))}
29    drows=function(tab){data.frame(rownames(tab))}
30
31 }
32
33
34
35 #####
36 #### 2. Load reads and metadata, processed ####
37 #####
38 {
39     reads=read.table("inputs/processed/reads_ok.txt")
40     metadata=read.table("inputs/processed/metadata_ok.txt")
41 }
42
43
44 #####
45 #### 3. Do log transformation, batch removal & variance stabilization ####
46 #####
47 {
48
49     ##Get voom tranformed data
50     design=model.matrix(~Batch,data=metadata)
51     dge <- DGEList(counts=reads)
52     dge <- calcNormFactors(dge)
53     v <- voom(dge,design,plot=FALSE)
54     write.table(design,paste0("outputs/2_coexpression/design.txt"))
55
56     ## Do the fit of the y values (do not use v as is as input of lmfit: we just want to run least
57     squares on the y matrix of log(cpm):
58
59     exp=v$E
60     write.table(exp,paste0("outputs/2_coexpression/exp.txt"))
61     fit2=lmFit(exp,design)
62     write.table(fit2,paste0("outputs/2_coexpression/fit2.txt"))
63 }
```

```

64
65 ## Define the voom plot axis:
66 voom_x=apply(log2(reads + 0.5),1,mean)
67 voom_y=(fit2$sigma)^0.25
68 voom_table=data.frame(x=voom_x,y=voom_y)
69
70
71 ## Get the table with the lowess fit.
72 span=0.5
73 l <- data.frame(lowess(voom_table$x, voom_table$y, f = span))
74 ## In some parts, you get two or three points identical until the last decimal figure printed:
75 l=l[!duplicated(l$x),]
76 write.table(l,paste0("outputs/2_coexpression/lowess_fit.txt"))
77
78 ## This defines a function-like object that represents the lowess fit.
79 tech_var_func=approxfun(l$x, l$y,rule=2)
80 voom_table$tech_var=tech_var_func(voom_table$x)
81 write.table(voom_table,paste0("outputs/2_coexpression/voom_table.txt"))
82 pl_voom_raw=ggplot(voom_table)+geom_point(aes(x=x,y=y))+geom_line(aes
    (x=x,y=tech_var),color="red")+theme_minimal()+xlab("Mean expression")+ylab("")
83
84 ## Now, we get the scaled residuals:
85 aux=as.matrix(log2(reads + 0.5))
86 scalings=aux
87 for(i in 1:ncol(scalings))scalings[,i]=tech_var_func(aux[,i])^4
88 res<-residuals(object=fit2, exp)
89 write.table(res,paste0("outputs/2_coexpression/res.txt"))
90 scaled_res=res/scalings
91 write.table(scaled_res,paste0("outputs/2_coexpression/scaled_res.txt"))
92
93
94 ## And reconstruct the expression where real residuals are substituted with scaled ones:
95 exp_prediction=exp-res
96 exp_stabilized=exp_prediction+scaled_res
97
98 ## Then, run the fit again, and retrieve the sigmas:
99 fit3=lmFit(exp_stabilized,design)
100 write.table(fit3,paste0("outputs/2_coexpression/fit3.txt"))
101 voom_y_stabilized=sqrt(fit3$sigma)
102
103 ## This is the same as before with the stabilized data
104 voom_table_stabilized=data.frame(x=voom_x,y=voom_y_stabilized)
105 #span=1
106 l_stabilized <- data.frame(lowess(voom_table_stabilized$x, voom_table_stabilized$y, f = span))
107 l_stabilized=l_stabilized[!duplicated(l_stabilized$x),]
108 write.table(l_stabilized,paste0("outputs/2_coexpression/lowess_fit_estabilized.txt"))
109 tech_var_func_stabilized=approxfun(l_stabilized$x, l_stabilized$y,rule=2)
110 voom_table_stabilized$tech_var=tech_var_func_stabilized(voom_table_stabilized$x)
111 write.table(voom_table_stabilized,paste0("outputs/2_coexpression/voom_table_stabilized.txt"))
112 pl_voom_stabilized=ggplot(voom_table_stabilized)+geom_point(aes(x=x,y=y))+geom_line(aes
    (x=x,y=tech_var),color="red")+xlim(6,16)+theme_minimal()+xlab("Mean expression")+ylab("")
113 voom_plots=plot_grid(pl_voom_raw,pl_voom_stabilized,ncol=2,align="h")
114
115 ## And here we have the expression matrix once Batch effects have been removed and the variance
    has been stabilized: optimal setup for coexpression network reconstruction.
116
117 exp=t(scaled_res) #Traspuesta
118 write.table(exp,paste0("outputs/2_coexpression/exp.txt"))
119
120 }

```

```

121 #####
122 #### 4. Get pairwise co-expression tests & network object, storing only the significant ones ####
123 #####
124
125 {
126     threshold_DC=0.01
127     dofs=nrow(exp)-2-(ncol(design)-1)
128
129     correlations_matrix <- cor(exp,method="pearson")
130     write.table(correlations_matrix,paste0("outputs/2_coexpression/correlations_matrix.txt"))
131     t_matrix=sqrt(dofs)*correlations_matrix/(sqrt(1-correlations_matrix*correlations_matrix))
132     write.table(t_matrix,paste0("outputs/2_coexpression/t_matrix.txt"))
133     p_value_matrix=2*pt(-abs(t_matrix),dofs)
134     write.table(p_value_matrix,paste0("outputs/2_coexpression/p_value_matrix.txt"))
135     p_value_vector=upperTriangle(p_value_matrix)
136     write.table(p_value_vector,paste0("outputs/2_coexpression/p_value_vector.txt"))
137
138     # From BH to ST
139     library(qvalue)
140     fdrs_vector=qvalue(p_value_vector)$qvalues
141     #fdrs_vector=p.adjust(p_value_vector,method="BH")
142     ## Benjamini-Hochberg
143
144     fdrs_matrix=p_value_matrix
145     upperTriangle(fdrs_matrix)=fdrs_vector
146     fdrs_matrix=t(fdrs_matrix)
147     upperTriangle(fdrs_matrix)=fdrs_vector
148     fdrs_matrix=t(fdrs_matrix)
149
150     rownames(fdrs_matrix)=rownames(p_value_matrix)
151     colnames(fdrs_matrix)=colnames(p_value_matrix)
152
153
154     network=data.frame(which(fdrs_matrix<threshold_DC,arr.ind=TRUE))
155     network=network[which(network$row<network$col),]
156
157     network$gen_row=rownames(fdrs_matrix)[network$row]
158     network$gen_col=rownames(fdrs_matrix)[network$col]
159
160     load_data=function(x,y){correlations_matrix[x,y]}
161     network$r=mapapply(load_data,x=network$row,y=network$col)
162
163     load_data=function(x,y){t_matrix[x,y]}
164     network$t=mapapply(load_data,x=network$row,y=network$col)
165
166     load_data=function(x,y){p_value_matrix[x,y]}
167     network$p=mapapply(load_data,x=network$row,y=network$col)
168
169     load_data=function(x,y){fdrs_matrix[x,y]}
170     network$fdr=mapapply(load_data,x=network$row,y=network$col)
171
172     network=network[order(network$p),]
173     rownames(network)=paste0(network$gen_row,"_",network$gen_col)
174     involved=unique(c(network$gen_row,network$gen_col))
175
176     length(involved)
177     # 10033
178
179     exp_sub=exp[,which(colnames(exp) %in% involved)]
180     scaled_res_sub=scaled_res[which(rownames(scaled_res) %in% involved),]
181
182     pairwise_matrixes=list(correlations_matrix=correlations_matrix,
183                           t_matrix=t_matrix, p_value_matrix=p_value_matrix,
184                           fdrs_matrix=fdrs_matrix,dofs=dofs)
185 }

```



```

186 #####
187 #### 5. Get WGCNA modules ####
188 #####
189
190 if(FALSE){
191   # Call the network topology analysis function
192   sft = pickSoftThreshold(exp_sub, powerVector = seq(2,30,by=2), blockSize = ncol(exp_sub),verbose =
193     5)
194   pow=25
195   # Plot the results:
196   power_tab=sft$fitIndices
197   power_tab$label=""
198   power_tab$label[which(power_tab$SFT.R.sq>0.85)[1]]=as.character(power_tab$Power
199     [which(power_tab$SFT.R.sq>0.85)[1]])
200
201   p1=ggplot(power_tab)+geom_point(aes(x=Power,y=SFT.R.sq))+geom_line(aes(x=Power,y=SFT.R.sq))+
202     geom_hline(yintercept=0.85)+geom_vline(xintercept=power_tab$Power
203       [which(power_tab$SFT.R.sq>0.85)[1]])+
204     geom_text_repel(aes(x=Power,y=SFT.R.sq,label=label))+
205     ggtitle(paste0("Scale freedom
206       (pow.est=",power_tab$Power[which(power_tab$SFT.R.sq>0.85)[1]],")"))
207
208   power_tab$label[which(power_tab$SFT.R.sq>0.85)[1]]=paste0("Power=",as.character(power_tab$Power
209     [which(power_tab$SFT.R.sq>0.85)[1]]),", <k>=",
210     as.character(round(power_tab$mean.k.[which(power_tab$SFT.R.sq>0.85)[1]])))
211
212
213
214   p2=ggplot(power_tab)+geom_point(aes(x=Power,y=mean.k.))+geom_line(aes(x=Power,y=mean.k.))+
215     geom_hline(yintercept=0.85)+geom_vline(xintercept=power_tab$Power
216       [which(power_tab$SFT.R.sq>0.85)[1]])+
217     geom_text_repel(aes(x=Power,y=mean.k.,label=label),vjust=1,hjust=1)+
218     ggtitle(paste0("Mean degree (pow.est=",power_tab$Power[which(power_tab$SFT.R.sq>0.85)[1]],")"))
219   p1_WGCNA=plot_grid(p1,p2,ncol=2,align="h")
220
221   ## This, with saveTOMs=TRUE will save the TOM file into the desired folder, which is this one:
222   system("mkdir -p outputs/2_coexpression/")
223   net = blockwiseModules(exp_sub, power = pow, TOMType = "unsigned", minModuleSize = 30,
224     reassignThreshold = 0, mergeCutHeight = 0.25, numericLabels = TRUE, pamRespectsDendro = FALSE,
225     maxBlockSize=length(involved)+1, saveTOMs = TRUE, saveTOMFileBase = "outputs/2_coexpression/",
226     verbose = 3)
227
228   ## Let's put exp_sub back to the more normal position:
229   exp_sub=data.frame(t(exp_sub))
230 }

```

```

224 #####
225 #### 6. Save results ####
226 #####
227
228 {
229     ### Variance stabilized and batch-free expression matrix.
230     write.table(scaled_res_sub,paste0("outputs/2_coexpression/stabilized_expression.txt"))
231
232     ## Pairwise objects and network table:
233     save(pairwise_matrixes,file="outputs/2_coexpression/pairwise_matrixes.Rdata")
234     write.table(network,paste0("outputs/2_coexpression/network",format(threshold_DC,digits=2),".txt"))
235
236     ## WGCNA net object containing the modules & some plots from the WGCNA analyses
237     # JQ lo pongo dentro del condicional
238     if(modules)
239     {
240         save(net,file="outputs/2_coexpression/WGCNA_net.Rdata")
241
242         pdf("outputs/2_coexpression/WGCNA_plots.pdf",width=7,height=4)
243         print(pl_WGCNA)
244         dev.off()
245     }
246
247     genes=read.table("inputs/processed/genes_ok.txt")
248     ### Get only the genes that appeared in the coexp. network:
249     genes=genes[which(rownames(genes) %in% colnames(exp_sub)),]
250     length(which(rownames(genes)!=colnames(exp_sub)))
251     length(which(rownames(genes)==colnames(exp_sub)))
252
253     write.table(genes,paste0("outputs/2_coexpression/coexpressed_genes.txt"))
254
255 }

```

7.3. Identify eQTLs

```

1 #####
2 #### 1. Load dependencies #####
3 #####
4
5 {
6     library(MatrixEQTL)
7     library(edgeR)
8     library(limma)
9     library(gdsfmt)
10    library(SNPRelate)
11    library(qvalue)
12    library(stats)
13    library(stringr)
14    library(data.table)
15
16    ul=function(tab,n=5)
17    {
18        rs=min(nrow(tab),n)
19        cs=min(ncol(tab),n)
20        return(tab[1:rs,1:cs])
21    }
22    ur=function(tab,n=5)
23    {
24        rs=min(nrow(tab),n)
25        cs=min(ncol(tab),n)
26        return(tab[1:rs,(ncol(tab)-cs):ncol(tab)])
27    }
28    dcols=function(tab){data.frame(colnames(tab))}
29    drows=function(tab){data.frame(rownames(tab))}
30
31 }
32

```



```

33 #####
34 #### 2. Load items needed for EQTL mapping #####
35 #####
36 {
37     SNPs_positions_file_name="inputs/processed/snpsspos_ok.txt"
38     expression_file_name="outputs/2_coexpression/stabilized_expression.txt"
39     genes_file_name="inputs/processed/genes_ok.txt"
40     Genotypes_file_name="inputs/processed/gtypes_ok.txt"
41 }
42
43 #####
44 #### 3. Check input items congruences #####
45 #####
46 {
47     snpsspos=read.table(SNPs_positions_file_name)
48     gtypes=read.table(Genotypes_file_name)
49     genes=read.table(genes_file_name)
50     expression=read.table(expression_file_name)
51     dim(expression)
52     # [1] 10033      90
53     dim(gtypes)
54     # [1] 6159725    90
55     dim(genes)
56     # [1] 10033      4
57     dim(snpsspos)
58     # [1] 6159725    3
59     ## Some checks:
60     ## Samples wise
61     length(which(colnames(expression)==colnames(gtypes)))/length(colnames(expression))
62     ## Genes wise:
63     length(which(rownames(expression)==rownames(genes)))/length(rownames(expression))
64     ## SNPs-wise (we could filter these out using the positions, but we will leave matrixEQTL do it)
65     length(which(rownames(snpsspos)==rownames(gtypes)))/length(rownames(snpsspos))
66 }

72 #####
73 #### 4. Build matrixEQTL inputs #####
74 #####
75
76 {
77     ## Matrix EQTL needs its inputs in .txt input tables saved in a particular format:
78     write.table(expression,expression_file_name, quote=F, sep="\t", row.names=TRUE)
79     write.table(genes,genes_file_name, quote=F, sep="\t", row.names=TRUE)
80     write.table(snpsspos,SNPs_positions_file_name, quote=F, sep="\t", row.names=TRUE)
81     write.table(gtypes,Genotypes_file_name, quote=F, sep="\t", row.names=TRUE)
82 }
83
84 #####
85 #### 5. Configure & Run Matrix for cisEQTL #####
86 #####
87
88 {
89     gene = SlicedData$new();
90     gene$fileDelimiter = "\t"; # the TAB character
91     gene$fileOmitCharacters = "NA"; # denote missing values;
92     gene$fileSkipRows = 1; # one row of column labels
93     gene$fileSkipColumns = 1; # one column of row labels
94     gene$fileSliceSize = 2000; # read file in slices of 2,000 rows
95     gene$LoadFile(expression_file_name);
96
97     ## Load covariates: there is none, so we declare an empty cov. file:
98

```

```

98
99 covariates_file_name=character()
100 cvrt = SlicedData$new();
101 cvrt$fileDelimiter = "\t";      # the TAB character
102 cvrt$fileOmitCharacters = "NA"; # denote missing values;
103 cvrt$fileSkipRows = 1;         # one row of column labels
104 cvrt$fileSkipColumns = 1;      # one column of row labels
105 if(length(covariates_file_name)>0) {
106   cvrt$LoadFile(covariates_file_name);
107 }
108
109 ## Load genotype data
110
111 snps = SlicedData$new();
112 snps$fileDelimiter = "\t";      # the TAB character
113 snps$fileOmitCharacters = "NA";
114 snps$fileOmitCharacters = "-9" # denote missing values;
115 snps$fileSkipRows = 1;         # one row of column labels
116 snps$fileSkipColumns = 1;      # one column of row labels
117 snps$fileSliceSize = 2000;     # read file in slices of 2,000 rows
118 snps$LoadFile(Genotypes_file_name)
119
120 ## Set up further program parameters
121 useModel = modelLINEAR
122 output_file_name_cis = tempfile()
123 pvOutputThreshold_cis = 1
124 pvOutputThreshold = 0;
125 errorCovariance = numeric()
126 cisDist = 1e5
127 output_file_name = tempfile()
128 output_file_name_cis = tempfile()
129
130 me = Matrix_eQTL_main(snps = snps, gene = gene, cvrt = cvrt, output_file_name = output_file_name,
131   useModel = useModel,
132   errorCovariance = errorCovariance, verbose = TRUE, output_file_name.cis = output_file_name_cis,
133   pvOutputThreshold = pvOutputThreshold, pvOutputThreshold.cis = pvOutputThreshold_cis, snpspos
134     = snpspos,
135   genepos = genes, cisDist = cisDist, pvalue.hist = "qqplot", min.pv.by.genesnp = TRUE,
136   noFDRsaveMemory = FALSE);
137 cis_eqtls=me$cis$eqtls
138 # 5766136      6
139
140 ## Filter out tests corresponding to SNPs in perfect linkage disequilibrium.
141
142 elemental=function(x){
143   paste0(x,collapse="")
144 }
145
146 gtypes_collapsed=apply(gtypes,1,elemental)
147 snpspos$gtypes_collapsed=gtypes_collapsed
148
149 revert=function(word){
150   word=str_replace_all(word,"0","3")
151   word=str_replace_all(word,"2","0")
152   word=str_replace_all(word,"3","2")
153   return(word)
154 }
155
156 snpspos$gtypes_collapsed_reversed=revert(snpspos$gtypes_collapsed)
157
158 aux=snpspos[,c(1,4,5)]
159 colnames(aux)[1]="snps"
160 result=merge(cis_eqtls, aux, by="snps")
161 result$test_ID=paste0(result$gtypes_collapsed,"_",result$gene)
162 result$test_ID_reversed=paste0(result$gtypes_collapsed_reversed,"_",result$gene)
163 # 5766136      10

```

```

161
162     result=result[which(!duplicated(result$test_ID)),]
163     # 2477579      10
164
165     result=result[which(!result$test_ID %in% result$test_ID_reversed),]
166     # 2409953      10
167
168     result$FDR_ok=qvalue(result$pvalue)$qvalues
169     cis_eqtls=result
170     # 2409953      11
171
172 }
173

```

```

174 #####
175 ### 6. Run again in trans for conjugated EQTLs: algorithm top #####
176 #####
177
178 {
179     ## First: select genes with at least one EQTL at 5% FDR
180
181     threshold=0.05
182     hits=cis_eqtls[which(cis_eqtls$FDR_ok<threshold),]
183     genes_with_at_least_one_eqtl=unique(hits$gene)
184     length(genes_with_at_least_one_eqtl)
185     # 815
186     cis_useful=cis_eqtls[which(cis_eqtls$gene %in% genes_with_at_least_one_eqtl),]
187     cis_useful=cis_useful[order(cis_useful$gene,cis_useful$FDR_ok),]
188     cis_useful=cis_useful[,c(1,2,3,4,11,6)]
189     # 280078      6
190     # cis_useful contains all cis-tests pointing to genes for which at least one EQTL at FDR=5% is
191       found (815 of such genes)
192
193     ## Let's pick up the top-3 most significant EQTLs per gene.
194     N=3
195
196     first_instances=cis_useful[which(!(duplicated(cis_useful$gene))),]
197     resto=cis_useful[which((duplicated(cis_useful$gene))),]
198
199     for(iter in 2:N){
200         appendix=resto[which(!(duplicated(resto$gene))),]
201         if(nrow(appendix)<length(genes_with_at_least_one_eqtl))
202         {
203             print(paste("In iter=",iter," some genes lack SNPs"))
204         }
205         first_instances=rbind(first_instances,appendix)
206         resto=resto[which((duplicated(resto$gene))),]
207     }

```

```

207
208 top_cis_snps=first_instances[order(first_instances$gene,first_instances$FDR_ok),]
209 dim(top_cis_snps)
210 # [1] 2445      6
211
212 length(unique(top_cis_snps$snps))
213 ## 2387: less bc. some SNPs are EQTL of more than 1 nearby gene
214
215 snpspos_trans_subset_top=snpspos[which(snpspos$snp %in% top_cis_snps$snps),c(1:3)]
216 dim(snpspos_trans_subset_top)
217 # [1] 2387      3
218
219 # Let us subset the genotype tables to contain only those 2387 SNPs
220 gtypes_trans_subset_top=gtypes[which(rownames(gtypes) %in% top_cis_snps$snps),]
221 # [1] 2387     90
222
223 # And the expression and gene-data tables to contain only those 815 genes.
224 genes_trans_subset_top=genes[which(rownames(genes) %in% top_cis_snps$gene),]
225 # [1] 815      4
226 expression_trans_subset_top=expression[which(rownames(expression) %in% top_cis_snps$gene),]
227 # [1] 815     90
228
229 ## Check the congruence of these new arguments:
230 ## Sample-wise
231 length(which(colnames(expression_trans_subset_top)==colnames(gtypes_trans_subset_top)))/length
  (colnames(expression_trans_subset_top))
232 ## Genes wise:
233 length(which(rownames(expression_trans_subset_top)==rownames(genes_trans_subset_top)))/length
  (rownames(expression_trans_subset_top))
234 ## SNPs-wise (we could filter these out using the positions, but we will leave matrixEQTL do it)
235 length(which(rownames(snpspos_trans_subset_top)==rownames(gtypes_trans_subset_top)))/length
  (rownames(snpspos_trans_subset_top))
236

```

```

237 ## Write them in matrix-EQTL-friendly format
238 system("mkdir -p outputs/3_EQTLs/trans_subset_inputs/top")
239
240 trans_subset_top_SNP_positions_file_name="outputs/3_EQTLs/trans_subset_inputs/top/snpspos_ok.txt"
241 trans_subset_top_expression_file_name="outputs/3_EQTLs/trans_subset_inputs/top/stabilized_expression.txt"
242 trans_subset_top_genes_file_name="outputs/3_EQTLs/trans_subset_inputs/top/coexpressed_genes.txt"
243 trans_subset_top_Genotypes_file_name="outputs/3_EQTLs/trans_subset_inputs/top/gtypes_ok.txt"
244
245
246 write.table(expression_trans_subset_top,trans_subset_top_expression_file_name, quote=F, sep="\t",
  row.names=TRUE)
247 write.table(genes_trans_subset_top,trans_subset_top_genes_file_name, quote=F, sep="\t",
  row.names=TRUE)
248 write.table(snpspos_trans_subset_top,trans_subset_top_SNP_positions_file_name, quote=F, sep="\t",
  row.names=TRUE)
249 write.table(gtypes_trans_subset_top,trans_subset_top_Genotypes_file_name, quote=F, sep="\t",
  row.names=TRUE)
250
251 ##### Now, Configure & Run Matrix EQTL
252
253 gene = SlicedData$new();
254 gene$fileDelimiter = "\t"; # the TAB character
255 gene$fileOmitCharacters = "NA"; # denote missing values;
256 gene$fileSkipRows = 1; # one row of column labels
257 gene$fileSkipColumns = 1; # one column of row labels
258 gene$fileSliceSize = 2000; # read file in slices of 2,000 rows
259 gene$LoadFile(trans_subset_top_expression_file_name);
260
261 ## Load covariates: there is none, so we declare an empty cov. file:
262

```

```

263 trans_subset_top_covariates_file_name=character()
264 cvrt = SlicedData$new();
265 cvrt$fileDelimiter = "\t"; # the TAB character
266 cvrt$fileOmitCharacters = "NA"; # denote missing values;
267 cvrt$fileSkipRows = 1; # one row of column labels
268 cvrt$fileSkipColumns = 1; # one column of row labels
269 if(length(covariates_file_name)>0) {
270   cvrt$LoadFile(trans_subset_top_covariates_file_name);
271 }
272
273 ## Load genotype data
274
275 snps = SlicedData$new();
276 snps$fileDelimiter = "\t"; # the TAB character
277 snps$fileOmitCharacters = "NA";
278 snps$fileOmitCharacters = "-9" # denote missing values;
279 snps$fileSkipRows = 1; # one row of column labels
280 snps$fileSkipColumns = 1; # one column of row labels
281 snps$fileSliceSize = 2000; # read file in slices of 2,000 rows
282 snps$LoadFile(trans_subset_top_Genotypes_file_name)
283
284 useModel = modelLINEAR
285 output_file_name_cis = tempfile()
286 pvOutputThreshold_cis = 1
287 pvOutputThreshold = 1;
288 errorCovariance = numeric()
289 cisDist = 1e5
290 output_file_name = tempfile()
291 output_file_name_cis = tempfile()
292
293
294
295 conjugated = Matrix_eQTL_main(snps = snps, gene = gene, cvrt = cvrt, output_file_name =
296   output_file_name, useModel = useModel, errorCovariance = errorCovariance, verbose = TRUE,
297   output_file_name.cis = output_file_name_cis, pvOutputThreshold = pvOutputThreshold,
298   pvOutputThreshold.cis = pvOutputThreshold_cis, snpspos = snpspos_trans_subset_top, genepos =
299   genes_trans_subset_top, cisDist = cisDist, pvalue.hist = "qqplot", min.pv.by.genesnp = TRUE,
300   noFDRsaveMemory = FALSE);
301
302 trans_eqtls_top=conjugated$trans$eqtls
303 #1942163 6
304 ## JQ por claridad
305 cis_eqtls_conjugated_top=conjugated$cis$eqtls
306 #[1] 3242 6
307 #so, we now have that 1942163+3242 (conjugated EQTLs labelled as trans plus cis (a minority,
308   obviously: cases here the SNP lies, by chance, nearby the conjugated gene)) = 815 * 2387
309   (genes times snps tested)
310
311 }
312
313

```

```

304 #####
305 ### 7. Run again in trans for conjugated EQTLs: algorithm balanced #####
306 #####
307
308 {
309
310   ## Load the table of cis-EQTLs involving genes with at least one hit (termed genes under genetic
      control).
311   cis=cis_eqtls[which(cis_eqtls$gene %in% top_cis_snps$gene),]
312   cis=cis[,c(1,2,11)]
313   colnames(cis)[3]="FDR"
314   cis=cis[order(cis$gene,cis$FDR),]
315
316   ## Load network
317   network=fread(paste0("outputs/2_coexpression/network0.01.txt"))
318   rownames(network)=network$V1
319   colnames(network)[1]="link"
320   red=network
321   red=red[,c("gen_row", "gen_col")]
322   N=3
323
324   ## Subset network links connecting pairs of genes under genetic control.
325   red=red[which(red$gen_row %in% top_cis_snps$gene),]
326   red=red[which(red$gen_col %in% top_cis_snps$gene),]
327
328
329   ## For each link, we will use 3 SNPs per gene that will be, in general link-specific (i.e. we will
      allow that the same gene uses different cis-EQTLs in different links if that helps to reach
      better significance-balance sets of SNPs). Here, I declare container table to store the
      identities and stats of these EQTLs.
330
331
332   final=data.frame(
333     gene_one=rep(NA,nrow(red)*N),gene_two=rep(NA,nrow(red)*N),
334     snp_one=rep(NA,nrow(red)*N),snp_two=rep(NA,nrow(red)*N),
335     FDR_one=rep(0,nrow(red)*N),FDR_two=rep(0,nrow(red)*N)
336   )
337
338   ## This loops runs the balancing algorithm
339
340   for(j in 1:nrow(red))
341   {
342     if(j%1000==0) print(j)
343     aux_row=cis[which(cis$gene == red$gen_row[j]),]
344     aux_row=aux_row[order(aux_row$FDR),]
345     aux_col=cis[which(cis$gene == red$gen_col[j]),]
346     aux_col=aux_col[order(aux_col$FDR),]
347     aux_row$stop=0
348     aux_col$stop=0
349
350     posicion_row=1
351     posicion_col=1
352     for(i in 1:N){
353       if(aux_row$FDR[posicion_row]<=aux_col$FDR[posicion_col])
354       {
355         guide_posicion=posicion_col
356         match_posicion=which.min(abs(-log10(aux_row$FDR[(posicion_row+1):nrow(aux_row)]) +
357           log10(aux_col$FDR[posicion_col])))+posicion_row
358         aux_col$stop[guide_posicion]=i
359         aux_row$stop[match_posicion]=i
360         posicion_col=guide_posicion+1
361         posicion_row=match_posicion+1
362       }else{

```



```

361         guide_posicion=posicion_row
362         match_posicion=which.min(abs(-log10(aux_col$FDR[(posicion_col+1):nrow(aux_col)]) +
363             log10(aux_row$FDR[posicion_row])))+posicion_col
364         aux_row$top[guide_posicion]=i
365         aux_col$top[match_posicion]=i
366         posicion_row=guide_posicion+1
367         posicion_col=match_posicion+1
368     }
369 }
370
371 aux_row=aux_row[which(aux_row$top>0),]
372 aux_col=aux_col[which(aux_col$top>0),]
373 append=cbind(aux_row,aux_col)[,c(2,6,1,5,3,7)]
374 colnames(append)=c("gene_one", "gene_two", "snp_one", "snp_two", "FDR_one", "FDR_two")
375
376 final[(N*j-N+1):(N*j)],<-append
377 }
378 system("mkdir -p outputs/3_EQTLs/trans_subset_inputs/match")
379 write.table(final,"outputs/3_EQTLs/trans_subset_inputs/match/SNPs_pairings.txt")
380
381 # Now, we gather all the SNPs appearing in the result table of the SNPs_balancing algorithm to then
382   run trans-EQTL mapping.
383
384 snpspos_trans_subset_match=snpspos[which(snp$pos %in%
385     unique(c(final$snp_one,final$snp_two))),c(1:3)]
386 gtypes_trans_subset_match=gtypes[which(rownames(gtypes) %in%
387     unique(c(final$snp_one,final$snp_two))),]
388 genes_trans_subset_match=genes[which(rownames(genes) %in% unique(c(final$gene_one,final$gene_two))),]
389 expression_trans_subset_match=expression[which(rownames(expression) %in%
390     unique(c(final$gene_one,final$gene_two))),]
391
392 ## 31571 unique SNPs involved (many appear several times), on 815 genes, 90 muestras
393
394
395 ## Check the congruence of these new arguments:
396 ## Sample-wise
397 length(which(colnames(expression_trans_subset_match)==colnames(gtypes_trans_subset_match)))/length
398   (colnames(expression_trans_subset_match))
399 ## Genes wise:
400 length(which(rownames(expression_trans_subset_match)==rownames(genes_trans_subset_match)))/length
401   (rownames(expression_trans_subset_match))
402 ## SNPs-wise (we could filter these out using the positions, but we will leave matrixEQTL do it)
403 length(which(rownames(snp$pos_trans_subset_match)==rownames(gtypes_trans_subset_match)))/length
404   (rownames(snp$pos_trans_subset_match))
405
406 system("mkdir -p outputs/3_EQTLs/trans_subset_inputs/match")
407
408 trans_subset_match_SNPs_positions_file_name="outputs/3_EQTLs/trans_subset_inputs/match/snp$pos_ok
409   .txt"
410 trans_subset_match_expression_file_name="outputs/3_EQTLs/trans_subset_inputs/match/stabilized_express
411   ion.txt"
412 trans_subset_match_genes_file_name="outputs/3_EQTLs/trans_subset_inputs/match/coexpressed_genes.txt"
413 trans_subset_match_Genotypes_file_name="outputs/3_EQTLs/trans_subset_inputs/match/gtypes_ok.txt"
414
415 write.table(expression_trans_subset_match,trans_subset_match_expression_file_name, quote=F,
416   sep="\t", row.names=TRUE)
417 write.table(genes_trans_subset_match,trans_subset_match_genes_file_name, quote=F, sep="\t",
418   row.names=TRUE)
419 write.table(snp$pos_trans_subset_match,trans_subset_match_SNPs_positions_file_name, quote=F,
420   sep="\t", row.names=TRUE)
421 write.table(gtypes_trans_subset_match,trans_subset_match_Genotypes_file_name, quote=F, sep="\t",
422   row.names=TRUE)

```

```

410 ##### Configure & Run Matrix EQTL #####
411
412 gene = SlicedData$new();
413 gene$fileDelimiter = "\t"; # the TAB character
414 gene$fileOmitCharacters = "NA"; # denote missing values;
415 gene$fileSkipRows = 1; # one row of column labels
416 gene$fileSkipColumns = 1; # one column of row labels
417 gene$fileSliceSize = 2000; # read file in slices of 2,000 rows
418 gene$LoadFile(trans_subset_match_expression_file_name);
419
420 ## Load covariates: there is none, so we declare an empty cov. file:
421
422 trans_subset_match_covariates_file_name=character()
423 cvrt = SlicedData$new();
424 cvrt$fileDelimiter = "\t"; # the TAB character
425 cvrt$fileOmitCharacters = "NA"; # denote missing values;
426 cvrt$fileSkipRows = 1; # one row of column labels
427 cvrt$fileSkipColumns = 1; # one column of row labels
428 if(length(covariates_file_name)>0) {
429   cvrt$LoadFile(trans_subset_match_covariates_file_name);
430 }
431
432 ## Load genotype data
433
434 snps = SlicedData$new();
435 snps$fileDelimiter = "\t"; # the TAB character
436 snps$fileOmitCharacters = "NA";
437 snps$fileOmitCharacters = "-9" # denote missing values;
438 snps$fileSkipRows = 1; # one row of column labels
439 snps$fileSkipColumns = 1; # one column of row labels
440 snps$fileSliceSize = 2000; # read file in slices of 2,000 rows
441 snps$LoadFile(trans_subset_match_Genotypes_file_name)
442
443
444 useModel = modelLINEAR
445 output_file_name_cis = tempfile()
446 pvOutputThreshold_cis = 1
447 pvOutputThreshold = 1;
448 errorCovariance = numeric()
449 cisDist = 1e5
450 output_file_name = tempfile()
451 output_file_name_cis = tempfile()
452
453 conjugated_match = Matrix_eQTL_main(snps = snps, gene = gene, cvrt = cvrt, output_file_name =
454   output_file_name, useModel = useModel, errorCovariance = errorCovariance, verbose = TRUE,
455   output_file_name.cis = output_file_name_cis, pvOutputThreshold = pvOutputThreshold,
456   pvOutputThreshold.cis = pvOutputThreshold_cis, snpspos = snpspos_trans_subset_match, genepos =
457   genes_trans_subset_match, cisDist = cisDist, pvalue.hist = "qqplot", min.pv.by.genesnp = TRUE,
458   noFDRsaveMemory = FALSE);
459
460 trans_eqtls_match=conjugated_match$trans$eqtls
461 # 25685872 6
462 cis_eqtls_conjugated_match=conjugated_match$cis$eqtls
463 # 44493 6
464
465 #so, we now have that 25685872 + 44493 (conjugated EQTLS labelled as trans plus cis (a minority,
466   obviously: cases here the SNP lies, by chance, nearby the conjugated gene)) = 815 * 31571 (genes
467   times snps tested)
468
469 }

```



```

462 #####
463 ### 8. Write output files #####
464 #####
465 {
466   system("mkdir -p outputs/3_EQTLs/results/match/")
467   system("mkdir -p outputs/3_EQTLs/results/top/")
468
469   write.table(cis_eqtls, file = "outputs/3_EQTLs/results/cis.txt")
470
471   write.table(top_cis_snps, file = "outputs/3_EQTLs/results/top/cis_top_N.txt")
472   write.table(trans_eqtls_top, file = "outputs/3_EQTLs/results/top/conjugated_trans_top.txt")
473   write.table(cis_eqtls_conjugated_top, file = "outputs/3_EQTLs/results/top/conjugated_cis_top.txt")
474
475   write.table(final, file = "outputs/3_EQTLs/results/match/cis_balanced_N.txt")
476   write.table(trans_eqtls_match, file = "outputs/3_EQTLs/results/match/conjugated_trans_match.txt")
477   write.table(cis_eqtls_conjugated_match, file =
478     "outputs/3_EQTLs/results/match/conjugated_cis_match.txt")
479
480 }
481 }

```

7.4. Mendelian Randomization for Top-SNPs method

```

1 #####
2 ### 1. Load dependencies #####
3 #####
4 {
5   library(edgeR)
6   library(limma)
7   library(gdsfmt)
8   library(SNPRelate)
9   library(qvalue)
10  library(stats)
11  library(MendelianRandomization)
12  library(data.table)
13  library(cobs)
14  library(ggplot2)
15  library(cowplot)
16
17  ul=function(tab,n=5)
18  {
19    rs=min(nrow(tab),n)
20    cs=min(ncol(tab),n)
21    return(tab[1:rs,1:cs])
22  }
23  ur=function(tab,n=5)
24  {
25    rs=min(nrow(tab),n)
26    cs=min(ncol(tab),n)
27    return(tab[1:rs,(ncol(tab)-cs):ncol(tab)])
28  }
29  dcols=function(tab){data.frame(colnames(tab))}
30  drows=function(tab){data.frame(rownames(tab))}
31 }

```

```

32 #####
33 ### 2. Load and pretty up input tables #####
34 #####
35
36 {
37     ## Input # 1: Network (here I just load it; the formatting requires an entire step, which also
38     needs the genes info table loaded here as well).
39     network=fread("outputs/2_coexpression/network0.01.txt")
40     colnames(network)[1]="link"
41     nrow(network)
42     # 25897510 pairs.
43     genes=read.table(paste0("outputs/2_coexpression/coexpressed_genes.txt"))
44
45     ## Input # 2: cis-EQTLs data. Load, rename FDR column, declare sd column
46
47     cis_hits_top=fread("outputs/3_EQTLs/results/top/cis_top_N.txt")
48     cis_hits_top=cis_hits_top[,c(2:7)]
49     cis_hits_top$sd=cis_hits_top$beta/cis_hits_top$statistic
50     colnames(cis_hits_top)[5]="FDR"
51     dim(cis_hits_top)
52     # 2445      6
53
54     ## Input # 3: conjugated-EQTLs data. Load cis and trans, rename FDR column, declare sd column,
55     type, and rbind
56
57     trans_hits_top=fread("outputs/3_EQTLs/results/top/conjugated_trans_top.txt")
58     trans_hits_top=trans_hits_top[,c(2:7)]
59     trans_hits_top=trans_hits_top[order(trans_hits_top$gene),]
60     dim(trans_hits_top)
61     # 1942163    6
62
63     trans_hits_top$type="trans"
64
65     conj_cis_hits_top=fread("outputs/3_EQTLs/results/top/conjugated_cis_top.txt")
66     conj_cis_hits_top=conj_cis_hits_top[,c(2:7)]
67     conj_cis_hits_top=conj_cis_hits_top[order(conj_cis_hits_top$gene),]
68     dim(conj_cis_hits_top)
69     # 3242      6
70     conj_cis_hits_top$type="cis"
71
72     conj_hits_top=rbind(trans_hits_top,conj_cis_hits_top)
73     conj_hits_top$sd=conj_hits_top$beta/conj_hits_top$statistic
74     colnames(conj_hits_top)[5]="FDR"
75     dim(conj_hits_top)
76
77     ## Input # 4: gtypes correlations.
78
79     gtypes=fread(paste0("inputs/processed/gtypes_ok.txt"))
80     gtypes=data.frame(gtypes)
81     rownames(gtypes)=gtypes$V1
82     gtypes=gtypes[,2:ncol(gtypes)]
83     gtypes_top=gtypes[which(rownames(gtypes) %in% unique(cis_hits_top$snps)),]
84     correlations_top=cor(t(gtypes_top))
85
86 }
87

```

```

88 #####
89 #### 3. Format network: subset only pairs under genetic control that are far enough from each other ####
90 #####
91
92 {
93   net_use_top=network[which(network$gen_row %in% cis_hits_top$gene & network$gen_col %in%
94     cis_hits_top$gene),]
95   dim(net_use_top)
96   ## [1] 131644      JQ: 9
97   rm(network)
98
99   vecinity_test=function(i,threshold=1000000){
100     gen_A=net_use_top$gen_row[i]
101     gen_B=net_use_top$gen_col[i]
102     chunk=genes[c(gen_A,gen_B),]
103     if(chunk$chr[1]!=chunk$chr[2])
104     {
105       return(0)
106     }else{
107       first_gene=which.min(chunk$start)
108       last_gene=which.max(chunk$start)
109       if(first_gene==last_gene){
110         return(1)
111       }else if(chunk$end[last_gene]<chunk$end[first_gene]){
112         return(1)
113       }else if(chunk$start[last_gene]-chunk$end[first_gene]<threshold){
114         return(1)
115       }else{
116         return(0)
117       }
118     }
119
120     net_use_top$vecinity=NA
121     net_use_top$vecinity=sapply(c(1:nrow(net_use_top)),vecinity_test)
122     dim(net_use_top)
123     #131644      10
124     net_use_top=net_use_top[which(net_use_top$vecinity!=1),]
125     dim(net_use_top)
126     #131230      10
127
128     ## Check to verify that the 815 genes survived:
129     length(unique(c(net_use_top$gen_row,net_use_top$gen_col)))
130     # 815
131
132   }
133
134   #####
135   #### 4. Checkpoint: save MR input objects ####
136   #####
137
138   {
139     system("mkdir -p outputs/4_MR/top")
140     MR_inputs=list(net_use_top=net_use_top,
141       cis_hits_top=cis_hits_top,conj_hits_top=conj_hits_top,correlations_top=correlations_top)
142     save(MR_inputs,file="outputs/4_MR/top/MR_inputs_checkpoint.Rdata")
143   }

```

```

144 #####
145 ### 5. Run MR #####
146 #####
147
148 {
149   net_infer_Mendel_Rand=function(coexp_top,cis_tests_top,conjugated_top,correl_top,i,N=3){
150     #if(i %% 200==0)
151     #print(i)
152     gen_a=coexp_top$gen_row[i]
153     gen_b=coexp_top$gen_col[i]
154
155     cis_chunk_fw=cis_tests_top[which(cis_tests_top$gene %in% gen_a),]
156     cis_chunk_bw=cis_tests_top[which(cis_tests_top$gene %in% gen_b),]
157
158     if(nrow(cis_chunk_fw)!=N)print(paste("Issue: cis_chunk_fw has", nrow(cis_chunk_fw)," tests"))
159     if(nrow(cis_chunk_bw)!=N)print(paste("Issue: cis_chunk_bw has", nrow(cis_chunk_bw)," tests"))
160
161     trans_chunk_fw=conjugated_top[which(conjugated_top$gene==gen_b & conjugated_top$snps %in%
162       cis_chunk_fw$snps),]
163     trans_chunk_bw=conjugated_top[which(conjugated_top$gene==gen_a & conjugated_top$snps %in%
164       cis_chunk_bw$snps),]
165
166     if(nrow(trans_chunk_fw)!=N)print(paste("Issue: trans_chunk_fw has", nrow(trans_chunk_fw),"
167       tests"))
168     if(nrow(trans_chunk_bw)!=N)print(paste("Issue: trans_chunk_bw has", nrow(trans_chunk_bw),"
169       tests"))
170
171     if(length(which(trans_chunk_fw$type=="cis")>0))print(paste("Issue: trans_chunk_fw includes
172       some cis tests"))
173
174     if(length(which(trans_chunk_bw$type=="cis")>0))print(paste("Issue: trans_chunk_bw includes
175       some cis tests"))
176
177     set_fw=which(colnames(correl_top) %in% cis_chunk_fw$snps)
178     set_bw=which(colnames(correl_top) %in% cis_chunk_bw$snps)
179
180     if(length(set_fw)!=N)print(paste("Issue:", length(set_fw)," SNPS in fw test in the corr.
181       matrix"))
182     if(length(set_bw)!=N)print(paste("Issue:", length(set_bw)," SNPS in bw test in the corr.
183       matrix"))
184
185     cors_fw=as.matrix(correl_top[set_fw,set_fw])
186     cors_bw=as.matrix(correl_top[set_bw,set_bw])
187
188     valid_cor_fw=length(cors_fw)>1
189     valid_cor_bw=length(cors_bw)>1
190
191     if(valid_cor_fw){
192       valid_cor_fw=abs(det(cors_fw))>1E-15
193     }
194     if(valid_cor_bw){
195       valid_cor_bw=abs(det(cors_bw))>1E-15
196     }
197
198     if(valid_cor_fw){
199       MR_input_fw <- mr_input(bx = cis_chunk_fw$beta,bxse = cis_chunk_fw$sd,by =
200         trans_chunk_fw$beta,byse = trans_chunk_fw$sd,corr=cors_fw)
201       MR_fw <- mr_ivw(MR_input_fw,correl = TRUE)
202     }else{
203       MR_input_fw <- mr_input(bx = cis_chunk_fw$beta,bxse = cis_chunk_fw$sd,by =
204         trans_chunk_fw$beta,byse = trans_chunk_fw$sd)

```

```

195     MR_fw <- mr_ivw(MR_input_fw,correl = FALSE)
196   }
197
198   if(valid_cor_bw){
199     MR_input_bw <- mr_input(bx = cis_chunk_bw$beta,bxse = cis_chunk_bw$sd,by =
      trans_chunk_bw$beta,byse = trans_chunk_bw$sd,corr=cors_bw)
200     MR_bw <- mr_ivw(MR_input_bw,correl = TRUE)
201   }else{
202     MR_input_bw <- mr_input(bx = cis_chunk_bw$beta,bxse = cis_chunk_bw$sd,by =
      trans_chunk_bw$beta,byse = trans_chunk_bw$sd)
203     MR_bw <- mr_ivw(MR_input_bw,correl = FALSE)
204   }
205
206   output_fw=c(MR_fw@Estimate,MR_fw@StdError,MR_fw@CILower,MR_fw@CIUpper,MR_fw@Pvalue)
207   output_bw=c(MR_bw@Estimate,MR_bw@StdError,MR_bw@CILower,MR_bw@CIUpper,MR_bw@Pvalue)
208   return(c(output_fw,output_bw))
209 }
210
211 #thing=load("outputs/4_MR/top/MR_inputs_checkpoint.Rdata")
212 #net_use_top=MR_inputs[["net_use_top"]]
213 #cis_hits_top=MR_inputs[["cis_hits_top"]]
214 #conj_hits_top=MR_inputs[["conj_hits_top"]]
215 #correlations_top=MR_inputs[["correlations_top"]]
216 #rm(MR_inputs)
217
218
219 system("mkdir -p outputs/4_MR/top/Result_chunks/")
220 for(j in 1:130){
221   print(paste("CHUNK: ",j,"\n"))
222   name=paste0("outputs/4_MR/top/Result_chunks/chunk_",j,".txt")
223   down=1+(j-1)*1000
224   up=j*1000
225   res=sapply(c(down:up),
226     net_infer_Mendel_Rand,
227     coexp_top=net_use_top,
228     cis_tests_top=cis_hits_top,
229     conjugated_top=conj_hits_top,
230     correl_top=correlations_top)
231   res=data.frame(t(res))
232   write.table(res,name)
233 }
234
235 name=paste0("outputs/4_MR/top/Result_chunks/chunk_131.txt")
236 down=130001
237 up=nrow(net_use_top)
238 res=sapply(c
  (down:up),net_infer_Mendel_Rand,coexp_top=net_use_top,cis_tests_top=cis_hits_top
  ,conjugated_top=conj_hits_top,correl_top=correlations_top)
239 res=data.frame(t(res))
240 write.table(res,name)
241
242 result_top=read.table(paste0("outputs/4_MR/top/Result_chunks/chunk_1.txt"))
243

```

```

243
244 for(j in 2:131)
245 {
246   print(paste("CHUNK: ",j,"\n"))
247   name=paste0("outputs/4_MR/top/Result_chunks/chunk_",j,".txt")
248   appendix_top=read.table(name)
249   result_top=rbind(result_top,appendix_top)
250 }
251
252 colnames(result_top)=c(
253   "Theta_FW", "Sd_FW", "LowerCI_FW", "UpperCI_FW", "P_FW",
254   "Theta_BW", "Sd_BW", "LowerCI_BW", "UpperCI_BW", "P_BW")
255
256 result_top$gen_one=net_use_top$gen_row
257 result_top$gen_two=net_use_top$gen_col
258 dim(result_top)
259 #131230      10
260
261 result_top$index=c(1:nrow(result_top))
262 result_top=result_top[order(-result_top$P_FW),]
263 result_top$n_larger=c(1:nrow(result_top))
264 result_top$pi_0_gorro_FW=result_top$n_larger/((nrow(result_top))*(1-result_top$P_FW))
265 result_top=result_top[order(result_top$P_FW),]
266 constraint_matrix=as.matrix(data.frame(c(0,2),c(0,1),c(1,0)))
267 result_top$pi_0_gorro_smooth_FW=cobs(result_top$P_FW,
268   result_top$pi_0_gorro_FW,
269   constraint="decrease",
270   nknots=14,method="quantile",
271   pointwise=constraint_matrix,maxiter=500,print.warn=TRUE,print.mesg=TRUE)$fitted
272
273
274 pi_0_FW=min(result_top$pi_0_gorro_smooth_FW)
275
276 result_top$Fdr_FW=p.adjust(result_top$P_FW,method="BH")*pi_0_FW
277 result_top$Fndr_FW=1-(pi_0_FW/result_top$pi_0_gorro_smooth_FW)
278
279 pl_FW=ggplot(result_top)+
280   geom_point(aes(x=P_FW,y=pi_0_gorro_FW),size=0.1,color="blue",alpha=0.3)+
281   geom_line(aes(x=P_FW,y=pi_0_gorro_smooth_FW),color="red")+
282   geom_line(aes(x=P_FW,y=Fndr_FW),color="green")+
283   geom_line(aes(x=P_FW,y=Fdr_FW),color="black")+
284   geom_hline(yintercept=pi_0_FW)+ylim(0,1)
285
286 result_top=result_top[order(-result_top$P_BW),]
287 result_top$n_larger=c(1:nrow(result_top))
288 result_top$pi_0_gorro_BW=result_top$n_larger/((nrow(result_top))*(1-result_top$P_BW))
289 result_top=result_top[order(result_top$P_BW),]
290 constraint_matrix=as.matrix(data.frame(c(0,2),c(0,1),c(1,0)))
291 result_top$pi_0_gorro_smooth_BW=cobs
292   (result_top$P_BW,result_top$pi_0_gorro_BW,constraint="decrease",nknots=14,method="quantile"
293   ,pointwise=constraint_matrix,maxiter=1000,print.warn=FALSE,print.mesg=FALSE)$fitted
294 pi_0_BW=min(result_top$pi_0_gorro_smooth_BW)
295
296 result_top$Fdr_BW=p.adjust(result_top$P_BW,method="BH")*pi_0_BW
297 result_top$Fndr_BW=1-(pi_0_BW/result_top$pi_0_gorro_smooth_BW)
298
299 pl_BW=ggplot(result_top)+ geom_point(aes(x=P_BW,y=pi_0_gorro_BW),size=0.1,color="blue",alpha=0.3)
300   + geom_line(aes(x=P_BW,y=pi_0_gorro_smooth_BW),color="red") +
301   geom_line(aes(x=P_BW,y=Fndr_BW),color="green") + geom_line(aes(x=P_BW,y=Fdr_BW),color="black")
302   + geom_hline(yintercept=pi_0_BW)+ylim(0,1)
303
304 pl_fdrs_fndrs=plot_grid(pl_FW,pl_BW,ncol=2)

```

```

299
300 pdf("outputs/4_MR/top/pl_fdrs_fndrs.pdf",height=4,width=8)
301 print(pl_fdrs_fndrs)
302 dev.off()
303
304 result_top=result_top[order(result_top$index),]
305 result_top=result_top[,c(11,12,1:5,15:18,6:10,19:22)]
306 result_top$link=paste0(result_top$gen_one,"_",result_top$gen_two)
307
308 }
309
310 #####
311 ### 6. Clasify link directions #####
312 #####
313
314 {
315   get_mean_log10Fdr=function(i){
316     if(i%%1000==0)print(i)
317     fdrs_one=cis_hits_top$FDR[which(cis_hits_top$gene %in% result_top$gen_one[i])]
318     fdrs_two=cis_hits_top$FDR[which(cis_hits_top$gene %in% result_top$gen_two[i])]
319     output=c(mean(-log10(fdrs_one)),mean(-log10(fdrs_two)))
320     return(output)
321   }
322
323   means_log10=t(sapply(1:nrow(result_top),get_mean_log10Fdr))
324
325   result_top$mean_log_fdr_one=means_log10[,1]
326   result_top$mean_log_fdr_two=means_log10[,2]
327
328
329   get_min_Fdr=function(i)
330   {
331     if(i%%1000==0)print(i)
332     fdrs_one=cis_hits_top$FDR[which(cis_hits_top$gene %in% result_top$gen_one[i])]
333     fdrs_two=cis_hits_top$FDR[which(cis_hits_top$gene %in% result_top$gen_two[i])]
334     output=c(max(-log10(fdrs_one)),max(-log10(fdrs_two)))
335     return(output)
336   }
337
338   maxs_log10=t(sapply(1:nrow(result_top),get_min_Fdr))
339
340   result_top$max_log_fdr_one=maxs_log10[,1]
341   result_top$max_log_fdr_two=maxs_log10[,2]
342
343   ## One_more_than_two (Two_more_than_one) means gene one (two) has more significant cis EQTLs than
344   gene two (one)
345   result_top$flag_direction_mean="One_more_than_two"
346   result_top$flag_direction_mean
347   [which(result_top$mean_log_fdr_one<result_top$mean_log_fdr_two)]= "Two_more_than_one"
348
349   ## One_more_than_two (Two_more_than_one) means gene one (two) more significant cis EQTLs than gene
350   two (one)
351   result_top$flag_direction_max="One_more_than_two"
352   result_top$flag_direction_max
353   [which(result_top$max_log_fdr_one<result_top$max_log_fdr_two)]= "Two_more_than_one"
354
355   th_fndr=0.1
356   th_fdr=0.1
357   ## Count FW links from one to two
358   FW_1=length(which(result_top$flag_direction_mean=="One_more_than_two" & result_top$Fdr_FW<th_fdr &
359     result_top$Fndr_BW<th_fndr))

```

```

354     ## Count FW links from two to one
355     FW_2=length(which(result_top$flag_direction_mean=="Two_more_than_one" & result_top$Fdr_BW<th_fdr &
      result_top$Fndr_FW<th_fndr))
356
357     ## Count BW links from one to two
358     BW_1=length(which(result_top$flag_direction_mean=="One_more_than_two" & result_top$Fdr_BW<th_fdr &
      result_top$Fndr_FW<th_fndr))
359     ## Count BW links from two to one
360     BW_2=length(which(result_top$flag_direction_mean=="Two_more_than_one" & result_top$Fdr_FW<th_fdr &
      result_top$Fndr_BW<th_fndr))
361
362     FW_links=FW_1+FW_2
363     BW_links=BW_1+BW_2
364
365     FW_links/(FW_links+BW_links)
366     # 0.4190195
367
368     ## Ahora, a ver el criterio max:
369
370     ## Count FW links from one to two
371     FW_1=length(which(result_top$flag_direction_max=="One_more_than_two" & result_top$Fdr_FW<th_fdr &
      result_top$Fndr_BW<th_fndr))
372     ## Count FW links from two to one
373     FW_2=length(which(result_top$flag_direction_max=="Two_more_than_one" & result_top$Fdr_BW<th_fdr &
      result_top$Fndr_FW<th_fndr))
374
375     ## Count BW links from one to two
376     BW_1=length(which(result_top$flag_direction_max=="One_more_than_two" & result_top$Fdr_BW<th_fdr &
      result_top$Fndr_FW<th_fndr))
377     ## Count BW links from two to one
378
379
380     BW_2=length(which(result_top$flag_direction_max=="Two_more_than_one" & result_top$Fdr_FW<th_fdr &
      result_top$Fndr_BW<th_fndr))
379
380     FW_links=FW_1+FW_2
381     BW_links=BW_1+BW_2
382
383     FW_links/(FW_links+BW_links)
384     # 0.4566005
385
386
387     write.table(result_top,"outputs/4_MR/top/result.txt")
388
389 }

```


7.5. Mendelian Randomization for Balanced-SNPs method

```

1 #####
2 #### 1. Load dependencies #####
3 #####
4 {
5   library(edgeR)
6   library(limma)
7   library(gdsfmt)
8   library(SNPrelate)
9   library(qvalue)
10  library(stats)
11  library(MendelianRandomization)
12  library(data.table)
13  library(cobs)
14  library(ggplot2)
15  library(cowplot)
16
17  ul=function(tab,n=5)
18  {
19    rs=min(nrow(tab),n)
20    cs=min(ncol(tab),n)
21    return(tab[1:rs,1:cs])
22  }
23  ur=function(tab,n=5)
24  {
25    rs=min(nrow(tab),n)
26    cs=min(ncol(tab),n)
27    return(tab[1:rs,(ncol(tab)-cs):ncol(tab)])
28  }
29  dcols=function(tab){data.frame(colnames(tab))}
30  drows=function(tab){data.frame(rownames(tab))}
31 }

32 #####
33 #### 2. Load and pretty up input tables #####
34 #####
35
36 {
37   ## Input # 1: Network (I just load it; the formatting requires an entire step, which also needs
38   ## the genes info table loaded here as well).
39   network=fread("outputs/2_coexpression/network0.01.txt")
40   colnames(network)[1]="link"
41   nrow(network)
42   # 25897510 pairs.
43   genes=read.table(paste0("outputs/2_coexpression/coexpressed_genes.txt"))
44
45   ## Input # 2: cis-EQTLs data. Here, in cis_hits_balanced we only stored
46   ## the tests of the IDs that will be paired.
47   ## But the tests stats are in cis.txt.
48   ## I merge them to build the sd columns, needed for MR
49   cis_hits_balanced=fread("outputs/3_EQTLs/results/match/cis_balanced_N.txt")
50   cis_hits_balanced=cis_hits_balanced[,c(2:7)]
51
52   cis_detailed_stats=fread("outputs/3_EQTLs/results/cis.txt")
53   cis_detailed_stats=cis_detailed_stats[,c(2:5,7)]
54   colnames(cis_detailed_stats)[1]="snp"
55
56   cis_detailed_stats$test=paste0(cis_detailed_stats$snp,"_",cis_detailed_stats$gene)
57
58   cis_hits_balanced$test=paste0(cis_hits_balanced$snp_one,"_",cis_hits_balanced$gene_one)
59   cis_hits_balanced_output=merge(cis_hits_balanced,cis_detailed_stats,by="test")
60   cis_hits_balanced_output=cis_hits_balanced_output[,c(2:7,10:12)]
61   colnames(cis_hits_balanced_output)[c(7:9)]=paste0(colnames(cis_hits_balanced_output)
62     [c(7:9)],"_one")

```

```

62
63 cis_hits_balanced_output$test=paste0
64   (cis_hits_balanced_output$snp_two,"_",cis_hits_balanced_output$gene_two)
65 cis_hits_balanced_output=merge(cis_hits_balanced_output,cis_detailed_stats,by="test")
66 cis_hits_balanced_output=cis_hits_balanced_output[,c(4,2,10,8,9,6,5,3,15,13,14,7)]
67
68 colnames(cis_hits_balanced_output)[c(9:11)]=paste0(colnames(cis_hits_balanced_output)
69   [c(9:11)],"_two")
70 cis_hits_balanced_output$link=paste0
71   (cis_hits_balanced_output$gene_one,"_",cis_hits_balanced_output$gene_two)
72 dim(cis_hits_balanced_output)
73 # 394932 7
74 cis_hits_balanced=cis_hits_balanced_output
75 rm(cis_hits_balanced_output)
76 cis_hits_balanced$sd_one=cis_hits_balanced$beta_one/cis_hits_balanced$statistic_one
77 cis_hits_balanced$sd_two=cis_hits_balanced$beta_two/cis_hits_balanced$statistic_two
78
79 ## Input # 3: conjugated-EQTLs data. Load cis and trans, rename FDR column, declare sd column,
80   type, and rbind
81 trans_hits_balanced=fread("outputs/3_EQTLs/results/match/conjugated_trans_match.txt")
82 trans_hits_balanced=trans_hits_balanced[,c(2:7)]
83 trans_hits_balanced=trans_hits_balanced[order(trans_hits_balanced$gene),]
84 dim(trans_hits_balanced)
85 # 25685872 6
86 trans_hits_balanced$type="trans"
87
88 conj_cis_hits_balanced=fread("outputs/3_EQTLs/results/match/conjugated_cis_match.txt")
89 conj_cis_hits_balanced=conj_cis_hits_balanced[,c(2:7)]
90 conj_cis_hits_balanced=conj_cis_hits_balanced[order(conj_cis_hits_balanced$gene),]
91 dim(conj_cis_hits_balanced)
92 # 44493 6
93 conj_cis_hits_balanced$type="cis"
94
95
96
97 conj_hits_balanced=rbind(trans_hits_balanced,conj_cis_hits_balanced)
98 conj_hits_balanced$sd=conj_hits_balanced$beta/conj_hits_balanced$statistic
99 colnames(conj_hits_balanced)[5]="FDR"
100 dim(conj_hits_balanced)
101 # 25730365 8
102
103 ## Input # 4: gtypes correlations.
104 gtypes=fread(paste0("inputs/processed/gtypes_ok.txt"))
105 gtypes=data.frame(gtypes)
106 rownames(gtypes)=gtypes$V1
107 gtypes=gtypes[,2:ncol(gtypes)]
108 gtypes_balanced=gtypes[which(rownames(gtypes) %in%
109   unique(c(cis_hits_balanced$snp_one,cis_hits_balanced$snp_two))),]
110 correlations_balanced=cor(t(gtypes_balanced))
111 }

```

```

105 #####
106 3. Format network: subset only pairs under genetic control that are far enough from each other ###
107 #####
108 {
109     ## This command to subset the network is different from the balanced case: different info in
110     dataframes
111
112     net_use_balanced=network[which(network$link %in% cis_hits_balanced$link),]
113     dim(net_use_balanced)
114     ## [1] 131644      JQ: 9
115     rm(network)
116
117     vecinity_test=function(i,threshold=1000000){
118         gen_A=net_use_balanced$gen_row[i]
119         gen_B=net_use_balanced$gen_col[i]
120         chunk=genes[c(gen_A,gen_B),]
121         if(chunk$chr[1]!=chunk$chr[2])
122         {
123             return(0)
124         }else{
125             first_gene=which.min(chunk$start)
126             last_gene=which.max(chunk$start)
127             if(first_gene==last_gene){
128                 return(1)
129             }else if(chunk$end[last_gene]<chunk$end[first_gene]){
130                 return(1)
131             }else if(chunk$start[last_gene]-chunk$end[first_gene]<threshold){
132                 return(1)
133             }else{
134                 return(0)
135             }
136         }
137     }
138
139     net_use_balanced$vecinity=NA
140     net_use_balanced$vecinity=apply(c(1:nrow(net_use_balanced)),vecinity_test)
141     dim(net_use_balanced)
142     #131644      10
143     net_use_balanced=net_use_balanced[which(net_use_balanced$vecinity!=1),]
144     dim(net_use_balanced)
145     #131230      10
146
147     ## Check to verify that the 815 genes survived:
148     length(unique(c(net_use_balanced$gen_row,net_use_balanced$gen_col)))
149     # 815
150 }

```

```

150 #####
151 ### 4. Filter elements from cis_hits_balanced and conj_hits_balanced that will not be used. ###
152 #####
153 {
154     cis_hits_balanced=cis_hits_balanced[which(cis_hits_balanced$link %in% net_use_balanced$link),]
155     nrow(cis_hits_balanced)
156     #393690
157
158     cis_hits_balanced$trans_test_fw=paste0(cis_hits_balanced$snp_one,"_",cis_hits_balanced$gene_two)
159     cis_hits_balanced$trans_test_bw=paste0(cis_hits_balanced$snp_two,"_",cis_hits_balanced$gene_one)
160
161     conj_hits_balanced$test=paste0(conj_hits_balanced$snps,"_",conj_hits_balanced$gene)
162     nrow(conj_hits_balanced)
163     # 25730365
164     conj_hits_balanced=conj_hits_balanced[which(conj_hits_balanced$test %in%
165         unique(c(cis_hits_balanced$trans_test_fw,cis_hits_balanced$trans_test_bw))),]
166     dim(conj_hits_balanced)
167     # [1] 776489      9
168     summary(factor(conj_hits_balanced$type))
169     ## trans
170     ## 776489
171     ## All remaining trans tests are trans (we erased the cis ones during the vicinity filtering)
172 }
173
174 #####
175 ### 5. Checkpoint: save MR input objects ###
176 #####
177
178 system("mkdir -p outputs/4_MR/match/")
179 MR_inputs=list(net_use_balanced=net_use_balanced,
180     cis_hits_balanced=cis_hits_balanced,conj_hits_balanced=conj_hits_balanced
181     ,correlations_balanced=correlations_balanced)
182 save(MR_inputs,file="outputs/4_MR/match/MR_inputs_checkpoint_match_fw.Rdata")
183
184 #####
185 ### 6. Run MR #####
186 #####
187
188 net_infer_Mendel_Rand=function
189 (coexp_balanced,cis_tests_balanced,conjugated_balanced,correl_balanced,i,N=3)
190 {
191     if(i%%100==0) print(i)
192     gen_a=coexp_balanced$gen_row[i]
193     gen_b=coexp_balanced$gen_col[i]
194     link_ref=paste0(gen_a,"_",gen_b)
195
196     cis_chunk=cis_tests_balanced[which(cis_tests_balanced$link == link_ref),]
197     if(nrow(cis_chunk_fw)!=N)print(paste("Issue: cis_chunk_fw has", nrow(cis_chunk_fw)," tests"))
198
199     cis_chunk_fw=cis_chunk[,c("snp_one","gene_one","beta_one","sd_one")]
200     cis_chunk_bw=cis_chunk[,c("snp_two","gene_two","beta_two","sd_two")]
201
202     colnames(cis_chunk_fw)=c("snps","gene","beta","sd")
203     colnames(cis_chunk_bw)=c("snps","gene","beta","sd")

```

```

203 trans_chunk_fw=conjugated_balanced[which(conjugated_balanced$gene==gen_b &
204 conjugated_balanced$snps %in% cis_chunk_fw$snps),]
205 trans_chunk_bw=conjugated_balanced[which(conjugated_balanced$gene==gen_a &
206 conjugated_balanced$snps %in% cis_chunk_bw$snps),]
207
208 if(nrow(trans_chunk_fw)!=N)print(paste("Issue: trans_chunk_fw has", nrow(trans_chunk_fw),"
209 tests"))
210 if(nrow(trans_chunk_bw)!=N)print(paste("Issue: trans_chunk_bw has", nrow(trans_chunk_bw),"
211 tests"))
212
213 set_fw=which(colnames(correl_balanced) %in% cis_chunk_fw$snps)
214 set_bw=which(colnames(correl_balanced) %in% cis_chunk_bw$snps)
215
216 if(length(set_fw)!=N)print(paste("Issue:", length(set_fw)," SNPS in fw test in the corr. matrix.
217 Iter",i))
218 if(length(set_bw)!=N)print(paste("Issue:", length(set_bw)," SNPS in bw test in the corr. matrix.
219 Iter",i))
220
221 cors_fw=as.matrix(correl_balanced[set_fw,set_fw])
222 cors_bw=as.matrix(correl_balanced[set_bw,set_bw])
223
224 valid_cor_fw=length(cors_fw)>1
225 valid_cor_bw=length(cors_bw)>1
226
227 if(valid_cor_fw){
228   valid_cor_fw=abs(det(cors_fw))>1E-15
229 }
230 if(valid_cor_bw){
231   valid_cor_bw=abs(det(cors_bw))>1E-15
232 }
233
234 if(valid_cor_fw){
235   MR_input_fw <- mr_input(bx = cis_chunk_fw$beta,bxse = cis_chunk_fw$sd,by =
236 trans_chunk_fw$beta,byse = trans_chunk_fw$sd,corr=cors_fw)
237 MR_fw <- mr_ivw(MR_input_fw,correl = TRUE)
238 }else{
239   MR_input_fw <- mr_input(bx = cis_chunk_fw$beta,bxse = cis_chunk_fw$sd,by =
240 trans_chunk_fw$beta,byse = trans_chunk_fw$sd)
241 MR_fw <- mr_ivw(MR_input_fw,correl = FALSE)
242 }
243
244 if(valid_cor_bw){
245   MR_input_bw <- mr_input(bx = cis_chunk_bw$beta,bxse = cis_chunk_bw$sd,by =
246 trans_chunk_bw$beta,byse = trans_chunk_bw$sd,corr=cors_bw)
247 MR_bw <- mr_ivw(MR_input_bw,correl = TRUE)
248 }else{
249   MR_input_bw <- mr_input(bx = cis_chunk_bw$beta,bxse = cis_chunk_bw$sd,by =
250 trans_chunk_bw$beta,byse = trans_chunk_bw$sd)
251 MR_bw <- mr_ivw(MR_input_bw,correl = FALSE)
252 }
253
254 output_fw=c(MR_fw@Estimate,MR_fw@StdError,MR_fw@CILower,MR_fw@CIUpper,MR_fw@Pvalue)
255 output_bw=c(MR_bw@Estimate,MR_bw@StdError,MR_bw@CILower,MR_bw@CIUpper,MR_bw@Pvalue)
256 return(c(output_fw,output_bw))
257 }
258
259 ##
260 #thing=load("outputs/4_MR/match/MR_inputs_checkpoint_match_fw.Rdata")
261 #net_use_balanced=MR_inputs[["net_use_balanced"]]
262 #cis_hits_balanced=MR_inputs[["cis_hits_balanced"]]
263 #conj_hits_balanced=MR_inputs[["conj_hits_balanced"]]
264 #correlations_balanced=MR_inputs[["correlations_balanced"]]
265 #rm(MR_inputs)

```

```

257
258 start_point=1
259 end_point=130
260
261 system("mkdir -p outputs/4_MR/match/Result_chunks_fv/")
262 for(j in start_point:end_point){
263   print(paste("CHUNK: ",j,"\n"))
264   name=paste0("outputs/4_MR/match/Result_chunks_fv/chunk_",j,".txt")
265   down=1+(j-1)*1000
266   up=j*1000
267   res=sapply(c(down:up),
268     net_infer_Mendel_Rand,
269     coexp_balanced=net_use_balanced,
270     cis_tests_balanced=cis_hits_balanced,
271     conjugated_balanced=conj_hits_balanced,
272     correl_balanced=correlations_balanced)
273
274   res=data.frame(t(res))
275   write.table(res,name)
276 }
277
278 name=paste0("outputs/4_MR/match/Result_chunks_fv/chunk_131.txt")
279 down=130001
280 up=nrow(net_use_balanced)
281 res=sapply(c
  (down:up),net_infer_Mendel_Rand,coexp_balanced=net_use_balanced
  ,cis_tests_balanced=cis_hits_balanced,conjugated_balanced=conj_hits_balanced
  ,correl_balanced=correlations_balanced)
282 res=data.frame(t(res))
283 write.table(res,name)
284
285 result_balanced=read.table(paste0("outputs/4_MR/match/Result_chunks_fv/chunk_1.txt"))
286
287 for(j in 2:131)
288 {
289   print(paste("CHUNK: ",j,"\n"))
290   name=paste0("outputs/4_MR/match/Result_chunks_fv/chunk_",j,".txt")
291   appendix_balanced=read.table(name)
292   result_balanced=rbind(result_balanced,appendix_balanced)
293 }
294
295 colnames(result_balanced)=c(
296 "Theta_FW", "Sd_FW", "LowerCI_FW", "UpperCI_FW", "P_FW",
297 "Theta_BW", "Sd_BW", "LowerCI_BW", "UpperCI_BW", "P_BW")
298
299 result_balanced$gen_one=net_use_balanced$gen_row
300 result_balanced$gen_two=net_use_balanced$gen_col
301 dim(result_balanced)
302 #131230      12
303
304 #####
305 #### 7. Clasify link directions #####
306 #####
307
308 result_balanced$index=c(1:nrow(result_balanced))
309 result_balanced=result_balanced[order(-result_balanced$P_FW),]
310 result_balanced$n_larger=c(1:nrow(result_balanced))
311 result_balanced$pi_0_gorro_FW=result_balanced$n_larger/((nrow(result_balanced))*
  (1-result_balanced$P_FW))
312 result_balanced=result_balanced[order(result_balanced$P_FW),]
313 constraint_matrix=as.matrix(data.frame(c(0,2),c(0,1),c(1,0)))
314 result_balanced$pi_0_gorro_smooth_FW=cobs
  (result_balanced$P_FW,result_balanced$pi_0_gorro_FW,constraint="decrease",nknots=14
  ,method="quantile",pointwise=constraint_matrix,maxiter=1000,print.warn=FALSE,print
  .msg=FALSE)$fitted
315

```

```

316 pi_0_FW=min(result_balanced$pi_0_gorro_smooth_FW)
317
318 result_balanced$Fdr_FW=p.adjust(result_balanced$P_FW,method="BH")*pi_0_FW
319 result_balanced$Fndr_FW=1-(pi_0_FW/result_balanced$pi_0_gorro_smooth_FW)
320
321 pl_FW=ggplot(result_balanced)+
  geom_point(aes(x=P_FW,y=pi_0_gorro_FW),size=0.1,color="blue",alpha=0.3) +
  geom_line(aes(x=P_FW,y=pi_0_gorro_smooth_FW),color="red") +
  geom_line(aes(x=P_FW,y=Fndr_FW),color="green")+ geom_line(aes(x=P_FW,y=Fdr_FW),color="black") +
  geom_hline(yintercept=pi_0_FW)+geom_hline(yintercept=0.1,color="green")+ylim(0,1)
322
323 result_balanced=result_balanced[order(-result_balanced$P_BW),]
324 result_balanced$n_larger=c(1:nrow(result_balanced))
325 result_balanced$pi_0_gorro_BW=result_balanced$n_larger/((nrow(result_balanced))*
  (1-result_balanced$P_BW))
326 result_balanced=result_balanced[order(result_balanced$P_BW),]
327 constraint_matrix=as.matrix(data.frame(c(0,2),c(0,1),c(1,0)))
328 result_balanced$pi_0_gorro_smooth_BW=cobs
  (result_balanced$P_BW,result_balanced$pi_0_gorro_BW,constraint="decrease",nknots=14
  ,method="quantile",pointwise=constraint_matrix,maxiter=1000,print.warn=FALSE,print
  .msg=FALSE)$fitted
329
330 pi_0_BW=min(result_balanced$pi_0_gorro_smooth_BW)
331
332 result_balanced$Fdr_BW=p.adjust(result_balanced$P_BW,method="BH")*pi_0_BW
333 result_balanced$Fndr_BW=1-(pi_0_BW/result_balanced$pi_0_gorro_smooth_BW)
334
335 pl_BW=ggplot(result_balanced)+
  geom_point(aes(x=P_BW,y=pi_0_gorro_BW),size=0.1,color="blue",alpha=0.3)+
  geom_line(aes(x=P_BW,y=pi_0_gorro_smooth_BW),color="red")+
  geom_line(aes(x=P_BW,y=Fndr_BW),color="green")+ geom_line(aes(x=P_BW,y=Fdr_BW),color="black") +
  geom_hline(yintercept=pi_0_BW)+geom_hline(yintercept=0.1,color="green")+ylim(0,1)
336
337
338 pl_fdrs_fndrs=plot_grid(pl_FW,pl_BW,ncol=2)
339
340 pdf("outputs/4_MR/match/pl_fdrs_fndrs_fv.pdf",height=4,width=8)
341 print(pl_fdrs_fndrs)
342 dev.off()
343
344 pl_corners=ggplot(result_balanced)+geom_point(aes(x=-log10(Fdr_BW),y=-log10(Fdr_FW)))
345 ## Este plot mola tanto que casi hace desconfiar :DD
346 result_balanced=result_balanced[order(result_balanced$index),]
347 result_balanced=result_balanced[,c(11,12,1:5,15:18,6:10,19:22)]
348 result_balanced$link=paste0(result_balanced$gen_one,"_",result_balanced$gen_two)
349
350 get_mean_log10Fdr=function(i)
351 {
352   if(i%%1000==0)print(i)
353   set=which(cis_hits_balanced$link %in% result_balanced$link[i])
354   fdrs_one=cis_hits_balanced$FDR_one[set]
355   fdrs_two=cis_hits_balanced$FDR_two[set]
356   output=c(mean(-log10(fdrs_one)),mean(-log10(fdrs_two)))
357   return(output)
358 }
359
360 means_log10=t(sapply(1:nrow(result_balanced),get_mean_log10Fdr))
361
362 result_balanced$mean_log_fdr_one=means_log10[,1]
363 result_balanced$mean_log_fdr_two=means_log10[,2]
364
365 logfdr1=-log10(cis_hits_balanced$FDR_one)
366 logfdr2=-log10(cis_hits_balanced$FDR_two)
367

```



```

368 get_min_Fdr=function(i)
369 {
370     if(i%%1000==0)print(i)
371     set=which(cis_hits_balanced$link %in% result_balanced$link[i])
372     output=c(max(logfdr1[set]),max(logfdr2[set]))
373     return(output)
374 }
375
376 maxs_log10=t(sapply(1:nrow(result_balanced),get_min_Fdr))
377
378 result_balanced$max_log_fdr_one=maxs_log10[,1]
379 result_balanced$max_log_fdr_two=maxs_log10[,2]
380
381
382 ## One_more_than_two (Two_more_than_one) means gene one (two) has more significant cis EQTLs than
    gene two (one)
383 result_balanced$flag_direction_mean="One_more_than_two"
384 result_balanced$flag_direction_mean
    [which(result_balanced$mean_log_fdr_one<result_balanced$mean_log_fdr_two)]= "Two_more_than_one"
385
386 ## One_more_than_two (Two_more_than_one) means gene one (two) more significant cis EQTLs than gene
    two (one)
387 result_balanced$flag_direction_max="One_more_than_two"
388 result_balanced$flag_direction_max
    [which(result_balanced$max_log_fdr_one<result_balanced$max_log_fdr_two)]= "Two_more_than_one"
389
390 th_fndr=0.1
391 th_fdr=0.1
392 ## Count FW links from one to two
393 FW_1=length(which(result_balanced$flag_direction_mean=="One_more_than_two" &
    result_balanced$Fdr_FW<th_fdr & result_balanced$Fndr_BW<th_fndr))
394 ## Count FW links from two to one
395
396
397 FW_2=length(which(result_balanced$flag_direction_mean=="Two_more_than_one" &
    result_balanced$Fdr_BW<th_fdr & result_balanced$Fndr_FW<th_fndr))
398
399 ## Count BW links from one to two
400 BW_1=length(which(result_balanced$flag_direction_mean=="One_more_than_two" &
    result_balanced$Fdr_BW<th_fdr & result_balanced$Fndr_FW<th_fndr))
401
402 ## Count BW links from two to one
403 BW_2=length(which(result_balanced$flag_direction_mean=="Two_more_than_one" &
    result_balanced$Fdr_FW<th_fdr & result_balanced$Fndr_BW<th_fndr))
404
405 FW_links=FW_1+FW_2
406 BW_links=BW_1+BW_2
407
408 FW_links/(FW_links+BW_links)
409 # 0.4955952
410
411 ## Ahora, a ver el criterio max:
412
413 ## Count FW links from one to two
414 FW_1=length(which(result_balanced$flag_direction_max=="One_more_than_two" &
    result_balanced$Fdr_FW<th_fdr & result_balanced$Fndr_BW<th_fndr))
415 ## Count FW links from two to one
416 FW_2=length(which(result_balanced$flag_direction_max=="Two_more_than_one" &
    result_balanced$Fdr_BW<th_fdr & result_balanced$Fndr_FW<th_fndr))
417
418 ## Count BW links from one to two
419 BW_1=length(which(result_balanced$flag_direction_max=="One_more_than_two" &
    result_balanced$Fdr_BW<th_fdr & result_balanced$Fndr_FW<th_fndr))
420 ## Count BW links from two to one
421 BW_2=length(which(result_balanced$flag_direction_max=="Two_more_than_one" &
    result_balanced$Fdr_FW<th_fdr & result_balanced$Fndr_BW<th_fndr))

```



```

419
420 FW_links=FW_1+FW_2
421 BW_links=BW_1+BW_2
422
423 FW_links/(FW_links+BW_links)
424 # 0.5002591
425
426 write.table(result_balanced,"outputs/4_MR/match/result_fw.txt")

```

7.6. Implementation of empirical null models to quantify directional bias in link prediction for Top-SNPs method

```

1 library(qvalue)
2
3 result_top=read.table("outputs/4_MR/top/result.txt")
4 forward_top=result_top[which(result_top$Fdr_FW<0.1 & result_top$Fndr_BW<0.1),]
5 backward_top=result_top[which(result_top$Fndr_FW<0.1 & result_top$Fdr_BW<0.1),]
6
7 iter=100000
8 N=nrow(result_top)
9 Fdr=0.1
10 Fndr=0.1
11 N1=nrow(forward_top)
12 N2=nrow(backward_top)
13
14 n_media=(N1+N2)/2
15 N1=n_media
16 N2=n_media
17
18 pi_0_FW=qvalue(result_top$P_FW)$pi0
19 pi_0_BW=qvalue(result_top$P_BW)$pi0
20
21 pi_medio=(pi_0_FW+pi_0_BW)/2
22 pi_0_FW=pi_medio
23 pi_0_BW=pi_medio
24
25 p1_FW=(N1*(1-Fdr))/(N*pi_0_FW*(2-Fdr))
26 p2_FW=(N2*Fndr)/(N*pi_0_FW)
27 p3_FW=(N2*(1-Fndr))/(N*(1-pi_0_FW))
28 p4_FW=(N1*Fdr)/(N*(1-pi_0_FW)*(2-Fdr))
29

```

```

30 p1_BW=(N1*(1-Fdr))/(N*pi_0_BW*(2-Fdr))
31 p2_BW=(N2*Fndr)/(N*pi_0_BW)
32 p3_BW=(N2*(1-Fndr))/(N*(1-pi_0_BW))
33 p4_BW=(N1*Fdr)/(N*(1-pi_0_BW)*(2-Fdr))
34
35 percentages=rep(NA,iter)
36
37 for(i in 1:iter){
38
39   ## Trato de declarar solo un dataframe con todo lo asignable para hacer la alocaación de memoria de
      todo de una vez:
40   tab=data.frame(
41     platonic_FW=rep("FALSE",N),empiric_FW=rep("UNKNOWN",N),
42     platonic_BW=rep("FALSE",N),empiric_BW=rep("UNKNOWN",N),
43     r_platonic_FW=runif(N),r_empiric_FW=runif(N),
44     r_platonic_BW=runif(N),r_empiric_BW=runif(N))
45
46   ## Declarar platonic labels:
47   tab$platonic_FW[which(tab$r_platonic_FW<pi_0_FW)]=TRUE"
48   tab$platonic_BW[which(tab$r_platonic_BW<pi_0_BW)]=TRUE"
49
50   ## Declarar empiric labels:
51   tab$empiric_FW[which(tab$platonic_FW==TRUE" & tab$r_empiric_FW< p1_FW)]=TRUE"
52   tab$empiric_FW[which(tab$platonic_FW==TRUE" & tab$r_empiric_FW>=p1_FW &
      tab$r_empiric_FW<(p1_FW+p2_FW))]=FALSE"
53   tab$empiric_FW[which(tab$platonic_FW==FALSE" & tab$r_empiric_FW< p4_FW)]=TRUE"
54   tab$empiric_FW[which(tab$platonic_FW==FALSE" & tab$r_empiric_FW>=p4_FW &
      tab$r_empiric_FW<(p4_FW+p3_FW))]=FALSE"
55
56   tab$empiric_BW[which(tab$platonic_BW==TRUE" & tab$r_empiric_BW< p1_BW)]=TRUE"
57   tab$empiric_BW[which(tab$platonic_BW==TRUE" & tab$r_empiric_BW>=p1_BW &
      tab$r_empiric_BW<(p1_BW+p2_BW))]=FALSE"
58   tab$empiric_BW[which(tab$platonic_BW==FALSE" & tab$r_empiric_BW< p4_BW)]=TRUE"
59   tab$empiric_BW[which(tab$platonic_BW==FALSE" & tab$r_empiric_BW>=p4_BW &
      tab$r_empiric_BW<(p4_BW+p3_BW))]=FALSE"
60
61   # Ahora calculo los porcentajes
62   number_links_only_FW=length(which(tab$empiric_FW==TRUE" & tab$empiric_BW==FALSE"))
63   number_links_only_BW=length(which(tab$empiric_FW==FALSE" & tab$empiric_BW==TRUE"))
64
65   percentages[i]=number_links_only_FW/(number_links_only_FW+number_links_only_BW)
66
67   if(i%%100==0) print(i)
68 }
69
70 summary(percentages)
71 #Min. 1st Qu. Median Mean 3rd Qu. Max.
72 #0.3591 0.4775 0.5000 0.5000 0.5227 0.6667
73 sd(percentages)
74 #0.0.03365883
75
76 save(percentages,file="outputs/4_MR/top/percentages_top_100K.Rdata")

```

7.7. Implementation of empirical null models to quantify directional bias in link prediction for Balanced-SNPs method

```

1 library(qvalue)
2
3 result_balanced=read.table("outputs/4_MR/match/result_fw.txt")
4 forward_balanced=result_balanced[which(result_balanced$Fdr_FW<0.1 & result_balanced$Fndr_BW<0.1),]
5 backward_balanced=result_balanced[which(result_balanced$Fndr_FW<0.1 & result_balanced$Fdr_BW<0.1),]
6
7 iter=100000
8 N=nrow(result_balanced)
9 Fdr=0.1
10 Fndr=0.1
11 N1=nrow(forward_balanced)
12 N2=nrow(backward_balanced)
13
14 n_media=(N1+N2)/2
15 N1=n_media
16 N2=n_media
17
18 pi_0_FW=qvalue(result_balanced$P_FW)$pi0
19 pi_0_BW=qvalue(result_balanced$P_BW)$pi0
20
21 pi_medio=(pi_0_FW+pi_0_BW)/2
22 pi_0_FW=pi_medio
23 pi_0_BW=pi_medio
24
25 p1_FW=(N1*(1-Fdr))/(N*pi_0_FW*(2-Fdr))
26 p2_FW=(N2*Fndr)/(N*pi_0_FW)
27 p3_FW=(N2*(1-Fndr))/(N*(1-pi_0_FW))
28 p4_FW=(N1*Fdr)/(N*(1-pi_0_FW)*(2-Fdr))
29
30
31 p1_BW=(N1*(1-Fdr))/(N*pi_0_BW*(2-Fdr))
32 p2_BW=(N2*Fndr)/(N*pi_0_BW)
33 p3_BW=(N2*(1-Fndr))/(N*(1-pi_0_BW))
34 p4_BW=(N1*Fdr)/(N*(1-pi_0_BW)*(2-Fdr))
35
36 percentages=rep(NA,iter)
37
38 for(i in 1:iter){
39     ## Trato de declarar solo un dataframe con todo lo asignable para hacer la alocaación de memoria de
40     ## todo de una vez:
41     tab=data.frame(
42         platonic_FW=rep("FALSE",N),empiric_FW=rep("UNKNOWN",N),
43         platonic_BW=rep("FALSE",N),empiric_BW=rep("UNKNOWN",N),
44         r_platonic_FW=runif(N),r_empiric_FW=runif(N),
45         r_platonic_BW=runif(N),r_empiric_BW=runif(N))
46
47     ## Declarar platonic labels:
48     tab$platonic_FW[which(tab$r_platonic_FW<pi_0_FW)]=TRUE
49     tab$platonic_BW[which(tab$r_platonic_BW<pi_0_BW)]=TRUE
50
51     ## Declarar empiric labels:
52     tab$empiric_FW[which(tab$platonic_FW==TRUE & tab$r_empiric_FW< p1_FW)]=TRUE
53     tab$empiric_FW[which(tab$platonic_FW==TRUE & tab$r_empiric_FW>=p1_FW &
54         tab$r_empiric_FW<(p1_FW+p2_FW))]=FALSE
55     tab$empiric_FW[which(tab$platonic_FW==FALSE & tab$r_empiric_FW< p4_FW)]=TRUE
56     tab$empiric_FW[which(tab$platonic_FW==FALSE & tab$r_empiric_FW>=p4_FW &
57         tab$r_empiric_FW<(p4_FW+p3_FW))]=FALSE
58 }

```

```

56   tab$empiric_BW[which(tab$platonc_BW=="TRUE" & tab$empiric_BW< p1_BW)]= "TRUE"
57   tab$empiric_BW[which(tab$platonc_BW=="TRUE" & tab$empiric_BW>=p1_BW &
      tab$empiric_BW<(p1_BW+p2_BW))]= "FALSE"
58   tab$empiric_BW[which(tab$platonc_BW=="FALSE" & tab$empiric_BW< p4_BW)]= "TRUE"
59   tab$empiric_BW[which(tab$platonc_BW=="FALSE" & tab$empiric_BW>=p4_BW &
      tab$empiric_BW<(p4_BW+p3_BW))]= "FALSE"
60
61   # Ahora calculo los porcentajes
62   number_links_only_FW=length(which(tab$empiric_FW=="TRUE" & tab$empiric_BW=="FALSE"))
63   number_links_only_BW=length(which(tab$empiric_FW=="FALSE" & tab$empiric_BW=="TRUE"))
64
65   percentages[i]=number_links_only_FW/(number_links_only_FW+number_links_only_BW)
66
67   if(i%%100==0) print(i)
68 }
69
70 summary(percentages)
71 #Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
72 #0.4404  0.4912  0.5000  0.5000  0.5088  0.5567
73 sd(percentages)
74 #0.01306218
75
76 save(percentages, file="outputs/4_MR/match/percentages_match_100K.Rdata")

```

Referencias

- [1] Getting started with RNA-sequencing. WANG Z., GERSTEIN M. & SNYDER M., "*RNA-seq: a revolutionary tool for transcriptomics*". Nat. Rev. Genet, 2009;10(1), 57–63. Retrieved in August 2022 from <https://www.technologynetworks.com/genomics/articles/rna-seq-basics-applications-and-protocol-299461>
- [2] Getting started with SNP-chips. "Use of SNP chips to detect rare pathogenic variants: retrospective, population based diagnostic evaluation". <https://www.bmj.com/content/372/bmj.n214>
- [3] THOMAS LAFRAMBOISE, "*Single nucleotide polymorphism arrays: a decade of biological, computational and technological advances, Nucleic Acids Research*"., Volume 37, Issue 13, 1 July 2009, Pages 4181–4193, <https://doi.org/10.1093/nar/gkp552>
- [4] STUART, JOSHUA M; SEGAL, ERAN; KOLLER, DAPHNE; KIM, STUART K. (2003) "*A gene-coexpression network for global discovery of conserved genetic modules*".
- [5] WEIRAUCH, MATTHEW T (2011). "*Gene coexpression networks for the analysis of DNA microarray data*". Applied Statistics for Network Biology: Methods in Systems Biology. pp. 215–250.
- [6] GORLOV, I., XIAO, X., MAYES, M. ET AL." *SNP eQTL status and eQTL density in the adjacent region of the SNP are associated with its statistical significance in GWA studies*"., <https://bmccgenomdata.biomedcentral.com/articles/10.1186/s12863-019-0786-0>
- [7] ALEXANDRA C. NICA EMMANOUIL T. DERMITZAKIS, "*Expression quantitative trait loci: present and future*".
- [8] Consulting about eQTLs in <https://es.wikipedia.org/wiki/eQTL>
- [9] JÉSSICA MOREIRA BATISTA DA SILVA (2020) "*Characterization of molecular heterogeneity in celiac disease patients*" p. 8-10.
- [10] ROBINSON MD, MCCARTHY DJ, SMYTH GK "*edgeR: A Bioconductor package for differential expression analysis of digital gene expression data*". Bioinformatics. 2009; 26(1):139-40
- [11] RITCHIE, M. E., PHIPSON, B., WU, D. I., HU, Y., LAW, C.W., SHI, W., SMYTH, G. K. (2015) "*limma powers differential expression analyses for RNA-sequencing and microarray studies*". Nucleic acids research, 43(7),e47-e47.
- [12] REGINA SANTESTEBAN AZANZA (2020) "*Métodos computacionales para la caracterización de relaciones causales entre genotipo y fenotipo: Heredabilidad de la expresión y coexpresión genética*".
- [13] BENJAMINI, Y., HOCHBERG, Y (1995) "*Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing*". Journal of the Royal statistical society: series B (Methodological), 57(1), 289-300
- [14] STOREY JD, TIBSHIRANI R (2003) "*Statistical significance for genomewide studies*". PNAS 100:9440–9445