

TFARM: Transcription Factor Association Rule Miner

Liuba Nausicaa Martino liuban.martino@gmail.com

Alice Parodi
Gaia Ceddia
Piercesare Secchi
Stefano Campaner
Marco Masseroli

April 26, 2023

Contents

`library(TFARM)`

1 Introduction

Looking for association rules between transcription factors in genomic regions of interest can be useful to find direct or indirect interactions among regulatory factors of DNA transcription. However, the results provided by the most recent algorithms for the search of association rules [?] [?] alone are often not intelligible enough, since they only provide a list of association rules. A novel method is proposed for subsequent mining of these results to evaluate the contribution of the items in each association rule. The *TFARM* package allows us to identify and extract the most relevant association rules with a given target transcription factor and compute the *Importance Index* of a transcription factor (or a combination of some of them) in the extracted rules. Such an index is useful to associate a numerical value to the contribution of one or more transcription factors to the co-regulation with a given target transcription factor.

2 Dataset

Association rules are extracted from a GRanges object in which metadata columns identify transcription factors and genomic coordinates are represented in the left-hand-side of the GRanges; thus, each row is a different genomic region. The element (i,j) (with $j > 4$) of the metadata section is equal to 0 if a binding site of transcription factor j is absent in region i, or to 1 (or any other value) if it is present. This dataset, called *matrix of presences*, should not have rows with only 0 values since we consider regions with no transcription factors as uninteresting regions. The first three columns of the GRanges contain the chromosome name, the genomic coordinates (i.e., *left* and *right* coordinate are the leftmost and rightmost bases of the DNA region), and the strand (encoded as "+", "-", or "*" if unknown), of each region respectively. The GRanges is obtained from the analysis of ENCODE ChIP-seq data:

it concerns the localization of transcription factors binding sites and histone modifications in DNA, as well as RefSeq data (<https://www.ncbi.nlm.nih.gov/refseq/>); specifically, here we focus on promotorial regions, but further analyses are possible on any region of interest. Such data have been processed and extracted with GMQL (GenoMetric Query Language [?], <http://www.bioinformatics.deib.polimi.it/GMQL/>) queries. In this example, the dataset we consider is the matrix of presences of 25 transcription factors' binding sites of the MCF-7 human breast adenocarcinoma cell line (i.e., all the transcription factors evaluated in ENCODE for this cell line), in the 2,944 promotorial regions of chromosome 1:

```
# Load and visualize the dataset:

data("MCF7_chr1")
length(MCF7_chr1)

## [1] 2944

MCF7_chr1

## GRanges object with 2944 ranges and 25 metadata columns:
##      seqnames      ranges strand |      PML
##      <Rle>        <IRanges> <Rle> | <numeric>
##      19      chr1      565116-568115 - |      0
##      20      chr1      565156-568155 - |      0
##      21      chr1      565264-568263 - |      0
##      22      chr1      566844-569843 + |      0
##      31      chr1      713069-716068 - |      0
##      ...      ...      ...      ... |      ...
##      4752     chr1 249152126-249155125 - |      1
##      4753     chr1 249152316-249155315 - |      1
##      4754     chr1 249166054-249169053 + |      0
##      4755     chr1 249166447-249169446 + |      0
##      4756     chr1 249198442-249201441 + |      1
##      SRF      CTCF      TCF7L2      FOSL2      SIN3AK20
##      <numeric> <numeric> <numeric> <numeric> <numeric>
##      19      0      1      1      0      0
##      20      0      1      1      0      0
##      21      0      1      1      0      0
##      22      0      1      1      0      0
##      31      1      1      0      0      1
##      ...      ...      ...      ...      ...
##      4752     0      1      0      0      1
##      4753     0      1      0      0      1
##      4754     1      1      0      1      1
##      4755     1      1      0      1      1
##      4756     0      1      0      0      1
##      HDAC2      EP300      GABPA      EGR1      HA.E2F1
##      <numeric> <numeric> <numeric> <numeric> <numeric>
##      19      0      0      0      0      0
##      20      0      0      0      0      0
##      21      0      0      0      0      0
##      22      0      0      0      0      0
##      31      0      0      0      0      1
```

```
##      ...      ...      ...      ...      ...      ...
## 4752      0      0      1      0      0
## 4753      0      0      1      0      0
## 4754      1      0      0      0      0
## 4755      1      0      0      0      0
## 4756      0      0      1      0      0
##          GATA3      REST      FOXM1      MYC      MAX
##      <numeric> <numeric> <numeric> <numeric> <numeric>
## 19      0      0      0      0      0
## 20      0      0      0      0      0
## 21      0      0      0      0      0
## 22      0      0      0      0      0
## 31      0      0      0      1      1
##      ...      ...      ...      ...      ...      ...
## 4752      0      0      0      1      1
## 4753      0      0      0      1      1
## 4754      0      0      0      1      1
## 4755      0      0      0      1      1
## 4756      0      0      0      1      1
##          TEAD4      CEBPB      JUND      RAD21      TAF1
##      <numeric> <numeric> <numeric> <numeric> <numeric>
## 19      0      0      0      0      0
## 20      0      0      0      0      0
## 21      0      0      0      0      0
## 22      0      0      0      0      0
## 31      0      0      0      0      1
##      ...      ...      ...      ...      ...      ...
## 4752      0      0      0      0      1
## 4753      0      0      0      0      1
## 4754      1      0      0      0      0
## 4755      1      0      0      0      0
## 4756      0      0      0      1      0
##          TCF12      ELF1      ZNF217      NR2F2
##      <numeric> <numeric> <numeric> <numeric>
## 19      0      0      0      0
## 20      0      0      0      0
## 21      0      0      0      0
## 22      0      0      0      0
## 31      0      1      0      0
##      ...      ...      ...      ...      ...
## 4752      0      1      0      0
## 4753      0      1      0      0
## 4754      0      1      0      1
## 4755      0      1      0      1
## 4756      0      1      0      0
## -----
## seqinfo: 24 sequences from an unspecified genome; no seqlengths
```

3 Extraction of the most relevant associations

We define a relevant association for the prediction of transcription factor TFt in the considered genomic regions as an association rule of the type:

$$\{TF1=1, TF2=1, TF3=1\} \rightarrow \{TFt=1\}$$

which means that the presence of the transcription factors TF1, TF2, and TF3 implies the presence of transcription factor TFt. Every association rule is characterized by a set of three measures: support, confidence and lift:

- *support*:

$$supp(X \rightarrow Y) = \frac{supp(X \cup Y)}{N} \quad 1$$

where N is the number of transactions, $X \cup Y$ is a set of items and $Supp(X \cup Y)$ is the support of the itemset $\{X, Y\}$, defined as

$$supp(X) = \frac{|\{t_i \in N; X \subseteq t_i\}|}{N} \quad 2$$

that is the proportion of transactions t_i in the dataset which contains the itemset X. The support of an association rule measures the frequency of a rule in the dataset and varies in the interval [0,1].

- *confidence*:

$$conf(X \rightarrow Y) = \frac{supp(X \cup Y)}{supp(X)} \quad 3$$

It gives an estimate of the conditioned probability $P(Y|X)$, that is the probability to find the right-hand-side (RHS) of the rule (i.e., the itemset Y) in a set of transactions, given that these transactions also contain the left-hand-side (LHS) of the rule (i.e., the itemset X). Therefore, it measures the reliability of the inference made by the rule $X \rightarrow Y$. The higher is the confidence of the rule, the higher is the probability to find the itemset Y in a transaction containing the itemset X. It varies in the interval [0,1].

- *lift*:

$$lift(X \rightarrow Y) = \frac{supp(X \cup Y)}{supp(X)supp(Y)} \quad 4$$

It measures the strength of the rule, and varies in the interval $[0, \infty]$.

To extract a set of relevant associations the user has to specify:

1. the presence/absence of the target transcription factor to be predicted, TFt;
2. the minimal support threshold of the rules to be extracted;
3. the minimal confidence threshold of the rules to be extracted.

Points 2. and 3. strongly depend on the dimensions of the dataset (i.e., number of rows - regions - and number of columns - transcription factors), the presence of the target transcription factor in the considered regions, the number of relevant associations that the user wants to find. Usually, the confidence threshold is set higher than 0.5, since it measures the posterior probability to have TFt given the presence of the pattern in the left-hand-side of the rule (e.g., $\{TF1=1, TF2=1, TF3=1\}$). The function `rulesGen` in the **TFARM** package extracts the association rules by calling the `apriori` function of the `arules` package [?] [?] [?]. It takes in input:

- the GRanges object in which the matrix of presences is represented;
- the target transcription factor;
- the minimum support threshold of the rules to be extracted;
- the minimum confidence threshold of the rules to be extracted;
- the logical parameter *type* that sets the type of left-hand-side of the rules to be extracted (i.e., containing only present transcription factors, or containing present and/or absent transcription factors).

The result of the `rulesGen` function is a data.frame containing:

- in the first column the left-hand-side of each extracted rule;
- in the second column the right-hand-side of each extracted rule (that is the presence/absence of the given target transcription factor);
- in the third column the support of each extracted rule;
- in the fourth column the confidence of each extracted rule;
- in the fifth column the lift of each extracted rule.

See *arulesViz* package for visualization tools of association rules.

```
# Coming back to the example on the transcription factors of cell line
# MCF-7, in the promotorial regions of chromosome 1.
# Suppose that the user wants to find the most relevant association
# rules for the prediction of the presence of TEAD4. This means extracting
# all the association rules with right-hand-side equal to {TEAD4=1}
# setting the parameter type = TRUE; the minimum support and minimum
# confidence thresholds are set, as an example, to 0.005 and 0.62,
# respectively:

r_TEAD4 <- rulesGen(MCF7_chr1, "TEAD4=1", 0.005, 0.62, TRUE)

## Apriori
##
## Parameter specification:
## confidence minval smax arem aval originalSupport maxtime
##      0.62    0.1    1 none FALSE          TRUE      5
## support minlen maxlen target ext
##    0.005      2    20  rules TRUE
##
## Algorithmic control:
## filter tree heap memopt load sort verbose
##    0.1 TRUE TRUE  FALSE TRUE    2    TRUE
##
## Absolute minimum support count: 14
##
## set item appearances ...[25 item(s)] done [0.00s].
## set transactions ...[25 item(s), 2944 transaction(s)] done [0.00s].
## sorting and recoding items ... [25 item(s)] done [0.00s].
## creating transaction tree ... done [0.00s].
## checking subsets of size 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 done [0.18s].
```

```
## writing ... [30 rule(s)] done [0.01s].
## creating S4 object ... done [0.01s].

dim(r_Tead4)

## [1] 30 7

head(r_Tead4)

##               lhs               rhs
## 1 {GABPA=1,TCF12=1,ZNF217=1,NR2F2=1} {TEAD4=1}
## 2 {FOSL2=1,GABPA=1,MYC=1,ZNF217=1} {TEAD4=1}
## 3 {GABPA=1,MAX=1,TCF12=1,ZNF217=1,NR2F2=1} {TEAD4=1}
## 4 {FOSL2=1,GABPA=1,MYC=1,ZNF217=1,NR2F2=1} {TEAD4=1}
## 5 {FOSL2=1,HDAC2=1,GABPA=1,MYC=1,ZNF217=1} {TEAD4=1}
## 6 {FOSL2=1,GABPA=1,MYC=1,MAX=1,ZNF217=1} {TEAD4=1}
##      support confidence coverage lift count
## 1 0.005095109 0.6521739 0.00781250 27.42857 15
## 2 0.007812500 0.6571429 0.01188859 27.63755 23
## 3 0.005095109 0.6521739 0.00781250 27.42857 15
## 4 0.006453804 0.6333333 0.01019022 26.63619 19
## 5 0.006114130 0.6428571 0.00951087 27.03673 18
## 6 0.007133152 0.6363636 0.01120924 26.76364 21
```

Once the set of the most relevant association rules (i.e., with support and confidence higher than the thresholds specified as parameters) is extracted, the user can look for *candidate co-regulator transcription factors* with the target transcription factor (in the example TEAD4), which are the transcription factors present in the LHS of the extracted rules. This is provided by the function `presAbs` of the **TFARM** package. The function `presAbs` takes in input:

- a string vector containing the names of all transcription factors present in the matrix of presences;
- the set of the most relevant association rules previously extracted with `rulesGen`;
- a logical parameter, *type*, which refers to the type of rules extracted with the `rulesGen` function. If *type* = *TRUE*, the LHS of the rules contains only items of the type TF=1, otherwise, if *type* = *FALSE*, the LHS of the rules can contain both items TF=1 and TF=0.

The `presAbs` function has two outputs:

- *pres*, which is a string vector containing all the items present in the LHSs of the considered set of rules;
- *abs*, which is a string vector containing all the items absent in the LHSs of the considered set of rules.

```
# Transcription factors present in at least one of the regions:

c <- names(mcols(MCF7_chr1))
c

## [1] "PML"      "SRF"      "CTCF"     "TCF7L2"   "FOSL2"
## [6] "SIN3AK20" "HDAC2"    "EP300"    "GABPA"    "EGR1"
```

```
## [11] "HA.E2F1" "GATA3" "REST" "FOXM1" "MYC"
## [16] "MAX" "TEAD4" "CEBPB" "JUND" "RAD21"
## [21] "TAF1" "TCF12" "ELF1" "ZNF217" "NR2F2"

lc <- length(c)

names(presAbs(c, r_TEAD4, TRUE))

## [1] "pres" "abs"

# Transcription factors present in at least one of the association rules:

p_TFs <- presAbs(c, r_TEAD4, TRUE)$pres
p_TFs

## [1] "FOSL2=1" "SIN3AK20=1" "HDAC2=1" "GABPA=1"
## [5] "HA.E2F1=1" "GATA3=1" "MYC=1" "MAX=1"
## [9] "TCF12=1" "ELF1=1" "ZNF217=1" "NR2F2=1"

# Transcription factors absent in all the association rules:

a <- presAbs(c[1:lc], r_TEAD4, TRUE)$abs
a

## [1] "PML=1" "SRF=1" "CTCF=1" "TCF7L2=1" "EP300=1"
## [6] "EGR1=1" "REST=1" "FOXM1=1" "TEAD4=1" "CEBPB=1"
## [11] "JUND=1" "RAD21=1" "TAF1=1"
```

All transcription factors in `p` are said to be *candidate co-regulators* of the TFt in the most relevant associations extracted with `rulesGen`.

4 Importance Index of a transcription factor

The extraction of candidate transcription factors for the interaction with a given target transcription factor (TFt) can be useful to provide a global vision of the possible associations of TFt. However, since the number of association rules and candidate co-regulators can be very high, this list does not provide an intelligible result, giving the lack of measures of how much each transcription factor contributes to the existence of a certain complex of transcription factors. Let us consider for example the rule

$$\{TF1=1, TF2=1, TF3=1\} \rightarrow \{TFt=1\}$$

Just looking at it, the user could not tell if the presences of TF1, TF2 and TF3 equally contribute to the prediction of the presence of TFt. A solution to this problem can be given by removing, alternatively, TF1, TF2, and TF3 from the rule and evaluate:

- 1) if the rule keeps on existing and being relevant
- 2) how the three quality measures of support, confidence, and lift of the rule change.

If a rule is not found as relevant after removing a transcription factor from its LHS, then the presence of that transcription factor in the pattern $\{TF1=1, TF2=1, TF3=1\}$ is fundamental for the existence of the association rule $\{TF1=1, TF2=1, TF3=1\} \rightarrow \{TFt=1\}$. Otherwise, if the rule keeps on existing as relevant, and its quality measures are similar to the ones

of the rule initially considered, then the presence of that transcription factor in the pattern $\{TF1=1, TF2=1, TF3=1\}$ is not fundamental for the existence of the association rule $\{TF1=1, TF2=1, TF3=1\} \rightarrow \{TFt=1\}$. Let us fix an item I (i.e., a candidate co-regulator for TFt) and extract the subset of the most relevant associations containing I , named $\{R^I\}$ (with J number of rules in $\{R^I\}$, $J=|\{R^I\}|$). Each element of $\{R^I\}_{j=1:J}$ is described by a set of quality measures of support, confidence and lift: $\{s_j^I, c_j^I, l_j^I\}_{j=1:J}$.

rule	support	confidence	lift
R_1^I	s_1^I	c_1^I	l_1^I
...
R_J^I	s_J^I	c_J^I	l_J^I

Table 1: Rules containing item I , and correspondent measures of support, confidence and lift.

Let then be $\{R_j^{I-}\}_{j=1:J}$ the set of rules obtained substituting the presence of item I with its absence in each element of $\{R_j^I\}_{j=1:J}$. For example, if I is $TF1$ and R_j^I is the rule $\{TF1=1, TF2=1, TF3=1\} \rightarrow \{TFt=1\}$, with measures $\{s_j^I, c_j^I, l_j^I\}$, then R_j^{I-} will be the rule $\{TF1=0, TF2=1, TF3=1\} \rightarrow \{TFt=1\}$ with measures $\{s_j^{I-}, c_j^{I-}, l_j^{I-}\}$. Thus, R_j^I and R_j^{I-} consider the same number of association rules for each item.

rule	support	confidence	lift
R_1^{I-}	s_1^{I-}	c_1^{I-}	l_1^{I-}
...
R_J^{I-}	s_J^{I-}	c_J^{I-}	l_J^{I-}

Table 2: Rules originally containing item I obtained by removing I , and correspondent support, confidence and lift measures.

To analyze the importance of a transcription factor, for example $TF1$, we can compare the two distributions $\{s_j^I, c_j^I\}_{j=1:J}$ and $\{s_j^{I-}, c_j^{I-}\}_{j=1:J}$ for each j in $\{1, \dots, J\}$. Since support and confidence vary in $[0, 1]$, while the lift is directly proportional to the confidence measure, we can define an index of the importance of item I in the rule R_j^I for j in $\{1, \dots, J\}$ as:

$$imp(I)_j = \Delta s_j + \Delta c_j \quad 5$$

with: $\Delta s_j = s_j^I - s_j^{I-} \quad \Delta c_j = c_j^I - c_j^{I-}$

The importance of I in its set of rules $\{R^I\}$ is obtained evaluating the mean of all its importances $imp(I)_j$ in the set of rules:

$$imp(I) = \frac{\sum_{j=1}^J imp(I)_j}{J} \quad 6$$

Then, evaluating the index $imp(I)$ for each item in the relevant association rules extracted can be useful to rank the transcription factors by their importance in the association with the target transcription factor, TFt . The presence of the transcription factors with highest mean Importance Index is assumed to be fundamental for the existence of some regulatory complexes (i.e., association rules assumed to be relevant); the transcription factors with lower mean Importance Index, instead, do not significantly influence the pattern of transcription factors associated with the target transcription factor. The definition of the Importance Index can be extended to couples of items, triplets, and so on. This can be easily done by

substituting the item I with a set of items (for example $I = \{TF1=1, TF2=1\}$), and applying the rest of the procedure in a completely analogous way. Thus, we identify as R^I the set of rules containing both TF1 and TF2 and R^{I^-} as the set of correspondent rules without the two transcription factors. This kind of approach allows the identification of interactions between transcription factors that would be unrevealed just looking at a list of association rules. The `rulesTF` function in **TFARM** package provides the subset of input rules containing a given transcription factor TF_i . It takes in input:

- a set of rules
- the transcription factor TF_i that the user wants to find in the LHSs of a subset of the considered rules
- a logical parameter, *verbose*: if *verbose* = *TRUE*, a console message is returned if the searched subset of rules is empty.

The output of the function is a data.frame containing the subset of rules whose LHSs contain TF_i , and the corresponding quality measures. Using the introduced notation, the output of the `rulesTF` function is $\{R_j^I\}_{j=1:J}$ with the quality measures $\{s_j^I, c_j^I, l_j^I\}_{j=1:J}$. The data.frame has J rows and five columns: the first column contains the LHS of the selected rules, the second one contains the RHS of the rules and the last three columns contain s_j^I, c_j^I, l_j^I (that is a data.frame like the one in Table ??).

```
# To find the subset of rules containing the transcription factor FOSL2:

r_FOSL2 <- rulesTF(TFi = 'FOSL2=1', rules = r_TEAD4, verbose = TRUE)
head(r_FOSL2)

##                               lhs
## 1          {FOSL2=1,GABPA=1,MYC=1,ZNF217=1}
## 2      {FOSL2=1,GABPA=1,MYC=1,ZNF217=1,NR2F2=1}
## 3      {FOSL2=1,HDAC2=1,GABPA=1,MYC=1,ZNF217=1}
## 4      {FOSL2=1,GABPA=1,MYC=1,MAX=1,ZNF217=1}
## 5  {FOSL2=1,HDAC2=1,GABPA=1,GATA3=1,MYC=1,ZNF217=1}
## 6 {FOSL2=1,GABPA=1,HA.E2F1=1,GATA3=1,MYC=1,ZNF217=1}
##                               rhs      support confidence      lift
## 1 {TEAD4=1} 0.007812500 0.6571429 0.011888587
## 2 {TEAD4=1} 0.006453804 0.6333333 0.010190217
## 3 {TEAD4=1} 0.006114130 0.6428571 0.009510870
## 4 {TEAD4=1} 0.007133152 0.6363636 0.011209239
## 5 {TEAD4=1} 0.005434783 0.6400000 0.008491848
## 6 {TEAD4=1} 0.005095109 0.6250000 0.008152174

dim(r_FOSL2)[1]

## [1] 28
```

```
# If none of the rules in input to rulesTF contains the given item TFi,
# and verbose = TRUE, a warnig message is reported to the user:
```

```
r_CTCF <- rulesTF(TFi = 'CTCF=1', rules = r_TEAD4, verbose = TRUE)

## Warning in rulesTF(TFi = "CTCF=1", rules = r_TEAD4, verbose = TRUE): None of
the rules contains CTCF=1
```

If the user wants to evaluate the importance of item I in a set of rules R^I , the user needs to substitute the presence of I in all the left-hand-side patterns of R^I with its absence: this is done using the function `rulesTF0` in **TFARM** package. This function takes in input:

- the transcription factor TF_i to be removed
- a set of rules containing TF_i
- the total set of rules
- the `GRanges` object containing the matrix of presences
- the target transcription factor.

It returns a `data.frame` with the rules obtained substituting the presence of TF_i with its absence and the correspondent measures. Using the introduced notation, the output of the `rulesTF0` function is $\{R^{I-}_j\}_{j=1:J}$ with the quality measures $\{s^{I-}_j, c^{I-}_j, l^{I-}_j\}_{j=1:J}$. The `data.frame` has J rows and five columns: the first column contains the LHS of the rules in R^I without TF_i , the second one contains the RHS of the rules and the last three columns contain $s^{I-}_j, c^{I-}_j, l^{I-}_j$ (that is a `data.frame` like the one in Table ??).

```
# For example to evaluate FOSL2 importance in the set of rules r_FOSL2:

r_noFOSL2 <- rulesTF0('FOSL2=1', r_FOSL2, r_TEAD4, MCF7_chr1, "TEAD4=1")
```

```
row.names(r_FOSL2) <- match(r_FOSL2$lhs, r_TEAD4$lhs)
row.names(r_noFOSL2) <- match(r_FOSL2$lhs, r_TEAD4$lhs)
head(r_noFOSL2)

##                               lhs
## 2                {FOSL2=0,GABPA=1,MYC=1,ZNF217=1}
## 4                {FOSL2=0,GABPA=1,MYC=1,ZNF217=1,NR2F2=1}
## 5                {FOSL2=0,HDAC2=1,GABPA=1,MYC=1,ZNF217=1}
## 6                {FOSL2=0,GABPA=1,MYC=1,MAX=1,ZNF217=1}
## 7    {FOSL2=0,HDAC2=1,GABPA=1,GATA3=1,MYC=1,ZNF217=1}
## 8 {FOSL2=0,GABPA=1,HA.E2F1=1,GATA3=1,MYC=1,ZNF217=1}
##      rhs      support confidence      lift
## 2 {TEAD4=1} 0.0027173913 0.09876543 0.02751359
## 4 {TEAD4=1} 0.0020380435 0.13043478 0.01562500
## 5 {TEAD4=1} 0.0010190217 0.06521739 0.01562500
## 6 {TEAD4=1} 0.0027173913 0.10126582 0.02683424
## 7 {TEAD4=1} 0.0006793478 0.05714286 0.01188859
## 8 {TEAD4=1} 0.0006793478 0.04166667 0.01630435
```

Now that the two sets of rules $\{R^I_j\}_{j=1:J}$ and $\{R^{I-}_j\}_{j=1:J}$, and the two sets of measures $\{s^I_j, c^I_j, l^I_j\}_{j=1:J}$ and $\{s^{I-}_j, c^{I-}_j, l^{I-}_j\}_{j=1:J}$ are obtained, the user can compute the Importance Index distribution for the chosen transcription factor TF_i . This can be done with the function `IComp` in the **TFARM** package which takes in input:

- the transcription factor TF_i
- the subset of rules `rules_TF` containing TF_i (provided by the function `rulesTF`) with their quality measures of support, confidence and lift
- the subset of rules `rules_noTF` obtained from `rules_TF` removing TF_i (provided by the function `rulesTF0`)

- a logical parameter (figures) to graphically rapresent $\{s_j^I, c_j^I, l_j^I\}_{j=1:J}$ and $\{s_j^{I-}, c_j^{I-}, l_j^{I-}\}_{j=1:J}$; set *figures* = *TRUE* to get it as an output.

The function has five outputs:

- *imp*, which is the set of importance index values of TFi in the given set of rules (*rules_TF*), one value for each rule.
- *delta*, which is the matrix of variations of standardized support, confidence, and lift, obtained removing TFi from *rules_TF*.
- *rwi*, which is a data.frame that contains rules from *rulesTF* associated with each candidate co-regulator transcription factor.
- *rwo*, which is a data.frame with rules in *rwi* obtained removing each transcription factor TFi.
- the plots of $\{s_j^I, c_j^I, l_j^I\}_{j=1:J}$ and $\{s_j^{I-}, c_j^{I-}, l_j^{I-}\}_{j=1:J}$ obtained if the user sets *figures* = *TRUE*.

```
# Perform the IComp function to compute the Importance Index distribution:

imp_FOSL2 <- IComp('FOSL2=1', r_FOSL2, r_noFOSL2, figures=TRUE)
names(imp_FOSL2)

## [1] "imp"      "delta"    "rwi"      "rwo"

imp_FOSL2$imp

## [1] 0.5634725 0.5073143 0.5827349 0.5395136 0.5876126
## [6] 0.5877491 0.5575295 0.6028063 0.5691677 0.5912939
## [11] 0.5261723 0.5722729 0.5876126 0.5876126 0.5859375
## [16] 0.5406774 0.5859375 0.5560171 0.5912939 0.5912939
## [21] 0.5261723 0.5722729 0.5722729 0.5876126 0.5849713
## [26] 0.5391110 0.5912939 0.5722729

head(imp_FOSL2$delta)

##      diff_supp diff_conf
## 1 0.005095109 0.5583774
## 2 0.004415761 0.5028986
## 3 0.005095109 0.5776398
## 4 0.004415761 0.5350978
## 5 0.004755435 0.5828571
## 6 0.004415761 0.5833333

head(imp_FOSL2$rwi)

##                                     lhs
## 2                      {FOSL2=1,GABPA=1,MYC=1,ZNF217=1}
## 4          {FOSL2=1,GABPA=1,MYC=1,ZNF217=1,NR2F2=1}
## 5          {FOSL2=1,HDAC2=1,GABPA=1,MYC=1,ZNF217=1}
## 6          {FOSL2=1,GABPA=1,MYC=1,MAX=1,ZNF217=1}
## 7  {FOSL2=1,HDAC2=1,GABPA=1,GATA3=1,MYC=1,ZNF217=1}
## 8 {FOSL2=1,GABPA=1,HA.E2F1=1,GATA3=1,MYC=1,ZNF217=1}
##                                     rhs      support confidence      lift
## 2 {TEAD4=1} 0.007812500 0.6571429 0.011888587
```

```
## 4 {TEAD4=1} 0.006453804 0.6333333 0.010190217
## 5 {TEAD4=1} 0.006114130 0.6428571 0.009510870
## 6 {TEAD4=1} 0.007133152 0.6363636 0.011209239
## 7 {TEAD4=1} 0.005434783 0.6400000 0.008491848
## 8 {TEAD4=1} 0.005095109 0.6250000 0.008152174

head(imp_FOSL2$rwo)

##                               lhs
## 2                {FOSL2=0,GABPA=1,MYC=1,ZNF217=1}
## 4      {FOSL2=0,GABPA=1,MYC=1,ZNF217=1,NR2F2=1}
## 5      {FOSL2=0,HDAC2=1,GABPA=1,MYC=1,ZNF217=1}
## 6      {FOSL2=0,GABPA=1,MYC=1,MAX=1,ZNF217=1}
## 7 {FOSL2=0,HDAC2=1,GABPA=1,GATA3=1,MYC=1,ZNF217=1}
## 8 {FOSL2=0,GABPA=1,HA.E2F1=1,GATA3=1,MYC=1,ZNF217=1}
##      rhs      support confidence      lift
## 2 {TEAD4=1} 0.0027173913 0.09876543 0.02751359
## 4 {TEAD4=1} 0.0020380435 0.13043478 0.01562500
## 5 {TEAD4=1} 0.0010190217 0.06521739 0.01562500
## 6 {TEAD4=1} 0.0027173913 0.10126582 0.02683424
## 7 {TEAD4=1} 0.0006793478 0.05714286 0.01188859
## 8 {TEAD4=1} 0.0006793478 0.04166667 0.01630435
```

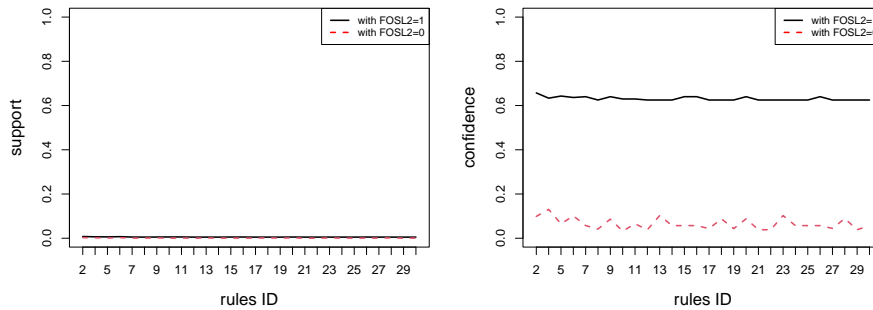


Figure 1: Support and Confidence for the extracted rules before and after the removal of item I . Left panel: Support distribution $\{s_j^I\}_{j=1:J}$, black thick line, and $\{s_j^{I-}\}_{j=1:J}$, red dotted line. Right panel: Confidence distribution $\{c_j^I\}_{j=1:J}$, black thick line, and $\{c_j^{I-}\}_{j=1:J}$, red dotted line.

The most useful application of the function `IComp` is the ranking of candidate co-regulators through their importances. As previously seen, the candidate co-regulators are returned by the function `presAbs`. The evaluation of the mean Importance of each co-regulator can be computed cycling the three functions `rulesTF`, `rulesTF0` and `IComp` over a string vector with all transcription factors present in the set of relevant association rules extracted, as returned by the function `presAbs`.

```
# For the considered example the user could run:
DELTA_mean_supp <- vector("list", length(p_TFs))
DELTA_mean_conf <- vector("list", length(p_TFs))
all <- lapply(p_TFs, function(pi) {
```

```

    A <- rulesTF(pi, r_Tead4, FALSE)
    B <- rulesTF0(pi, A, r_Tead4, MCF7_chr1, "TEAD4=1")
    IComp(pi, A, B, figures=FALSE)
  })

for (i in 1:length(p_TFs)) {
  IMP_Z[[i]] <- all[[i]]$imp
  # Extract the delta variations of support and confidence:
  DELTA_mean_supp[[i]] <- apply(all[[i]]$delta[1], 2, mean)
  DELTA_mean_conf[[i]] <- apply(all[[i]]$delta[2], 2, mean)
}

IMP <- data.frame(
  TF = p_TFs,
  imp = sapply(IMP_Z, mean),
  sd = sapply(IMP_Z, sd),
  delta_support = as.numeric(DELTA_mean_supp),
  delta_confidence = as.numeric(DELTA_mean_conf),
  nrules = sapply(IMP_Z, length),
  stringsAsFactors=FALSE
)

library(plyr)

##
## Attaching package: 'plyr'
##
## The following object is masked from 'package:IRanges':
##
## desc
##
## The following object is masked from 'package:S4Vectors':
##
## rename

```

```

# Sort by imp column of IMP

IMP.ord <- arrange(IMP, desc(imp))
IMP.ord

##      TF      imp      sd delta_support
## 1  ZNF217=1 0.6337003 0.0137307484 0.005423460
## 2   FOSL2=1 0.5700001 0.0245475607 0.004440023
## 3   GABPA=1 0.5630153 0.0759261948 0.005151721
## 4    MYC=1 0.5556017 0.1562668117 0.005349864
## 5    MAX=1 0.4937297 0.3622483586 0.005289208
## 6   GATA3=1 0.4826400 0.2562831652 0.005064229
## 7   TCF12=1 0.4607986 0.0205092581 0.001358696
## 8 SIN3AK20=1 0.4128163 0.2408282842 0.005038496
## 9   HDAC2=1 0.3214813 0.2852705613 0.004959239
## 10 HA.E2F1=1 0.1295856 0.0001961108 0.004585598
## 11   ELF1=1 -0.1201596 0.4630672884 0.004840353
## 12  NR2F2=1 -0.1433436 0.4269140814 0.004925272

```

```
##      delta_confidence nrules
## 1      0.6282768      30
## 2      0.5655600      28
## 3      0.5578636      30
## 4      0.5502518      28
## 5      0.4884405      14
## 6      0.4775758      22
## 7      0.4594399       2
## 8      0.4077778      12
## 9      0.3165220      15
## 10     0.1250000       4
## 11     -0.1250000       8
## 12     -0.1482689      10
```

In this way we get, besides the mean Importance Index of each candidate co-regulator of TFt (TFt = TEAD4 in the example), the standard deviation of the distribution of the Importance Index of each candidate co-regulator of TFt, and the number of rules in which each candidate co-regulator of TFt is present. The function `IComp` can be easily generalized for the computation of the mean Importance Index of combinations of transcription factors (see the example used for the `heatI` function in the following section).

4.1 Validation of the Importance Index formula

We defined the Importance Index of an item in an association rule as the linear combination of the variations of the support and confidence of the rule obtained substituting the presence of the item in the left-hand-side of the association rule, with its absence (as in Formula ??). In this way, we assume that each of the two variations equally contributes to the evaluation of the contribution of the item to the prediction of the presence of another item in the right-hand-side of the considered association rule. Nevertheless, one of the two quality measures might be more or less sensitive than the other to the removal of the item from the rule, leading to a greater or smaller variation of one or more of the values of support and confidence.

TF	Δs_s	Δc_c	
TF_1	$\Delta s_{s,1}$	$\Delta c_{c,1}$	
...
TF_1	$\Delta s_{s,n_1}$	$\Delta c_{c,n_1}$	
...
TF_M	$\Delta s_{s,K-n_M+1}$	$\Delta c_{c,K-n_M+1}$	
...
TF_M	$\Delta s_{s,K}$	$\Delta c_{c,K}$	

Table 3: Matrix with the variations of the support and confidence, obtained removing each transcription factor from the subset of rules in which it is present. M is the total number of transcription factors, K is the total number of rules and n_i is the number of rules for transcription factor TF_i .

Thanks to the Principal Components Analysis [?] [?], computed by the function `IPCA` in the **TFARM** package, we can evaluate if it is possible to find a subspace of \mathbb{R}^2 in which the most variability of the dataset containing the variations of the measures (Table ??) is captured.

This can be easily done by extracting the delta variations of support and confidence, using the function `IComp`, simply getting its *delta* output, as well as a matrix containing the candidate co-regulators found, and the number of rules in which each of them appears.

A principal component is a combination of the original variables after a linear transformation; the set of principal components defines a new reference system. The new coordinates of data represented in the reference system defined by principal components are called *scores*, and the coefficients of the linear combination that define each principal component are called *loadings* (so, loadings give a measure of the contribution of every observation to each principal component).

The `IPCA` function takes in input:

- the list of variations of distributions of support and confidence measures, obtained from the `IComp` function
- a matrix with the mean importance of all candidate co-regulators and the number of rules in which each of them appears.

It returns:

- a summary, containing: the standard deviation on each principal component, the proportion of variance explained by each principal component, and the cumulative proportion of variance described by each principal component;
- the scores of each principal component
- the loadings of each principal component
- a plot with the variability and the cumulate percentage of variance explained by each principal component
- a plot with the loadings of the principal components

```
# Select the candidate co-regulators and the number of rules
# associated with them, then perform the Principal Component Analysis:

colnames(IMP)

## [1] "TF"          "imp"          "sd"
## [4] "delta_support" "delta_confidence" "nrules"

TF_Imp <- data.frame(IMP$TF, IMP$imp, IMP$nrules)
i.pc <- IPKA(DELTA, TF_Imp)
names(i.pc)

## [1] "summary" "scores" "loadings"

i.pc$summary

## Importance of components:
##              Comp.1      Comp.2
## Standard deviation    0.4110147 0.1435518
## Proportion of Variance 0.8912784 0.1087216
## Cumulative Proportion 0.8912784 1.0000000

head(i.pc$loadings)
```

```
##          Comp.1      Comp.2
## delta s 0.7009374  0.7132228
## delta c 0.7132228 -0.7009374
```

```
head(i.pc$scores)
```

```
##          Comp.1      Comp.2
## [1,] -0.3199744  0.006961379
## [2,] -0.3199744  0.006961379
## [3,] -0.3202125  0.006719116
## [4,] -0.3202000  0.005278056
## [5,] -0.3199703  0.006481026
## [6,] -0.3209268  0.005992327
```

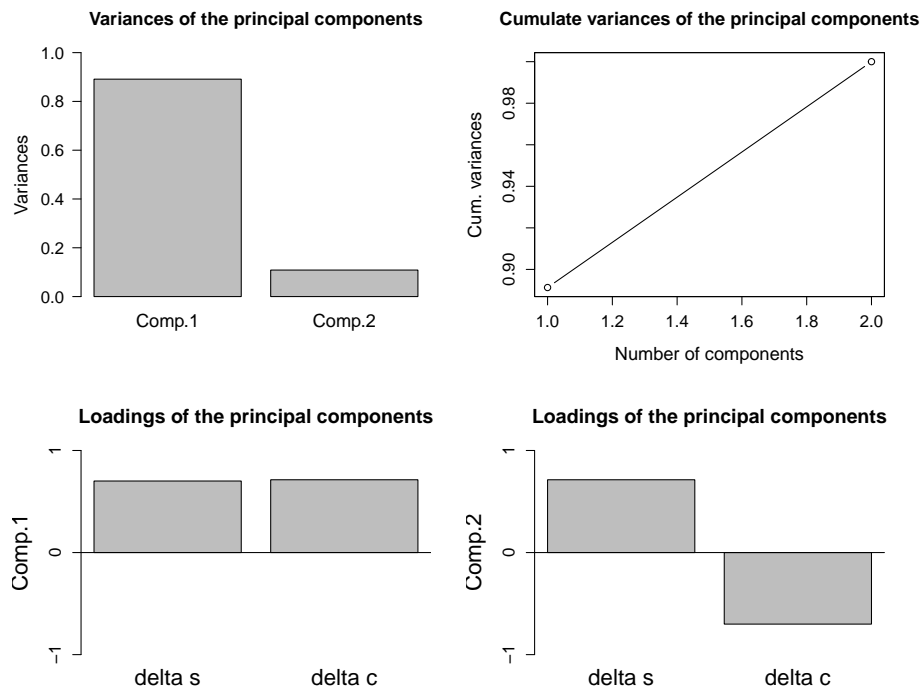


Figure 2: Principal Component Analysis of Importance Index Variances of each of the two principal components (on the top left), the cumulate proportion of variance explained by each principal component (on the top right), and loadings of the two principal components.

Looking at the value of the variance associated with the first principal component in Figure ??, this value explains 89.13% of the variability of the DELTA dataset. Moreover, from the plot of the loadings in Figure ??, it is easy to note that the first principal component is a linear combination of the variations of support and confidence, that equally contribute to the combination. So, it is reasonable to define the Importance Index as in Formula ??.

5 Visualization tools

The function `distribViz` in the *TFARM* package provides a boxplot visualization of the Importance Index distributions of a set of transcription factors (or of combinations of transcription factors).

```
# Considering for example the candidate co-regulators
# found in the set of rules r_TEAD4:

distribViz(IMP_Z, p_TFs)

## $stats
##           [,1]      [,2]      [,3]      [,4]      [,5]
## [1,] -0.3702446 0.5261723 0.5297554 -0.03725091 0.1294158
## [2,] -0.3702446 0.5567733 0.5388463 0.14509511 0.1294158
## [3,] -0.3702446 0.5775039 0.5617618 0.63009511 0.1295856
## [4,] 0.1299253 0.5876126 0.6300951 0.63009511 0.1297554
## [5,] 0.6300951 0.6028063 0.6300951 0.64543478 0.1297554
##           [,6]      [,7]      [,8]      [,9]     [,10]
## [1,] -0.06701281 0.6300951 0.6300951 -0.3702446 0.1297554
## [2,] 0.12975543 0.6300951 0.6300951 -0.3702446 0.1297554
## [3,] 0.31176178 0.6300951 0.6300951 -0.3702446 0.4709284
## [4,] 0.63009511 0.6454348 0.6300951 -0.1615716 0.6300951
## [5,] 0.64543478 0.6572690 0.6300951 -0.1615716 0.6454348
##           [,11]     [,12]
## [1,] 0.4462964 0.6300951
## [2,] 0.4462964 0.6300951
## [3,] 0.4607986 0.6300951
## [4,] 0.4753009 0.6454348
## [5,] 0.4753009 0.6649554
##
## $n
## [1] 8 28 30 22 4 15 14 28 10 12 2 30
##
## $conf
##           [,1]      [,2]      [,3]      [,4]      [,5]
## [1,] -0.64964662 0.5682955 0.5354395 0.4667194 0.1293173
## [2,] -0.09084251 0.5867122 0.5880841 0.7934708 0.1298539
##           [,6]      [,7]      [,8]      [,9]     [,10]
## [1,] 0.1076461 0.6236176 0.6300951 -0.4745059 0.2427202
## [2,] 0.5158775 0.6365726 0.6300951 -0.2659832 0.6991367
##           [,11]     [,12]
## [1,] 0.4283940 0.6256701
## [2,] 0.4932032 0.6345201
##
## $out
## [1] 0.5073143 0.3225770 0.3225770 -0.3571826 -0.3649356
## [6] 0.2353649 0.3061141 0.3152983 0.1424778 0.6454348
## [11] 0.6454348 0.3017311 0.3017311 0.6454348 0.6454348
## [16] 0.6454348 0.6454348 0.6572690 0.6572690 0.5944293
## [21] 0.5944293
```

```
##
## $group
## [1] 2 3 3 7 7 8 8 8 8 8 8 8 8 8 8 8 8 9
## [19] 9 12 12
##
## $names
## [1] "ELF1=1"      "FOSL2=1"      "GABPA=1"      "GATA3=1"
## [5] "HA.E2F1=1"   "HDAC2=1"      "MAX=1"        "MYC=1"
## [9] "NR2F2=1"     "SIN3AK20=1"   "TCF12=1"      "ZNF217=1"
```

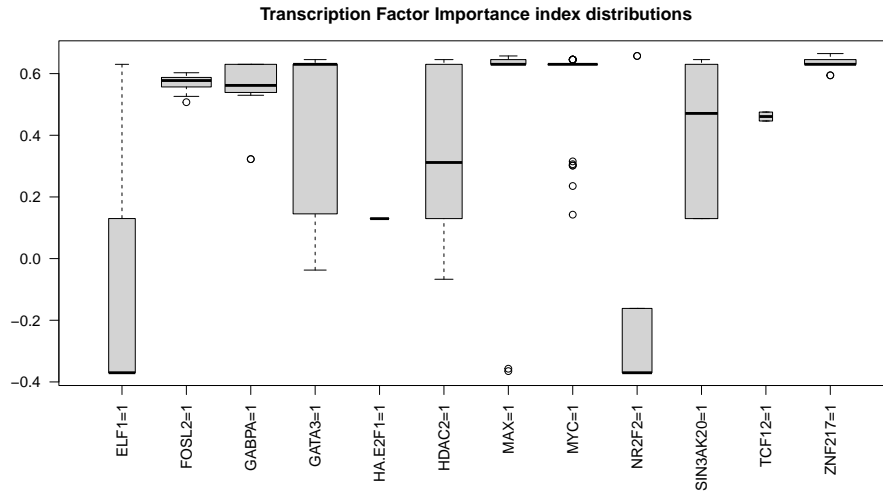


Figure 3: Importance Index distribution. Importance Index distributions of candidate co-regulators of TEAD4 in the set of the 30 most relevant associations for the prediction of TEAD4 in promotorial regions of chromosome 1 in MCF7 cell line.

The shape of boxplots changes as follows:

- The higher the number of rules containing the candidate co-regulator I , the larger the boxplot for I is
- The higher the variability of the Importance Index of I , the longer the boxplot for I is
- The higher the median of the Importance Index distribution of I , the higher the boxplot for I is aligned with respect to the y-axis.

Moreover, named q_1 and q_3 the first and third quartiles of the Importance Index distribution for a given item I , all the rules where I has importance $x \leq q_1 - 1.5 * (q_3 - q_1)$ or $x \geq q_1 + 1.5 * (q_3 - q_1)$ are considered outlier rules, and represented with a circle outside the boxplot.

For example, in the boxplots in Figure ?? it is easy to notice that:

1. SIN3AK20, HDAC2, GATA3, and GABPA have the highest median Importance Index, and they are present in a high number of relevant association rules
2. HA.E2F1 and TCF12 have intermediate median Importance Index and the lowest variability Importance Index distribution

3. ELF1 and NR2F2 are present in a high number of relevant association rules, but they have low median Importance Index and high variability of the Importance Index distribution.

It can also be noticed that for the transcription factors GABPA, MYC, MAX, NR2F2, FOSL2, and ZNF217 there are some outlier rules, that are rules in which the Importance Index of the candidate co-regulators is a lot higher or lower than the rest of the distribution. These outliers can be extracted as reported in the following text:

```
# Select the index of the list of importances IMP_Z
# containing importance distributions of transcription factor ZNF217
ZNF217_index <- which(p_TFs == 'ZNF217=1')

# Select outlier rules where ZNF217 has importance greater than 0
o <- IMP_Z[[ZNF217_index]] > 0
rule_o <- all[[ZNF217_index]]$rwi[o,]
# Display the one rule for example the sixth
rule_o[6,]

##                                lhs      rhs
## 6 {FOSL2=1,GABPA=1,MYC=1,MAX=1,ZNF217=1} {TEAD4=1}
##      support confidence      lift
## 6 0.007133152  0.6363636 0.01120924

# So, ZNF217 is very relevant in the pattern of transcription factors
# {FOSL2=1,GABPA=1,MYC=1,MAX=1,ZNF217=1}
# for the prediction of the presence of TEAD4.

# To extract support, confidence and lift of the corresponding rule
# without ZNF217:
all <- all[[ZNF217_index]]$rwo[o,]
all[6,]

##                                lhs      rhs support
## 6 {FOSL2=1,GABPA=1,MYC=1,MAX=1,ZNF217=0} {TEAD4=1}      0
##      confidence lift
## 6          0      0

# Since the measure of the rule obtained removing ZNF217 is equal to zero,
# the rule {FOSL2=1,GABPA=1,MYC=1,MAX=1,ZNF217=0} -> {TEAD4=1},
# obtained removing ZNF217, is found in the relevant rules for the prediction
# of the presence of TEAD4.
```

The function `heatI` is another useful visualization tool of the package [TFARM](#); it takes in input:

- a string vector with names of transcription factors
- a vector of mean importances of pairs of transcription factors in the previous input.

It returns a heatmap visualization of the mean importances of transcription factors in the considered string vector.

Evaluating importances of combinations of transcription factors, the number of Importance Index distribution grows combinatorially. This makes it more difficult to see which are the most critical combinations (even sorting them by their mean importances). For pairs of transcription factors, the function `heatI` gives an heatmap visualization of a square matrix whose elements are as follows (Table ??): called M the number of candidate co-regulators, the element (i,j) of such matrix is the mean Importance Index of a couple of transcription factors (TF_i, TF_j) . This matrix is symmetric with respect to the main diagonal.

	TF_1	TF_2	...	TF_{M-1}	TF_M
TF_1	$\text{imp}(TF_1)$	$\text{imp}(\{TF_1, TF_2\})$...	$\text{imp}(\{TF_1, TF_{M-1}\})$	$\text{imp}(\{TF_1, TF_M\})$
TF_2	$\text{imp}(\{TF_2, TF_1\})$	$\text{imp}(TF_2)$...	$\text{imp}(\{TF_2, TF_{M-1}\})$	$\text{imp}(\{TF_2, TF_M\})$
...					
TF_{M-1}	$\text{imp}(\{TF_{M-1}, TF_1\})$	$\text{imp}(\{TF_{M-1}, TF_2\})$...	$\text{imp}(TF_{M-1})$	$\text{imp}(\{TF_{M-1}, TF_M\})$
TF_M	$\text{imp}(\{TF_M, TF_1\})$	$\text{imp}(\{TF_M, TF_2\})$...	$\text{imp}(\{TF_M, TF_{M-1}\})$	$\text{imp}(TF_M)$

Table 4: Mean importance matrix of couples of transcription factors

To get this matrix, we need to build all possible combinations of pair of candidate co-regulators. It can be easily computed with the function `combn` in the package `combinat`. This function takes as input a vector (which is a string vector of transcription factors) and the number of required elements in the combinations. Function `combn(p, 2)` generates all pair combinations of p elements. The elements of each combination are then combined in the form TF_1, TF_2 .

```
# Construct couples as a vector in which all possible combinations of
# transcription factors (present in at least one association rules)
# are included:

couples_0 <- combn(p_TFs, 2)
couples <- paste(couples_0[1,], couples_0[2,], sep=',')
head(couples)

## [1] "FOSL2=1,SIN3AK20=1" "FOSL2=1,HDAC2=1"
## [3] "FOSL2=1,GABPA=1"    "FOSL2=1,HA.E2F1=1"
## [5] "FOSL2=1,GATA3=1"    "FOSL2=1,MYC=1"

# The evaluation of the mean Importance Index of each pair is
# then computed similarly as previously done for single transcription factors:

# Compute rulesTF, rulesTF0 and IComp for each pair, avoiding pairs not
# found in the r_TEAD4 set of rules

IMP_c <- lapply(couples, function(ci) {
  A_c <- rulesTF(ci, r_TEAD4, FALSE)
  if (all(!is.na(A_c[[1]][1]))) {
    B_c <- rulesTF0(ci, A_c, r_TEAD4, MCF7_chr1, "TEAD4=1")
    IComp(ci, A_c, B_c, figures=FALSE)$imp
  }
})

# Delete all NULL elements and compute the mean Importance Index of each pair
```

```

I_c <- matrix(0, length(couples), 2)
I_c <- data.frame(I_c)
I_c[,1] <- paste(couples)

null.indexes <- vapply(IMP_c, is.null, numeric(1))
IMP_c <- IMP_c[!null.indexes]
I_c <- I_c[!null.indexes,]

I_c[,2] <- vapply(IMP_c, mean, numeric(1))
colnames(I_c) <- colnames(IMP[,1:2])

# Select rows with mean Importance Index different from NaN, then order I_c:

I_c <- I_c[!is.na(I_c[,2]),]
I_c_ord <- arrange(I_c, desc(imp))
head(I_c_ord)

##           TF           imp
## 1 TCF12=1,ZNF217=1 0.6479683
## 2 GABPA=1,NR2F2=1 0.6370300
## 3 ZNF217=1,NR2F2=1 0.6370300
## 4 MAX=1,NR2F2=1 0.6368886
## 5 GABPA=1,MAX=1 0.6366597
## 6 MAX=1,ZNF217=1 0.6366597

# Construction of a vector in which mean Importance Index values of pairs
# of transcription factors are represented.
# These transcription factors are taken from the output of presAbs as
# present in at least one association rules.

# The function rbind is used to combine IMP columns and I_c_ord columns and
# then the function arrange orders the data frame by the imp column.

I_c_2 <- arrange(rbind(IMP[,1:2], I_c_ord), desc(imp))
p_TFs <- sub("=1", "", p_TFs)
I_c_2$TF <- sub("=1", "", I_c_2$TF)

i.heat <- heatI(p_TFs, I_c_2)

```

To build the heatmap, the user must also consider the single transcription factor mean importances (since the heatmap diagonal elements are the mean importances of single transcription factors).

The obtained heatmap is represented in Figure ???. The color scale indicates that the lowest mean importances are represented in dark red, whereas the highest ones are represented in light white.

This representation is useful to notice that, for example:

- ZNF127 has high mean Importance Index alone and in couple with all other candidate co-regulator transcription factors;

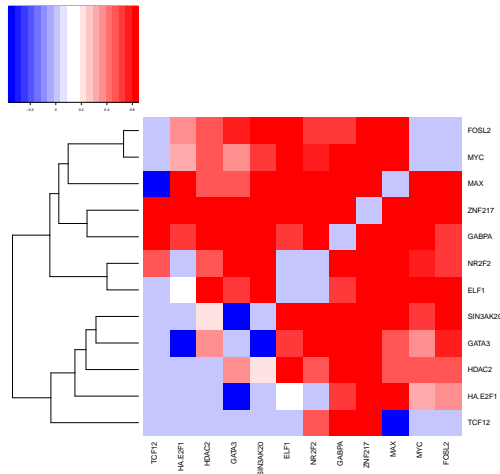


Figure 4: Heatmap. Mean importance of couples of candidate co-regulator transcription factors in the set of the 30 most relevant rules for the prediction of the presence of TEAD4 in promotorial regions of chromosome 1 in cell line MCF-7. The mean importances of single transcription factors are represented in the main diagonal as in Table ??.

- TCF12 has low mean Importance Index alone and in couple with all other candidate co-regulator transcription factors, except with GABPA, ZNF127, and NR2F2.