

Users' Guide MORE

Sonia Tarazona Campos
Maider Agueralde Martin
Blanca Tomás Riquelme

March 8, 2023

Contents

1	Introduction	2
2	Getting started	3
3	Input data	3
4	Generating the regression models with MORE	5
4.1	GetGLM input parameters	5
4.2	GetGLM output	7
4.3	Running an example	8
5	Retrieving significant regulations from MORE results	11
5.1	RegulationPerCondition input parameters	11
5.2	Interpreting RegulationPerCondition output with an example	11
6	Plotting MORE results	12
6.1	plotGLM input parameters	12
6.2	Interpretation of MORE plots	13
7	How to use MORE with R Shiny	17
7.1	MORE Shiny application example	19
8	How to cite MORE package	20

1 Introduction

One of the most common questions to be addressed when performing a multi-omics experiment is how the levels of given biological entities are being regulated by other biological entities under certain conditions. An example of this type of study would be understanding the regulatory mechanisms behind the changes in gene expression.

Potential regulators of a given gene such as miRNAs, transcription factors (TF), methylation sites, etc. can be either retrieved or predicted from public databases or obtained by a combination of experimental and computational procedures. However, a methodology for selecting the specific regulators of a particular biological system studied under certain experimental conditions is required. This is the goal of the MORE (Multi-Omics REgulation) method: modelling gene expression as a function of experimental variables, such as diseases or treatments, and the potential regulators of a given gene. The idea is to obtain more specific candidate regulators for the biological system under study by applying regression models, specifically, generalized linear models (GLM). MORE facilitates the application of GLM to multi-omic data and, although it was originally conceived to study gene expression regulation, its usage can be extended to protein or metabolite levels, for instance.

MORE requires several data inputs: gene expression data, regulators omic data, experimental design and potential associations between genes and regulators. With this input data, MORE generates the initial model equation, which is different for each gene because each one of them has different potential regulators. MORE admits numerical omic data (continuous or discrete data) or binary data.

It is strongly recommended to fit MORE models only to genes that present significant changes in any of the experimental conditions studied, that is, to differentially expressed genes (DEGs). DEGs can be selected with the standard procedures depending on the experimental design but DEGs selection is not included in the MORE algorithm and must be done by the user.

This idea can be extended to potential regulators, since regulators that do not change across conditions are not good candidates to regulate gene expression. Removing non-DE regulators will also help to reduce the number of predictors in the model since an excess of them would prevent the estimation of regression coefficients. Even so, MORE has several functionalities to filter regulators with missing values or with low variation, highly correlated regulators, and to perform variable selection.

MORE package also includes a function to retrieve the significant regulations and the magnitude of the regulatory effect under each experimental condition considered, and an additional function to graphically investigate the relationship between genes and regulators.

2 Getting started

The MORE method is available as an R package from <https://github.com/ConesaLab/MORE.git>. As for other packages in GitHub, it can be installed from R with the following instructions:

```
> install.packages("devtools")
> devtools::install_github("ConesaLab/MORE")
```

3 Input data

This section describes the main data files required by MORE to generate the regression models.

Gene expression data Expression values for each gene, in rows, under each experimental condition or replicate, in columns. MORE accepts either a **matrix** or a **data frame**. See an example below:

```
> head(TestData$GeneExpressionDE)
```

	Group1.Time1.Rep1	Group1.Time2.Rep1	Group1.Time3.Rep1	Group1.Time4.Rep1
ENSMUSG00000000078	7366	7408	6081	7045
ENSMUSG000000056999	1494	2742	3749	4070
ENSMUSG000000024873	294	323	660	757
ENSMUSG000000015461	6122	6289	6914	7633
ENSMUSG000000058135	733	3927	7317	8337
ENSMUSG000000038208	1042	1069	1052	1021
	Group1.Time5.Rep1	Group1.Time6.Rep1	Group1.Time7.Rep1	Group1.Time8.Rep1
ENSMUSG00000000078	6717	7893	7696	8072
ENSMUSG000000056999	4750	5595	7721	6896
ENSMUSG000000024873	714	665	866	950
ENSMUSG000000015461	6944	5459	4752	2783
ENSMUSG000000058135	8853	7306	5892	3549
ENSMUSG000000038208	1088	1055	1024	1098
	Group2.Time1.Rep1	Group2.Time2.Rep1	Group2.Time3.Rep1	Group2.Time4.Rep1
ENSMUSG00000000078	5920	4956	5434	4357
ENSMUSG000000056999	1618	4081	5749	7849
ENSMUSG000000024873	243	486	440	425
ENSMUSG000000015461	2304	724	108	1
ENSMUSG000000058135	1947	4318	7141	7096
ENSMUSG000000038208	1138	1205	1039	1372
	Group2.Time5.Rep1	Group2.Time6.Rep1	Group2.Time7.Rep1	Group2.Time8.Rep1
ENSMUSG00000000078	5405	4964	6762	5500
ENSMUSG000000056999	7133	7744	6763	6095
ENSMUSG000000024873	522	592	899	848
ENSMUSG000000015461	276	1315	2755	4159
ENSMUSG000000058135	7717	6573	7212	3675
ENSMUSG000000038208	1444	1111	1046	885

Experimental design Matrix or data frame containing the experimental covariates, such as treatments, diseases, strains, dose of a drug, etc. The rows of the object must be the same as the columns in **Gene expression data** and in the same order, as shown below. There is no restriction for the number of columns but it must be taken into account that MORE will combine all the experimental covariates into a

single variable. For instance, in the example below, the new single covariate would combine Time and Group2 to obtain the values: 1_0, 2_0, ..., 7_1, 8_1. Therefore, in this case, it makes more sense excluding Time from the experimental design and just including the covariate Group2.

```
> TestData$design
      Time Group2
Group1.Time1.Rep1 1      0
Group1.Time2.Rep1 2      0
Group1.Time3.Rep1 3      0
Group1.Time4.Rep1 4      0
Group1.Time5.Rep1 5      0
Group1.Time6.Rep1 6      0
Group1.Time7.Rep1 7      0
Group1.Time8.Rep1 8      0
Group2.Time1.Rep1 1      1
Group2.Time2.Rep1 2      1
Group2.Time3.Rep1 3      1
Group2.Time4.Rep1 4      1
Group2.Time5.Rep1 5      1
Group2.Time6.Rep1 6      1
Group2.Time7.Rep1 7      1
Group2.Time8.Rep1 8      1
```

Regulatory omic data This object must be a list where each element is a matrix or data frame containing the data for each “regulatory” omic (miRNA expression, transcription factor expression, etc.), with structure similar to gene expression data: regulators in rows and experimental conditions in columns (the columns must be the same as in gene expression and in the same order). See example below (TestData\$data.omics\$‘miRNA-seq’).

```
> head(TestData$data.omics$`miRNA-seq`)
      Group1.Time1.Rep1 Group1.Time2.Rep1 Group1.Time3.Rep1 Group1.Time4.Rep1 Group1.Time5.Rep1 Group1.Time6.Rep1
mmu-miR-125a-5p      7100          7229          10066          8868          12649          16321
mmu-miR-141-5p        19           49           40           107           176           183
mmu-miR-145a-3p     117329       85793       73670       57547       43926       21426
mmu-miR-148b-3p     125403       136517       137890       190417       150718       217291
mmu-miR-150-5p       1256        1098        707        1025        599        253
mmu-miR-152-3p       344         445         322         529         655        1082
      Group1.Time7.Rep1 Group1.Time8.Rep1 Group2.Time1.Rep1 Group2.Time2.Rep1 Group2.Time3.Rep1 Group2.Time4.Rep1
mmu-miR-125a-5p     14760       15164       1655       166       8       562
mmu-miR-141-5p       168         343         58         8         2         9
mmu-miR-145a-3p     11769         2       11593       1539       125       7075
mmu-miR-148b-3p     203966       203880       108489       116911       102194       94160
mmu-miR-150-5p       190           0       1289       844       617       517
mmu-miR-152-3p     1645       1441         63         17         0       121
      Group2.Time5.Rep1 Group2.Time6.Rep1 Group2.Time7.Rep1 Group2.Time8.Rep1
mmu-miR-125a-5p     2818       4997       7528       9159
mmu-miR-141-5p       13         117       121       283
mmu-miR-145a-3p     19253       32778       61536       106987
mmu-miR-148b-3p     126998       150972       171944       205767
mmu-miR-150-5p       374         370       164         0
mmu-miR-152-3p       94         270       665       687
```

Associations For each regulatory omic, associations between regulators and genes which indicate which are the potential regulators of each gene that will be consequently incorporated into the initial equation of the regression model. The association objects must be data frames and stored in a single **list** (attached bellow the example of miRNA-seq, TestData\$associations\$‘miRNA-seq’). The names of the elements of this list must be the names of the list collecting regulatory omic data and must be in the same order.

If the user wants to consider all regulators of an omic as potential regulators they must set to NULL the object of this omic in the **associations** list. Moreover, if the user does not provide the list of **associations**, all regulators of all omics in **data.omics** will be considered as potential regulators for all genes. However, this option is very time consuming. By default, NULL.

```
> head(TestData$associations$`miRNA-seq`)
      Gene      ID
1 ENSMUSG00000024873 mmu-miR-335-3p
2 ENSMUSG00000024873 mmu-miR-1912-3p
3 ENSMUSG00000024873 mmu-miR-615-5p
4 ENSMUSG00000024873 mmu-miR-322-5p
5 ENSMUSG00000024873 mmu-miR-1894-3p
6 ENSMUSG00000024873 mmu-miR-7082-3p
```

4 Generating the regression models with MORE

The **GetGLM** function in MORE adjusts a generalized linear model (GLM) for each gene (protein, metabolite etc.) in the *GeneExpression* object, to determine which regulators and experimental covariates have a significant effect on the response variable (gene expression, protein levels, etc.). These are the arguments accepted by the function, which are described in detail in Section 4.1.

```
GetGLM(GeneExpression, associations, data.omics, edesign = NULL,
center = TRUE, scale = FALSE, Res.df = 5, epsilon = 0.00001,
alfa = 0.05, MT.adjust = "none", family = negative.binomial(theta=10),
elasticnet = 0.5, stepwise = "backward", interactions.reg = TRUE,
min.variation = 0, correlation = 0.9, min.obs = 10, omic.type = 0)
```

4.1 GetGLM input parameters

GeneExpression Matrix or data frame containing gene expression data with genes in rows and experimental samples in columns. The row names must be the gene IDs.

associations List where each element corresponds to a different omic data type (miRNAs, transcription factors, methylation, etc.). The names of the elements of the list will be the omics. Each element is a data frame with two columns (optionally three) describing the potential interactions between genes and regulators for that omic. First column must contain the regulators, second the gene IDs, and an additional column can be added to describe the type of interaction (for example, in methylation data, if a CpG site is located in the promoter region of the gene, in the first exon, etc.).

data.omics List where each element corresponds to a different omic data type (miRNAs, transcription factors, methylation, etc.). The names of this list will be the omics (in

the same order as in **associations**) and each element of the list is a matrix or data frame with omic regulators in rows and samples in columns.

edesign Data frame or matrix describing the experimental design. Rows must be the samples, that is, the columns in the **GeneExpression**, and columns must be the experimental covariates to be included in the model, such as disease, treatment, etc.

center If TRUE (default), the omic data are centered.

scale If TRUE, the omic data are scaled. Default value is set to FALSE.

Res.df Number of degrees of freedom in the residuals. By default, 5. Increasing **Res.df** will increase the power of the statistical model and decrease the number of significant predictors.

epsilon A threshold for the positive convergence tolerance in the GLM model. By default, 0.00001.

alfa Significance level. By default, 0.05.

MT.adjust Multiple testing correction method to be used within the stepwise variable selection procedure. By default, "none". You can see the different options in `?p.adjust`.

family Error distribution and link function to be used in the model (see `glm` for more information). By default, `negative.binomial(theta = 10)`.

elasticnet ElasticNet mixing parameter. By default, 0.5. These are the values that can be passed to this argument:

NULL No ElasticNet variable selection is performed.

Value between 0 and 1 ElasticNet is applied with this number being the combination between ridge and lasso penalization.

Value 0 The ridge penalty.

Value 1 The lasso penalty.

stepwise Stepwise variable selection method to be applied. It can be one of: "none", "backward" (by default), "forward", "two.ways.backward" or "two.ways.forward".

interactions.reg If TRUE (default), MORE allows for interactions between each regulator and the experimental covariate.

min.variation Vector with as many elements as number of omics (names of this vector will be the omics), indicating the minimum change in standard deviation that a regulator must show across conditions in order not to be considered as having low variation and be removed from the regression models, for numerical regulators. Or the minimum change in the proportion a regulator must show across conditions, for

binary regulators. When a single value is given, the minimum change will be considered the same for all omics. The user has the option to set this value to NULL if they do not want to provide a value but are sure that they want to filter more than constant regulators across conditions. In this case, the value will be calculated as the 10% of the maximum observed variability across conditions for continuous regulators and as the 10% of the maximum observed proportion difference across conditions for binary regulators. By default, its value is 0.

correlation Correlation threshold (in absolute value) to decide which regulators are correlated, in which case, a representative of the group of correlated regulators is chosen to enter the model. By default, 0.9.

min.obs Minimum number of observations a gene must have to compute the GLM model. By default, 10.

omic.type Vector with as many elements as the number of omics, which indicates if the omic values are numeric (0, default) or binary (1). When a single value is given, the type for all the omics is set to that value. By default, 0.

4.2 GetGLM output

The object returned by the GetGLM function is a list that contains the following elements:

ResultsPerGene is a list with as many elements as genes in GeneExpression object. For each gene, there is a list containing the following information:

Y Data frame with the response variable values for that gene (y), the values fitted by the model (fitted.y), and the residuals of the GLM (residuals).

X Data frame with all the predictors included in the final model.

coefficients Matrix with the estimated coefficients for significant regulators and their p-values.

allRegulators Data frame with all the initial potential regulators in rows and the following information in columns: gene, regulator, omic, area (the third optional column in associations), filter (if the regulator has been filtered out of the model, this column indicates the reason), and Sig (1 if the regulator is significant and 0 if not). Regarding the filter column, several values are possible:

- *MissingValue*: If the regulator has been filtered out of the study because it has missing values.
- *LowVariation*: If the regulator has been filtered out of the study because it has lower variability than the threshold set by the user in min.variation parameter.

- *ElasticNet*: If the regulator has been excluded of the model by the Elastic-Net variable selection procedure.
- *Model*: When the regulator is included in the initial equation model.
- *omic_mcX_X*: For instance, *TF_mc1_R*. This notation is related to highly correlated regulators and how they are treated to avoid the multicollinearity problem. Following the *TF_mc1_R* example, it means that there are two or more transcription factors (TF omic) potentially regulating the gene that are highly correlated (in absolute value). In such cases, one of them is chosen as the representative and indicated with *_R*. The rest of them are labeled with *_P* if positively correlated with the representative and with *_N* if negatively correlated. An additional row is then added to this table, with the regulator *TF_mc1_R* and the filter label being *Model*, since only this representative is considered in the model. When there are several groups of correlated regulators for the same omic, it is indicated with *_mc1_*, *_mc2_*, etc.

significantRegulators A character vector containing the significant regulators.

GlobalSummary List that contains the following elements:

GoodnessOfFit Matrix that collects the model p-value, the final number of degrees of freedom in the residuals, the R-squared value (which for GLMs is defined as the percentage of deviance explained by the model) and the Akaike information criterion (AIC) value for all the genes that had at least a significant regulator.

ReguPerGene Matrix containing, for each omic and gene, the number of initial regulators, the number of regulators included in the initial model and the number of significant regulators.

GenesNOmodel List of genes for which the final GLM model could not be obtained. There are four possible reasons for that and they are indicated: “Too many missing values”, “No predictors after EN”, “-Inf/Inf values” and “GLM error”, where EN refers to the ElasticNet variable selection procedure.

GenesNoregulators List of genes for which there were no initial regulators, only generated in case any gene was under this condition.

Arguments List with the arguments used to generate the model: experimental design matrix, minimum degrees of freedom in the residuals, significance level, family distribution, variable selection, etc.

4.3 Running an example

An example of the execution of **GetGLM** function is shown next by using simulated data. The data file `TestData.RData` is available in the package, so the user can test and visual-

ize these results.

```
> data("TestData")
```

In this file, the gene expression matrix corresponds to the omic RNA-seq (**GeneExpressionDE**) and there is a list with three matrices of regulators in the **data.omics** object:

ChIP-seq Omic that contains binary values (1 if the TF analysed in the ChIP-seq experiment binds to the gene and 0 if not).

miRNA-seq miRNA expression data (count data).

TF TF expression data (count data).

Therefore, the **OmicType** parameter is defined as

```
> OmicType = c(1, 0, 0)
> names(OmicType) = names(data.omics)
```

The experimental design matrix (**edesign**) consists of 8 time points in two conditions, which results in a total of 16 experimental samples, but time is not to be considered as an experimental covariate since we are interested in comparing temporal profiles for the two experimental groups.

We can run the following **GetGLM** code to obtain the regression models for our genes:

```
> set.seed(123)
> SimGLM = GetGLM(GeneExpression = GeneExpressionDE,
associations = associations, data.omics = data.omics,
edesign = edesign[,-1, drop = FALSE], Res.df = 7,
epsilon = 0.00001, alfa = 0.05, MT.adjust = "fdr",
family = negative.binomial(theta = 10), elasticnet = 0.3,
stepwise = "backward", interactions.reg = TRUE,
correlation = 0.9, min.variation = NULL,
min.obs = 10, omic.type = OmicType)
```

The estimated coefficients of the significant regulators in the final model GLM for the gene ENSMUSG00000000078 and their associated p-values are

```
> SimGLM$ResultsPerGene$ENSMUSG00000000078$coefficients
               coefficient      p-value
(Intercept)    8.768718e+00 2.043440e-23
Mef2d          -2.381339e-05 2.018984e-04
TF_mc1_R       -1.339656e-04 4.201164e-03
`Group1:mmu-miR-431-5p` 3.675612e-06 1.273422e-03
`mmu-miR-18a-5p`    2.549854e-06 4.122175e-04
Hmga1          -1.099662e-05 5.564550e-03
```

The **allRegulators** table shows, for each gene, their regulators, omic, area, the kind of filter applied and the significance of the regulator. In this case (see **filter** column), in ChIP-seq, the regulators are correlated, so a random representative regulator (R) has been

chosen. In addition, it is indicated if the correlation with the representative is negative (N) or positive (P). Other possible filters are also indicated as ElasticNet, LowVariation... while Model means that the regulator was included in the model. On the other hand, the **Sig** column returns 1 if the regulator was significant in the final model and 0 if not.

```
> SimGLM$ResultsPerGene$ENSMUSG00000000078$allRegulators[1:10,]
```

	gene	regulator	omic	area	filter	Sig
13_5789384_5789613	ENSMUSG000000000078	13_5789384_5789613	ChIP-seq		ChIP-seq_mc1_P	0
13_5803679_5803880	ENSMUSG000000000078	13_5803679_5803880	ChIP-seq		ChIP-seq_mc1_P	0
13_5804696_5804875	ENSMUSG000000000078	13_5804696_5804875	ChIP-seq		ChIP-seq_mc1_R	0
13_5860779_5861047	ENSMUSG000000000078	13_5860779_5861047	ChIP-seq		ChIP-seq_mc1_N	0
13_5861542_5861793	ENSMUSG000000000078	13_5861542_5861793	ChIP-seq		ChIP-seq_mc1_P	0
13_5926774_5926959	ENSMUSG000000000078	13_5926774_5926959	ChIP-seq		ChIP-seq_mc1_P	0
mmu-miR-431-5p	ENSMUSG000000000078	mmu-miR-431-5p	miRNA-seq		Model	1
mmu-miR-18a-5p	ENSMUSG000000000078	mmu-miR-18a-5p	miRNA-seq		Model	1
mmu-miR-18b-5p	ENSMUSG000000000078	mmu-miR-18b-5p	miRNA-seq		ElasticNet	0
mmu-miR-25-3p	ENSMUSG000000000078	mmu-miR-25-3p	miRNA-seq		ElasticNet	0

```
> SimGLM$ResultsPerGene$ENSMUSG00000000078$allRegulators[90:105,]
```

	gene	regulator	omic	area	filter	Sig
Satb1	ENSMUSG000000000078	Satb1	TF		Model	0
Tbp	ENSMUSG000000000078	Tbp	TF		TF_mc1_N	1
Zfp628	ENSMUSG000000000078	Zfp628	TF		TF_mc1_N	1
Zfp64	ENSMUSG000000000078	Zfp64	TF		ElasticNet	0
Cux1	ENSMUSG000000000078	Cux1	TF		LowVariation	0
Gtf2b	ENSMUSG000000000078	Gtf2b	TF		LowVariation	0
Mafk	ENSMUSG000000000078	Mafk	TF		LowVariation	0
Pou6f1	ENSMUSG000000000078	Pou6f1	TF		LowVariation	0
Rara	ENSMUSG000000000078	Rara	TF		LowVariation	0
Rfx1	ENSMUSG000000000078	Rfx1	TF		LowVariation	0
Rreb1	ENSMUSG000000000078	Rreb1	TF		LowVariation	0
Rxra	ENSMUSG000000000078	Rxra	TF		LowVariation	0
Zfp692	ENSMUSG000000000078	Zfp692	TF		LowVariation	0
Zfp775	ENSMUSG000000000078	Zfp775	TF		LowVariation	0
ChIP-seq_mc1_R	ENSMUSG000000000078	ChIP-seq_mc1_R	ChIP-seq		ElasticNet	0
TF_mc1_R	ENSMUSG000000000078	TF_mc1_R	TF		Model	1

5 Retrieving significant regulations from MORE results

The function **RegulationPerCondition** is applied to the **GetGLM** output and returns a summary table that contains all the significant regulations, that is, all the pairs gene-regulator that were significant in MORE models. Moreover, it provides the regression coefficient that relates the gene and the regulator for each experimental condition after testing if this coefficient is significant or not.

```
RegulationPerCondition(getGLMoutput, betaTest = TRUE)
```

5.1 RegulationPerCondition input parameters

getGLMoutput Object containing the output of **getGLM** function.

betaTest If TRUE (default), a hypothesis test is performed on the coefficients of the final GLM model. If the coefficients are significant, their values are displayed. If not, it is indicated with a 0.

5.2 Interpreting RegulationPerCondition output with an example

Following with the previous example, we can run the **RegulationPerCondition** function with the argument **betaTest** equal to TRUE:

```
> myresults = RegulationPerCondition(SimGLM)
```

The output is the following table where only the last nine pairs gene-regulator are shown:

```
> myresults[1:20,]
```

	gene	regulator	omic	area	representative	Group0	Group1
Mef2d	ENSMUSG000000000078	Mef2d	TF			-2.381e-05	-2.381e-05
mmu-miR-431-5p	ENSMUSG000000000078	mmu-miR-431-5p	miRNA-seq			0.000e+00	3.676e-06
mmu-miR-18a-5p	ENSMUSG000000000078	mmu-miR-18a-5p	miRNA-seq			2.550e-06	2.550e-06
Hmga1	ENSMUSG000000000078	Hmga1	TF			-1.100e-05	-1.100e-05
Arid5a	ENSMUSG000000000078	Arid5a	TF		Rfxank	-1.340e-04	-1.340e-04
Rfxank	ENSMUSG000000000078	Rfxank	TF		Rfxank	-1.340e-04	-1.340e-04
Tbp	ENSMUSG000000000078	Tbp	TF		Rfxank	1.340e-04	1.340e-04
Zfp628	ENSMUSG000000000078	Zfp628	TF		Rfxank	1.340e-04	1.340e-04
mmu-miR-668-5p	ENSMUSG000000056999	mmu-miR-668-5p	miRNA-seq			1.453e-05	1.453e-05
mmu-miR-6931-5p	ENSMUSG000000056999	mmu-miR-6931-5p	miRNA-seq			0.000e+00	-3.536e-05
mmu-miR-543-5p	ENSMUSG000000056999	mmu-miR-543-5p	miRNA-seq			0.000e+00	-2.643e-05
mmu-miR-7648-3p	ENSMUSG000000024873	mmu-miR-7648-3p	miRNA-seq			-1.276e-05	-1.276e-05
mmu-miR-6931-5p1	ENSMUSG000000024873	mmu-miR-6931-5p	miRNA-seq		mmu-miR-6931-5p	-2.059e-05	0.000e+00
mmu-miR-145a-3p	ENSMUSG000000024873	mmu-miR-145a-3p	miRNA-seq		mmu-miR-6931-5p	-2.059e-05	0.000e+00
Tgif1	ENSMUSG000000015461	Tgif1	TF			0.000e+00	-6.945e-04
Sp4	ENSMUSG000000058135	Sp4	TF			-1.122e-04	-1.122e-04
Ep300	ENSMUSG000000058135	Ep300	TF		Ep300	-2.972e-04	-2.972e-04
Zfp513	ENSMUSG000000058135	Zfp513	TF		Ep300	-2.972e-04	-2.972e-04
Gfi1b	ENSMUSG000000038208	Gfi1b	TF			3.529e-05	3.529e-05
Tcf3	ENSMUSG000000038208	Tcf3	TF			0.000e+00	4.089e-05

This table shows the significant regulators for each gene. The **representative** column indicates if the regulator was chosen as the random representative of a correlated group

of regulators or, otherwise, which regulator was taken as the representative of the group. When no information is provided in this column, it means that the regulator was not part of a correlated group of regulators. Regulators correlated positively with the representative will have the same coefficients (same sign) than the representative, while negatively correlated regulators will have the same coefficients of the representative but with opposite sign.

The final columns correspond to the regression coefficients of each regulator for each experimental group. In this case, the experimental design matrix (**edesign**) contained two conditions, so the column `Group0` corresponds to the first condition and `Group1` corresponds to the second one. These are the conclusions we can draw from the coefficients:

- If two experimental groups have the same coefficients, it means that the regulator has the same effect on the gene in both groups.
- If one of the coefficients is 0, it means that the regulator has no effect on the gene under this experimental condition.
- Experimental groups with different non-zero coefficients indicate that the regulator affects the gene in all these experimental groups but the magnitude of the effect is not the same for all these groups.

6 Plotting MORE results

MORE package includes the function **plotGLM** to graphically represent the relationship between genes and regulators: for a given pair gene-regulator, to explore the regulators of a given gene, or to analyze which genes are regulated by a specific regulator.

```
plotGLM(GLMoutput, gene, regulator = NULL, reguValues = NULL,
plotPerOmic = FALSE, gene.col = 1, regu.col = NULL, order = TRUE,
xlab = "", cont.var = NULL, cond2plot = NULL)
```

6.1 plotGLM input parameters

GLMoutput Object generated by the function **getGLM**.

gene ID of the gene to be plotted.

regulator ID of the regulator to be plotted. If `NULL` (default value), all the regulators of the gene are plotted.

reguValues Vector containing the values of a regulator that the user can optionally provide. If `NULL` (default value), these values are taken from **GLMoutput** as long as they are available.

plotPerOmic If TRUE, all the significant regulators of the given gene and the same omic are plotted in the same graph. If FALSE (default value), each regulator is plotted in a separate plot.

gene.col Color to plot the gene. By default, 1 (black).

regu.col Color to plot the regulator. If NULL (default), a color will be assigned by the function, that will be different for each regulatory omic.

order If TRUE (default), the values in X axis are ordered.

xlab Label for X axis.

cont.var Vector with length equal to the number of observations in data, which optionally may contain the values of the numerical variable (e.g. time) to be plotted in X axis. By default, NULL.

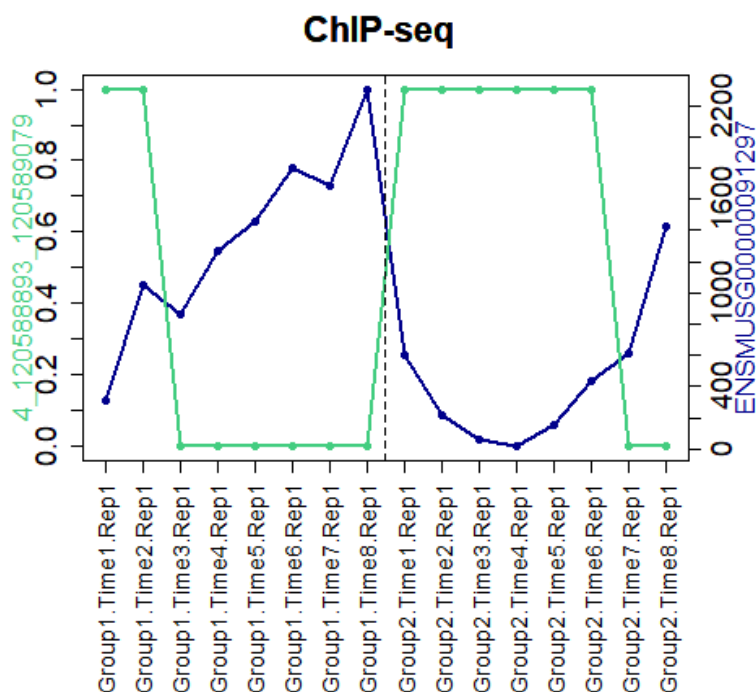
cond2plot Vector or factor indicating the experimental group of each value to represent. If NULL (default), the labels are taken from the experimental design matrix.

6.2 Interpretation of MORE plots

Following the previous example, the MORE graphic below represents the expression profile of a given gene (ENSMUSG00000091297) and the values for a significant regulator of this gene (ChIP-seq regulator 4_120588893_120589079), and can be generated with the following code:

```
> plotGLM(GLMoutput = SimGLM,
  gene = "ENSMUSG00000091297",
  regulator = "4_120588893_120589079",
  plotPerOmic = FALSE,
  gene.col = "blue4",
  regu.col = "seagreen3")
```

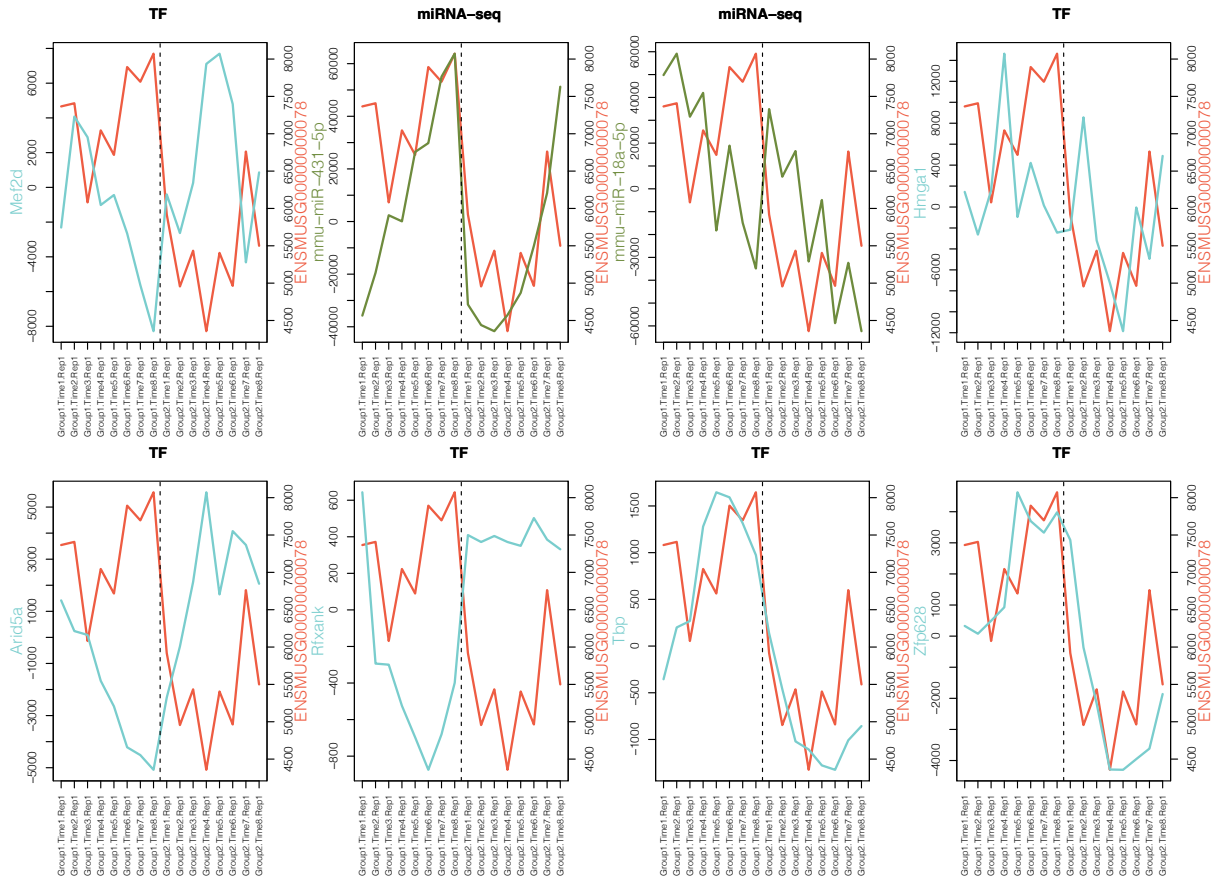
The X axis is divided in the two conditions (1 or 2) and within each condition, the observations are displayed, which correspond to different time points in this case. The right Y axis shows the expression values for the gene (plotted in blue), while the left Y axis indicates the values for the regulator (plotted in green).



If we set the regulator argument to NULL, all the significant regulators of gene ENSMUSG00000000078 will be plotted (8 regulators from two different omics: miRNA-seq and TF).

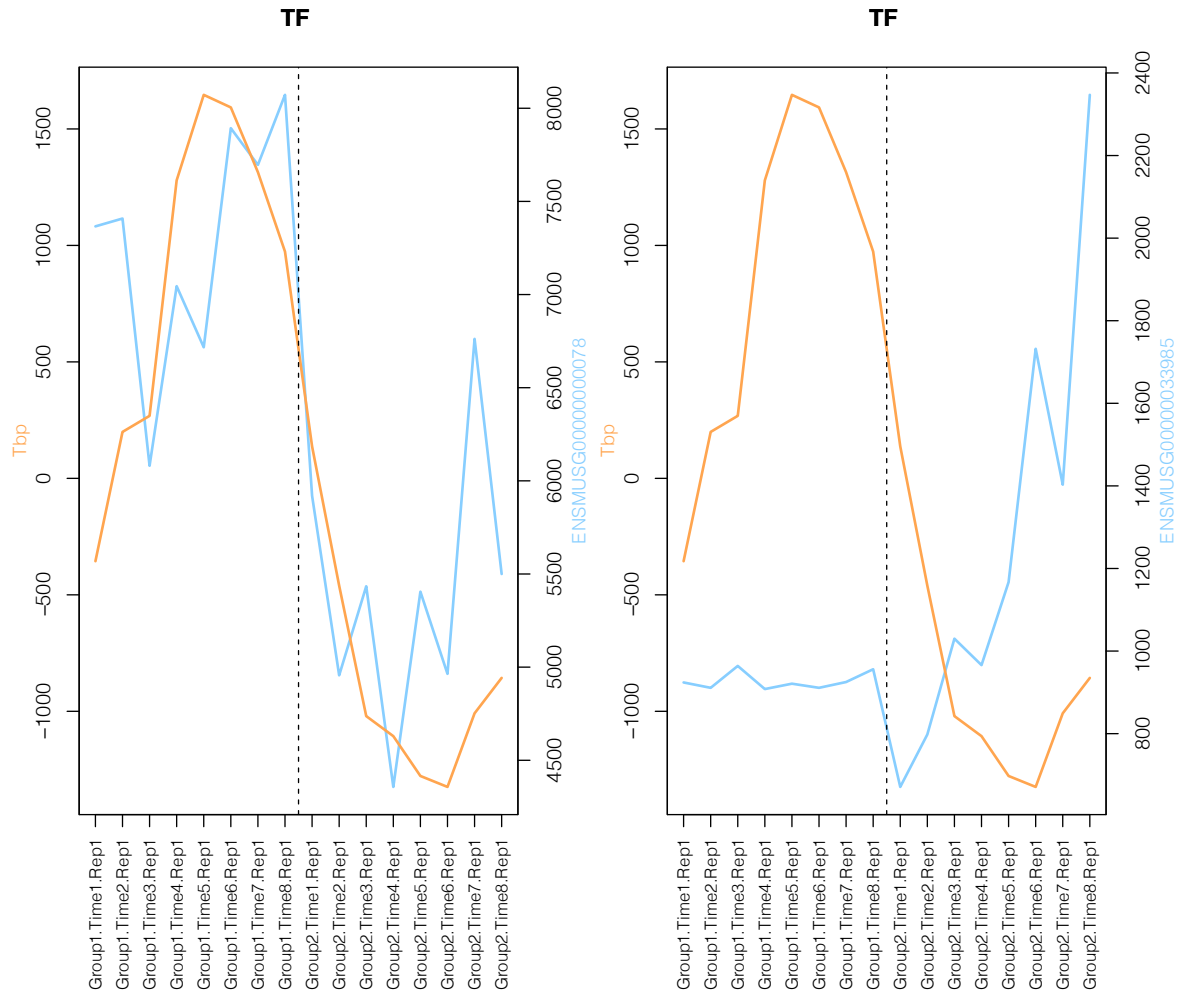
```
> par(mfrow = c(2,4))
> plotGLM(GLMoutput = SimGLM,
  gene = "ENSMUSG00000000078",
  regulator = NULL,
  plotPerOmic = FALSE,
  gene.col = "tomato2")
```

The title of each plot indicates the omic represented in that plot. The values for significant regulators are plotted in different colors according to the omic. The values for the gene are plotted in red, as indicated in the previous code.



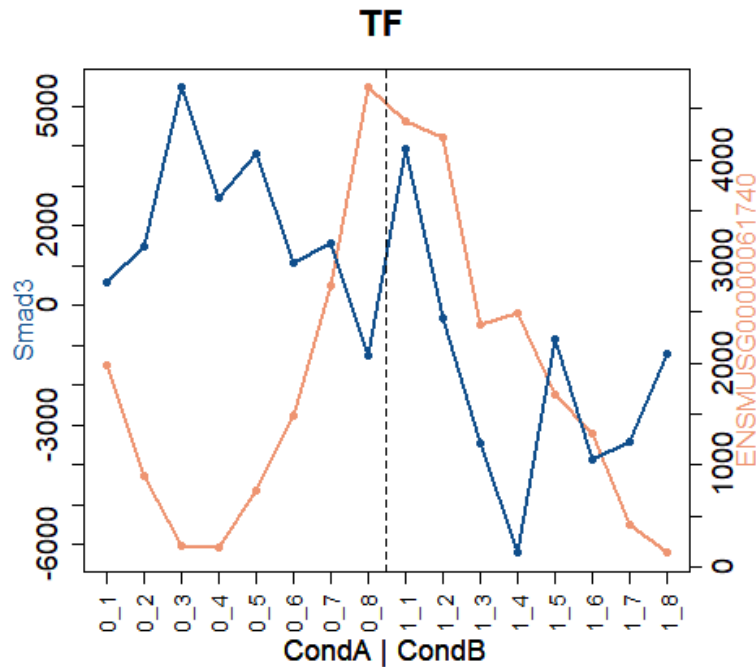
If we want to plot all the genes that are being significantly regulated by a given regulator (e.g. Tbp), we must set the gene argument to NULL as follows. In this case, the TF regulates two genes: ENSMUSG00000000078 and ENSMUSG00000033985.

```
> par(mfrow = c(1,2))
> plotGLM(GLMoutput = SimGLM,
  gene = NULL,
  regulator = "Tbp",
  plotPerOmic = FALSE,
  gene.col = "skyblue1",
  regu.col = "tan1")
```



Users can also define their own values for time points with the vector **cont.var**. In addition, they can assign a label for axis X to differentiate between two conditions, A or B, which is the **xlab** parameter. The code and the resulting graph, where the gene is plotted in light orange and the regulator is plotted in blue, can be found below.

```
> plotGLM(GLMoutput = SimGLM,
  gene = "ENSMUSG000000061740",
  regulator = "Smad3",
  plotPerOmic = FALSE,
  gene.col = "lightsalmon2",
  regu.col = "dodgerblue4",
  cont.var = c(1,2,3,4,5,6,7,8),
  xlab = "CondA | CondB")
```

7 How to use MORE with R Shiny

Shiny is an R package that allows for building web applications from R packages or scripts so users that are not familiar with R language can still easily use R packages. We have generated a MORE web application with R Shiny with this purpose. To use the MORE shiny tool, users must first install the Shiny and the Shiny themes packages from CRAN repository with `install.packages()` function.

```
> install.packages("shiny")
> install.packages("shinythemes")
```

Moreover, users must also have previously installed the MORE package described in section 2.

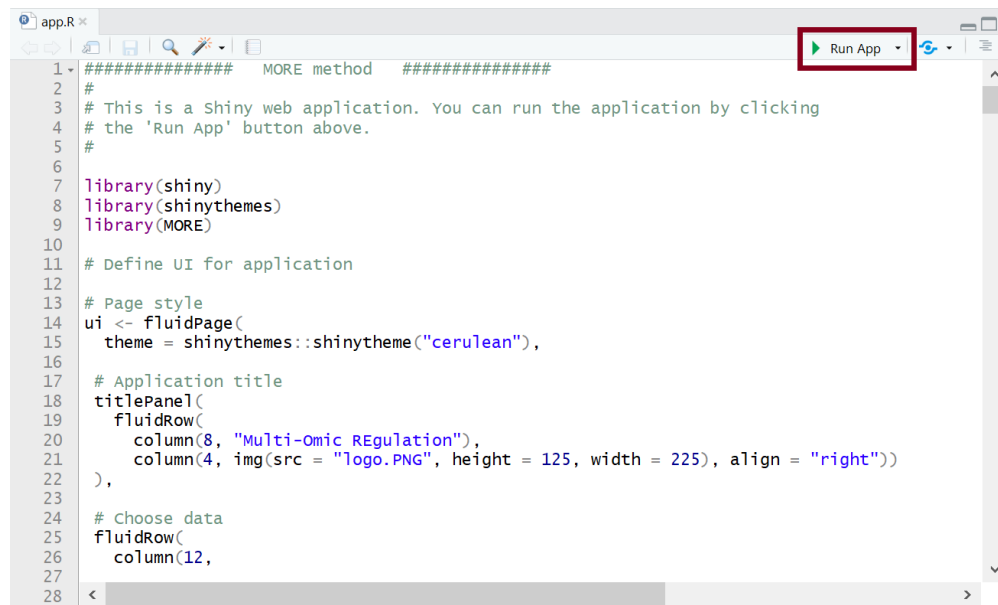
The MORE shiny scripts needed to run the tool are available in the Downloads bitbucket folder for MORE package (<https://bitbucket.org/ConesaLab/more/downloads/>). There is a file called `app.R`, an example dataset stored in the file `TestDataShiny.RData` and a folder called `www`.

app.R Script to run the web application for MORE method. Users must open this file from RStudio to start using the application .

TestDataShiny.RData Example data file to test the application. We used it as an example of an execution of MORE Shiny.

www Folder containing the style options. Please, do not delete anything in this folder.

Therefore, in order to run the application, please open the `app.R` file, where the MORE method application is located, and execute it using the button `Run App` (the red box in the following picture).



```
1 ##### MORE method #####
2 #
3 # This is a Shiny web application. You can run the application by clicking
4 # the 'Run App' button above.
5 #
6
7 library(shiny)
8 library(shinythemes)
9 library(MORE)
10
11 # Define UI for application
12
13 # Page style
14 ui <- fluidPage(
15   theme = shinythemes::shinytheme("cerulean"),
16
17   # Application title
18   titlePanel(
19     fluidRow(
20       column(8, "Multi-Omic REGulation"),
21       column(4, img(src = "logo.PNG", height = 125, width = 225), align = "right"))
22     ),
23
24   # Choose data
25   fluidRow(
26     column(12,
27
28
```

Figure 1: Run MORE application, click `Run App` button

A window will open with the MORE application, as shown in the following figure.

Multi-Omic REGulation

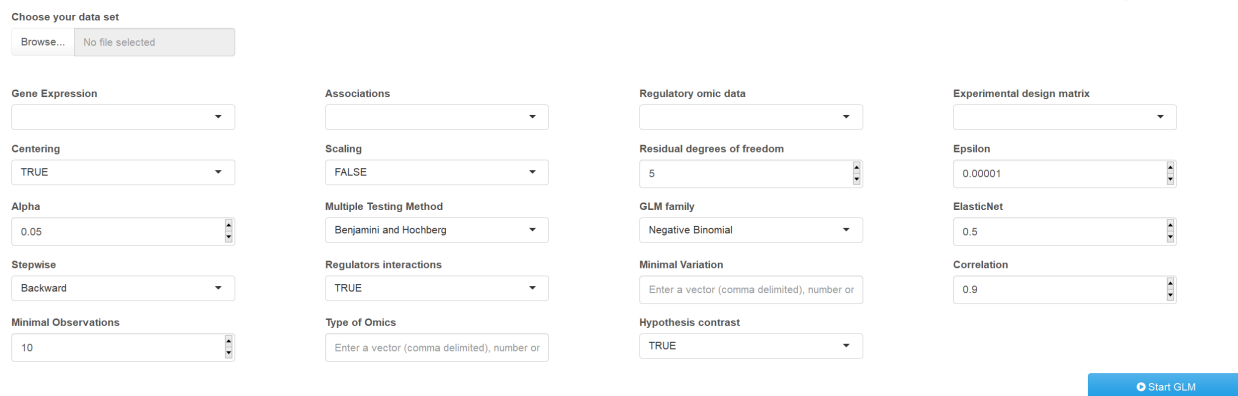


Figure 2: MORE Shiny

The different options have been explained in the previous sections but, next, we run an

example to clarify how to use them. This example can also be visualized in this video: <https://youtu.be/SSIaeFRNsXg>.

7.1 MORE Shiny application example

Please take into account that MORE Shiny only supports `.RData` files at this time. Once the `RData` is loaded, the user must indicate the names of the data files in the `RData` that corresponds to (see the input data defined in the **section 3**): gene expression matrix, experimental design matrix, regulators matrix and association matrix. In summary (for more details see **section 3**):

Gene expression matrix Matrix or data frame that contains expression values for each gene, in rows, under each experimental condition or replicate, in columns.

Experimental design matrix Matrix or data frame that contains the experimental covariates, such as treatments, points of time...

Regulatory omic data List that contains matrices or data frames containing the data for each regulatory omic, e.g. miRNA expression. The data frames structure is similar to gene expression data.

Association matrix List that contains data frames with the potential regulators for each regulatory omic considered. The association objects must be data frames and stored in a single list.

In this case, the file `TestDataShiny.RData` contains the objects described in **section 4.3**, unlike the experimental design matrix, that contains only one column with two conditions. By clicking on the button `Browse...`, users can choose their own data file (see Figure 3 blue box). Once the data is loaded, the user must enter the same input parameters of the `GetGLM` and `RegulationPerConditon` functions as defined in **section 4.1** and **section 5.1**.

It should be taken into account that if users want to enter a `NULL` value for a given parameter, they must leave the box empty.

In the example of the application, we will consider the input parameters shown in the following figure, leaving blank those we want to be `NULL`.

Now, clicking the button `Start GLM`, we will obtain a summary table (see Figure 4). Specifically, this is the table defined in the **section 5.2**, the output of the **RegulationPerConditon()** function. The user can download the table in `csv` format by pressing the button `Download` (see Figure 4 orange box). It is necessary to save the file that will contain the table with name and extensión **.csv**.

The button `MORE plots` (see Figure 4) will generate plots to visualize the relationship between genes and regulators. The user can change the different parameters without re-executing the application to tune the plots or plot new elements.

Figure 3: MORE Shiny: inputs for running example

Here we show the example for the gene ENSMUSG00000000078 (orange) and TF Mef2d (blue). Pressing the button `Generate Plot`, Shiny generates the first graph. However, if it is expected to obtain more than one graphic, the user can see all of them in a pdf file. This pdf file is generated by pressing the button `Download` (box orange in Figure 5) and saving the document, for example, as [plotsGLM.pdf](#). It is essential to save the file with name and extension **.pdf**.

Show entries

gene	regulator	omic	area	representative	Group0	Group1
ENSMUSG00000000078	mmu-miR-378a-5p	miRNA-seq			-1.224e-04	-1.224e-04
ENSMUSG00000000078	mmu-miR-139-3p	miRNA-seq			-3.021e-02	-3.021e-02
ENSMUSG00000000078	Mef2d	TF			-2.105e-05	-2.105e-05
ENSMUSG000000024873	mmu-miR-7081-5p	miRNA-seq			-3.052e-01	-3.052e-01
ENSMUSG000000024873	mmu-miR-34a-5p	miRNA-seq			-1.509e-04	-1.509e-04
ENSMUSG000000024873	mmu-miR-291a-5p	miRNA-seq			-8.512e-06	-8.512e-06
ENSMUSG000000024873	mmu-miR-879-5p	miRNA-seq			6.577e-04	6.577e-04
ENSMUSG000000015461	mmu-miR-376c-3p	miRNA-seq			3.343e-05	3.343e-05
ENSMUSG000000015461	Amt	TF			3.965e-04	3.965e-04
ENSMUSG000000058135	mmu-miR-665-3p	miRNA-seq			-5.274e-01	-5.274e-01

gene regulator omic area representative Group0 Group1

Showing 1 to 10 of 69 entries

[MORE plots](#)

Previous 1 2 3 4 5 6 7 Next

[Download](#)

Figure 4: Summary table and MORE plots button. The user can download the whole table by pressing the button `Download` (orange box)

8 How to cite MORE package

Tarazona, S., Tomás-Riquelme, B., Martínez-Mira, C., Clemente-Císcar, M., Conesa, A. (2018). MORE: Multi-Omics REgulation by regression models. R package version 0.1.0.

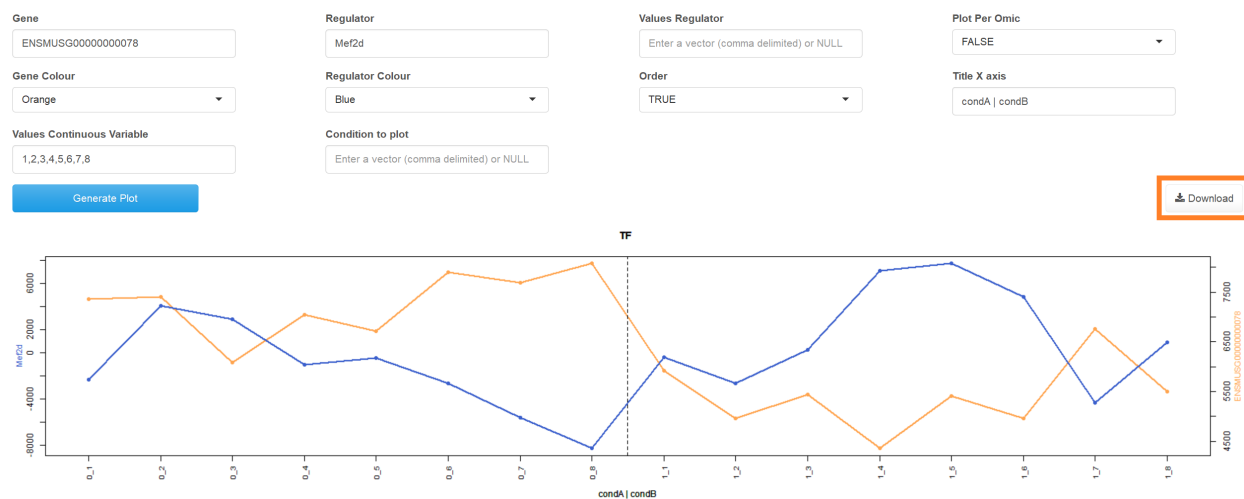


Figure 5: Inputs and plot for the pair gen–regulator (ENSMUSG00000000078–Mef2d). The user can download all plots by pressing the button `Download` (orange box)