

# MS-CleanR tutorial:

Peak list cleaning, data concatenation and peak annotation

03/03/2020

guillaume.marti@univ-tlse3.fr

## 1. Prerequisite :

### Software installation

#### 1.1. Downloading

**MS-DIAL** version up to 4.16:

[http://prime.psc.riken.jp/Metabolomics\\_Software/MS-DIAL/index2.html](http://prime.psc.riken.jp/Metabolomics_Software/MS-DIAL/index2.html)

**MS-FINDER** version up to 3.30:

[http://prime.psc.riken.jp/Metabolomics\\_Software/MS-FINDER/index2.html](http://prime.psc.riken.jp/Metabolomics_Software/MS-FINDER/index2.html)

**R version up to 3.6.1** : <https://cran.r-project.org/>

**R studio**: <https://rstudio.com/products/rstudio/>

#### 1.2. Installation

- In **R**, copy and paste the following command to update R version if necessary

```
if(!require(installr)) {  
install.packages("installr"); require(installr)}  
updateR()
```

- In **R studio**, update all your packages with the command

```
SetRepositories()
```

Select 1 and 2 for CRAN and BIOCONDUCTOR packages

Select the command Update on the right windows in the Package part

- Install MS-cleanR by copying and pasting the command :

```
devtools::install_github("SyrupType/mscleanr", auth_token =  
"a29ed88006a073473a47cab074f9bc7e8dfacbac")
```

## 2. MS-CleanR workflow

Within your project directory, create one subfolder for each ionization mode namely “pos” and “neg”. In each of this new directory, create another subfolder named “peaks”.

### 2.1. Process the data with MS-DIAL



Process data with MS-DIAL in either pos or neg mode or both according to the tutorial <https://mtbinfo-team.github.io/mtbinfo.github.io/>

#### Important notices:

- A) During data importation, it is important to note the type (Blank, QC or Sample) and class of every sample in **Class ID column** (blank, sample class, QC)

File name	File type	Class ID	Batch	Analytical order	Injection volume (µL)	Y variable	Included
BLANC-M-POSN	Blank	blank	1	7	1	0	<input checked="" type="checkbox"/>
BLANC-P-POSN	Blank	blank	1	8	1	0	<input checked="" type="checkbox"/>
BLANC-Q-POSN	Blank	blank	1	9	1	0	<input checked="" type="checkbox"/>
BLANC-T-POSN	Blank	blank	1	10	1	0	<input checked="" type="checkbox"/>
BLANC-U-POSN	Blank	blank	1	11	1	0	<input checked="" type="checkbox"/>
BLANC-X-POSN	Blank	blank	1	12	1	0	<input checked="" type="checkbox"/>
BLANC-Y-NEG	Blank	blank	1	13	1	0	<input checked="" type="checkbox"/>
blc-neg-1	Blank	blank	1	14	1	0	<input checked="" type="checkbox"/>
CAM1-POS-1N	Sample	CAM1	1	15	1	0	<input checked="" type="checkbox"/>
CAM1-POS-2N	Sample	CAM1	1	16	1	0	<input checked="" type="checkbox"/>
CAM1-POS-3N	Sample	CAM1	1	17	1	0	<input checked="" type="checkbox"/>
CAM1-POS-4N	Sample	CAM1	1	18	1	0	<input checked="" type="checkbox"/>
CAM1-POS-5N	Sample	CAM1	1	19	1	0	<input checked="" type="checkbox"/>
CAM1-POS-6N	Sample	CAM1	1	20	1	0	<input checked="" type="checkbox"/>
CAM1-POS-7N	Sample	CAM1	1	21	1	0	<input checked="" type="checkbox"/>
CAM1-POS-8N	Sample	CAM1	1	22	1	0	<input checked="" type="checkbox"/>
CAM1-POS-9N	Sample	CAM1	1	23	1	0	<input checked="" type="checkbox"/>
CAM2-POS-1N	Sample	CAM2	1	24	1	0	<input checked="" type="checkbox"/>
CAM2-POS-2N	Sample	CAM2	1	25	1	0	<input checked="" type="checkbox"/>
CAM2-POS-3N	Sample	CAM2	1	26	1	0	<input checked="" type="checkbox"/>
CAM2-POS-4N	Sample	CAM2	1	27	1	0	<input checked="" type="checkbox"/>
CAM2-POS-5N	Sample	CAM2	1	28	1	0	<input checked="" type="checkbox"/>
CAM2-POS-6N	Sample	CAM2	1	29	1	0	<input checked="" type="checkbox"/>
CAM2-POS-7N	Sample	CAM2	1	30	1	0	<input checked="" type="checkbox"/>
CAM2-POS-8N	Sample	CAM2	1	31	1	0	<input checked="" type="checkbox"/>
CAM2-POS-9N	Sample	CAM2	1	32	1	0	<input checked="" type="checkbox"/>
QC-ALL-POS-1N	QC	QC	1	33	1	0	<input checked="" type="checkbox"/>
QC-ALL-POS-2N	QC	QC	1	34	1	0	<input checked="" type="checkbox"/>
QC-ALL-POS-3N	QC	QC	1	35	1	0	<input checked="" type="checkbox"/>
QC-ALL-POS-4N	QC	QC	1	36	1	0	<input checked="" type="checkbox"/>
QC-ALL-POS-5N	QC	QC	1	37	1	0	<input checked="" type="checkbox"/>
QC-ALL-POS-6N	QC	QC	1	38	1	0	<input checked="" type="checkbox"/>
TAK1-NEG-1	Sample	TAK1	1	39	1	0	<input checked="" type="checkbox"/>
TAK1-POS-2-N	Sample	TAK1	1	40	1	0	<input checked="" type="checkbox"/>
TAK1-POS-3N	Sample	TAK1	1	41	1	0	<input checked="" type="checkbox"/>
TAK1-POS-4N	Sample	TAK1	1	42	1	0	<input checked="" type="checkbox"/>
TAK1-POS-5N	Sample	TAK1	1	43	1	0	<input checked="" type="checkbox"/>
TAK1-POS-6N	Sample	TAK1	1	44	1	0	<input checked="" type="checkbox"/>

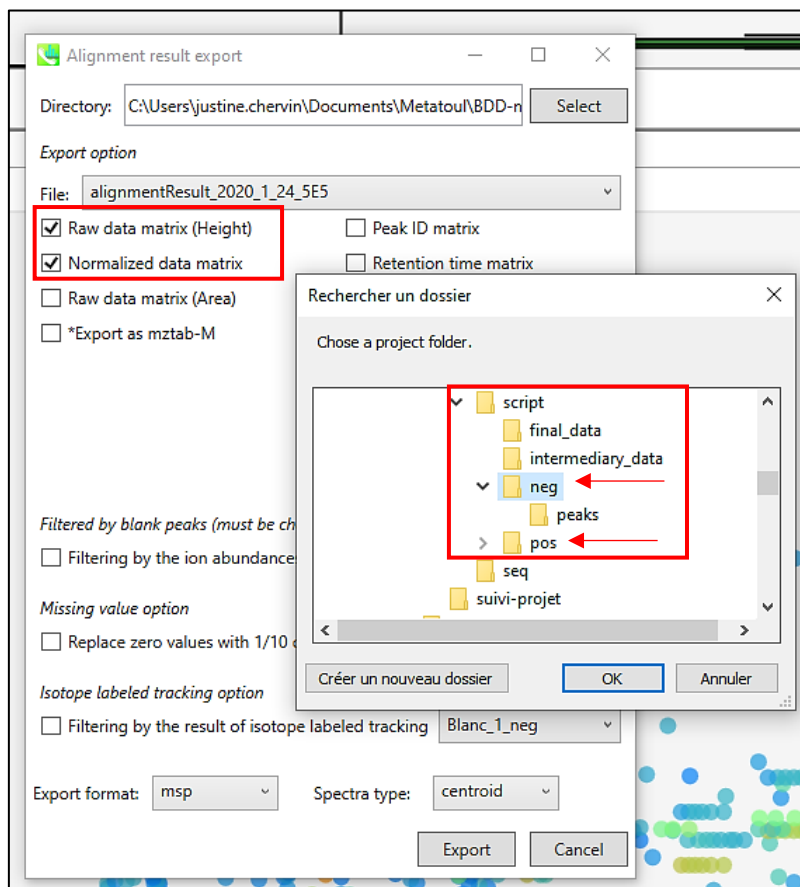
- B) Be careful to have the **same number of samples** between pos and neg mode and in the **same order**.

### 2.2. Export peak list

After alignment process:

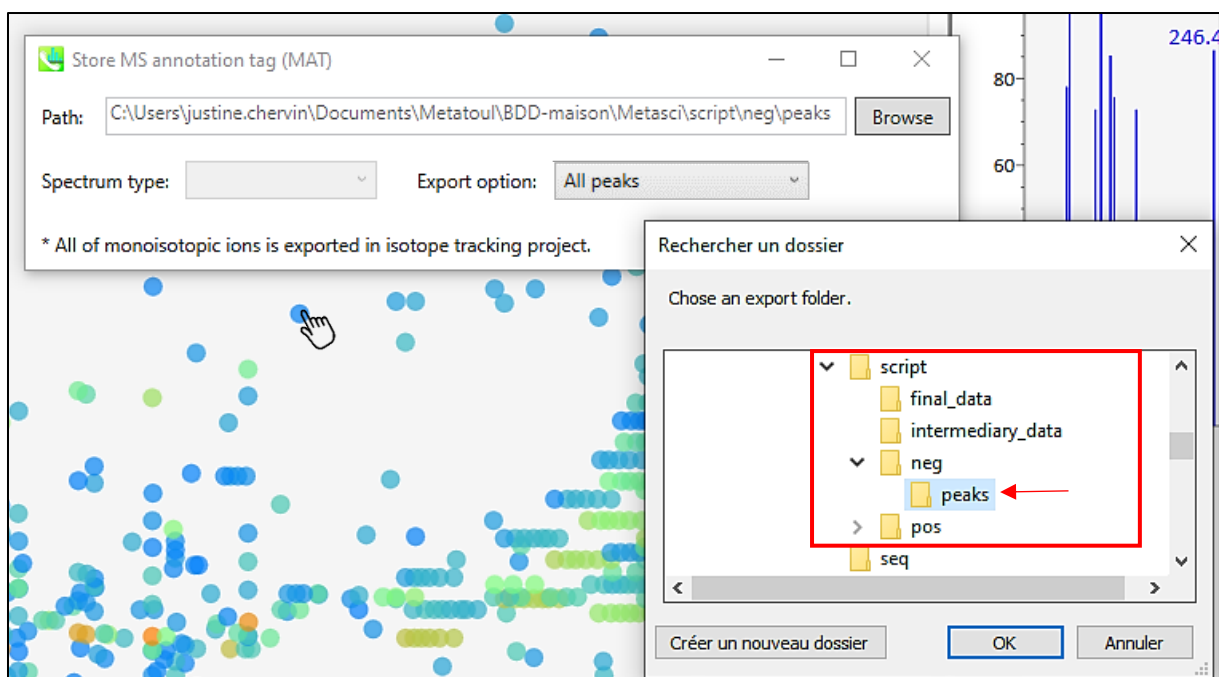
- **Normalized data** by Total ion chromatogram (TIC) or another method

- Export alignment results: both **Raw data matrix (Height)** and **Normalized data matrix** respectively in previously created folders named “pos” and “neg”.



### 2.3. Export all peaks

By clicking on one point, export « **all peaks** » to the “peaks” directory respectively created in “pos” and “neg” folders.



## 2.4. Open the shiny interface of MS-CleanR



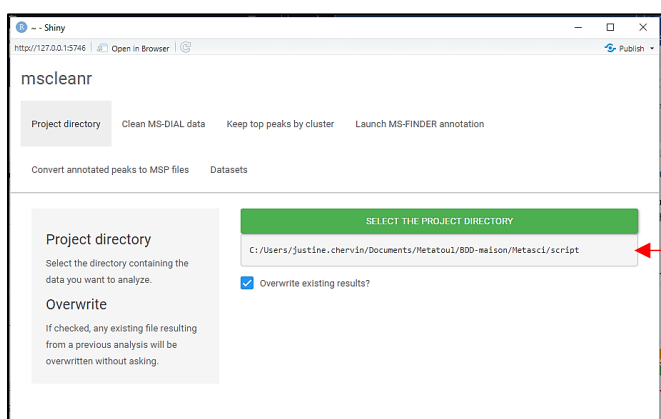
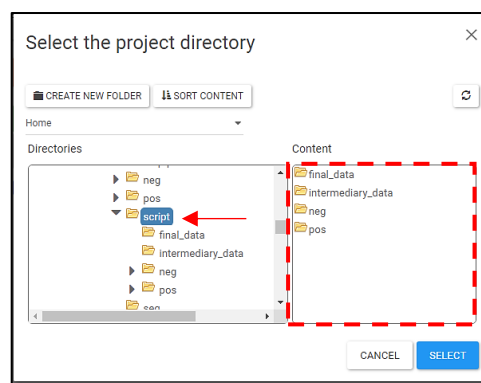
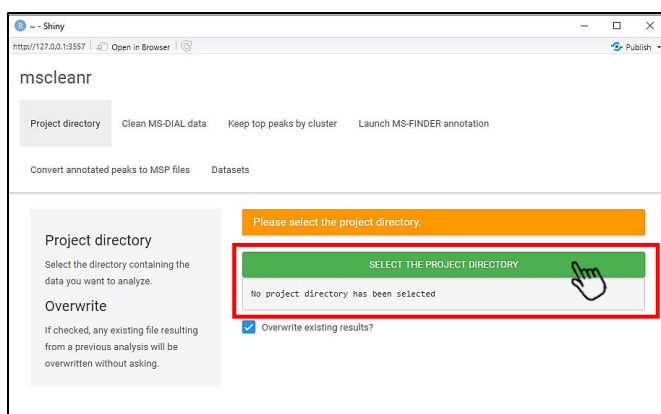
Select the MScleanR package in **Rstudio** and open the shiny interface using the following command:

- *Note that if you encounter some issues, try to open the Shiny interface in internet browser.*
- *Sometimes Windows block file writing, close the shiny or R studio and run it again to solve the problem.*

```
runGUI()
```

### 2.4.1. Select the project directory

First step is to define the project directory on the first index called “**Project Directory**” by clicking on the green rectangle “*Select the project directory*” and by selecting the parent folder containing “pos and “neg” directories.



When your project directory is selected, it is written in the grey rectangle.

#### 2.4.2. Define your parameter of filtration and Clean your data

In the second index called “**Clean MS-DIAL data**” various parameters can be personalized to filter your data. You can decide to select any filter according to your goal and experimental design.

Command	Description
<b>Blank ratio</b>	Subtract blank peaks to samples based on the indicated “ <b>Minimum blank ratio</b> ” by default at 0.8. This operation is done on the <b>Height files</b> between Blanks and QCs.
<b>Incorrect Mass</b>	Delete all peaks with a mass defect in X.8 and X.9 which appear to be artifacts.
<b>Relative standard Deviation (RSD)</b>	Filter based on the <b>Maximum RSD</b> value set at 30 by default. The RSD is calculated on each defined class. If RSD of one feature is under 30 for all class, it is removed from the peak list.
<b>Relative Mass Defect (RMD)<sup>1</sup></b>	RMD is calculated in ppm as ((mass defect/measured monoisotopic mass) × 10e6) Analysis of natural products from the DNP shows that 95 % of RMD are comprised between 50 and 3000 (values by default).
<b>Delete ghost peaks</b>	Delete variables with <i>m/z</i> values corresponding to blank peaks but with a different RT in samples.

<sup>1</sup> Ekanayaka EA, Celiz MD, Jones AD. Relative mass defect filtering of mass spectra: a path to discovery of plant specialized metabolites. *Plant Physiol.* 2015;167(4):1221–1232. doi:10.1104/pp.114.251165

<b>Maximum mass difference</b>	m/z value tolerance set by default to 0.005 for Pearson correlation and pos/neg merging
<b>Maximum retention time difference</b>	RT value tolerance set by default to 0.025 (absolute value) for Pearson and pos/neg merging
<b>Use Pearson correlation to compute clusters?</b>	Extend MS-DIAL clusters with Pearson correlation. <b>Minimum correlation</b> and <b>maximum p-value</b> are respectively set by default to 0.8 and 0.05

Once your parameters are fixed, click on the green rectangle named “*Clean MS-DIAL data*”. A green window appears with the writing “*Cleaning data...*”.

Cleaning data...

During the cleaning:

- Clusters are formed based on MS-DIAL “post curation column”, Pearson correlation, links such as adducts, neutral losses, dimers, ...;
- Adducts are corrected based on previous found links;
- Pos and Neg clusters are concatenated if relational links are found (adducts mass difference)
- Once the cleaning is done, one new folder is created named “intermediary\_data”. Different information is obtained at the bottom of the index “Clean MS-DIAL data”.

Shiny

http://127.0.0.1:4872 | Open in Browser

Check what filters you want to use to clean your ms data.

**Deltas**

Indicates the acceptable retention time and mass differences to consider that peaks are related.

**Clusterisation options**

You can choose to use the Pearson correlation between peaks as a supplementary data used during clusterisation

**(Optional) Reference files**

Optionally, you can use your own files for adducts and neutral losses. See the documentation for more information.

Maximum mass difference  
0.005

Maximum retention time difference  
0.025

☒ Use Pearson correlation to compute clusters?

Minimum correlation  
0.8

Maximum p-value  
0.05

You can optionally import personal reference files for adducts and neutral losses. By default, data displayed in the Datasets tab will be used.

☐ Use personal reference files?

**CLEAN MS-DIAL DATA**

```

/!\ Deleting C:/Users/justine.chervin/Documents/Metatoul/BD0-maison/Metasci/script/final_data
/!\ Deleting C:/Users/justine.chervin/Documents/Metatoul/BD0-maison/Metasci/script/intermediary_data
*** Treating C:/Users/justine.chervin/Documents/Metatoul/BD0-maison/Metasci/script ***
MSDial peaks after filtering: 621 positive, 266 negative, 0 NA, 887 total
MSDial links: 1978
73 peaks identified by MSDial
Correlation links found in pos : 1939
Correlation links found in neg : 323
Clusters detected with MSDial data: 180
Using package neutral losses for positive mode
Adducts/Neutral losses detection (cluster pos 91 )
Adducts/Neutral losses detection (cluster pos 92 )
Adducts/Neutral losses detection (cluster pos 93 )
Adducts/Neutral losses detection (cluster pos 94 )
Adducts/Neutral losses detection (cluster pos 95 )
Adducts/Neutral losses detection (cluster pos 96 )
Adducts/Neutral losses detection (cluster pos 97 )
Adducts/Neutral losses detection (cluster pos 98 )
Adducts/Neutral losses detection (cluster pos 99 )
Adducts/Neutral losses detection (cluster pos 100 )
Adducts/Neutral losses detection (cluster pos 101 )

```

Delete previous results if necessary

Number of final peaks  
Number of MS-DIAL links  
Number of MS-DIAL identification

Adduct / neutral loss relations

Adduct correction if necessary

At this step, several files are created in the folder “intermediary\_data”.

Files	Description
Adducts_massdiff_filtered	Reference file for mass difference between regular adducts
Adducts_massdiff_total	Reference file for mass difference between all possible adducts
Adducts_detected_by_MS-DIAL	Reference file for adduct ponderation of regular adducts found by MS-DIAL
Adducts_filtered.graphml	A graph to display feature clusters based on adducts links
Adducts_final_selection	Final adducts resulting from MSdial and modified after pos/neg concatenation
Adducts_initial.graphml	A graph to display feature clusters based on MSdial data
Annotated_MS-peaks-MSDial	List of annotated peaks based on the database (msp file) imported in MS-DIAL
Deleted_blank_ghosts	List of peaks deleted with “delete ghost peaks”
Deleted_blanks	List of peaks deleted with the filter “blank ratio”
Deleted_mz	List of peaks deleted with the filter “incorrect mass”
Deleted_rmd	List of peaks deleted with the filter “RMD”

Deleted_rsd	List of peaks deleted with the filter “RSD”
Links_clusters_final	List of correlation (adduct, neutral loss, msdial) between peaks in neg and pos
Links_post_selection	Feature links after adduct prioritization process
Links_pre_selection	Feature links with all adducts possibilities
MS_peaks-clusters.graphml	A graph of final clusters (MS-DIAL + Pearson)
MS_peaks-clusters_final	List of final clusters (MS-DIAL + Pearson) in both pos and neg ionization
MS_peaks-clusters_msdial	List of MS-DIAL clusters in both pos and neg ionization
parameters	List of parameter used for the cleaning
samples	List of samples with indication of sample name, class, file type, script class and column name

### 2.4.3. Select number of peaks per cluster

In the third index “**Keep top peaks by cluster**” you can select the number of features you want to keep in each cluster.

This step is based on the hypothesis that in one cluster, only one true metabolite is present. The other variables used to come from feature degeneration. Generally, this metabolite appears to be the **most intense** and/or **the most connected within the graph** (adducts, neutral loss, dimers...).

You can then choose to select as many peaks as you want and either the most intense(s) by clicking “**Intensity**”, the most connected by clicking “**Degree**” or both. We advise to select both criteria and keep 2 top peaks by cluster for further MS-finder request.

mscleanr Project directory Clean MS-DIAL data **Keep top peaks by cluster** Launch MS-FINDER annotation Convert annotated peaks to MSP files Datasets

Keeps only the top peaks by cluster and by method.

**Selection mode**  
Peaks selected can be the most intense (intensity), most connected (degree), or both. If there are ties, the number of peaks selected can be greater than the number requested by the user.

**Exporting filtered peaks**  
Copy MAT files corresponding to the selected peaks in a new folder for an faster analysis in MS-FINDER.

Selection criteria  
☒ Intensity (most intense peaks)  
☐ Degree (most connected peaks)

Number of peaks to keep (by cluster and by method)  
 1

☒ Export final peaks in a new folder?

**KEEP TOP PEAKS BY CLUSTER**

Keeping only selected peaks...



At this step, a new folder is created in both “pos” and “neg” folders named “**filtered peaks**”. All .MAT files corresponding to kept peaks are copied from “peaks” folder and pasted in this new folder “filtered peaks”.

The screenshot shows the mscleanr Shiny application interface. The top navigation bar includes 'Project directory', 'Clean MS-DIAL data', 'Keep top peaks by cluster' (the active step), 'Launch MS-FINDER annotation', and 'Convert annotated peaks to MSP files'. Below this, the 'Datasets' section is visible. The main content area is divided into three columns. The left column contains a description: 'Keeps only the top peaks by cluster and by method.' and a 'Selection mode' section explaining that peaks can be selected by intensity, degree, or both. The middle column shows 'Selection criteria' with checkboxes for 'Intensity (most intense peaks)' and 'Degree (most connected peaks)', both of which are checked. The 'Number of peaks to keep (by cluster and by method)' is set to 2. A checkbox for 'Export final peaks in a new folder?' is also checked. The right column displays a green header 'KEEP TOP PEAKS BY CLUSTER' followed by a list of operations in a terminal-like font, including deleting files and adding adduct modifications for various peaks (120, 214, 266, 322, 268, 328, 434, 441, 444, 456, 465, 466).

Number of kept peaks

Modification of adduct annotation directly in .MAT file for further MS-FINDER annotation

## 2.5. Interrogation of MS-FINDER

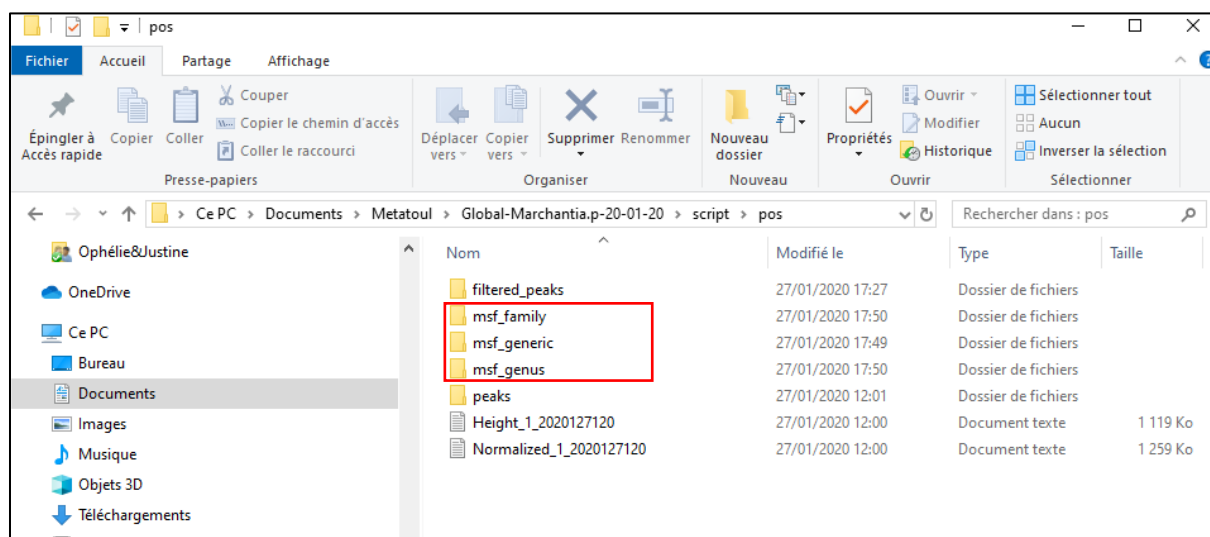


From « **filtered peaks** » folder, interrogate MS-FINDER based on several databases of your choice (for example plant genus, plant family, generic databases from MS-FINDER, ...).

*Optional: Add a “Compound\_level” column within your in-house database for MSfinder. This level will be used for annotation ranking in the next step.*

The most **important thing** to do is to create respectively in “pos” and “neg” directories, new folders named “**msf\_X**” (for example msf\_genus) which correspond to the name of each database used for feature annotation. The msf\_generic is mandatory and correspond to internal database in MS-Finder.

For each database used, export “structure” and “formula” as a single file in the corresponding folder.



## 2.6. Launch MS-FINDER annotation



Once all your MS-FINDER interrogations are done and your folder “msf\_X” filled with “**structure**” and “**formula**” files, go to the fourth index called “**Launch MS-FINDER annotation**”.

This step will permit the annotation of the variables based either only on the score of MS-FINDER or on the prioritization of the different databases, used to indicate the more pertinent annotation.

mscleanr

Project directory Clean MS-DIAL data Keep top peaks by cluster **Launch MS-FINDER annotation** Convert annotated peaks to MSP files

Datasets

☐ Select the best annotation for each peak based only on MSFINDER scores?

Indicate the compound levels in your annotation files, separated by commas (leave blank if none).

1a,1b (A)

Indicate the biosource levels in your annotation process, separated by commas.

genus, family, generic (B)

Indicate the scores multipliers associated to your compound or biosource levels, separated by commas (leave blank if none).

genus:2,family:1.5,generic:1| (C)

LAUNCH MS-FINDER ANNOTATION

This option is used to report the identification with the best MS-FINDER score

This option is used when you want to prioritize some databases.

In (A) you have to indicate the compound level within your database

In (B) you have to order your database

In (C) you can dedicate to your database levels a multiplier to calculate new scores from MS-FINDER ones.

Annotating peaks with MS-FINDER data...

Project directory Clean MS-DIAL data Keep top peaks by cluster **Launch MS-FINDER annotation** Convert annotated peaks to MSP files Datasets

☐ Select the best annotation for each peak based only on MSFINDER scores?

Indicate the compound levels in your annotation files, separated by commas (leave blank if none).

1a,1b

Indicate the biosource levels in your annotation process, separated by commas.

genus,family,generic

Indicate the scores multipliers associated to your compound or biosource levels, separated by commas (leave blank if none).

1a:2,b:1.5,genus:2,family:1.5,generic:1

LAUNCH MS-FINDER ANNOTATION

```

/!\ Level b present in scores but not in biosource or compound levels.
/!\ Deleting C:/Users/justine.chervin/Documents/Metatoul/Global-Marchantia.p-20-01-20/script/fina
*** Treating C:/Users/justine.chervin/Documents/Metatoul/Global-Marchantia.p-20-01-20/script ***
Annotating with 2 compound levels ( 1a, 1b ) and 3 biosource levels ( genus, family, generic ).
*** Annotating clusters with [M+H]+ / [M-H]- couples ***
Annotating cluster 41
Annotating cluster 110
Annotating cluster 124
Annotating cluster 146

```

Summary of compounds and biosource levels used

Paste of annotation in the final peak list

Two files are created in the “final-data” folder:

- **Annotated MS peaks cleaned** = the final peak list with annotation from MS-FINDER

- **Annotated MS peaks normalized** = the final peak list renormalized based on total peak area

The final peak list looks like as follow. Different information are available such as:

- The average m/z value;
- The average RT value;
- The annotation based on MS-FINDER interrogation on the “**Structure**” column with the associated **Total score** of MS-FINDER and **Final score** calculated from the indicated multipliers.
- The source of the annotation in the “**level**” column;
- The ontology of the compound; ...

The variable are also identified as:

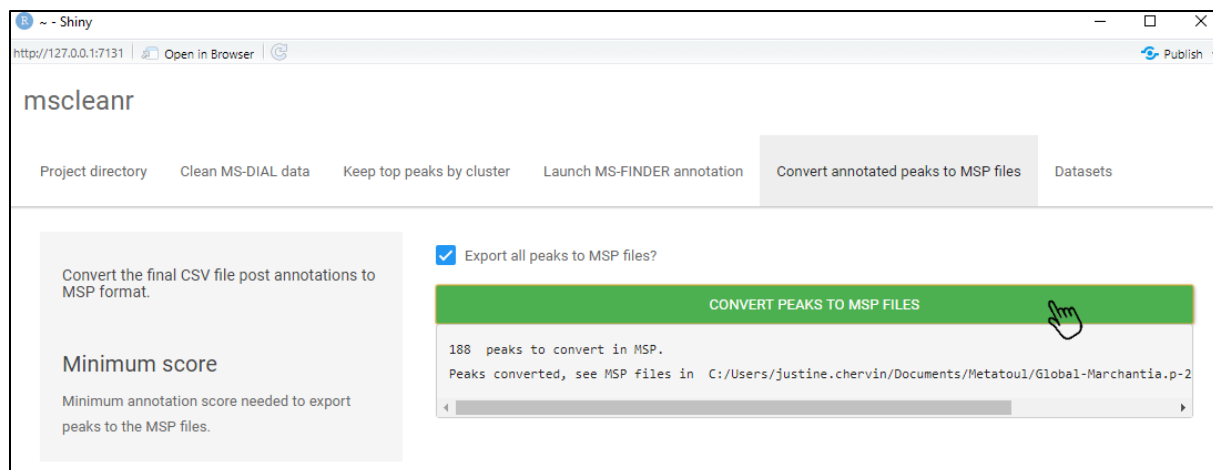
- Unknown compound = variable with no annotation
- Simple ID = based on a single feature in pos or neg mode
- Double ID =based on annotation retrieve in pos and neg mode

J	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
1	Annotation result	annotation_wamir	source	Alignment	Average RT	Average MS	Adapt type	level	Formula	Structure	Total score	Final score	Title	MS1 count	MSMS count	PRECURSORMZ	PRECURSORTYPE	Theoretical mass	Mass error	Formula score	Ontology
2	Unknown compound	neg	108	1.346	316.7895	[M-H] <sup>+</sup>															
3	Unknown compound	neg	173	6.558	427.0308	[M-H] <sup>+</sup>															
4	Unknown compound	neg	56	1.303	242.05179	[M-H] <sup>+</sup>															
5	Unknown compound	pos	10	3.463	111.0761	[M+H] <sup>+</sup>															
6	Unknown compound	pos	105	7.532	179.0466	[M+H] <sup>+</sup>															
7	Unknown compound	pos	424	6.896	401.07574	[M+H] <sup>+</sup>															
8	Unknown compound	pos	159	21.543	206.61177	[M+H] <sup>+</sup>															
9	Unknown compound	pos	202	1.402	236.1484	[M+H] <sup>+</sup>															
25	Simple ID	neg	11	2.818	128.05514	[M-H] <sup>+</sup>		genus	3H7NO3	5(5)-5-carboxy-4,5-dihydro	7.0089	28.0356	(ROA) PYROGLUTAN	3	13	128.0551	[M-H] <sup>+</sup>	129.0425931	0.0002166	4.495	Alpha amino acids and derivatives
26	Simple ID	neg	111	2.014	323.02853	[M-H] <sup>+</sup>		genus	29H13N2O9P	5-(2S,3R,4S,5R)-3,4-dihydro	7.3562	29.4248	(ROA) URIDINE MO	3	37	323.0285	[M-H] <sup>+</sup>	324.0358666	9.02E-05	4.314	Pentose phosphates
27	Simple ID	neg	128	7.55	345.1874	[M-H] <sup>+</sup>		generic	21H12O29	3,4,5-Trimethoxyphenyl g	5.7246	5.7246	Unknown	3	26	345.1877	[M-H] <sup>+</sup>	346.1269223	0.0004058	2.949	Phenolic glycosides
28	Simple ID	neg	13	3.904	129.01901	[M-H] <sup>+</sup>		generic	3H6O4	Itaconic acid	6.4729	6.4729	(ROA) ITACONATE	3	5	129.019	[M-H] <sup>+</sup>	130.0266087	0.0003322	3.471	Branched fatty acids
29	Simple ID	neg	132	1.286	353.08966	[M-H] <sup>+</sup>		genus	21H18O9	6-Hydroxy-2-methyl-7-[(2	5.8621	11.7242	Unknown	3	72	353.0897	[M-H] <sup>+</sup>	354.0950822	-0.001894	3.344	Phenolic glycosides
30	Simple ID	neg	140	21.533	367.20425	[M-H] <sup>+</sup>		generic	22H28N2O3	17-O-Acetylalmitine	3.9909	3.9909	Unknown	3	12	367.2043	[M-H] <sup>+</sup>	368.2099928	-0.001384	2.205	Alkaline-sarcosine alkaloids
31	Simple ID	neg	148	6.887	380.15028	[M-H] <sup>+</sup>		genus	21H12N5O6	2-(3S,4S,5R,6S)-2-Hydroxy	6.5026	26.0104	Unknown	3	62	380.1503	[M-H] <sup>+</sup>	381.1448355	0.001257	3.868	Fatty acid glycosides of mono- and disaccharides
32	Simple ID	neg	155	13.242	392.20822	[M-H] <sup>+</sup>		generic	22H13N1O6	Pulchellamine G-(+)-Pulci	6.5573	6.5573	Unknown	3	60	392.2082	[M-H] <sup>+</sup>	393.2151377	-0.000339	3.752	Guananilides and derivatives
33	Simple ID	neg	164	11.837	408.20285	[M-H] <sup>+</sup>		generic	22H13N1O7	Isariotrin C(-)-Isariotrin C	5.7896	5.7896	Unknown	3	86	408.2029	[M-H] <sup>+</sup>	409.2100523	-0.000124	3.514	Neopanes
34	Simple ID	neg	171	18.496	423.36013	[M-H] <sup>+</sup>		generic	22H12N2O8S	2-(1S,3R,4R,9aS)-1-Hydro	5.3463	5.3463	Unknown	3	77	423.3601	[M-H] <sup>+</sup>	424.1668076	-0.000599	3.345	Sulfenilides
35	Simple ID	neg	174	11.272	435.1293	[M-H] <sup>+</sup>		genus	22H24O10	4-(1S,2,4-Dihydroxy-6-[(2S	6.1769	24.7076	Unknown	3	72	435.1293	[M-H] <sup>+</sup>	436.136947	0.0003705	3.631	Flavonoid O-glycosides
36	Simple ID	neg	175	17.228	437.13907	[M-H] <sup>+</sup>		genus	22H22O5	Psychantol A, 6'-Hydroxy	6.7501	13.5002	Unknown	3	78	437.1391	[M-H] <sup>+</sup>	438.1467238	0.0003474	3.688	Diarylethers
37	Simple ID	neg	176	18.368	437.17554	[M-H] <sup>+</sup>		genus	29H28O4	Marchantin C, 1'-Me ethe	6.0043	24.0172	Unknown	3	81	437.1755	[M-H] <sup>+</sup>	438.1831093	0.0003329	3.277	Lignans, neolignans and related compounds
38	Simple ID	neg	18	1.722	133.01407	[M-H] <sup>+</sup>		genus	24H6O5	L-maleate	7.4137	29.6548	(ROA) MALATE	3	10	133.0141	[M-H] <sup>+</sup>	134.0215333	0.0001468	4.345	Beta hydroxy acids and derivatives
39	Simple ID	neg	163	8.576	407.13458	[M-H] <sup>+</sup>		genus	22H24O9	2S,3R,4S,5S,6R)-2-[2-(2-13	6.663	13.326	Unknown	3	64	407.1346	[M-H] <sup>+</sup>	408.1420323	0.0001559	3.776	Stilbene glycosides
40	Simple ID	neg	189	12.962	449.10275	[M-H] <sup>+</sup>		genus	28H18O6	2,2',3,3',7,7'-Hexahydroxy-	5.4968	21.9872	Unknown	3	30	449.1028	[M-H] <sup>+</sup>	450.1103383	0.0002618	3.418	Phenanthrois
41	Simple ID	neg	19	1.344	135.0298	[M-H] <sup>+</sup>		generic	24H2O5	Erythronic acid	7.0979	7.0979	Unknown	3	27	135.0298	[M-H] <sup>+</sup>	136.0371734	9.69E-05	3.559	Sugar acids and derivatives
42	Simple ID	neg	190	12.714	449.10278	[M-H] <sup>+</sup>		genus	28H18O6	2,2',3,3',7,7'-Hexahydroxy-	5.0015	20.006	Unknown	3	45	449.1028	[M-H] <sup>+</sup>	450.1103383	0.0002618	3.261	Phenanthrois
43	Simple ID	neg	191	17.84	452.27817	[M-H] <sup>+</sup>		generic	22H44NO7P	LysolPE(0/16:0)	6.2888	6.2888	Unknown	3	30	452.2782	[M-H] <sup>+</sup>	453.2855394	6.29E-05	3.643	2-acyl-sn-glycero-3-phosphoethanolamines
44	Double ID	pos	465	10.908	463.08771	[M+H] <sup>+</sup>		genus	22H18O12	2R,3R,4S,5S,6S)-6-[(1S,4	7.4323	29.7292	Unknown	3	54	463.0877	[M+H] <sup>+</sup>	462.079826	-0.000598	4.628	Flavonoid-7-O-glucuronides
45	Double ID	pos	466	10.005	463.08774	[M+H] <sup>+</sup>		genus	22H18O12	2R,3R,4S,5S,6S)-6-[(1S,4	7.5668	30.2672	Unknown	3	39	463.0877	[M+H] <sup>+</sup>	462.079826	-0.000598	4.631	Flavonoid-7-O-glucuronides
46	Double ID	pos	488	17.402	471.18042	[M+H] <sup>+</sup>		generic	22H30N2O8S	Doxapamine	5.4081	5.4081	Unknown	3	237	471.1804	[M+H] <sup>+</sup>	470.1722869	-0.000837	3.534	Methionine and derivatives
47	Double ID	pos	207	10.179	473.20294	[M-H] <sup>+</sup>		generic	22H34O11	UNPD79762	5.6695	5.6695	Unknown	3	61	473.2029	[M-H] <sup>+</sup>	474.2101119	-6.45E-05	3.305	Terpene glycosides
48	Double ID	pos	517	9.531	639.11963	[M+H] <sup>+</sup>		genus	22H28O18	2R,3R,4S,5R,6S)-6-[(1S,4	7.244	28.976	Unknown	3	47	639.1196	[M+H] <sup>+</sup>	638.113314	-0.00041	4.641	Flavonoid-7-O-glucuronides
49	Double ID	neg	181	17.298	439.15469	[M-H] <sup>+</sup>		genus	22H24O5	Marchantin C, 12-Hydroxy	6.9413	25.3652	Unknown	3	32	439.1547	[M-H] <sup>+</sup>	440.1623739	0.0003974	3.493	Lignans, neolignans and related compounds
49	Double ID	pos	436	17.098	425.17465	[M+H] <sup>+</sup>		generic	20H28N2O6S	2-(1S,3R,4R,9aS)-1-Hydro	5.8585	5.8585	Unknown	3	115	425.1747	[M+H] <sup>+</sup>	424.1668076	-0.000616	3.618	Sulfenilides
49	Double ID	pos	450	17.064	441.16974	[M+H] <sup>+</sup>		genus	22H24O5	Marchantin C, 2'-Hydroxy	6.2361	24.9444	Unknown	3	194	441.1697	[M+H] <sup>+</sup>	440.1623739	-0.97E-05	3.5	Lignans, neolignans and related compounds
49	Double ID	pos	45	7.123	293.0234	[M-H] <sup>+</sup>		generic	21H12N2O2	L-Tryptophan	8.0971	8.0971	(ROA) TRYPTOPHAN	3	41	293.0824	[M-H] <sup>+</sup>	294.0998776	0.0002012	4.082	Indolyl carboxylic acids and derivatives

## 2.7. Export peaks as msp files



In the fifth index “**Convert annotated peaks to MSP files**”, you will be able to create two msp files named “peaks-neg” and “peaks-pos” in the folder “final\_data”. All peaks can be converted, or user can choose a scoring threshold based on multiplied MSfinder score.



These two files could then be imported in MetGem software or GNPS facility to create mass spectral similarity networks.