

Data-Dependent-Assisted Data-Independent Acquisition (DaDIA.R) User Manual

(Version 2, 2021-01-04)

Jian Guo¹, Sam Shen¹, Shipei Xing¹, Tao Huan^{1,*}

¹ Department of Chemistry, Faculty of Science, University of British Columbia, Vancouver Campus, 2036 Main Mall, Vancouver, V6T 1Z1, BC, Canada

* Author to whom correspondence should be addressed:

Dr. Tao Huan

Tel: (+1)-604-822-4891

E-mail: thuan@chem.ubc.ca

Internet: <https://huan.chem.ubc.ca/>

- DaDIA.R is an R script for performing DaDIA workflow of metabolic feature extraction and annotation.
- The program is written in the language ‘R’ and is publicly available at <https://github.com/HuanLab/DaDIA.git>
- Please see below for detailed instructions on using the DaDIA.R code:

- 1) **File preparation. (Important: the number of samples in line 38 has to agree with the real number of DIA mzXML files).** User needs to create two folders to store DDA and DIA files separately. All mzXML files from DDA analyses need to be put in the DDA folder. All mzXML files from DIA analyses need to be put in the DIA folder. The library file in the format of .msp should also be put in the DIA folder. In addition, a .txt file containing the information about the DIA m/z range should also be put into the DIA folder. **Figure 1** illustrates the details of how the files should be organized in the corresponding folders. The values in the .txt file are separated by tab. If it is DIA(SWATH) data, the txt file should contain the information about the m/z range for the survey scan and SWATH windows. If it is DIA(AIF) data, the txt file should contain the information about the m/z range for the survey scan and AIF window. The m/z range file examples for both DIA(SWATH) and DIA(AIF) are illustrated in **Figure 2**. Note that the column headers should be kept the same as the examples shown.

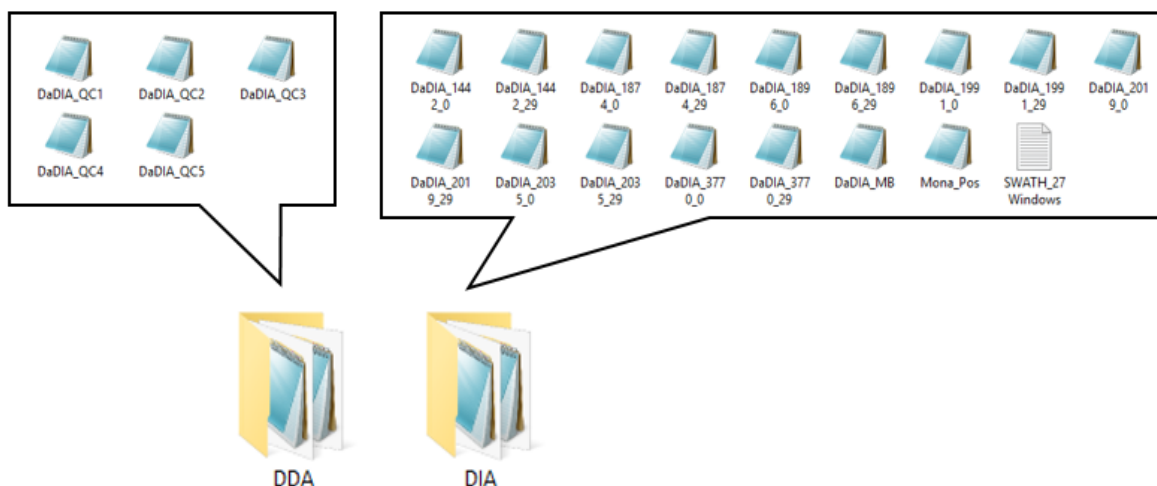


Figure 1. m/z range file sample format for DIA(SWATH) (left) and DIA(AIF) (right) files.

Experiment	MSType	Minmz	Maxmz
0	SCAN	100	1000
1	MSMS	100	126
2	MSMS	125	151
3	MSMS	150	176
4	MSMS	175	201
5	MSMS	200	226
6	MSMS	225	251
7	MSMS	250	276
8	MSMS	275	301
9	MSMS	300	326
10	MSMS	325	351
11	MSMS	350	376
12	MSMS	375	401
13	MSMS	400	426
14	MSMS	425	451
15	MSMS	450	476
16	MSMS	475	501
17	MSMS	500	526
18	MSMS	525	551
19	MSMS	550	576
20	MSMS	575	601
21	MSMS	600	626
22	MSMS	625	651
23	MSMS	650	676
24	MSMS	675	701
25	MSMS	700	801
26	MSMS	800	901
27	MSMS	900	1000

Figure 2. m/z range file sample format for DIA(SWATH) (left) and DIA(AIF) (right) files.

2) Download the R-scrip “DaDIA.R” from Github (<https://github.com/HuanLab/DaDIA.git>).

R package installation. In R-studio, user needs to first install libraries “xcms”, “MSnbase”, “dplyr”, “doParallel”, “foreach”, “metaMS”, and “CAMERA” if they are previously not installed. **R Version 4.0 or above, XCMS Development Version 3.11.4 or above, and metaMS Version 1.25.1 are required; all other packages should be updated to the newest available version.**

3) **Parameter setting.** After all the required libraries are successfully installed. User needs to set the parameters to their desired values. All the parameters available for customized setting are in line 17 – 60, as shown in **Figure 3**. The function of each parameter is described in **Table 1**.

```

17 #####
18 #Part 1: Parameters for feature extraction
19 DDA.directory <- "E:/DDA/"
20 DIA.directory <- "E:/DIA/"
21 cwpDDA <- CentWaveParam(ppm=10,
22                         peakwidth=c(5,60),
23                         mzdiff = 0.01,
24                         snthresh = 6,
25                         integrate = 1,
26                         prefilter = c(3,100),
27                         noise = 100) #XCMS parameters for DDA feature extraction
28 cwpDIA <- CentWaveParam(ppm=10,
29                         peakwidth=c(5,60),
30                         mzdiff = 0.01,
31                         snthresh = 6,
32                         integrate = 1,
33                         prefilter = c(3,100),
34                         noise = 100) #XCMS parameters for DIA feature extraction
35 mass.tol <- 10 #mz tolerance in ppm: used in feature dereplication and MS2 matching
36 mass.const.tol <- 0.05 #mz tolerance in constant value: used in feature rescue
37 rt.tol <- 60 #rt tolerance in seconds
38 num.samples <- 1 #enter how many DIA samples here
39 plot.DaDIA <- TRUE #plot DaDIA features
40 plot.DaDIA.mztol <- 0.5 #DaDIA feature plotting mz window width
41 plot.DaDIA.rttol <- 30 #DaDIA feature plotting rt window width
42 #Parameters for alignment
43 bw <- 5
44 minfrac <- 0.5
45 mzwid <- 0.015
46 max <- 100
47 quantitative.method <- "maxo"
48 # "maxo" = peak height
49 # "into" = peak area
50 #####
51 #Part 2: Parameters for database search (dot product)
52 feature.annotation <- TRUE #annotate DaDIA features
53 db.name <- "Library.msp" #annotation library name
54 ms1.tol <- 0.01 #dot product calculation ms1 tolerance
55 ms2.tol <- 0.02 #dot product calculation ms2 tolerance
56 dot.product.threshold <- 0.1 #dot product annotation threshold
57 match.number.threshold <- 1 #annotation match number threshold
58 adduct_isotope.annotation <- TRUE #perform CAMERA annotation
59 export.mgf <- TRUE #export individual MS2 spectra as .mgf
60 combine.mgf <- TRUE #combine all exported .mgf files
61 MS2mirrorplot <- TRUE #plot mirror plots for features with dot product larger than dot product threshold
62 #####

```

Figure 3. Parameter settings of DaDIA.R.

Table 1. The functions of all DaDIA parameters.

Line #	Parameter Name	Parameter Function
-----------	----------------	--------------------

19	<i>DDA.directory</i>	Set the directory containing all DDA .mzxml files
20	<i>DIA.directory</i>	Set the directory containing all DIA .mzxml files, <i>m/z</i> window .txt file (no specific name is required as the program recognizes it by its file type), and annotation library .msp file
21	<i>cwpDDA</i>	Set XCMS parameters for DDA feature extraction
28	<i>cwpDIA</i>	Set XCMS parameters for DIA feature extraction
35	<i>mass.tol</i>	Set <i>m/z</i> tolerance (\pm ppm) for MS ¹ feature dereplication and MS ² matching
36	<i>mass.const.tol</i>	Set <i>m/z</i> tolerance (\pm constant value) for rescuing DIA features using DDA data
37	<i>rt.tol</i>	Set retention time tolerance (\pm sec) for identifying the same features
38	<i>num.samples</i>	Set number of DIA samples to run
39	<i>plot.DaDIA</i>	Set whether to plot EIC for DaDIA features
40	<i>plot.DaDIA.mztol</i>	Set <i>m/z</i> window width for DaDIA feature EIC plotting
41	<i>plot.DaDIA.rttol</i>	Set RT window width for DaDIA feature EIC plotting
43	<i>bw</i>	Set XCMS feature alignment bandwidth
44	<i>minfrac</i>	Set XCMS feature alignment minimum sample fraction
45	<i>mzwid</i>	Set XCMS feature alignment <i>m/z</i> slice width
46	<i>max</i>	Set XCMS feature alignment maximum # of groups / slice
47	<i>quantitative.method</i>	Set whether to use peak height or peak area for quantitative calculations
52	<i>feature.annotation</i>	Set whether to perform MS ² extraction and DaDIA feature annotation
53	<i>db.name</i>	Set the name of the library used for metabolite annotation
54	<i>ms1.tol</i>	Set MS ¹ tolerance in dot product calculation for metabolite annotation
55	<i>ms2.tol</i>	Set MS ² tolerance in dot product calculation for metabolite annotation
56	<i>dot.product.threshold</i>	Set annotation dot product score threshold
57	<i>match.number.threshold</i>	Set annotation match number threshold
58	<i>adduct_isotope.annotation</i>	Set whether to perform CAMERA adduct and isotope annotation
59	<i>export.mgf</i>	Set whether to export MS ² spectra as individual .mgf files
60	<i>combine.mgf</i>	Set whether to concatenate all exported .mgf files into a single .mgf file
61	<i>MS2mirrorplot</i>	Set whether to plot MS ² mirror plot

- 4) Note: user needs to set the directory in the user's computer that contains all DDA samples in line 18. User needs to set the directory in the user's computer that contains all DIA samples, *m/z* range definitions in .txt format (for DIA(SWATH) or DIA(AIF)), and the annotation library in .msp format in line 19. If you want the warning messages to show, change the options in line 7 from -1 to 0.
- 5) In R-studio, click on "→Source" in the top right corner of the R-studio interface to begin the DaDIA data processing.
- 6) After running the script for single DDA and single DIA sample, one csv file "DaDIATable.csv" containing all metabolic features extracted and one csv file "annotated_output.csv" containing all feature annotation results will be generated in the DIA folder. After running the script for multiple DDA and DIA samples, multiple csv files "n_DaDIATable.csv" (n is the number of DIA samples) containing all metabolic features extracted for each sample, one csv file "alignedDaDIATable.csv" containing aligned features, and one csv file "annotated_output.csv" containing the annotation results for the aligned features will be generated in the DIA folder.
- 7) Notably, in "annotated_output.csv" file, the columns with the header "MS²-Available" contains either TRUE or FALSE values. TRUE means there are MS² spectra assigned to the features, while FALSE means there are no MS² spectra assigned. If the user performs "CAMERA", there will be three additional columns shown up in the file showing the isotopic, adduct, and pcgroup information. The features with the same number in "pcgroup" are actually the same metabolite as they are highly correlated peaks.

Specific Notes

- a) Note: if users wish to use their own in-house library in .csv format for annotation, they must first convert their library from .csv file to an .msp file using the R script "convertMSP.R" at the provided website on GitHub (<https://github.com/HuanLab/DaDIA.git>).
- b) Note: user can choose to plot the EIC of all the metabolic features by switching on the plot function in line 38 and setting the *m/z* and retention time tolerance in line 39, 40. If the user chooses to plot DaDIA features, a folder named "DaDIA_EIC" will be generated in the DIA folder containing all the EIC plots. The name of the EIC plots is composed of feature retention time and *m/z* values.
- c) Note: user can choose to output individual .mgf files for each feature by switching on the function in line 58. Two folders named "DDAmgf" and "DIAmgf" will be generated in the DIA folder containing all individual mgf files from either DDA or DIA data. If the user wishes to output one additional mgf file combining all the MS² information into one file, they can switch on the function in line 59. An mgf file named "combined_mgf" will be generated in the DIA folder. The names of all the mgf files are composed of precursor mass, retention time and the source of the MS² spectrum. (whether it is from DDA or DIA data)

- d) Note: user can choose to generate the MS² mirror plot by switching on the function in line 60. A folder named “MS2mirrorplot” will be generated in the DIA folder containing all the mirror plots in .png file for the features with dot products larger than dot product threshold set by user. In each mirror plot, the corresponding metabolite name is shown on the top of the plot. The name of each plot file is composed of the feature ID and the dot product of the feature.