# Guideline on running Citrullination_Diagnostic_Ion_Analysis

## Requirements

**Python requirements**

- Python version 3.6 or greater
- Dependencies: pandas, numpy, pyteomics, itertools, collections, statistics, re, io

**File requirements**

- Input search result file: The input file must be a .csv file(s) containing the following three columns (see next page for the details):

  'Title'    MS2 spectrum title as written in MGF file
  'Peptide'  Peptide sequence with modification delta mass rounded up to third decimal places
  'Charge'   Charge state of the peptide

- Spectrum file: The spectrum file must be a Mascot Generic Format (MGF) file(s) containing MS2 spectra corresponding to those matched to the PSMs in the input search result file. If MS2 spectra in the input search file and spectrum file are not equivalent, only the common MS2 spectra will be retained and subsequently processed.

# Input search result file

## Example input search result file:

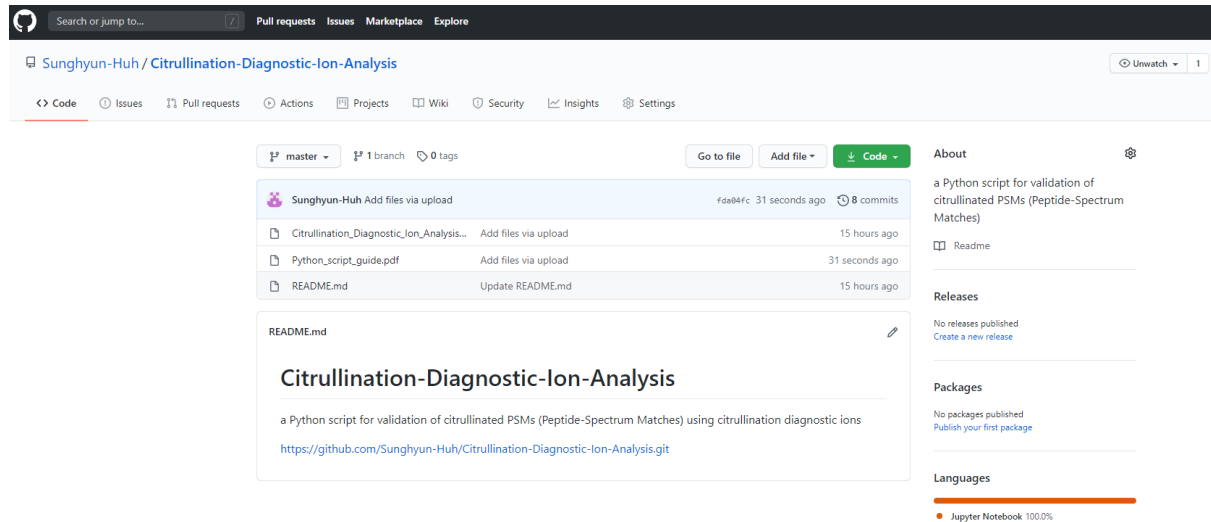| A | B | C |
|---|---|---|
| Title | Peptide | Charge |
| 20160312_02_A1.10012.10012.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1001 | YETSGIGEAR+0.984VK | 2 |
| 20160312_02_A1.10045.10045.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1004 | NIVTPR+0.984TPPPSQGK | 2 |
| 20160312_02_A1.10116.10116.3 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1011 | NIVTPR+0.984TPPPSQGK | 3 |
| 20160312_02_A1.10222.10222.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1022 | NIVTPR+0.984TPPPSQGK | 2 |
| 20160312_02_A1.10334.10334.3 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1033 | NIVTPR+0.984TPPPSQGK | 3 |
| 20160312_02_A1.10362.10362.3 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1036 | TPSTAHLR+0.984VPK | 3 |
| 20160312_02_A1.10418.10418.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1041 | NIVTPR+0.984TPPPSQGK | 2 |
| 20160312_02_A1.10479.10479.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1047 | AQSR+0.984EQLAALK | 2 |
| 20160312_02_A1.1054.1054.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1054" | DSR+0.984SGSPM+15.995AR | 2 |
| 20160312_02_A1.10602.10602.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1060 | NIVTPR+0.984TPPPSQGK | 2 |
| 20160312_02_A1.10646.10646.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1064 | Q+0.984KR+0.984LQ+0.984AM+15.995Q+0.984K | 2 |
| 20160312_02_A1.10671.10671.3 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1067 | SGSEAGSPRR+0.984PRRQR | 3 |
| 20160312_02_A1.1073.1073.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1073" | R+0.984GGGGGRR+0.984SK | 2 |
| 20160312_02_A1.10764.10764.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1076 | R+0.984FIN+0.984DMVK | 2 |
| 20160312_02_A1.10769.10769.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1076 | NIVTPR+0.984TPPPSQ+0.984GK | 2 |
| 20160312_02_A1.10874.10874.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1087 | MAR+0.984EAEFEAEQER | 2 |
| 20160312_02_A1.11026.11026.3 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1102 | RGR+0.984PPKDEK | 3 |
| 20160312_02_A1.11286.11286.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1128 | N+0.984R+0.984Q+0.984VIC+57.021VTLK | 2 |
| 20160312_02_A1.11398.11398.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1139 | GGTSR+0.984ALAAASSVK | 2 |
| 20160312_02_A1.11489.11489.3 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1148 | EEFER+0.984Q+0.984N+0.984KQLR | 3 |
| 20160312_02_A1.11557.11557.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1155 | TVEMR+0.984DGEVIK | 2 |
| 20160312_02_A1.11735.11735.2 File:"20160312_02_A1.raw", NativeID:"controllerType=0 controllerNumber=1 scan=1173 | Q+0.984Q+0.984IADLR+0.984EDLKR | 2 |

search_result_file

Format for peptide sequences should follow that of MS-GF+ search result. Specifically, modification delta masses should be rounded up to third decimal places. Currently allowed modifications are as follows:

| Modification | Mod on peptide |
|---|---|
| Carbamidomethyl Cys | C+57.021 |
| Oxidation Met | M+15.995 |
| Deamidated Asn | N+0.984 |
| Deamidated Gln | Q+0.984 |
| Citrullinated Arg | R+0.984 |
| Pyro-Glu from Glu | E-17.027 |
| Pyro-Glu from Gln | Q-18.011 |
| iTRAQ 4plex Lys | K+144.102 |
| iTRAQ 8plex Lys | K+304.205 |
| TMT Lys | K+229.163 |
| iTRAQ 4plex N-term | +144.102 |
| iTRAQ 8plex N-term | +304.205 |
| TMT N-term | +229.163 |
| Acetyl N-term | +42.011 |

# Downloading the Python script

The Python script can be downloaded via following GitHub page:

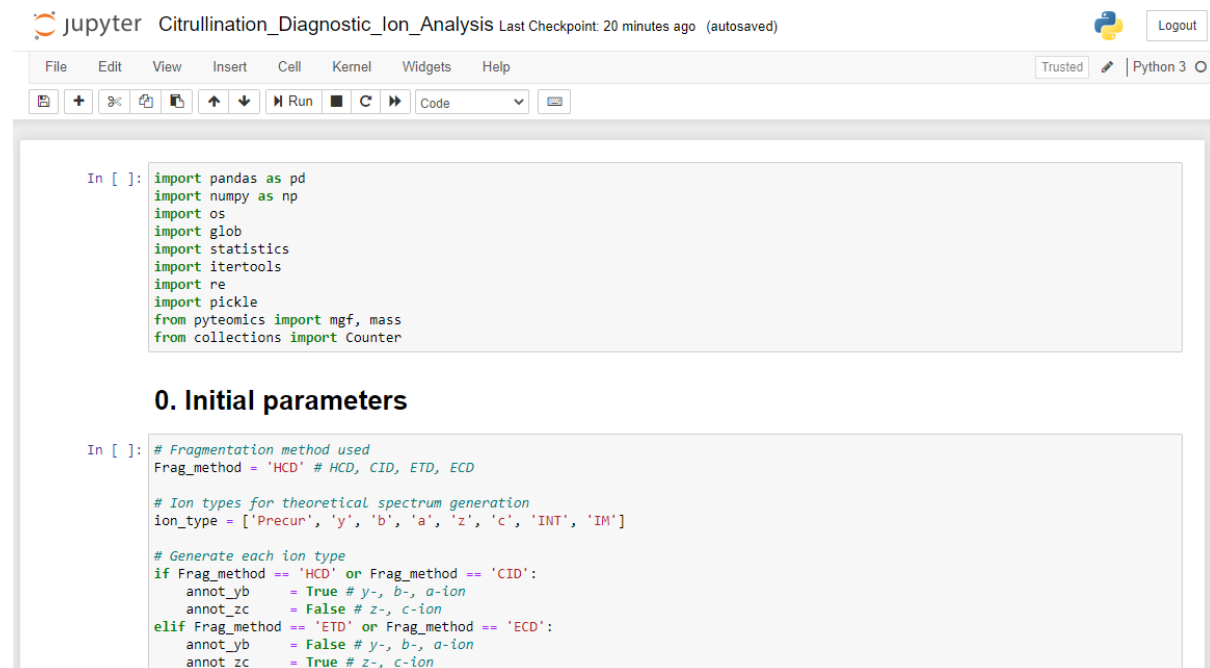https://github.com/Sunghyun-Huh/Citrullination-Diagnostic-Ion-Analysis



Users can download the Jupyter notebook via the following steps:

- Click on the 'Citrullination_Diagnostic_Ion_Analysis.ipynb'.

- Click on the 'Raw'.

- Press ctrl+s and manually type 'ipynb' after the filename to download as a .ipynb file.

# Running the Python script

A snapshot of the Jupyter notebook:



```python
import pandas as pd
import numpy as np
import os
import glob
import statistics
import itertools
import re
import pickle
from pyteomics import mgf, mass
from collections import Counter
```

## 0. Initial parameters

```python
# Fragmentation method used
Frag_method = 'HCD' # HCD, CID, ETD, ECD

# Ion types for theoretical spectrum generation
ion_type = ['Precur', 'y', 'b', 'a', 'z', 'c', 'INT', 'IM']

# Generate each ion type
if Frag_method == 'HCD' or Frag_method == 'CID':
    annot_yb     = True # y-, b-, a-ion
    annot_zc     = False # z-, c-ion
elif Frag_method == 'ETD' or Frag_method == 'ECD':
    annot_yb     = False # y-, b-, a-ion
    annot_zc     = True # z-, c-ion
```

# Running the Python script

A snapshot of initial parameters settings:

## 0. Initial parameters

```
In [ ]: # Fragmentation method used
        Frag_method = 'HCD' # HCD, CID, ETD, ECD

        # Ion types for theoretical spectrum generation
        ion_type = ['Precur', 'y', 'b', 'a', 'z', 'c', 'INT', 'IM']

        # Generate each ion type
        if Frag_method == 'HCD' or Frag_method == 'CID':
            annot_yb    = True # y-, b-, a-ion
            annot_zc    = False # z-, c-ion
        elif Frag_method == 'ETD' or Frag_method == 'ECD':
            annot_yb    = False # y-, b-, a-ion
            annot_zc    = True # z-, c-ion
        annot_precur = True # precursor ion
        annot_INT    = True # internal ion
        annot_IM     = True # immonium ion
        annot_dict = {
            'Precur' : annot_precur,
            'y'      : annot_yb,
            'b'      : annot_yb,
            'a'      : annot_yb,
            'z'      : annot_zc,
            'c'      : annot_zc,
            'INT'    : annot_INT,
            'IM'     : annot_IM
        }

        # MS2 mass tolerance (ppm)
        ms2_ppm = 15

        # Signal-to-noise (SNR) filter for MS2 spectrum
        apply_SNR = True # Apply SNR filter
        SNR       = 2 # SNR threshold
        low       = 0.05 # Define low x% intensity as baseline noise level

        # Maximum charge state of sequence ions
        max_charge = 2 # 2, 3, ... 'max'

        # Maximum number of neutral loss from a single ion
        max_NL = 3
```

Explanations of initial parameters are as follows:

| | |
|---|---|
| Frag_method | Fragmentation method used in the input data (value = 'HCD', 'CID', 'ETD', or 'ECD'). If set as 'HCD' or 'CID', $y$-ion, $b$-ion, and $a$-ion will be generated for theoretical spectrum. If set as 'ETD' or 'ECD', $z$-ion, $c$-ion will be generated for theoretical spectrum. Commonly, precursor, internal, and immonium ion will be generated for all fragmentation method used. |
| ms2_ppm | MS2 level mass tolerance in ppm (default = 15 ppm) |
| apply_SNR | Determine whether to apply signal-to-noise filter to remove noise peaks (value = True or False; default = True) |
| SNR | Signal threshold level. The average intensity of noise peaks (as defined in 'low') multiplied by this signal threshold level will be the final signal-to-noise filter. If 'apply_SNR' = True, all peaks below the signal-to-noise filter will be removed (default = 2) |
| low | Proportion of MS2 peaks regarded as noise. If 'apply_SNR' = True, all peaks below this noise level will be treated as noise (default = 0.05) |
| max_charge | Maximum charge state of fragment ions (default = 2) |
| max_NL | Maximum number of neutral losses from a single ion (default = 3) |

# Running the Python script

A snapshot of codes for loading input files:

## 1. Input files

```
In [ ]: # Set current working directory
        PATH = "F:/Project/"
        os.chdir(PATH)
```

```
In [ ]: # Input files
        spec_files   = glob.glob('spectrum_file.mgf') # MGF file(s)
        search_files = glob.glob('search_result_file.csv') # Search result file(s)
```

Users can upload local input files via the following steps:

- Set the directory in which the input files are located.

- Copy and paste the input filenames. In case of multiple MGF or search result files in the same directory, type in '*.mgf' or '*.csv'.

# Output result file

**Example output file:**



Explanations of output columns are as follows:

| Column | Explanation |
|---|---|
| mod_peptide | Simplified peptide with a predefined set of symbols for modifications |
| Pep_length | Peptide length |
| mz_Precursor | Theoretical precursor $m/z$ |
| Cit_Count | Number of citrullinated sites |
| Total_NL_label | Annotations of all diagnostic neutral loss ions |
| precNL_label | Annotations of precursor neutral losses |
| seqNL_label | Annotations of sequence ion neutral losses |
| intNL_label | Annotations of internal ion neutral losses |
| Total_INT_label | Annotations of all diagnostic internal ions |
| Dipeptide_label | Annotations of diagnostic dipeptides |
| Tripeptide_label | Annotations of diagnostic tripeptides |
| IM_NH3_label | Annotation of IM(Cit)-NH$_3$ |
| Total_NL_count | Number of all diagnostic neutral loss ions |
| precNL_count | Number of precursor neutral losses |
| seqNL_count | Number of sequence ion neutral losses |
| intNL_count | Number of internal ion neutral losses |
| Total_INT_count | Number of all diagnostic internal ions |
| Dipeptide_count | Number of diagnostic dipeptides |
| Tripeptide_count | Number of diagnostic tripeptides |
| IM_NH3_count | Number of IM(Cit)-NH$_3$ |
| Cit_probability | Probability (P) of citrullination status calculated by the EN model (HCD data only) |
| Cit_prediction | Classification of citrullination status using a P cutoff >0.5 (HCD data only) |