

Utilities for Mass Spectrometry Analysis of Proteins

User's Manual

Version 2.2.0

May 2022

To download Utilities for Mass Spectrometry Analysis of Proteins visit:

www.umsap.nl

Utilities for Mass Spectrometry Analysis of Proteins

Copyright © 2017 Kenny Bravo Rodriguez.

All Rights Reserved.

Contents

List of Figures	IV
List of Tables	V
1 Introduction	1
1.1 Citing Utilities for Mass Spectrometry Analysis of Proteins	1
1.2 Acknowledgments	2
2 Obtaining and Installing Utilities for Mass Spectrometry Analysis of Proteins	3
2.1 Obtaining Utilities for Mass Spectrometry Analysis of Proteins	3
2.2 Installing Utilities for Mass Spectrometry Analysis of Proteins	3
2.3 Uninstalling Utilities for Mass Spectrometry Analysis of Proteins	4
3 Workflow in Utilities for Mass Spectrometry Analysis of Proteins	5
3.1 The input files	6
3.2 The output files	6
3.3 Using Utilities for Mass Spectrometry Analysis of Proteins	7
3.4 Navigating through Utilities for Mass Spectrometry Analysis of Proteins	8
3.5 Backward compatibility	8
4 UMSAP Control	10
4.1 The interface	10
4.2 The Tools menu	10
5 Correlation Analysis	13
5.1 The interface	13
5.2 The analysis	15

5.3	The result window	15
5.4	The Tools menu	15
6	Data Preparation	17
6.1	The interface	17
6.2	The analysis	18
6.3	The result window	19
6.4	The Tools menu	20
7	Limited Proteolysis	21
7.1	Definitions	21
7.2	The input files	22
7.3	The interface	22
7.4	The analysis	25
7.5	The result window	26
7.6	The Tools menu	28
8	Proteome Profiling	29
8.1	Definitions	29
8.2	The input files	29
8.3	The interface	29
8.3.1	The Tools menu	33
8.4	The analysis	33
8.5	The output files	34
8.6	Visualizing the output files	35
8.6.1	The Tools menu	36
8.6.1.1	Filters	36
	Bibliography	38

List of Figures

3.1	The main window of UMSAP	5
3.2	Structure of the output generated by UMSAP	7
4.1	The UMSAP Control window	12
5.1	The Correlation Analysis tab	13
5.2	The Correlation Analysis result window	16
6.1	The Data Preparation tab	17
6.2	The Data Preparation result window	19
7.1	The Limited Proteolysis module tab	23
7.2	The Result - Control experiments helper window for the Limited Proteolysis module	25
7.3	The Limited Proteolysis result window	27
8.1	The Proteome Profiling module tab	30
8.2	The Result - Control experiments helper window for the Proteome Profiling module	32
8.3	The structure of the Output folder from the Proteome Profiling module	34
8.4	The Proteome Profiling analysis window	35

List of Tables

3.1	List of built-in keyboard shortcuts	9
-----	---	---

Chapter 1

Introduction

Utilities for Mass Spectrometry Analysis of Proteins (UMSAP) is a graphical user interface (GUI) designed to speed up the post-processing of data obtained during mass spectrometry studies involving proteins. The program is not intended to analyze a mass spectrum or a mass chromatogram, neither to identify the peaks in a mass spectrum. The main objective is the fast post-processing of the vast amount of data generated in mass spectrometry experiments involving proteins after peak identification have been performed.

The program is organized in modules with each module performing a single type of data post-processing. The reason for this clear separation is the high dependency between the type of mass spectrometry experiment performed and the way in which the resulting data must be post-processed. The modules are designed in such a way that the required user input is minimized but still users can control every aspect of the analysis. Currently, the software contains three modules, but several others are already planned.

1.1 Citing Utilities for Mass Spectrometry Analysis of Proteins

If results obtained with UMSAP are published in any way, please acknowledge the use of UMSAP by including the following sentence:

”Utilities for Mass Spectrometry Analysis of Proteins was created by Kenny Bravo Rodriguez at the University of Duisburg-Essen and is currently developed at the Max Planck Institute of Molecular Physiology.”

Any published work, which uses UMSAP, should include the following reference:

Kenny Bravo-Rodriguez, Birte Hagemeier, Lea Drescher, Marian Lorenz, Michael Meltzer, Farnusch Kaschani, Markus Kaiser and Michael Ehrmann. (2018). Utilities for Mass Spectrometry Analysis of Proteins (UMSAP): Fast post-processing of mass spectrometry data. [Rapid Communications in Mass Spectrometry](#), 32(19), 1659–1667.

Electronic documents should include a direct link to the official web page of UMSAP at: www.umsap.nl

1.2 Acknowledgments

I would like to thank all the persons that have contributed to the development of UMSAP, either by contributing ideas and suggestions or by testing the code. Special thanks go to: Dr. Farnusch Kaschani, Dr. Juliana Rey, Dr. Petra Janning and Prof. Dr. Daniel Hoffmann.

In particular, I would like to thank Prof. Dr. Michael Ehrmann.

Chapter 2

Obtaining and Installing Utilities for Mass Spectrometry Analysis of Proteins

2.1 Obtaining Utilities for Mass Spectrometry Analysis of Proteins

UMSAP is distributed free of charge for anyone interested in using it. To obtain a copy of the software just register at www.umsap.nl and go to the Download page.

No extra software or packages are needed for UMSAP to properly work. So far, UMSAP have been tested in macOS 10.14.6 and 12.3 and Windows 10. Support for some Linux distributions will be available in the future.

2.2 Installing Utilities for Mass Spectrometry Analysis of Proteins

Windows

Unzip the file you just downloaded from www.umsap.nl. Then, copy the folder UMSAP to the location in your file system where you want to keep it. Finally, create a shortcut to the executable file UMSAP.exe found inside the main folder UMSAP. That is all. You are now ready to use UMSAP.

macOS

Unzip the file you just downloaded from www.umsap.nl. Then, just move the UMSAP.app folder to /Applications/. That is all. You are now ready to use UMSAP.

Depending on the security settings in macOS, it may be needed to explicitly allow UMSAP to be opened the first time the app is used.

2.3 Uninstalling Utilities for Mass Spectrometry Analysis of Proteins

UMSAP will not create any installation file in your computer. Therefore, the only thing you need to do, to completely uninstall UMSAP, is to delete the folder UMSAP.app in macOS or UMSAP in Windows. You should also delete any shortcut pointing to the executable file of UMSAP and the configuration file `.umsap_config.json` in your home folder. That is all.

Chapter 3

Workflow in Utilities for Mass Spectrometry Analysis of Proteins

When you start UMSAP, the program will display the main window (Figure 3.1). From this window you can access all the modules and utilities either by the menu entries: Modules and Utilities or by the corresponding buttons on the right side list. A complete description of each module and utility is given in the following chapters.

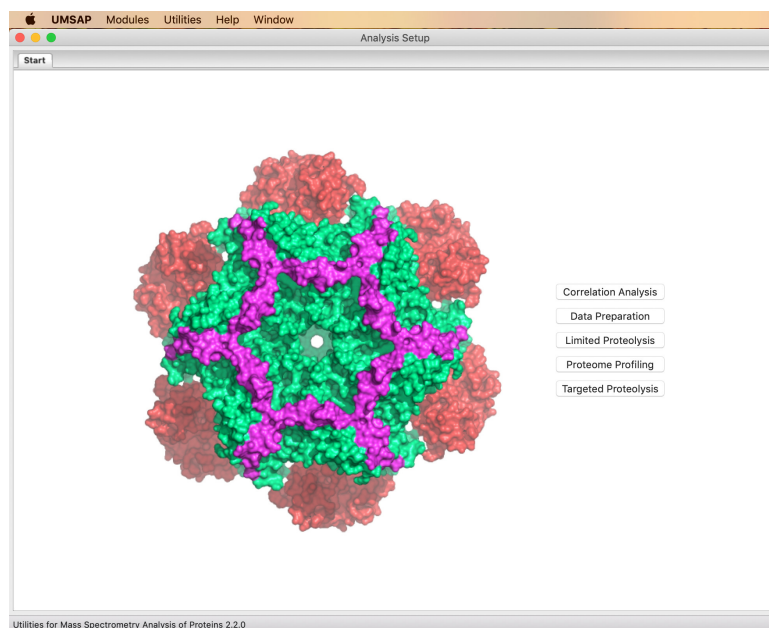


Figure 3.1: The main window of UMSAP. From this window users can access all the available Modules and Utilities.

3.1 The input files

UMSAP has two main input files. One file contains the detected peptide sequences after all peak assignments have been completed, and the other file contains the detected proteins. The program expects these files to be plain text files containing a table with the data. Columns in the files are expected to be tab separated. The first row in the files is expected to contain only the names of the columns. There is no limit in the amount and type of data present in the Data files. However, each module will expect certain columns to be present. Columns not needed by the modules will simply be ignored.

In addition, certain modules use other input files as well. The modules Targeted Proteolysis and Limited Proteolysis use fasta files containing the sequences of the recombinant and native proteins used in the experiments. The first sequence found in the fasta file is assumed to be the sequence of the recombinant protein. The second sequence found in the fasta file is assumed to be the sequence of the native protein. All other sequences found in the fasta file are discarded. If the sequence of the native protein is given UMSAP performs a sequence alignment between the native and recombinant sequences. The alignment allows UMSAP to translate the results obtained with the residue numbers of the recombinant protein to the residue numbers of the native protein. This is done to facilitate future comparison of results between different recombinant proteins of the same native protein. However, when analyzing the results of the alignment UMSAP assumes that the recombinant and native sequences differs only in the N and C-terminal tags while the sequence between the tags is identical. If this is not the case, e.g. there are point mutations or insertion/deletion in the sequence of the recombinant protein no native sequence file should be given to UMSAP.

The Targeted Proteolysis module may also use a local PDB file.

3.2 The output files

Results generated by UMSAP will be saved in two folders and a file with extension .umsap (Figure 3.2). Direct manipulation of the umsap file and files within these folders should be avoided. UMSAP provides a way to manage them through the UMSAP Ctrl window (Chapter 4). Nevertheless, all the files created by UMSAP are plain text files with json or cvs (tab separated) format, in order for users to be able to read their content. Changing the content of the files is highly discouraged as this will lead to errors in the reliability and visualization of the results with UMSAP.

The folder Input_Data_Files contains a copy of the input files used for the analysis in the project. When adding a new analysis to the project, the new input files used will be copied to the Input_Data_Files folder. The date and time of the analysis will be added to the name of the file to avoid overwriting existing files inside the folder.

The folder Steps_Data_Files contains a folder for each analysis in the project. These folders contain the main results for the analysis as well as a step by step account of the calculations and any further analysis performed after the main results were created.

The .umsap file contains information about all the analysis in the project and allows

managing the project and the visualization of the results. An unlimited number of analysis can be added to any given .umsap file. UMSAP will never overwrite or replace an .umsap file, instead new analysis will be added to the selected .umsap file.

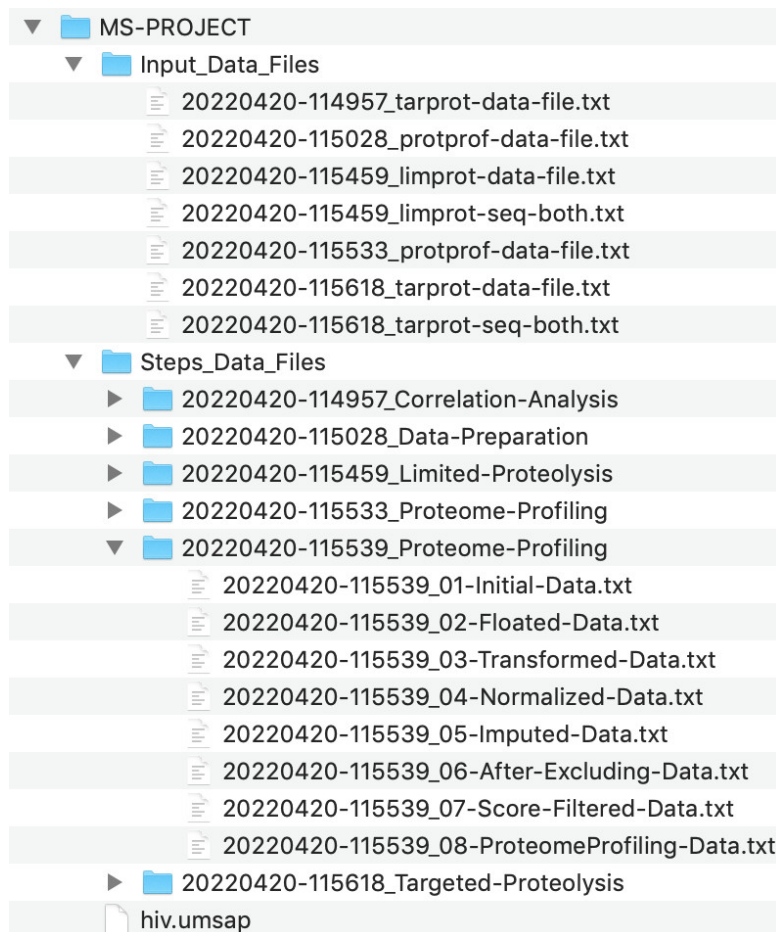


Figure 3.2: Structure of the output generated by UMSAP. Results are saved in the Steps_Data_Files folder. The .umsap file allows managing and visualizing the results.

3.3 Using Utilities for Mass Spectrometry Analysis of Proteins

Once the input files are ready to be analyzed, using UMSAP is straightforward. Just open the program and select a module or utility. In the new tab, fill in the needed information and hit the Start Analysis button at the bottom of the tab. Depending on the amount of data and the complexity of the analysis to perform it may take a few minutes for the program to complete the task at hand. While the analysis is running, a window, containing a progress bar, will appear. This window will give a rough guess of the remaining time needed to complete the current analysis and will report any error

encountered. It will be helpful if users send a crash report to umsap@umsap.nl, so we can correct them.

In order to make the program as user-friendly as possible help messages will pop up from buttons and labels. The help messages will contain a brief description of what is the button or label for and what input is expected from the user. In this way, users can find basic information about a particular element of the interface without needing to go to the manual or online tutorials. If more information is needed, users may consult the manual or click the Help button at the bottom of the module/utility tab to read an online tutorial.

Depending on the module or utility just run, new windows will be created to show a graphical representation of the results. All plots support to zoom into a rectangular selection of the plot and to reset the zoom level.

3.4 Navigating through Utilities for Mass Spectrometry Analysis of Proteins

The entries Modules and Utilities will be available in the menu of every window. The Modules entry in the menu gives direct access to all modules. The same is true for the Utilities entry. These menu entries are the fastest way to access all the functions in UMSAP. In a typical UMSAP session, users will work with different independent windows simultaneously. The windows have descriptive names, so users can quickly guess the content of any window. The scheme of the windows name is *File Name - Utilities or Module Name - ID of the Analysis*. For example, the window with name *hiv.umpap - Target Proteolysis - 20220420-115618 - Cleavage Sites* will be displaying the Targeted Proteolysis analysis with ID 20220420-115618 - Cleavage Sites from file *hiv.umsap*.

A list of current shortcuts is given in Table 3.1.

3.5 Backward compatibility

Unfortunately, UMSAP 2.2.0 is not capable to read any file generated with previous versions of UMSAP.

Shortcut	Action	Window
Alt+Cmd+L	Create the Limited Proteolysis tab	All
Alt+Cmd+P	Create the Proteome Profiling tab	All
Alt+Cmd+T	Create the Targeted Proteolysis tab	All
Cmd+R	Read umsap file	All
Cmd+C	Copy	Text and List boxes
Cmd+X	Cut	Text and List boxes
Cmd+V	Paste	Text and List boxes
Cmd+A	Select all	Text and List boxes
Cmd+P	Show Data Preparation results	Results plot
Cmd+D	Duplicate result window	Results plot
Cmd+E	Export data	Results plot
Cmd+I	Export image	Results plot
Cmd+K	Clear all selections	Results plot
Cmd+A	Add analysis	UMSAP Ctrl
Cmd+X	Delete analysis	UMSAP Ctrl
Cmd+E	Export analysis	UMSAP Ctrl
Cmd+U	Reload file	UMSAP Ctrl
Cmd+Z	Reset the zoom on a plot	Selected plot
Alt+Shift+I	Export all images	Multiple plots
Alt+Shift+Z	Reset all zooms	Multiple plots
Shift+I	Export main plot image	Multiple plots
Shift+Z	Reset main plot zoom	Multiple plots
Alt+I	Export secondary plot image	Multiple plots
Alt+Z	Reset secondary plot zoom	Multiple plots
Cmd+A	Show all peptides	Limited Proteolysis
Cmd+L	Toggle Band/Lane selection mode	Limited Proteolysis
Cmd+S	Export sequence alignments	Limited Proteolysis
Shift+A	Add label to Volcano plot	Proteome Profiling
Shift+P	Toggle Pick label / Select protein	Proteome Profiling
Shift+Cmd+A	Apply all Filters	Proteome Profiling
Shift+Cmd+F	Auto apply all Filters	Proteome Profiling
Shift+Cmd+R	Remove selected Filters	Proteome Profiling
Shift+Cmd+Z	Remove last applied Filter	Proteome Profiling
Shift+Cmd+X	Remove all Filters	Proteome Profiling
Shift+Cmd+C	Copy Filters	Proteome Profiling
Shift+Cmd+V	Paste Filters	Proteome Profiling
Shift+Cmd+S	Save Filters	Proteome Profiling
Shift+Cmd+L	Load Filters	Proteome Profiling
Shift+Cmd+E	Export filtered data	Proteome Profiling
Cmd+S	Export sequence alignments	Targeted Proteolysis

Table 3.1: List of built-in keyboard shortcuts. Windows users should replace Cmd with Ctrl.

Chapter 4

UMSAP Control

The UMSAP Control windows shows the content of an .umsap file (Figure 4.1).

4.1 The interface

The analysis contained in the selected .umsap file are displayed in alphabetical order and grouped by the analysis type. The checkboxes to the left of the names of the Utilities and Modules allow creating the corresponding window showing the results available for the selected Utility or Module.

Each analysis in the file is represented by the user-provided Analysis ID. Unfolding any ID will display all the configuration values provided by the user prior to running the analysis. In addition, a left click over any Analysis ID will create the corresponding tab in the Analysis Setup window (Figure 3.1) and populate all fields with the values in the selected analysis. This is the fastest way to configure the analysis tab to rerun an analysis with slight changes in the configuration options. After rerunning an analysis or simply adding a new analysis to the .umsap file, the window will be automatically updated to display the new results.

4.2 The Tools menu

The UMSAP Control windows allows also to manage the content of the selected .umsap file. Currently, it is possible to Add (Cmd+A) analysis from a different .umsap file, to Delete (Cmd+X) analysis from an existing .umsap file and to Export (Cmd+E) the analysis in an .umsap file to a new .umsap file.

Adding analysis from an .umsap file to the already opened .umsap file will result in the addition of the new information to the already opened .umsap file and in the copy of the necessary files and folders to folders Input_Data_Files and Steps_Data_Files. During this process there is a small chance to end up with duplicated file and/or folder names or Analysis ID. In this case, UMSAP will rename the file/folder/Analysis ID to avoid any overwriting and will update any reference in the .umsap file to the files/folders that

were renamed.

Deleting any analysis from an .umsap file will also result in the removal of the files and/or folders referenced in the deleted analysis. Files in `Input_Data_Files` are only deleted if they are not referenced by any remaining analysis. Deleting all analysis in an .umsap file will result in the removal of the .umsap file and folders `Input_Data_Files` and `Steps_Data_Files`. If the folder containing the project is empty after deleting all UMSAP files and folders the project folder is also deleted.

Exporting some or all analysis in an opened .umsap file to an already existing .umsap file is not possible. When exporting the selected analysis to a project folder containing an `Input_Data_Files` and/or `Steps_Data_Files` folder, UMSAP will create a new folder in the selected project folder and export all the information to this empty folder.

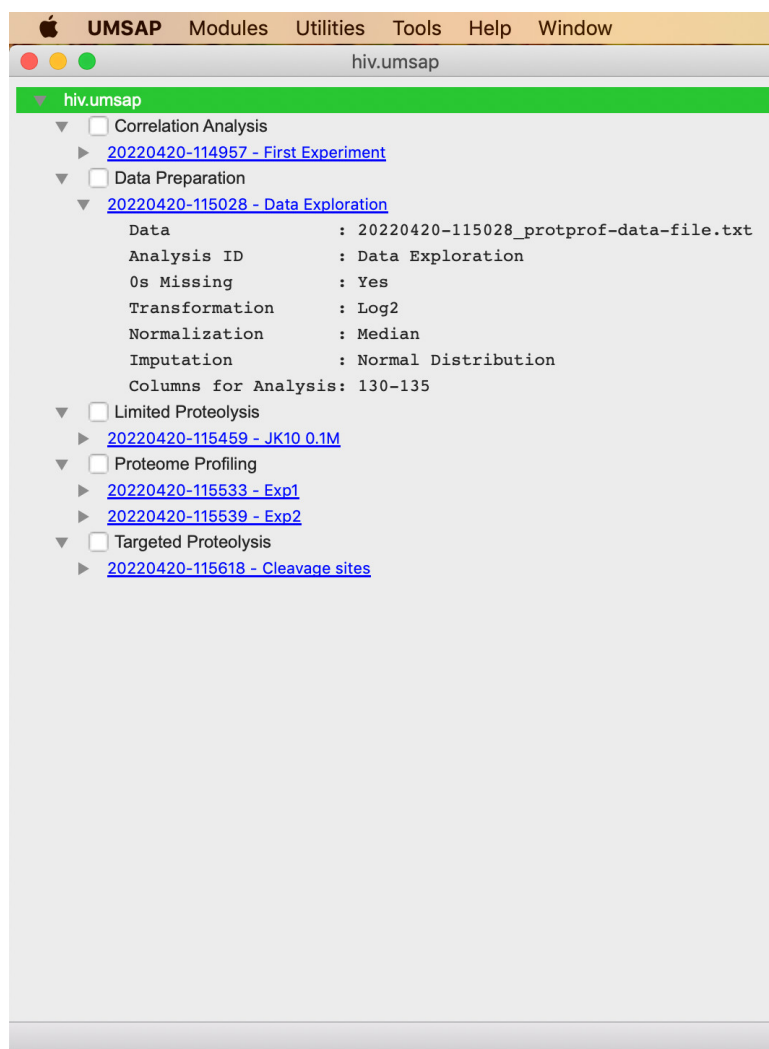


Figure 4.1: The UMSAP Control window. The content of the selected .umsap file is shown in alphabetical order. The window allows managing the content of the .umsap file and to visualize the results of the analysis in the file.

Chapter 5

Correlation Analysis

The Correlation Analysis utility calculates the correlation in the MS data used as input for UMSAP.

5.1 The interface

The Correlation Analysis tab is divided in four regions (Figure 5.1).

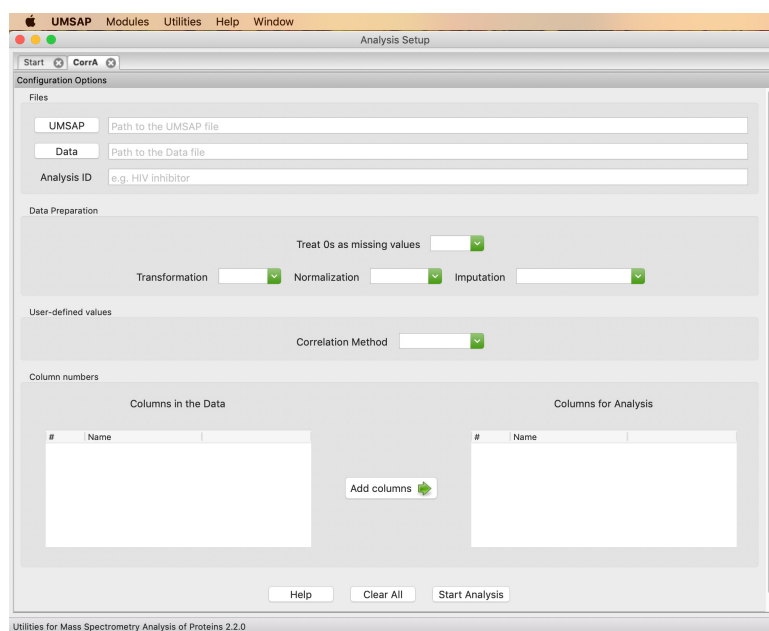


Figure 5.1: The Correlation Analysis tab. This tab allows to perform a correlation analysis of the data contained in a given Data file.

Region Files contains two buttons and a text field. Here users select the input and output files for the analysis.

1. The button UMSAP allows users to browse the file system to select the location and

name of the .umsap file. When selecting an already existing .umsap file the operating system will ask if it is ok to replace the file, the answer can be yes since UMSAP will never overwrite or replace an .umsap file, instead the new analysis will be added to the already existing file. Only .umsap files can be selected here.

2. The button Data allows users to browse the file system to select the input data file that will be used for the analysis. The Data file is expected to be a plain text file with tab separated columns and the name of the columns in the first row of the file. In addition, columns to be analyzed must contain only numbers and must be of the same length. Only .txt files can be selected here.

3. The text field Analysis ID allows users to provide an ID for the analysis to be run. The date and time of the analysis will be automatically added to the beginning of the name.

Region Data Preparation contains four dropdown boxes. Here users select how the data in the Data file should be prepared before starting the analysis.

1. The dropdown Treat 0s as missing values allows user to define how to handle 0 values present in the Data file.

2. The dropdown Transformation allows user to select the Transformation method to be applied to the data.

3. The dropdown Normalization allows users to select the Normalization method to be applied to the data.

4. The dropdown Imputation allows user to select the Imputation method used to replace missing values in the data.

Region User-defined values contains one dropdown box.

1. The dropdown Correlation Method allows users to select the correlation method to use.

Region Column numbers contains two lists and a button. Here users select the columns in the Data file to be used in the Correlation Analysis.

1. The list to the left will display the names of the columns present in the selected Data file. The list is automatically filled once the Data file is selected. Rows in the list can not be deleted, except in the case of loading a different Data file or using the Clear All button at the bottom of the tab. Selected rows can be copied with the Cmd+C shortcut.

2. The list to the right will contain the columns in the Data file that will be used for the Correlation Analysis. This list must contain at least two rows for the analysis to proceed. Selected rows in this list can be deleted with the Cmd+X shortcut and rows already copied from the list in the left can be pasted with the Cmd+V shortcut. While pasting the rows, duplicate rows will be silently discarded. Importantly, the order of the rows and columns in the matrix containing the correlation coefficients will be the same as the order of the columns in this list. Therefore, users are advised to fill the list in such a way that replicates of the same experiment are consecutive to each other in the list.

3. The button Add columns will add the selected rows in the left list to the right list. The rows will be added to the right list in the same order as they are selected in the left list. Duplicate rows will be silently discarded.

The bottom of the tab contains three buttons.

1. The button Help leads to an online tutorial about Correlation Analysis in UMSAP.
2. The button Clear All will delete all user input from the tab.
3. The button Start Analysis starts the Correlation Analysis.

5.2 The analysis

First, UMSAP will check the validity of the user provided input. Then, columns in the right list are read from the Data file. The columns must contain only numbers and the same amount of rows must be found in all columns. Failing to comply with this will result in the program aborting the analysis. After this, all steps selected in the Data Preparation region are carried out (Chapter 6). Finally, the correlation coefficients are calculated using the selected method. If any of the coefficients cannot be calculated, then the corresponding coefficient is set to NA.

5.3 The result window

The correlation coefficients resulting from a Correlation Analysis will be shown as a color coded matrix (Figure 5.2). Values between -1 to 0 will be shown in shades of red, 0 will be shown as white and values between 0 to 1 will be shown in shades of blue. NA values will be shown in green. The columns and rows of the matrix are the column names used to calculate the correlation coefficients. Information about a specific matrix element can be obtained by simply placing the mouse pointer over the matrix element.

5.4 The Tools menu

The Tools menu in the window showing the correlation coefficients allows user to view any of the Correlation Analysis contained in the selected .umsap file or to modify the appearance of the displayed plot. For example, the column numbers can be displayed instead of the column names or the color bar can be hidden. In addition, only a subset of the columns can be shown using the Select Columns entry.

The Tools menu also allows duplicating the window (Cmd+D) for easier comparison of two or more analysis, checking the Data Preparation steps of the analysis (Cmd+P), creating an image of the plot (Cmd+I), exporting the correlation coefficient matrix to a tab separated CSV file (Cmd+E) and resetting the zoom level of the plot (Cmd+Z).

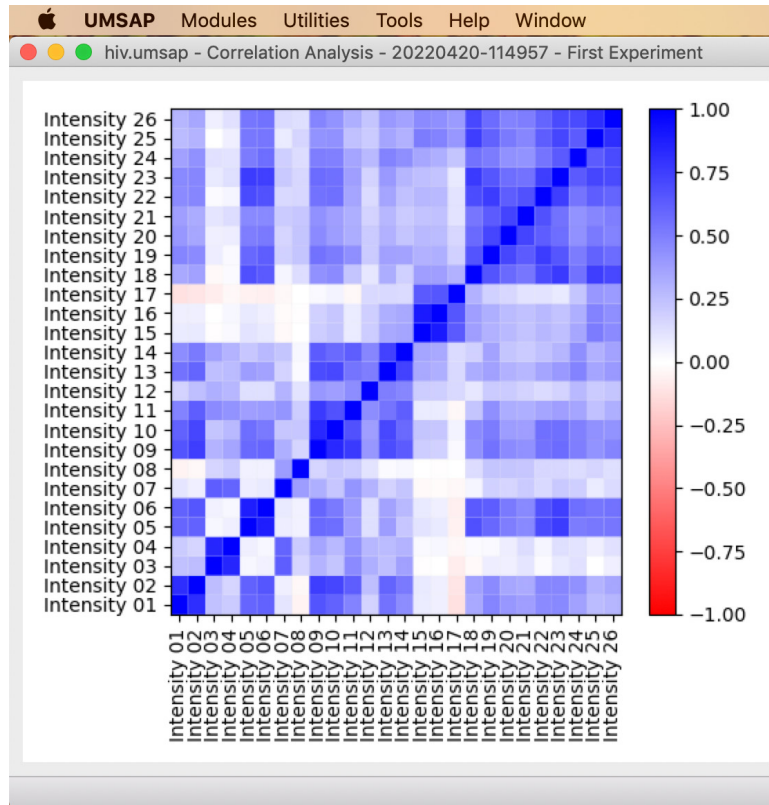


Figure 5.2: The Correlation Analysis result window. The correlation coefficients are shown as a color coded matrix. Values between -1 to 0 are shown in shades of red, 0 is shown in white and values between 0 to 1 in shades of blue. NA values are shown in green.

Chapter 6

Data Preparation

The Data Preparation utility allows exploring the distribution of the data in the selected Data File and the impact that different Data Preparation options have over the data.

6.1 The interface

The Data Preparation tab is divided in two main regions (Figure 6.1).

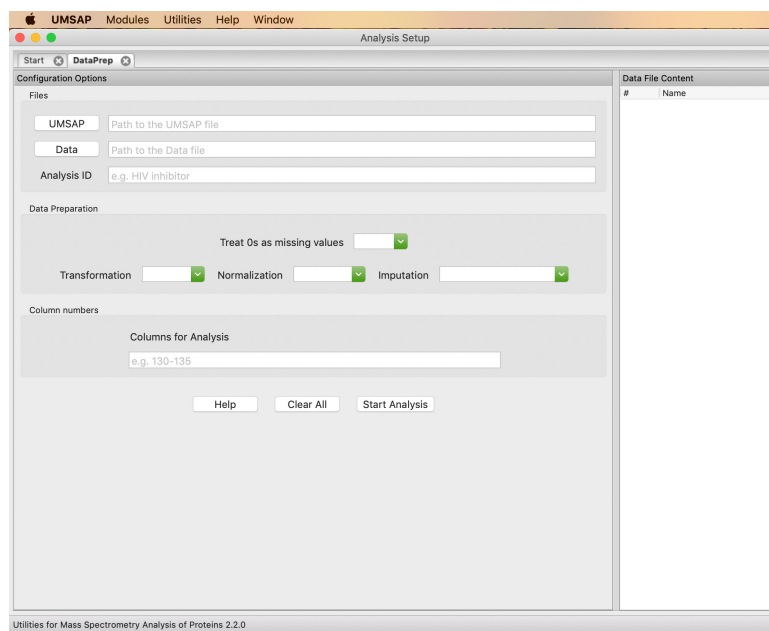


Figure 6.1: The Data Preparation tab. This tab allows to perform a statistical exploration of the data contained in a given Data file.

The Data File Content region holds only a list to show the name of the columns in the selected Data File. The list will be automatically filled after selecting the file.

The Configuration Options region contains all the fields needed to configure and run

the analysis.

Section Files contains two buttons and a text field. Here users select the input and output files for the analysis.

1. The button UMSAP allows users to browse the file system to select the location and name of the .umsap file. When selecting an already existing .umsap file the operating system will ask if it is ok to replace the file, the answer can be yes since UMSAP will never overwrite or replace an .umsap file, instead the new analysis will be added to the already existing file. Only .umsap files can be selected here.
2. The button Data allows users to browse the file system to select the input data file that will be used for the analysis. The Data file is expected to be a plain text file with tab separated columns and the name of the columns in the first row of the file. In addition, columns to be analyzed must contain only numbers and must be of the same length. Only .txt files can be selected here.
3. The text field Analysis ID allows users to provide an ID for the analysis to be run. The date and time of the analysis will be automatically added to the beginning of the name.

Section Data Preparation contains four dropdown boxes. Here users select how the data in the Data file should be prepared before starting an analysis.

1. The dropdown Treat 0s as missing values allows user to define how to handle 0 values present in the Data file.
2. The dropdown Transformation allows user to select the Transformation method to be applied to the data.
3. The dropdown Normalization allows users to select the Normalization method to be applied to the data.
4. The dropdown Imputation allows user to select the Imputation method used to replace missing values in the data.

Section Column numbers contains a text field. Here users specify the Columns in the Data File to be used during the Data Preparation steps. Only integers can be accepted here. Column numbers can be copied (Cmd+C) and paste (Cmd+V) from the selected rows in the list on region Data File Content or just type the numbers.

The bottom of the region contains three buttons.

1. The button Help leads to an online tutorial about Correlation Analysis in UMSAP.
2. The button Clear All will delete all user input from the tab.
3. The button Start Analysis starts the Correlation Analysis.

6.2 The analysis

First, UMSAP will check the validity of the user provided input. After this, the selected Data File is read and the following steps are taken:

1. The content of all specified columns in Data File is checked to make sure only numbers are found in them. 0 values present in the columns are left or remove depending on the selected value for field Treat 0s as missing values.
2. The indicated Transformation method is applied to the selected columns.
3. The indicated Normalization method is applied to the transformed data.
4. The indicated Imputation method is applied to the normalized data.

The results from the four steps is saved, so users can check the effect of the selected workflow over the data. Currently, only one method is implemented for the Transformation, Normalization and Imputation of the Data, respectively. The only alternative is to skip the corresponding step. The methods available will be expanded in the near future. All steps are column wise applied.

6.3 The result window

The window showing the results from a Data Preparation workflow is divided in three regions (Figure 6.2).

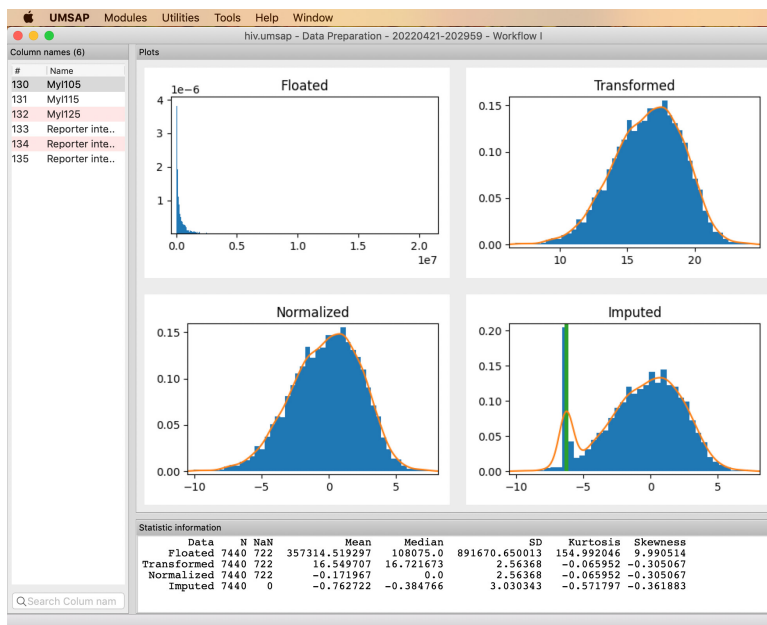


Figure 6.2: The Data Preparation result window. Histograms for the initial, transformed, normalized and imputed data are shown.

Region Column names shows a table with the number (0 based) and name of the analyzed columns.

Region Plots shows the results from the Data Preparation workflow in four histograms for the selected column in region Column names. The histograms are created for the initial, transformed, normalized and imputed data. They show the probability density as blue bars and the calculated probability density function in orange. The green bars

in the Imputed histogram represent the imputed values.

Region Statistic information shows a description of the data for the selected column in region Column names.

6.4 The Tools menu

The Tools menu in the window showing the results from a Data Preparation workflow allows user to view any of the Data Preparation analyzes contained in the selected .umsap file. The Tools menu also allows duplicating the window (Cmd+D) for easier comparison of two or more set of results, creating an image of the plots (Alt+Shift+I), exporting the data shown to a tab separated CSV file (Cmd+E) and resetting the zoom level of the plots (Alt+Shift+Z).

Chapter 7

Limited Proteolysis

The Limited Proteolysis module is designed to post-process the results from an enzymatic digestion performed in two steps. The first step is assumed to be a limited proteolysis in which a large protein is split in smaller fragments. The fragments are then separated using a SDS-PAGE electrophoresis. Finally, selected bands from the gel are submitted to a full enzymatic digestion and the generated peptides are analyzed using mass spectrometry.

The main objective of the module is to identify the protein fragments generated in the initial limited proteolysis from the peptides found in the MS analyzed gel spots. This is achieved by performing an equivalence test([1](#)) between the peptides in the selected gel spots and a control spot containing the full length target protein. In this way, peptides leaked from one gel spot to another can be eliminated. Several replicates of the experiment are expected.

7.1 Definitions

Before explaining in detail the interface of the module and how does the module work, let's make clear the meaning of some terms that will be used in the following paragraphs.

- *Recombinant protein*: actual amino acid sequence used in the mass spectrometry experiments. It may be identical to the native sequence of the Target protein under study or not.
- *Native protein*: full amino acid sequence expressed in wild type cells.
- *Detected peptide*: any peptide detected in any of the mass spectrometry experiments including the control experiments.
- *Relevant peptide*: a detected peptide with a Score value above a user defined threshold, see page 24.
- *Filtered peptide*: a relevant peptide with equivalent intensities in the control and a given gel spot at the chosen significance level.

- *Fragment*: group of filtered peptides with no gaps when their sequences are aligned to the sequence of the recombinant/native protein.

For example, there are three fragments in the alignment shown below. The first fragment is formed by sequences 1 to 3 since there is no gap in the sequence MKKTAIAIAVAL. SEQ4 forms the second fragment because there is a gap between the last residue in SEQ3 and the first residue in SEQ4 and another gap between the last residue in SEQ4 and the first residue in SEQ5. For the same reason SEQ5 forms the third fragment.

REC.PROT	MKKTAIAIAVALAGFATVAQAASWSHPQFEKIEGRRDRGQKTQSAPGTL	50
SEQ1	MKKTAIAIAV.....	10
SEQ2	..KTAIAIAV.....	8
SEQ3IAIAVAL.....	7
SEQ4ATVAQAASWS.....	10
SEQ5DRGQKTQSAPG...	11

7.2 The input files

The Limited Proteolysis module requires a Data file containing the detected peptides and a sequence file containing the amino acid sequence of the recombinant protein used in the study. Both files must follow the guidelines specified in Section 3.1. In short, the Data file must have a tabular format with tab separated columns and the name of the columns are expected as first row. The Sequence file is expected to contain at least one sequence and to be FASTA formatted. If more than one sequence is found in the Sequence file the first sequence will be taken as the sequence of the Recombinant protein and the second sequence will be taken as the sequence of the Native protein. All other sequences are discarded.

7.3 The interface

The tab of the Limited Proteolysis module is divided in two regions (Figure 7.1).

The Data File Content region holds only a table to show the name of the columns in the selected Data File. The table will be automatically filled after selecting the file.

The Configuration Options region contains all the fields needed to configure and run the analysis.

Section Files contains three buttons and a text field. Here users select the input and output files for the analysis.

1. The button UMSAP allows users to browse the file system to select the location and name of the .umsap file. When selecting an already existing .umsap file the operating system will ask if it is ok to replace the file, the answer can be yes since UMSAP will never overwrite or replace an .umsap file, instead the new analysis will be added to the already existing file. Only .umsap files can be selected here.

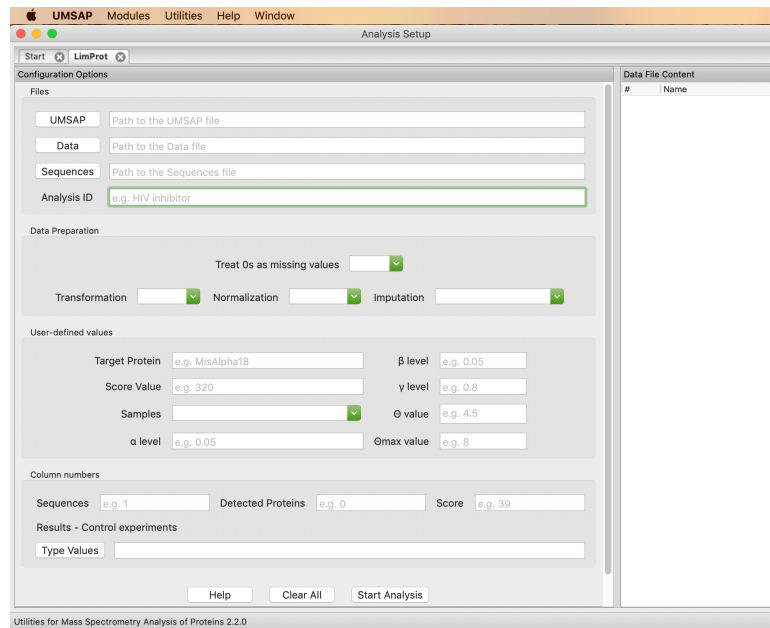


Figure 7.1: The Limited Proteolysis module tab. This tab allows users to perform the analysis of the results obtained during a two steps enzymatic proteolysis experiment where the products from the first limited digestions are separated using SDS-PAGE electrophoresis.

2. The button Data allows users to browse the file system to select the input data file that will be used for the analysis. The Data file is expected to be a plain text file with tab separated columns and the name of the columns in the first row of the file. In addition, columns to be analyzed must contain only numbers and must be of the same length. Only .txt files can be selected here.

3. The button Sequences allows users to browse the file system to select the FASTA file containing the sequence of the Recombinant protein and the Native protein. The FASTA file must contain at least one sequence.

4. The text field Analysis ID allows users to provide an ID for the analysis to be run. The date and time of the analysis will be automatically added to the beginning of the name.

Section Data Preparation contains four dropdown boxes. Here users select how the data in the Data file should be prepared before starting the analysis.

1. The dropdown Treat 0s as missing values allows user to define how to handle 0 values present in the Data file.

2. The dropdown Transformation allows user to select the Transformation method to be applied to the data.

3. The dropdown Normalization allows users to select the Normalization method to be applied to the data.

4. The dropdown Imputation allows user to select the Imputation method used to

replace missing values in the data.

Section User-defined values contains seven text fields and one dropdown box. Here users configure the Limited Proteolysis to be run.

1. The text field Target Protein allows users to specify the protein of interest. Users may type here any unique protein identifier present in the Data file. The search for the Target Protein is case-sensitive, meaning that eFeB is not the same as efeb.
2. The text field Score Value allows users to define a threshold value above which the detected peptides will be considered as relevant. The Score Value is an indicator of how reliable was the detection of the peptide during the MS experiments. The value given to UMSAP depends on the program generating the Data file. Only one real number equal or greater than zero will be accepted here. A value of zero means all detected peptides belonging to the Target Protein will be treated as relevant peptides.
3. The dropdown Samples allows users to specify whether samples are independent or paired. For example, samples are paired when the same Petri dish is used for the control and experiments.
- 4–8. The parameters α , β , γ , Θ and Θ_{\max} are used to configure the equivalence test⁽¹⁾ performed to identify peptides in the selected gel spots with equivalent intensity values to the control spots (Section 7.4). α , β and γ must be between 0 and 1. The value of Θ is optional. If left blank UMSAP will calculate a value for each peptide based on the intensity values found in the Data file. If given then the given value will be used for each peptide. Θ_{\max} is the maximum value to consider the intensity values in the gel spot and control as equivalent.

Section Column numbers contains four text fields. Here, users provide the column numbers in the Data file from where UMSAP will get the information needed to perform the Limited Proteolysis analysis. All columns specified in this section must be present in the Data file. Column numbers start at 0. The column numbers are shown in the table of Region Data File Content after the Data file is selected.

1. The text field Sequences allows users to specify the column in the Data file containing the sequences of the peptides identified in the MS experiments. Only one integer number equal or greater than zero will be accepted here.
2. The text field Detected Proteins allows users to specify the column in the Data file containing the unique protein identifier for the proteins detected in the MS experiments. It is in this column where the program will look for the Target Protein value given in Section User-defined values. It is important that in this column the Target Protein value is used to refer to only one protein. Only one integer number equal or greater than zero will be accepted here.
3. The text field Score allows users to specify the column in the Data file containing the Score values. It is in this column where the program will look for the values to be compared against the Score threshold given in Section User-defined values.
4. The text field Results - Control experiments allows users to specify the columns in the Data file containing the results of the control and experiments. The button Type Values call a helper window (Figure 7.2) where users can type the information needed.

Figure 7.2: The Result - Control experiments helper window for the Limited Proteolysis module. This window allows users to specify the column numbers in the Data file containing the MS results for the selected gel spots.

The helper window is divided in two Regions. Region Data File Content will show the column numbers and names of the columns present in the selected Data file. Region Configuration Options has two sections. The upper section allows defining the number of bands and lanes of interest in the gel as well as the label for lanes, bands and control spot. The button Setup Fields creates the corresponding text fields in the bottom section to type the column numbers. Each text field should contain the column numbers with the MS results for the given gel spot. The values for the text fields should be positive integer numbers or a range of integers, e.g. 60–62 or left blank for empty gel spots. Selected rows in the table can be copied (Cmd+C) and then pasted (Cmd+V) in the text fields. Duplicate column numbers are not allowed.

7.4 The analysis

First, UMSAP will check the validity of the user-provided input and then the selected Data file is read. The columns specified in section Column numbers are extracted from the Data file. All other columns present in the Data file are discarded. After this, all steps selected in the Data Preparation section are applied to the columns specified in the text field Result - Control experiments (Chapter 6). Then, the following actions are performed.

All rows in the prepared data containing peptides that do not belong to the Target protein are removed. Then, all rows containing peptides from the Target protein but with Score values lower than the user-defined Score threshold are removed. These steps leave only relevant peptides, this means, peptides with a Score value higher than the user-defined threshold that belong to the Target protein. For each one of these relevant peptides the equivalence test is performed (1).

The implementation of the equivalence test is based on the following equations:

$$s^* = s \sqrt{\frac{n-1}{\chi^2_{(\gamma, n-1)}}} \quad (7.1)$$

$$\Theta = \delta + s^* [t_{(1-\alpha, 2n-2)} + t_{(1-\beta/2, 2n-2)}] \sqrt{\frac{2}{n}} \quad (7.2)$$

$$(\bar{y}_1 - \bar{y}_2) \pm t_{(1-\alpha, n_1+n_2-2)} \cdot \sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \quad (7.3)$$

where s^* is an estimate of the upper confidence limit of the standard deviation, $\chi^2_{(\gamma, n-1)}$ is the (100γ) th percentile of the chi-squared distribution with $n-1$ degrees of freedom, Θ is the acceptance criterion, δ is the absolute value of the true difference between the group's mean values, t is the Student's t value, \bar{y} is the measurement mean and s_p is the pooled standard deviation of the measurements calculated with:

$$s_p = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}} \quad (7.4)$$

α , β , γ and Θ are the parameters defined in section User-defined values of region Configuration Options in the tab of the module.

In essence, for each relevant peptide, the control experiments are used to estimate the upper confidence limit for the standard deviation using Equation 7.1 and then the acceptance criterion is calculated with Equation 7.2. Finally, the confidence interval for the mean difference for the gel spot and the control is calculated with Equation 7.3 and compare to Θ . Peptides with equivalent mean intensity in at least one gel spot and the control are retained while not equivalent peptides are discarded.

If the value of Θ is given in section User-defined values of the module's tab then only the confidence interval for the mean difference is calculated, and the value is directly compared to the given Θ value. The maximum possible Θ value must always be provided. The reason for this is that when only a few replicates of the experiments are performed the calculated Θ value may be too large and then the equivalence test is easily past by all relevant peptides.

After the filtered peptide (FP) are identified the modules creates the output files.

7.5 The result window

The window showing the results from a Limited Proteolysis workflow is divided in four regions (Figure 7.3).

Region Peptide List contains a table with all FP detected during the analysis. Selecting a peptide in the table will highlight with a thick black border all gel spot in Region Gel Representation and all fragments in region Protein Fragments where the selected

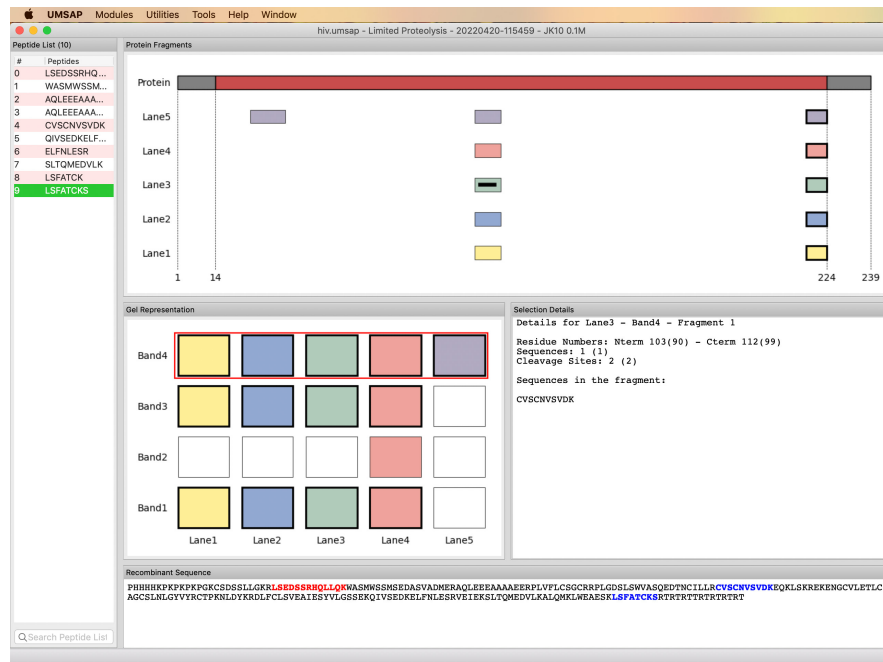


Figure 7.3: The Limited Proteolysis result window. Users can perform here the analysis of the fragments obtained in the limited proteolysis experiments.

peptide was found. In addition, region Recombinant Sequence will show the selected peptide in blue. The search box at the bottom allows searching for a sequence in the list of FP.

Region Gel Representation contains a representation of the analyzed gel. Here, each gel spot is represented with a square. White squares represent gel spot for which no peptide from the Target protein was detected with intensity values equivalent to the controls or that were not analyzed because no column number information was given when configuring the analysis. The rest of the square will be colored according to the band/lane they belong to.

There are two selection modes available for Region Gel Representation. In the Lane selection mode a left click over an empty space in the gel representation will select the closer lane. In this mode the gel spot will be colored according to the band they belong to. The selected lane will be highlighted with a red rectangle. The band selection mode works similarly, but users can select bands and the gel spot are colored according to the lane they belong to. The selection mode can be toggled through the keyboard shortcut Cmd+L or the Tools menu. In addition, the entire gel can be selected (Cmd+A) or a single gel spot can be selected with a left click. Any selection on the gel will update the content of Regions Protein Fragments, Selection Details and Recombinant Sequence.

Region Protein Fragments will display a graphical representation of the fragments found in each gel spot for the selected band/lane. The first fragment in this region represents the full length of the recombinant sequence of the Target protein. Here the central red section represents the sequence in the recombinant protein that is identical to the native protein sequence while gray sections represent the sequences in the recombinant protein

that are different to the native protein sequence. If the sequence of the native protein was not given then the fragment is shown in gray. The fragments are color coded using the same colors of the band/lane they belong to. Selecting a fragment will update the information shown in regions Selection Details and Recombinant Sequence.

Region Selection Details will show information about the selected lane/band, gel spot in region Gel Representation or the selected fragment in region Protein Fragments. The displayed information for a selected band/lane includes the number of non-empty lanes/bands, the number of fragments identified in each non-empty gel spot in the band/lane and the protein regions identified. Selecting a gel spot will display this information only for the gel spot. Selecting a fragment in region Protein Fragments will display the following information: number of cleavage sites, first and last residue number for the selected fragment and a sequence alignment of all peptides forming the fragment.

7.6 The Tools menu

The Tools menu in the window showing the results from a Limited Proteolysis analysis allows user to view any of the analyzes contained in the selected .umsap file. Users can toggle the band/lane selection mode (Cmd+L), select all gel spot in the analysis (Cmd+A) and clear all selections made in the window (Cmd+K). In addition, the zoom level in the plots can be reset and an image of the plots can be created.

The Tools menu also allows duplicating the window (Cmd+D) for easier comparison of two or more analysis, checking the Data Preparation steps of the analysis (Cmd+P), exporting the results of the analysis to a tab separated CSV file (Cmd+E) and to export the sequence alignments (Cmd+S) between the peptides found in the analysis and the sequence of the recombinant protein.

Chapter 8

Proteome Profiling

The Proteome Profiling module is designed to identify differentially expressed protein under various experimental conditions. A typical example is to compare the effect of two substance over protein expression using whole cell lysates.

8.1 Definitions

Before explaining in detail the interface of the module and how does the module work, let's make clear the meaning of some terms that will be used in the following paragraphs.

- *Detected protein*: any protein detected in any of the mass spectrometry experiments including the control experiments.
- *Relevant proteins*: a detected protein with a Score value above a user defined thresholds (page 31).

8.2 The input files

The Proteome Profiling module requires only one input file. This Data file must follow the guidelines specified in Section 3.1. In short, the Data file must have a tabular format with tab separated columns and the name of the columns are expected as first row. All columns given as input in section Column numbers of Region Configuration Options of the interface (Figure 8.1) must be present in the Data file.

8.3 The interface

The tab of the Proteome Profiling module is divided in two regions (Figure 8.1).

Region 1 contains four buttons allowing users to quickly delete all provided input and start a new analysis. The Clear all button will delete all user provided input and will

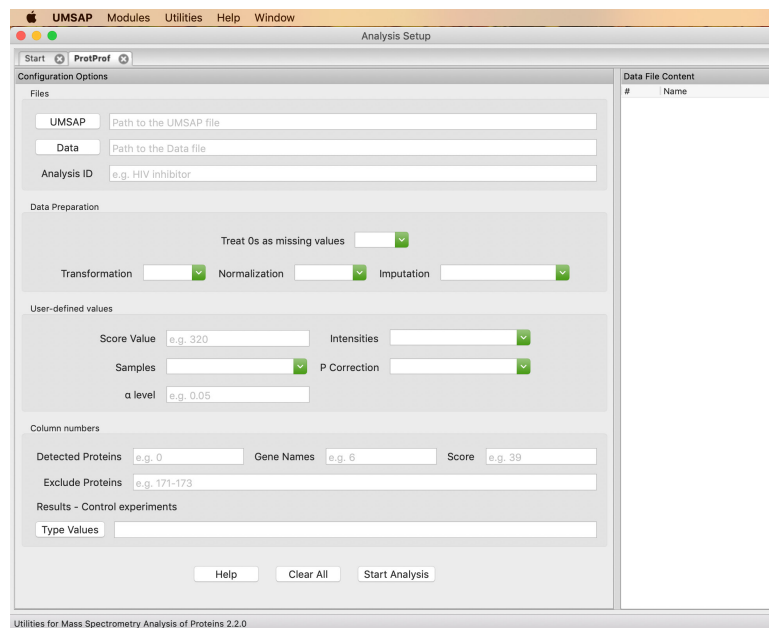


Figure 8.1: The Proteome Profiling module tab. This tab allows users to perform a proteome profiling analysis.

empty the list box in Region 3. The Clear files button will delete all values in section Files of Region 2 and will empty the list box in Region 3. The Clear values button will delete all values in section Values of Region 2. Finally, the Clear columns button will delete all values in section Column numbers of Region 2.

Region 2 contains the fields where users provide the information needed in order to perform the post-processing of the Data file. The section Files of Region 2 will provide the path to the Data and Output files. It contains three buttons.

1.- The Data file button allows users to browse the file system and select a Data file. Only .txt files can be selected here. Once the Data file is selected, the name of the columns in the file will be shown in the list box of Region 3. If the path to the Data file is typed in, the display of the name of the columns in Region 3 can be triggered by pressing the Enter key in the keyboard while the Data file entry box has the focus of the keyboard.

2.- The Output folder button allows users to browse the file system and select the location of the folder that will contain the output. By default, UMSAP will create a ProtProf folder inside the selected Output folder to save all the generated results. If only the name of the output folder is given, the output folder will be created in the same folder containing the Data file. If this field is left empty, then the ProtProf folder will be created in the same folder containing the Data file. If the selected Output folder already contains a ProtProf folder, then the current date and time to the seconds will be added to the name in order to avoid overwriting the files from previous analyses.

3.- The Output name button does nothing but the text box to its right allows users to specify the name of the files that will be generated during the analysis. If this field is

left empty, then the name `protprof` will be used for the output files.

The section Values of Region 2 of the interface contains four parameters. Here, users provide information about how the Data file should be processed.

1.- The parameter Score value allows users to define a threshold value above which the detected proteins will be considered as relevant. The Score value is an indicator of how reliable was the detection of the protein during the MS experiments. The value given to UMSAP depends on the program generating the Data file. Only one real number equal or greater than zero will be accepted as a valid input here. A value of zero means all detected proteins will be treated as relevant proteins.

2.- The parameter Data normalization allows selecting the normalization procedure to be performed before running the analysis of the data in the Data file. Currently, only a \log_2 normalization is possible but this will be expanded to include quantile, variance stabilization and local regression normalization, among other methods.

3.- The parameter Median correction indicates whether to apply a median correction to protein intensities in each experiment. The main advantage of this correction is to get symmetric volcano plots.

4.- The parameters P correction allows selecting the correction method for the p values calculated during the analysis.

The section Column numbers of Region 2 contains six parameters. Here, users provide the column numbers in the Data file from where UMSAP will get the information needed to perform the analysis of the module. All columns specified in this section must be present in the Data file. Users must be aware that Python starts counting from 0. Therefore, the number of the columns in the Data file starts from 0 and not from 1. The column numbers displayed in the list box of Region 3 after the Data file is selected can be directly used for the values of the parameters.

1.- The parameter Detected proteins allows users to specify the column in the Data file containing the protein identifiers found in the Data file. Only one integer number equal or greater than zero will be accepted here.

2.- The parameter Gene names allows users to specify the column in the Data file containing the gene names of the proteins found during the MS experiments. Only one integer number equal or greater than zero will be accepted here.

3.- The parameter Score allows to specify the column in the Data file containing the Score values. It is in this column where the program will look for the values to be compared against the Score threshold given in section Values of Region 2 of the interface. Only one integer number equal or greater than zero will be accepted here.

4.- The parameter Exclude proteins allows to specify several columns in the Data file. Proteins found in these columns will be excluded from the analysis. The module assumes that these columns contains numeric values and values greater than zero indicate that the respective protein must be eliminated from the analysis. Only integer numbers equal or greater than zero will be accepted here. A value of NA means that all proteins will be considered during the analysis.

5.- The parameter Columns to extract allows users to specify which columns in the Data

file will be copy to the shorter version of the Data file, see page ?? for more details. A range of columns may be specified as 4–10 with both numbers included in the range. Any number of columns may be specified here. Only integer numbers equal or greater than zero will be accepted. A value of NA means no shorter version of the Data files will be created.

6.- The parameter Results - Control experiments allows users to specify the columns in the Data file containing the results of the experiments. There are two ways to provide the information for this parameter. Users can load the values from a .txt file using the Load values button or use the Type values button to call a helper window, see Figure 8.2. Duplicate column numbers are not allowed here.

The helper window is divided in four Regions. Region 1 allows to define the number of conditions and relevant points analyzed, to define the kind of control experiment performed and to create the matrix in Region 2. The fields Names allow to input the names for the conditions and relevant points. A comma separated list of names is expected here. In the case of the name of the control experiment only one name is expected. If left empty default name values will be used. Each text field in Region 2 should contain the column numbers containing the MS results for the given experiment. The values for the text fields should be positive integer numbers or a range of integers, e.g. 60–62 or NA for empty experiments. Fields left empty are set to NA when the values are exported to the window of the module. The column numbers can be seen in the list box in Region 3. Selected entries in the list box can be copied and then pasted to the text fields using the right mouse button or the Tools menu. Region 4 contains two buttons to Cancel or to Export the values to the window of the module.

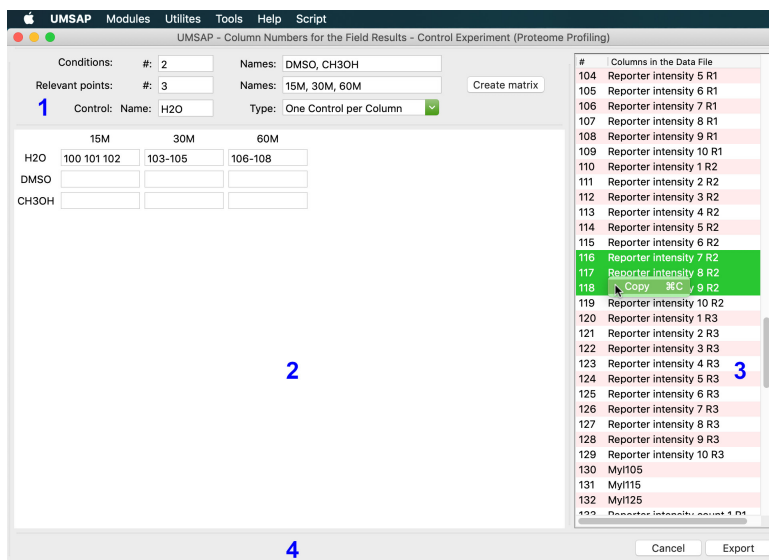


Figure 8.2: The Result - Control experiments helper window for the Proteome Profiling module. This window allows users to specify the column numbers in the Data file containing the MS results for the selected conditions, relevant points and control experiments.

The column numbers, the labels for conditions, relevant points and control experiments

and the information for the controls can also be loaded from a .txt file. The format of the file content is very simple. The first four lines give the values for the labels and control and the rest of the lines specify the column numbers with a comma (,) separating the values for a different condition - relevant point. The following is an example for two conditions, two relevant points and one control for each condition:

```
Control type : One Control per Row
Control name : MyControl
Condition names: DMSO, H2O
Relevant point names: 30min, 1D
```

```
105 115 125, 106 116 126, 101 111 121
130 131 132, 108 118 128, 103 113 123
```

The values for Control type are the same as in the helper window.

Region 3 of the Proteome Profiling module main window contains a list box that will display the number and name of the columns found in the Data file. The list box is automatically filled when the Data file is selected. Selected columns in the list box can be directly added to any field in section Column numbers of Region 2 of the interface using the right mouse button over the list box or the Tools menu.

Region 4 contains two buttons and the progress bar. The Help button leads to an online tutorial while the Start analysis button will trigger the processing of the Data file. The progress bar will give users a rough idea of the remaining processing time after the analysis is started.

8.3.1 The Tools menu

The tools menu in the module window allows to copy the selected columns in the list box in Region 3 of the interface to the fields of section Column numbers of Region 2 of the interface. The list box in Region 3 of the interface can also be cleared. In addition, through this menu users can create a .usr file with the given options to the module before running the analysis. If something goes wrong during the analysis having the .usr file means that users do not have to type the values of all the parameters again, see ?? for more details.

8.4 The analysis

First, UMSAP will check the validity of the user provided input. In particular, all experiments need to have the same number of replicates as the respective control. Then, the Data file is processed as follow. All proteins found in the Exclude proteins columns are discarded. Proteins that were not identified in all conditions are discarded. Finally, all proteins with a Score value lower than the defined threshold are removed. The intensity values of the remaining proteins are normalized and then a median correction is applied to each experiment if Median correction was selected in section Values. With

the resulting intensity values the fold change for each protein and experiment as well as the intensity ratios with respect to the control experiments are calculated and two different analysis are performed.

The fold change is calculated as:

$$FC = ave(I_{C,RP})/ave(I_{Control}) \quad (8.1)$$

The first analysis is a t-test to determine if each experiment is significantly different to the corresponding control. The second analysis is a t-test, or ANOVA test if more than two conditions are studied, to determine if the values for the studied conditions are significantly different for each selected relevant point.

Finally, the corrected p values are calculated.

8.5 The output files

The Proteome Profiling module generates up to three files and a folder named Data.Steps, see Figure 8.3. The folder Data.Steps contains a step by step account of all the calculations performed so users can check the accuracy of the calculation or perform further analysis. The files inside Data.Steps are plain text file with tab separated columns. The first line contains the name of the columns in the file.

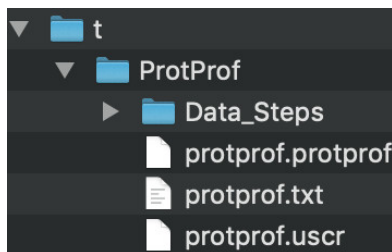


Figure 8.3: The structure of the Output folder from the Proteome profiling module. The file `protprof.txt` will be created only if requested.

The three files created have extension `.txt`, `.uscr` and `.protprof`. The file with extension `.txt` contains all the lines in the Data files but only the columns specified with the parameter Columns to extract in section Column numbers of Region 2 of the Proteome Profiling module. This file is generated only if the value of the parameter Columns to extract is not NA.

The file with extension `.uscr` contains all the input given by the user so a new analysis may be performed without typing all the option values again, see ??.

The file with extension `.protprof` is the main output of the module. This file contains all the results and can be used to generate the graphical representation of the results.

8.6 Visualizing the output files

After creating the .protprof at the end of the analysis, the Proteome Profiling module will automatically load the file and create a windows to display the results, see Figure 8.4. This window is divided in four Regions.

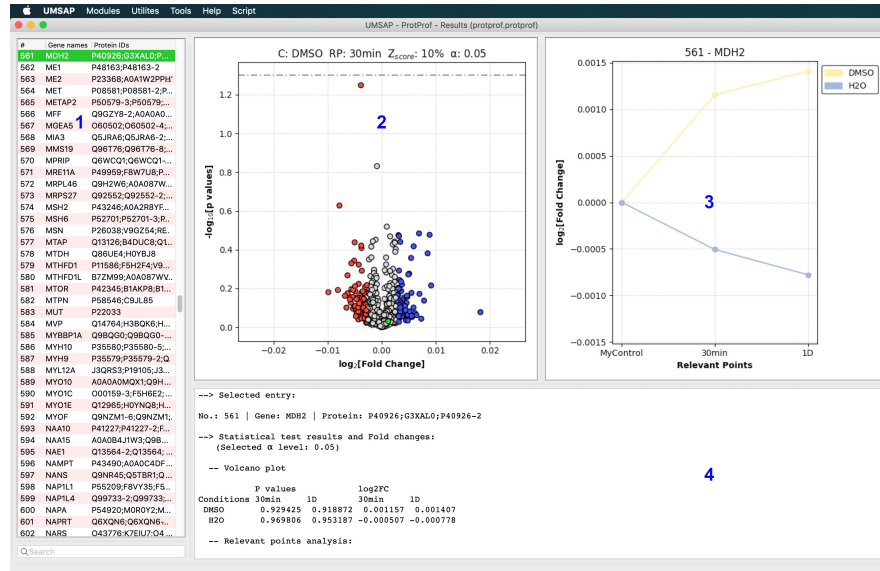


Figure 8.4: The Proteome Profiling analysis window. Users can performed here the analysis of the proteome profiling.

Region 1 contains a list of all protein IDs and Gene names contained in the .protprof file being shown. The search box at the bottom allows to search for a protein in the list. Selecting a protein in the list will highlight the protein in Region 2 and display information about it in Regions 3 and 4.

Region 2 contains a volcano plot showing the results for the t-test comparing the condition (C), relevant point (RP) to the corresponding control. The volcano plot has a horizontal line indicating the chosen significance level. In addition, the points in the plot can be colored by Z-score allowing to quickly identified the top up (blue) or down (red) regulated proteins. In Figure 8.4, the top 10 % up and down regulated proteins are colored and α is set to 0.05. Selecting a protein in the plot will highlight the selected protein in Region 1 and display information about it in Regions 3 and 4. See subsection 8.6.1 for more options.

Region 3 contains a plot of $\log_2[FC]$ vs Relevant points. The plot allows to see the behavior of the FC along the relevant points for each condition tested in the experiments. See subsection 8.6.1 for more options.

Region 4 shows a summary of the results for a selected proteins. Proteins can be selected in the listbox in Region 1 or in the volcano plot of Region 2. The information includes a summary of the selected protein including the number of the protein in the listbox, the gene name and the protein id. Calculated p and $\log_2[FC]$ values as well as averages and standard deviations for intensities and ratios.

8.6.1 The Tools menu

The Tools menu for the results window of the Proteome Profiling module allows to further customize the plots and to apply different filters to the protein list shown in the window.

Under the menu entry Volcano Plot, users can change the condition and relevant point shown in the volcano plot, the Z score value used to color the points in the plot and the α value. An image of the volcano plot can also be created. If the condition or relevant point displayed is changed the Apply Filters menu entry allows to recalculate the filters for the current condition, relevant point displayed.

Under the menu entry Relevant Points, users can show all the proteins at once with different or the same colors and create an image of the Relevant Points graph.

The menu entry Export Data can be used to export the data shown in the window to a plain text file (see ??) while the Corrected P values entry will display the information in the .protprof file using the corrected p values instead of the regular p values.

8.6.1.1 Filters

The menu entry Filters allows to Add or Remove the filters applied to the protein list. The idea behind Filters is to identify proteins with a desired behavior and discard the rest of the proteins from the listbox in Region 1 and the plots in Regions 2 and 3. Filters are applied to the current Condition and Relevant point shown in the Volcano plot in Region 2. If the Condition or the Relevant point shown is changed, the new plot will show the proteins obtained after the filter was applied. This allows to follow the behavior of the filtered proteins in all Conditions and Relevant points. The menu entry Applied Filters in the Volcano plot submenu allows to recalculate the filters based on the new Conditions and Relevant point shown. Any number of filters can be applied. The applied filters are shown in the bottom left corner of the results window.

Filter can be removed in any given order using the menu entry Any in the Remove filter submenu. Additionally, the last applied filter can be removed with the menu entry Last Added or the shortcut Ctrl/Cmd + Z.

Currently, the implemented filters are:

Z score

This menu entry allows to filter proteins by the Z score value of the Fold change.

Log2FC

This menu entry allows to filter proteins by the absolute value of the $\log_2[FC]$.

P value

This menu entry allows to filter proteins by the p value calculated for the comparison of the currently displayed Condition and Relevant point to the control experiment. The threshold p value can be given in the 0 to 1 range or as a $-\log_{10}$ value. Regular or corrected P values can be used in the filter.

α value

This menu entry allows to filter proteins by the p values calculated for the comparison of the relevant points. The returned list of proteins consists of proteins for which the calculated p value is less than the selected α value, for at least one relevant point.

Monotonic

This menu entry allows to filter proteins by the behavior of the $\log_2[FC]$ along the Relevant points. The filter searches for proteins that have a monotonically increasing or decreasing (or both) behavior for the condition shown in the Volcano plot.

Divergent

This menu entry allows to filter proteins by the behavior of the $\log_2[FC]$ along the Relevant points. In this case, the filter searches for proteins that have a monotonically increasing and decreasing behavior in at least two of the conditions tested.

Bibliography

1. G. B. Limentani, M. C. Ringo, F. Ye, M. L. Bergquist, E. O. MCSorley, *Analytical Chemistry* **77**, 221 A–226 A, ISSN: 0003-2700 (June 2005).