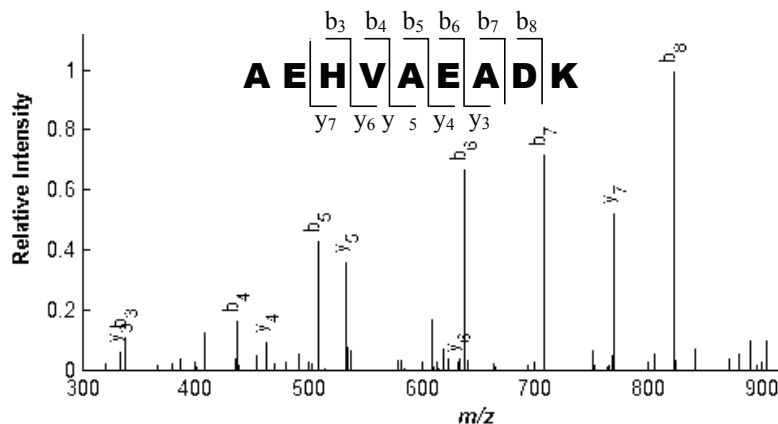


Figure 2. An example of a MS/MS spectrum and its annotation



Peptide identification is the core of protein sequencing in that a protein sequence is identified from its constituent peptides with the help of protein databases. There are three different approaches for PFF peptide identification: database searching, *de novo* sequencing, and sequence tagging.

Database searching is the most widely used method in high-throughput proteomics. It correlates an experimental MS/MS spectrum with theoretical ones generated from a list of peptides digested *in silico* from proteins in a database. As annotated MS/MS spectra increase, a spectrum library searching method which compares an experimental spectrum to previously identified ones in the reference library has recently proposed (Craig et al., 2006; Frewen et al., 2006). *De novo* sequencing tries to infer the complete peptide sequence directly from the m/z differences between peaks in a MS/MS spectrum without any help of databases. Sequence tagging yields one or more partial sequences by *de novo* sequencing, then finds homologous sequences in the protein database and finally scores the homologous candidate sequences to obtain the true peptide.

Robust data preprocessing, scoring scheme, validation algorithm, and fragmentation model are the keys to all of the above three approaches.

MAIN FOCUS

Preprocessing of MS/MS Spectra

The purpose of data preprocessing is to retain the potentially contributing data while removing the noisy or

misleading ones as many as possible to improve the reliability and efficiency of protein identification. In general, only a relatively small fraction (say, <30%) of a large MS/MS spectra set is identified confidently due to some reasons. First, there exist some chemical or electrical noises in a spectrum resulting in random matches to peptides. Second, the peptide ion charge state of a spectrum is usually ambiguous due to the inability of instruments. Third, there are considerable duplicate spectra most likely representing a unique peptide. Four, many poor-quality spectra often have no results which can be filtered in advance using classification or regression methods. Clustering algorithms (Tabb et al., 2003; Beer et al., 2004), linear discriminant analysis (LDA) (Nesvizhskii et al., 2006), quadratic discriminant function (Xu et al., 2005), SVM (Bern et al., 2004; Klammer et al., 2005), linear regression (Bern et al., 2004), Bayesian rules (Colinge et al., 2003a; Flikka et al., 2006), decision tree and random forest (Salmi et al., 2006) have been widely used to address the above problems.

Scoring Scheme for Peptide Candidates

The goal of scoring is to find the true peptide-spectrum match (PSM). In principle there are two implementation frameworks for this goal: descriptive or probabilistic.

In descriptive frameworks, an experimental and a theoretical MS/MS spectra are represented as vectors $\mathbf{S} = (s_1, s_2, \dots, s_n)$ and $\mathbf{T} = (t_1, t_2, \dots, t_n)$, respectively, where n denotes the number of predicted fragments, s_i and t_i are binary values or the observed and predicted intensity

5 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/data-mining-protein-identification-tandem/10862

Related Content

A Multi-Agent System for Handling Adaptive E-Services

Pasquale De Meo, Giovanni Quattrone, Giorgio Terracina and Domenico Ursino (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition* (pp. 1346-1351).

www.irma-international.org/chapter/multi-agent-system-handling-adaptive/10996

Bitmap Join Indexes vs. Data Partitioning

Ladjel Bellatreche (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition* (pp. 171-177).

www.irma-international.org/chapter/bitmap-join-indexes-data-partitioning/10816

On Interactive Data Mining

Yan Zhao (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition* (pp. 1085-1090).

www.irma-international.org/chapter/interactive-data-mining/10956

Realistic Data for Testing Rule Mining Algorithms

Colin Cooper and Michele Zito (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition* (pp. 1653-1658).

www.irma-international.org/chapter/realistic-data-testing-rule-mining/11040

Text Mining Methods for Hierarchical Document Indexing

Han-Joon Kim (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition* (pp. 1957-1965).

www.irma-international.org/chapter/text-mining-methods-hierarchical-document/11087