

# THE PSEUDO-GRADIENT ALGORITHM FOR RESIDUAL GAS ANALYSIS

## PSEVDOGRADIENTNA METODA ZA ANALIZO MASNIH SPEKTROV

Igor Belič

Institute of metals and technology, Lepi pot 11, 1000 Ljubljana, Slovenia  
igor.belic@imt.si

*Prejem rokopisa – received: 2008-12-19; sprejem za objavo – accepted for publication: 2009-01-12*

The article focuses on a special approach to the qualitative and quantitative residual gas analysis of a vacuum-chamber atmosphere. The main outline of mass-spectrometry usage is given together with a brief comparison of two completely different scientific fields. The new approach, where the mass spectra are formalized in a vector annotation, to residual gas analysis is proposed. The main problem with residual gas analysis based on a mass spectrum is its ambiguity, which cannot be easily resolved. One way to deal with the problem is practical knowledge about the interpretation of mass spectra. The pseudo-gradient algorithm for mass-spectra analysis was developed and tested. The general testing platform was developed and the algorithm was tested within it. The presented work is the foundation for a comprehensive study of mass-spectrum analysis techniques. An important concept is the virtual environment that provides the mass-spectrum generator, the standard fragmentation patterns' database, the data space for the various algorithms that can test various approaches to the mass-spectrum analysis, and the database of the achieved results, which is necessary for the comparison of different algorithms and represents the backbone for the dynamic mass-spectra generation and analysis.

Key words: mass spectra, residual gas analysis, gradient descent, optimization

V prispevku je opisan nov način kvalitativne in kvantitativne analize rezidualnih plinov atmosfere v vakuumskem sistemu. Uvodoma je podan pregled različnih vrst uporabe masne spektrometrije, skupaj s pregledom dveh povsem različnih znanstvenih področij. Prikazan je nov način, ki pojmuje masne spektre v formalni vektorski obliki. Osnovni problem pri analizi masnih spektrov, ki ga ni mogoče enostavno rešiti, je njihova večličnost. Ena od možnosti je uvedba praktičnih izkušenj, ki algoritmu pomaga pri razrešitvi nekaterih dvoumij. Razvita in analizirana je psevdo-gradientna metoda za rezidualno analizo plinov iz masnega spektra. Razvito je bilo preizkusno okolje, v katerem je bil analiziran algoritem.

Predstavljeno delo je temelj za široko študijo tehnik za analizo masnih spektrov. Pomemben koncept je virtualno okolje, ki vsebuje generator masnih spektrov, podatkovno bazo standardnih masnih spektrov, podatkovni prostor za hranjenje rezultatov preizkusov algoritmov, ki so nujni za kasnejšo primerjavo delovanja posameznih algoritmov. Virtualno okolje je tudi temelj za kasnejšo dinamično analizo masnih spektrov.

Gljučne besede: masna spektrometrija, rezidualna analiza plinov, gradientna metoda, optimizacija

## 1 INTRODUCTION

The use of mass spectrometry is very widespread and versatile. Many types and various configurations of mass spectrometers are in use. Basically, there are two different ways that mass spectrometry is used: one is for the detection of substances through the fragmentation patterns detected by the mass spectrometers and the other is the use of mass spectrometry in vacuum science. The first is a very wide field that includes chemical, pharmaceutical, biological (from now on bio-chemical) and other purposes, mainly for the study of the structure of substances. The main feature of mass spectrometry that deals with the structures of numerous organic substances (combined with other methods – gas chromatography, etc.) is that one substance (or a small number of them) of a very complex and often unknown structure is introduced to the mass-spectrometer analyzer and the obtained mass spectrum is then analysed.<sup>1,2</sup> Finding the fragmentation patterns of organic substances is a very complicated process, but necessary in order to implement the automatic mass-spectra identification.<sup>3</sup> For this

purpose several software products were developed and are widely used. To illustrate the unbalance in software support that utilizes the mass-spectra analysis in bio-chemical and vacuum fields, several products are listed.

The bio-chemical field is very well supported by numerous commercial software packages. One such software product is the Mascot mass-database search program.<sup>4</sup> This is a powerful search engine that uses mass spectrometry data to identify proteins from primary-sequence databases.<sup>5,6</sup> Another one is Analyst software.<sup>7</sup> A Microsoft Windows-based data system that provides instrument control and data analysis for the entire family of Thermo Scientific mass spectrometers and related instruments is the Xcalibur.<sup>8,9</sup> Several specialized software modules have been designed to work with Xcalibur to meet the needs of specific applications. MassLynx<sup>10,11</sup> is a fundamental platform for acquiring, analysing, managing and sharing mass-spectrometry information. The Sample List is a key feature of MassLynx, keeping everything about the samples together in one place. It is also the central place for initiating any activities relating to the sample. It provides

general purpose and specialised application managers dedicated to providing information for specific types of mass-spectrometry analyses and data. Agilent Technologies MassHunter Workstation software<sup>12</sup> is designed to streamline mass-spectrometry analysis, from instrument tuning through to the final report. It enables users to find all the compounds in samples, find differences between samples and find sample sets. The MassHunter workstation uses data-mining techniques to search the mass databases in public and private domains. Geneva Bioinformatics (GeneBio) has launched a software platform known as SmileMS for the identification and analysis of small molecules by mass spectrometry.<sup>13</sup> SmileMS provides machine-independent software for the analysis of MS spectra from small molecules. It allows each user to add both public and private databases. Mascot Distiller<sup>5</sup> is used to identify the mass spectra corresponding to proteins from the Information Non-Redundant Protein Sequences Database. The list of commercial software dealing with bio-chemical mass spectrometry is long and versatile<sup>14</sup> (4000 Series Explorer, Analyst QS, BioAnalyst, Cliquid, DiscoveryQuant, GPS Explorer, LightSight, MALDI Imaging, MarkerView, Metabolite ID, MRMPilot, MultiQuant, Pro ICAT, Pro ID, Pro QUANT, ProteinPilot, SimGlycan, TissueView, etc.).

These mass spectra are obtained by various types of instrumentation layouts, such as tandem mass spectrometers MS/MS<sup>3,15,16,17,18,19</sup>, collision-induced dissociation (CID) MS/MS<sup>20</sup> time-of-flight mass spectrometers (TOF)<sup>16,18,21,22,23,24</sup>, ion-trap mass spectrometers<sup>19</sup>, quadrupole/ion-trap mass spectrometers<sup>15,25</sup>, and the MALDI (Matrix-assisted laser desorption/ionization) technique used for mass spectrometry<sup>6,23</sup>, etc.

In the bio-chemical field mass spectrometry is widely used in order to identify the peaks resulting from a chromatographic separation.<sup>26</sup> The most common approach to solve the problem for unknowns on whom very little other structural information is available is the use of a retrieval algorithm and a reference mass-spectra database.

There are large reference databases of mass spectra at the NIST and Wiley libraries. The NIST/EPA/NIH mass spectral databases<sup>27</sup> contain, in the 2008 version, 220,460 spectra of 192,108 unique compounds (for electron-impact mass spectrometry), and 14,802 spectra of 5308 precursor ions for tandem mass spectrometry.

The largest database is the Wiley's Registry of Mass Spectral Data<sup>28</sup>, in 2008 containing 560,000 different spectra, over 348,000 spectra with chemical structures. It is one of the most comprehensive mass spectral libraries ever published. There are also specially designed mass-spectra search programs to make use of these vast databases.

In vacuum science there are again two major categories of how mass spectrometry is used. The first, which is in terms of operation very close to the bio-chemical

uses, is the vacuum-system leak detection.<sup>29</sup> The tracer gas is introduced to the system and again one gas fragmentation pattern is monitored by a mass spectrometer.

Another and yet much more complicated category is residual gas analysis, which is a general term for the analysis of gas and vapour species in vacuum chambers and vacuum processes. The main feature of all the previously mentioned methods is that the main objective of it is the detection and analysis of various species that compose the vacuum-chamber atmosphere simultaneously. It is no more the case where one or a very small number of species is to be detected and compared to the known spectra in a database. Here, the number of "expected" vacuum-chamber atmosphere constituent gases is to be combined in such a way as to provide the partial pressures of each and every one of them.

A search over the internet reveals that the field is far less equipped with software than is the case with bio-chemical mass spectrometry. Among the features found in the literature<sup>30</sup>, the software (MASsoft<sup>30</sup>, Merlin Research<sup>31</sup>, Questor 5 Process Analysis software<sup>31</sup>) provides the following interesting features:

- Histogram, Trend Analysis and Analog peak displays.
- Mixed mode scanning, e.g., Trend, Histogram and Analog peaks in multiple-windows
- Simultaneous real-time display of graphical and tabular trend analysis data.
- Peak-height identification.
- Real-time background subtraction.
- Automatic mass-scale alignment.
- Statistical analysis of data in real time
- Partial pressure ratios
- Quantitative or normalised composition analysis
- Visualisations

In this article we focus on a special approach to qualitative and quantitative residual gas analysis.

### 1.1 Mass spectra as vectors

The fragment patterns of atoms and molecules that may coexist in a vacuum system are composed of a sequence of number pairs, as shown by the CH<sub>4</sub> example in **Table 1**.<sup>32</sup>

Looking at this annotation from another perspective, one can see that the sequence of numbers (**Table 1** left) represents only a fraction of all possible  $m/u$  components (only from  $m/u = 12$  to  $m/u = 17$ ). The complete "picture" of the CH<sub>4</sub> standard fragmentation pattern is presented in **Table 1**, on the right.

The same is true for all the other constituents of a vacuum-chamber atmosphere. The fragmentation pattern can be rewritten in vector form as follows:

$$\mathbf{G}_{\text{CH}_4} = (0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 2.4, 7.7, 15.6, 85.8, 100, 1.2) \quad (1)$$

**Table 1:** CH<sub>4</sub> fragmentation pattern**Tabela 1:** Fragmentacijski vzorec CH<sub>4</sub>

$m/u$	Normalized amplitude	$m/u$	Normalized amplitude
12	2.4	1	0
13	7.7	2	0
14	15.6	3	0
15	85.8	4	0
16	100	5	0
17	1.2	6	0
		7	0
		8	0
		9	0
		10	0
		11	0
		12	2.4
		13	7.7
		14	15.6
		15	85.8
		16	100
		17	1.2

This is a classical vector annotation, where the number of all the used  $m/u$  components defines the dimensionality of the vector space.

Any peak in a mass spectrum of a mixture of gases may consist of a combination of molecular ions and/or fragment ions. The contributions add linearly and the peak height for the discrete  $m/u$  is equal to the sum of the individual peak heights that would be produced if each constituent were alone in the system. This is the classical presumption of superposition.<sup>32</sup>

$$S_i = \sum_j s_{ij} \quad (2)$$

The sum in (2) is taken over all gases currently present in the vacuum chamber. The consecutive gas is marked by the index  $j$ ;  $S_i$  is the total peak height at the mass number  $i$ ;  $s_{ij}$  is the peak height contribution from gas  $j$  at the  $m/u$  number  $i$ .

$s_{ij}$  is related to the fragmentation pattern, the analyzer sensitivity, and the partial pressure of gas  $j$  by the equation:

$$s_{ij} = A_j P_j g_{ij} \quad (3)$$

Where  $A_j$  is the analyzer's sensitivity to the gas  $j$ ;  $P_j$  is the partial pressure of gas  $j$ ; the principal peak of the standard fragment pattern is always set to 1. Other peaks at  $m/u = i$  represent the ratio to the principal peak and are denoted by  $g_{ij}$ .

In a formal vector annotation the mass spectrum is a weighted sum of the constituent fragment patterns, such as:

$$S = \sum_j w_j G_j; \quad w_j = a_j p_j \quad (4)$$

$S$  is the vector of the total peak heights;  $w_j$  is the weight, valid for the  $j$ -th constituent fragment pattern; and vector  $G_j$ , is the  $j$ -th gas-standard fragment pattern.

However, the equation exactly models the ideal situation, where all the constituent components of the atmosphere are known a priori. In practice one must always expect the unknown substances to be present in the vacuum system and contribute to the mass spectra. Another problem that additionally complicates the interpretation of mass spectra is the inevitable noise produced by the instrumentation equipment. Both components add to the so-called noise. The equation can be therefore rewritten as:

$$S = \sum_j w_j G_j + N; \quad w_j = a_j p_j; \quad N = S_u + N_e \quad (5)$$

Where  $N$  is the vector representing the noise. The overall noise  $N$  is composed of two components:  $S_u$  – the spectral components of unknown constituents, and  $N_e$  – the electrical noise.

### 1.2 The ambiguity of mass spectra

The equations (4) and (5) formally represent the synthesis of the mass spectra. On the other hand, from the practical point of view, we are interested in a reverse process, i.e., analysis. The basic question of mass spectrometry is which constituents are composing the given mass spectra and in what quantities.

From the mathematical theory, it is very well known that such a task is only achievable under the very strict condition of orthogonality, which must be fulfilled for all the constituent vectors. If this is not the case, the result can still be achieved, but it is never unique.

**Table 2:** The example of mass-spectra ambiguity

**Tabela 2:** Primer večličnosti masnega spektra

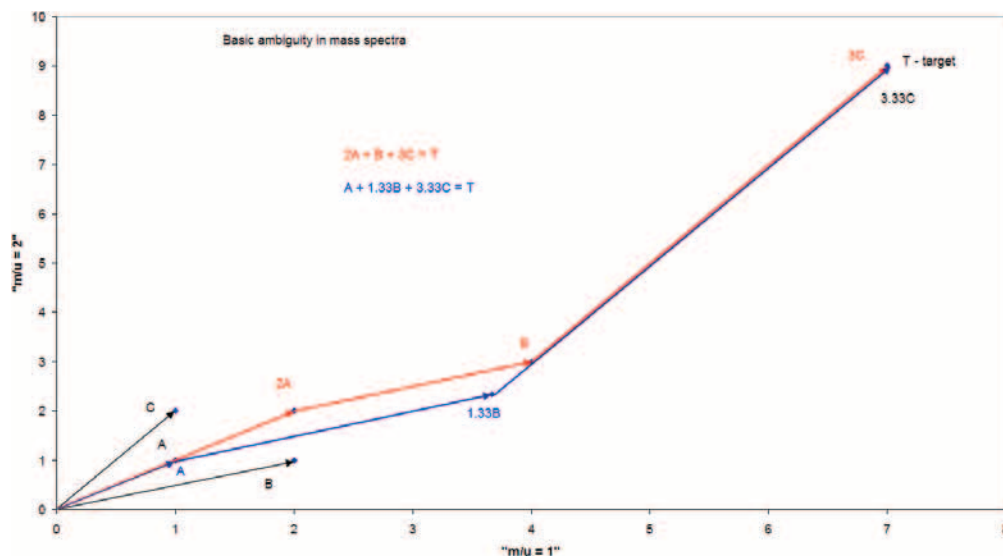
Standard Fragmentation Patterns	m/u	Constituent A	Constituent B	Constituent C	
	1	1	2	1	
	2	1	1	2	
Case 1	Multiplication factors 1	2	1	3	
	m/u	Constituent A	Constituent B	Constituent C	Point T
	1	2	2	3	7
	2	2	1	6	9
Case 2	Multiplication factors 2	1	1.333	3.333	
	m/u	Constituent A	Constituent B	Constituent C	Point T
	1	1	2.666	3.333	7
	2	1	1.333	6.666	9

Let us suppose we have three constituent gases, and that we observe only two  $m/u$  values:  $m/u = 1$  and  $m/u = 2$ .

The experiment to demonstrate the basic ambiguity of the mass spectra is shown in **Table 2**. For this purpose, the standard fragmentation patterns for three virtual gases were set. The next step is to generate the mass spectra that contains the three gases. For this, the multiplication factor for each constituent must be set and the mass spectra is created, according to equation (4), taking into account that all the sensitivity factors  $a_j$  are

set to 1. The result of the mass-spectra synthesis is the actual spectra, or "target" point T, which represents the superposition of all three constituent gases, multiplied by their factors.

The next step is the inverse process, meaning that we have the mass spectra – point T, and we seek the multiplication factors that produce that mass spectra. The figures in **Table 2** clearly show that the same mass spectra can be achieved by the factors 2, 1, 3, and 1, 1.333, 3.333 as well. The experiment is graphically illustrated in **Figure 1**.



**Figure 1:** The illustration of mass-spectra ambiguity

**Slika 1:** Ilustracija večličnosti masnega spektra

When the mass spectra are considered to be vectors, the ambiguity in its interpretation can be clearly demonstrated.

### 1.3 Interpretation of mass spectra for residual gas analysis – general considerations

It is obvious that from knowing only the resulting mass spectra it is next to impossible to reconstruct the proper contribution for each constituent spectrum.

The first step in the mass-spectra analysis is the identification of the residual gases that constitute the observed spectrum.

It is to be expected that the major residual constituents in vacuum systems are inorganic gases, such as  $\text{H}_2\text{O}$ ,  $\text{H}_2$ ,  $\text{CO}_2$ ,  $\text{CO}$ ,  $\text{O}_2$ ,  $\text{N}_2$ ,  $\text{Ar}$ ,  $\text{Ne}$ , etc. Most inorganic gas molecules found in vacuum systems are composed of 2, 3, or 4 atoms. Their fragment patterns are therefore simple. These gases are relatively easy to identify in the overall mass spectrum.

There are some simple rules at the basis of the experimental work with mass spectrometry. These rules are extremely important for the understanding of mass-spectra properties. As such they can provide the valuable additional information needed for a more accurate mass-spectra identification.<sup>32</sup>

- The spectra of simple inorganic gas molecules generally have even mass numbers. For these inorganic gas molecules a simple ionization is more probable than fragmentation.
- The parent peaks of almost all inorganic gases are the largest. They are found at an even mass number.
- Mass spectra should not contain large quantities of highly reactive gases ( $\text{F}_2$ ,  $\text{O}_2$ , etc.). Generally these gases are most easily adsorbed by the vacuum system walls. In cases where the mass-spectra analysis gives large quantities of such gases, the algorithm should be capable of suppressing their quantities (or to signal possible erroneous events, e.g., vacuum-system leaks)
- If the mass peak at  $m/u = 14$  (mostly  $\text{N}^+$ ) is larger than the mass peaks at  $m/u = 12$  and  $16$  ( $\text{C}^+$  and  $\text{O}^+$  from  $\text{CO}$ ) then an air leak is very possible.
- The noble gases,  $\text{He}$ ,  $\text{Ne}$ ,  $\text{Ar}$ , etc., are usually not the dominant residual gases. They have a highly unreactive nature and they can scarcely be found in the atmosphere. Like with the atmosphere, in vacuum systems where gas reactivity is important for efficient pumping, noble gases (especially  $\text{Ar}$  and  $\text{Ne}$ ) can be observed in appreciable quantities.
- A  $m/u = 1$  peak greater than a few percent of  $m/u = 2$  is often a sure indication of significant amounts of water vapor in the system. Please note that except for water,  $\text{H}^+$  is not a significant fragment of any gas usually seen in a vacuum system. Even with the standard spectra for  $\text{H}_2\text{O}$ , the  $m/u = 1$  is not plotted. The partial pressures of gases like  $\text{H}_2\text{O}$ ,  $\text{CO}$ ,  $\text{CO}_2$ ,

etc. can provide useful information about the progress of vacuum-system bakeouts.

- A mass spectrum with a series of peaks that are separated by  $\Delta m/u = 14$  or  $\Delta m/u = 15$  indicates that there are hydrocarbons present in the vacuum system.
- Organic gases are unstable, thus the parent peaks may not appear in the spectrum. Fragments often represent the dominant mass peaks.
- Fragments with odd  $m/u$  numbers are expected to be the most populated.
- The  $m/u$  peaks 57, 55, 53; or 43, 41, 39; or 29, 27; or combinations of these are a sure indication of organic species in the vacuum system. If the highest  $m/u$  of each of these series is the largest peak then the organic is of the saturated type, (forepump oil). If a lower mass number in the series is the largest, this is generally an indication of some degree of unsaturation (Multi-bonded carbon-carbon) and can be caused by some polymeric substance.

### 1.4 Descent techniques

One of the numerous methods available to solve the problem of mass-spectrum identification is the use of the descent method. There are, in fact, several descent methods. In our research, the classical gradient method has been altered to the so-called pseudo-gradient descent method. We are well aware that any used method cannot overcome the basic problem of mass spectra: the ambiguity, therefore, is that one cannot expect the results of the proposed method to be flawless. Our aim is to find a method that will provide results that are good enough for practical use.

Descent techniques<sup>33</sup> are generally used for the solution of unconstrained minimization problems. In our case we are trying to find the combination of standard fragment patterns that yield the target, i.e., measured spectra. The standard fragment patterns are formalized as vectors, and we are seeking the set of multiplication factors, i.e., the weight that in the weighted sum forms the target vector. The distance from the calculated weighted sum to the target vector is measured by means of the classical Euclidean distance. The distance function is also often called the error function. This is the classical optimization problem. The unconstrained minimization problem is  $\min D(\mathbf{w}, \mathbf{x})$ . Here, a value of the variable vector  $\mathbf{w} = (w_1, \dots, w_n)^T$  is sought that minimizes the objective, in our case the distance function  $D(\mathbf{w}, \mathbf{x})$ . The distance is calculated in  $m$ -dimensional Euclidean space, spanned over the unit vector  $\mathbf{x} = (x_1, \dots, x_m)^T$ . The problem  $\min D(\mathbf{w}, \mathbf{x})$  is a special case of the general nonlinear programming or optimization problem. For the sake of simplicity, the distance function  $D(\mathbf{w}, \mathbf{x})$  is often denoted by  $D(\mathbf{w})$ , bearing in mind that the distance function is primarily changed by  $\mathbf{w}$ , but is calculated in the space defined by  $\mathbf{x}$ .

A new descent direction is generated for each iteration. The descent iteration may involve the evaluation of the first and possibly the higher order derivatives of the objective distance function. Each step of the descent process yields a considerable improvement of the objective function.

The descent techniques involve a series of iterations that generally consist of three parts:

- 1) Determination of the direction of descent  $s^k$ ,
- 2) Determination of the descent length  $\lambda^k$ ,
- 3) Calculation of the descent step  $w^{k+1} = w^k + \lambda^k w^k$ .

The descent direction is an  $n$ -dimensional vector  $s = (s_1, \dots, s_n)^T$ . It exists in the same space as the weights do and represents the direction in which weights should be changed in order to decrease the distance function. At the  $k$ -th iteration the direction vector  $s^k$  originates at the current point  $w^k$ . It points in a descent or "downhill" direction, i.e., the value of the objective function decreases from  $w^k$  to a point at some distance in that direction. A unit vector  $s^k$  is said to be the descent direction, with respect to the objective function  $D(w)$  at  $w^k$ , if there is a  $\lambda_p > 0$ , such that for all  $\lambda$  satisfying  $\lambda_p \geq \lambda > 0$  we have

$$D(w^{k+1}) = D(w^k + \lambda s^k) < D(w^k) \quad (6)$$

If  $D(w)$  is differentiable,  $s^k$  is a descent direction if

$$\lim_{\lambda \rightarrow 0} \frac{D(w^k + \lambda s^k) - D(w^k)}{\lambda} = \left. \frac{dD(w^k + \lambda s^k)}{d\lambda} \right|_{\lambda=0} = (s^k)^T \nabla D(w^k) < 0 \quad (7)$$

Where the  $\nabla D(w)$  denotes the gradient of the objective function  $D(w)$  evaluated at the point  $w^k$ . If  $D(w)$  is differentiable, the product  $(s_k)^T \nabla D(w)$  of the directions  $s^k$  and the gradient  $\nabla D(w)$  is, by definition, the directional derivative of  $D(w)$  in the direction  $s^k$  evaluated at  $w^k$ . If this directional derivative exists and is negative, then  $s^k$  is a descent direction.

The descent step-length  $\lambda$  is a scalar. It is a measure of the distance along the descent direction  $s^k$  between two successive iteration points  $w^k$  and  $w^{k+1}$ . In other words, at the  $k$ -th iteration a step of length  $\lambda^k$  is taken from point  $w^k$  to the point  $w^{k+1}$ .

### The Descent Iteration

A typical descent iteration can be summarized in the following steps:

- 1) Compute a descent direction  $s^k = (s_1^k, \dots, s_n^k)^T$
- 2) Compute a descent step-length  $\lambda^k$
- 3) Perform a descent step to obtain a new point

$$w^{k+1} = w^k + \lambda^k s^k \quad (8)$$

The  $k$ -th descent step is defined as follows

$$\Delta w^k = w^{k+1} - w^k = \lambda^k s^k \quad (9)$$

A sequence of  $k$  descent steps leads from a starting point  $w^0$  to a point  $w^k$  given by

$$w^k = w^0 + \sum_{l=0}^{k-1} \lambda^l s^l = w^0 + \sum_{l=0}^{k-1} \Delta w^l \quad (10)$$

At the  $k$ -th iteration a matrix  $\Delta \bar{W}_k$  is defined by

$$\Delta \bar{W}_k = [\Delta w^0, \Delta w^1, \dots, \Delta w^{k-1}] \quad (11)$$

That is, the columns of  $\Delta \bar{W}_k$  are the  $k$  descent steps  $\Delta w^0, \Delta w^1, \Delta w^{k-1}$  preceding  $\Delta w^k$ .

The choice of the locally steepest direction as a descent direction leads to the steepest descent techniques. A locally steepest direction is obtained if the descent step  $\Delta w^k$  minimizes

$$\Delta D = \sum_{i=1}^n \frac{\partial D(w^k, x^k)}{\partial w_i} \Delta w_i^k \quad (12)$$

The distance between two points  $x^1$  and  $x^2$  in the  $m$ -dimensional space is defined as

$$D(x^1, x^2) = [(x^1 - x^2)^T A (x^1 - x^2)]^{1/2} \quad (13)$$

$A$  is a positive definite  $m \times m$  symmetric metric matrix. The definition of positive definiteness ensures that  $D(x^1, x^2) > 0$  for any nonzero  $x^1$  and  $x^2$ . In general, it is not necessary to use the same unit of distance along the different coordinate axes.

The choice of  $I$ ,  $m \times m$  identity matrix as a metric, i.e.,  $A = I$ , leads to the first-order steepest descent technique. Rescaling of the variables, e.g.,  $\bar{x} = Bx$ , is equivalent to introducing a new metric matrix relative to the old coordinate system  $x$ . It is not necessary to use the same metric matrix throughout the whole iterative process of a descent technique.

## 2 EXPERIMENTAL WORK

Our aim is to design an algorithm that is able to perform qualitative and quantitative mass-spectra analysis. It is not possible to undertake the task without the formation of a virtual, simulated environment where various aspects of mass-spectra analysis can be addressed. In real vacuum systems, it is very expensive and time-consuming work to create a large amount of different situations to test primarily the algorithms. We want to establish the situation where the creation of analyzed spectra is completely under our control. Undivided attention can therefore be placed on the studied algorithms. All the problems are then concentrated on algorithms alone. When in a situation where the behaviour of the algorithm is very thorough, it can be implemented on real vacuum systems, giving the chance to assess its behaviour in real circumstances. At this stage of the implementation of the algorithm in "real life" is a process where the majority of problems is due to the used equipment, while the algorithm's behaviour is known, tested and understood.

Such an approach is not a very common one, and it is almost general practise that two problems are studied at

once, i.e., the mass spectra in a real situation and the algorithm.

2.1 The environment

The environment that binds together the generation of the mass spectra, the various algorithms for the mass-spectra identification, the experiments database and the evaluation procedures was established. The first step is to create the environment where static mass-spectra analysis can be performed.

By static, it is meant that the mass spectrum is created and the algorithm to analyze the spectrum is activated. During the analysis time, the mass spectrum is fixed. Our intention is to create the environment where mass spectrum is generated as the time variant spectrum, i.e., dynamic.

The environment (Figure 2) enables the study of behaviour of different algorithms in noisy conditions.

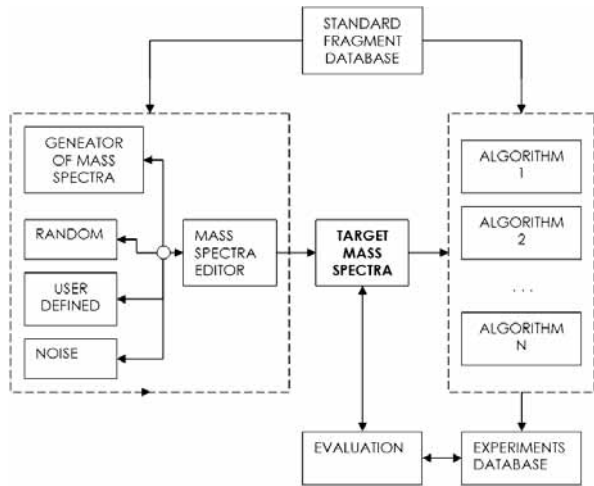


Figure 2: The environment for the generation of mass spectra, algorithms for the mass-spectra constituents analysis, the experiments database and evaluation

Slika 2: Okolje za generiranje masnih spektrov, algoritmi za analizo, podatkovna baza za shranjevanje podatkov poskusov in njihovih obdelav

m/u	1	2	3	4	5	6	7	8	9	10	11	12	13
	H2	He	CH4	NH3	H2O	Ne	C2H2	C2H4	N2	CO	C2H6	NO	CH4O
1	0.0500	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
2	1.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
3	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
4	0.0000	1.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
5	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
6	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
7	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
8	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
9	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
10	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
11	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
12	0.0000	0.0000	0.0240	0.0000	0.0000	0.0000	0.0250	0.0210	0.0000	0.0450	0.0000	0.0000	0.0000
13	0.0000	0.0000	0.0770	0.0000	0.0000	0.0000	0.0560	0.0350	0.0000	0.0000	0.0000	0.0000	0.0000
14	0.0000	0.0000	0.1550	0.0220	0.0000	0.0000	0.0020	0.0630	0.0720	0.0050	0.0340	0.0750	0.0000
15	0.0000	0.0000	0.8580	0.0750	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0460	0.0240	0.0000
16	0.0000	0.0000	1.0000	0.8000	0.0110	0.0000	0.0000	0.0000	0.0000	0.0090	0.0000	0.0150	0.0000
17	0.0000	0.0000	0.0120	1.0000	0.2300	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
18	0.0000	0.0000	0.0000	0.0040	1.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0190
19	0.0000	0.0000	0.0000	0.0000	0.0010	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
20	0.0000	0.0000	0.0000	0.0000	0.0030	1.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
21	0.0000	0.0000	0.0000	0.0000	0.0030	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
22	0.0000	0.0000	0.0000	0.0000	0.0000	0.9990	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
23	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
24	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0560	0.0370	0.0000	0.0000	0.0000	0.0000	0.0000
25	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.2010	0.1170	0.0000	0.0000	0.0420	0.0000	0.0000

Figure 3: The standard mass-spectra pattern database (section)

Slika 3: Podatkovna baza standardnih masnih spektrov (izrez)

Name	Size
Exp1_1	42 KB
Exp1_2	42 KB
Exp1_3	42 KB
Exp1_4	42 KB
Exp1_5	42 KB
Exp1_6	42 KB
Exp1_7	42 KB
Exp1_8	42 KB
Exp1_9	42 KB
Exp1_10	42 KB
Exp1_11	42 KB
Exp1_12	42 KB
Exp1_13	42 KB
Exp1_14	42 KB
Exp1_15	42 KB
Exp1_16	42 KB
Exp1_17	42 KB
Exp1_18	42 KB
Exp1_19	42 KB
Exp1_20	42 KB
Exp1_21	42 KB
Exp1_22	42 KB
Exp1_23	42 KB
Exp1_24	42 KB
Exp1_25	42 KB

Figure 4: The structure of the experiments database

Slika 4: Struktura podatkovne baze poskusov

The core of the environment is the database, including the standard fragment patterns for molecular ions and/or fragment ions (Figure 3).

In our experiments, 47 gases, with  $m/u$  ratios from 1 to 47 formed the database. However, by no means is the system limited to these figures. The database leaves the user to freely add or remove the standard fragmentation patterns. The database is a vital element of the environment since it provides the standard fragment patterns as vectors. The data is used for the generation of the mass spectra, which is done in two different ways. One is random generation, where the multiplication factors  $w_j$  from equation (4) are chosen randomly. Another option is the user-defined selection of  $w_j$ , where the user can make the selection of gases to compose the vacuum-system atmosphere. The third option makes it possible to add noise to the generated spectrum – equation (5). Once the spectrum is generated it can be edited by the user.

The spectrum-generation process results in the vector annotation of mass spectra. As such it is prepared to be analyzed by various algorithms.

Each experiment is saved in the so-called experiments database, which consists of separate files, each for a separate experiment. Figure 4 depicts the structure of the experiments database.

Each experiment file holds the data of the generated spectrum and the results of its analysis by the algorithm.

The following data and the graphical representations describe the experiment:

- 1) Each experiment starts with the random generation (Figure 2) of virtual vacuum system atmosphere composition in terms of multiplication factors  $w$ . Multiplication factors are stored, and shown in

ITERATION	DISTANCE TO TARGET	Min	COMPOSITION	CALCULATION	GAS	COMPOSITION	CALCULATION	NUMBER OF ITERATIONS
1	43.3708193	1	0.3528	0.3071	H2	7.0555	7.0555	2031
15	0.333843511	3	0.3065	0.0000	H4	5.3342	5.3337	
102	0.511781417	4	5.3342	5.3337	3CH4	5.7952	5.7231	RMS error 0.1728
153	0.418172299	5	0.0000	0.0000	4RH3	2.8958	2.9618	
204	0.382576178	6	0.0000	0.0000	5HD	3.9195	1.6038	0.1728
255	0.302156221	7	0.0000	0.0000	6Ma	7.7474	7.7437	
306	0.296292137	8	0.0000	0.0000	7CH2	0.1402	0.4475	0.1728
357	0.265289587	9	0.0000	0.0000	8C2H4	7.6072	6.2495	
408	0.248817983	10	0.0000	0.0000	9N2	8.1443	5.8928	0.1728
459	0.228797754	11	0.0000	0.0000	10CO	7.0904	8.4584	
510	0.219876477	12	1.3689	1.3215	11C2H6	0.6535	2.5953	0.1728
561	0.205376054	13	0.9045	0.8931	12NO	4.1403	2.6926	
612	0.200082926	14	3.2771	3.2439	13CH4O	8.6262	8.2227	0.1728
663	0.195043284	15	6.3935	6.1152	14O2	7.9548	8.1767	
714	0.187456842	16	10.5885	10.5842	15H2S	3.7354	3.7354	0.1728
765	0.178423953	17	4.8728	4.6755	16A	9.6195	9.6208	
816	0.178249957	18	3.8748	3.8165	17C2H2	6.7445	8.7023	0.1728
867	0.168873905	19	0.1794	0.1888	18C3H8	0.5624	0.6353	
918	0.167444196	20	8.7857	8.7882	19CO2	9.4954	7.8174	0.1728
969	0.158982987	21	0.0252	0.0252	20N2O	3.6402	5.3183	
1020	0.158987269	22	0.8838	0.8617	21C2H4O	5.2487	5.1583	0.1728
1071	0.151980317	23	0.0000	0.0000	22C2H6O	7.6111	8.0771	
1122	0.151520278	24	0.3731	0.3388	23H2O	0.5368	0.6036	0.1728
1173	0.144489602	25	1.1931	1.1865	24CH2O2	5.5048	5.6398	
1224	0.143716743	26	7.5891	7.5947	25C2H10	4.6870	4.7023	0.1728
1275	0.137419367	27	12.9908	12.8891	26C3H8O	2.8811	2.8362	
1326	0.136706161	28	26.8744	26.8691	27C2H3	6.2273	6.2263	0.1728
1377	0.130718878	29	21.7336	21.7257				
1428	0.120053873	30	6.5761	6.4763				0.1728
1479	0.124355156	31	16.3175	16.3159				
1530	0.120689152	32	15.3253	15.2923				0.1728
1581	0.118603488	33	1.8632	1.8535				
1632	0.117724758	34	3.7670	3.7681				0.1728
1683	0.113612133	35	2.5779	2.5777				
1734	0.112084044	36	0.1868	0.1885				0.1728
1785	0.109204348	37	2.0448	2.0441				
1836	0.107681023	38	1.8748	1.8749				0.1728
1887	0.104848188	39	7.2396	7.2422				
1938	0.103727511	40	12.1625	12.1622				0.1728
1989	0.100912234	41	10.3645	10.3514				
		42	7.9788	7.9785				0.1728
		43	9.7785	9.7754				
		44	16.2742	16.2763				0.1728
		45	5.9104	5.8077				
		46	5.1170	5.1251				0.1728
		47	1.6071	1.6072				

Figure 5: The data stored in the experiment file  
Slika 5: Podatki, shranjeni v datoteki, ki pripadajo enemu poskusu

Figure 5 in rightmost table (the column entitled COMPOSITION). The same data is graphically represented in Figure 6, lower graph. Please note that the randomly generated data is presented as the left (blue) bars in a bar graph.

- Once the randomly generated composition is known, the mass spectrum is formed (Equation 4), where standard fragmentation patterns (Figure 3) and the corresponding weights are combined to form the spectra (the column COMPOSITION in the middle table of Figure 5). The same data is graphically shown in Figure 6 middle bar graph. Again the generated spectrum is shown by the left (blue) bars. These data form the so called "target" for the mass spectrum analysis algorithm.
- At this stage the spectrum generation process is complete the data is stored and graphically presented. The reverse process of spectrum analysis can commence. The pseudo gradient algorithm (or any other) works on a spectrum (column COMPOSITION in middle table of Figure 5) and seeks the multiplication factors  $w$  that best fit the target spectrum. The calculated weights are stored in column CALCULATION in rightmost table of Figure 5. The same data is graphically presented in lower bar graph (the right (red) bars) of Figure 6. Similarly the calculated spectrum is stored in column CALCULATION in middle table of Figure 5, and graphically presented in middle bar graph (again the right (red) bars).
- The data that describe the convergence of the used algorithm (the leftmost table of Figure 5) consisting of consequent iteration counter (column ITERATION), and the current distance to the given target spectrum (column DISTANCE TO TARGET) is stored. The same data is graphically represented in Figure 6 – upper graph.

- The lower right table in Figure 5 summarizes the five mass spectra constituents where the errors produced by the tested algorithm were absolutely the largest.
- In addition the number of iterations required to fulfil the preset algorithm termination criteria, and the overall RMS error for all mass spectra constituents (Figure 5 upper right part) are presented.

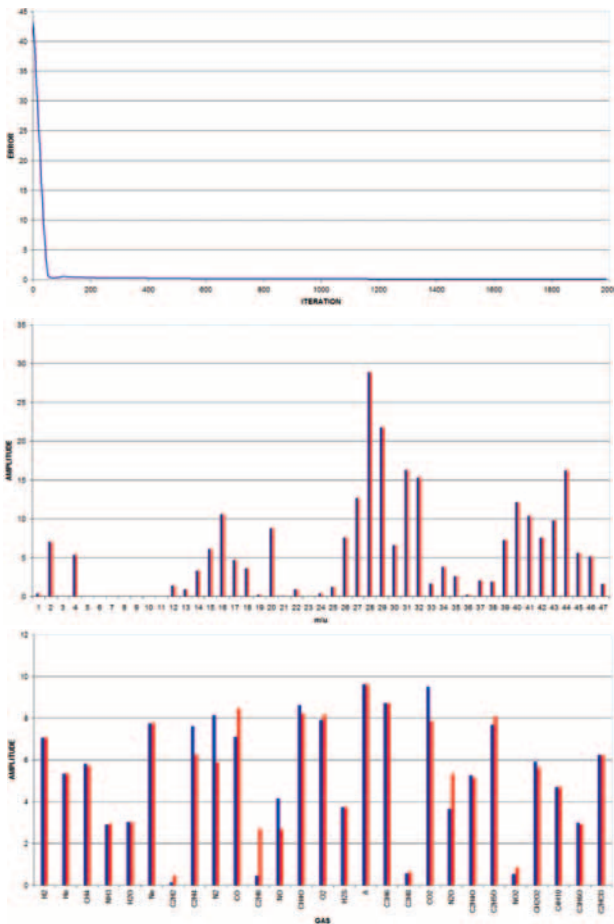
In this stage, the system is prepared to generate and analyze static mass-spectra vectors. The design of the environment is such that a dynamic module can easily be added in order to generate and simultaneously analyze the time profiles of the mass spectra.

### 2.2 The pseudo-gradient descent algorithm

The classical gradient approach requires the calculation of the first derivative of the distance to the target function  $D(w, x)$ . The distance function in the  $D(. , x)$  is the standard Euclidean distance (13), unfortunately the dependence  $D(w, .)$  is more complicated. In our approach we have decided to calculate the differences produced by the small changes of  $w$  rather than produce the necessary step directions from the formal first derivative of the function. Therefore, the formal gradient descent algorithm becomes the pseudo-gradient algorithm.

- Read the data from the standard fragment pattern database.
- Calculate increments for all gases (These are the differences used for the calculation of pseudo gradients. Increments mean that each  $m/u$  component of the standard fragment pattern is multiplied by a small value – 0.001).
- Define the memory data space for the error function. For each iteration, the error value represents the distance from the current mass spectrum produced as the result of current weight values, and the target spectrum. The error function is the dependence of the error and the iteration counter.
- Read the target spectrum.
- Set the initial weights for all possible gases included in the database to 0.01 – this is the starting point where the gradient descent starts.
- Calculate the first distance from the spectra that uses the initial weights and the target spectrum.
- Calculate the step-length  $\lambda^k$  for the current iteration. The algorithm uses the variable step-length  $\lambda^k$ , which is set to 1/20 of the current distance (Equation 13). This step is necessary to ensure the smooth convergence of the algorithm.
- Start the iteration loop, which is executed until the actual distance remains larger than 0.1 or the repetition counter remains under 3000.
  - Calculate the differences in distance that are caused by making very small differences – increments in the gas factors – weights. Both directions are





**Figure 6:** The graphical representation of mass-spectra analysis data  
**Slika 6:** Grafična predstavitev podatkov analize masnih spektrov

probed and the change that produces the reduction of distance to the target is maintained. Special care must be taken during reduction, while the weights for each constituent gas must always remain positive. The lowest possible weight value is therefore 0.

- o Keep only incremental changes that produce a reduction of the distance to the target.
- o Find the direction of downhill change – the pseudo gradient (12).
- o Use the step-length and direction vector to form the new, closer set of weights.
- o Calculate the new distance (6).
- o Calculate the new step-length  $\lambda^{k+1}$ , again as 1/20 of the actual distance to the target.
- As the algorithm executes iterations, the data structures (**Figures 5 and 6**) are consequently filled with data.

### 3 RESULTS AND DISCUSSION

The mass spectrum is the linear superposition of the individual peak heights of the constituent gases. The

standard fragmentation patterns of the constituent gases are not orthogonal; therefore, the algorithms that are used for the mass-spectra identification produce ambiguous results.

The experiment was designed to probe the pseudo-gradient descent algorithm for 1000 randomly generated spectra.

The pseudo gradient algorithm shows excellent convergence. In all 1000 examples the convergence was without any detected instabilities. For all the performed tests the error function is a monotonously decreasing one. It is important to note that the algorithm seeks the values for weights that produce the mass spectrum as close as possible to the given spectra. This means that the end spectrum is always as near as possible to the preset error tolerance in the  $m/u$  "space", while the actual weight values can sometimes produce serious errors in the weight space.

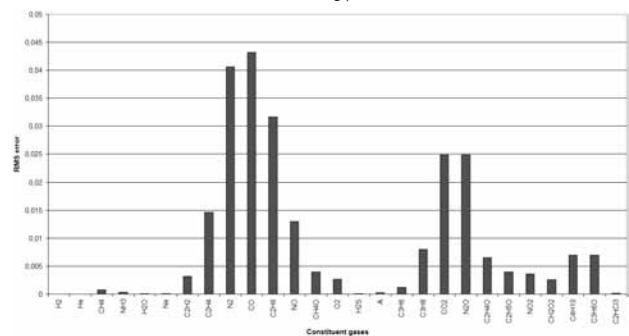
The average starting error for all 1000 tests was 23.84, and the error tolerance to stop the algorithm was arbitrarily set to 0.1.

A total of 142 tests out of 1000 did not fulfil the pre-set condition of error tolerance in 3000 iterations. All the others did the job in an average 802 iterations. The average error for the 142 failed tests was 0.14. The problem of 142 unfinished tests could be easily solved by moving the number of allowed iterations from 3000 to some higher value. The arbitrarily set error tolerance is also very easy to change – according to the needs posed by the concrete problem.

**Figure 7** plots the root-mean-square errors (of the calculated weights) for all the constituent gases and for all 1000 tests. Since the mass spectrum analysed by the tested algorithm is always exactly known (it is generated by the mass-spectrum generator), the error of the performed analysis can always be calculated (not only assessed, as is the case during the analysis of real data).

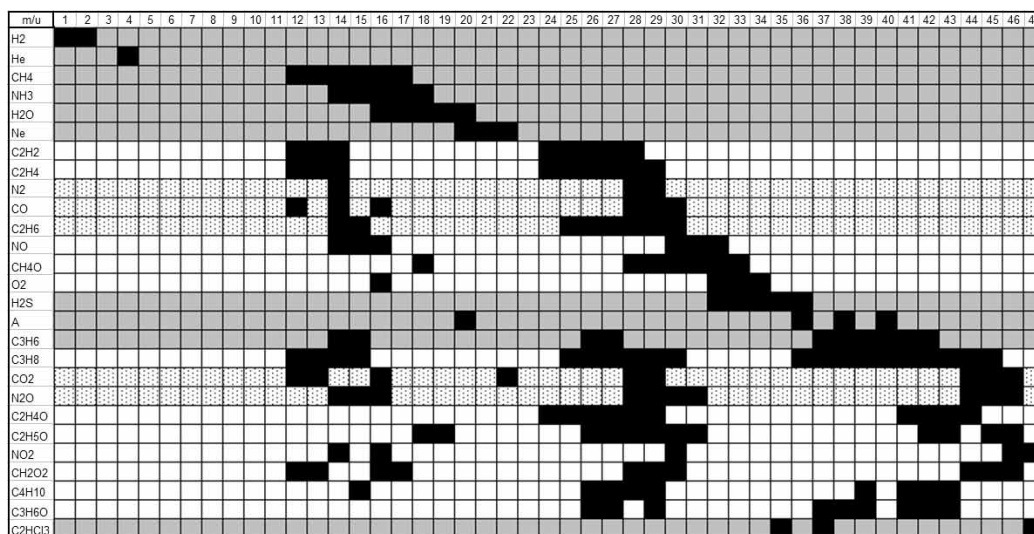
The root-mean-square error (RMS) is calculated by equation 14.

$$E_j = \frac{\sqrt{\sum_{i=1}^N (w_{ij} - wt_{ij})^2}}{N} \quad (14)$$



**Figure 7:** Root-mean-square error for the calculated weights of the constituent gases

**Slika 7:** Srednjekvadratna napaka izračunanih uteži plinov, ki sestavljajo atmosfero v vakuumskem sistemu



**Figure 8:** Mass-spectra standard fragment patterns – the black squares represent the presence of a peak  
**Slika 8:** Standardni masni spektri – črni kvadrati predstavljajo prisotnost posameznega vrha v spektru

Where  $E_j$  represents the RMS error for the  $j$ -th constituent gas, the index  $i$  is the test counter that runs from 1 to 1000,  $w_{ij}$  is the weight of the  $j$ -th constituent gas (the result of the analysis) during the  $i$ -th test,  $wt_{ij}$  is the target weight (of the generated spectrum) of the  $j$ -th constituent gas during the  $i$ -th test. The value for  $N$  is, in our case, 1000.

From **Figure 7** we can conclude that there is a group of constituent gases where analysis gives excellent results, almost without noticeable errors. Such gases are:  $H_2$ , He,  $CH_4$ ,  $NH_3$ ,  $H_2O$ , Ne,  $H_2S$ , A,  $C_3H_6$ ,  $C_2HCl_3$ .

The group of constituent gases with a "medium" error would consist of  $C_2H_2$ ,  $C_2H_4$ , NO,  $CH_4O$ ,  $O_2$ ,  $C_3H_8$ ,  $C_2H_4O$ ,  $C_2H_6O$ ,  $NO_2$ ,  $CH_2O_2$ ,  $C_4H_{10}$ ,  $C_3H_6O$ .

The largest errors are detected for  $N_2$ , CO,  $C_2H_6$ ,  $CO_2$ ,  $N_2O$ .

The reason for the errors is mainly in the mentioned mass-spectra ambiguity. Putting it another way: the standard fragmentation patterns for the gases overlap, so it is very hard to distinguish which is which.

In **Figure 8** the existence of the peak in the fragment pattern is symbolized as a black square. The constituent gasses that can be recognised precisely are shown with a grey background, while the dotted area represents the gasses with the poorest results.

#### 4 CONCLUSION

The presented work is a foundation for the comprehensive study of mass-spectrum analysis techniques. It introduces several new ideas in the field of mass spectrometry. The most important concept is the virtual environment that provides the mass-spectrum generator, the space for the various algorithms that can test various approaches to the mass-spectrum identification, the database of the achieved results, which is extremely important for the comparison of different algorithms, the

backbone for the dynamic mass-spectra generation and analysis. The environment also makes it possible to analyse the immunity to noise for all algorithms.

The main purpose of the virtual environment is the controlled mass-spectra data, which allows an exact evaluation of the errors produced by the studied algorithm. Normally, such algorithms are developed and tested directly on "live" data, which combines the problems that originate from the measurements environment with those produced by the algorithm, often without a real chance of dividing the two. The approach with the virtual environment makes it possible to study problems regarding the algorithms alone, prior to applying them to the "live" environment.

Such an approach also makes it possible to run numerous tests in a relatively short time, which would never be possible with real vacuum systems.

The pseudo-gradient descent algorithm for mass-spectrum identification was proposed and tested. The results are promising, taking into account that additional work will be needed to overcome the mass-spectrum ambiguity. The steps that will be tested are:

- 1) Introduction of "common knowledge" regarding the mass spectrometry and identification techniques.
- 2) Finding out how does the selection of the initial point influence the algorithm's performance.
- 3) Introduction of several steps in the identification process – after the first run of the algorithm the constituent gases that can be detected with the highest precision should be removed and the algorithm should be run again for the remaining gases. This step should be repeated several times.

Another problem is the study of a dynamic mass spectrum, i.e., the analysis of mass-spectrum time

profiles. The virtual environment is prepared to enable such studies.

## 5 REFERENCES

- <sup>1</sup> K. Schulz, K. Schlenz, S. Malt, R. Metasch, W. Römhild, J. Dreßler, D. W. Lachenmeier, *Journal of Chromatography A*, 1211 (2008), 1–2, 113–119
- <sup>2</sup> I. S. Sheoran, A. R. S. Ross, D. J. H. Olson, V. K. Sawhney, *Plant Science*, 176 (2008), 99–104
- <sup>3</sup> J. Chena, X. Li, C. Suna, Y. Pana, U. P. r Schlunegger, *Talanta*, 77 (2008), 152–159
- <sup>4</sup> A. Cingöz, F. Hugon-Chapuis, V. Pichon, *Journal of Chromatography A*, 1209 (2008), 95–103
- <sup>5</sup> <http://www.matrixscience.com/>
- <sup>6</sup> H. M. Santosa, C. Mota, C. Lodeiroa, I. Mouraa, I. Isaacb, J.L. Capelo, *Talanta*, 77 (2008), 870–875
- <sup>7</sup> <https://products.appliedbiosystems.com/>
- <sup>8</sup> <http://www.thermo.com/com/cda/product/detail/0,1055,1000001009250,00.html>
- <sup>9</sup> T. Murakamia, T. Kawasakia, A. Takemuraa, N. Fukutsua, N. Kishia, F. Kusud, *Journal of Chromatography A*, 1208 (2008), 164–174
- <sup>10</sup> [http://www.scientific-computing.com/products/product\\_details.php?product\\_id=207](http://www.scientific-computing.com/products/product_details.php?product_id=207)
- <sup>11</sup> S. Shia, Y. Zhaob, H. Zhouc, Y. Zhanga, X. Jianga, K. Huangd, *Journal of Chromatography A*, 1209 (2008), 145–152
- <sup>12</sup> <http://www.chem.agilent.com>
- <sup>13</sup> <http://scientific-computing.com>
- <sup>14</sup> <http://www3.appliedbiosystems.com>
- <sup>15</sup> L. Lionetto, A. M. Lostia, A. Stigliano, P. Cardelli, M. Simmaco, *Clinica Chimica Acta*, 398 (2008), 53–56
- <sup>16</sup> M. Mezcuca, C. Ferrer, J. F. García-Reyes, M.J. Martínez-Bueno, M. Sigrist, A. R. Fernández-Alba, *Food Chemistry*, 112 (2009), 221–225
- <sup>17</sup> K. M. Robinson, J. T. Morr, J.S. Beckman, *Archives of Biochemistry and Biophysics*, 423 (2004), 213–217
- <sup>18</sup> M. Y. Zhanga, N. Kagana, M. A. Sungb, M. M. Zaleskab, M. Monaghanb, *Journal of Chromatography B*, 874 (2008), 51–56
- <sup>19</sup> W. C. Chau, J. Wu, Z. Cai, *Chemosphere*, 73 (2008), S13–S17
- <sup>20</sup> X. Deng, G. Gao, S. Zheng, F. Li, *Journal of Pharmaceutical and Biomedical Analysis*, 48 (2008), 562–567
- <sup>21</sup> L. Qi, J. Cao, P. Li, Q. Yu, X. Wen, Y. Wang, C. Li, K. Bao, X. Ge, X. Cheng, *Journal of Chromatography A*, 1203 (2008), 27–35
- <sup>22</sup> H. Qu, B. Li, X. Li, G. Tu, J. Lü, W. Sun, *Microchemical Journal*, 89 (2008), 159–164
- <sup>23</sup> E. Vanlaere, K. Sergeant, P. Dawyndt, W. Kallow, M. Erhard, H. Sutton, D. Dae, B. Devreese, B. Samyn, P. Vandamme, *Journal of Microbiological Methods*, 75 (2008), 279–286
- <sup>24</sup> R. Wihlborg, D. Pippitt, R. Marsili, *Journal of Microbiological Methods*, 75 (2008), 244–250
- <sup>25</sup> J. Radjenović, S. Péreza, M. Petrović, D. Barcelóa, *Journal of Chromatography A*, 1210 (2008), 142–153
- <sup>26</sup> R. Oprean, L. Oprean, M. Tamas, R. Sandulescu, L. Roman, *Journal of Pharmaceutical and Biomedical Analysis*, 24 (2001), 1163–1168
- <sup>27</sup> <http://www.sisweb.com/software/ms/nist.htm>
- <sup>28</sup> <http://eu.wiley.com/WileyCDA/Section/id-301546.htm>
- <sup>29</sup> A. Calcatelli, M. Bergoglio, D. Mari, *Vacuum*, 81 (2007), 1538–1544
- <sup>30</sup> <http://www.hiddenanalyticalinc.com>
- <sup>31</sup> <http://www.extrel.com>
- <sup>32</sup> M. J. Drinkwine, D. Lichtman, *Partial pressure analyzers and analysis*, American Vacuum Society, Milwaukee, 1979
- <sup>33</sup> J. L. S. Samuel, *Iterative methods for nonlinear optimization problems*. Prentice HALL, Englewood Cliffs, 1972