

Law of Localised Fine Structure^{*}

with application in mass spectrometry proteomic studies.

Mateusz Łacki and Anna Gambin

Faculty of Mathematics, Informatics and Mechanics
University of Warsaw, Banacha 2, 02-097 Warszawa, Poland
mateusz.lacki@biol.uw.edu.pl
aniag@mimuw.edu.pl

Abstract. The distribution of isotopic envelope is approximated indirectly by Poisson laws. We use the devised approximations to design alternative algorithms for calculating the isotopic fine structure with significant coverage.

Keywords: Isotopic Fine Structure, Poisson Approximation, Stable Isotopes, Avergin Model.

1 Introduction

There are many reasons why mass spectrometry analysis is hard. It is hard in that there are potentially many many sources of interferences that can distort the information about the actual composition of a sample. The study of the nature of these interferences is needed to achieve the goal of making out of mass spectrometers yet more reliable an identification tool.

Part of the noise in the mass to charge domain is innately related to the elements themselves and stems from the existence of isotopes. It is because of them that a given analyte is represented as a series of peaks, a spectrum, rather than only one peak. The theoretical underpinnings of how to mathematically model the impact of isotopes are already well established, see [19]. The main idea behind the model is to abstract from the exact positionings of the extra neutrons on a particular chemical compound and thus concentrate only on their relative amounts among all atoms of a given element. Assuming that the isotopic configurations are independent and follow the element dependent distribution, one arrives to the conclusion that the correct law describing occurrence of different isotopes in a chemical compound is the product of multinomial distributions.

There is one huge problem with that law: together with the growth of molecule one observes an exponential growth in the number of possible isotope configurations, which precludes their direct enumeration. To solve this problem, different simplifications were proposed, amounting to different ways of binning configurations together explicitly [4], by hiding them under the guise of Fourier Transform [14], or by ...

^{*} This research was partially supported by Polish National Science Center grant n° 2011/01/B/NZ2/00864.

However, it is considered of paramount importance in Mass Spectrometry to develop machines with still higher resolution powers and it is very likely that this trend will continue. Even today there are machines that already can distinguish peaks attributed to different configurations with the same number of extra neutrons.

Add Olson and others but Olson above all.

Here we propose to approach the problem of fine structure so that it overcomes the shortcomings of the aggregate model, as used in [4]. That particular model bins together configurations having the same number of extra neutrons distributed on different atoms. For instance, if one considers water molecule H_2O , the model would glue together configuration with one extra neutron only on the first hydrogen together with that having it on the second together with that on oxygen atom. We devise an algorithm to deaggregate these probability clusters. We call the peaks obtained via that algorithm a *localised fine structure*.

What motivates the solution to this problem is a search for better molecule fingerprints. The development of new mass spectrometers capable of distinguishing differences in masses of neutrons is proceeding at a vigorous pace. Soon, scientists will face the need of more detailed models than those abstracting from mass defects. It is also common for chemists to search for presence of specific substance in the sample. Usually this is done by looking at some highly specific range in the mass domain of the gathered spectra. Our model provides deeper insight to what might happen while focussing on that particular bit of collected data: conditioning on configurations with the same number of extra neutrons translates directly into focussing in a specific region of the mass-to-charge domain.

The algorithm assumes that one can easily find a peak not far from the most probable one and that the distribution is close to what one would call unimodal¹ and that the most of distributions probability lies in a rather small neighbourhood of the mode. Both the guess about the starting point and the way the neighbourhood gets explored depend on the Poisson approximation to the distribution under study. To our best knowledge this type of approximation have not yet been used for algorithmic purposes. It has been used however in the context of proteomic and peptide research: in [3] it is being used for high throughput protein identification; its use was revalidated in [18] in case of peptides.

We also observe that the use of Poisson approximation gives a theoretical explanation for the equatransneutronic binning used in [13] and actually helps deaggregating results obtained using that approach as well.

Diophantine equations.

¹ We provide a precise definition of unimodality for discrete probability distributions in Section ...

2 Approximations

By an isotopic configuration we understand information on numbers of different isotopes a chemical compound in the sample is made of. For the purpose of simplicity, we focus here on chemical compounds composed of carbon, hydrogen, nitrogen, oxygen, and sulfur; still, results of this section generalize to any compound whatsoever. Thus, we concentrate on compounds like $C_c H_h O_o N_n S_s$, where the low case letters describe the numbers of atoms of particular element type. Among such compounds one can already find peptides and proteins. An isotopic configuration could be represented by an extended empirical formula,

$$^{12}C_{c_0} \ ^{13}C_{c_1} \ ^1H_{h_0} \ ^2H_{h_1} \ ^{14}N_{n_0} \ ^{15}N_{n_1} \ ^{16}O_{o_0} \ ^{17}O_{o_1} \ ^{18}O_{o_2} \ ^{32}S_{s_0} \ ^{33}S_{s_1} \ ^{34}S_{s_2} \ ^{36}S_{s_4}. \quad (1)$$

In the above representation, small letters with indices represent counts of different atoms with indices displaying the number of additional neutrons an isotope has with respect to the highest possible isotopic variant.

Rather than (1), we shall be using an equivalent probabilistic notation, treating upper case letters, like ^{12}C , as random variables and considering small case letters, c_0 , to be their realisations. An expression like $A = \{^{13}C = c_1, \ ^2H = h_1\}$ is shorthand for saying: let us focus on all configurations (1) that have c_1 heavy carbons and h_1 deuters in total.

Following [11], one assumes that the law of vector

$$(^{12}C, \ ^{13}C, \ ^1H, \ ^2H, \ ^{14}N, \ ^{15}N, \ ^{16}O, \ ^{17}O, \ ^{18}O, \ ^{32}S, \ ^{33}S, \ ^{34}S, \ ^{36}S), \quad (2)$$

given $C_c H_h O_o N_n S_s$, is a product of independent multinomial distributions,

$$\mathbb{M} = \text{Multi}\left(\mathbb{P}(^{12}C), \mathbb{P}(^{13}C); c\right) \otimes \cdots \otimes \text{Multi}\left(\mathbb{P}(^{32}S), \mathbb{P}(^{33}S), \mathbb{P}(^{34}S), \mathbb{P}(^{36}S); s\right), \quad (3)$$

where the probabilities of observing particular isotopes, $\mathbb{P}(^{12}C), \dots, \mathbb{P}(^{36}S)$, are established in independent experiments, cf. Table 1. For instance, the probability of a given carbons configuration (c_0, c_1) equals

$$\text{Multi}\left(\mathbb{P}(^{12}C), \mathbb{P}(^{13}C); c\right) \left((c_0, c_1)\right) = \binom{c}{c_0, c_1} \mathbb{P}(^{12}C)^{c_0} \mathbb{P}(^{13}C)^{c_1}$$

and it should be multiplied by similar expression for hydrogen, nitrogen, oxygen and sulfur to obtain probability for expression like (1).

Observe, that given $C_c H_h O_o N_n S_s$, part of the information in (2) is redundant and can be shortened by neglecting counts of the lightest isotope variants, leaving us with

$$(^{13}C, \ ^2H, \ ^{15}N, \ ^{17}O, \ ^{18}O, \ ^{33}S, \ ^{34}S, \ ^{36}S). \quad (4)$$

Missing therms can be retrieved from relationships $^{12}C + ^{13}C = c$, $^1H + ^2H = h$, and so on, that occur with probability one.

Definition 1 *We call the set of configurations*

$$LFS_K = \{^{13}C + ^2H + ^{15}N + ^{17}O + 2 \times ^{18}O + ^{33}S + 2 \times ^{34}S + 4 \times ^{36}S = K\} \quad (5)$$

a localised fine structure with K extra neutrons.

The reason for numbers 2 and 4 appearing above is that ^{18}O and ^{34}S have two additional neutrons, and ^{36}S – four; cf. Table 1.

The problem of enumerating all elements of LFS_K is known as the money exchange problem. In general, it corresponds to finding all integer solutions (x_1, \dots, x_k) of a *Linear Diophantine Equation*

$$d_1x_1 + \dots + d_kx_k = K, \quad (6)$$

where (d_1, \dots, d_k) are integer coefficients. According to [1], if the greatest common divisor of (d_1, \dots, d_k) is equal to one, then the number of solutions to (6) is approximately $\frac{K^{k-1}}{(k-1)!d_1\dots d_k}$. Carbon has only one additional isotope, so $\exists_i d_i = 1$ in (6). The above estimate encompasses therefore all of organic chemistry.

Nonetheless, since configurations in LFS_K are naturally prioritized by probability (3) one would be satisfied with enumerating only the most probable ones.

Problem 1 *For a given K , find a small set $B \subset LFS_K$ of configurations s.t.*

$$\mathbb{M}_K(B) := \frac{\mathbb{M}(B)}{\mathbb{M}(LFS_K)} \approx 1, \quad (7)$$

where \mathbb{M}_K is the product of multinomial laws (3) conditional on the set of configurations in LFS_K and is referred to as **The Law of Localised Fine Structure**.

In statistical terms, we are interested in approximating some critical set of large probability, as measured by the *Law of Localised Fine Structure*.

Why should one study law described by (7) in the first place? Simply because the masses of different configurations in LFS_K concentrate around the compound’s monoisotopic mass shifted to the right by K Daltons; c.f [9]. For medium sized compounds, LFS_K ’s for different K should in principle form disjoint clusters in the mass to charge domain, with some interference for bigger compounds. Studying LFS_K guarantees exploration of a precised place in the mass to charge domain.

To solve Problem 1 we approximate measure \mathbb{M}_K by a more analytically tractable measure \mathbb{Q}_K defined on the LFS_K . We then devise an algorithm to find a possibly small set of configurations $B^* \subset LFS_K$, s.t. $\mathbb{Q}_K(B^*) \approx 1$. Since $\mathbb{Q}_K \approx \mathbb{M}_K$, so $\mathbb{M}_K(B^*) \approx 1$ and B^* solves Problem 1, possibly suboptimally.

A natural way to define proper \mathbb{Q}_K is to first approximate \mathbb{M} by some \mathbb{Q} and then pose $\mathbb{Q}_K(\circ) := \frac{\mathbb{Q}(\circ \cap LFS_K)}{\mathbb{Q}(LFS_K)}$, i.e. condition \mathbb{Q} on the occurrence of configurations from LFS_K . To prove it works, we have to first mention, that by approximation we understand convergence in distribution, as described in [10]. Then, we make use of the following lemma:

Lemma 1. *Let $\mu^{[n]}, \mu$ be discrete measures. If $\mu^{[n]}$ converges in distribution to μ , $\mu^{[n]} \rightharpoonup \mu$, and an event A has nonzero probability under any of that measures, $\forall_n \mu^{[n]}(A), \mu(A) > 0$, then measures conditioned by A , $\mu_A^{[n]}(\circ) := \frac{\mu^{[n]}(\circ \cap A)}{\mu^{[n]}(A)}$ converge in distribution to $\mu_A(\circ) := \frac{\mu(\circ \cap A)}{\mu(A)}$; or $\mu_A^{[n]} \rightharpoonup \mu_A$ for short.*

Proof is to be found in **Appendix**.

Let us now unveil the usefulness of Lemma 1. There is an entire family of measures mentioned in it, $\mu^{[n]}$. We assume, that one of them is simply our initial measure: there exists n^* s.t. $\mathbb{M} = \mu^{[n^*]}$. Also, we assume the approximation of $\mu^{[n^*]}$ by measure μ is already a good one. Our choice for μ is to be the product of independent Poisson measures, which is stimulated by the following, well known lemma.

Lemma 2. *If all $\lim_{n \rightarrow \infty} np_{k,n} = \lambda_k$ exist for $k \in \{1, \dots, w\}$, then*

$$\text{Multi}\left(p_0^{[n]}, p_1^{[n]}, \dots, p_w^{[n]}; n\right) \rightarrow \text{Poiss}(\lambda_1) \otimes \dots \otimes \text{Poiss}(\lambda_w), \quad (8)$$

where Poiss stands for the Poisson distribution, $\text{Poiss}(\lambda)(k) = \frac{\lambda^k}{k!} e^{-\lambda}$.

In Lemma 2 one assumes that the number of trials n goes to infinity. In our model this corresponds to an infinite enlargement of the compound. The existence of limits assumes that this enlargement is done so that on such an idealized compound only the lightest isotopes would appear infinitely often. Moreover, since the support of any Poisson distribution is equal to the set of all integer numbers, the state space of configurations gets significantly enlarged and contains configurations that are nonphysical for any real chemical compound. For instance, positive probabilities would be prescribed to configurations with numbers of isotopes greater than the number of possible places for them on any finite compound. Observe also, that the probabilities $p_k^{[n]}$ are pending towards zero: for good approximation one would expect therefore the probabilities of observing heavier isotopes, e.g. quantities like $\mathbb{P}(^{13}\text{C})$, $\mathbb{P}(^2\text{H})$, \dots , $\mathbb{P}(^{36}\text{S})$, to be relatively small. That is the case – cf. Table 1.

Observe, that Lemma 2 defines a proper limit for just one multinomial distribution, whereas \mathbb{M} is a product thereof. The problem is other than what to do with products: one can approximate independently each multinomial. However, the quality of such approximation depends on all the counts of different elements in a molecule. For instance, in case of $\text{C}_c\text{H}_h\text{O}_o\text{N}_n\text{S}_s$ the better the approximation² the bigger the smallest among numbers (c, h, n, o, s) . Due to the polymer structure, one would expect some more information could be revealed on that matter for proteins and peptides. Indeed, empirical research by Senko et al. [17] established the concept of avergine, i.e. an averaged protein: any protein composed of m amino acids should have its mass approximately equal to the mass of the idealised compound

$$\text{C}_{\lfloor m \times 4.9384 \rfloor} \text{H}_{\lfloor m \times 7.7583 \rfloor} \text{O}_{\lfloor m \times 1.4773 \rfloor} \text{N}_{\lfloor m \times 1.3577 \rfloor} \text{S}_{\lfloor m \times 0.0417 \rfloor}.$$

The weakest link in the approximation might result from small numbers of sulfur. This is an acknowledged problem in empirical studies, as exposed in [18]. The longer the polymers however, the smaller the differences should be.

² The *goodness* of approximation is expressed in the total variance distance; see [15].

The final question is: what values should be used as λ 's in Lemma 2? We *calibrate* those values by equating them to the averages of the multinomial distributions from (3): in case of carbon we set $\lambda_{13C} \approx c \times \mathbb{P}(^{13}C)$. In contrast to our method, λ 's in [3, 18] are chosen to be the minimisers in a free parameter optimisation scheme with χ^2 penalty³.

All in all, the probability assigned to event

$$\{^{13}C = c_1, ^2H = h_1, ^{15}N = n_1, ^{17}O = o_1, ^{18}O = o_2, ^{33}S = s_1, ^{34}S = s_2, ^{36}S = s_4\}$$

is given by

$$\frac{\lambda_{13C}^{c_1}}{c_1!} \frac{\lambda_{2H}^{h_1}}{h_1!} \frac{\lambda_{15N}^{n_1}}{n_1!} \frac{\lambda_{17O}^{o_1}}{o_1!} \frac{\lambda_{33S}^{s_1}}{s_1!} e^{-\mu} \frac{\lambda_{18O}^{o_2}}{o_2!} \frac{\lambda_{34S}^{s_2}}{s_2!} e^{-\eta} \frac{\lambda_{36S}^{s_4}}{s_4!} e^{-\gamma}, \quad (9)$$

where

$$\begin{aligned} \mu &= \lambda_{13C} + \lambda_{2H} + \lambda_{15N} + \lambda_{17O} + \lambda_{33S} \\ \eta &= \lambda_{18O} + \lambda_{34S} \\ \gamma &= \lambda_{36S}. \end{aligned}$$

The usefulness of approximation by a product of independent Poisson lies in two important properties, as summarised in the following lemmas.

Lemma 3. *Suppose we have a collection of m independent Poisson-distributed random variables, $X_i \sim \text{Poiss}(\kappa_i)$. Then $X_1 + \dots + X_m \sim \text{Poiss}(\kappa_1 + \dots + \kappa_m)$.*

Lemma 4. *Suppose we have a collection of m independent Poisson-distributed random variables, $X_i \sim \text{Poiss}(\kappa_i)$. Then X_1, \dots, X_m given that $X_1 + \dots + X_m = K$ is multinomially distributed,*

$$(X_1, \dots, X_m | X_1 + \dots + X_m = K) \sim \text{Multi}\left(\frac{\kappa_1}{\sigma}, \dots, \frac{\kappa_m}{\sigma}; K\right),$$

where $\sigma = \sum_{i=1}^m \kappa_i$.

Both lemmas are proved in [12]. Lemma 3 shows how to simplify calculations for a Diophantine equations with all parameters set to one, $a_i \equiv 1$. Lemma 4 describes the law resulting from conditioning independent Poisson variables by such an expression.

Suppose that we concentrated on molecules composed entirely of elements that can have only one additional neutron, e.g. $C_c H_h N_n$. By Lemma 4 we get:

Result 1 *For $C_c H_h N_n$, let $\tilde{\mu} := \lambda_{13C} + \lambda_{2H} + \lambda_{15N}$. Then*

$$\mathbb{Q}_K = \text{Multi}\left(\frac{\lambda_{13C}}{\tilde{\mu}}, \frac{\lambda_{2H}}{\tilde{\mu}}, \frac{\lambda_{15N}}{\tilde{\mu}}; K\right).$$

³ Note however, that these two solutions should not differ too much for larger compounds, for it is known that both the Poisson and Multinomial distributions are concentrated near their means, see [2].

Proof. The corresponding Diophantine equation is $^{13}\text{C} + ^2\text{H} + ^{15}\text{N}$.

It is valuable to see, how Lemma 4 generalizes while conditioning on a more complex Diophantine equation. Observe, that (5) can be rewritten as

$$LFS_K = \left\{ \underbrace{^{13}\text{C} + ^2\text{H} + ^{15}\text{N} + ^{17}\text{O} + ^{33}\text{S}}_{G_1} + 2 \times \underbrace{(^{18}\text{O} + ^{34}\text{S})}_{G_2} + 4 \times \underbrace{^{36}\text{S}}_{G_4} = K \right\},$$

so that in light of Lemma 3, $\mathbb{Q}(A)$ can be calculated in an easier way:

$$\mathbb{Q}(LFS_K) = \sum_{k_1+2k_2+4k_4=K} \mathbb{P}(G_1 = k_1, G_2 = k_2, G_4 = k_4),$$

where $G_1 \sim \text{Pois}(\mu)$, $G_2 \sim \text{Pois}(\eta)$, and $G_4 \sim \text{Pois}(\gamma)$ are mutually independent. There is a strict link between G_i and the concept of *equatransneutronic groups* described in [13]: it is equal to the total number of atoms bearing exactly i additional neutrons.

To calculate \mathbb{Q}_K it remains to divide (9) by $\mathbb{Q}(LFS_K)$. Observe however that a more significant expression is to be obtained, if additionally we multiply both the nominator and the denominator of that expression by $\frac{\mu^{k_1}}{k_1!} \frac{\eta^{k_2}}{k_2!} \frac{\gamma^{k_4}}{k_4!}$:

Result 2 *The approximative fine structure law with K additional neutrons for $C_e H_h O_o N_n S_s$ is equal to*

$$\text{Multi} \left(\frac{\lambda_{13C}}{\mu}, \frac{\lambda_{2H}}{\mu}, \frac{\lambda_{15N}}{\mu}, \frac{\lambda_{17O}}{\mu}, \frac{\lambda_{33S}}{\mu}; k_1 \right) \otimes \text{Multi} \left(\frac{\lambda_{18O}}{\eta}, \frac{\lambda_{34S}}{\eta}; k_2 \right) \otimes \mathbb{L}(k_1, k_2, k_4),$$

where

$$\mathbb{L}(k_1, k_2, k_4) = \frac{\frac{\mu^{k_1}}{k_1!} \frac{\eta^{k_2}}{k_2!} \frac{\gamma^{k_4}}{k_4!}}{\sum_{k'_1+2k'_2+4k'_4=K} \frac{\mu^{k'_1}}{k'_1!} \frac{\eta^{k'_2}}{k'_2!} \frac{\gamma^{k'_4}}{k'_4!}}. \quad (10)$$

Otherwise stated, the approximative distribution is a mixture of independent multinomial distributions weighted by the \mathbb{L} distribution, which, for lack of name, we shall call the *lucky law*. Under the Poisson approximation, the *lucky law* is the resulting law on the *equatransneutronic configurations*. General expression is to be found in the **Appendix**.

3 Algorithms

Result 2 opens up a new way to do calculations: by using the approximation one reduces the complexity of Problem 1 to that of studying \mathbb{L} . Usually *Lucky Law* is less dimensional, and so the reduction is significant. In proteomics, one easily establishes the set all possible configurations S in a double *for loop* and calculates their probabilities. In general, the problem could be approached using a tailored MCMC algorithm defined on the space of Linear Diophantine Equations.

Having enumerated the *lucky* configurations, we order them by descending probability and select the critical $L\%$ -set S^* . For each configuration (k_1, k_2, k_4) in S^* one can then independently find the $M\%$ and $B\%$ -critical sets of the underlying multinomial distributions. We achieve this by controlled *breadth first search*: the configurations of the multinomial distribution can be thought of vertices V of an underlying graph, $G = (V, E)$. Two configurations $\mathbf{v}, \mathbf{w} \in V$ define an edge $(\mathbf{v}, \mathbf{w}) \in E$ if and only if $\exists_{i \neq j} v_i = w_i + 1$ and $v_j = w_j - 1$. One then starts the algorithm in the vicinity of the mode of current Multi $(p_1, \dots, p_w; n)$: as proxy, we use the point with coordinates equal to the floor of $np_i + 1$. More elaborate set of candidates can be used, see [8]. One then enlists all the neighbours of the initial node and puts them on a *max-priority queue*, see [5]. One then recursively looks at neighbours of the top-priority configuration, checks their probability and enqueues them. In the same time, using a hash-table, one must store information on the visited configurations to avoid multiple visits to the same node. Observe that in case of molecules containing elements with only one isotope, e.g. $\text{C}_c\text{H}_h\text{N}_n$, this step alone would suffice to solve the problem, as showed in Result 1.

Having obtained the critical sets we calculate their exterior product and obtain a set of valid configurations in $LF S_K$. We calculate then their true probability under \mathbb{M}_K and their mass. Finally, we merge all obtained solutions.

We call the above algorithm DEFINE. A prototype of it has been implemented in **R**. Fig. 1 in shows how well the prototype manages in solving Problem 1. Observe also, that the *for loop* can be carried out in parallel.

Algorithm 1 DEFINE

- 1: **require:** $\text{C}_c\text{H}_h\text{O}_o\text{N}_n\text{S}_s$, K, L, M, B
 - 2: Establish $\lambda_{13\text{C}}, \lambda_{2\text{H}}, \lambda_{15\text{N}}, \lambda_{17\text{O}}, \lambda_{33\text{S}}, \mu, \gamma$
 - 3: Find $S = \{(k_1, k_2, k_4) : k_1 + 2k_2 + 4k_4 = K\}$.
 - 4: Find $P = \{\mathbb{P}(\mathbf{k}) : \mathbf{k} \in S\}$
 - 5: Order S using P and select the top $L\%$. Call result S^* .
 - 6: **for all** $\mathbf{k} \in S^*$
 - 7: $\mathfrak{M} :=$ Critical $M\%$ set of Multi $\left(\frac{\lambda_{13\text{C}}}{\mu}, \frac{\lambda_{2\text{H}}}{\mu}, \frac{\lambda_{15\text{N}}}{\mu}, \frac{\lambda_{17\text{O}}}{\mu}, \frac{\lambda_{33\text{S}}}{\mu}; k_1\right)$
 - 8: $\mathfrak{B} :=$ Critical $B\%$ set of Multi $\left(\frac{\lambda_{18\text{O}}}{\eta}, \frac{\lambda_{34\text{S}}}{\eta}; k_2\right)$
 - 9: Partial Result $:= \left\{ \left(\mathbb{M}(\mathbf{x}, \mathbf{y}, z), M(\mathbf{x}, \mathbf{y}, z) \right) : \mathbf{x} \in \mathfrak{M}, \mathbf{y} \in \mathfrak{B}, z = k_4 \right\}$
 - 10: **end for**
 - 11: Result $:= \bigcup$ Partial Results.
-

4 Discussion and Conclusions

The presented algorithm is by far a suboptimal way of handling Problem 1. In fact, more elaborate algorithms could come into being by more careful explo-

ration of the space of configuration of the approximative distribution. However, we judge that all such algorithms could share the idea of using, in one way or another, the approach developed in DEFINE: namely, start by choosing a configuration presumed to be in vicinity of the mode of \mathbb{M}_K , and proceed by a controlled *breadth first search* until either a certain number of configurations is reached, or they already gathered ones already have enough of probability upon them.

A different problem to those mentioned before could be solved in this way to, namely:

Problem 2 Find a small set C among all possible configurations, s.t. $\mathbb{M}(C) \approx 1$.

The candidate for the biggest peak would by then be the product of modes of each multinomial model in (3).

This problem has been more efficiently be solved by different approaches, e.g. by the use of Fourier transform methods, see .

Note also that, at least in case of *Time of Flight* analyzers, there is an additional advantage of studying the LFS_K configurations over those gathered in the above mentioned set C : it is known, that in these instruments the resolution depends on the mass of analyte, see [7]. It is more difficult to differentiate correctly between molecules with similar masses, when both of them are big. Models solving Problem 2 would have to add some sort of binning procedure with bin width being a function of mass, that not being straightforward to model. Thanks to the localisation in the mass to charge domain, while studying LFS_K we simply neglect that sort of problem.

In general, modelling probabilistically the fine structure of the isotopic envelope could serve in an automatic peptide identification procedure. Differences in the fine structure with K^* s.t. $\mathbb{M}(LSF_K^*) = \max_K \mathbb{M}(LSF_K)$ could be particularly informative. However, the design of an appropriate scheme is way beyond the scope of this article.

Note also, that as a possible application of finding a critical A set amounting to, such that $\mathbb{M}(A) \approx 95\%$, one might envisage the problem of finding an optimal binning procedure to match real data resulting from a particular mass spectrometer. In this way, one could measure the machine's resolution without any need to refer to somewhat underdefined notions of *p percent valley* and *peak width*, see [7].

Acknowledgments. We would like to thank Piotr Dittwald for thorough introduction to the problem of fine isotopic structure and many productive brawls over proper definitions. We also thank Dirk Valkenborg for pointing out the existence of the concept of equatransneutronic isotopes. Finally, a huge thanks goes to prof. Alan Rockwood, for the time we spent together discussing the mass spec related issues and also for showing me America.

References

1. Agnarsson, G.: On the sylvester denumerants for general restricted partitions. (154), 49–60 (2002)

Find Rockwood's publication.

2. Bobkov, S., Ledoux, M.: On modified logarithmic sobolev inequalities for bernoulli and poisson measures. *Journal of Functional Analysis* 156(2), 347 – 365 (1998), <http://www.sciencedirect.com/science/article/pii/S0022123697931876>
3. Breen, E.J., Hopwood, F.G., Williams, K.L., Wilkins, M.R.: Automatic poisson peak harvesting for high throughput protein identification. *Electrophoresis* 21(11), 2243–2251 (2000), [http://dx.doi.org/10.1002/1522-2683\(20000601\)21:11<2243::AID-ELPS2243>3.0.CO;2-K](http://dx.doi.org/10.1002/1522-2683(20000601)21:11<2243::AID-ELPS2243>3.0.CO;2-K)
4. Claesen, J., Dittwald, P., Burzykowski, T., Valkenborg, D.: An Efficient Method to Calculate the Aggregated Isotopic Distribution and Exact Center-Masses. *Journal of The American Society for Mass Spectrometry* 23(4), 753–763 (Apr 2012), <http://dx.doi.org/10.1007/s13361-011-0326-2>
5. Cormen, T.H., Stein, C., Rivest, R.L., Leiserson, C.E.: *Introduction to Algorithms*. McGraw-Hill Higher Education, 2nd edn. (2001)
6. Dugundji, J.: *Topology*. Allyn and Bacon series in advanced mathematics, Allyn and Bacon, Boston (1966), <http://opac.inria.fr/record=b1079460>
7. Eidhammer, I., Flikka, K., Martens, L., Mikalsen, S.O.: *Computational Methods for Mass Spectrometry Proteomics*. Wiley-Interscience (2007)
8. Gall, F.L.: Determination of the modes of a multinomial distribution. *Statistics & Probability Letters* 62(4), 325 – 333 (2003), <http://www.sciencedirect.com/science/article/pii/S0167715202004303>
9. Hughey, C.A., Hendrickson, C.L., Rodgers, R.P., Marshall, A.G., Qian, K.: Kendrick mass defect spectrum: a compact visual analysis for ultrahigh-resolution broadband mass spectra. *Anal. Chem.* 73(19), 4676–4681 (Oct 2001)
10. Kallenberg, O.: *Foundations of modern probability*. Probability and its applications, Springer, New York (2002)
11. Kienitz, H.: Mass Spectrometry and its Applications to Organic Chemistry. *Ange wandte Chemie* 73(17-18) (1961)
12. Kingman, J.F.C.: *Poisson processes*, Oxford Studies in Probability, vol. 3. The Clarendon Press Oxford University Press, New York (1993), Oxford Science Publications
13. Olson, M., Yergey, A.: Calculation of the isotope cluster for polypeptides by probability grouping. *Journal of the American Society for Mass Spectrometry* 20(2), 295–302 (2009), <http://dx.doi.org/10.1016/j.jasms.2008.10.007>
14. Rockwood, A.L.: Relationship of Fourier transforms to isotope distribution calculations. *Rapid Commun. Mass Spectrom.* 9(1), 103–105 (jan 1995), <http://dx.doi.org/10.1002/rcm.1290090122>
15. Roos, B.: On the rate of multivariate poisson convergence. *Journal of Multivariate Analysis* 69(1), 120 – 134 (1999), <http://www.sciencedirect.com/science/article/pii/S0047259X98917894>
16. Rosman, K., Taylor, P.: *Isotopic Compositions of the Elements 1997*. J. Phys. Chem. Ref. Data 27(6) (1998)
17. Senko, M.W., Beu, S.C., McLaffertycor, F.W.: Determination of monoisotopic masses and ion populations for large biomolecules from resolved isotopic distributions. *J. Am. Soc. Mass Spectrom.* 6(4), 229–233 (Apr 1995)
18. Valkenborg, D., Assam, P., Thomas, G., Krols, L., Kas, K., Burzykowski, T.: Using a poisson approximation to predict the isotopic distribution of sulphur-containing peptides in a peptide-centric proteomic approach. *Rapid Communications in Mass Spectrometry* 21(20), 3387–3391 (2007), <http://www.scopus.com/inward/record.url?eid=2-s2.0-35349002401&partnerID=40&md5=d68ac3d300587b84f183dcd1c7f508cd>, cited By (since 1996)15

19. Valkenborg, D., Mertens, I., Lemière, F., Witters, E., Burzykowski, T.: The isotopic distribution conundrum. *Mass spectrometry reviews* 31(1), 96–109 (2012), <http://dx.doi.org/10.1002/mas.20339>

Tables

Table 1. Basic Information on Stable Isotopes, as found in [16].

Element	Isotope	Extra Neutrons	Mass [Da]	Probability
Carbon	^{12}C	0	12	0.9893
	^{13}C	1	13.0033	0.0107
Hydrogen	^1H	0	1.0078	0.999885
	^2H	1	2.0141	0.000115
Nitrogen	^{14}N	0	14.0031	0.99632
	^{15}N	1	15.0001	0.00368
Oxygen	^{16}O	0	15.9949	0.99757
	^{17}O	1	16.9991	0.00038
	^{18}O	2	17.9992	0.00205
Sulfur	^{32}S	0	31.9721	0.9493
	^{33}S	1	32.9714	0.0076
	^{34}S	2	33.9679	0.0429
	^{36}S	4	35.9671	0.0002

Figures

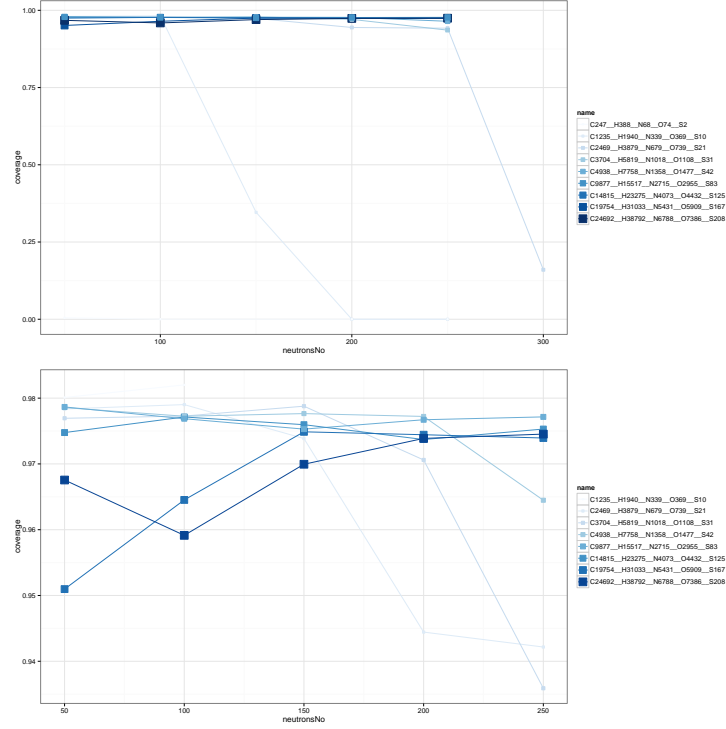


Fig. 1. Coverage obtained using DEFINE algorithm. The image on the bottom zooms into the upper reaches of the top picture. Both show the coverage of distribution original \mathbb{M}_K for $K \in \{50, 100, 150, 200, 250, 300\}$ for several chemical compounds. The bigger the compound (empirical formulas in the legend) the bigger the squares and the more intense the colour. Observe that for lighter compounds the results do not seem promising: we attribute this to the overall quality of conditional distributions \mathbb{M}_K . Simply, all the multinomial distribution in (3) are unimodal and for larger K the solutions to Diophantine equation (5) do not encompass the region next to the mode, where the distribution is centered. For the reasons exposed in **Discussion and Conclusions**, it is impractical to look at these distribution in the first place.

Appendix

Proof of Lemma 1

We want to prove that if $\mu^{[n]} \rightharpoonup \mu$ and $\mu^{[n]}(A), \mu(A) > 0$, then also $\mu_A^{[n]} \rightharpoonup \mu_A$. We do this under the assumption that both $\mu^{[n]}$ and μ are discrete measures on probability space E .

By the *Portmanteau Lemma*, see [10], $\mu^{[n]} \rightharpoonup \mu$ implies that for any set A with boundry ∂A subject to $\mu(\partial A) = 0$, one should observe

$$\lim_{n \rightarrow \infty} \mu^{[n]}(A) = \mu(A). \quad (11)$$

The notion of boundry requires the notion of topology: thus, we decide on the discrete topology, which is natural in this context ⁴. In this topology however, $\partial A = \emptyset$, for it is a set theoretical difference of the closure and the interior, both of which are equal to A . Hence, $\mu(\partial A) = 0$. Thus, (11) always holds.

Ex definitione, $\mu^{[n]} \rightharpoonup \mu$ means, that for any bounded function $f : E \rightarrow \mathbb{R}$ one observes

$$\int f d\mu^{[n]} \xrightarrow{n \rightarrow \infty} \int f d\mu. \quad (12)$$

A simple calculation using both (11) and (12) completes the proof:

$$\int f d\mu_A^{[n]} = \frac{\int f d\mu^{[n]}}{\mu^{[n]}(A)} \xrightarrow{n \rightarrow \infty} \frac{\int f d\mu}{\mu(A)} = \int f d\mu.$$

General form of the *Lucky Law*

If the compound contains elements with their *additional neutron acceptances* in set $I = \{1, 2, 4\}$, formula (10) generalizes to

$$\mathbb{L}(\mathbf{k}) = \frac{\prod_{i \in I} \frac{\mu_i^{k_i}}{k_i!}}{\sum_{\{\mathbf{k}^* : \sum_{i \in I} i k_i^* = K\}} \prod_{i \in I} \frac{\mu_i^{k_i^*}}{k_i^*!}},$$

where \mathbf{k} is an ordered tuple indexed by I . Nature poses a natural limit on the complexity of the *lucky law*, as at most $\#I \leq 10$.

Ascertain that asking Frederik.

⁴ For appropriate topological notions consult [6].