



UvA-DARE (Digital Academic Repository)

Naive Bayes classification model for isotopologue detection in LC-HRMS data

van Herwerden, D.; O'Brien, J.W.; Choi, P.M.; Thomas, K.V.; Schoenmakers, P.; Samanipour, S.

DOI

[10.1016/j.chemolab.2022.104515](https://doi.org/10.1016/j.chemolab.2022.104515)

Publication date

2022

Document Version

Final published version

Published in

Chemometrics and Intelligent Laboratory Systems

License

CC BY

[Link to publication](#)

Citation for published version (APA):

van Herwerden, D., O'Brien, J. W., Choi, P. M., Thomas, K. V., Schoenmakers, P., & Samanipour, S. (2022). Naive Bayes classification model for isotopologue detection in LC-HRMS data. *Chemometrics and Intelligent Laboratory Systems*, 223, Article 104515. <https://doi.org/10.1016/j.chemolab.2022.104515>

General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

Disclaimer/Complaints regulations

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, P.O. Box 19185, 1000 GD Amsterdam, The Netherlands. You will be contacted as soon as possible.



Naive Bayes classification model for isotopologue detection in LC-HRMS data



Denice van Herwerden^{a,*}, Jake W. O'Brien^b, Phil M. Choi^{b,c}, Kevin V. Thomas^b, Peter J. Schoenmakers^a, Saer Samanipour^{a,b,d,**}

^a Van 't Hoff Institute for Molecular Sciences (HIMS), University of Amsterdam, Amsterdam, the Netherlands

^b Queensland Alliance for Environmental Health Sciences (QAEHS), The University of Queensland, Australia

^c Water Unit, Health Protection Branch, Prevention Division, Queensland Department of Health, Australia

^d Norwegian Institute for Water Research (NIVA), Oslo, Norway

ARTICLE INFO

Keywords:

Isotopologue identification
High-resolution mass spectrometry
Naive Bayes classification

ABSTRACT

Isotopologue identification or removal is a necessary step to reduce the number of features that need to be identified in samples analyzed with non-targeted analysis. Currently available approaches rely on either predicted isotopic patterns or an arbitrary mass tolerance, requiring information on the molecular formula or instrumental error, respectively. Therefore, a Naive Bayes isotopologue classification model was developed that does not depend on any thresholds or molecular formula information. This classification model uses the elemental mass defects of six elemental ratios and successfully identified isotopologues for both theoretical isotopic patterns and wastewater influent samples, outperforming one of the most commonly used approaches (i.e., 1.0033 Da mass difference method - CAMERA). For the theoretical isotopologues, the classification model outperformed an "in-house" mass difference method with a true positive rate (TP_r) of 99.0% and false positive rate (FP_r) of 1.8% compared to a TP_r of 16.2% and an FP_r of 0.02%, assuming no error. As for the wastewater influent samples, the classification model, with a TP_r of 99.8% and false detection rate (FD_r) of 0.5%, again performed better than the mass difference method, with a TP_r of 96.3% and FD_r of 4.8%. Therefore, it can be concluded that the classification model can be used for isotopologue identification, requiring no thresholds or information on the molecular formula.

1. Introduction

Non-targeted analysis (NTA) in combination with liquid chromatography high-resolution mass spectrometry (LC-HRMS) is a comprehensive approach for the characterization of unknown chemicals in complex sample matrices, originating from, for example, environmental or biological backgrounds [1–6]. These samples can contain thousands of structurally known and unknown chemicals. To identify these chemicals, the raw data files need to be processed to extract and group information that belongs to unique chemical constituents (i.e., parent, isotopologue, adduct, and (in-source) fragment ions) [1]. During this step, one approach for reducing the number of individual unidentified features is the detection or removal of isotopologues (i.e., heavier versions of the same monoisotopic peak).

For LC-HRMS data, two main approaches have been used to detect

isotopologues [7,8]. The first strategy relies on predicting the molecular formula, which can be translated to a predicted isotopic pattern [7,9]. The main shortcoming of this approach is the difficulties associated with accurate and reliable molecular formula prediction for unknown chemical constituents. The wrong molecular formula could be assigned to a feature either due to instrumental error or the absence of a chemical constituent in a database. These wrongly assigned molecular formulas could lead to identifying the potential isotopologues of a feature with the wrong isotopic pattern, resulting in higher false positive and false negative identification rates.

On the other hand, a theoretical mass difference has been used and implemented in, for example, the open source-software package CAMERA [8,10]. Besides isotopologue detection, CAMERA also performs other filtering steps, such as a retention time comparison, shape correlation, and intensity ratio check for isotopic patterns [8]. The

* Corresponding author.

** Corresponding author. Van 't Hoff Institute for Molecular Sciences (HIMS), University of Amsterdam, Amsterdam, the Netherlands.

E-mail addresses: d.vanherwerden@uva.nl (D. van Herwerden), s.samanipour@uva.nl (S. Samanipour).

<https://doi.org/10.1016/j.chemolab.2022.104515>

Received 22 November 2021; Received in revised form 17 January 2022; Accepted 1 February 2022

Available online 12 February 2022

0169-7439/© 2022 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

implemented mass difference method uses a mass of $n \times 1.0033$ Da to assign isotopologues to their corresponding monoisotopic masses [8,10]. Here n equals the depth of the isotopologue mass. For example, an isotopologue mass depth of four corresponds to the mass range of the monoisotopic peak plus three isotopologues. This approach, even though elegant given that it does not require information on the molecular formula, does require an arbitrary mass tolerance as input. This means that the mass tolerance changes, depending on the resolution during acquisition, and needs to be correctly provided by the user.

In this manuscript, an isotopologue classification model is proposed that requires no prior knowledge of the molecular composition or arbitrary tolerances. The Naive Bayes classification model was generated using elemental mass defects, for which the potential in isotopologue detection was explored. For performance evaluation of the classification model, a comparison was made with an “in-house” developed mass difference method. This comparison was performed for both theoretical isotopic patterns and wastewater influent samples.

2. Experimental section

2.1. LC-HRMS analysis

The 44 wastewater influent and three quality control samples were analyzed with LC-HRMS. Briefly, samples were collected over a time window of 24 h, using on-site autosamplers set to use the optimized conditions described by Ort et al. [11]. These samples were filtered, spiked with 10 ng L^{-1} of 19 labeled internal standards, and stored frozen until analysis. For analysis, $10 \mu\text{L}$ of the sample was injected on a biphenyl column at 45°C and separated using a 10-min gradient from 5 to 100% methanol with 0.1% formic acid. The eluent was analyzed using a QToF with a nominal resolution of 30 000 to 35 000 in positive ion mode with a mass range of 50–600 Da and collision energy of 10 eV. Further details on the analysis are provided elsewhere [12].

2.2. Data processing

The raw data files were converted to mzXML file format, using

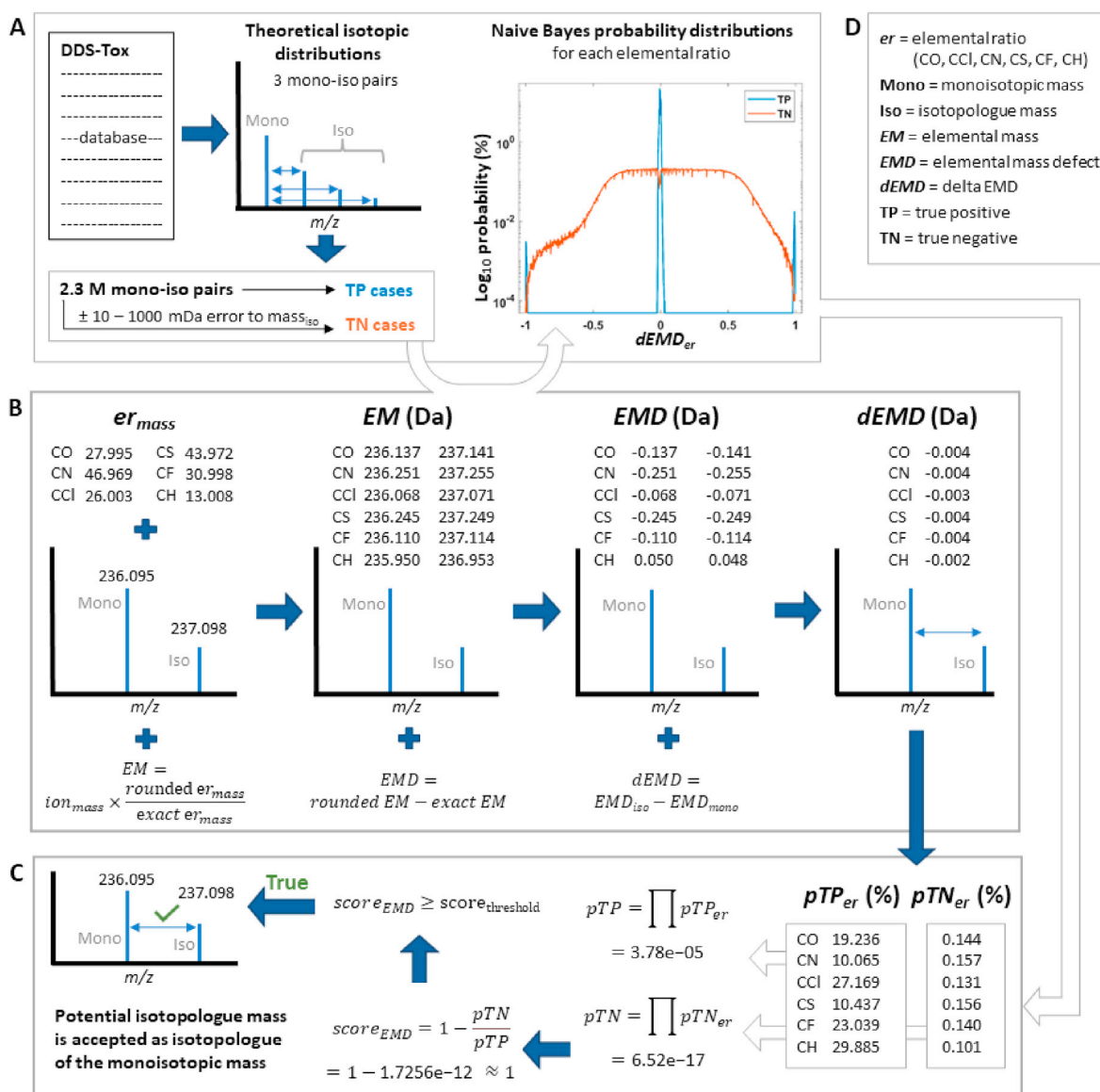


Fig. 1. Section A shows the Workflow for the construction of the Naive Bayes isotopologue classification model, which requires calculations of the $dEMD$ values (section B) for the mono-iso pairs. The workflow for the use of the classification model for the example mono-iso pair in B is shown in section C. Finally, D contains a list of abbreviations.

MSConvertGUI (64-bit, ProteoWizard [13]). Feature lists were generated with the self adjusting feature detection (SAFD) algorithm, using the following settings: 10 000 maximum number of iterations, a minimum intensity of 500, resolution of 20 000, 0.02 m/z minimum window size in the mass domain, 0.75 minimum regression coefficient, a maximum signal increment of 5, a signal to noise ratio of 2, and a minimum and maximum peak width in the time domain of 3 and 200 s, respectively [12]. These feature lists were used for the performance evaluation of the classification model on real samples.

2.3. Theoretical isotopic patterns

The isotopic patterns used for setting up the probabilistic isotopologue classification model were calculated for 737 594 chemicals from the DDS-TOX database [14]. These chemicals consist of a curated list of compounds relevant to environmental and human health. The isotopic patterns were obtained using pyOpenMS [9] (v2.6.0), combining both the isotopic masses from the fine [15,16] and coarse [9] isotope pattern generator (Fig. 1A). The fine isotope pattern generator calculates the hyperfine isotopic pattern that is obtained when the mass defect of the individual isotopes is taken into account [9]. This mass defect equals the difference between the actual mass of an atom and the sum of the building blocks (e.g., neutrons) the atom is comprised of. For this method, a maximum unexplained probability of 0.01% was used. On the other hand, the coarse isotope pattern generator calculates the unit mass isotopic patterns, using the summed probability for each isotopologue peak, ignoring the hyperfine structures. For this, a maximum isotopic tree depth was required that corresponds to one plus the maximum number of isotopes that could be present in a single molecule [16]. Considering the fact that an increasing number of isotopes within a molecule results in a lower occurrence probability (i.e., intensity), a maximum isotopic tree depth of 6 was chosen.

The full isotopic pattern for a compound was comprised of the fine and coarse isotopic patterns, excluding duplicate isotopologues from the coarse isotopic pattern that had a mass difference of ≤ 0.003 Da with any of the other isotopologues, which is the typical mass error observed in LC-HRMS experiments [17]. In this manuscript, a monoisotopic parent ion with one of its isotopologues is referred to as a mono-iso pair. For example, if a monoisotopic parent ion has 5 theoretical isotopologues, 5 mono-iso pairs are obtained. In total, 2 691 244 mono-iso pairs were generated, which were employed for training (85% of the mono-iso pairs) and testing (15% of the mono-iso pairs) of the probabilistic isotopologue classification model (available on figshare) [18].

2.4. Elemental ratio calculations

To construct the probabilistic isotopologue classification model, elemental mass defects (EMDs) were used. The assumption here is that the monoisotopic and isotopologue mass have the same EMD because they have the same molecular structure with the isotopologue having one or more of its atoms being replaced with heavier versions (i.e., isotopes) of the same elements. To calculate the EMD for both the monoisotopic and isotopologue mass, the elemental mass (EM) needs to be calculated according to equation (1). Here, the ion_{mass} can either be the mono-isotopic or the isotopologue mass and the er_{mass} (i.e., elemental ratio mass) depends on the elemental ratio used. For the classification model the elemental ratios CO, CCl, CN, CS, CF, and CH were used, which have an er_{mass} of 27.995, 46.969, 26.003, 43.972, 30.998, and 13.008, respectively. These values are the sum of the elemental masses of each element for a single elemental ratio. For example, the er_{mass} of CO equals the monoisotopic mass of a carbon atom plus that of an oxygen atom (i.e., $12.000 + 15.995 = 27.995$). The selected elemental ratios were chosen based on both the frequency they were encountered in the DDS-Tox database (Table S1) and the fact that only 0.007% of the database entries contain none of the selected elements.

After the EM is calculated, the EMD for the monoisotopic and

isotopologue mass can be obtained according to equation (2) (i.e., EMD_{mono} and EMD_{iso} , respectively). These EMD values are used to calculate the delta EMD ($dEMD$) for a mono-iso pair (Equation (3)). It is important to note that the EMD_{mono} should always be subtracted from the EMD_{iso} and not vice versa when using the probabilistic isotopologue classification model described in this paper. An example case for calculating the $dEMD$ value can be found in Fig. 1B. The full set of isotopologue and monoisotopic EMD values for the DDS-Tox database can be found on figshare [18].

$$EM = ion_{mass} \times \frac{\text{rounded } er_{mass}}{\text{exact } er_{mass}} \quad (1)$$

$$EMD = \text{rounded } EM - \text{exact } EM \quad (2)$$

$$dEMD = EMD_{iso} - EMD_{mono} \quad (3)$$

2.5. EMD probability distributions

To generate the EMD probability distributions for the classification model, both true positive (TP) and true negative (TN) mono-iso pairs were required (Fig. 1A). The mono-iso pairs in the training set were used as the TP cases and TN cases were generated based on the mono-iso pairs from the training set with a randomly added mass error between 0.01 and 1 Da to the isotopologue mass. For all mono-iso pairs in the TP and TN training set, the $dEMD$ s were calculated for the selected elemental ratios (Equation (3)). These $dEMD$ values were used to construct the TP and TN probability distributions for each of the six elemental ratios. To build these probability distributions, the generated $dEMD$ values were binned, using a range between -1 and 1 Da with a 0.002 Da step size. For each $dEMD$ bin, the number of occurrences plus one was used. This prevented that a $dEMD$ range could have a probability equal to zero, in case no occurrences for that specific $dEMD$ were found in the training set. Finally, the probability distributions were calculated by dividing the occurrence distribution values by the total number of occurrences.

2.6. Naive Bayes classification

Naive Bayes classification was used to develop the probabilistic isotopologue detection model, using the TP and TN $dEMD$ probability distributions obtained for the selected elemental ratios (i.e. CO, CCl, CN, CS, CF, and CH). To calculate the posterior probabilities (i.e., $P(A|B)$) for classifying a potential mono-iso pair as TP or TN, Bayes theorem is used (Equation (4)) [19]. Here, $P(A)$ is the probability of a mono-iso pair being TP or TN, $P(B)$ is the occurrence likelihood for a specific $dEMD$ value, $P(B|A)$ is the probability for a $dEMD$ value in case of A, and e equals the number of elemental ratios used, which would be six for our model (i.e., CO, CN, CCl, CS, CF, and CH).

$$P(A|B) = \prod_{i=1}^e \frac{P(B|A)_i \times P(A)_i}{P(B)_i} \quad (4)$$

Since $P(B)$ is a marginal probability (i.e., constant probability normalizing factor), equation (4) can be rewritten to equation (5). Additionally, a uniform distribution is assumed for the prior $P(A)$, further reducing the formula to equation (6).

$$P(A|B) \propto \prod_{i=1}^e P(B|A)_i \times P(A)_i \quad (5)$$

$$P(A|B) \propto \prod_{i=1}^e P(B|A)_i \quad (6)$$

Lastly, for the classification of the potential mono-iso pair, the TP and TN probabilities are obtained using equation (6). These probabilities were converted to probability percentages (i.e., on a scale of 0–100). Due to the wide range of values that can be obtained for the TP and TN probabilities, a score_{EMD} is used instead for the evaluation (Equation (7)).

Here $P(TP)$ and $P(TN)$ equal the true positive and true negative probabilities, respectively. This $score_{EMD}$ ranges between 1 and minus infinity. In case the potential mono-iso pair has a $score_{EMD}$ above a set threshold, the potential isotopologue is classified as a correct isotopologue of the monoisotopic ion. An example for the calculation of the $score_{EMD}$ can be found in Fig. 1C.

$$score_{EMD} = 1 - \frac{P(TN)}{P(TP)} \quad (7)$$

2.7. Performance assessment

For the performance assessment, the test set was used. In this instance, TN cases were also generated based on the mono-iso pairs from the test set with a random mass error added to the isotopologue mass of 0.01–1 Da. For both the TP and TN cases, the $score_{EMD}$ were calculated (Equation (7)). To select a suitable $score_{EMD}$ cut-off value and assess the performance of the classification model, the TP and false positive (FP) rates were calculated for a range of $score_{EMD}$ values. The $score_{EMD}$ from 0.7–1 Da with a step size of 0.002 Da were employed to calculate the TP_r and FP_r (Equations (8) and (9), respectively). Here, the TP s equal the number of cases from the test set that were correctly classified as an isotopologue, FN s are the number of cases that were incorrectly classified as not an isotopologue, TN s are cases that were correctly classified as not an isotopologue, and FP s are the cases that were wrongly classified as isotopologues.

$$TP_r = \frac{TP}{TP + FN} * 100 \quad (8)$$

$$FP_r = \frac{FP}{TN + FP} * 100 \quad (9)$$

2.8. Mass difference method

The mass difference method is a commonly used approach for automated isotopologue detection in LC-HRMS data. This method has already been implemented in different open access algorithms such as CAMERA and MZmine [8,10]. Since these packages often perform more steps than just isotope detection and require specifically formatted input data, an “in-house” mass difference strategy was developed [8]. This “in-house” method is a julia implementation of the isotopologue detection implemented in CAMERA, which was used to benchmark our classification model against. For the mass difference method, to assess if a signal is an isotopologue of a monoisotopic peak, first the mass difference between the signal and monoisotopic ion was calculated. Then, the residue of the division of the mass difference by 1.0033 Da is obtained. For example, if the mass difference is 2.0081 Da, the residue would be 0.0015 Da. In case the residue is lower than the specified mass tolerance, the signal is accepted as an isotopologue of the monoisotopic mass. For the mass difference method, when dealing with the training set, a mass tolerance of ± 0.0001 Da was used based on the assumption that the theoretical isotopologues do not contain any mass error. On the other hand for the wastewater samples, this mass tolerance was increased to ± 0.01 Da to better reflect the inherent mass error in such data caused by background signal and instrumental fluctuations.

2.9. Isotopologue detection performance for wastewater samples

To test the isotopologue classification model on real samples, the isotopologue detection performance was evaluated for the feature lists obtained from 44 wastewater influent samples and three quality control samples. Additionally, a reference compound list comprised of 45 chemicals was used, containing the monoisotopic masses (i.e., protonated molecular mass), retention times, and parent isotopologue distributions (Table S3). The isotopologue distributions for these chemicals

were obtained from the isotope pattern preview tool in MZmine2 (v2.53), using the protonated molecular formula, a minimum intensity of 0.01%, a merge width of 0.0001 Da, and a charge of 1, which showed to cover an isotopologue mass depth of six [10].

The presence of a reference compound was confirmed based on the reference retention time ± 0.1 min and the monoisotopic parent mass with a mass tolerance of 0.01 Da. When a reference compound monoisotopic parent mass was present, all features within a time range of ± 0.1 min were extracted. If a feature's mass was higher than the monoisotopic mass and lower than the monoisotopic mass plus 1.0033×6 (i.e., isotopologue mass depth of six), it was evaluated as a potential isotopologue with both the classification model and the mass difference method. When a model correctly identifies an isotopologue according to the reference parent isotopologue distribution, it is considered a TP case. Whereas the FP cases are incorrectly identified isotopologues and the FN cases are the TP cases that were not detected by a model. With these cases, the TP_r and FD_r were calculated for the classification model and mass difference method (Equations (8) and (10), respectively), which were used to compare the two isotopologue identification methods.

$$FD_r = \frac{FP}{TP + FP} * 100 \quad (10)$$

2.10. Calculations and code availability

All calculations were performed using a personal computer running Windows 10 Education with 12 cores and 32 GB of memory. For obtaining the theoretical isotopologues of the DDS-Tox database Python (v3.9.4) was used and for calculations related to the classification model Julia (v1.6.0) was used. The mzXML files were imported in julia using the MS_Import package, which is available at https://bitbucket.org/SSamani_pour/ms_import.jl/src/master/. The code for the probabilistic isotopologue classification model is available at https://bitbucket.org/Denice_van_Herwerden/emdforiso/src/master/. This package includes both the probabilistic isotopologue classification model and functions to use the model with feature lists obtained either from SAFD [12] or other algorithms. The code for SAFD is available at https://bitbucket.org/SSamani_pour/safd.jl/src/master/.

3. Results and discussion

3.1. Exploring the EMD probability distributions

Calculating the EMD values for the theoretical isotopologues showed that the EMD values for the monoisotopic and isotopologue masses were similar. Fig. 2 shows the EMD values for the theoretical isotopic distribution of carbamazepine. In this example, a minimum and maximum absolute difference in EMD (i.e., $dEMD$) of 0.003 and 0.020 Da were found, respectively. Additionally, an increase in $dEMD$ between the EMD_{mono} and EMD_{iso} was observed for isotopologues with a higher isotopologue mass depth. Even though the elements S and F are not present in the molecular formula of carbamazepine, a similar EMD trend is observed as for the elements O, N, Cl, and H. On the other hand, Figs. S1 and S2 show that the presence of other elements (e.g., Br and P) in the molecular formula also do not influence the EMD values.

Overall, similar trends were observed for all theoretical isotopologue distributions with EMD values ranging from -0.5 to 0.5 Da for all six elemental ratios. To evaluate this trend, the Pearson correlation coefficients between the EMD_{mono} and EMD_{iso} values were obtained [20]. These coefficients were calculated separately for each elemental ratio and isotopologue mass depth of 1 till 6 (Table S2). The highest correlation of 1.00 was found for the elemental ratio CN with an isotopologue mass depth of 1 and the lowest value was 0.86 for both the elemental ratios CCl and CS with an isotopologue mass depth of 5 (Figs. S3 and S4, respectively). Overall, the Pearson correlation coefficient decreases with a higher isotopologue mass depth except for an isotopologue mass depth

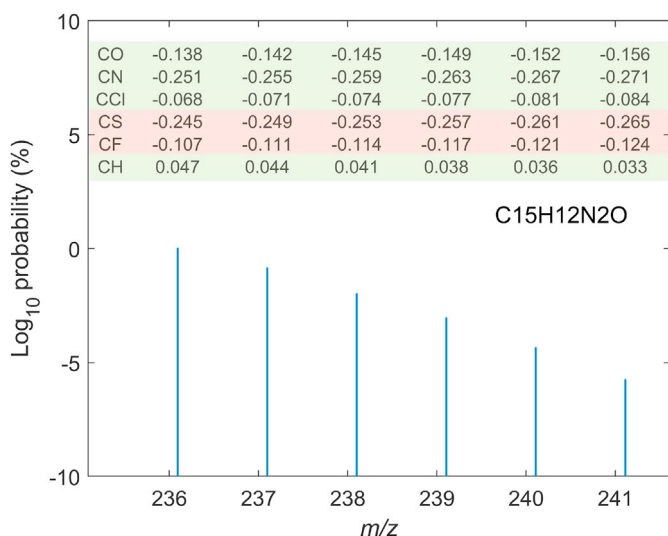


Fig. 2. Isotopic distribution of carbamazepine with the corresponding \log_{10} probability percentages. For the monoisotopic (236.095 Da) and each isotopologue peak (237.098, 238.102, 239.105, 240.108, and 241.112 Da), the EMD values are shown above in Da for the elemental ratios CO, CN, CCl, CS, CF, and CH. Additionally, the elemental ratios that are present in the molecule are marked in green and the ones that are not are marked in red. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

of 6. It is expected that this was due to a relatively low number of mono-iso pairs with a depth of 6 (Table S2). These results showed that similar EMD values for mono-iso pairs were obtained throughout the theoretical dataset.

After calculating all $dEMD$ values for the mono-iso pairs of both the TP and TN cases, the TP and TN probability percentage distributions were obtained for the selected elemental ratios (Fig. 3). For the TP probability distributions, there were 2 regions for which the TP probabilities were

higher than the TN probabilities. The first region being around a $dEMD$ of 0, which is in accordance with the hypothesis that the monoisotopic and isotopologue mass of the same compound obtain similar EMD values. As for the second region, $dEMD$ values close to 1 and -1 Da were found. For the TN probability distributions, a small decrease in probability was observed around a $dEMD$ of 0 Da, which was caused by the minimum added mass error to the isotopologue mass of the TN mono-iso pairs (i.e., 0.01 Da). In case the minimum added mass error would be set to a higher or lower value than 0.010 Da, the decrease in TN probability around a $dEMD$ value of 0 would become broader or narrower, respectively. For example, if a higher value than 0.010 Da is used for the minimum randomly added mass error, the broader decrease in TN probabilities will result in more $dEMD$ values having a higher TP probability than TN probability, potentially leading to more FP isotopologue identifications. Overall, these plots showed that the $dEMD$ could be used to differentiate between isotopologue and non-isotopologue masses.

3.2. Classification model performance

A receiver operator curve was generated for selection of the $score_{EMD}$ threshold. This curve showed the TP_r versus the TN rate for scores EMD between 0.7 and 1 (Fig. S5). Based on this plot a $score_{EMD}$ threshold of 0.9997 was selected. This corresponded with a TP_r and FP_r of 99.0 and 1.8%, respectively.

3.2.1. Comparison with existing method

To evaluate the performance of the classification model with that of the existing mass difference method, the performance for the “in-house” mass difference method was evaluated for a mass tolerance of 0.0001 Da. The mass tolerance was selected based on the assumption that there is no error present in the theoretical mono-iso pairs and the full receiver operator curve can be found in section S4. For a mass tolerance of 0.0001 Da, a TP_r and FP_r of 16.2 and 0.02% was found, respectively. The relation between the isotopologue probability and TP detected isotopologues for both the mass difference method and classification model, showed that for the lowest probability isotopologues (i.e., 78% of the total mono-iso

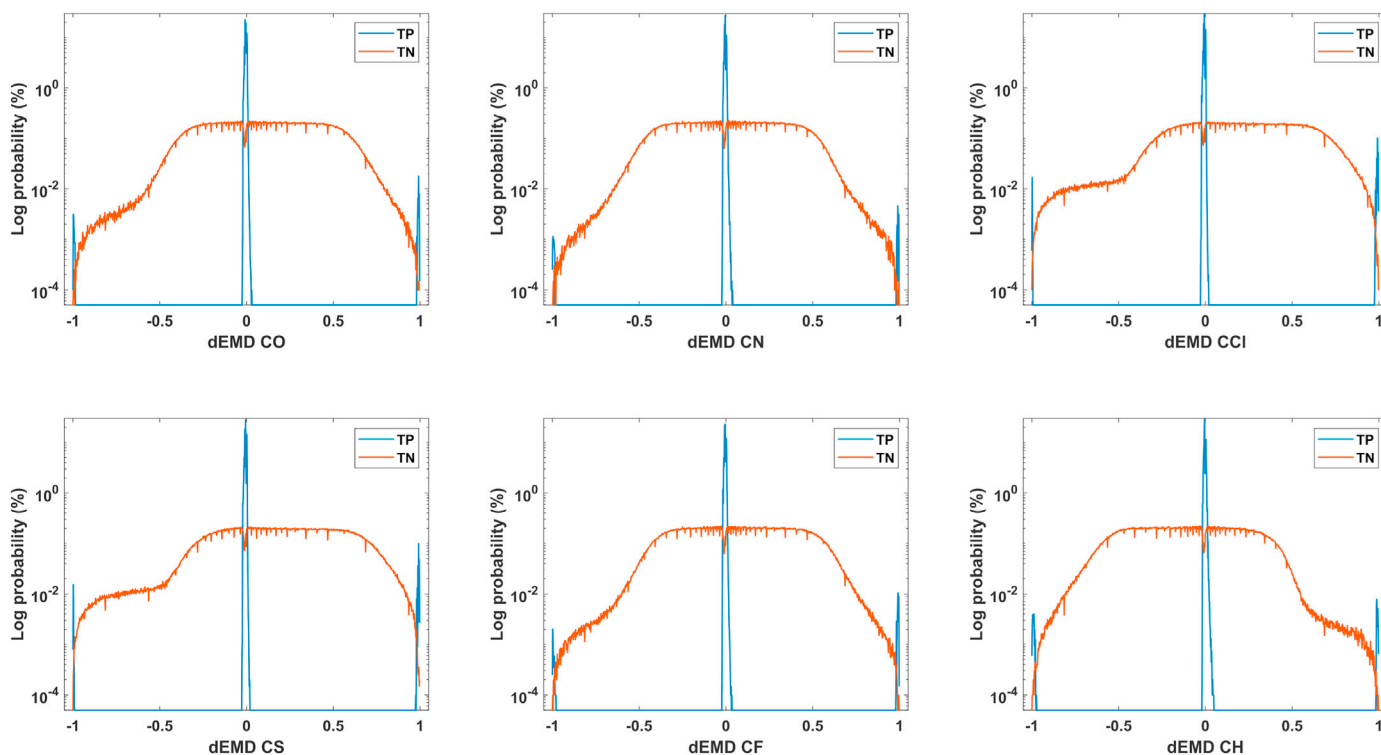


Fig. 3. TP and TN probability distributions for the $dEMD$ values for the selected elemental ratios CN, CCl, CO, CS, CF, and CH.

pairs) only 5.2% was detected for the mass difference method and 98.8% for the classification model (Table S4). Fig. S7 shows an example isotopic distribution for 3-thiomethylparacetamol, where only two out of seven isotopologues actually follow the 1.0033 Da mass trend. Finally, when the mass difference results are compared to the results of the classification model (i.e., TP_r of 99.0% and FP_r of 1.8%), both methods performed well with regard to the FP_r (i.e., $\leq 5\%$). However, the classification model outperformed the mass difference method for the TP_r .

3.3. Model implementation for real samples

To evaluate the model performance for real samples, isotopologue detection was performed for 44 wastewater influent and three quality control samples. A total of 391 features were evaluated as potential isotopologues from the 45 reference compounds in question. Overall, 212 TP cases, one FN case, and one FP case were found for the classification model, resulting in an average TP_r of 99.8% and an FD_r of 0.5%. The FN case was caused by an 0.011 Da mass error between the mono-isotopic and isotopologue mass, which is larger than the minimum mass error (i.e., 0.010 Da) assumed for the true negative cases that are used for training the model. As for the FP case, the detected isotopologue mass was 155.068 m/z and the monoisotopic parent ion mass was 152.072 m/z . If the decreasing intensity for less likely isotopologues would have been taken into account, this ion would not have been included due to the absence of the isotopologues with a higher probability (e.g., 153.068 and 154.075 m/z , Fig. S8). However, when dealing with non-targeted data, the molecular composition is often not known, meaning that a consistent assumption of decreasing isotopologue intensity cannot be made. For example, halogenated compounds do not follow this trend. Overall, it can be concluded that the classification model can also be applied to real data.

For the mass difference method, a total of 203 TP , 10 FN , and 13 FP cases were found, corresponding to an average TP_r and FD_r of 96.3 and 4.8%. For these cases, all FN s were caused by a mass error larger than 0.01 Da and all FP s were caused by the same reason as the FP of the classification model. Across multiple datasets a signal at 304.182 m/z was identified as an isotopologue of codeine, for which the monoisotopic mass was 300.159 m/z . Only in some cases, an isotopologue at 301.163 m/z was detected, which would still mean that there were no isotopologues with an isotopologue mass depth of 2 or 3 present with higher intensities than the signal at 304.182 m/z . To conclude, the classification model had a higher TP_r and lower FD_r than the mass difference method. However, if the decreasing intensity with lower isotopologue probabilities would have been taken into account, the methods would both have had an FD_r of 0.0%.

4. Potentials and limitations

The classification model provides a good alternative approach for the detection of isotopologues, requiring no information on the molecular formula or arbitrary thresholds. However, it should be noted that the classification model is unable to distinguish between isotopologues coming from different chemicals or signals with the same monoisotopic mass. This would require additional separation, such as chromatography, although, the classification model would still not be able to distinguish between isotopologues from overlapping isobaric compounds. Besides reducing the total number of features for identification, correct isotopologue identification can also assist in accurate molecular formula assignment. When multiple formulas are possible for a monoisotopic mass, the isotopic patterns can be predicted and compared with the detected isotopologues masses to eliminate less likely candidates. Lastly, the model was built based on isotopic distributions with a tree depth of six, meaning that it might not be able to correctly classify ions with more than 6 isotopologues if these ions would be detected at all due to their low occurrence probabilities. However, if required, the EMDforIso.jl package enables the user to retrain the classification model using

different training sets and parameters.

5. Conclusion

This manuscript demonstrated the potential of using elemental ratios for the detection of isotopologues. The classification model that was constructed based on the elemental ratios CO, CN, CCl, CS, CF, and CH, showed overall good performance and even outperformed the commonly used mass difference method. For the theoretical mono-iso pairs, when assuming no error, the classification model outperformed the mass difference method with a TP_r of 99.0% and FP_r of 1.8% compared to a TP_r of 16.2% and an FP_r of 0.02%. As for the wastewater influent samples, the classification model, with a TP_r of 99.8% and FD_r of 0.5%, performed better than the mass difference method, with a TP_r of 96.3% and FD_r of 4.8%.

Author statement

All the authors have read the manuscript "Naïve Bayes classification model for isotopologue detection in LC-HRMS data" and approved its content prior to submission.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgement

The authors thank the wastewater treatment plants that assisted in the collection of the wastewater influent samples, the Chemometrics and Advances Separation Team (CAST) for their insights, and the authors gratefully acknowledge the financial support from the Australian Research Council ARC Discovery Project (DP190102476) and Linkage Project (LP150100364). Additionally, Jake W. O'Brien is the recipient of an NHMRC Emerging Leadership Fellowship (EL1 2009209). Finally, the Queensland Alliance for Environmental Health Sciences, The University of Queensland, gratefully acknowledges the financial support of the Queensland Department of Health.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.chemolab.2022.104515>.

References

- [1] S. Bastian, J. Youngjoon, K. Sarit, H.L. Amy, D. Pradeep, O. Jake, G.R. Maria Jose, G.G. Sara, M.F. Jochen, T.V. Kevin, S. Saer, Trac. Trends Anal. Chem. 133 (2020) 116063.
- [2] E.L. Schymanski, et al., Anal. Bioanal. Chem. 407 (2015) 6237–6255.
- [3] E. Werner, J.-F. Heilier, C. Ducruix, E. Ezan, C. Junot, J.-C. Tabet, J. Chromatogr. B 871 (2008) 143–163 (Hyphenated Techniques for Global Metabolite Profiling).
- [4] S. Samanipour, M.J. Reid, K. Bæk, K.V. Thomas, Environ. Sci. Technol. 52 (2018) 4694–4701.
- [5] S. Samanipour, S. Kaserzon, S. Vijayasathya, H. Jiang, P. Choi, M.J. Reid, J.F. Mueller, K.V. Thomas, Talanta 195 (2019) 426–432.
- [6] W. Brack, J. Hollender, M.L. de Alda, C. Müller, T. Schulze, E. Schymanski, J. Slobodnik, M. Krauss, Environ. Sci. Eur. 31 (2019) 62.
- [7] M.A. Stravs, J. Hollender, Anal. Chem. 86 (2017) 6812–6817.
- [8] C. Kuhl, R. Tautenhahn, C. Bötcher, T.R. Larson, S. Neumann, Anal. Chem. 84 (2012) 283–289.
- [9] H.L. Röst, et al., Nat. Methods 13 (2016) 741–748.
- [10] T. Pluskal, S. Castillo, A. Villar-Briones, M. Orešič, BMC Bioinf. 11 (2010).
- [11] C. Ort, M.G. Lawrence, J. Rieckermann, A. Joss, Environ. Sci. Technol. 44 (2010) 6024–6035.
- [12] S. Samanipour, J.W. O'Brien, M.J. Reid, K.V. Thomas, Anal. Chem. 91 (2019) 10800–10807.
- [13] M.C. Chambers, et al., Nat. Biotechnol. 30 (2012) 918–920.

- [14] C.M. Grulke, A.J. Williams, I. Thillanadarajah, A.M. Richard, *Comput. Toxicol.* 12 (2019) 100096.
- [15] M.K. Lacki, M. Startek, D. Valkenburg, A. Gambin, *Anal. Chem.* 89 (2017) 3272–3277.
- [16] M. Loos, C. Gerber, F. Corona, J. Hollender, H. Singer, *Anal. Chem.* 87 (2015) 5738–5744.
- [17] N.A. Alygizakis, S. Samanipour, J. Hollender, M. Ibáñez, S. Kaserzon, V. Kokkali, J.A. Van Leerdam, J.F. Mueller, M. Pijnappels, M.J. Reid, E.L. Schymanski, J. Slobodnik, N.S. Thomaidis, K.V. Thomas, *Environ. Sci. Technol.* 52 (2018) 5135–5144.
- [18] D. van Herwerden, S. Samanipour, Dataset for: probabilistic classification model for isotope detection in high-resolution mass spectrometry. https://figshare.com/articles/dataset/DDS-Tox_isotope_distributions/16559517/2, 2021.
- [19] C. Albon, in: R. Roumeliotis, J. Bleiel (Eds.), *Machine Learning with Python Cookbook*, O'Reilly Media, Inc., 2018.
- [20] R.G. Brereton, *Applied Chemometrics for Scientists*, John Wiley & Sons, 2007.