

1 **Training and evaluating machine learning algorithms for ocean** 2 **microplastics classification through vibrational spectroscopy**

4 **Abstract**

5 Microplastics are contaminants of emerging concern - not only environmental, but also to
6 human health. Characterizing them is of fundamental importance to evaluate their potential
7 impacts and target specific actions aiming to reduce potential harming effects. This study
8 extends the exploration of machine learning classification algorithms applied to FTIR spectra
9 of microplastics collected at sea. A comparison of successful classification models was
10 made in order to evaluate prediction performance for 13 classes of polymers. A rigorous
11 methodology was applied using a pipeline scheme to avoid bias in the training and selection
12 phases. The application of an oversampling technique also contributed by compensating
13 unbalanceness in the dataset. The log-loss was used as the minimization function target and
14 to assess performance. In our analysis, Support Vector Machine Classifier provides a good
15 relationship between simplicity and performance, for a fast and useful automatic
16 characterization of microplastics.

18 **Keywords**

19 Microplastics; Marine Pollution; Chemical Identification; Vibrational Spectroscopy; FTIR;
20 Machine Learning

23 **1. Introduction**

24 Plastic debris is today found in virtually every habitat on earth (Free et al., 2014; Lebreton et
25 al., 2018; Saito et al., 2018; Wang et al., 2019; Zhang et al., 2020). Yet, major scientific
26 concern has been given to plastics at sea, where they are almost omnipresent and when
27 found in large quantities and/or concentrations their negative impacts may be serious
28 (Barnes, 2002; Chae and An, 2017). Originated mainly on land as a result of large
29 consumption of disposable items and poor waste management, plastics compose most
30 marine debris (Galgani et al., 2015).

31 According to the International Organization for Standardization (ISO), microplastics are any
32 solid plastic particle insoluble in water with any dimension from 1 μm to 1 000 μm (1 mm) and
33 large microplastics between 1 mm and 5 mm. (ISO/TR 21960:2020). They can originate from
34 the degradation and fragmentation of larger plastic debris when exposed to environmental
35 conditions or can be directly emitted in their form as, for example, microbeads, pellets, or
36 textile fiber (Fendall and Sewell, 2009; GESAMP, 2019; Hidalgo-Ruz et al., 2012; Thompson
37 et al., 2009). Their wide occurrence and physical characteristics, such as density and
38 chemical composition, have contributed to consider these as emerging environmental
39 contaminants (Sauvé and Desrosiers, 2014).

40 An adequate characterization of these contaminants can provide substantial information on
41 the inputs and transport in the oceans, rates of degradation and fragmentation, interaction
42 with biota, consequences of their presence in natural habitats, and, finally, risk assessment
43 and management. Multiple parameters are relevant for these analyses, such as size, mass,
44 sampling site, and DNA from the plastisphere (Zettler et al., 2013). Determination of the
45 chemical composition of sampled plastics has been used to suppose where they originated
46 from, what human use they might have had, and even estimate their age (Song et al., 2015;
47 Turner and Holmes, 2011).

48 While diversity is wanted in the production and use phases of polymer products, it adds
49 complexity to the identification process. For macroplastics, characterization may still be done

50 by recycling codes in the product or by physical characterization. Yet, for microplastics, these
51 procedures have been shown to result in high rates of false positives when not followed by
52 any chemical analysis (GESAMP, 2019; Hidalgo-Ruz et al., 2012) Thus, a chemical
53 characterization of microplastics is fundamental to more adequately assess the sources, fate,
54 and impacts of microplastics.

55 Vibrational Spectroscopy, namely Raman and Fourier Transform Infrared Spectroscopy
56 (FTIR), are established techniques to assess chemical composition (Köppler et al., 2016).
57 Both analyses provide information on specific chemical bonds and functional groups, albeit by
58 different methods (Kuptsov and Zhizhin, 1998). They allow differentiation between synthetic
59 and natural polymers, identification of polymer type, and degree of weathering, hence their
60 wide utilization in microplastic identification methods (Andrady, 2017).

61 FTIR analysis can be done visually by an expert, but it is often more convenient to use a peak
62 matching dedicated software with the available databases (Li et al., 2006). They perform what
63 is called a *library search*, going through all spectra in the database and showing those with
64 better ranking. This strategy allows a faster polymer identification, but since sampled
65 microplastics are somehow dirty and weathered, peaks absent in the same non-weathered
66 polymers can be seen in microplastics. (Jung et al., 2018; Xu et al., 2019). Therefore, their
67 chemical signature (spectrum) can be significantly different from virgin polymers that
68 commonly constitute databases (GESAMP, 2019; Primpke et al., 2018). Procedural factors
69 can also influence the quality of measurements and impact identification accuracy and time.
70 The size of identifiable microplastics, for example, is limited by the FTIR setup and should not
71 be less than 500 μm (unless the spectrometer is coupled with a microscope) (Köppler et al.,
72 2016; Shim et al., 2017). Weathered microplastics could be fragile and shatter while being
73 manipulated, which also hinders spectrum acquisition (Shim et al., 2017).

74 A few alternatives have been proposed and tested to surpass these limitations. Pre-
75 processing of the raw data, the usage of specifically designed databases, and dual database

76 searches are some of the approaches that have been demonstrated to improve match results
77 and allow library searches to be used with good confidence (Primpke et al., 2018; Primpke et
78 al.2020; Renner et al., 2019).

79 Most notably, μ -FTIR and μ -Raman techniques are being used to address common
80 limitations on microplastics analysis. Coupling a spectrometer and a microscope, allows the
81 characterization of smaller particles and with an additional Focal Plane Array (FPA) detector,
82 several particles can be characterized at once (Cabernard et al. 2018; Brandt et al. 2020). In
83 this sense, these techniques can provide a more accurate picture of the evaluated
84 microplastic environment. Yet, this kind of equipment may be less available mainly due to
85 cost issues.

86 Improving methodologies for conventional vibrational spectroscopy could provide better
87 alternatives to marine pollution researchers in places where financial support is insufficient,
88 mainly the global south. Having said that, the proposed methodology could be used to
89 evaluate the performance of classification algorithms using data from different spectroscopic
90 techniques, especially FTIR imaging where data volumes are significantly larger ($\sim 10^6$
91 spectra per image) and spectra quality tends to be lower (Primpke et al. 2017).

92 Sampling microplastics at sea produces hundreds or even thousands of individual particles,
93 which by manual methodologies would take far too long to be characterized. Increased
94 automation of the characterization process of microplastic samples with higher confidence in
95 the results could quicken information acquisition on this emerging contaminant, filling multiple
96 knowledge gaps and significantly advancing understanding on the field. Indeed, that is the
97 main objective of the study presented in this paper. A machine learning pipeline was
98 proposed for the selection of the best among a few machine learning algorithms to classify
99 microplastics spectra, then discuss the main findings.

100 ML algorithms have recently been proposed as faster and more accurate methods to analyze
101 spectra from marine microplastics and have been successfully used in chemometrics (and

102 many other fields) for more than a decade (Conroy et al. 2005). Most efforts have been
103 applied to imaging techniques, probably due to high computational demands associated with
104 higher data volumes - thus, the necessity to improve efficiency (Hufnagl et al., 2019).
105 Nevertheless, conventional spectroscopic techniques could benefit from these improvements.

106 In machine learning, classification algorithms are a category of prediction models used to
107 attribute unidentified samples to a given class based on a set of variables. These algorithms
108 have been applied by Hufnagl (Random Forest Classifier) and Kedzierski (K-nearest
109 neighbors) to spectral data of microplastics and resulted in expressive classification
110 accuracies (Hufnagl et al., 2019; Kedzierski et al., 2019). However, a comparison of different
111 algorithms hasn't been made. Given the already demonstrated potential of these techniques,
112 further investigation could lead to even better results. Since the learning process highly
113 depends on the database available for training, comparable results should ideally derive from
114 the same dataset. In the present study, data previously published by Kiedzierski were used
115 for this purpose (Kedzierski et al., 2019).

116 Methodologies focusing on the selection of relevant attributes on microplastics spectra have
117 been proposed by Renner 2017, Hufnagl 2019, and da Silva 2020. Indeed, it is the approach
118 experts take when visually interpreting a spectrum. In this paper, a *dimensionality reduction*
119 (*DR*) technique called Principal Component Analysis (PCA) was applied to extract the most
120 relevant features in the whole dataset. Some machine learning models lose performance
121 when the data is represented in a high dimension space, falling into the so-called curse of
122 dimensionality (Trunk, 1979). Feature selection or extraction is used to reduce the number of
123 variables that describe a certain set of instances (samples) while retaining most of the
124 information. This can, sometimes, reduce the time of implementation without significant loss
125 of information and, when the curse of dimensionality is observed, even improve prediction
126 performance.

127 The main disadvantage of the machine learning approach to this classification task is that it is
128 limited to analyzing classes of polymers that are well represented in the training database.
129 Since there is an enormous variety of plastics, setting up rather complete databases would be
130 laborious. However, from all the known range of polymers, there are only a few commonly
131 found as marine microplastics, making the setup of representative databases a lot more
132 feasible. The method described in this paper goes further in the process of setting up this
133 type of database by using an artificial oversampling technique to mitigate the effects of
134 imbalanced datasets in the learning process, which can produce models that are biased
135 towards the majority class.

136 **2. Materials and Methods**

137 This study extends previous efforts by combining several machine learning techniques and
138 comparing algorithms, to find, among them, the best method to automatically classify
139 microplastic spectra. It differs from other approaches by proposing a pipeline methodology to
140 train, evaluate and select classification models. The methodology described next was
141 rigorously defined to justify choosing one model over all others. Once the model is chosen, it
142 can be trained and called to make a prediction. A researcher interested in reproducing the
143 experiment must follow the entire methodology, whereas a potential user of the proposed
144 predictive tool can skip to section 2.9 for a short explanation on how it can be used.

145 All programming was made in a Core i5-7200U with 16 GB ram, using only 1 core of the
146 processor to avoid issues related to parallel programming. We used Python programming
147 language and some of its data analysis and machine learning libraries, namely: scikit-learn
148 (Pedregosa et al., 2011), scipy (Virtanen et al., 2020), and imblearn (Lemaître et al., 2017).
149 All programming is available in a supplementary file (Appendix D) so that the methodology
150 described in this paper can be audited, reproduced and even improved by peers. It is also
151 available online at GitHub (https://github.com/EdsonCilos/mp_classification).

152 **2.1. The dataset**

153 The data used to train the algorithms were previously published by Kedzierski et al., 2019,
154 and were generated from the Attenuated Total Reflection FTIR spectroscopy of samples
155 collected during expeditions in the Mediterranean Sea. Spectra were recorded in
156 absorbance mode in the range of 4000 to 600 cm^{-1} with 4 cm^{-1} resolution and 16 scans. The
157 labeled dataset was constructed using the raw data consisting of 958¹ spectra previously
158 identified and assigned to 17 different classes.

159 **2.2. Dataset pre-processing pipeline**

160 Typically in the spectra evaluation by an expert, the analysis of some frequencies is not
161 considered due to the presence of noise or just because the knowledge in the field
162 prescribes that samples can be distinguished by some specific peaks. Our perspective,
163 however, stands in the Machine Learning background, assigning to the algorithm the task to
164 find the most successful pattern matching. As pointed by Vapnik: “In a wide philosophical
165 sense, predictive models do not necessarily connect prediction of an event with an
166 understanding of the law that governs the event; they are just looking for a function that
167 explains the data best.” (Vapnik, 2006)

168 In other words, our goal is to look for a good explanation of the available data, and therefore,
169 no feature selection based on expert knowledge was done, nor seemed to be necessary. As
170 a consequence, all recorded wavelengths were used for the analysis. Additionally, a major
171 drawback of a priori feature selection is that it may be suboptimal and should be revised
172 when new classes are introduced, while PCA is embedded in the pipeline.

173 Since the original dataset was highly imbalanced, the first step in pre-processing consisted
174 of renaming the classes that were underrepresented in the dataset, assigning them to a
175 generic class called “unknown”. Samples from the 4 least represented classes, namely,
176 Polyurethane, Animal fiber-like, Poly(vinyl chloride), and acrylic (PMMA) were moved to the
177 unknown class, effectively removing 4 classes. The latter having only 3 samples, while the
178 others only 1 sample, making it impossible to train a model using the methodology adopted

179 in this paper. The dataset then consisted of 958 samples from 14 different classes (13
180 microplastic classes + unknown).

181 It is important to note that apart from relabeling underrepresented classes and correcting the
182 baseline, all preprocessing steps must be done after the train/test split, to avoid the so-called
183 “snooping bias/data leakage” (see “Results and Discussion” for further details). It is a
184 common mistake when coding machine learning algorithms. It happens when some
185 information that typically would be available only in subsequent steps (like in a “production
186 environment”), is introduced in the learning process.

187 The remaining pre-processing steps were assembled in a pipeline, which we will call, for
188 short, a Pre-pipeline (Later on in Section 2.6 we shall discuss the full Pipeline of the
189 methodology, which includes the Pre-pipeline). All steps were tested in the “turn on/ turn off”
190 configuration, to evaluate which of those was able to improve the final model’s performance.
191 It consisted of the application of the following techniques in the raw data (in this order):

192

- 193 1. Baseline correction using Asymmetric Least Squares (Eilers and Boelens, 2005).
194 Parameters were: lambda (2nd derivative constraint) = $1e^5$; p (weighting deviations) =
195 0.05; itermax (number of iterations to perform) = 10
- 196 2. Standard scaler (Z-transform), since the features are arranged at different scales,
197 which could affect the model performance (Géron, 2017);
- 198 3. Principal component analysis (PCA), as a method to verify if higher dimensions could
199 impact the algorithm’s performance (Jolliffe, 2002). PCA is an unsupervised learning
200 technique that creates a new linear space with orthogonal variables, called principal
201 components (PCs), which are the directions of most variance in the dataset. If the
202 original features are uncorrelated, the new space contains the same number of
203 dimensions as the original space, but, generally, a much smaller number of variables

204 is sufficient to describe the data without much loss of information. In our analysis, the
205 first n PCs containing 99% of the explained variance were used.

206 4. Oversampling to mitigate the effects of imbalance in the dataset. Despite having
207 assembled underrepresented classes, the dataset was still highly imbalanced, as the
208 2 most populated classes contained 45% of all samples in the dataset. This can
209 cause distortions in the learning process, as the algorithm would be more likely to
210 assign a new sample to these classes (Prati et al., 2009). This was done using the
211 Imbalanced Learning API inside a Pipeline scheme (Figures 1 and 3). It generates
212 new instances by randomly sampling the under-represented classes, effectively
213 copying existing spectra (Lemaître et al., 2017).

214

215 The four pre-processing techniques described (Baseline correction, Standard Scaler, PCA,
216 and Oversampling) were subjected to analysis to evaluate if including them would improve
217 the performances of the classifiers. Several combinations were tried for each classifier (e.g.
218 kNN without PCA, with Standard Scaler, and with oversampling). The best combination was
219 forwarded to another trial to examine if a frequency smoothing strategy could improve the
220 performance likely by suppressing features that arise from non-ideal instrument and sample
221 conditions (Renner et al., 2019; Zimmermann and Kohler, 2013). Ergo, a Savitzky-Golay
222 filter was included before the standard scaler step, in which an evaluation of 115
223 combinations of windows and degrees was performed. This step was done separately in
224 order to speed the evaluation, considering it already took 4 days.

225 **2.3. The holdout strategy**

226 Initially, the dataset was split into training and test sets (holdout strategy), stratifying
227 proportionally to the number of classes, to make sure every class is present both in the train
228 and test sets. The test set ratio choice was 25%. To allow reproducibility, we fixed the seed,
229 as well as all random states, equal to 0. It is worth mentioning that the authors used the

230 holdout strategy to be able to evaluate, under mathematical assumptions, an “unbiased”²
 231 estimate of the generalization error. The holdout strategy allows the trained algorithm to be
 232 tested on an “external” dataset, completely unknown to the algorithm. This can be done
 233 using mathematical tools like Hoeffding's inequality (Hoeffding, 1963), Vapnik-Chervonenkis
 234 theory (Vapnik, 2000) or usual statistical inference techniques. Thereby, a unique and final
 235 evaluation in the test set is capable of providing a confidence interval for the proposed
 236 algorithm, as we shall see in the Results/Discussion section.

237 **2.4. The classification algorithms**

238 According to the "No Free Lunch Theorem" (Wolpert, 1996): there is no a priori reason to
 239 prefer one machine learning model over another without making any assumption about the
 240 data. Therefore, in practice, we must select some models and test them on the problem that
 241 we aim to solve. For this reason, we chose some of the most popular and relevant models in
 242 machine learning: Decision Tree (DT) (Moore, 1987; Murphy, 2012, p. 544), Gaussian Naive
 243 Bayes (GNB) (Murphy, 2012, p. 82), k-Nearest Neighbor (kNN) (Nordhausen, 2009, p. 14),
 244 Random Forest (RF) (Nordhausen, 2009, p. 587), Logistic Regression (LR) (Murphy, 2012,
 245 p. 14; Nordhausen, 2009, p. 119), Support Vector Machine Classifier (SVC) (Murphy, 2012,
 246 p. 496; Nordhausen, 2009, p. 417), and Neural Networks (Goodfellow et. al, 2016).

247 Each model has general characteristics that may be more suited to one application than the
 248 other. A comparison of 5 of the models proposed here is presented in Table 1, adapted from
 249 Kotsiantis (Kotsiantis, 2007). For instance, GNB has a high speed of learning and
 250 classification but is highly affected by interdependent attributes. SVC, on the other hand, has
 251 a high speed of classification and is very tolerant to irrelevant attributes, but is slow to learn
 252 and falls short in explainability (Kotsiantis, 2007).

	Decision Trees	Neural Networks	Naive Bayes	kNN	SVC
Accuracy in general	**	***	*	**	****

Speed of learning	***	*	****	****	*
Speed of classification	****	****	****	*	****
Tolerance to missing values	***	*	****	*	**
Tolerance to irrelevant attributes	***	*	**	**	****
Tolerance to redundant attributes	**	**	*	**	***
Tolerance to highly interdependent attributes	**	***	*	*	***
Dealing with discrete/binary/continuous attributes	****	***	***	***	**
Tolerance to noise	**	**	***	*	**
Dealing with danger of overfitting	**	*	***	***	**
Attempts for incremental learning	**	***	****	****	**
Explanation ability	****	*	****	**	*
Model parameter handling	***	*	****	***	*

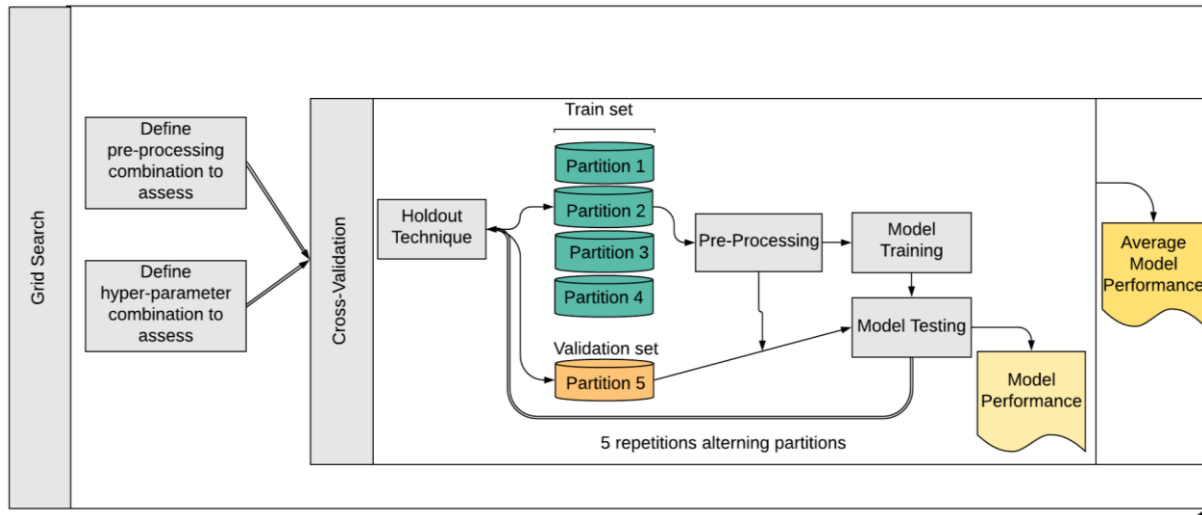
253 Table 1 - A comparison of models' characteristics (Adapted from Kotsiantis, 2007)

254 The algorithms have been configured to return probabilities (provided by Sklearn;
255 *sklearn.calibration*) rather than just a deterministic output. This seems to be suitable for the
256 current task, allowing the researcher to be more confident when the probability score is
257 higher. In this context, the cross-entropy (or log-loss) of a multinoulli distribution was chosen
258 as an objective function to be minimized, to make the probability distribution of the model as
259 close as possible to the empirical distribution (see Appendix B for further details)
260 (Goodfellow et al., 2016).

261 **2.5. Gridsearch**

262 Every machine learning model has a set of parameters that must be predefined by the user
263 and are not learned during training. These are called "hyperparameters". Grid Search is an
264 efficient tool to find the best combination of hyperparameters, given a predefined grid.

265 All combinations of hyperparameters can be found in the project's file (Appendix D) (refer to
 266 param_grid.py). Just to mention a few, Logistic Regression, for example, was tested using
 267 diverse penalty schemes (l1, l2 and Elastic Net with different l1 ratios ranging from 0.1,
 268 0.2,..., 0.9), different regularization parameter C ($C = 10^j$ for $j = -2, -1, \dots, 3$) and several
 269 solvers depending on the penalty (newton-cg, sag, saga, lbfgs and/or liblinear).



270

271 *Figure 1 - Graphical representation of the pipeline used with Grid Search for assessing*
 272 *multiple classifiers.*

273 In a cross-validation we split the data in k folds (or partitions), for each fold we train the
 274 model in the remaining data (all data except the fold) and evaluate the model in the fold; in
 275 the end the k results are averaged producing a score. Such a strategy was used to evaluate
 276 the grid search, in a 5-fold cross-validation scheme (Figure 1). We fixed the seed that
 277 generated the partition in the cross-validation, therefore a reproduction of this experiment is
 278 likely to achieve the same result.

279 Throughout this section whenever we refer to a “training set” we are referring to a new
 280 training set created by cross-validation or another validation mechanism. Otherwise, when
 281 referring to the initial training set, we will refer to it as the "original training set".

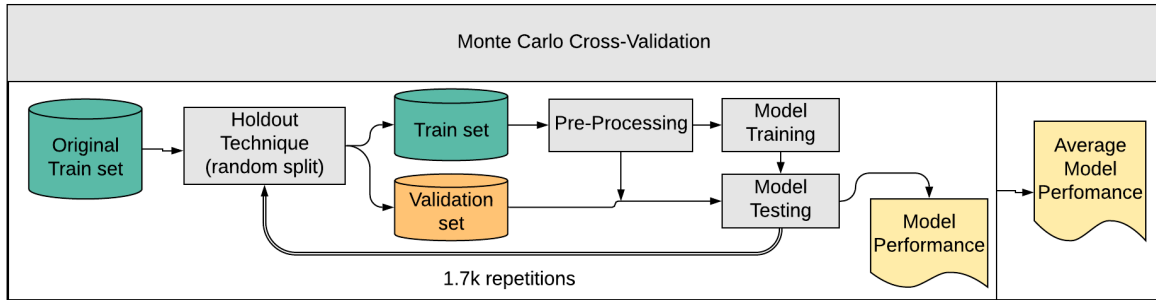
282 Cross-validation was performed within Grid Search to avoid overfitting on the validation set -
 283 because these hyperparameters can be tuned optimally for a specific validation set. That is,

284 for a certain combination of hyperparameters, 5 different training and test sets are used to
285 train and test the model's performance. The final result for that specific combination is the
286 average of these tests (standard deviations were also computed).

287 **2.6. Monte Carlo Cross-Validation**

288 After the GridSearch, we looked for the best scores. Decision Tree e Gaussian Naive Bayes
289 classifiers performed very poorly. Among the remaining models, we selected the most
290 promising configuration in each case: the top 2 Neural Networks (with less or more neurons),
291 top 2 Support Vector Machine (linear kernel and 'rbf' kernel), the best Random Forest, the
292 best kNN and the best Logistic Regression. These models were subjected to subsequent
293 analysis considering a Monte Carlo Cross-Validation (MCCV) (briefly described next). The
294 remaining classifiers showed a considerably higher log-loss and/or its parameters were
295 unlikely to improve the performance; therefore, they were not forwarded to the next step.
296 The selected model will be referred to as "the final hypothesis".

297 An MCCV consists in randomly splitting the data in train and validation sets, several times.
298 It's a kind of holdout strategy with many trials, also known as "repeated learning-test". This
299 technique allows us to draw a "monte carlo picture" and check a model's performance in a
300 histogram, allowing us to compute some statistics like mean and standard variation. In our
301 case, the MCCV consisted of 1700 trials. Since the grid search was done with a fixed seed,
302 it can produce models that perform well in rare cases. By introducing randomness through
303 MCCV, we can make a second filter in the models, in an attempt to avoid "rare best
304 performing models".

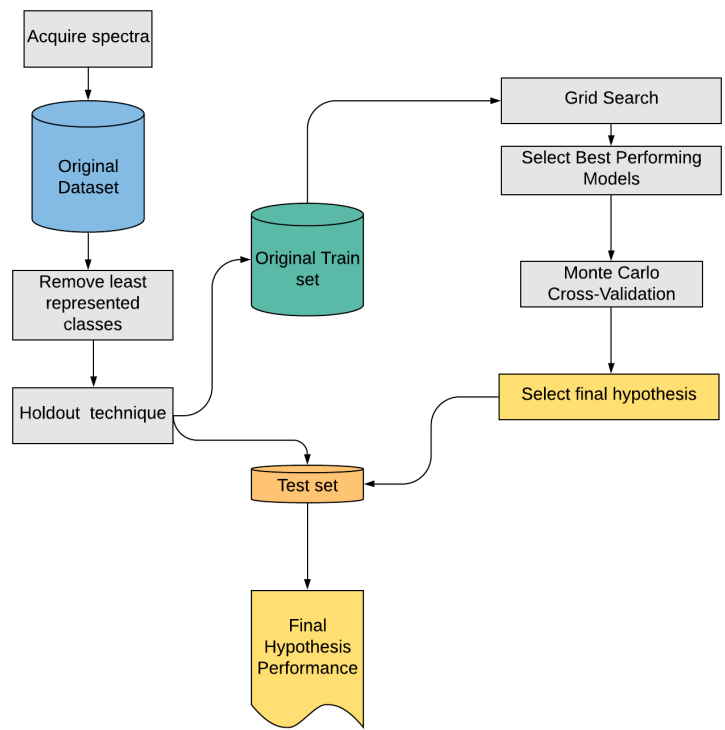


305

306 *Figure 2 - Graphical representation of the pipeline used in the Monte Carlo Cross-Validation.*

307 *After 1,700 repetitions the results are averaged. The process was repeated on the seven*
 308 *most promising models.*

309 **2.7. The Pipeline**



310

311 *Figure 3 - Graphical representation of the full pipeline methodology.*

312 The application of the tools described until now was done by means of a pipeline
 313 methodology. In total, 878 models under 5 combinations of pre-processing techniques were
 314 subjected to the same workflow. Thus, using a pipeline was useful to automate the
 315 sequence of operations that we proposed (Figure 3).

316 One important step to mention is the pre-processing of the data, which happens after the
317 train/test split, something that is commonly overlooked, but that avoids the snooping bias.
318 Notice, as well, that all steps concerning the selection of the best model are done using the
319 “original training set”, not the “original dataset”. Keeping part of the original dataset aside, we
320 can use it in the end to test the selected model on a dataset completely new to the algorithm.
321 We suggest this pipeline to be used to evaluate and compare the performances of different
322 models in the task of microplastics classification with ATR-FTIR spectroscopy. It can be
323 equally used with different datasets, containing, for example, other polymer classes.

324 **2.8. Comparability and External Validation**

325 In order to be able to compare our methodology with previous works, we applied, to the final
326 hypothesis, the training methodology described by Kedzierski et al., 2019 (which consisted
327 essentially of a MCCV with 1000 slightly different splits) with three differences: Firstly, we
328 used a stratified split. This is important because our dataset is highly unbalanced with some
329 classes underrepresented, therefore it is possible that in certain splits some classes remain
330 only in the training set (more probable) or only in the test set (less probable), which can, in
331 any case, introduce bias in our model. Secondly, Kedzierski et al. assign all unclear
332 predictions (less than 3 votes in kNN) to the “unknown” class, which the present
333 methodology does not. Finally, and most importantly, it seems that the authors included a
334 standardization pre-process in the entire dataset, which introduces a “snooping” bias in the
335 model. This means, for example, that the model “already knows” the mean and standard
336 deviation of the features of the validation set before even being trained. On the other hand,
337 our pre-processing methodology is built inside a pipeline that uses, for example, the
338 standardization learned in the training set and applies it in the validation set, which avoids
339 such bias.

340 A further verification was included to prove that the approach proposed here is usable
341 outside the Kedzierski dataset. To do so, a different environmental dataset was used to train

342 and evaluate the final hypothesis, which was done by MCCV, as described in section 2.6.
343 There was no tuning of hyperparameters at this step. The dataset contained 800 FTIR
344 spectra of polymers ingested by turtles and was published in a previous study by Jung et al.
345 (2018). There were initially 9 classes, including a differentiation between high and low
346 density poly(ethylene). As previously, the less representative classes were suppressed and
347 spectra moved to the “unknown” class, which resulted in 5 classes.

348 **2.9. Application of the classifier**

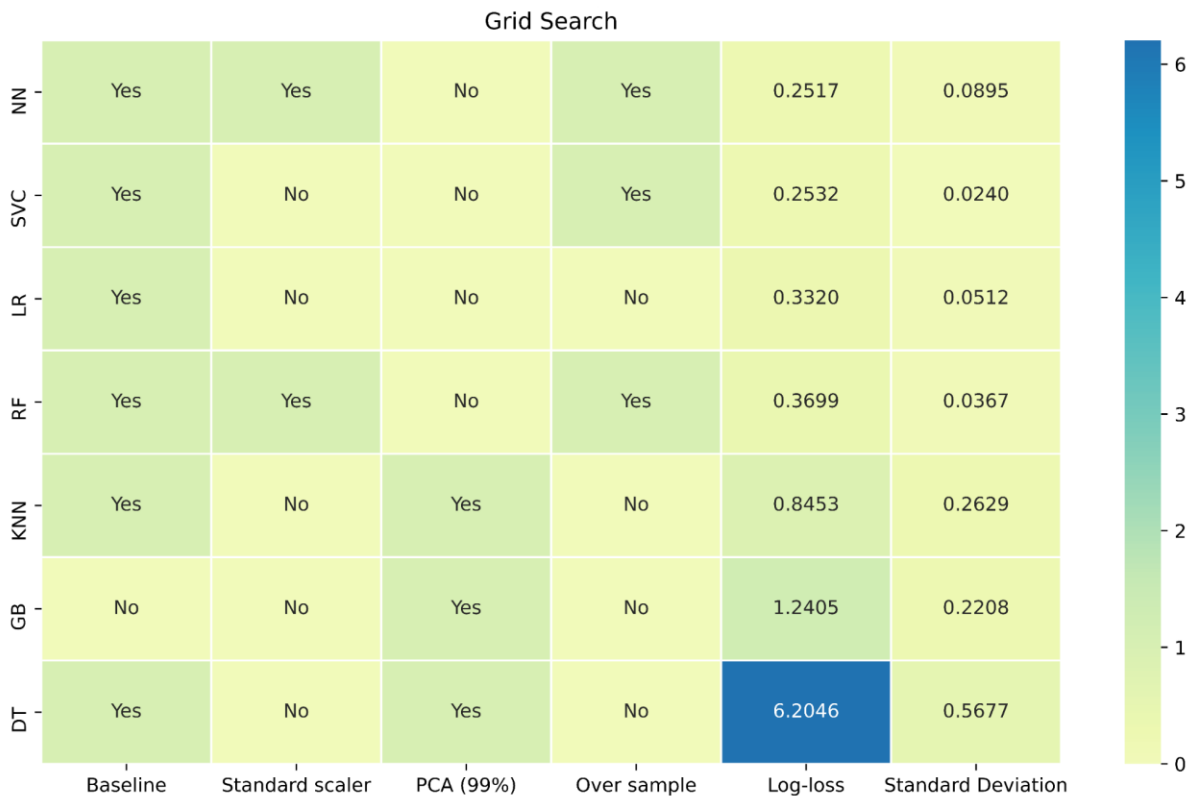
349 If the reader is particularly interested in using the proposed final hypothesis (classifier), they
350 can do so by running the code “*final_classifier.py*” in a python console, which is located in
351 the root folder of the project repository.

352 It should be stressed that the data provided for classification should not have been pre-
353 processed at all, as the appropriate pre-processing will be performed by the algorithm
354 considering the results presented in the next section of this paper. Additionally, the spectral
355 data should have been collected at the same frequency range described in section 2.1.

356 Given the limitation of the classifier to recognize only the classes it was trained on, the
357 classification result will be given as a probability of that sample belonging to the assigned
358 class. If the user is then unsure of the assignment, it can be visually confirmed. An example
359 of this application and output is given in Appendix C.

360 **3. Results/Discussion**

361 Using the methodology described in the previous section, 28096 models were evaluated in
362 the 5-fold cross-validation. Results (log-loss and standard deviation) for the best combination
363 of hyperparameters and pre-processing techniques for each different model are presented in
364 Figure 4.



365

366

367 *Figure 4 - Best performance found within the gridsearch for each evaluated algorithm,*
 368 *namely Random Forests (RF), Support Vector Machine Classifier (SVC), Logistic*
 369 *Regression (LR), K-Nearest Neighbours (kNN), Decision Trees (DT), Gaussian Naive Bayes*
 370 *(GB), considering the best combination of hyperparameters and pre-processing techniques*
 371 *- the baseline correction, Standard scaler, Principal Component analysis (PCA) and (Naive)*
 372 *oversampling. The log-loss shown here is the mean of the five results (cross-validation). The*
 373 *correspondent standard deviation is also presented.*

374

375 For each simulation, a more detailed table can be found in Appendix D (paths can be found
 376 in Appendix A). Analyzing the GridSearch, the seven most promising models, based on the
 377 proposed metrics, were: LR with baseline correction; Sigmoid NN with baseline correction,
 378 standard scaler, and oversampling; Sigmoid NN with baseline correction and standard
 379 scaler; SVC with linear kernel, baseline correction, and oversampling; SVC with rbf kernel,
 380 baseline correction, and oversampling; RF with baseline correction, standard scaler, and

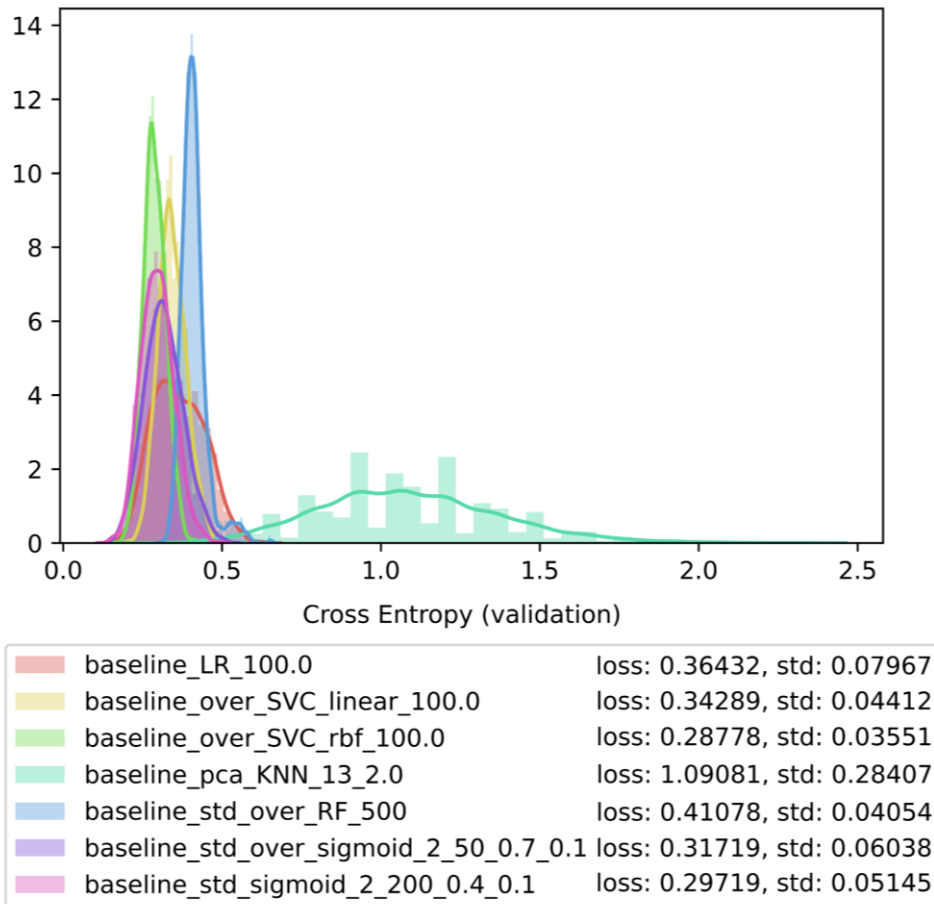
381 oversampling; kNN with baseline correction and PCA. The reader can refer to the
382 “*gridsearch*” subfolder in the project repository for further details.

383 The pre-processing combination that improved the performance was different for each
384 model. For instance, SVC did not benefit from PCA, while kNN, DT, and GB did.

385 When comparing several models with and without PCA, most of them present little scoring
386 difference when PCA is applied. Meanwhile, we verified significant differences in the learning
387 time. For example, baseline + oversample took around 5 hour to run, whereas baseline +
388 pca + oversample took about 17 min. It should be noticed, however, that this is not a
389 detailed picture of which algorithm could benefited from PCA, but an experimental evidence
390 to affirm that future analyses, mainly when they involve the comparison of several
391 algorithms, can be done with PCA (99%) without significant loss of information but with a
392 saving time bonus. All grid search and MCCV (next) simulations have their time computed
393 and stored in the file “*time.csv*” in the “*results*” project’s folder.

394 Previous results suggested the use of peak smoothing techniques over the spectra (da Silva
395 et al. 2020; Zimmermann et al. 2013). In our case, however, a Savitzky-Golay filter did not
396 improve the performance of any model. For this reason, the procedure was not included in
397 the project. Such a phenomenon probably happens because smoothing peaks favours
398 experts’ visual inspection, while in a pattern recognition this is likely to be irrelevant.

399 The next step was to subject the most promising models to an MCCV in order to validate the
400 consistency of their performances when randomness is allowed. The resulting histogram of
401 each model is presented in Figure 5.



402

403 *Figure 5 - Cross-entropy log-loss for the seven most promising models: Logistic Regression*
 404 *with baseline correction; SVC with linear kernel, baseline correction, and oversampling; SVC*
 405 *with rbf kernel, baseline correction, and oversampling; kNN with baseline correction and*
 406 *PCA; RF with baseline correction, standard scaler, and oversampling; Sigmoid NN with*
 407 *baseline correction, standard scaler, and oversampling; Sigmoid NN with baseline correction*
 408 *and standard scaler*

409 Overall, the seven models had a decrease in performance (increased log-loss) but it could
 410 be said that they are, in fact, robust against randomness. Even though the results indicate
 411 that kNN and RF are the worst models, only by looking at these results it is not possible to
 412 choose the best model. To extend the analysis, the authors examined their accuracies and
 413 the results were very similar between models, so no further conclusion could be made.
 414 Considering the methodology proposed by (Renner et al., 2019), which was also used by
 415 (Kedzierski et al. 2019), we evaluated, for each model, how many classes were “working

416 well” (sensitivity greater than 75%) and equal conclusions were possible. Therefore, since
417 none of these metrics provided clear boundaries for a final decision, we shall take into
418 account Occam's Razor, which states that “the simplest model that fits the data is also the
419 most plausible” (Abu-Mostafa et al., 2012, p. 167). Appropriately, SVC with a linear kernel
420 was chosen as the final hypothesis. The reader can refer to the “*results/graphics*” subfolder
421 in the project repository for further graphs in this matter.

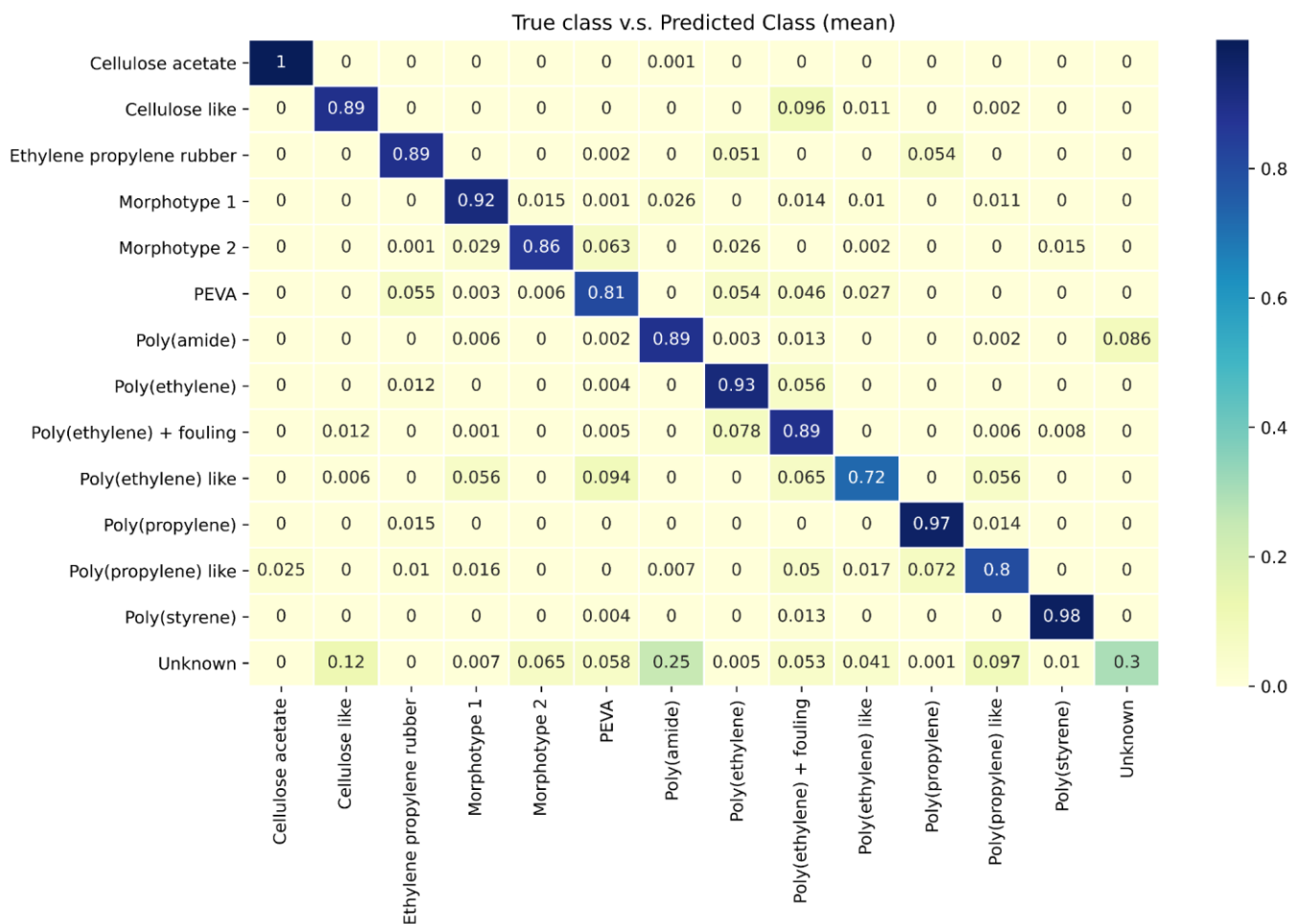
422 More complex, non-linear models, such as neural networks, can be tuned to fit the data well.
423 Nonetheless, if the data were collected under different circumstances, or other pre-
424 processing techniques or parameters were changed, this would cause minor changes in the
425 final performance of every model. The neural networks, for instance, would probably still
426 perform well, but there would be many more hyper-parameters to tune again in order to
427 obtain optimal results. Linear SVC has only one such parameter (C) that could be readily
428 tuned in case new reference samples or a whole new training set was provided. Having
429 considered this, SVC is expected to perform efficiently with low tuning and prediction times.

430 Additionally, regarding the similarity in the performance of distinct models under various
431 criteria, it is the authors’ opinion that the “ideal theoretical structure” of the data is indeed
432 linearly separable. By “ideal theoretical structure” of the data, we mean: no error in data
433 labeling, materials without wear, and all micro-plastic polymers well represented, thus we
434 conjecture that under FTIR the microplastics exhibits an almost linear pattern, with linearity
435 suffering small violations due to the presence of these “non-ideal elements”. Another
436 important remark is that our model improved performance when combined with oversample
437 technique, agreeing with the fact that in an imbalanced problem, a SVC classifier can
438 produce suboptimal models that are biased toward the majority class (He, H., Ma, Y., 2013,
439 p. 83). Lastly, Linear SVC can be implemented using the LIBLINEAR library, which is
440 capable of handling large data sets (Fan et al., 2008). Therefore, we expect that the model
441 could be implemented in a massive dataset, in which we expect similar results³.

442 Since the methodology describes supervised learning procedures, it can be expected to
443 observe where the algorithm makes mistakes when tested on a validation set. This is

444 possible by looking at the confusion matrix shown in Figure 6, generated by averaging
 445 results obtained for every split in the MCCV for the selected model.

446



448 *Figure 6 - Confusion Matrix for the results of MCCV for the final hypothesis. Each row*
 449 *represents the true class, meanwhile each column is the class predicted proportion, on*
 450 *average, by baseline_over_SVC_linear_100.0. Ideally, the matrix should be an identity*
 451 *matrix (perfect match).*

452 A confusion matrix compares the true class in the validation set with the class predicted by
 453 the algorithm for the same sample (Ballabio et al., 2018). Considering this, it can be
 454 observed, for instance, that samples from the “Cellulose Acetate” class were predicted to
 455 belong to their true class every time, whereas for the “Cellulose Like” class, 9.6% of all

456 predictions were incorrectly assigned to the “Poly(ethylene) + fouling” class and 1.1% were
457 incorrectly assigned to the “Poly(ethylene) like” class. The same kind of error was observed
458 by Kedzierski et al. (2019) using a kNN classifier. Since the algorithm is trained to distinguish
459 between classes by identifying patterns in the provided dataset, which only contains
460 information on spectral features and polymer classes, this could indicate similarities between
461 spectra of both classes. As pertinently pointed out by the referred authors, spectral bands
462 associated with aging and biofouling of microplastics could be the cause, with a
463 preponderance of the latter. A visual examination of spectra from the “Poly(ethylene) +
464 fouling” and “Cellulose Like” classes was done to check for specific bands related to these
465 changes in their molecular fingerprint. However, it was not possible to visually establish any
466 specific relations that could be causing the algorithm to mistake both classes, as spectra had
467 too many overlapping peaks. We emphasize that the results are promising, nonetheless.

468 It may be clearer to depict faults made by the algorithm by observing the “Ethylene
469 Propylene Rubber” (EPR) class, where the algorithm wrongly predicted the “Poly(ethylene)”
470 (PE) and “Poly(propylene)” (PP) classes. This is comprehensible since EPR is made of the
471 same monomers as PE and PP and share their chemical characteristics, hence, their
472 spectral features. For classes such as “Poly(ethylene) like”, “Poly(ethylene) + fouling” and
473 “Poly(propylene) like”, the main error made by the algorithm was to assign samples truly
474 belonging to these classes to their regular counterparts (i.e. “Poly(ethylene)” and
475 “Poly(propylene)”.

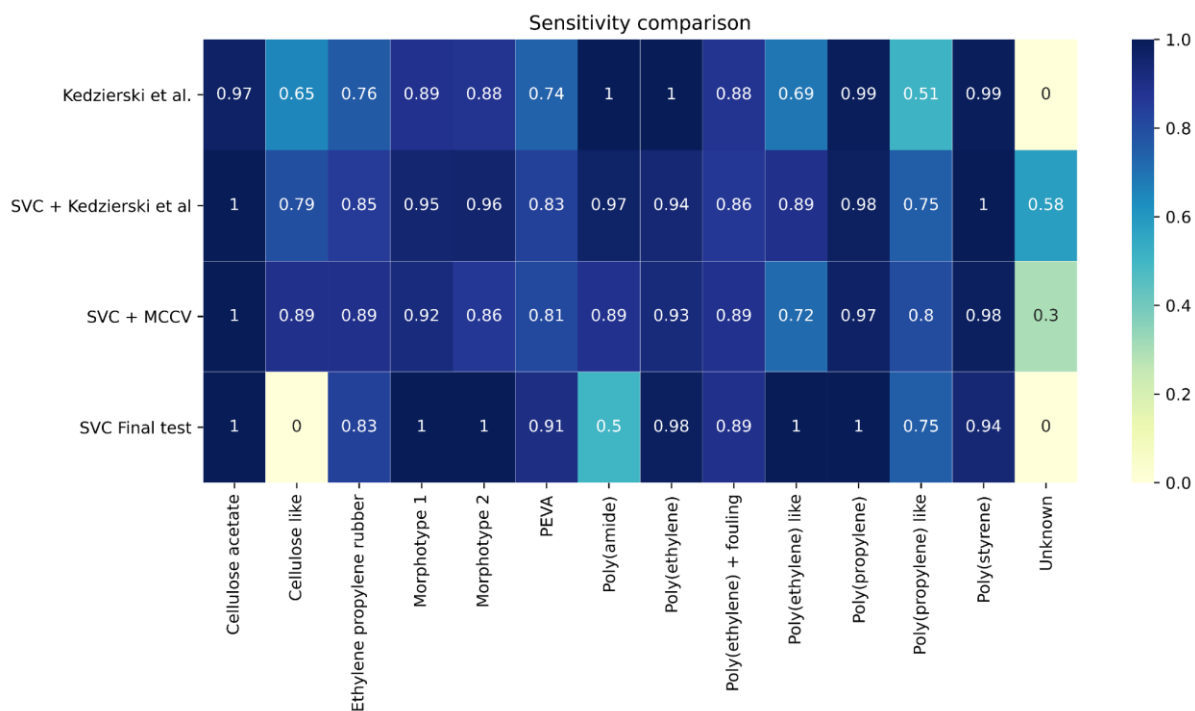
476 In the case of “Poly(propylene) like”, which had one of the worst results, most incorrect
477 predictions (7,2%) were made to the “Poly(propylene) class. Since they were originally
478 assigned to these classes by an expert due to their spectral similarities, this could be
479 expected - and, in fact, was also observed by Kedzierski et al (2019). However, despite
480 having similar spectra, the algorithm was able to differentiate them with considerable
481 accuracy, showing better class results than the k-nearest neighbour (k-NN) method originally
482 proposed by the referred authors. Further research on the quality of the data can respond

483 more assertively if errors of the algorithm are indeed mistakes or indicative of a human error
 484 in labelling the data. The class “Morphotype 2” was mistaken mostly with “Poly(amide)”,
 485 however, it returned excellent results meaning this grouping is probably very concise and
 486 should contain samples from remarkably similar polymers, or even, mostly samples from a
 487 single polymer type. It is also unlikely that this polymer is one of the others already in this
 488 database.

489 Concerning the oversample methodology, adding new samples in the distribution of the
 490 training set can “impose non-uniform error costs, causing the learner to be biased in favour
 491 of predicting the rare class” (He and Ma, 2013, p. 37). As we can see in Figure 6, the
 492 proposed model does not suffer from this problem, since the less representative classes are
 493 not “stealing” sensitivity of the more representative ones, like (PE) or (PP). This indicates,
 494 within the results shown, that the model properly handled the unbalanced dataset problem.

495 After these considerations, the sensitivity calculated from results shown in Figure 6 is
 496 presented in Figure 7, along with results from a previous paper and the final test results.

497



498

499 *Figure 7 - Sensitivity for every class in the dataset according to results published by*
500 *Kedzierski et al. (2019) using a k-Nearest Neighbour algorithm, Support Vector Machine*
501 *Classifier (SVC) in the methodology proposed by Kedzierski et al. (2019), SVC in the Monte*
502 *Carlo Cross-Validation and SVC in the final test.*

503 Assuming that the dataset is statistically representative, the Final test results are what could
504 be approximately expected of the performance of the algorithm in “real” conditions, that is,
505 with the application of the deployed model to non-labeled or unknown spectra. Contrarily to
506 the results given by the MCCV and the methodology adapted from Kedzierski et al. (2019)
507 (refer to section 2.7. for a comparison between methodologies), which used the entire
508 dataset for training, having test samples used for training even if in different
509 splits/simulations, the test set in this case was completely new to the algorithm, giving an
510 “unbiased” prediction.

511 All methodologies had comparable sensitivities, with those from Kedzierski et al. (2019)
512 being lower for some classes (Figure 7). Nonetheless, the referred author had shown that
513 conventional machine learning algorithms could be powerful tools for classifying
514 microplastics spectra. Interestingly, classes that did not perform as well with the kNN
515 algorithm, like “PEVA”, “Poly(ethylene) like” and Poly(propylene) like” also had worse
516 sensitivities with SVC in the MCCV, but presented a slight improvement. In general, this
517 comparison indicates that Support Vector Machine Classifiers offer more appropriate
518 classification boundaries for this specific task with even better scalability, since kNN
519 algorithms are highly computationally demanding with large datasets (Gutierrez et al., 2016).

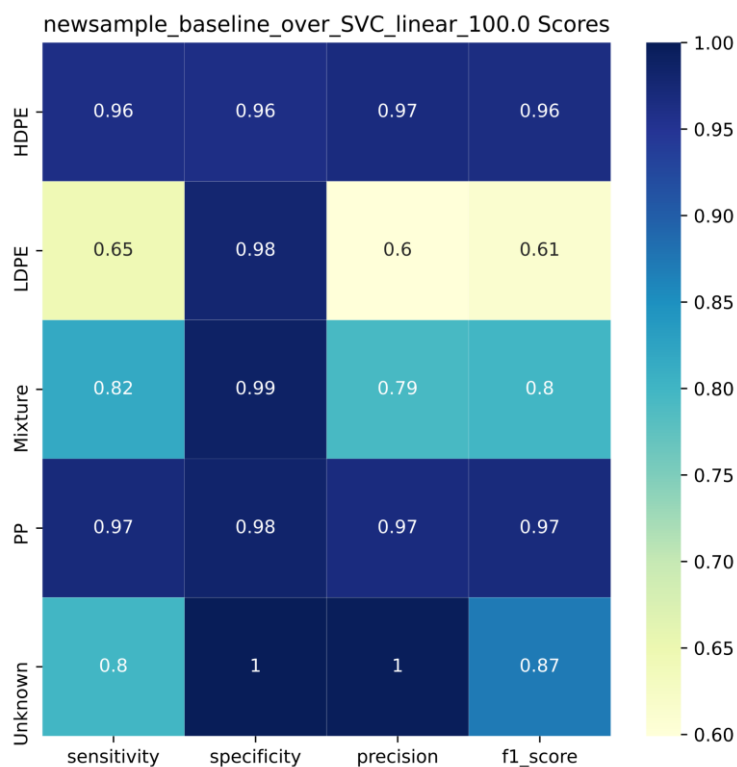
520 The “Cellulose like” class presented no sensitivity, however. Among 240 samples, 4 of them
521 belong to this class, therefore this undesirable event is plausible to happen and should not
522 be interpreted as the actual performance for that class. In this regard, we expect that more
523 samples in this class could improve the model performance.

524 Assuming the data is large enough, independent and identically distributed (i.i.d), the overall
525 accuracy (non class-specific) of the model (which can be given with a 95% probability) is
526 approximately 91.25% +/- 3.6% (Mitchell, 1997, p. 132, formula 5.1). The condition “large
527 enough data” ensures that the sample mean will be approximately normal, due to the Central
528 Limit Theorem (Mitchell, 1997, p. 142). A more conservative confidence interval without such
529 assumption can be computed by using Hoeffding's Inequality (Hoeffding, 1963), which can
530 establish the following bound for the generalization accuracy (Abu-Mostafa, 2012, p. 40):

$$531 \quad |A_{out}(g) - A_{in}(g)| = |E_{out}(g) - E_{in}(g)| \leq \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}},$$

532 with probability greater than $1 - \delta$ where M is the number of final models to be evaluated in
533 the test set, $A_{out}(g)$ is the generalization accuracy of the model g , $A_{in}(g)$ is the in-sample
534 accuracy, similarly $E_{out}(g)$ the generalization error (the ratio of misclassified labels) and
535 $E_{in}(g)$ the in-sample error. With $N = 240$, $M = 1$ and $\delta = 0.05$ and assuming the data is i.i.d,
536 the confidence interval is 91.25% +/- 8.8% with 95% of probability.

537 In regard to the model tested on the Jung dataset, a detailed score is presented in Figure 8.
538 The lower sensitivities for “LDPE” and “Mixture” are attributed to errors similar to those found
539 in the confusion matrix in Figure 6.



540

541 *Figure 8 - Detailed score for the MCCV applied to the Jung dataset considering five classes:*
 542 *High-density Poly(ethylene) (HDPE), Low-density Poly(ethylene) (LDPE), a mixture of PE*
 543 *and PP, Poly(propylene) (PP) and Unknown*

544 The model's cross-entropy log loss was 0.24554 (with a standard deviation of 0.0662) and
 545 the accuracy was 94% (with a standard deviation of 0.01209). This result is in accordance
 546 with those from the Kedzierski dataset in the MCCV (shown in Figure 5), although their
 547 comparison is not straightforward since the number of classes is different. The excellent
 548 performance of the Linear SVC classifier on a different dataset is another evidence of the
 549 linear separability of MP spectra.

550

551 **4. Conclusion**

552 The present study was able to demonstrate the performance of different machine learning
 553 classification algorithms to the classification of ocean microplastics using previously
 554 identified samples' ATR-FTIR spectra as input data. For this purpose, we proposed log-loss

555 to measure models' multi-class probabilities, an approach that has not been used previously
556 in machine learning methods for microplastics characterization.

557 We presented a machine learning combination not yet proposed for classification of
558 microplastics, choosing linear SVC as the final classifier after thoroughly evaluating multiple
559 conventional classification algorithms. The rigorous pipeline methodology described in this
560 paper is essential to avoid introducing bias in the model training and evaluation, which
561 supports the selection of the algorithm and substantiates the performance obtained.

562 The best pre-processing of the raw spectra was also evaluated for each classifier. We
563 identified that for the selected model, only a baseline correction and a naive oversampling
564 was more effective.

565 Linear SVC is well suited for scaling up. After being trained, the classifier can be directly
566 deployed and applied to classify unknown spectra. The user doesn't have to upload the
567 database or train the SVC. Given a spectra, the proposed model would return a probability
568 score for each class, rather than simply a deterministic output. This procedure seems to be
569 more realistic with practice, since the researcher may have greater or lesser confidence in
570 the resulting model's evaluation, depending on the probabilities returned.

571 This study attests for the use of this methodology applied to ATR-FTIR data. In many cases,
572 other vibrational spectroscopic techniques may have been used, namely Raman and
573 microspectroscopy. Despite the new challenges these techniques impose, such as particle
574 morphology, moisture, and blank spectra, in the case of an FPA detector, due to the
575 similarities between them, similar results could be expected. This hypothesis could be
576 verified by applying the same evaluation methodology described in this paper.

577 Machine learning models work within the learned database, which means that the database
578 must be representative of the required task. Thus, it should be applied with caution to
579 spectra collected under different circumstances. Beyond that, the major limitation of the
580 learning approach is that the final model cannot predict samples for a class it does not know,

581 or doesn't know well, which could restrict the full deployment of this procedure, seeing that
582 polymers can be of many kinds.

583 The disproportionate amount of samples in some classes in the dataset was addressed in
584 this study, but not solved. Regardless of the method adopted, "artificially balancing the
585 training distribution may help with the effects of class imbalance, but does not remove the
586 underlying problem" (He and Ma, 2013). In other words, the most appropriate solution to the
587 problem needs to be solved through some specialized algorithm, which does not seem to be
588 the case, since the statistical support will be missing anyway, or through the acquisition of
589 new data, which seems to be the critical point to solve the issue of the unbalanced dataset
590 and even improve the results in this and future works.
591 Data sharing is necessary for the improvement of machine learning algorithms as it
592 increases the size and diversity of training databases. Nonetheless, the authors came
593 across some obstacles, such as: different file formats, huge files with lots of redundant
594 metadata, missing class labels on identified spectra and data that were already pre-
595 processed. We recommend the publication of raw spectra, as the prior pre-processing may
596 introduce bias. The JCAMP-DX file format is a standard defined by the The International
597 Union of Pure and Applied Chemistry (IUPAC) and is considered the optimal format for
598 sharing spectroscopy data. Also, providing the polymer class labels of every spectra should
599 reduce redundant work and speed the collective learning process.

600 Since identification is a critical step in the study of these contaminants, improving confidence
601 and speeding up the process are crucial to the advancement of the area. Despite not yet
602 being able to stand alone as a method to automatically classify every conceivable sample,
603 given its limitations, this study presents robust statistics to support the utilization of machine
604 learning methods to the problem of automatic classification of microplastics.

605

606 ¹ Comparing with Kedzierski et al., 2019 dataset, we deleted 12 spectra, since they were not
607 available as raw data.. Spectras related to Cellulose Acetate are also unavailable in the raw
608 data, but we requested to Kedzierski, which gently provided us the data.

609 ² Quotes were used here because, despite every effort to minimize bias, machine learning
610 models cannot really be unbiased.

611 ³ As long as this massive dataset adequately represents the distribution of MP classes in a
612 way that the imbalance is less deleterious. Even in a scenario with little data, Linear SVC
613 without oversample is still highly competitive.

614 **Acknowledgements**

615 This work was funded by the Brazilian government through the National Council for Scientific
616 and Technological Development (CNPq). Thanks are also due to Mikael Kedzierski, author
617 of a previous paper, who was considerate and transparent in regard to his own work.

618

619 **References**

620 Abu-Mostafa, Y.S., Magdon-Ismail, M., Lin, H.-T., 2012. Learning From Data.

621 Andrad, A.L., 2017. The plastic in microplastics: A review. Marine Pollution Bulletin 119, 12–22.

622 <https://doi.org/10.1016/J.MARPOLBUL.2017.01.082>

623 Barnes, D.K.A., 2002. Biodiversity: Invasions by marine life on plastic debris. Nature 416.

624 <https://doi.org/10.1038/416808a>

625 Ballabio, D., Grisoni, F., & Todeschini, R. 2018. Multivariate comparison of classification
626 performance measures. Chemometrics and Intelligent Laboratory Systems, 174, 33-44.

627 Brandt J, Bittrich L, Fischer F, et al. High-Throughput Analyses of Microplastic Samples Using

628 Fourier Transform Infrared and Raman Spectrometry. Applied Spectroscopy. 2020;

629 74(9):1185-1197. <https://doi.org/10.1177/0003702820932926>

630 Burman, P., 1989. A comparative study of ordinary cross-validation, v-fold cross-validation and

631 the repeated learning-testing methods. Biometrika 76.

632 <https://doi.org/10.1093/biomet/76.3.503>

633 Cabernard, L., Roscher, L., Lorenz, C., Gerdt, G., Primpke, S., 2017. Comparison of Raman and

634 Fourier Transform Infrared Spectroscopy for the Quantification of Microplastics in the

635 Aquatic Environment. *Environmental Science & Technology* 2018 52 (22), 13279-13288.
636 <https://doi.org/10.1021/acs.est.8b03438>

637 Chae, Y., An, Y.J., 2017. Effects of micro- and nanoplastics on aquatic ecosystems: Current
638 research trends and perspectives. *Marine Pollution Bulletin* 124.
639 <https://doi.org/10.1016/j.marpolbul.2017.01.070>

640 Classifier comparison [WWW Document], n.d. . [https://scikit-](https://scikit-learn.org/stable/auto_examples/classification/plot_classifier_comparison.html)
641 [learn.org/stable/auto_examples/classification/plot_classifier_comparison.html](https://scikit-learn.org/stable/auto_examples/classification/plot_classifier_comparison.html).

642 Conroy, J., Ryder, A. G., Leger, M. N., Hennessey, K., Madden, M. G., 2005. Qualitative and
643 quantitative analysis of chlorinated solvents using Raman spectroscopy and machine
644 learning", *Proc. SPIE 5826, Opto-Ireland 2005: Optical Sensing and Spectroscopy*.
645 <https://doi.org/10.1117/12.605056>

646a Silva, V.H., Murphy, F., Amigo, J.M., Stedmon, C., Strand, J. 2020. Classification and
647 Quantification of Microplastics (<100 µm) Using a Focal Plane Array-Fourier Transform
648 Infrared Imaging System and Machine Learning. *Anal. Chem.* 92(20), 13724-13733.

649 Fan, R., Chang, K., Hsieh, C., Wang, X., Lin, C. 2008. LIBLINEAR: A Library for Large Linear
650 Classification. *Journal of Machine Learning Research*

651 Fendall, L.S., Sewell, M.A., 2009. Contributing to marine pollution by washing your face:
652 Microplastics in facial cleansers. *Marine Pollution Bulletin* 58.
653 <https://doi.org/10.1016/j.marpolbul.2009.04.025>

654 Forests of randomized trees [WWW Document], n.d. . [https://scikit-](https://scikit-learn.org/stable/modules/ensemble.html#forests-of-randomized-trees)
655 [learn.org/stable/modules/ensemble.html#forests-of-randomized-trees](https://scikit-learn.org/stable/modules/ensemble.html#forests-of-randomized-trees).

656 Free, C.M., Jensen, O.P., Mason, S.A., Eriksen, M., Williamson, N.J., 2014. High-levels of
657 microplastic pollution in a large , remote , mountain lake. *Marine Pollution Bulletin* 85, 156–
658 163. <https://doi.org/10.1016/j.marpolbul.2014.06.001>

659 Galgani, F., Hanke, G., Maes, T., 2015. Global distribution, composition and abundance of
660 marine litter, in: *Marine Anthropogenic Litter*. https://doi.org/10.1007/978-3-319-16510-3_2

661 Géron, A. 2017. *Hands-on Machine Learning with Scikit-Learn and TensorFlow*. O'Reilly Media

662 GESAMP, 2019. Guidelines for the monitoring and assessment of plastic litter in the ocean.

663 Ghosal, S., Chen, M., Wagner, J., Wang, Z., Wall, S., 2018. Molecular identification of polymers
664 and anthropogenic particles extracted from oceanic water and fish stomach: A Raman micro-
665 spectroscopy study *. *Environmental Pollution* 233, 1113–1124.
666 <https://doi.org/10.1016/j.envpol.2017.10.014>

667 Goodfellow, I., Bengio, Y., Courville, A., 2016. *Deep Learning*. MIT Press.

668 Gutierrez, P., Lastra, M., Bacardit, J., Benitez, J., Herrera, F., 2016. GPU-SME-kNN: scalable
669 and memory efficient kNN and lazy learning using GPUs. *Inf. Sci.*, 373 (2016), pp. 165-182

670 He, H., Ma, Y., 2013. *Imbalanced learning: Foundations, algorithms, and applications*,
671 *Imbalanced Learning: Foundations, Algorithms, and Applications*.
672 <https://doi.org/10.1002/9781118646106>

673 Hoeffding, W., 1963. Probability inequalities for sums of bounded random variables. *Journal of*
674 *the American Statistical Association*. **58** (301): 13–30.
675 [doi:10.1080/01621459.1963.10500830](https://doi.org/10.1080/01621459.1963.10500830).

676 Hidalgo-Ruz, V., Gutow, L., Thompson, R.C., Thiel, M., 2012. Microplastics in the Marine
677 Environment: A Review of the Methods Used for Identification and Quantification.
678 *Environmental Science Technology* 46, 3060–3075. <https://doi.org/10.1021/es2031505>

679 Lufnagl, B., Steiner, D., Renner, E., Löder, M.G.J., Laforsch, C., Lohninger, H., 2019. A
680 methodology for the fast identification and monitoring of microplastics in environmental
681 samples using random decision forest classifiers. *Analytical Methods* 11, 2277–2285.
682 <https://doi.org/10.1039/C9AY00252A>

683 Jolliffe, I. T., 2002. *Principal Component Analysis*. Springer

684 Melissa R. Jung, George H. Balazs, Thierry M. Work, T. Todd Jones, Sara V. Orski, Viviana
685 Rodriguez C., Kathryn L. Beers, Kayla C. Brignac, K. David Hyrenbach, Brenda A. Jensen,
686 and Jennifer M. Lynch., 2018. Polymer Identification of Plastic Debris Ingested by Pelagic-
687 Phase Sea Turtles in the Central Pacific. *Environmental Science & Technology* 52 (20),
688 11535-11544 DOI: 10.1021/acs.est.8b03118

689

690ung, M.R., Horgen, F.D., Orski, S. v., Rodriguez C., V., Beers, K.L., Balazs, G.H., Jones, T.T.,
691 Work, T.M., Brignac, K.C., Royer, S.J., Hyrenbach, K.D., Jensen, B.A., Lynch, J.M., 2018.
692 Validation of ATR FT-IR to identify polymers of plastic marine debris, including those
693 ingested by marine organisms. *Marine Pollution Bulletin* 127, 704–716.
694 <https://doi.org/10.1016/j.marpolbul.2017.12.061>

695Käppler, A., Fischer, D., Oberbeckmann, S., Schernewski, G., Labrenz, M., Eichhorn, K.-J., Voit,
696 B., 2016. Analysis of environmental microplastics by vibrational microspectroscopy: FTIR,
697 Raman or both? *Analytical and Bioanalytical Chemistry* 408, 8377–8391.
698 <https://doi.org/10.1007/s00216-016-9956-3>

699Kedzierski, M., Falcou-Préfol, M., Kerros, M.E., Henry, M., Pedrotti, M.L., Bruzaud, S., 2019. A
700 machine learning algorithm for high throughput identification of FTIR spectra. *Chemosphere*.

701Kotsiantis, S.B., 2007, Supervised machine learning: A review of classification techniques.
702 *Informatica (Ljubljana)*, 31 (3), 249-268.

703Kuptsov, a. h., Zhizhin, g. n., 1998. Handbook of Fourier Transform Raman and Infrared Spectra
704 of Polymers, Bioseparation.

705Young, Y., Hee, S., Jang, M., Myung, G., Rani, M., Lee, J., Joon, W., 2015. A comparison of
706 microscopic and spectroscopic identification methods for analysis of microplastics in
707 environmental samples. *Marine Pollution Bulletin* 93, 202–209.
708 <https://doi.org/10.1016/j.marpolbul.2015.01.015>

709Debreton, L., Slat, B., Ferrari, F., Aitken, J., Marthouse, R., Hajbane, S., 2018. Evidence that the
710 Great Pacific Garbage Patch is rapidly accumulating plastic. *Scientific Reports* 1–15.
711 <https://doi.org/10.1038/s41598-018-22939-w>

712Lemaître, G., Nogueira, F., Aridas, C.K., 2017. Imbalanced-learn: A python toolbox to tackle the
713 curse of imbalanced datasets in machine learning. *Journal of Machine Learning Research*
714 18.

715Bi, J., Hibbert, D.B., Fuller, S., Vaughn, G., 2006. A comparative study of point-to-point algorithms
716 for matching spectra. *Chemometrics and Intelligent Laboratory Systems* 82.
717 <https://doi.org/10.1016/j.chemolab.2005.05.015>

71 Mitchell, T. M., 1997. Machine Learning. McGraw-Hill Science

71 Moore, D.H., 1987. Classification and regression trees, by Leo Breiman, Jerome H. Friedman,
720 Richard A. Olshen, and Charles J. Stone. Brooks/Cole Publishing, Monterey, 1984,358
721 pages, Cytometry 8. <https://doi.org/10.1002/cyto.990080516>

72 Murphy, K.P., 2012. Machine learning: a probabilistic perspective (adaptive computation and
723 machine learning series), Mit Press. ISBN.

72 Nordhausen, K., 2009. The Elements of Statistical Learning: Data Mining, Inference, and
725 Prediction, Second Edition by Trevor Hastie, Robert Tibshirani, Jerome Friedman.
726 International Statistical Review 77. https://doi.org/10.1111/j.1751-5823.2009.00095_18.x

72 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
728 Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D.,
729 Brucher, M., Perrot, M., Duchesnay, É., 2011. Scikit-learn: Machine learning in Python.
730 Journal of Machine Learning Research 12.

73 Prati, R.C., Batista, G.E.A.P.A., Monard, M.C., 2009. Data mining with unbalanced class
732 distributions: Concepts and methods, in: Proceedings of the 4th Indian International
733 Conference on Artificial Intelligence, IICAI 2009.

73 Primpke, S., Lorenz, C., Gerdt, G., 2017. An automated approach for microplastics analysis
735 using focal plane array (FPA) FTIR microscopy and image analysis. Analytical Methods
736 1499–1511. <https://doi.org/10.1039/c6ay02476a>

73 Primpke, S., Wirth, M., Lorenz, C., Gerdt, G., 2018. Reference database design for the
738 automated analysis of microplastic samples based on Fourier transform infrared (FTIR)
739 spectroscopy. Analytical and Bioanalytical Chemistry 5131–5141.

74 Primpke, S., Cross, R.K., Mintenig, S.M., Simon, M., Vianello, A., Gerdt, G. and Vollertsen, J.
741 2020. Toward the Systematic Identification of Microplastics in the Environment: Evaluation of
742 a New Independent Software Tool (siMPle) for Spectroscopic Analysis. Appl. Spectrosc.
743 74(9), 1127-1138.

74 RandomForestClassifier [WWW Document], n.d. . [https://scikit-](https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html)
745 [learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html](https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html).

746 Renner, G., Nellessen, A., Schwiers, A., Wenzel, M., Schmidt, T.C., Schram, J., 2019. Data
747 preprocessing & evaluation used in the microplastics identification process : A critical review
748 & practical guide. Trends in Analytical Chemistry 111, 229–238.
749 <https://doi.org/10.1016/j.trac.2018.12.004>

750 Rocha-Santos, T., Duarte, A.C., 2015. A critical overview of the analytical approaches to the
751 occurrence, the fate and the behavior of microplastics in the environment. TrAC - Trends in
752 Analytical Chemistry. <https://doi.org/10.1016/j.trac.2014.10.011>

753 Saito, H., Fletcher, R., Yogi, T., Kayo, M., Miyagi, S., Ogido, M., Fujikura, K., 2018. Human
754 footprint in the abyss : 30 year records of deep-sea plastic debris Sanae Chiba. Marine
755 Policy 96, 204–212. <https://doi.org/10.1016/j.marpol.2018.03.022>

756 Sauv e, S., Desrosiers, M., 2014. A review of what is an emerging contaminant. Chemistry Central
757 Journal. <https://doi.org/10.1186/1752-153X-8-15>

758 Shim, W.J., Hong, S.H., Eo, S.E., 2017. Identification methods in microplastic analysis: A review.
759 Analytical Methods 9, 1384–1391. <https://doi.org/10.1039/c6ay02558g>

760 sklearn.calibration.CalibratedClassifierCV. viewed 08 June 2021, <[https://scikit-](https://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html#sklearn.calibration.CalibratedClassifierCV)
761 [learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html#sklearn.](https://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html#sklearn.calibration.CalibratedClassifierCV)
762 [calibration.CalibratedClassifierCV](https://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html#sklearn.calibration.CalibratedClassifierCV)>

763 Song, Y.K., Hong, S.H., Jang, M., Han, G.M., Shim, W.J., 2015. Occurrence and Distribution of
764 Microplastics in the Sea Surface Microlayer in Jinhae Bay, South Korea. Archives of
765 Environmental Contamination and Toxicology 69. <https://doi.org/10.1007/s00244-015-0209-9>

766 Support Vector Machines [WWW Document], n.d. . [https://scikit-](https://scikit-learn.org/stable/modules/svm.html)
767 [learn.org/stable/modules/svm.html](https://scikit-learn.org/stable/modules/svm.html).

768 Thompson, R.C., Swan, S.H., Moore, C.J., vom Saal, F.S., 2009. Our plastic age. Philosophical
769 Transactions of the Royal Society B: Biological Sciences.
770 <https://doi.org/10.1098/rstb.2009.0054>

771 Trunk, G. v, 1979. A Problem of Dimensionality: A Simple Example given. IEEE Transactions on
772 Pattern Analysis and Machine Intelligence 306–307.

773 Turner, A., Holmes, L., 2011. Occurrence, distribution and characteristics of beached plastic
774 production pellets on the island of Malta (central Mediterranean). *Marine Pollution Bulletin*
775 62. <https://doi.org/10.1016/j.marpolbul.2010.09.027>

776 Vapnik, V., *Estimation of Dependencies Based on Empirical Data*. Empirical Inference
777 Science: Afterword of 2006, New York: Springer, 2006.

778 Vapnik, V. *The Nature of Statistical Learning*. Information Science and Statistics. New York:
779 Springer-Verlag, 2000.

780 Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski,
781 E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S.J., Brett, M., Wilson, J., Millman,
782 K.J., Mayorov, N., Nelson, A.R.J., Jones, E., Kern, R., Larson, E., Carey, C.J., Polat, İ.,
783 Feng, Y., Moore, E.W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I.,
784 Quintero, E.A., Harris, C.R., Archibald, A.M., Ribeiro, A.H., Pedregosa, F., van Mulbregt, P.,
785 Vijaykumar, A., Bardelli, A. pietro, Rothberg, A., Hilboll, A., Kloeckner, A., Scopatz, A., Lee,
786 A., Rokem, A., Woods, C.N., Fulton, C., Masson, C., Häggström, C., Fitzgerald, C.,
787 Nicholson, D.A., Hagen, D.R., Pasechnik, D. v., Olivetti, E., Martin, E., Wieser, E., Silva, F.,
788 Lenders, F., Wilhelm, F., Young, G., Price, G.A., Ingold, G.L., Allen, G.E., Lee, G.R., Audren,
789 H., Probst, I., Dietrich, J.P., Silterra, J., Webber, J.T., Slavič, J., Nothman, J., Buchner, J.,
790 Kulick, J., Schönberger, J.L., de Miranda Cardoso, J.V., Reimer, J., Harrington, J.,
791 Rodríguez, J.L.C., Nunez-Iglesias, J., Kuczynski, J., Tritz, K., Thoma, M., Newville, M.,
792 Kümmerer, M., Bolingbroke, M., Tartre, M., Pak, M., Smith, N.J., Nowaczyk, N., Shebanov,
793 N., Pavlyk, O., Brodtkorb, P.A., Lee, P., McGibbon, R.T., Feldbauer, R., Lewis, S., Tygier,
794 S., Sievert, S., Vigna, S., Peterson, S., More, S., Pudlik, T., Oshima, T., Pingel, T.J.,
795 Robitaille, T.P., Spura, T., Jones, T.R., Cera, T., Leslie, T., Zito, T., Krauss, T., Upadhyay,
796 U., Halchenko, Y.O., Vázquez-Baeza, Y., 2020. SciPy 1.0: Fundamental algorithms for
797 scientific computing in Python. *Nature Methods* 17. <https://doi.org/10.1038/s41592-019->
798 0686-2

799 Wang, J., Liu, X., Li, Y., Powell, T., Wang, X., Wang, G., Zhang, P., 2019. Microplastics as
800 contaminants in the soil environment : A mini-review. *Science of the Total Environment* 691,
801 848–857. <https://doi.org/10.1016/j.scitotenv.2019.07.209>

802 Volpert, D. H., 1996. The Lack of A Priori Distinctions between Learning Algorithms, *Neural*
803 *Computation* 8(7), 1341–1390.

804 Ku, Q.S., Liang, Y.Z., 2001. Monte Carlo cross-validation. *Chemometrics and Intelligent*
805 *Laboratory Systems* 56. [https://doi.org/10.1016/S0169-7439\(00\)00122-2](https://doi.org/10.1016/S0169-7439(00)00122-2)

806 Zettler, E.R., Mincer, T.J., Amaral-Zettler, L.A., 2013. Life in the “plastisphere”: Microbial
807 communities on plastic marine debris. *Environmental Science and Technology* 47.
808 <https://doi.org/10.1021/es401288x>

809 Zhang, Y., Kang, S., Allen, S., Allen, D., Gao, T., Sillanpää, M., 2020. Atmospheric microplastics:
810 A review on current status and perspectives. *Earth-Science Reviews*.
811 <https://doi.org/10.1016/j.earscirev.2020.103118>

812 Zimmermann, B., Kohler, A., 2013. Optimizing savitzky-golay parameters for improving spectral
813 resolution and quantification in infrared spectroscopy. *Applied Spectroscopy* 67, 892–902.
814 <https://doi.org/10.1366/12-06723>

815