

Automated Classification of Undegraded and Aged Polyethylene Terephthalate Microplastics from ATR-FTIR Spectroscopy using Machine Learning Algorithms

Christian Ebere Enyoh (✉ cenyoh@gmail.com)

Saitama University, Japan <https://orcid.org/0000-0003-4132-8988>

Qingyue Wang (✉ seiyo@mail.saitama-u.ac.jp)

Saitama University, Japan <https://orcid.org/0000-0002-7673-2836>

Research Article

Keywords: Automated classification, PET microplastics, ATR-FTIR spectroscopy, Machine learning algorithms, Environmental analysis, Method standardization

DOI: <https://doi.org/>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Automated Classification of Undegraded and Aged Polyethylene**
2 **Terephthalate Microplastics from ATR-FTIR Spectroscopy using**
3 **Machine Learning Algorithms**

4 **Christian Ebere Enyoh and Qingyue Wang**

5 Graduate School of Science and Engineering, Saitama University, 255
6 Shimo-Okubo, Sakura-ku, Saitama City, Saitama 338-8570, Japan.

7 Correspondence: cenyoh@gmail.com, ORCID ID: 0000-0003-4132-
8 8988(C.E.E); seiyo@mail.saitama-u.ac.jp, ORCID ID: 0000-0002-7673-2836

9 (Q.W.)

10
11
12
13
14
15
16
17
18
19
20
21
22
23

24

25

26 **Abstract**

27 Automated analysis of microplastics is essential due to the labor-intensive,
28 time-consuming, and error-prone nature of manual methods. Attenuated
29 Total Reflectance Fourier Transform Infrared (ATR-FTIR) spectroscopy offers
30 valuable molecular information about microplastic composition. However,
31 efficient data analysis tools are required to effectively differentiate between
32 various types of microplastics due to the large volume of spectral data
33 generated by ATR-FTIR. In this study, we propose a machine learning (ML)
34 approach utilizing ATR-FTIR spectroscopy data for accurate and efficient
35 classification of undegraded and aged polyethylene terephthalate (PET)
36 microplastics (MPs). We evaluate seven ML algorithms, including Random
37 Forest (RF), Gradient Boosting (GB), Decision Tree (DT), k-Nearest
38 Neighbors (k-NN), Logistic Regression (LR), Support Vector Machine (SVM),
39 and Multi-Layer Perceptron (MLP), to assess their performance. The models
40 were optimized using 5-fold cross-validation and evaluated using multiple
41 metrics such as confusion matrix, accuracy, precision, recall (sensitivity), and
42 F1-score. The experimental results demonstrate exceptional performance by
43 RF, GB, DT, and k-NN models, achieving an accuracy of 99% in correctly
44 classifying undegraded and aged PET MPs. The proposed approach
45 capitalizes on the potential of ATR-FTIR spectra to discern distinct chemical
46 signatures of undegraded and aged PET particles, enabling precise and
47 reliable classification. Furthermore, the method offers the benefit of
48 automating the classification process, streamlining the analysis of
49 environmental samples. It also presents the advantage of providing an
50 effective means for method standardization, facilitating more automated and
51 optimized extraction of information from spectral data. The method's
52 versatility and potential for large-scale application make it a valuable
53 contribution to the field of MP environmental research.

54

55 **Keywords:** *Automated classification; PET microplastics; ATR-FTIR*
56 *spectroscopy; Machine learning algorithms; Environmental analysis;*
57 *Method standardization*

58

59

60

61

62

63 **1. Introduction**

64 Microplastics (MPs), defined as plastic particles smaller than 5mm in size,
65 have become a pervasive environmental concern due to their widespread
66 distribution and potential ecological impacts (Verla et al, 2019; Enyoh et al,
67 2021). Among the various types of MPs, polyethylene terephthalate (PET)
68 MPs are particularly prevalent, originating from commonly used items such
69 as single-use plastic bottles and polyester textiles (Enyoh et al, 2023;
70 Chowdury et al, 2022). While in the environment PET MPs can degrade and
71 thus can serve as an efficient vector for toxic pollutants to ecosystems (Verla
72 et al, 2019a). Therefore, understanding the abundance and distribution of
73 both undegraded and aged PET microplastics is crucial for assessing their
74 environmental risks and formulating effective mitigation strategies.

75 Analyzing microplastics manually is labor-intensive, time-consuming, and
76 error-prone, necessitating the development of automated approaches to
77 streamline the process. One promising technique for microplastic analysis is

78 Attenuated Total Reflectance Fourier Transform Infrared (ATR-FTIR)
79 spectroscopy, which provides valuable molecular information about the
80 composition of microplastics (Ioakeimidis et al, 2016; Chowdhury et al, 2022).
81 However, the vast amount of spectral data generated by ATR-FTIR
82 necessitates the utilization of powerful data analysis tools to accurately
83 distinguish between different MPs types.

84 Recently researchers are now developing automated analytical method for
85 the characterization of MPs. Hufnagl et al. (2019) proposed a Radom Forest
86 Classifier method utilizing micro Fourier Transform Infrared (μ -FTIR)
87 hyperspectral images to identify various types of microplastics (MPs),
88 including polyethylene, polypropylene, poly(methyl methacrylate),
89 polyacrylonitrile, and polystyrene, in environmental samples. They developed
90 a model for four plastic types using spectral descriptors determined by
91 spectroscopy experts for polymer characterization. Kedzierski et al. (2019)
92 developed automated methods for identifying the chemical nature of
93 microplastics (MPs) through FTIR-ATR spectra, using k-nearest neighbors'
94 classification. The spectra were collected during the Tara Expedition in the
95 Mediterranean Sea, and a learning database containing 969 microplastic
96 spectra was created for testing. The results demonstrated the effectiveness
97 of machine learning in identifying spectra of common polymers, such as
98 poly(ethylene). However, it was noted that the learning database would
99 benefit from enhancement with less common microplastic spectra. The
100 method was further applied to over 4,000 spectra of unidentified

101 microplastics. The verification protocol revealed less than a 10% difference
102 in the results between the proposed automated method and human expertise.
103 Notably, 75% of the discrepancies could be easily corrected with minimal
104 intervention, indicating the reliability and efficiency of the automated
105 approach in identifying the chemical nature of microplastics. These findings
106 highlight the potential of machine learning in large-scale microplastic
107 characterization studies and underscore the importance of continuously
108 updating the learning database for enhanced performance.

109 On the other hand, Wander et al. (2020) conducted an exploratory analysis of
110 μ -FTIR imaging by employing Principal Component Analysis (PCA) and
111 Uniform Manifold Approximation and Projection (UMAP) to reduce data
112 dimensionality and visualize particle similarity. Although this strategy
113 significantly reduced the analyzed data and removed background information
114 from the images, further analysis was necessary for spectra characterization.

115 Da Silva et al. (2020) presented an automated analytical method for
116 characterizing small microplastics ($< 100 \mu\text{m}$) using μ -FTIR hyperspectral
117 imaging and machine learning (ML) tools. They evaluated Partial Least
118 Squares Discriminant Analysis (PLS-DA) and Soft Independent Modelling of
119 Class Analogy (SIMCA) models with different data pre-processing strategies
120 for classifying nine of the most common polymers worldwide. Additionally,
121 they analyzed the hyperspectral images to automatically quantify particle
122 abundance and size. PLS-DA demonstrated superior analytical performance,
123 exhibiting higher sensitivity, sensibility, and lower misclassification error

124 compared to SIMCA models. Moreover, PLS-DA was less sensitive to edge
125 effects on spectra and poorly focused regions of particles. The approach was
126 successfully tested on a seabed sediment sample from Roskilde Fjord,
127 Denmark, showcasing the method's efficiency. This proposed method offers
128 an automated and efficient approach for microplastic polymer
129 characterization, abundance enumeration, and size distribution, thereby
130 contributing to methods standardization with significant benefits.

131 In a recent study, Yan et al (2022) developed an ensemble model comprising
132 of 7 ML models including support vector machine, K nearest neighbor, least
133 discriminant analysis etc to identify MPs types from ATR-FTIR spectra. The
134 Kedzierski and Jung Datasets have been used to assess the suggested
135 ensemble learning approach. The findings demonstrate that, in terms of 5
136 metrics (Kappa score, F1 score, accuracy, recall, and precision). Their
137 technique showed excellent performance with Each MP receiving a higher-
138 class report and a clearer confusion matrix (i.e., less muddled categories).
139 Moses et al. (2023) utilized a focal plane array (FPA) based micro-FTIR
140 (μ FTIR) to compare two widely used data analysis algorithms concerning the
141 abundance, polymer composition, and size distributions of MPs derived from
142 selected environmental water samples in the size range of 11–500 μ m. The
143 two algorithms under investigation were: (a) the siMPle analysis tool
144 (systematic identification of MicroPlastics in the environment) in combination
145 with MPAPP (MicroPlastic Automated Particle/fibre analysis Pipeline), and
146 (b) the BPF (Bayreuth Particle Finder). The findings of the study revealed a

147 generally good agreement between the two algorithms, but certain
148 discrepancies were observed, particularly concerning specific polymer
149 types/clusters and the smallest MP size classes. This highlights the
150 importance of conducting a detailed comparison of MP algorithms, as it is
151 crucial for ensuring better comparability of MP data. By addressing these
152 differences and potential limitations, researchers can enhance the accuracy
153 and reliability of MP characterization, thus advancing the understanding of
154 MP pollution in aquatic environments.

155 These studies have demonstrated the use of ML models for MPs
156 identification, however, there are needs to develop systems that could
157 identify MPs at different states (undegraded or aged). As a first step, in this
158 research, we present a novel approach that harnesses the potential of
159 unsupervised machine learning algorithms to automate the classification of
160 undegraded and aged PET MPs using ATR-FTIR spectroscopic data. By
161 employing unsupervised methods, our approach reduces the need for manual
162 labeling of data, making it both cost-effective and adaptable to a wide range
163 of microplastic samples. The ultimate goal of this study is to develop a robust
164 and efficient tool that can accurately differentiate between undegraded and
165 aged PET MPs, facilitating a comprehensive understanding of their
166 prevalence in various environmental matrices.

167 **2. Methodology**

168 **2.1. Samples preparation**

169 Commercial plastic products such disposable water bottles made of PET, was
170 used to obtain the MPs tested in this work (excluding the caps). The PET MPs
171 preparation has been described in detailed in previous studies (Enyoh et al,
172 2022; 2023; Enyoh and Wang, 2022; 2023). The plastic waste was first cut
173 into small pieces using stainless-steel scissors. Subsequently, the resulting
174 pieces were further ground into fine particles using a High-Speed Blender.
175 The grinding process involved cycles of blending with cooling intervals. The
176 number of cycles varied based on the amount of plastic to be ground, with
177 more than 10 cycles required for PET. After grinding, the particles were
178 sieved to obtain PET MPs within a specific size range. The prepared PET MPs
179 (500 μm) were then cleaned by soaking in methanol overnight, followed by
180 drying in an oven and rinsing with ultrapure water.

181 The cleaned MPs were labeled as undegraded PET MPs. Further undegraded
182 PET MPs were prepared by just putting the PET MPs in deionized water for
183 24 hrs at room temperature and using a pristine PET MPs. To age the PET
184 MPs, they were treated with sulfuric acid (H_2SO_4) at an elevated temperature
185 (60 $^\circ\text{C}$). Additionally, artificial thermal aging was performed on PET MPs,
186 which involved exposing the PET MPs to hydrogen peroxide (H_2O_2) and
187 elevated temperature (60 $^\circ\text{C}$) in an oven. The resulting aged MPs were
188 washed and dried, and they were categorized as Aged PET MPs. To confirm
189 the presence of undegraded and aged PET microplastics (MPs), a scanning
190 electron microscopy (SEM) technique was employed using a Variable
191 Pressure Scanning Electron Microscope (VP-SEM) SU-1510 with an

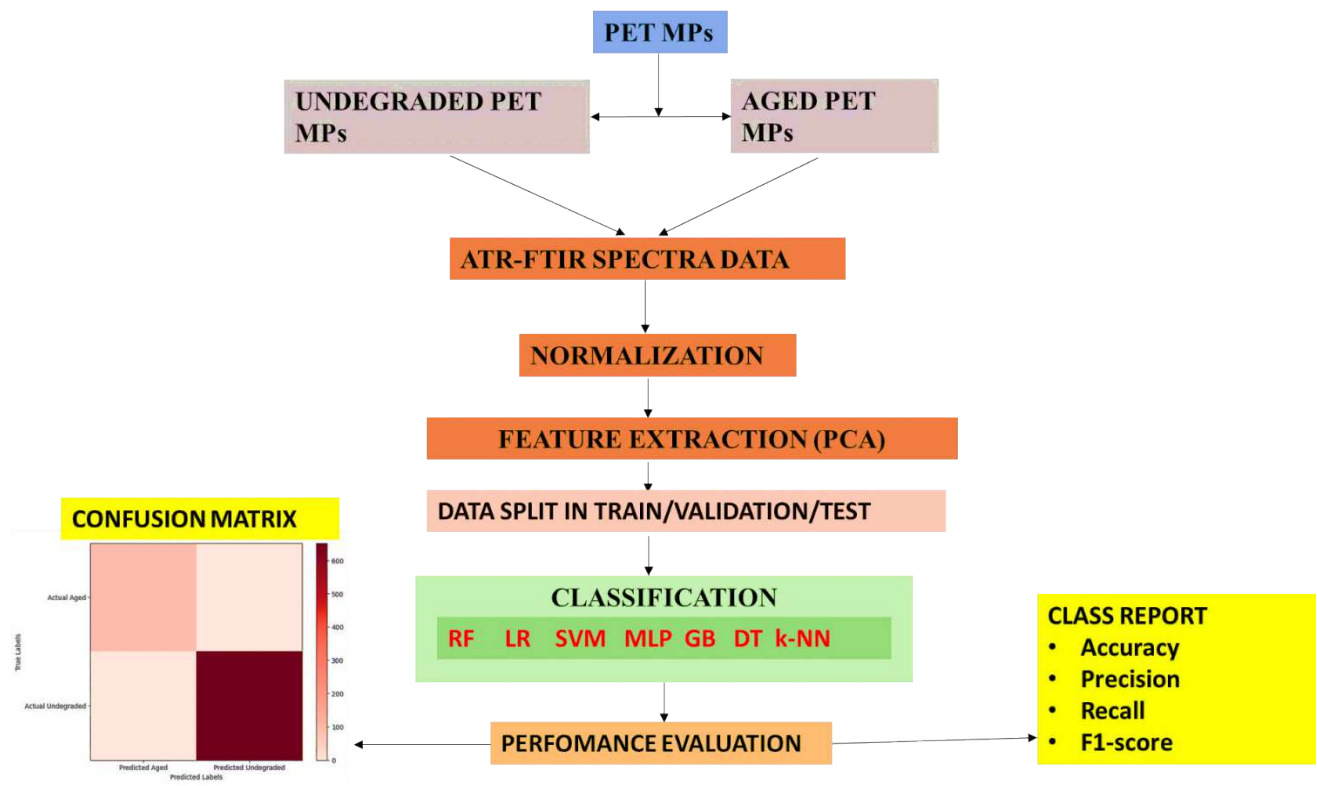
192 accelerating voltage of 15kV (Hitachi Ltd, Tokyo, Japan). The samples
193 underwent preparation and were placed on a stud before undergoing sputter
194 coating. Sputter coating involved applying a film approximately 15 nm thick
195 using argon gas at a pressure of around 4 psi and a current of approximately
196 16 mA for a duration of 4 minutes. The sputter coating process was carried
197 out using the E102 Ion Sputter (Hitachi Ltd, Tokyo, Japan). By employing
198 SEM and sputter coating, a high-resolution images and surface information
199 was obtained, allowing the verification of the presence and characteristics of
200 undegraded and aged PET MPs with a greater level of detail.

201 **2.2. Sample analysis- ATR-FTIR spectral acquisition**

202 Attenuated Total Reflectance-Fourier-Transform Infrared (ATR-FTIR)
203 spectroscopy is a widely employed method to assess changes in the functional
204 groups of polyethylene terephthalate (PET) during environmental
205 degradation (Chowdhury et al., 2022). For the analysis, the functional groups
206 on the surface of the prepared microplastic particles (MPs) were determined
207 using an ATR-FTIR system (JASCO FTIR-6100). Before analyzing the MPs, the
208 instrument was blanked to ensure accurate measurements. The MP samples
209 were securely attached to a KBr disc and placed in the FTIR instrument for
210 measurement. The infrared spectrum was recorded in the range of 400–4000
211 cm^{-1} by averaging 64 scans at a resolution of 4 cm^{-1} . This process provided
212 valuable molecular information about the composition of the PET
213 microplastics, aiding in the investigation of environmental deterioration
214 effects.

215 **2.3. Machine learning (ML) development**

216 The ML sequence of workflow for the classification of undegraded and aged
217 PET MPs is illustrated in Figure 1. The entire process, including data
218 processing, model training, and classification, was executed using the Python
219 programming language within a Jupyter Notebook environment. The
220 computations were conducted on a system equipped with a 64-bit Intel Core
221 i5 vPro processor and 4 GB of RAM.



222
223 **Figure 1: ML sequence of workflow for the classification of PET MPs**

224 **2.3.1. Data normalization**

225 The ATR-FTIR spectra data typically consist of multiple columns representing
226 different features, such as undegraded and aged spectra. However, these
227 features may have different scales, which could potentially introduce bias

228 towards columns with larger values during the modeling process (Yan et al,
229 2022). To address this issue, the data was normalized using min-max scaling,
230 also known as feature scaling or data normalization. Min-max scaling
231 transforms each feature column in the dataset to a common range, typically
232 between 0 and 1. The formula for min-max scaling is represented as follows:
233

$$234 \text{ Scaled_value} = \frac{\text{original_value} - \text{min_value}}{\text{max_value} - \text{min_value}} \quad (1)$$

235 Where: original_value is the value of a data point in the feature column;
236 min_value is the minimum value of the feature column; max_value is the
237 maximum value of the feature column and scaled_value is the resulting
238 normalized value for the data point, which lies between 0 and 1.

239 By applying this transformation, all the feature columns, including
240 undegraded (0) and aged spectra (1), were scaled to the same range, and
241 further splitting the data into features X and target labels y. The X DataFrame
242 contains the numerical features (0 and 1), and the y Series contains the
243 corresponding numerical class labels (undegraded and aged). This makes
244 them comparable and preventing any bias due to different scales.
245 Normalizing the data in this manner ensures that each feature contributes
246 equally to the machine learning model and improves the model's performance
247 and convergence during training.

248 **2.3.2. Feature Selection**

249 For the feature selection step, the principal component analysis (PCA) an
250 unsupervised ML technique was employed on the previously scaled data. PCA

251 is a widely used dimensionality reduction method that aims to transform the
252 original features in a multivariate dataset (in this case, ATR-FTIR spectral
253 data) into a new set of uncorrelated variables called principal components
254 (PCs) (Enyoh et al, 2023b). These PCs are ordered in such a way that the first
255 component captures the most significant variance in the data, the second
256 component captures the second most significant variance, and so on. The
257 primary objective of using PCA for feature selection is to reduce the number
258 of dimensions (features) while preserving the most important information in
259 the data (Da Silva et al., 2020). The mathematical model for PCA
260 decomposition is shown in Equation (2):

261

$$262 \quad X = TP^T + E \quad (2)$$

263 In this equation, X represents the measured spectral data (sample by
264 wavenumber), T is the score matrix (sample by component), P^T is the loading
265 matrix (component by wavenumber), and E is the residuals (unexplained
266 data, sample by wavenumber).

267 By applying PCA to the scaled data, the number of dimensions is reduced to
268 a specified number of principal components (in this case, 2 components). The
269 selection of the number of components is based on the analyzed explained
270 variance ratios provided by PCA. Once the PCs are obtained, they serve as
271 the new feature set for constructing ML classification models.

272 **2.3.4. Data Splitting**

273 After performing data normalization and feature selection using PCA, the
274 next step in the modeling process involved splitting the spectral dataset into
275 two separate parts by using a `train_test_split` function from Sklearn: the
276 training data and the testing data. This division was done to create a clear
277 distinction between the data used to train the machine learning model and
278 the data used to evaluate its performance. The split was done in a specific
279 ratio to ensure an effective and reliable evaluation of the model's
280 generalization capabilities.

281 The dataset was divided into two subsets:

- 282 1. Training Data: This subset comprised 80% of the original spectral
283 dataset. The training data was used to train the ML model, allowing it
284 to learn the underlying patterns and relationships present in the data.
285 During training, the model adjusted its internal parameters to minimize
286 the prediction errors and optimize its performance on the training data.
 - 287 2. Testing Data: The remaining 20% of the spectral dataset formed the
288 testing data. This independent subset served as a previously unseen
289 sample for the model. After training, the model was evaluated on this
290 testing data to assess its performance in predicting the target labels
291 (e.g., undegraded or aged spectra) for new, unseen samples. The
292 testing data acted as a simulation of real-world scenarios where the
293 model encounters new observations that it has not seen during training.
- 294 By splitting the dataset into training and testing subsets, the model's ability
295 to generalize and make accurate predictions on unseen data was assessed.

296 This process is crucial to determine if the model has learned patterns that
297 can be applied to new, unseen data without overfitting or underfitting.

298 **2.3.5. ML data training and testing Methodology for PET MPs** 299 **classification**

300 Seven (7) ML models Random Forest classifier (RF), Logistic Regression (LR),
301 Support Vector Machines classifier (SVM), Neural Networks based on
302 multilayer perceptron classifier (MLP), Gradient Boost (GB), Decision Trees
303 (DT) and k-Nearest Neighbor (k-NN) were evaluated for classifying the
304 undegraded and aged PET MPs in this study.

305 ***RF***

306 The RF classifier is composed of multiple DTs. When making a new
307 classification, each DT independently provides a classification for the input
308 data. The RF algorithm then evaluates these classifications and selects the
309 final prediction based on the class that receives the most votes from the
310 individual trees (Mao and Wang, 2012; Cinar and Koklu, 2019). RF is
311 particularly efficient in handling datasets with a large number of variables
312 (Enyoh et al., 2023a). The simplified equation for the RF, as represented by
313 equation 3, is as follows:

$$314 \text{RF}_{(x)} = \text{mode}(\text{DT}_1(x), \text{DT}_2(x), \dots, \text{DT}_n(x)) \quad (3)$$

315 Here, RF(x) represents the class prediction made by the RF for a given input
316 instance x. The mode function selects the most frequently occurring class
317 prediction from the individual decision trees DT₁, DT₂, ..., DT_n, where n is the
318 number of trees in the forest. Based on a randomly selected portion of the

319 training data, each decision tree in the RF is built individually. A random
320 selection of predictor variables is also taken into account for partitioning the
321 data at each node of the tree.

322 ***LR***

323 LR primary purpose is to elucidate the relationship between these dependent
324 and independent variables. To achieve this, LR fits the weights of the input
325 variables to the training data, aiming to minimize the discrepancy between
326 the predicted probabilities and the actual class labels (Cruyff et al., 2016).
327 The simplified equation for logistic regression, represented as equation (4),
328 is as follows:

$$329 \quad y = \frac{1}{(1 + e^{(-z)})} \quad (4)$$

330 where the variable "y" denotes the predicted output or the probability of a
331 specific class. This probability is obtained by passing the linear combination
332 of the input variables and their respective weights, represented by "z,"
333 through the sigmoid function. The sigmoid function transforms any real-
334 valued number to a value within the range of 0 to 1, enabling the
335 interpretation of the output as a probability. This property makes logistic
336 regression suitable for tasks where the prediction is associated with a
337 probability score, allowing for a more nuanced understanding of the model's
338 predictions.

339 ***SVM***

340 SVM is a fundamental technique used for both classification and regression
341 tasks. It creates a hyperplane that aids in distinguishing between different

342 classes or predicting numerical values. In two-dimensional space, SVM
343 achieves linear separation, while in three-dimensional space, it uses a planar
344 separation. In multidimensional space, it relies on a hyperplane for effective
345 separation of data points (Schölkopf et al., 2001). The classification process
346 in SVM involves identifying the optimal hyperplane that maximizes the
347 margin between different classes. The larger the margin, the better the
348 separation and generalization of the model (Cinar and Koklu, 2019). The
349 simplified form for the predicted output from SVM, represented by equation
350 (5), is as follows (Enyoh et al, 2023):

$$351 \quad y(x)_{pre} = \sum_{i=1}^n \alpha_i \cdot K(x_i, x_j) + b \quad (5)$$

352 where $K(x_i, x_j)$ is the radial basis function kernel. α_i and b denote Lagrange
353 multiplier and threshold parameter, respectively.

354 ***MLP***

355 In this study, we further utilized a popular artificial neural network (ANN)
356 known as the Multilayer Perceptron (MLP). The MLP learns through a
357 technique called backpropagation, where weights are adjusted either after
358 analyzing the entire dataset or after each individual data point. The
359 architecture of the MLP involves organizing neurons into layers, with a
360 hidden layer situated between the input and output layers. Depending on the
361 complexity of the problem, an MLP can consist of multiple hidden layers. The
362 input layer captures information about the problem to be addressed, while
363 the output layer produces the final results or predictions. The study's findings

364 and data processing within the network are conveyed through this output
365 layer (Enyoh et al., 2023). The equation (6), in which f is the activation
366 function, N is the number of inputs per neuron, and k is the layer (hidden,
367 output), may be used to represent the ANN system in its simplest form (Enyoh
368 et al, 2023a).

$$369 \quad Y_j^{k+1} = f(\sum_{i=1}^N X_i^k w_{ij}^k + b_i^k) \quad (6)$$

370 In this research, the model was configured with 100 hidden layers, and the
371 activation function used was the Rectified Linear Unit (ReLU). ReLU,
372 represented by the function $f(x) = \max(0, x)$, introduces non-linearity to the
373 model and effectively addresses the issue of vanishing gradients. It is a widely
374 adopted activation function in deep learning due to its popularity and
375 effectiveness. For optimizing the training process, the Adam optimization
376 algorithm was employed, and a random state of 42 was set. Adam is known
377 for its adaptive learning rate strategy, which dynamically adjusts the learning
378 rate during training. This adaptive approach scales the gradients based on
379 their estimated first and second moments, resulting in faster convergence
380 and improved performance when compared to traditional gradient descent
381 algorithms. By using Adam, the model achieves faster convergence, allowing
382 for better generalization to unseen data.

383 **GB**

384 The Gradient Boosting (GB) classifier is an ensemble learning technique that
385 combines multiple weak learners, represented by Decision Trees (DTs), to

386 create a robust and accurate model (Hastie et al, 2009). The algorithm follows
387 an iterative process, where it gradually adds new DTs to the ensemble. Each
388 subsequent tree focuses on reducing the errors made by the previous trees
389 (Hastie et al, 2009). During each iteration, the algorithm calculates the
390 gradient of the loss function concerning the predicted values and constructs
391 a new tree to minimize this gradient (Piryonesi et al, 2021). The predictions
392 from all the trees in the ensemble are then combined to make the final
393 prediction. The simplified equation for the Gradient Boosting classifier is the
394 sum of weak learners, where each weak learner compensates for the errors
395 made by the preceding learner. It can be expressed as shown in equation (7).

$$396 \quad y(x) = y_0(x) + \eta * g_1(x) + \eta * g_2(x) + \dots + \eta * g_n(x) \quad (7)$$

397 Where $y(x)$ is the predicted output (whether undegraded or aged), $y_0(x)$ is
398 the initial prediction, $g_1(x)$, $g_2(x)$, ..., $g_n(x)$ are the weak learners (usually
399 decision trees), and η is the learning rate (in this case = 0.1). At the
400 beginning, $y(x)$ is initialized with $y_0(x)$, which is the mean or median value of
401 the target variable.

402 ***DT***

403 DT is often visualized as a tree diagram, where each branch and node
404 represent a classification query. The root node stands for an attribute, and
405 the inner nodes indicate tests or evaluations of properties. The branches
406 depict the outcomes of these evaluations, leading to the final decision
407 represented by the leaf nodes, which correspond to the classes (Enyoh et al.,
408 2023; Rokach and Maimon, 2005). DT offers several advantages, making it

409 well-suited for handling complex problems and providing inferences in the
410 form of logical classification rules (Cinar and Koklu, 2019). Its distinct
411 advantages include ease of implementation, seamless integration into
412 databases, and high reliability (Wu et al., 2008). In its simplified form, a
413 Decision Tree can be expressed as shown in equation (8).

$$414 \quad y(x) = (x_1, x_2, x_3, \dots, x_n) \quad (8)$$

415 Where y is the target variable for classifying (undegraded or aged). The
416 vector x is composed of the features, x_1, x_2, \dots etc., that are used for that task.

417 ***k-NN***

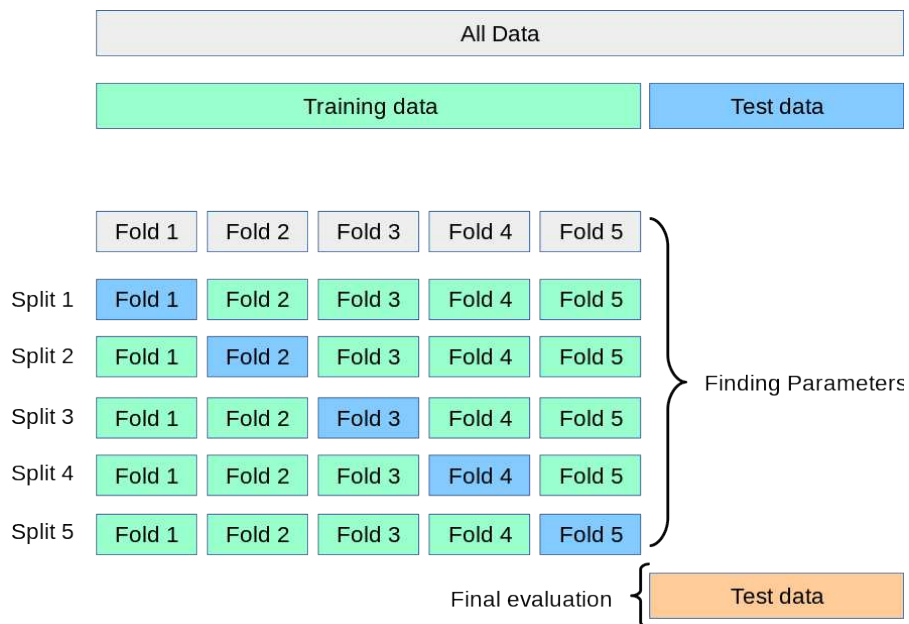
418 k -NN is a popular and widely used machine learning model, especially for
419 large-scale training datasets. It operates based on a distance metric to
420 identify the most similar data points in the training set (Ibeto et al., 2021). In
421 the k -NN algorithm, each data point is conceptually plotted in a multi-
422 dimensional space, where each axis represents a different variable or feature.
423 When a new data point needs to be classified (the test data), the algorithm
424 compares it with all the available data points in the training set. The test data
425 will have several neighbors that are close to it in terms of all the measured
426 characteristics. To determine the class of the test data, the algorithm selects
427 the k nearest data points based on the distance metric. The class with the
428 majority of data points among these selected neighbors is assigned to the test
429 data (Richman, 2011).

430 In this specific study, the k value, representing the number of nearest
431 neighbors to consider, was chosen as 5. This means that when classifying new

432 data points, the algorithm will look at the class labels of the 5 nearest
433 neighbors to make the final prediction.

434 **2.4. Model Optimization**

435 The ML model was optimized using cross-validation to avoid overfitting and
436 improve its performance. In this specific case, the k-fold cross-validation
437 process is performed with 5 folds, meaning the dataset is divided into 5 equal
438 parts, and the model is trained and tested 5 times, each time using a different
439 fold as the test set and the remaining four folds as the training set (Figure 2).



440
441 **Figure 2. 5-fold cross-validation applied in this study (Adapted from**
442 **https://scikit-learn.org/stable/modules/cross_validation.html,**
443 **assessed 28/07/2023)**

444 The accuracy values range from approximately 93.33 % to 100 % (Table 1).
445 The accuracy measures the proportion of correctly classified samples over
446 the total number of samples in each fold. The cross-validation process helps

447 to assess the model's performance on different subsets of the data, providing
 448 an estimate of how well the model generalizes to unseen data. In this case,
 449 the model's performance is consistently high in Folds 2, 3, and 4, achieving
 450 perfect accuracy (100 %), meaning it correctly classified all samples in these
 451 folds. In Folds 1 and 5, the accuracy is slightly lower (93.33 % and 99.47 %,
 452 respectively), but still relatively high, suggesting that the model is performing
 453 well on different subsets of the data.

454 **Table 1. The reported cross-validation (CV) scores obtained in each**
 455 **of the 5 folds for all ML algorithms**

Fold	Accuracy (%)						
	RF	LR	SVM	MLP	GB	DT	KNN
1	93.3	100	97.19	99.47	93.45	93.45	97.99
3							
2	100	100	100	100	100	100	100
3	100	100	100	100	100	100	100
4	100	100	100	100	100	100	100
5	94.8	94.7	94.91	96.52	94.64	94.64	94.78
1	8						

456

457 2.5. Model Evaluation metrics

458 To evaluate the model's performance on the testing data to see how well it
 459 generalizes to unseen samples, common evaluation metrics for classification

460 tasks such as confusion matrix, accuracy, precision, recall, and F1-score was
461 computed. By utilizing these metrics, a comprehensive understanding of a
462 model's strengths and weaknesses in classifying undegraded and aged PET
463 MPs accurately is obtained. Additionally, a learning curve of the different ML
464 were also evaluated.

465 **2.5.1. Confusion Matrix (CM)**

466 A confusion matrix (CM) is a table that summarizes the performance of a
467 classification model. It presents the actual class labels (undegraded and
468 aged) against the predicted class labels (undegraded and aged). The
469 confusion matrix includes four key terms:

- 470 □ True Positives (TP): The number of positive instances correctly
471 classified as positive.
- 472 □ False Positives (FP): The number of negative instances
473 incorrectly classified as positive.
- 474 □ True Negatives (TN): The number of negative instances correctly
475 classified as negative.
- 476 □ False Negatives (FN): The number of positive instances
477 incorrectly classified as negative.

478 The confusion matrix helps to visualize the model's performance across
479 different classes and serves as the foundation for calculating accuracy,
480 precision, recall, and F1-score.

481 **2.5.2. Accuracy**

482 Accuracy measures the proportion of correctly classified instances over the
483 total number of instances in the dataset. Mathematically, accuracy is defined
484 as in equation (10):

$$485 \text{ Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}} \quad (10)$$

486 **2.5.3. Precision**

487 Precision quantifies the ability of the model to correctly identify positive
488 instances among the instances predicted as positive. It focuses on minimizing
489 false positives. The precision is calculated as:

$$490 \text{ Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (11)$$

491

492 **2.5.4. Recall (Sensitivity or True Positive Rate)**

493 Recall evaluates the model's ability to correctly identify positive instances out
494 of all the actual positive instances in the dataset. It focuses on minimizing
495 false negatives. The recall is calculated as:

$$496 \text{ Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (12)$$

497 **2.5.5. F1-score**

498 The F1-score is the harmonic mean of precision and recall and provides a
499 balanced assessment of the model's performance. It takes into account both
500 false positives and false negatives. The F1-score is calculated as:

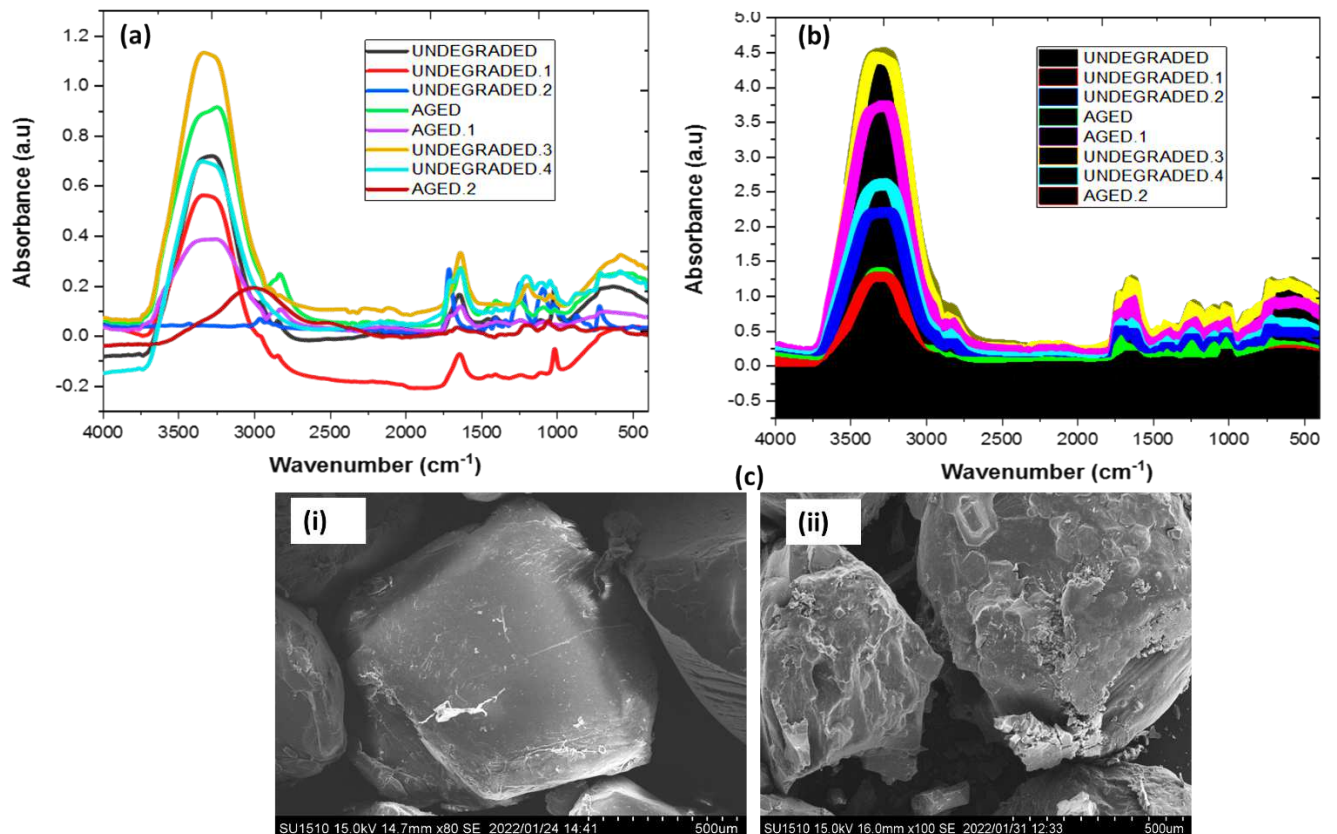
$$501 \text{ F1 - score} = \frac{2 \times \text{Precision} + \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \times \text{TP}}{2 \times \text{TP} + \text{FP} + \text{FN}} \quad (13)$$

502 **3. Results and discussion**

503 **3.1. Description of the ATR-FTIR spectral and SEM**

504 Figure 3 presents the ATR-FTIR spectra for both undegraded and aged PET
505 microplastics (MPs). The spectra indicate the principal bands corresponding
506 to various functional groups in each material. According to standard PET
507 spectroscopy, the bands at 1000 cm^{-1} , 1099 cm^{-1} , 1701 cm^{-1} , 2925 cm^{-1} , and
508 3400 cm^{-1} in the undegraded PET MPs correspond to aromatic -CH vibrations,
509 O-C-C, C = O (carbonyl), -C-C- (alkyl), and -OH groups in the PET structure.
510 These bands provide evidence of the chemical structure specific to
511 undegraded PET (Chowdhury et al, 2022). However, it is noteworthy that the
512 undegraded 2 (pristine PET MPs) showed no -OH group in its spectra. The
513 presence of broad -OH peaks in the spectra of other undegraded PET MPs
514 could be attributed to treatment processes, such as using methanol to remove
515 additives or exposure to water. After aging, the spectra indicate that the
516 major functional groups in the PET are retained, but there is a slight change
517 in the spectrum for both undegraded and aged PET MPs. The band strength
518 at 1099 cm^{-1} decreased in the aged PET MPs due to the thermal aging process
519 that the pristine PET MPs underwent. The thermal treatment dispersed the
520 PET's long chain backbone into smaller fragments, leading to the formation
521 of an interactive cross-linked reactive PET product, primarily associated with
522 the -O-C-C- group present in undegraded PET MPs. The cross-linking process
523 generates radicals, and when these radicals' peroxide, peroxy radicals are
524 formed. During the termination stage, the generated peroxy radical interacts
525 with additional free radical PET monomers. As a result, the band at 1099 cm^{-1} ,

526 related to pristine PET MPs, is nearly lost. This observation confirms that the
527 FTIR-ATR spectra support the alteration and aging of the PET MPs.
528 To provide additional confirmation of the aging process in comparison to
529 undegraded PET MPs, SEM was performed, as depicted in Figure 3c. The
530 SEM images distinctly display a noticeable difference in the surface
531 characteristics of the undegraded and aged PET MPs. In the SEM images,
532 the gaps between the strands of the aged PET MPs appear expanded when
533 compared to the undegraded PET MPs. This observation serves as strong
534 evidence confirming the occurrence of aging in the PET MPs. The SEM
535 analysis visually confirms the structural changes that have taken place on the
536 surface of the aged PET MPs, further corroborating the findings of the aging
537 process.



538

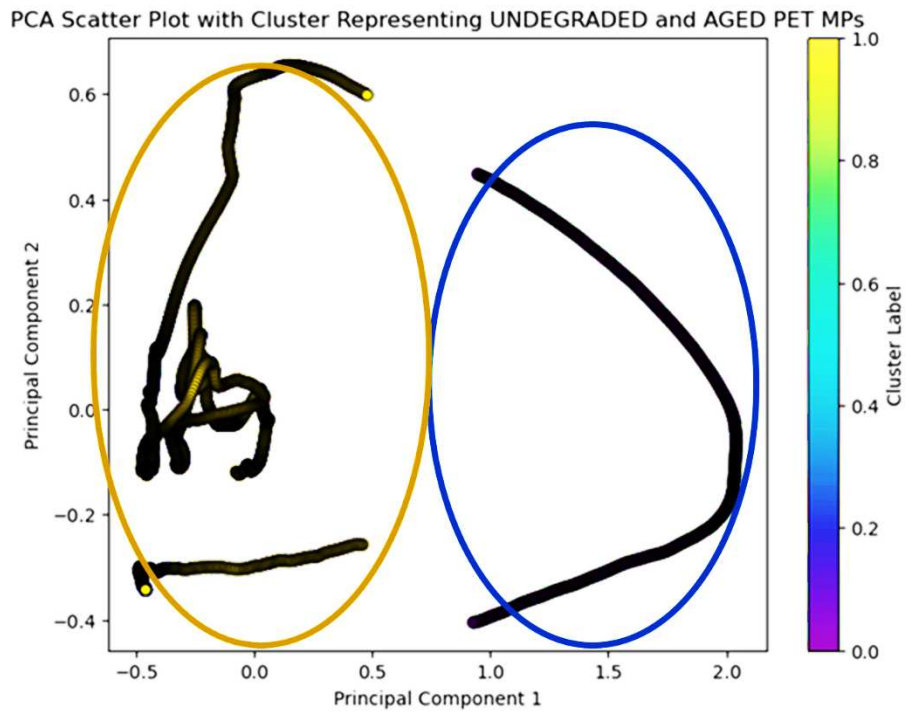
539 **Figure 3. (a) Pre-processed spectra ATR-FTIR spectra of the**
540 **undegraded and aged PET MPs, (b) the dataset used for PCA and (c)**
541 **surface morphology of the degraded and aged PET MPs**

542 **3.2. Exploratory analysis (PCA) of the normalized ATR-FTIR spectra**

543 The application of Principal Component Analysis (PCA) in this study served
544 as a feature selection technique to reduce the dimensionality of the dataset.
545 As a result, two principal components (PC 1 and PC 2) were extracted from
546 the original data. These principal components captured the most important
547 information from the dataset and allowed for a more concise representation
548 of the samples. In the PCA plot (Figure 4), the samples were visualized based
549 on their scores in the PC 1 and PC 2 axes. It was observed that the
550 undegraded PET MPs were predominantly clustered together in one group,
551 while the aged PET MPs formed a separate cluster. This segregation of
552 samples based on their respective clusters indicates that the PCA successfully
553 identified underlying patterns and distinctive characteristics of each group.
554 Furthermore, the explained variance plot revealed that 80% of the data's
555 variability was captured in PC 1, while PC 2 accounted for 10% of the
556 variability (Figure 4b). This indicates that PC 1 carries a substantial amount
557 of information, making it the most significant component for discriminating
558 between undegraded and aged PET MPs. PC 2, although explaining less
559 variance, still contributes valuable information for differentiation between
560 the two groups. In general, the PCA analysis with the two extracted principal
561 components demonstrated its efficacy in distinguishing undegraded and aged

562 PET MPs, providing an informative and concise representation of the dataset
563 with a substantial proportion of the data's variability retained.

564

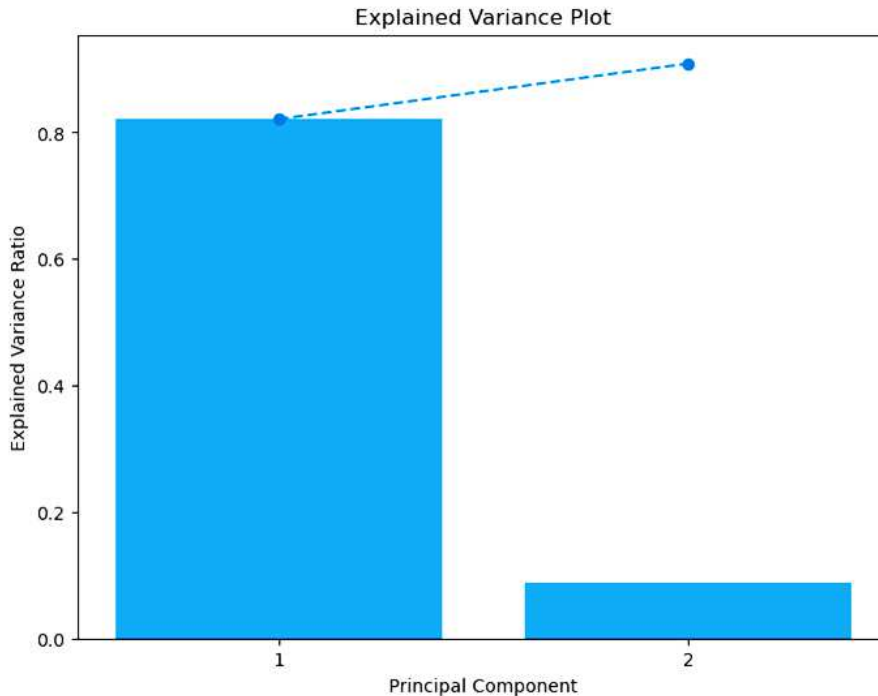


565

566

(a)

567



(b)

Figure 4. (a) The PCA plots in space from the normalized spectral for undegraded and aged PET MPs and (b) the explained variance plot for the principal components extracted

3.3. ML classification models evaluation

3.3.1 Confusion matrix (CM)

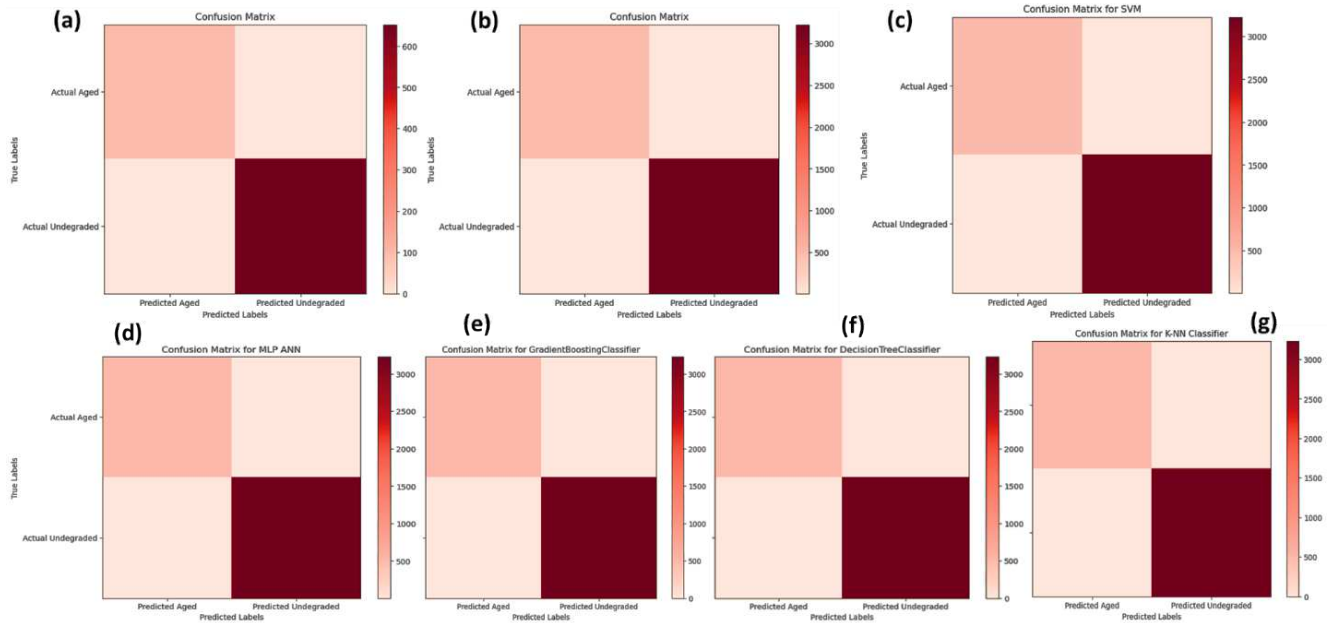
The confusion matrix (CM) is a matrix of numbers that provides valuable insights into how a ML model performs in classifying undegraded and aged PET MPs dataset. In Figure 5 and summarized in Table 2, the CM shows the predicted class labels on the x-axis and the true class labels on the y-axis for each ML algorithm.

By analyzing the CM, we can observe the percentage of correctly classified undegraded and aged PET MPs for each ML algorithm. It also reveals

582 instances where the model gets confused and misidentifies certain PET MPs
583 as others. These misclassifications can be crucial in understanding the
584 strengths and limitations of each ML approach in accurately classifying the
585 MPs.

586 Among the ML models such as Random Forest (RF), Gradient Boosting (GB),
587 Decision Tree (DT), and k-Nearest Neighbors (k-NN), most of the samples are
588 correctly classified with no false positives (FP) expected (Table 2). However,
589 Logistic Regression (LR), Support Vector Machine (SVM), and Multi-Layer
590 Perceptron (MLP) had a few samples i.e. 22, 7 and 4 respectively classified
591 as false positives (FP). Additionally, these models had some samples i.e. 10,
592 6 and 2 respectively, that were expected to be of the positive class but were
593 classified as the negative class, resulting in false negatives (FN).

594 Understanding the FP and FN rates is crucial as it helps to identify the areas
595 where the ML models may have challenges in classification. By examining
596 these misclassifications and the impacts of Principal Component Analysis
597 (PCA) on the ML methods, informed decisions on model improvement and
598 identify strategies to enhance the classification accuracy of undegraded and
599 aged PET MPs can be made.



600

601 **Figure 5. Confusion matrix plot for all ML classifiers (a) RF (b) LR**

602 **(c) SVM (d) MLP (e) GB (f) DT (g) K-NN**

603 **Table 2. Complexity matrix of algorithms used**

Algorithm	Confusion Matrix		
	Positive	Negative	
RF	92	0	Positive
	0	656	Negative
LR	481	22	Positive
	10	3223	Negative
SVM	496	7	Positive
	6	3227	Negative
MLP	499	4	Positive

	2	3231	Negative
GB	503	0	Positive
	0	3233	Negative
DT	503	0	Positive
	0	3233	Negative
k-NN	503	0	Positive
	0	3233	Negative

604

605 3.3.2. Class reports for all ML algorithms

606 In this study, a more comprehensive analysis of the ML algorithms'
607 performance in classifying PET microplastics (MPs) is conducted. To evaluate
608 the specific performance of different methods for undegraded and aged PET
609 MPs classification, accuracy, precision, recall, and F1 score are used as
610 performance metrics. The results of these classification performance
611 measurements are presented in Table 3, providing valuable insights into the
612 strengths and weaknesses of each ML approach in accurately classifying
613 undegraded and aged PET MPs. Based on the evaluation of various
614 performance metrics, it is evident that the ML models generally performed
615 exceptionally well in classifying undegraded and aged PET microplastics
616 (MPs). The accuracy of the models, ranging from 0.98 to 0.99 following 5-fold
617 cross-validation, indicates their high ability to predict the correct classes.
618 Among the models, LR, MLP, and k-NN demonstrated the highest accuracy,
619 achieving a remarkable 0.99. Additionally, RF, GB, DT, and k-NN showed

620 perfect precision, recall, and F1 scores (all equal to 1.00) for correctly
621 classifying both undegraded and aged PET MPs. These models showcased
622 excellent performance and robustness in distinguishing between the two
623 classes, resulting in no false positives or false negatives. However, other ML
624 models such as LR, SVM, and MLP, while still achieving high accuracy,
625 showed metrics slightly below 1.00, especially when classifying aged PET
626 MPs. This indicates that these models may have some challenges in
627 accurately identifying aged PET MPs, as they had some misclassifications.
628 Based on the overall performance and considering the precision, recall, and
629 F1 scores, RF, GB, DT (Figure 6a), and k-NN (Figure 6b) are the most
630 favorable models for the classification of undegraded and aged PET MPs.
631 These models demonstrated superior accuracy and robustness in correctly
632 identifying both classes with no misclassifications. However, further
633 examination and potential fine-tuning of LR, SVM, and MLP may be needed
634 to improve their performance, particularly in classifying aged PET MPs.

635
636
637

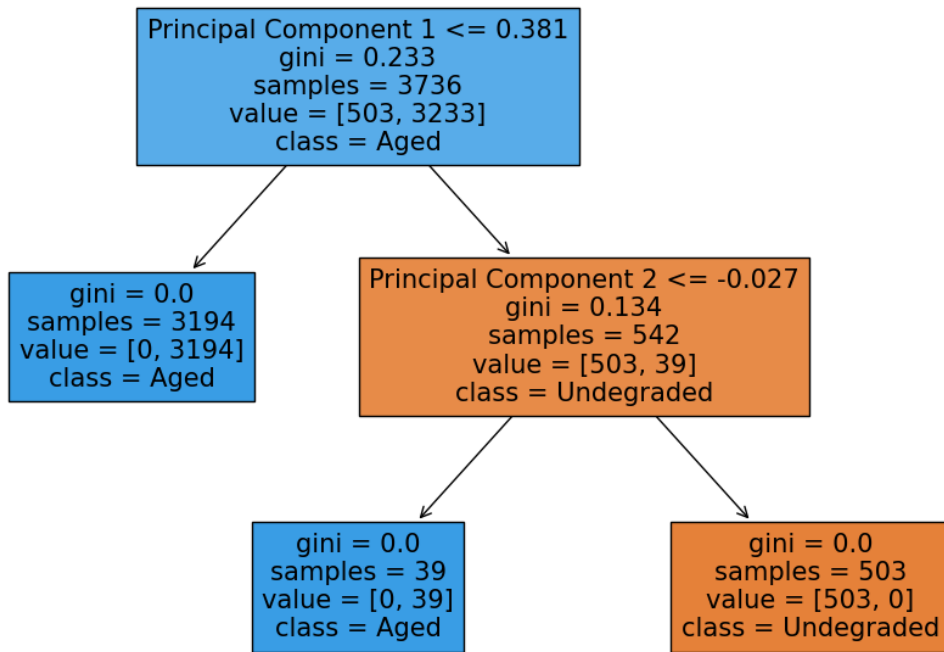
638 **Table 2. Classification report for the different ML algorithms**

Measure	Class (0 =undegraded and 1 = aged)	RF	LR	SVM	MLP	GB	DT	k-NN
Accuracy		0.98	0.99	0.98	0.99	0.98	0.98	0.99
Precision	0	1.00	0.98	0.99	1.00	1.00	1.00	1.00

	1	1.00	0.99	1.00	1.00	1.00	1.00	1.00
Recall	0	1.00	0.96	0.99	0.99	1.00	1.00	1.00
(Sensitivity)	1	1.00	1.00	1.00	1.00	1.00	1.00	1.00
F1-Score	0	1.00	0.97	0.99	0.99	1.00	1.00	1.00
	1	1.00	1.00	1.00	1.00	1.00	1.00	1.00

639

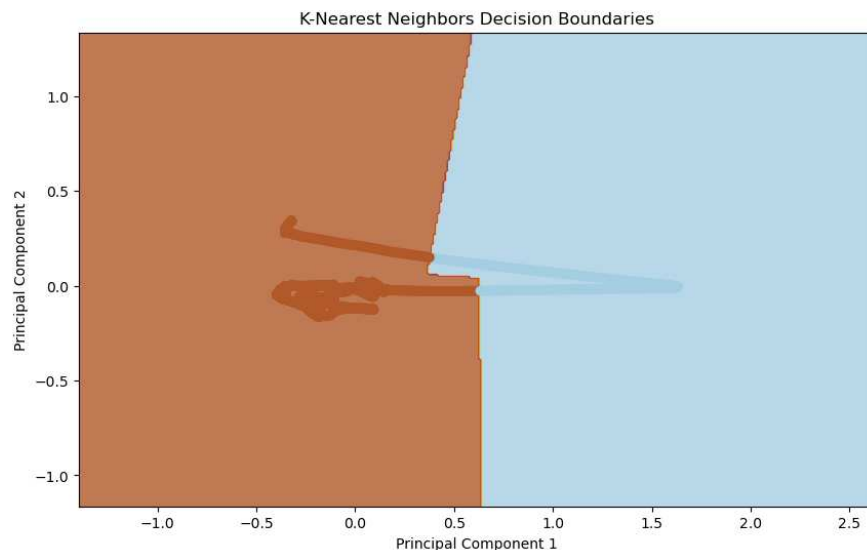
Decision Tree Visualization



640

641

(a)



(b)

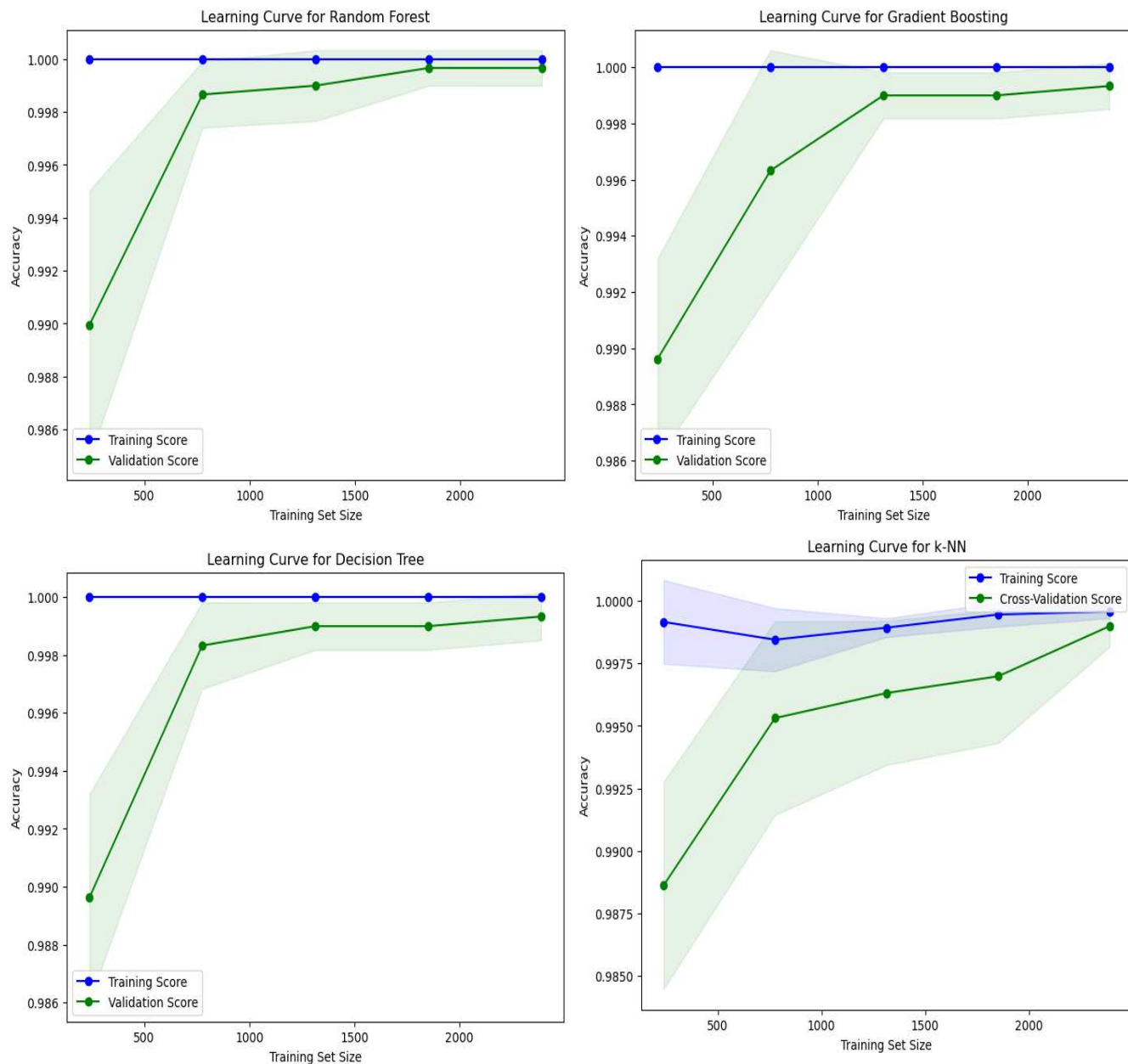
Figure 6. Visualization of (a) DT and (b) k-NN for classifying undegraded and aged PET MPs using ATR-FTIR spectral

3.3.3. Learning curves for ML algorithms

Interpreting learning curves is crucial in assessing the performance and behavior of a machine learning model during the training process. Learning curves depict the model's training and validation (or cross-validation) performance as a function of the training set size or the number of training iterations (Perlich, 2011).

The learning curves for best ML models is presented in Figure 7, which typically show the training set loss(accuracy) and cross-validation set loss on the y-axis and the number of training samples on the x-axis. The training set loss measures how well the model is fitting the training data, while the validation set loss evaluates the model's generalization performance on unseen data. The learning curves play a vital role in assessing the

658 performance and behavior of a machine learning model as the training set
659 size increases. In this case, the learning curves clearly demonstrated positive
660 signs of a well-behaved model. The decreasing and converging nature of both
661 the training and cross-validation curves indicates that the model is learning
662 from the data and improving its performance as more training samples are
663 added. The fact that both curves decrease implies that the model is fitting the
664 training data well, capturing the underlying patterns and trends. The
665 convergence of the training and cross-validation curves is a positive
666 indicator, as it suggests that the model is not suffering from significant
667 overfitting. Overfitting occurs when a model memorizes the training data too
668 closely and fails to generalize to new, unseen data (Ying, 2019). A small gap
669 between the training and cross-validation curves implies that the model's
670 performance on new data is similar to its performance on the training data,
671 which is a desirable outcome. This finding indicates that the model is
672 exhibiting good generalization capabilities, meaning it can make accurate
673 predictions on data it has never seen before. The convergence of the learning
674 curves suggests that the model is not just memorizing the training data but
675 is learning to capture the underlying patterns that can be applied to new data.
676 These ML models are likely to make accurate predictions on new, unseen
677 data, making it a reliable and effective tool for the task at hand.



678

679 **Figure 7. Learning curves for the best ML algorithms for the**

680 **classification of undegraded and aged PET MPs.**

681 **4. Conclusion**

682 In this study, we developed a machine learning approach utilizing ATR-FTIR

683 spectral data to classify undegraded and aged PET microplastics (MPs)

684 particles. Among the seven ML models evaluated, Random Forest (RF),

685 Gradient Boosting (GB), Decision Tree (DT), and k-Nearest Neighbors (k-NN)
686 demonstrated the best performance with an impressive accuracy of 99% in
687 classifying undegraded and aged PET MPs. These results showcase the
688 significant potential of ATR-FTIR spectra in accurately distinguishing
689 between undegraded and aged PET MPs particles. The proposed strategy not
690 only enables effective classification but can also be adapted for use with
691 various environmental samples. Furthermore, our method offers the
692 advantage of automating the sorting process for microplastic particles,
693 making it a valuable tool for standardizing methods. By optimizing spectra
694 and extracting essential information from the data, our approach streamlines
695 and enhances the classification process, providing more reliable and efficient
696 results. Finally, the method's versatility and potential for method
697 standardization make it a valuable contribution to the field of microplastic
698 analysis in environmental research.

699 **Contributions**

700 **C.E.E.:** Conceptualization, Methodology, Software, Formal analysis, Validation,
701 Visualization, Investigation, Data curation, Project administration, Writing- Original draft
702 preparation, Writing- Reviewing and Editing. **W.Q:** Supervision, Project administration,
703 Funding acquisition, Resources, Writing- Reviewing and Editing.

704 **Funding**

705 This study was partially supported by the Special Funds for Basic Research (B)
706 (No.22H03747, FY2022-FY2024) of Grant-in-Aid for Scientific Research of Japanese Ministry
707 of Education, Culture, Sports, Science and Technology (MEXT).

708 **References**

709 Chowdhury, T., Wang, Q. & Enyoh, C.E. Degradation of Polyethylene
710 Terephthalate Microplastics by Mineral Acids: Experimental, Molecular
711 Modelling and Optimization Studies. J Polym Environ (2022).
712 <https://doi.org/10.1007/s10924-022-02578-z>

713
714 Cinar, I, and Koklu, M. (2019). Classification of Rice Varieties Using Artificial
715 Intelligence Methods. International Journal of Intelligent Systems and
716 Applications in Engineering IJISAE, 2019, 7(3), 188-194.

717
718 Cruyff, M.; Böckenholt, U.; van der Heijden, P.G.M.; Frank, L.E.; Chaudhuri,
719 A.; Christofides, C.T.; Rao, C.R. (2016). Handbook of Statistics, Volume 34.
720 Data Gathering, Analysis and Protection of Privacy through Randomized
721 Response Techniques: Qualitative and Quantitative Human Traits, pp. 287 -
722 315, Elsevier. p. 287-315.

723
724 Da Silva, V. H., Murphy, F., Amigo, J. M., Stedmon, C., & Strand, J. (2020).
725 Classification and Quantification of Microplastics (<100 µm) Using a Focal
726 Plane Array-Fourier Transform Infrared Imaging System and Machine
727 Learning. Analytical Chemistry, 92(20), 13724-13733.
728 doi:10.1021/acs.analchem.0c01324

729 Enyoh C.E. and Wan, Q. (2022). Combined experimental and molecular
730 dynamics removal processes of contaminant phenol from simulated

731 wastewater by polyethylene terephthalate microplastics, Environmental
732 Technology, DOI: 10.1080/09593330.2022.2139636

733

734 Enyoh C.E., A.W. Verla, F.O. Ohiagu & E.C. Enyoh (2021). Progress and
735 future perspectives of microplastic research in Nigeria. International Journal
736 of Environmental Analytical Chemistry,
737 <https://doi.org/10.1080/03067319.2021.1887161>

738

739 Enyoh C.E., Duru C.E., Prosper E., Wang Q. (2023). Evaluation of
740 Nanoplastics Toxicity to the Human Placenta in Systems. Journal of
741 Hazardous Materials 446:130600. DOI: 10.1016/j.jhazmat.2022.130600

742

743 Enyoh, C.E., Qingyue Wang , Prosper O. (2022). Response Surface
744 Methodology for modeling the Adsorptive uptake of Phenol from Aqueous
745 solution Using Adsorbent Polyethylene Terephthalate Microplastics.
746 Chemical Engineering Journal Advances. DOI: 10.1016/j.ceja.2022.100370

747

748 Enyoh, C.E.; Wang, Q.; (2023) Adsorption and toxicity characteristics of
749 ciprofloxacin on differently prepared polyethylene terephthalate
750 microplastics from both experimental and theoretical perspectives. Journal of
751 Water Process Engineering, 53, 103909. DOI:
752 <https://doi.org/10.1016/j.jwpe.2023.103909>

753

754 Enyoh, C.E.; Wang, Q.; Momimul, R.H.; Senlin, L.; (2023b) Preliminary
755 characterization and probabilistic risk assessment of microplastics and
756 potentially toxic elements (PTEs) in garri (cassava flake), a common staple
757 food consumed in West Africa Environmental Analysis Health and Toxicology
758 2023; 38(1): e2023005. <https://doi.org/10.5620/eaht.2023005>

759
760 Enyoh, C.E.; Wang, Q.; Senlin, L.(2023) Optimizing the Efficient Removal of
761 Ciprofloxacin from Aqueous Solutions by Polyethylene Terephthalate
762 Microplastics using Multivariate Statistical Approach. Chemical Engineering
763 Science 278(12):118917:. DOI: 10.1016/j.ces.2023.118917

764
765 Hastie, T.; Tibshirani, R.; Friedman, J. H. (2009). 10. Boosting and Additive
766 Trees. The Elements of Statistical Learning (2nd ed.). New York: Springer.
767 pp. 337-384.

768
769 Hufnagl, B.; Steiner, D.; Renner, E.; Löder, M. G. J.; Laforsch, C.; Lohninger,
770 H. A Methodology for the Fast Identification and Monitoring of Microplastics
771 in Environmental Samples Using Random Decision Forest Classifiers. Anal.
772 Methods. 2019, 11 (17), 2277-2285.

773
774 Ibeto C.N., C.E. Enyoh, A.C. Ofomatah, L.A. Oguejiofor, T. Okafocha & V.
775 Okanya (2021): Microplastics pollution indices of bottled water from South

776 Eastern Nigeria, International Journal of Environmental Analytical
777 Chemistry, DOI: 10.1080/03067319.2021.1982926

778

779 Ioakeimidis, C., Fotopoulou, K. N., Karapanagioti, H. K., Geraga, M., Zeri, C.,
780 Papathanassiou, E., ... Papatheodorou, G. (2016). The degradation potential
781 of PET bottles in the marine environment: An ATR-FTIR based approach.
782 Scientific Reports, 6(1). doi:10.1038/srep23501

783

784 Kedzierski, M., Falcou-Préfol, M., Kerros, M. E., Henry, M., Pedrotti, M. L., &
785 Bruzard, S. (2019). *A machine learning algorithm for high throughput*
786 *identification of FTIR spectra: Application on microplastics collected in the*
787 *Mediterranean Sea. Chemosphere, 234, 242-*
788 *251.* doi:10.1016/j.chemosphere.2019.05.113

789

790 Mao, W. and F. Wang, *New advances in intelligence and security informatics.*
791 2012: Academic Press.

792

793 Moses, S.R., Roscher, L., Primpke, S. et al. Comparison of two rapid
794 automated analysis tools for large FTIR microplastic datasets. Anal Bioanal
795 Chem 415, 2975-2987 (2023). <https://doi.org/10.1007/s00216-023-04630-w>

796

797 Perlich, C. (2011). Learning Curves in Machine Learning. In: Sammut, C.,
798 Webb, G.I. (eds) Encyclopedia of Machine Learning. Springer, Boston, MA.
799 https://doi.org/10.1007/978-0-387-30164-8_452

800

801 Piryonesi, S. Madeh; El-Diraby, Tamer E. (2021). Using Machine Learning to
802 Examine Impact of Type of Performance Indicator on Flexible Pavement
803 Deterioration Modeling. Journal of Infrastructure Systems. 27 (2): 04021005.
804 doi:10.1061/(ASCE)IS.1943-555X.0000602.

805

806 Richman, J.S. (2011), Multivariate neighborhood sample entropy: a method
807 for data reduction and prediction of complex data, in Methods in enzymology.
808 Elsevier. p. 397-408.

809

810 Rokach, L.; Maimon, O. (2005). Top-down induction of decision trees
811 classifiers-a survey. *IEEE Transactions on Systems, Man, and Cybernetics -*
812 *Part C: Applications and Reviews*. **35** (4): 476-487.

813

814 Schölkopf, B., Platt, J. C., Shawe-Taylor, J., Smola, A. J., & Williamson, R. C.
815 (2001). Estimating the support of a high-dimensional distribution. *Neural*
816 *computation*, 13(7), 1443-1471.

817 <https://doi.org/10.1162/089976601750264965>

818

819 Verla AW, Enyoh CE, Verla EN (2019). Microplastics, an emerging concern:
820 a review of analytical techniques for detecting and quantifying microplastics.
821 Anal Methods Environ Chem J. 12:15-32.

822
823 Verla, A.W., Enyoh, C.E., Verla, E.N. et al. (2019a). Microplastic-toxic
824 chemical interaction: a review study on quantified levels, mechanism and
825 implication. SN Appl. Sci. 1, 1400. [https://doi.org/10.1007/s42452-019-1352-](https://doi.org/10.1007/s42452-019-1352-0)
826 [0](https://doi.org/10.1007/s42452-019-1352-0)

827
828 Wander, L.; Vianello, A.; Vollertsen, J.; Westad, F.; Braun, U.; Paul, A.
829 Exploratory Analysis of Hyperspectral FTIR Data Obtained from
830 Environmental Microplastics Samples. Anal. Methods. 2020, 12 (6), 781-791.

831
832 Wu, Xindong; Kumar, Vipin; Ross Quinlan, J.; Ghosh, Joydeep; Yang, Qiang;
833 Motoda, Hiroshi; McLachlan, Geoffrey J.; Ng, Angus; Liu, Bing; Yu, Philip S.;
834 Zhou, Zhi-Hua (2008). Top 10 algorithms in data mining. Knowledge and
835 Information Systems. 14 (1): 1-37. doi:10.1007/s10115-007-0114-2

836
837 Yan X., Zhi, C., Alan, M., Yuansong, Q. (2022). An ensemble machine
838 learning method for microplastics identification with FTIR spectrum. Journal
839 of Environmental Chemical Engineering, 10, 4, 108130.
840 <https://doi.org/10.1016/j.jece.2022.108130>

841

- 842 Ying X. (2019). An Overview of Overfitting and its Solutions. *J. Phys.: Conf.*
843 *Ser. 1168* 022022. **DOI** 10.1088/1742-6596/1168/2/022022

Figures

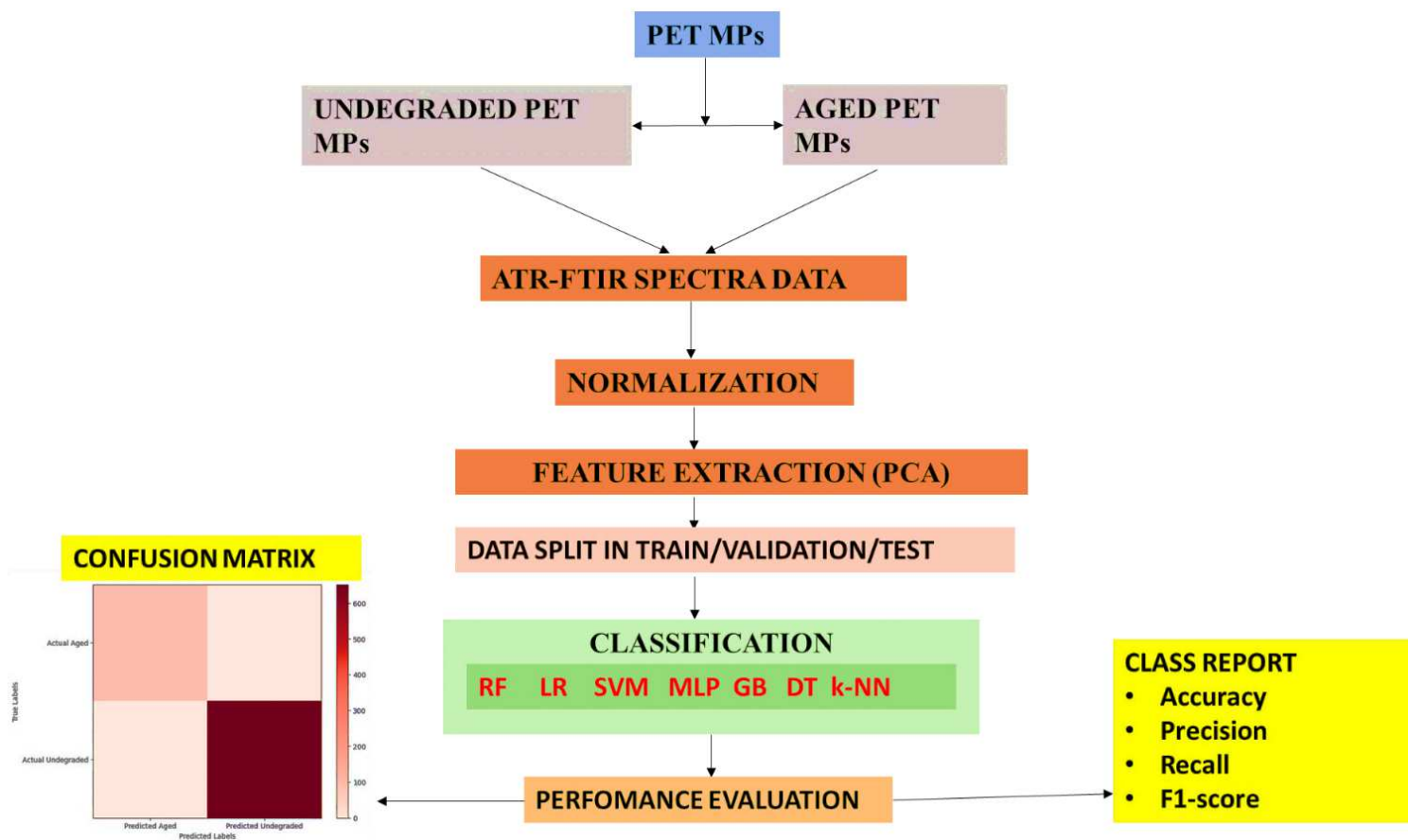


Figure 1

ML sequence of workflow for the classification of PET MPs

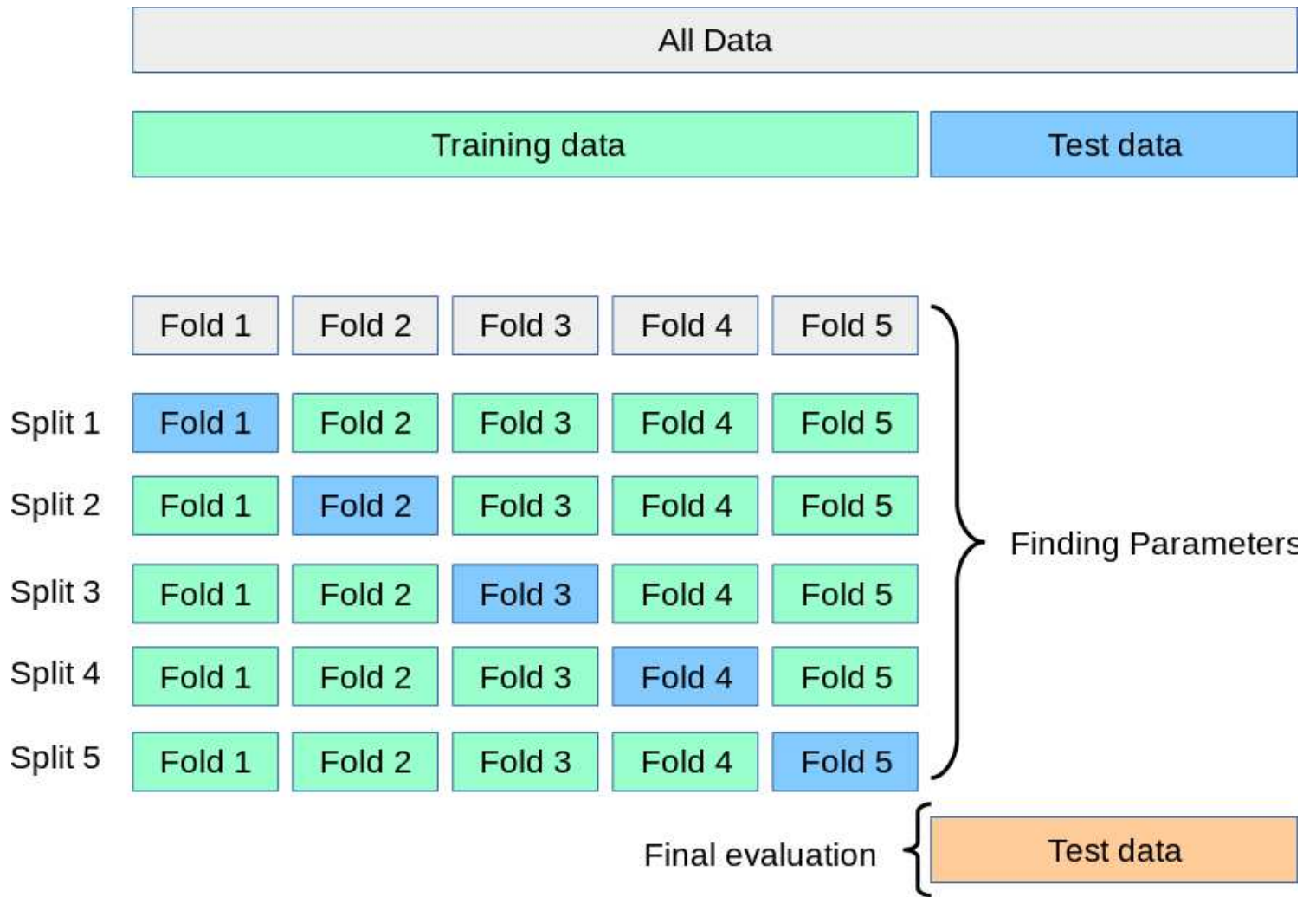


Figure 2

5-fold cross-validation applied in this study (Adapted from https://scikit-learn.org/stable/modules/cross_validation.html, assessed 28/07/2023)

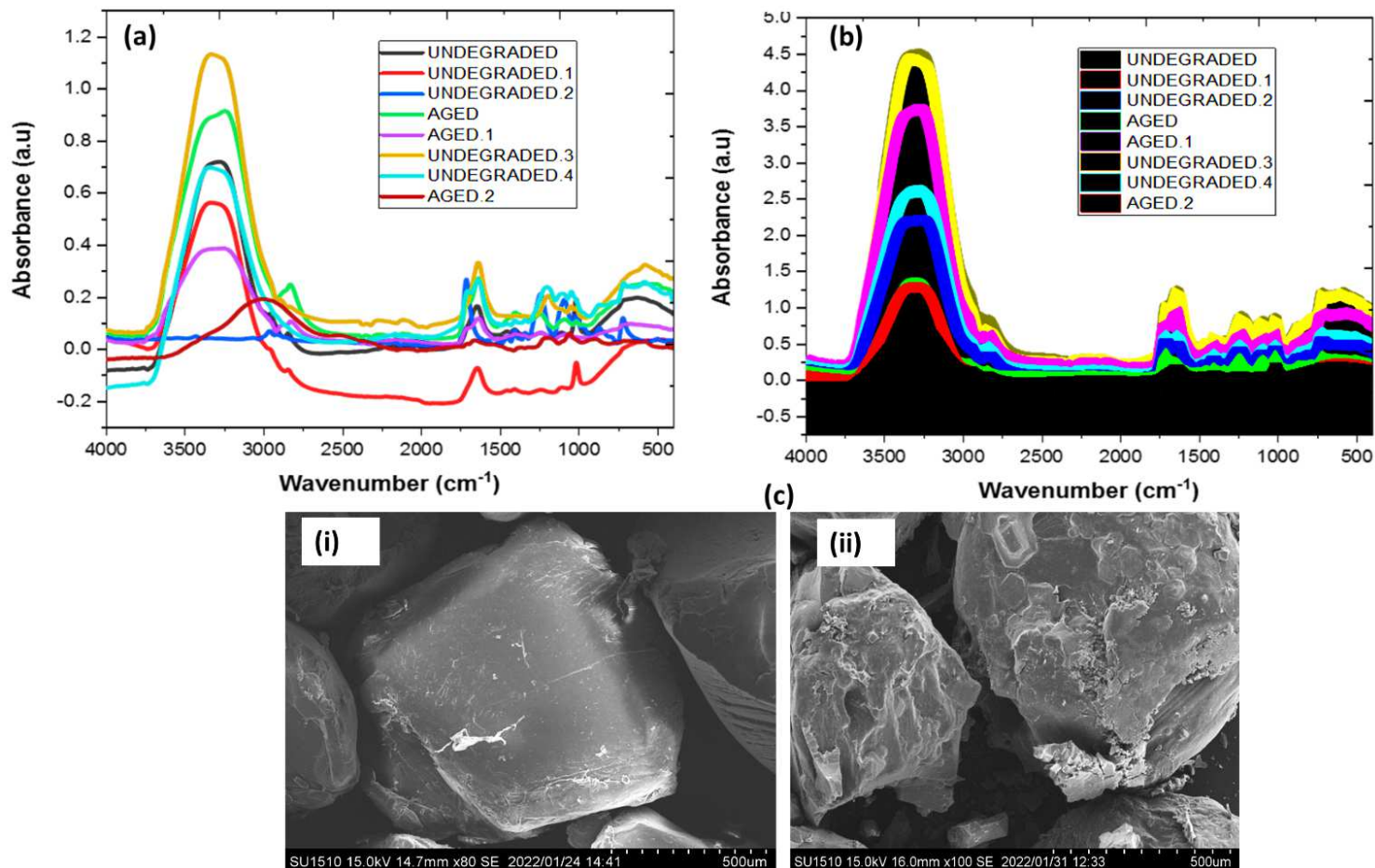
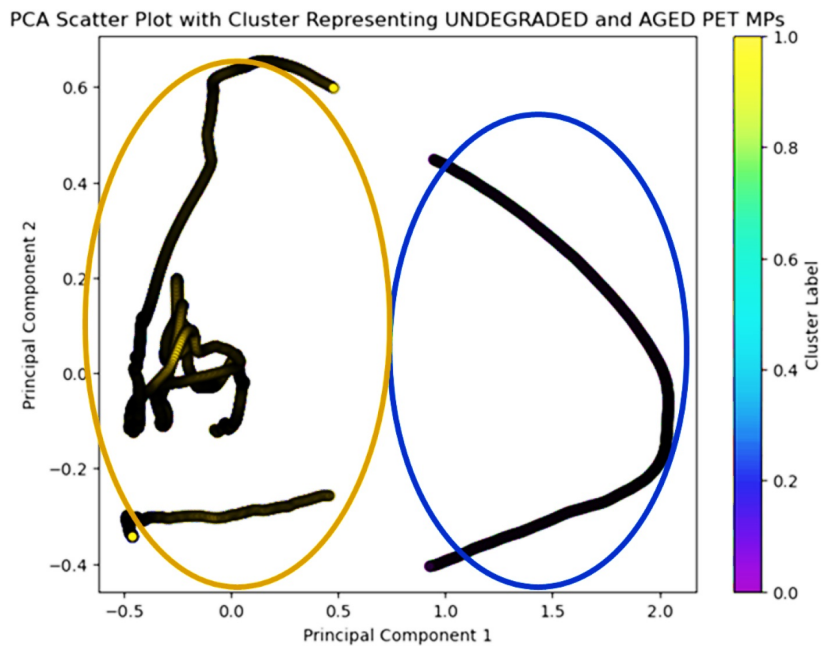
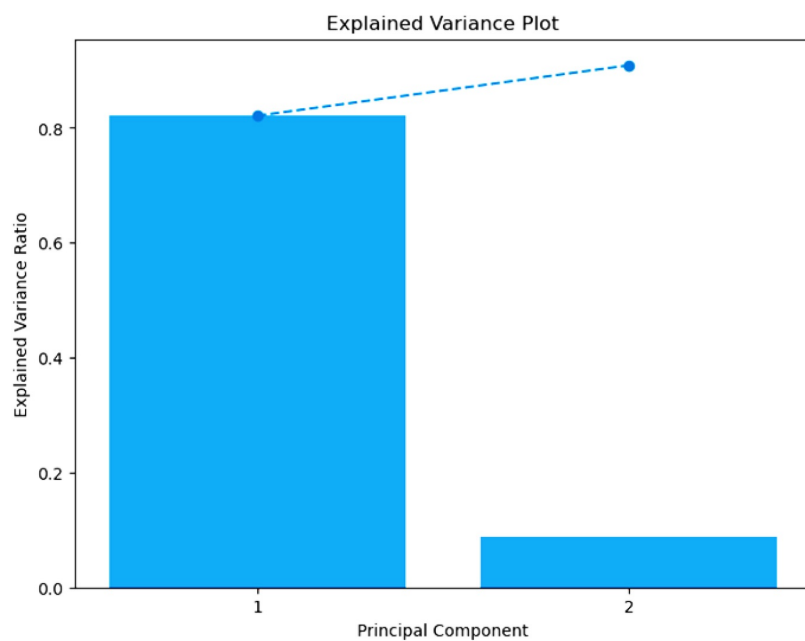


Figure 3

(a) Pre-processed spectra ATR-FTIR spectra of the undegraded and aged PET MPs, (b) the dataset used for PCA and (c) surface morphology of the degraded and aged PET MPs



(a)



(b)

Figure 4

(a) The PCA plots in space from the normalized spectral for undegraded and aged PET MPs and (b) the explained variance plot for the principal components extracted

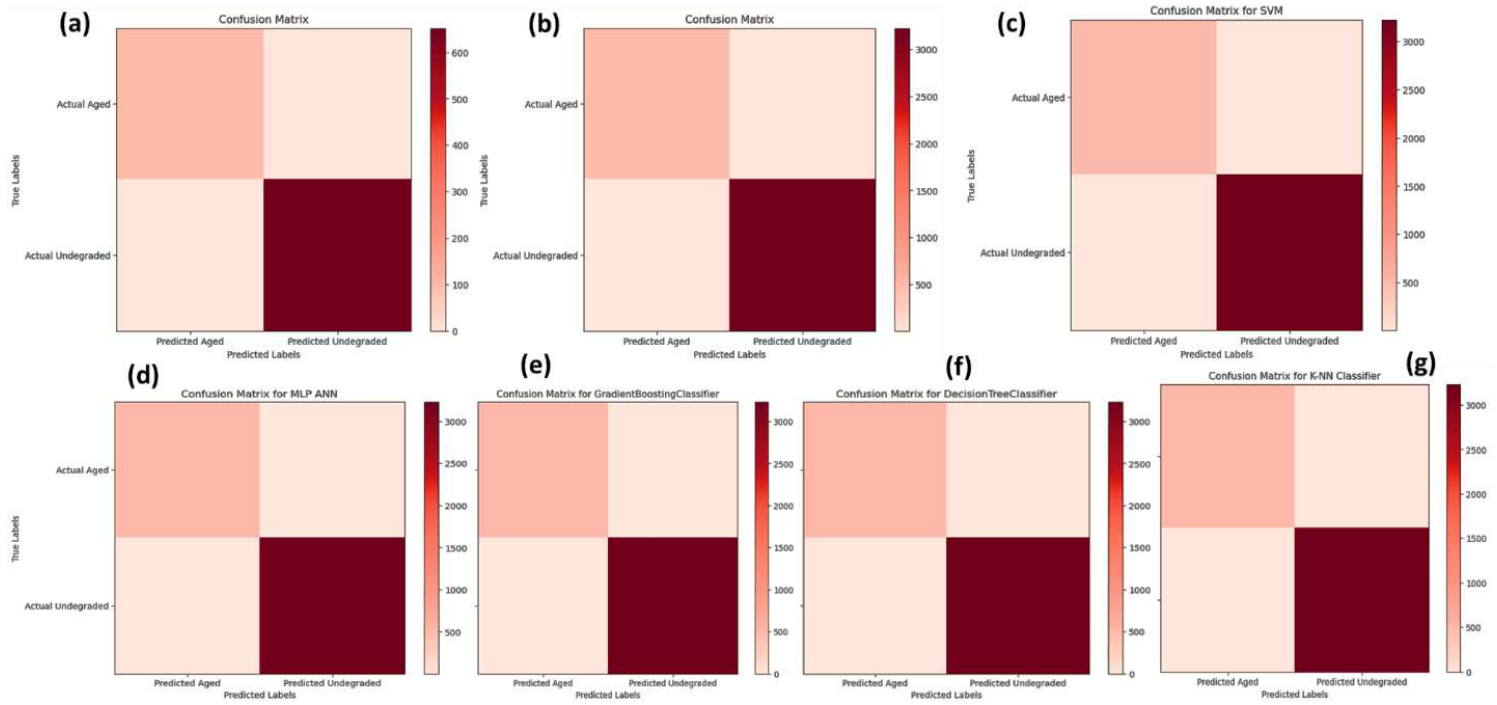
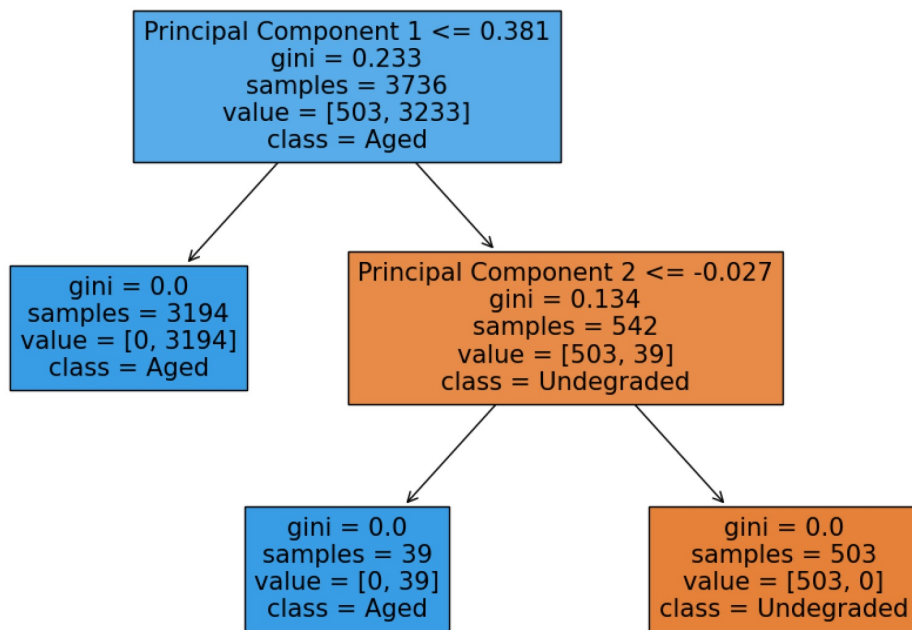


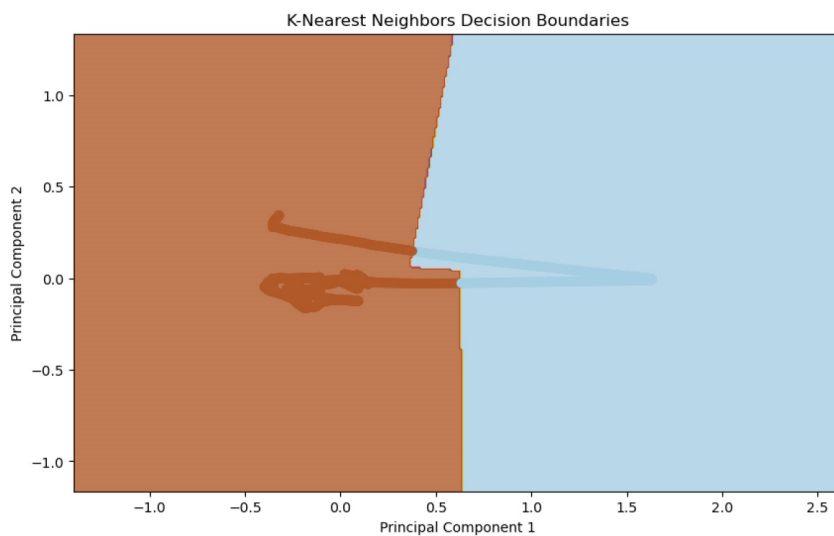
Figure 5

Confusion matrix plot for all ML classifiers (a) RF (b) LR (c) SVM (d) MLP (e) GB (f) DT (g) K-NN

Decision Tree Visualization



(a)



(b)

Figure 6

Visualization of (a) DT and (b) k-NN for classifying undegraded and aged PET MPs using ATR-FTIR spectral

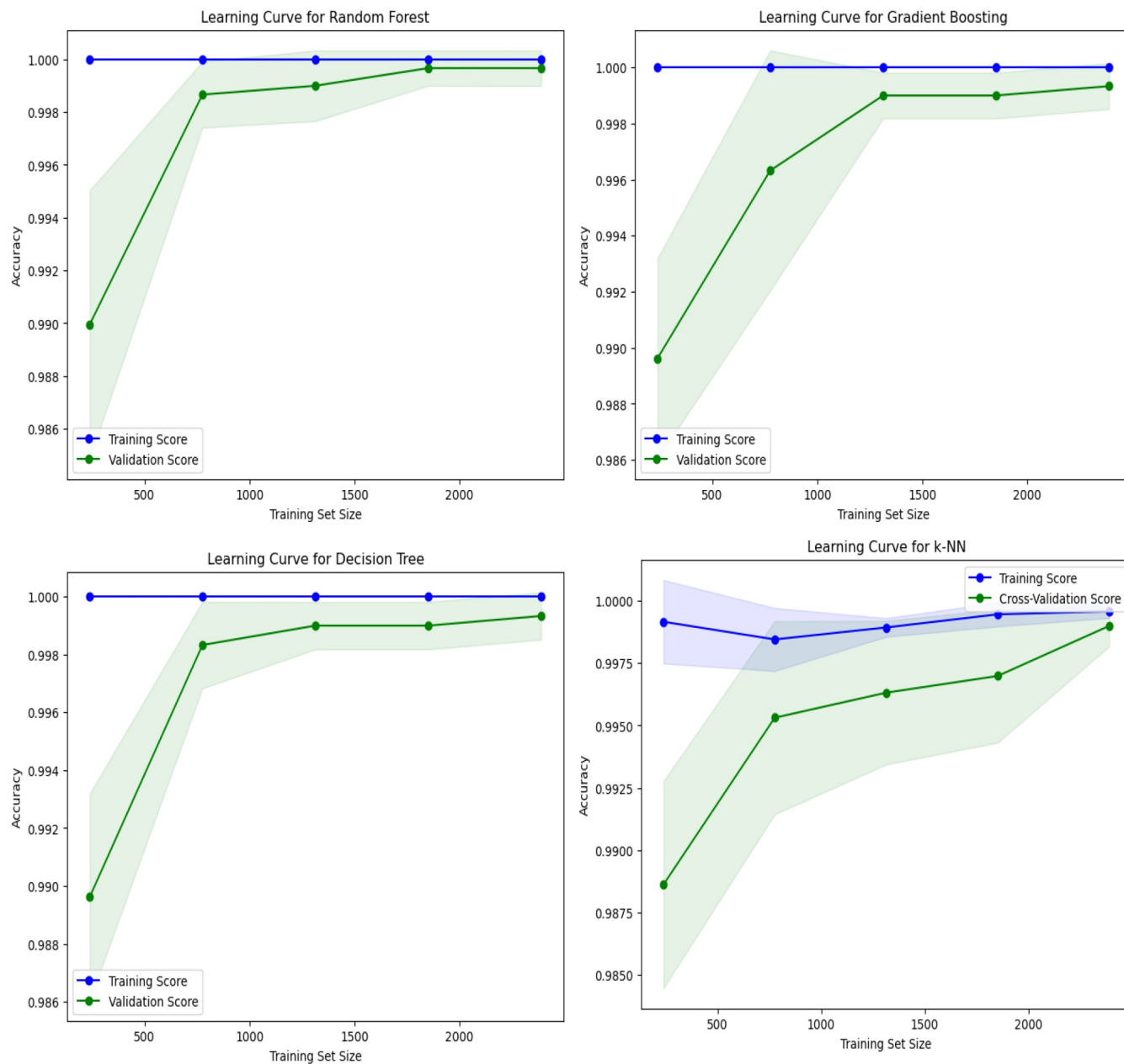


Figure 7

Learning curves for the best ML algorithms for the classification of undegraded and aged PET MPs.