

pubs.acs.org/JPCA

Article

¹ Machine Learning Identification of Organic Compounds Using ² Visible Light

³ Published as part of The Journal of Physical Chemistry virtual special issue "Early-Career and Emerging ⁴ Researchers in Physical Chemistry Volume 2".

s Thulasi Bikku, Rubén A. Fritz, Yamil J. Colón, and Felipe Herrera*



6 ABSTRACT: Identifying chemical compounds is essential in 7 several areas of science and engineering. Laser-based techniques 8 are promising for autonomous compound detection because the 9 optical response of materials encodes enough electronic and 10 vibrational information for remote chemical identification. This has 11 been exploited using the fingerprint region of infrared absorption 12 spectra, which involves a dense set of absorption peaks that are 13 unique to individual molecules, thus facilitating chemical 14 identification. However, optical identification using visible light 15 has not been realized. Using decades of experimental refractive 16 index data in the scientific literature of pure organic compounds 17 and polymers over a broad range of frequencies from the ultraviolet 18 to the far-infrared, we develop a machine learning classifier that can



19 accurately identify organic species based on a single-wavelength dispersive measurement in the visible spectral region, away from 20 absorption resonances. The optical classifier proposed here could be applied to autonomous material identification protocols and 21 applications.

I. INTRODUCTION

²² Scientific data analysis has been accelerated by machine ²³ learning by training models that allow a rapid interpretation of ²⁴ complex data patterns and the automated control of measure-²⁵ ment devices. ^{1,2} The autonomous identification and discovery ²⁶ of chemical compounds for applications in science and ²⁷ industry can therefore benefit from the development of ²⁸ compact, portable, and highly accurate sensors powered by ²⁹ machine learning.³ Remote molecular sensing based on light ³⁰ exploits the dispersive and absorptive response of a material ³¹ system to electromagnetic radiation.^{4,5} Although chemical ³² methods can be very specific,^{6,7} optical sensing techniques can ³³ be beneficial because light-matter interaction is nondestructive ³⁴ and can be processed remotely.

The refractive index is an optical property of materials fundamentally related to microscopic physicochemical characteristics such as the dynamic polarizability as well as macroscopic variables such as concentration, temperature, and pressure.^{8,9} Refractive indices are commonly used for uptroverse and pressure.^{8,9} Refractive indices are commonly used for an quantifying the content of target molecules in agricultural,^{10,11} therefore, a refractive index database over a broad range of frequencies as can serve as useful training data for a machine learning workflow that enables the autonomous identification of 44 chemical compounds using light. 45

Molecular spectroscopy databases have already been used 46 for training machine learning algorithms in chemical 47 identification problems.¹³⁻¹⁵ Efforts have focused on training $_{48}$ classifiers with infrared (IR) absorption and Raman scattering 49 databases, in the mid-infrared (mid-IR) spectral range ($\lambda \approx 3-50$ 50 μ m).^{11–32} In general, the information in the mid-IR is so 51 rich and complex that it would be very unlikely for different 52 molecules to have the same peak structure, particularly in the 53 so-called "fingerprint" region of the IR spectrum.³³⁻³⁵ 54 Compounds with the same chemical formula but different 55 spatial conformation (isomers) can thus be discriminated by 56 analyzing the position of their Raman peaks, for example.³⁶ Machine learning classifiers trained on vibrational spectroscopy 58 data can therefore be very accurate, with identification errors of 59 a few percent or less.^{13,16,28-32} 60

Received: November 11, 2022 Revised: February 3, 2023





Figure 1. Illustration of the proposed machine learning classification scheme for chemical identification. The left panel shows a sample of the refractive index $n(\lambda)$ and extinction $k(\lambda)$ spectra used for model training. The spectral data is public domain and can be downloaded at https://refractiveindex.info.⁵²

Alternative machine learning strategies for molecular 61 62 classification that are not trained with infrared absorption or 63 scattering spectra have also been reported.²⁸⁻³² Away from 64 absorption resonances, molecules and materials experience a 65 dispersive response in the presence of external electromagnetic 66 fields, which is determined by the real part of the dielectric 67 function of the medium $\epsilon(\lambda)$. The imaginary part of $\epsilon(\lambda)$ gives 68 the extinction coefficient $k(\lambda)^{1}$, which quantifies electro-69 magnetic energy loss due to absorption and scattering. 70 Absorption and dispersion are fundamentally related to each 71 other via the Kramer-Kronig relation.⁹ In frequency regions 72 with negligible extinction, the dispersive response is entirely determined by the refractive index $n(\lambda) = \sqrt{\text{Re}[\epsilon]}$, which 73 provides information on the electronic structure of materials (λ 74 $75 \approx 400-800$ nm) or their dielectric properties ($\lambda > 50 \ \mu$ m). As 76 mentioned above, the refractive index $n(\omega)$ is also an extensive property and thus correlates with the molecular density. 77

The refractive index has been used for detecting tissue 79 damage in biomedical samples using terahertz transmission 80 images ($\lambda > 100 \ \mu m$) and deep learning, reaching recall 81 performances of up to 93% with feature engineering 82 techniques.^{15,37} Databases with the static refractive indices ($\lambda = \infty$) of small organic compounds, solvents, and polymers 84 have also been developed for training regression models that 85 predict the refractive index given a set of structural and 86 quantum chemical descriptors.^{38–43} Predictive models can be 87 useful for the automated discovery of organic optoelectronic 88 materials.⁴⁴ Experimental databases with organic refractive 99 indices at the sodium wavelength (589 nm) have also been 90 used for constructing quantitative structure-property relation-91 ships (QSPR) that can be used, for example, in the chemical 92 analysis of organic mixtures using visible light.⁴⁵

⁹³ For molecular classification tasks based on IR spectral data, ⁹⁴ the target feature of interest for a classifier ("molecule ⁹⁵ identity") is learned from a high-dimensional feature vector ⁹⁶ that encodes the spectral peaks of interest.¹³ In contrast, the ⁹⁷ static refractive index of organic compounds is a single-valued ⁹⁸ feature that does not have the same information content.⁴⁶ ⁹⁹ The same applies to refractive index databases at a single wavelength away from ultraviolet and infrared absorption 100 resonances. New databases and analysis techniques are thus 101 needed for applications that can benefit from accurate organic 102 molecule classification protocols based on the dispersive 103 response of materials over a tunable range of wavelengths in 104 the visible. 105

Motivated by the growing interest in developing autono- 106 mous tools for organic materials discovery, $^{47-49}$ here we 107 demonstrate a machine learning classifier of organic com- 108 pounds based on the measurement of the refractive index $n(\lambda_0)$ 109 at a single optical wavelength λ_0 anywhere in the visible 110 spectral region (400–750 nm), where most organic com- 111 pounds are fully transparent.⁵⁰ The classifier is trained with a 112 publicly available materials science database containing the 113 optical constants of 61 organic molecules and polymers over 114 the spectral range spanning from ultraviolet to far-infrared 115 wavelengths. The classification scheme is illustrated in Figure 116 fil 1. We envision the proposed classifier for data analysis at the 117 fil output of a liquid or gas-phase chromatograph that can 118 separate complex chemical mixtures into multiple single- 119 component fractions.⁵¹

The rest of the article is organized as follows: After 121 describing the structure and information content of the original 122 spectral database in Sec. II.1, we describe the data 123 preprocessing strategies in Sec. II.2 and demonstrate 124 classification errors smaller than 1% in the visible range in 125 Sec. III. Comparisons with recent Raman-based classifiers are 126 given in Sec. III.4. We conclude and suggest future directions 127 and applications in Sec. IV. 128

II. METHODS

II.1. Experimental Refractive Index Database. We do 129 web scraping on a public domain database of experimental 130 optical constants available at https://refractiveindex.info.⁵² 131 The Web site is a repository of published data from the 132 scientific literature since 1940. The site is organized into 133 categories that resemble a virtual bookshelf: The "Shelf" 134 category groups materials into inorganic, organic, glasses, 135 others, and 3D; the "Book" category contains the chemical 136



Figure 2. Distribution of the raw data from the https://refractiveindex.info/⁵² over a broad range of wavelengths and refractive indices. Side panels show cuts of the data record distribution over wavelengths (top) and refractive indices (right).



Figure 3. Number of data records per molecular compound (61 compound class labels) available in the raw database from ref 52. The record distribution after data augmentation via Sellmeier fitting in the UV/vis is also shown. The correspondence between the class labels and organic compound names is given Table S6 of the Supporting Information.

137 compound name, which will be the "class" label predicted by 138 the machine learning classifier on output; the "Page" subsection refers to the source where the optical data was 139 140 first reported in the literature as well as comments and other 141 information such as the group velocity and group velocity 142 dispersion, the measurement wavelength range, and the state of 143 matter of the sample (gas, liquid, or solid). Each "Page" record 144 has a csv file with the spectral data for n (refractive index) and 145 k (extinction coefficient) over a range of wavelengths (λ). 146 From the general data set, we build a smaller set by selecting 'Organic Materials" in the "Shelf" category, which contains 61 147 148 organic molecules and polymers. The compiled file has 149 194 816 data records, sorted in columns with the features 150 "Shelf", "Book", "Page", " λ " (wavelength), "n" (refractive 151 index), and "k" (extinction coefficient). The data has 418 152 missing n values and 60 944 missing k values. Example data

records are shown for the "Acetone" class in Table S1 of the 153 Supporting Information (SM).

In Figure 2, we show a visualization of the organic data set 155 f2 over a grid of wavelengths (λ) and refractive indices (n). The 156 color code indicates the number of data records in each (λ , n) 157 region [The records involve measurements up to 25 μ m and 158 refractive indices in the range (0.0–2.3)]. The upper and right- 159 side panels show the number of organic compounds (counted 160 as a label on the "Book" category) with data records in each 161 wavelength and refractive index range. 162

II.2. Preprocessing the Training Data Set. The original 163 database is naturally heterogeneous and imbalanced, as some 164 compounds and frequency regions have been studied more 165 intensely than others in the literature. This represents a 166 challenge for the implementation of machine learning 167 classification algorithms.⁵³ Data sets with unequal data records 168

f3

169 per target feature lead to majority and minority classes, thus 170 affecting the overall predictive accuracy of classification 171 models.⁵⁴ Figure 3 shows the number of data records per 172 organic molecule class for the raw database (RD) and after 173 preprocessing through data augmentation (DA). The number 174 of records in the original data varies from only a few spectral 175 points (e.g., Polymethyl Pentene, Class Label = 58) to a few 176 thousand (e.g. Ethanol, Class Label = 1), which illustrates the 177 class imbalance problem. In addition, different compounds in 178 the refractive index database can have different types of 179 spectroscopic features over different wavelengths, which is a 180 form of type imbalance.⁵⁵

We tried several preprocessing strategies to overcome class and type imbalance in the raw data set, including oversampling (OS), undersampling (US), and physics-based DA. We also used feature engineering (FE-1 and FE-2) strategies, spectralbased binning (SBB), and spectral-based binning with feature engineering (SBB-FE1 and SBB-FE2), attempting to increase the prediction accuracies with the original imbalanced data set. The details of these preprocessing strategies are given in the Results section.

II.3. Random Forest Classifier. The random forest (RF) 190 191 algorithm was chosen as the default method for classification. 192 Random forest is a supervised machine learning algorithm 193 based on an ensemble of decision trees.⁵⁶ It uses bagging and 194 feature randomness when building each individual decision 195 tree to create an uncorrelated forest of trees whose prediction 196 is more accurate by committee than individually. The data set 197 was divided into 75% for training and 25% for testing. The results below were obtained with a single instance of the split, 198 199 but due to the intrinsic overfitting evidenced by the relatively 200 large difference between our training and testing accuracies, we 201 carried out additional cross-validation tests over random splits. 202 The cross-validated accuracies do not vary significantly among 203 different splits (see Table S3 in the SM). RF was implemented 204 using Python's Scikit-Learn library with default hyper-205 parameters.⁵⁷ In early stages of this study, we obtained 206 classification accuracies with alternative models such as 207 gradient boosting, support vector classification, and logistic 208 regression on the raw data set, and random forest performed 209 best (see test results Table S5 in the SI). The code and the 210 data sets used in this work are publicly available at https:// 211 github.com/fherreralab/organic optical classifier.

III. RESULTS AND DISCUSSION

212 We now discuss the random forest classification performance 213 obtained by training with the original imbalanced data set and 214 after preprocessing the data. We tested several preprocessing 215 techniques to gain insight on the optimal experimental setup 216 necessary for achieving reasonable identification outcomes 217 using minimal optical measurements.

III.1. Classification Accuracies with Imbalanced Class 219 **Sets.** Incorporating domain knowledge into the features of the 220 training set is known to increase the prediction performance of 221 classification models. Based on this intuition and inspired by 222 the feature structure of Raman-based classifiers^{18,24,26} where 223 the target class is associated with a high-dimensional feature 224 vector containing the spectral peaks, we tested whether 225 grouping small fractions of the curves $n(\lambda)$ and $k(\lambda)$ into 226 feature vectors $\mathbf{x} = [\lambda_j, n_j, k_j]$ improved the classifier 227 performance. As explained above, the original data set has a 228 three-dimensional feature vector per target class (j = 1)229 representing a single evaluation of the *n* and *k* curves at a given value of λ . We then build six-dimensional (j = 1,2) and nine- 230 dimensional feature vectors (j = 1,2,3) containing the 231 information from two and three consecutive points in the n 232 and k curves. We refer to the six- and nine-dimensional feature 233 vector schemes as feature engineering 1 (FE1) and 2 (FE2), 234 respectively. As a result of the increase in the number of 235 features per target class, the number of records is reduced. For 236 classification problems with large data sets, binning strategies 237 can prove useful for improving the overall class prediction 238 accuracy.^{55,56,58} In addition to the FE1 and FE2 strategies, we 239 adopt a spectral-based binning (SBB) strategy based on the 240 division of the wavelength domain into five spectroscopic 241 regions: UV [$\lambda < 0.40 \ \mu m$] containing 1773 data records, 242 Visible (VIS) $[0.40 < \lambda < 0.75 \ \mu m]$ with 5979 records; Near- 243 Infrared (Near-IR) $[0.75 < \lambda < 1.50 \ \mu m]$ with 35 445 records; 244 Infrared (IR) $[1.50 < \lambda < 4.0 \ \mu m]$ with 135 407 records; and 245 Far-Infrared (Far-IR) $[\lambda > 4.0 \ \mu m]$ with 66 678 records. Our 246 splitting of the IR spectrum into subregions is not necessarily 247 standard.5 248

In Figure 4, we compare the training and testing accuracies $_{249}$ f4 of the random forest classification training with the raw $_{250}$



Figure 4. Overall training and testing accuracy for imbalanced data, not separated by wavelength range for the different data preprocessing strategies.

database (RD) as well as FE and SBB preprocessing strategies. 251 Accuracies are also provided in Tables S3 and S4 of the SM. 252 The overall testing accuracies are about 80% with and without 253 preprocessing, although feature engineering and spectral-based 254 binning tend to reduce the prediction accuracy. Combining FE 255 and SBB did not improve performance relative to the raw data. 256 The low performance of the classifier over the entire range of 257 wavelengths (UV to Far-IR) comes roughly speaking from an 258 average of spectral regions of very high accuracies (IR) and 259 regions with very low accuracies (UV/vis). In what follows, we 260 separately studied the performance in different spectral regions. 261

III.2. Addressing Class Imbalance in the Database. In 262 Figures 3 and 5, we illustrate that the original data set contains 263 f5 a disproportionate number of data records per organic 264 compound in the IR spectral bin ($1.50-4.0 \ \mu m$), relative to 265 the UV and VIS bins. This type of record distribution 266 generates a class imbalance problem⁶⁰ that we address using 267 the following strategies: undersampling (US), oversampling 268 (OS), and data augmentation (DA). 269

Resampling strategies aim to balance classes in the training 270 data by reshaping the data set such that the numbers of records 271 in the different classes becomes comparable.⁶⁰ Undersampling 272



Figure 5. Distribution of data records per spectral bin in the raw database (RD), after oversampling (OS), undersampling (US), and data augmentation (DA).

273 is a method that randomly reduces the number of data records 274 in the majority class while oversampling duplicates randomly 275 chosen records in the minority class. Both resampling 276 strategies can be effective when used independently or 277 combined. Figure 5 shows the record distribution used for 278 training the random forest model after OS and US are carried 279 out. In the resampled data set, the number of records in the 280 UV/vis is comparable with the IR.

In addition to resampling (OS and US), we augment the data set using physics-based modeling. Specifically, we generate refractive index data using the Sellmeier equation⁹

$$n^{2}(\lambda) = A + \frac{B_{1}\lambda^{2}}{\lambda^{2} - C_{1}} + \frac{B_{2}\lambda^{2}}{\lambda^{2} - C_{2}}$$
(1)

28

f6

285 where $(A_i B_i C_i)$ are phenomenological coefficients. For each 286 molecule in the original database that has refractive index 287 measurements in the UV/visible, we fit the experimental curves 288 for n^2 using eq 1 to obtain Sellmeier coefficients. These are 289 then used to interpolate the index data in the region λ = 290 [200,750] nm. Therefore, we exclude the near-infrared 291 wavelengths from the fitting, despite the Sellmeier equation 292 still being valid in this region for most organic materials. This data-augmentation procedure increases the number of records 293 294 in the UV and VIS bins to about 3000 additional points per chemical compound. Figure 5 shows that the distribution of 295 296 training records set after DA changed significantly with respect to RD. The augmented data set has about 400 000 records, 297 with comparable numbers of records in all spectral bins. DA 298 299 has also been used to resolve the imbalance problem in 300 Raman-based classification problems.^{28–32}

After addressing class imbalance in the raw data set, we 301 tested the prediction performance of the random forest 302 303 classifier on the different spectral bins. Figure 6 shows that the accuracies after US and OS improved significantly relative 304 to the imbalanced data set, reaching 97% for UV and 99% for 305 VIS regions (see also Table S4 in the SM). The accuracies after 306 US and OS however did not change significantly in the Near-307 308 IR (2% improvement) and IR (0% improvement) bins. There 309 is a decrease of 25% in accuracy for Far-IR, likely due to the 310 reduction of data records on this region during undersampling. 311 On the other hand, data augmentation (DA) significantly 312 improved the performance of the classifier (99% for UV and 313 98% for VIS) without affecting the accuracies in the infrared 314 bins (see Table S4 in the SM). The accuracies for FE1 and 315 FE2 significantly decreased as expected. The reduced number



Figure 6. Testing accuracy after spectral-based binning with different data preprocessing strategies: Feature engineering (FE1 and FE2), oversampling (OS), undersampling (US), and data augmentation (DA). The accuracies of the raw database are also shown for comparison.

of data records resulting from the grouping of the features 316 dramatically affects the model performance in the less 317 represented spectral bins. 318

The results in Figure 6 suggest that measuring the refractive 319 index in the UV/vis region is in principle sufficient for a precise 320 classification of chemical compounds. In comparison with the 321 accuracies reached in the far-infrared spectrum ($\lambda > 5 \ \mu$ m, 322 fingerprint region), the classification accuracies in the visible 323 region after data augmentation are equally good, without 324 additional steps of hyperparameter optimization for the 325 random forest model.

III.3. Estimating the Required Measurement Precision. The precision with which the refractive index is 328 measured in experiments should affect the classifier perform-329 ance. Compounds with similar refractive indices at a given 330 wavelength in the visible would not be reliably distinguished if 331 the index measurement does not have enough significant digits. 332 We quantify this intuition by computing the classification 333 accuracy achieved by the random forest model using training 334 data whose precision was manually truncated to a finite 335 number of significant figures. In Figure 7, we show the testing 336 f7 accuracies obtained using different decimal digits for $n(\lambda)$ in 337 the training set, separating the results by preprocessing 338 strategy. While the classification accuracies are in general 339 poor (<82%) for single-wavelength index measurements with 340



Figure 7. Testing accuracies for refractive index measurements with different precision. Results are shown for different preprocessing strategies: Featured engineering (FE1 and FE2), oversampling (OS), undersampling (US), and data augmentation (DA). The accuracies for the raw database (RD) are also shown for comparison.

Table 1. Prediction Accuracies U	Using Raman Spectra	l Databases"
----------------------------------	---------------------	--------------

Method	Spectral Database	Testing Accuracies (%)	Data Augmentation	Ref.
1D-CNN	Minerals and organic compounds	100	Standard augmentation transformations in vibrational spectroscopy	31
CRL	72 organic compounds	97.5	Gaussian noise and linear combination	28
DNN	72 organic compounds	92.6* 96.4**	Shifting, Gaussian noise, and interpolation	30
CNN	72 organic compounds	81.9* 86.0**	Shifting, Gaussian noise, and interpolation	30
DRCNN	72 organic compounds	98.1	Shifting and Gaussian noise	29
1D-CNN and KNN	620 mineral and 211 synthetic organic pigments	97.38	Shifting, Gaussian noise, Savitzky-Golay smoothing, spline interpolation, and polynomial reconstruction	32
RF	61 organic compounds	99 (UV) 98.1 (vis) 99.2 (Near-IR) 83.1 (IR) 94.8 (Far-IR)	Sellmeier equation fitting on UV/vis optical ranges	This work

^{*a*}The number of molecules present in the database is specified when the data is available. We also list the data-augmentation methodologies used in these works. Model acronyms: 1D-CNN = 1D convolutional neural network, CRL = contrastive representation learning, DNN = deep neural network, CNN = convolutional neural network, DRCNN = deeply recursive convolutional neural network, KNN = K-nearest neighbor classifier, RF = random forest. *Without transfer learning, **with transfer learning.

341 two significant figures or less, the raw database already gives 342 better classification accuracies (~95%) when trained with 343 index data of at least three decimal places. Increasing the 344 number of significant figures beyond four decimals does not 345 significantly improve the accuracy, regardless of the prepro-346 cessing strategy used.

These high testing accuracies for single-wavelength measure-347 348 ments are only possible because the trained random forest 349 model has learned the molecular spectra over a broad range of 350 wavelengths spanning hundreds of nanometers. We confirm 351 this by training the same random forest model described above 352 (see Methods Section) but with a reduced training set with 353 only 89 refractive index records at 1000 nm (near IR). Only 42 354 compounds have measurements reported around this target 355 wavelength. The testing accuracies obtained with such a 356 restricted one-wavelength data set are very poor (30% or less, 357 see Table S3 in the SM). This confirms the limited 358 classification ability expected for conventional single-point 359 index measurements. Our work extends this intuition by 360 exploiting the correlations learned by the decision trees in the 361 refractive index values at different wavelengths. Figure 7 shows 362 that these correlations are more effectively learned when the 363 index measurements have interferometric precision (4 digits or 364 more).

III.4. Comparison with Raman Classifiers. Class 365 366 imbalance has also been addressed in other molecular 367 classification studies based on Raman spectral databases, which also tend to be heterogeneous with a data record 368 $_{369}$ distribution that over-represents compounds of particular $_{370}$ interest in chemistry. $^{28-32}$ To address class imbalance, 371 researchers have explored data-augmentation strategies such 372 as peak shifting, noise addition, smoothing, spline interpola-373 tion, and polynomial reconstruction. These preprocessing 374 strategies were implemented to reconstruct the data set before 375 training deep learning classifier models. In our analysis of the 376 refractive index.info database, the Sellmeier fitting procedure is 377 a valid augmentation strategy that can be used to reconstruct 378 the part of the training set corresponding to ultraviolet and 379 visible wavelengths, without introducing data leakage.⁹ In 380 Table 1, we compare the refractive index approach with the 381 recent Raman-based classifiers in the literature, showing that the dispersive method can also give high classification 382 performance, using a similar a preprocessing strategy 383 (interpolation) on data sets with a comparable volume of 384 data records and molecular classes.^{28–32} 385

IV. CONCLUSIONS

We have built a machine learning classification scheme to 386 identify organic compounds based on refractive index 387 measurements in the visible spectral region, in which most 388 organic compounds are highly transparent. We trained a 389 random forest classifier using decades of experimental data 390 from the scientific literature. The database contains 194 816 391 spectral records of refractive index and extinction curves of 61 392 organic compounds and polymers over a broad range of 393 wavelengths from the UV to the far-infrared. There is a class 394 imbalance problem in the experimental data that restricts the 395 classification accuracy for refractive index inputs at visible 396 wavelengths (400-750 nm) to approximately 80%. This 397 imbalance is primarily due to the disproportionate number of 398 infrared absorption records reported in the mid- and far- 399 infrared regions. Imbalance is a common problem when 400 working with spectroscopic databases as the experimental data 401 is deposited from different sources.²⁸⁻³² 402

We addressed this class imbalance issue by preprocessing the 403 raw data before training the classifier using resampling and 404 physics-based data-augmentation strategies analogous to those 405 employed by other machine learning classifiers trained with 406 Raman spectra.²⁸⁻³² By training the random forest model with 407 preprocessed balanced data, we achieve molecular classification 408 testing accuracies in the UV and visible regions better than 409 98%. Additional improvements can be expected with additional 410 steps of model hyperparameter optimization. Such high 411 accuracies are comparable to those obtained using only 412 Raman spectroscopy databases (see Table 1), thus demon- 413 strating the feasibility of using machine learning tools for 414 enhancing the capabilities of laser-based chemical sensing 415 devices. The high classification precision reached via data 416 augmentation for chemical identification suggests that a similar 417 technique would be applicable with other problems that 418 involve continuous-variable data sets that can be interpolated 419 using physics models. Further research is needed to generalize 420 421 the proposed molecular classifier to identify the structural and 422 other chemical features of the molecules that are present in the 423 Refractive Index Database.⁵² Our work thus serves as a starting 424 point for the development of remote chemical sensors based 425 on laser light.

426 **ASSOCIATED CONTENT**

427 Data Availability Statement

428 The data that support the findings of this study are openly 429 available at 10.5281/zenodo.6419970.⁶¹

430 **Supporting Information**

431 The Supporting Information is available free of charge at 432 https://pubs.acs.org/doi/10.1021/acs.jpca.2c07955.

List of organic compound classes in the training
database, additional random forest classification tests, a
visualization of the trained trees, and classification

436 metrics using other machine learning models (PDF)

437 **AUTHOR INFORMATION**

438 Corresponding Author

439 Felipe Herrera – Department of Physics, Universidad de

440 Santiago de Chile, 3493 Santiago, Chile; Millennium

441 Institute for Research in Optics, https://www.miroptics.cl/

442 eng/; ^(b) orcid.org/0000-0001-8121-1931

443 Authors

444 Thulasi Bikku – Department of Physics, Universidad de

445 Santiago de Chile, 3493 Santiago, Chile; Computer Science

446 and Engineering, Vignan's Nirula Institute of Technology and

447 Science for Women, Guntur, Andhra Pradesh 522009, India

448 **Rubén A. Fritz** – Department of Physics, Universidad de

449 Santiago de Chile, 3493 Santiago, Chile

450 Yamil J. Colón – Department of Chemical and Biomolecular

451 Engineering, University of Notre Dame, Notre Dame, Indiana

452 46556, United States;
^(a) orcid.org/0000-0001-5316-9692

453 Complete contact information is available at:

454 https://pubs.acs.org/10.1021/acs.jpca.2c07955

455 Notes

456 The authors declare no competing financial interest.

457 **ACKNOWLEDGMENTS**

458 R.A.F. is supported by DICYT-USACH grant POSTDOC 459 USA1956_DICYT. F.H. and T.B. are supported by ANID 460 through grants FONDECYT Regular No. 1181743 and 461 Millennium Science Initiative Program ICN17_012. Y.J.C. 462 thanks the University of Notre Dame for financial support 463 through start-up funds.

464 **ADDITIONAL NOTE**

⁴⁶⁵ ¹We follow photonics notation and denote the extinction ⁴⁶⁶ coefficient by k.

467 **REFERENCES**

468 (1) Doucet, M.; Archibald, R.; Heller, W. T. Machine Learning for
469 Neutron Reflectometry Data Analysis of Two-Layer Thin Films.
470 Machine Learning: Science and Technology 2021, 2, 035001.

471 (2) Ratner, D.; Sumpter, B.; Alexander, F.; Billings, J. J.; Coffee, R.; 472 Cousineau, S.; Denes, P.; Doucet, M.; Foster, I.; Hexemer, A.; et al. 473 Office of Basic Energy Sciences (BES) Roundtable on Producing and 474 Managing Large Scientific Data with Artificial Intelligence and Machine Learning, 2019. https://www.osti.gov/biblio/1630823/ (accessed 475 November 15, 2022). 476

(3) Ballard, Z.; Brown, C.; Madni, A. M.; Ozcan, A. Machine 477 Learning and Computation-Enabled Intelligent Sensor Design. *Nature* 478 *Machine Intelligence* 2021, 3 (7), 556–565. 479

(4) Wang, X.; Wolfbeis, O. S. Fiber-Optic Chemical Sensors and 480 Biosensors (2015–2019. Anal. Chem. 2020, 92 (1), 397–430. 481

(5) McDonagh, C.; Burke, C. S.; MacCraith, B. D. Optical Chemical 482 Sensors. *Chem. Rev.* **2008**, *108* (2), 400–422. 483

(6) Gründler, P.. Chemical Sensors: An Introduction for Scientists and 484 Engineers; Springer: Berlin, 2007. 485

(7) Swager, T. M.; Mirica, K. A. Introduction: Chemical Sensors. 486 Chem. Rev. 2019, 119 (1), 1–2. 487

(8) Mohan, S.; Kato, E.; Drennen, J. K.; Anderson, C. A. Refractive 488 Index Measurement of Pharmaceutical Solids: A Review of Measurement Methods and Pharmaceutical Applications. *J. Pharm. Sci.* **2019**, 490 *108* (11), 3478–3495. 491

(9) Boyd, R. W. Nonlinear Optics, 4th ed.; Academic Press: 492 Amsterdam, 2020. 493

(10) Pereira, F. M. V.; de Souza Carvalho, A.; Cabeça, L. F.; 494 Colnago, L. A. Classification of Intact Fresh Plums according to 495 Sweetness Using Time-Domain Nuclear Magnetic Resonance and 496 Chemometrics. *Microchemical Journal* **2013**, *108*, 14–17. 497

(11) Magwaza, L. S.; Opara, U. L. Analytical Methods for 498
 Determination of Sugars and Sweetness of Horticultural Products— 499
 a Review. Scientia Horticulturae 2015, 184, 179–192. 500

(12) Choudhary, V.; Rönnow, D. A Nondestructive Testing Method 501 for the Determination of the Complex Refractive Index Using Ultra 502 Wideband Radar in Industrial Applications. *Sensors* **2020**, *20*, 3161. 503

(13) Madden, M. G.; Howley, T. A Machine Learning Application 504 for Classification of Chemical Spectra. *Applications and Innovations in* 505 *Intelligent Systems XVI* **2009**, 77–90. 506

(14) Park, H.; Son, J.-H. Machine Learning Techniques for THz 507 Imaging and Time-Domain Spectroscopy. *Sensors* **2021**, *21* (4), 1186. 508 (15) Cao, Y.; Huang, P.; Chen, J.; Ge, W.; Hou, D.; Zhang, G. 509 Qualitative and Quantitative Detection of Liver Injury with Terahertz 510 Time-Domain Spectroscopy. *Biomedical Optics Express* **2020**, *11* (2), 511 982.

(16) Madden, M. G.; Ryder, A. G. Machine Learning Methods for 513
 Quantitative Analysis of Raman Spectroscopy Data. Proc. SPIE 4876, 514
 Opto-Ireland 2002: Optics and Photonics Technologies and Applications 515
 2003. 516

(17) Zhang, L.; Li, C.; Peng, D.; Yi, X.; He, S.; Liu, F.; Zheng, X.; 517 Huang, W. E.; Zhao, L.; Huang, X. Raman Spectroscopy and Machine 518 Learning for the Classification of Breast Cancers. *Spectrochimica Acta* 519 *Part A: Molecular and Biomolecular Spectroscopy* **2022**, 264, 120300. 520

(18) Ryzhikova, E.; Ralbovsky, N. M.; Sikirzhytski, V.; Kazakov, O.; 521 Halamkova, L.; Quinn, J.; Zimmerman, E. A.; Lednev, I. K. Raman 522 Spectroscopy and Machine Learning for Biomedical Applications: 523 Alzheimer's Disease Diagnosis Based on the Analysis of Cerebrospinal 524 Fluid. Spectrochimica Acta Part A: Molecular and Biomolecular 525 Spectroscopy **2021**, 248, 119188. 526

(19) Cui, A.; Jiang, K.; Jiang, M.; Shang, L.; Zhu, L.; Hu, Z.; Xu, G.; 527 Chu, J. Decoding Phases of Matter by Machine-Learning Raman 528 Spectroscopy. *Physical Review Applied* **2019**, *12* (5), 054049. 529

(20) Guo, S.; Popp, J.; Bocklitz, T. Chemometric Analysis in Raman 530 Spectroscopy from Experimental Design to Machine Learning–Based 531 Modeling. *Nat. Protoc.* **2021**, *16* (12), 5426–5459. 532

(21) Martinez, J. C.; Guzmán-Sepúlveda, J. R.; Bolañoz Evia, G. R.; 533 Córdova, T.; Guzmán-Cabrera, R. Enhanced Quality Control in 534 Pharmaceutical Applications by Combining Raman Spectroscopy and 535 Machine Learning Techniques. *Int. J. Thermophys.* **2018**, 39 (6), 79. 536

(22) Berghian-Grosan, C.; Magdas, D. A. Application of Raman 537 Spectroscopy and Machine Learning Algorithms for Fruit Distillates 538 Discrimination. *Sci. Rep.* **2020**, *10* (1), 21152. 539

(23) Zhang, Y.; Cong, Q.; Xie, Y.; JingxiuYang; Zhao, B. 540 Quantitative Analysis of Routine Chemical Constituents in Tobacco 541 by Near-Infrared Spectroscopy and Support Vector Machine. 542 543 Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy 544 **2008**, 71 (4), 1408–1413.

- 545 (24) Hu, W.; Ye, S.; Zhang, Y.; Li, T.; Zhang, G.; Luo, Y.; Mukamel, 546 S.; Jiang, J. Machine Learning Protocol for Surface-Enhanced Raman 547 Spectroscopy. *J. Phys. Chem. Lett.* **2019**, *10* (20), 6026–6031.
- (25) Ralbovsky, N. M.; Lednev, I. K. Towards Development of a

549 Novel Universal Medical Diagnostic Method: Raman Spectroscopy

550 and Machine Learning. Chem. Soc. Rev. **2020**, 49 (20), 7428–7453.

551 (26) Guevara, E.; Torres-Galván, J. C.; Ramírez-Elías, M. G.; 552 Luevano-Contreras, C.; González, F. J. Use of Raman Spectroscopy to

553 Screen Diabetes Mellitus with Machine Learning Tools. *Biomedical* 554 Optics Express **2018**, 9 (10), 4998–5010.

555 (27) Zou, T.; Dou, Y.; Mi, H.; Zou, J.; Ren, Y. Support Vector 566 Regression for Determination of Component of Compound Oxy-557 tetracycline Powder on Near-Infrared Spectroscopy. *Anal. Biochem.* 558 **2006**, 355 (1), 1–7.

559 (28) Li, B.; Schmidt, M. N.; Alstrøm, T. S. Raman Spectrum 560 Matching with Contrastive Representation Learning. *Analyst* **2022**, 561 147, 2238.

562 (29) Zhou, W.; Tang, Y.; Qian, Z.; Wang, J.; Guo, H. Deeply-563 Recursive Convolutional Neural Network for Raman Spectra 564 Identification. *RSC Adv.* **2022**, *12* (8), 5053–5061.

(30) Zhang, R.; Xie, H.; Cai, S.; Hu, Y.; Liu, G.; Hong, W.; Tian, Z.
566 Transfer-Learning-Based Raman Spectra Identification. *J. Raman*567 Spectrosc. 2020, 51 (1), 176–186.

(31) Hamed Mozaffari, M.; Tay, L.-L. Overfitting One-Dimensional
Convolutional Neural Networks for Raman Spectra Identification.
Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy
2022, 272, 120961.

572 (32) Shang, L.; Bao, Y.; Tang, J.; Ma, D.; Fu, J.; Zhao, Y.; Wang, X.;

573 Yin, J. A Novel Polynomial Reconstruction Algorithm-Based 1D 574 Convolutional Neural Network Used for Transfer Learning in Raman 575 Spectroscopy Application. J. Raman Spectrosc. **2022**, 53 (2), 237–246.

576 (33) López-López, M.; García-Ruiz, C. Infrared and Raman 577 Spectroscopy Techniques Applied to Identification of Explosives. 578 *TrAC Trends in Analytical Chemistry* **2014**, *54*, 36–44.

579 (34) Baker, M. J.; Byrne, H. J.; Chalmers, J.; Gardner, P.; Goodacre, 580 R.; Henderson, A.; Kazarian, S. G.; Martin, F. L.; Moger, J.; Stone, N.; 581 et al. Correction: Clinical Applications of Infrared and Raman 582 Spectroscopy: State of Play and Future Challenges. *Analyst* **2018**, *143* 583 (8), 1934–1934.

584 (35) Tatsch, E.; Schrader, B. Near-Infrared Fourier Transform 585 Raman Spectroscopy of Indigoids. *J. Raman Spectrosc.* **1995**, *26* (6), 586 467–473.

587 (36) Liu, C.; Xu, L.; He, H.; Jia, W.; Hua, Z. Discrimination of 588 Phenethylamine Regioisomers and Structural Analogues by Raman 589 Spectroscopy. *Journal of Forensic Sciences* **2021**, *66* (1), 365–374.

(37) Park, H.; Son, J.-H. Machine Learning Techniques for THz
Imaging and Time-Domain Spectroscopy. *Sensors* 2021, 21 (4), 1186.
(38) Katritzky, A. R.; Sild, S.; Karelson, M. Correlation and
Prediction of the Refractive Indices of Polymers by QSPR. *J. Chem.*Inf. Comput. Sci 1998, 38 (6), 1171–1176.

(39) Duchowicz, P. R.; Fioressi, S. E.; Bacelo, D. E.; Saavedra, L. M.;
Toropova, A. P.; Toropov, A. A. QSPR Studies on Refractive Indices
of Structurally Heterogeneous Polymers. *Chemometrics and Intelligent Laboratory Systems* 2015, 140, 86–91.

(40) Bouteloup, R.; Mathieu, D. Improved Model for the Refractive
Index: Application to Potential Components of Ambient Aerosol. *Phys. Chem. Chem. Phys.* 2018, 20 (34), 22017–22026.

602 (41) Erickson, M. E.; Ngongang, M.; Rasulev, B. A Refractive Index 603 Study of a Diverse Set of Polymeric Materials by QSPR with 604 Quantum-Chemical and Additive Descriptors. *Molecules* **2020**, 25 605 (17), 3772.

606 (42) Khan, P. M.; Rasulev, B.; Roy, K. QSPR Modeling of the 607 Refractive Index for Diverse Polymers Using 2D Descriptors. *ACS* 608 *Omega* **2018**, 3 (10), 13374–13386.

609 (43) Astray, G.; Cid, A.; Moldes, O.; Ferreiro-Lage, J. A.; Gálvez, J. 610 F.; Mejuto, J. C. Prediction of Refractive Index of Polymers Using Artificial Neural Networks. Journal of Chemical & Engineering Data 611 2010, 55 (11), 5388–5393. 612

(44) Liu, J.; Ueda, M. High Refractive Index Polymers: Fundamental 613 Research and Practical Applications. *J. Mater. Chem.* **2009**, *19* (47), 614 8907. 615

(45) Redmond, H.; Thompson, J. E. Evaluation of a Quantitative 616 Structure–Property Relationship (QSPR) for Predicting Mid-Visible 617 Refractive Index of Secondary Organic Aerosol (SOA. *Phys. Chem.* 618 *Chem. Phys.* **2011**, *13* (15), 6872. 619

(46) Chen, R.-C.; Dewi, C.; Huang, S.-W.; Caraka, R. E. Selecting 620 Critical Features for Data Classification Based on Machine Learning 621 Methods. *Journal of Big Data* **2020**, 7 (1), 52. 622

(47) Flores-Leonar, M. M.; Mejía-Mendoza, L. M.; Aguilar-Granda, 623 A.; Sanchez-Lengeling, B.; Tribukait, H.; Amador-Bedolla, C.; 624 Aspuru-Guzik, A. Materials Acceleration Platforms: On the Way to 625 Autonomous Experimentation. *Current Opinion in Green and* 626 *Sustainable Chemistry* **2020**, 25, 100370. 627

(48) Cai, W.; Abudurusuli, A.; Xie, C.; Tikhonov, E.; Li, J.; Pan, S.; 628 Yang, Z. Toward the Rational Design of Mid-Infrared Nonlinear 629 Optical Materials with Targeted Properties via a Multi-Level Data- 630 Driven Approach. *Adv. Funct. Mater.* **2022**, *32*, 2200231. 631

(49) Stein, H. S.; Gregoire, J. M. Progress and Prospects for 632 Accelerating Materials Science with Automated and Autonomous 633 Workflows. *Chemical Science* **2019**, *10* (42), 9640–9649. 634

(50) Rocha, F. S.; Gomes, A. J.; Lunardi, C. N.; Kaliaguine, S.; 635 Patience, G. S. Experimental Methods in Chemical Engineering: 636 Ultraviolet Visible Spectroscopy-UV-Vis. *Canadian Journal of* 637 *Chemical Engineering* **2018**, 96 (12), 2512–2517. 638

(51) Wade, J. H.; Bailey, R. C. Refractive Index-Based Detection of 639 Gradient Elution Liquid Chromatography Using Chip-Integrated 640 Microring Resonator Arrays. *Anal. Chem.* **2014**, *86* (1), 913–919. 641

(52) RefractiveIndex.INFO. Refractive index database. https:// 642 refractiveindex.info (accessed January 15, 2023). 643

(53) Fernández, A.; García, S.; Galar, M.; Prati, R. C.; Krawczyk, B.; 644 Herrera, F. *Learning from Imbalanced Data Sets*; Springer Science 645 +Business Media: New York, 2018. 646

(54) Cai, J.; Luo, J.; Wang, S.; Yang, S. Feature Selection in Machine 647 Learning: A New Perspective. *Neurocomputing* **2018**, 300, 70–79. 648

(55) More, A. S.; Rana, D. P. Review of random forest classification 649 techniques to resolve data imbalance. *Proceedings - 1st International* 650 *Conference on Intelligent Systems and Information Management, ICISIM* 651 2017, Aurangabad, India, 2017; pp 72–78. . 652

(56) Parmar, A.; Katariya, R.; Patel, V. A Review on Random Forest: 653 An Ensemble Classifier. *International Conference on Intelligent Data* 654 *Communication Technologies and Internet of Things (ICICI) 2018* **2019**, 655 26, 758–763. 656

(57) Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; 657 Thirion, B.; Grisel, O.; et al. Scikit-learn: Machine learning in Python. 658 *Journal of Machine Learning Research* **2011**, *12*, 2825–2830. 659

(58) Lustgarten, J. L.; Gopalakrishnan, V.; Grover, H.; Visweswaran, 660 S. Improving Classification Performance with Discretization on 661 Biomedical Datasets. *AMIA Annu Symp Proc.* **2008**, 2008, 445–449. 662

(59) NASA ipac. *Near, Mid and Far-Infrared*. https://archive.ph/ 663 20120529003352/http://www.ipac.caltech.edu/Outreach/Edu/ 664 Regions/irregions.html (accessed June 15, 2022). 665

(60) Ali, A.; Shamsuddin, S. M.; Ralescu, A. L. Classification with 666 Class Imbalance Problem: A Review. *Int. J. Adv. Soft Comput. Its Appl.* 667 **2015**, 7, 176. 668

(61) Bikku, T.; Fritz, R. A.; Colón, Y.; Herrera, F. *fherreralab*/ 669 organic_optical_classifier: v1.0.0, 2022. DOI: 10.5281/zeno- 670 do.6419970 (accessed November 20, 2022). 671