Modern Approaches in Classification of Covalent Organic Frameworks by Textural Properties Using JSL

Michael Nazarkovsky¹ and Felipe Lopes de Oliveira²

¹Chemistry Department, Pontifical Catholic University of Rio de Janeiro (DQ-PUC Rio), Rio de Janeiro, Brazil

²Institute of Chemistry, Federal University of Rio de Janeiro (IQ-UFRJ), Rio de Janeiro, Brazil

Introduction

Covalent organic frameworks (COFs) are an emerging class of materials composed only of light elements such as carbon, nitrogen, oxygen, hydrogen, and boron.^[1] COFs are designed in a bottom-up approach by the covalent bonding of one or more building blocks, following a well-defined geometric pattern dictated by the topology of the network resulting in organic, crystalline and nanoporous reticular material^[2] (Fig.1).



Fig.1. The sketch over the COFs structures: geometry, building blocks and the bonds types

There are various molecular geometries that can be used as building blocks and several types of chemical bonds that can be used to connect them.^[3] In this way, it is possible to form structures with extremely specific characteristics, such as pore size that can be controlled in the order of angstroms and the physicochemical characteristics of the pore interior, only adjusting the building blocks used. Therefore, it is possible to think in the building blocks as Lego pieces, where different pieces can be used to build several structures with unique characteristics (Fig.2).

Due to this incredible modulation capacity, COFs can be used for several extremely interesting applications such as heterogeneous catalysis,^[4–6] energy production, and storage, organic semiconductors,^[7] chemical sensing,^[8] thermal insulators,^[9] and gas capture and storage,^[10,11] among others^[12] (Fig.3).



Fig.2. The concept of the building blocks for 2D and 3D structures



Fig.3. The variety of the COFs application

However, this large chemical diversity can also be a problem. The development of these materials usually follows a combination of chemical intuition and the use of available molecules. With several steps of synthesis and characterization that are repeated until a material with high crystallinity is obtained. This makes its development very time-consuming, expensive, and requires several trained professionals. And often, even with all this effort, it is not possible to obtain a material with the desired characteristics.

To overcome this complex experimental development, we propose a simple and fast computational approach based on the textural properties such as specific surface area, pore size, pore volume, and dimensionality to study these materials in the aim to accelerate their development.

In the presentation reported for the Summit, the structures were classified by means of unsupervised (Hierarchical Clustering, PCA) and supervised (Multiple Logistic Regression, Naive Bayes, KNN, SVM) methods coded in JSL (JMP Pro 15). The COFs were separated in two main groups, named 2D and 3D, based on their geometrical characteristics. The largest clusters were selected to perform supervised machine

learning stratifying the data by the structure (62.9% 2D and 37.1% 3D) to avoid disproportion in training/validation ratio (80%/20%). The 2D/3D classification by textural properties was successfully accomplished with the 100% accuracy (validation) for all models. Other metrics, such as Entropy R2, Generalized R2, Mean -log(p) were compared and discussed to select the optimal model. It was shown in the presentation to the present Summit 2021.

The paper basically focuses on the treatment of the whole dataset – it means that the data will pass the path from the cleaning until tuning the models. The presentation rested in the extraction of the clusters utilizing the typical machine learning for 100% Accuracy. Herewith, we claim the role of imbalanced output and how to resolve this issue dealing with all the data. Also, in the paper we used a deep learning algorithm, implementing a boosted Neural Network.

Methodology

For this work 590 COF structures with diverse building blocks, topologies, and dimensionalities were selected from a database of known as CURATED COFs,^[13] which compiles several synthesized structures reported in the literature. The textural properties were calculated using the Zeo++^[14] software, which performs geometric analysis of porous and crystalline materials based on the formation of the Voronoi networks and determines a set of structural descriptors for the selected COFs. The set of descriptors selected can be divides in three classes: i) pore size, ii) specific area and iii) channels number (Fig.4).



Fig.4. The computational pipeline for the COFs classification by machine learning

The pore size related descriptors were the large included sphere (LIS), which gives us information on the pore diameter of the material, the Large Free Sphere (LFS) and the Large free sphere path (LSFP). These last two are generally similar to the pore diameter, but for structures that have functional groups within the pore it can be

considerably different, so it is interesting to evaluate these three descriptors. Also related to the pore size the pore volume (V_p) in cm³/g is used to classify the structures.

The specific area related descriptors were the specific gravimetric area (SSA) and specific volumetric area (SSAV), in square meters per gram and in square meters per volume unit respectively.

The channel related descriptors were the number of accessible channels, the dimensionality of the channels, the size of the channels, and the accessible volume fraction (AVF), which is the percentage of the material composed by pores.

All the data are listed and defined: $COF_name - Name of the structure$ LIS - Large included sphere (Pore Diameter) [Å] LFS - Large free sphere [Å] LSFP - large sphere free path [Å] SSA - Specific surface area [m²/g] SSAV - Specific volumetric surface area m²/cm³ AVF - Accessible volume fraction $V_p - Pore volume cm³/g$ N(chan) - Number of accessible channels DimChan - Dimensionality of the channels (1D, 2D ou 3D)ChanSize - Channel size [Å]

Results

The dataset was screened for missing, negative or zero values, whereas the set shrank to 580 structures. Thereupon, the outliers two outliers by criterion of quantile range were excluded from SSA (2 outlier) and Vp (2 outlier) belonging to the same two points at Q = 3 and the tail quantile of 0.1. Hence, the final number of the structures has become 576 with no other outliers (Table 1). Now, the cleaned dataset can be subjected to the machine learning experiments.

The scatter plots demonstrate strong collinearity within two sets of the inputs: between LIS, LFS and LSFP and between ChanSize X, Y and Z, which is critical while linear modeling. SSAV has a weak negative correlation with other parameters in 2D and it is a bit stronger for 3D keeping its negative character. Generally, the 2D dataset is characterized by more pronounced correlations, than 3D is.

 Table 1. The final screening analysis after excluding outliers for 2D (A) and 3D (B) structures

			•	•	•		•	,	• •	
10%	90%	Low	High	Number of Outlier	rs	10	% <mark>90</mark> %	Low	High	Number of Outliers
Quantile	Quantile	Threshold	Threshold	Outliers (Count) Colum	n Quanti	e Quantile	Threshold	Threshold	Outliers (Count)
9.50553	29.506	-50.496	89.5074	0	LIS, Å	5.7825	1 26.7402	-57.09	89.6132	0
8.60495	29.106	-52.898	90.6091	0	LFS, Å	4.3313	9 21.1113	-46.008	71.4509	0
9.48433	29.506	-50.581	89.571	0	LSFP, Å	5.6650	9 26.2119	-55.975	87.8524	0
1607.67	2767.96	-1873.2	6248.83	0	SSA, ma	2/g 1036.3	1 7863.62	-19446	28345.6	0
985.999	1556.34	-725.01	3267.35	0	SSAV, n	n2/cm3 835.43	4 2339.99	-3678.2	6853.67	0
0.25414	0.6182	-0.838	1.71037	0	AVF	0.1003	6 0.81916	-2.0561	2.97559	0
0.31535	1.56549	-3.4351	5.31593	0	Vp, cm3	3/g 0.0991	2 5.66305	-16.593	22.3549	0
8.80503	29.1802	-52.32	90.3055	0	ChanSiz	e X, Å 5.0636	8 26.7402	-59.966	91.7698	0
7.88985	28.91	-55.171	91.9706	0	ChanSiz	e Y, Å 3.6919	9 21.1113	-48.566	73.3692	0
8.80503	29.1802	-52.32	90.3055	0 (A) ChanSiz	e Z, Å 5.0353	6 26.6137	-59.7	91.3488	0 (B)
	10% Quantile 9.50553 8.60495 9.48433 1607.67 985.999 0.25414 0.31535 8.80503 7.88985 8.80503	10% 90% Quantile Quantile 9.50553 29.506 8.60495 29.106 1607.67 2767.96 985.999 1556.34 0.25414 0.6182 0.31535 1.56549 8.80503 29.1802 7.88985 28.31 8.80503 29.1802	10% 90% Low Quantile Quantile Threshold 9.50553 29.506 -50.496 8.60495 29.506 -52.898 9.48433 29.506 -1873.2 985.999 1556.34 -725.01 0.25414 0.6182 -0.838 0.31535 1.56544 -3.4351 8.80503 29.1802 -52.32 7.88985 2.8.91 -55.171 8.80503 29.1802 -52.32	10% 90% Low High Quantile Quantile Threshold Threshold 9.50553 29.506 -50.496 89.5074 8.60495 29.506 -52.898 90.6091 9.48433 29.506 -1873.2 6248.83 985.999 1556.34 -725.01 3267.35 0.25414 0.6182 -0.838 1.71037 0.31535 1.56549 -3.4351 5.31593 8.80503 29.1802 -52.32 90.3055 7.8986 2.891 -55.171 91.9706 8.80503 29.1802 -52.32 90.3055	10% 90% Low High Threshold Number of Outlier Quantile Variantia Threshold Threshold Outliers (Countiantiantiantiantiantiantiantiantiantia	10% 90% Low High Threshold Number of Outliers Outliers (Count) 9.50553 29.506 -50.496 89.5074 0 LlS, Å 8.60495 29.106 -52.898 90.6091 0 LFS, Å 9.48433 29.506 -50.581 89.571 0 LSFP, Å 1607.67 2767.96 -1873.2 6248.83 0 SSA, mi 985.999 1556.34 -725.01 3267.35 0 SSA, mi 0.25414 0.6182 -0.838 1.71037 0 AVF 0.31535 1.56549 -3.4351 5.31593 0 Vp, cmi 8.80503 29.1802 -52.32 90.3055 0 ChanSiz 7.88985 2.810 -55.717 19.19706 0 ChanSiz 8.80503 29.1802 -52.32 90.3055 (A) ChanSiz	10% 90% Low High Threshold Number of Outliers Outliers (Count) 109 Column 9.50553 29.506 -50.496 89.5074 0 LIS, Å 5.7825 8.60495 29.106 -52.898 90.6091 0 LFS, Å 4.3313 9.48433 29.506 -50.496 89.5074 0 LIS, Å 5.7825 1607.67 2767.96 -1873.2 6248.83 0 SSA, m2/g 10363 985.999 1556.34 -725.01 3267.35 0 SSAV, m2/cm3 835.43 0.25414 0.6182 -0.838 1.71037 0 AVF 0.1003 0.31535 1.56549 -3.4351 5.31593 0 Vp, cm3/g 0.0991 8.80503 29.1802 -52.32 90.3055 0 ChanSize Y, Å 3.6919 8.80503 29.1802 -52.32 90.3055 0 (A) ChanSize Y, Å 3.6919	10% 90% Low High Number of Outliers Outliers (Count) 10% 90% Quantile 9.50553 29.506 -50.496 89.5074 0 Uls, Å 5.78251 26.7402 8.60495 29.506 -50.496 89.5074 0 Uls, Å 5.78251 26.7402 1607.67 2767.96 -1873.2 6248.83 0 SSA, m2/g 1036.31 7863.62 985.999 1556.34 -725.01 3267.35 0 SSAV, m2/cm3 835.434 2339.99 0.25414 0.6182 -0.838 1.71037 0 AVF 0.10036 0.81916 0.31535 1.56549 -3.4351 5.31593 0 Vp, cm3/g 0.0912 5.66305 8.80503 29.1802 -52.32 90.3055 0 ChanSize X, Å 5.06368 26.7402 7.88985 28.91 -55.717 91.9706 ChanSize Y, Å 3.66399 21.1113 8.80503 29.1802 -52.32 90.3055 (A) ChanSize Y, Å	10% 90% Low High Number of Outliers Outliers (Count) 10% 90% Low 9.50553 29.506 -50.496 89.5074 0 Uls, Å 5.78251 26.7402 -57.09 8.60495 29.106 -52.898 90.6091 0 LS, Å 5.78251 26.7402 -57.09 9.48433 29.506 -50.581 89.571 0 LSFP, Å 5.66509 26.2119 -55.975 1607.67 2767.96 -1873.2 6248.83 0 SSA, m2/g 1036.01 17863.62 -19446 985.999 1556.34 -725.01 3267.35 0 SSAV, m2/cm3 835.434 233.99 -3678.2 0.25414 0.6182 -0.838 1.71037 0 AVF 0.10036 0.81916 -2.0561 0.31535 1.56549 -3.4351 5.31593 0 Vp, cm3/g 0.0912 5.66305 16.593 8.80503 29.1802 -52.32 90.3055 0 ChanSize X, Å 5.06368	10% 90% Low High Pupuntile Number of Outliers Outliers (Count) 10% 90% Low High Pupuntile 9.50553 29.506 -50.496 89.5074 0 LIS, Å 5.78251 26.7402 -57.09 89.6132 8.60495 29.506 -50.496 89.5074 0 LIS, Å 5.78251 26.7402 -57.09 89.6132 9.48433 29.506 -50.581 89.571 0 LSFP, Å 5.66509 26.2119 -55.975 87.8524 1607.67 2767.96 -187.32 6248.83 0 SSA, m2/g 1036.31 7863.62 -19446 28345.6 985.999 1556.34 -725.01 3267.35 0 SSAV, m2/cm3 835.434 239.99 -3678.2 6853.67 0.25414 0.6182 -0.838 1.71037 0 AVF 0.10036 0.81916 -2.0561 2.97559 0.31535 1.56549 -3.4351 5.31593 0 Vp, cm3/g 0.09912 5.66305

Table. 2. The scatterplot matrix with correlations for over the input parameters for 2D





Table. 3. The scatterplot matrix with correlations for over the input parameters for 3D

The distribution of the structures to be classified has shown an unbalanced profile, where the 85%-portion is attributed to 2D. Thus, in order to build a correct predictive modeling, the training-validation split must be performed under stratification by "type", as a binary response of 2D/3D (Fig.5, 6). Also, the profit table for "type" was adjusted (3D - 0.15 and 2D - 0.85).





Fig.6. The training-validation split stratified by "type"

It is worth noting that non-continuous parameters, such a ChanDim and N(chan) can play an essential role as variables, having some issues in the regression (linear or logistic) with well-known instability during the parameters estimation. In other cases, for non-linear or non-parametric models, this feature is not encountered (Fig.7). It is seen that absolute majority of monodimensional channels (91.0%) is typical for 2D, whereas the 3D class possesses tridimensional and two-dimensional at the ratio of > 2.3.



Fig.7. The distribution of the number of channels over the channels dimensionalities for 2D and 3D structures.

Starting with KNN (K_{max} = 100, Euclidean distance), the minimal misclassification rate or (1-Accuracy)× 0.01 of 8.7% with F1 = 0.995 are the results at K = 3 (Fig.8). The mean values of the neighbors' distances and the modes for 3D are shifted to higher values both for training and validation. To develop this model, all the variables (continuous and nominal) were involved. This leads us to conclusion that both types of

the structures can be considered as mean clusters, whose coordinates are the principal components.



Fig.8. The summary of KNN: the model selection, the confusion matrices of training and validation, the distributions of the neighbors distances over the training-validation split

As for the multiple logistic regression, after removing all the variables at p > 0.05 and excluding those with unstable status in the parameters estimation, the model has got three main factors: LIS, SSA and SSAV (Fig.9). In other words, specific surface area and pore diameter are the main characteristics in the linear classification with a non-significant lack of fit (p = 1.000). The profit matrix has given 7.83% of balanced misclassification, whereas F1 = 0.952. In the contrast to KNN, this model has given less accuracy, but it involved only three parameters to distinguish 2D/3D COFs. In the other hand, the accuracy for KNN during validation has turned out to be worse than at training - as opposed to the case of multiple logistic regression, where the training has higher misclassification (8.89%) than validation.



Fig.9. The summary on the Multiple Logistic Regression

The Quadratic Discriminant Analysis involving only three main variables from Multiple Logistic Regression (SSA, LIS and SSAV) basically demonstrates quite weak classification – the improved decision matrix on validation gives MR 13.91% at F1 = 0.912 (Fig.10).

Decision Tree after 14 splits showed that ChanSize of all three directions (X, Y and Z) and LSFP did not contribute to the model (Fig. 11) and, hence, the model was relaunched and tuned excluding these parameters (Fig. 12). Simplified model has not lost the accuracy (F1 = 0.979) and the profit value in the respective matrix – its validation misclassification rate is 3.48% with balanced distribution of the misclassified data, and less than for training (4.34%) and F1 = 0.974. The model has shown to be more accurate than multiple logistic regression and does not overfit the data.



Fig.10. The Discriminant Analysis summary



Fig.11. The overall model of Decision Tree to classify 2D/3D structures

Similar situation is encountered with Bootstrap Forest (10 terms sampled per split, 12 terms, the minimum size split = 5) – two variables related to the ChanSize, namely X and Z, contribute poorly to the model (< 0.01), giving quite high MR (6.09%) and F1 (0.963). (Fig.13). Thus, the tuned Bootstrap Forest (5 terms per split, 10 terms, the minimal split = 5) revealed a decrease of the trees (only two trees), the misclassification rate in validation (3.48%) and F1 = 0.979 (Fig.14) due to the redistribution in 2D classification: from 6.1% of misclassified after tuning the MR was reduced down to 3.1%.

The Boosted Tree (140 layers, 5 splits per tree, the learning rate = 0.1) let us to exclude two last variables (ChanSize Z and LSFP) for more adequate fitting (Fig.15) and the resulted tuned (adjusted) Boosted Tree gave more reduced number of the layers (88 layers, 10 splits per tree, learning rate = 0.1) (Fig.16, 17). Moreover, the metrics of the model underwent the significant improvement fur to better attribution of 2D structures:

MR for 2D in the initial full model was 3.1%, whereas after tuning, MR became 1.0%, which does not change after the 36th layer.



Fig.12. Tuned Decision Tree model's summary: the split tree structure, split history, the variable contribution, decision and confusions matrices

Column Con	tribution	IS			Decisio	on Ma	trix						
	Number				Т	raining		Va	alidatio	n	Specifie	d Profit	Matrix
Term	of Splits	G^2		Portion	Decision			Decis	ion		Decis	sion	
SSA, m2/g	94	88.9562796		0.4259	Actual	Cou	nt	Actual	Cou	nt	Actual	2D	30
LFS, Å	48	26.3203781		0.1260	type	2D	3D	type	2D	3D	2D	0	-0.176
N(chan)	75	23.4579892		0.1123	2D	370	23	2D	92	6	3D	-1	(
Vp, cm3/g	30	14.8287573		0.0710	3D	0	68	3D	1	16			
AVF	40	13.2701419		0.0635		Dec	ision		Deci	ision	1		
LIS, Å	43	11.7627856		0.0563	Actual	Ra	ate	Actual	Ra	te			
SSAV, m2/cm3	38	9.04431493		0.0433	type	2D	3D	type	2D	3D			
ChanSize Y, Å	8	7.82321515		0.0375	20	0.9/1	0.059	20	0.939	0.061			
LSFP, Å	12	6.21571379		0.0298	30	0.000	1,000	3D	0.059	0.001			
ChanDim	32	5.08790944		0.0244	Micelay	cificat	ion	Micela	cificat	ion			
ChanSize X, Å	6	1.05734272		0.0051	Rate			IVIISCIA	D	ato			
ChanSize Z, Å	4	1.0177201		0.0049		0.0	100		0.0	609			
						0.0	455		0.0	005			
Cumulative	Validation				Confus	ion M	atrix						
0.05	_				Tr	aining		Vali	dation				
5 0.04 -						Predi	cted		Predic	ted			
- E0 0					Actual	Cou	Int	Actual	Cour	nt			
Suc			MR		type	2D	3D	type	2D	3D			
0.02-					2D	392	1	2D	98	0			
.01-					3D	15	53	3D	3	14			
0.00													
0.001	5 10	15 20 25	30										
	Numb	er of Trees											

Fig.13. The overall Bootstrap Forest model's summary: the variable contribution, cumulative validation progress, decision and confusions matrices.



Fig.14. Tuned Bootstrap Forest model's summary: the variable contribution, cumulative validation progress, decision and confusions matrices.

Column Con	tribution	s		Decisio	on Ma	trix										
	Number			Training		Validation		Specified Profit Matrix			atrix					
Term	of Splits	G^2	Portion		Decis	ion		Decis	sion		1	Decisio	on			
ChanDim	266	38566.3038	0.6160	Actual	Cou	nt	Actual	Cou	nt	Act	ual	2D	3D			
N(chan)	182	12141.5987	0.1939	type	2D	3D	type	2D	3D	2D		0 -	0.176			
SSA, m2/g	49	2960.58252	0.0473	2D	389	4	2D	95	3	3D		-1	0			
SSAV, m2/cm3	61	2613.02466	0.0417	3D	1	67	3D	1	16					1		
AVF	46	2386.14396	0.0381		Dec	ision		Dec	ision		onfus	on M	atriv			
LFS, Å	34	1496.32315	0.0239	Actual	R	te	Actual	R	ate		.omus		uuin			
Vp, cm3/g	22	846.574883	0.0135	type	2D	3D	type	20	3D		Ir	aining		Validation		
LIS, Å	19	626.325979	0.0100	20	0.000	0.010	20	0.960	0.021			Predi	cted		Predie	cted
ChanSize X, Å	9	537.024361	0.0086	20	0.990	0.010	20	0.505	0.031		Actual	Cou	int	Actual	Cou	nt
ChanSize Y, Å	9	323.122722	0.0052	50	0.015		30	0.055	0.541		type	2D	3D	type	2D	3D
ChanSize Z. Å	2	83.4297327	0.0013	Misclas	sificat	ion	Misclas	ssificat	tion		2D	392	1	2D	97	1
LSFP, Å	1	23.9562445	0.0004		R	ate		F	Rate		3D	4	64	3D	3	14
LSFP, A	1	23.9562445	0.0004		0.0	108		0.0	348	8	50	4	04	30	2	14

Fig.15. The Boosted Tree summary for the full model: the variable contribution, decision and confusions matrices

Column Con	tribution	s		Decisio	on Ma	trix										
	Number			Т	raining		Va	lidatio	n	Spe	ecified P	rofit M	atrix			
Term	of Splits	G^2	Portion		Decis	ion		Decis	ion		1	Decisio	n			
ChanDim	173	14788.2565	0.3976	Actual	Cou	nt	Actual	Cou	nt	Act	ual	2D	3D			
N(chan)	306	10700.569	0.2877	type	2D	3D	type	2D	3D	2D		0 -1	0.176			
SSA, m2/g	87	2708.25816	0.0728	2D	392	1	2D	97	1	3D		-1	0			
AVF	73	2267.8157	0.0610	3D	0	68	3D	1	16							
SSAV, m2/cm3	68	2101.46698	0.0565		Decision Decision Confusio				ion M	atrix						
LFS, Å	36	1497.87456	0.0403	Actual	Ra	te	Actual	Ra	te		т.			N/-1	dest en	_
Vp, cm3/g	58	1186.98528	0.0319	type	20	30	type	20	31	•		aining		Val	dation	1
LIS, Å	45	980.941258	0.0264	20	0.007	0.003	20	0.000	0.010			Predi	cted		Predi	cted
ChanSize X, Å	15	499.499875	0.0134	3D	0.000	1,000	3D	0.550	0.01	1	Actual	Cou	int	Actual	Cou	int
ChanSize Y. Å	19	466.658052	0.0125	50	0.000	1.000	50	0.039	0.94		type	2D	3D	type	2D	3D
				Misclas	sificat	ion	Misclas	sificat	ion		2D	392	1	2D	97	1
					R	ate		R	ate		3D	0	68	3D	2	15
					0.0	022		0.0	174					L		

Fig.16. The Boosted Tree summary after tuning: the variable contribution, decision and confusions matrices



Fig.17. The cumulative MR progress with the layers for the full model – MR(Full) and for the tuned model – MR (Adj)

Regarding the Support Vector Machine which comprises all the variables, this model misclassified the structures at the same MR (<1%), as KNN did with quite balanced distribution in the confusion matrix at validation (Fig.18). But a huge difference between validation and training may speak for underfitting, although such a difference in the case of KNN is more pronounced. Thus, in spite of the equal F1 at validation (0.995), the model underfits less than KNN.



Fig.18. The Support Vector Machine results: the parallel plots of the variables vs. real distribution and the predicted 2D/3D distribution - over training and validation. The confusion matrices for training and validation.

After arranging all the models by their metrics into Tables 4, we can conclude that during the training, the Tree family has shown itself more efficient by the key KPIs. Support Vector Machines stays on the third position in training, but is a leader model at the validation. Comparing the ROCs (Fig. 19), the main three models reveal their effectiveness: Boosted Tree, SVM and Bootstrap Forest.

Model	Entropy R ²	Generalized R ²	Mean -Log p	MR	Average Profit	AUC
Boosted Tree	0.9690	0.9799	0.013	0.0022	-4·10 ⁻⁴	0.9999
Classification Tree	0.7916	0.8545	0.0872	0.0304	-0.013	0.9828
Support Vector Machines	0.7807	0.8461	0.0917	0.0282	-0.015	0.9920
Bootstrap Forest	0.7391	0.8136	0.1091	0.0412	-0.02	0.9865
Logistic Regression	0.6043	0.7001	0.1655	0.0672	-0.028	0.9533
Discriminant Analysis	0.0254	0.0371	0.4077	0.0998	-0.036	0.9456

Table 4. The model comparison by performance at training



The situation is a bit different at the validation: Classification Tree, SVM and Boosted Tree keep leadership as on ROCs (Fig.20), as by tabulated KPIs (Table 5). The only significant difference between the AUCs for SVM is demonstrated by the

Discriminant Analysis, however, *p* is not small enough (0.0481) to reject H_o about the non-significant difference between the models by their AUC (Table 6). Thus, for the sufficient classification of the COFs structures by the textural properties we can recommend KNN (K = 3). SVM (v = 0.5 cost = 1, pumber of the

properties, we can recommend KNN (K = 3), SVM (γ = 0.5, cost = 1, number of the support vectors = 147), Classification Tree (14 splits).

Table 5	. The mode	el comparison b	y perform	nance at v	alidation		
Model	Entropy R ²	Generalized R ²	Mean -Log p	MR	Average Profit	AUC	F1
Support Vector Machines	0.8469	0.8956	0.0641	0.0174	-0.01	0.9940	0.995
Classification Tree	0.8439	0.8934	0.0654	0.0174	-0.013	0.9964	0.979
Boosted Tree	0.8169	0.8735	0.0767	0.0261	-0.01	0.9934	0.990
Bootstrap Forest	0.8130	0.8707	0.0783	0.0261	-0.013	0.9928	0.979
Logistic Regression	0.6928	0.7762	0.1287	0.0609	-0.021	0.9838	0.952
Discriminant Analysis	0.3596	0.4585	0.2683	0.0870	-0.032	0.9706	0.912



Fig.20. The ROC curves for validation over all the involved models

Model	vs. Model	AUC Difference	Std Error	Lower 95%	Upper 95%	X ²	p>χ²
Bootstrap Forest	SVM	-0.001	0.0096	-0.020	0.0177	0.0155	0.9009
Boosted Tree	SVM	-0.001	0.0023	-0.005	0.0038	0.0708	0.7902
Logistic Regression	SVM	-0.010	0.0084	-0.027	0.0063	1.4723	0.2250
Discriminant Analysis	SVM	-0.023	0.0118	-0.047	-0.000	3.9061	0.0481*
Classification Tree	SVM	0.0024	0.0072	-0.012	0.0164	0.1127	0.7371
Bootstrap Forest	Logistic Regression	0.0090	0.0112	-0.013	0.0310	0.6423	0.4229
Boosted Tree	Logistic Regression	0.0096	0.0081	-0.006	0.0254	1.4186	0.2336
Bootstrap Forest	Discriminant Analysis	0.0222	0.0159	-0.009	0.0534	1.9472	0.1629
Boosted Tree	Discriminant Analysis	0.0228	0.0115	0.0002	0.0454	3.9178	0.0478*
Logistic Regression	Discriminant Analysis	0.0132	0.0082	-0.003	0.0293	2.6001	0.1069
Bootstrap Forest	Classification Tree	-0.004	0.0050	-0.013	0.0062	0.5203	0.4707
Boosted Tree	Classification Tree	-0.003	0.0068	-0.016	0.0104	0.1925	0.6608
Logistic Regression	Classification Tree	-0.013	0.0094	-0.031	0.0059	1.7836	0.1817
Discriminant Analysis	Classification Tree	-0.026	0.0146	-0.055	0.0029	3.107	0.0780
Bootstrap Forest	Boosted Tree	-0.001	0.0092	-0.019	0.0174	0.0043	0.9480

 Table 6. The AUC difference analysis

Finally, all the available variables served for the boosted Neural Network of 39 sigmoidal activation functions incorporated in a single layer (Fig.21) after one tour of iterations with squared penalty method. Comparing training and validation matrices, one can easily see that neither over- nor underfitting occurred – this also is proven by minor difference between the Entropy R² and Generalized R². The respective interactive html-coded calculator (profiler) and the estimates for each activation function summary are uploaded to a Git-Hub repository: *https://github.com/Nazarkovsky/COFs-dimensionality-prediction.-Boosted-Neural-Network*.

Frainin	g			Validati	on		
type				type			
Measur	es		Value	Measur	es		Value
General	ized RS	quare	0.9963629	General	ized R	Square	0.9974741
Entropy	RSqua	re	0.9943244	Entropy	RSqu	are	0.996053
Misclass	sificatio	n Rate	0	Misclass	ificati	on Rate	0
Confus	ion Ma	trix		Confus	ion M	atrix	
	Predic	ted			Predi	cted	
Actual	Cou	nt		Actual	Cou	int	
type	2D	3D		type	2D	3D	
2D	393	0		2D	98	0	
3D	0	68		3D	0	17	
Confu	usion R	ates		Confu	usion F	Rates	
	Predi	icted			Prec	licted	
Actual	Ra	te		Actual	R	ate	
type	2D	3D		type	20	3D	
2D	1.000	0.000		2D	1.000	0.000	
3D	0.000	1.000		3D	0.000	1.000	

Fig.21. The Neural Network summary, confusion matrix and the diagram

Conclusion

As a conclusion, we recommend to test the most efficient models, SVM, KNN and Neural Network on new data to improve the performance of the proposed techniques, if necessary. After testing the models with thousands of the structures, we can take the final decision on the selection of the preferred machine learning algorithm (or an ensemble of the models) to be deployed and utilized as an application or interactive table to predict the dimensionality of the covalent organic frameworks. This study is helpful for the future predictive modeling of the physicochemical properties for each type of COFs.

References

- A.P. Cote, A.P. Côte, A.I. Benin, N.W. Ockwig, M. O'Keeffe, A.J. Matzger, O.M. Yaghi, *Science (80-.).* 2005, *310*, 1166–1170.
- [2] C.S. Diercks, O.M. Yaghi, The atom, the molecule, and the covalent organic framework., *American Association for the Advancement of Science*, [2017,].
- [3] X. Feng, X. Ding, D. Jiang, *Chem. Soc. Rev.* **2012**, *41*, 6010–6022.
- [4] R.S.B. Gonçalves, A.B.V. Deoliveira, H.C. Sindra, B.S. Archanjo, M.E. Mendoza, L.S.A. Carneiro, C.D. Buarque, P.M. Esteves, *ChemCatChem.* 2016, *8*, 743– 750.
- [5] R.A. Maia, F. Berg, V. Ritleng, B. Louis, P.M. Esteves, Chem. A Eur. J. 2020, 26, 2051–2059.
- [6] F.L. Oliveira, A. S. França, A.M. Castro, R.O.M. Alves de Souza, P.M. Esteves, R.S.B. Gonçalves, *Chempluschem.* **2020**, *85*, 2051–2066.
- [7] A.K. Mandal, J. Mahmood, J.B. Baek, *ChemNanoMat.* **2017**, *3*, 373–391.
- [8] R. Xue, H. Guo, T. Wang, L. Gong, Y. Wang, J. Ai, D. Huang, H. Chen, W. Yang, *Anal. Methods.* **2017**, *9*, 3737–3750.
- [9] S.K.S. Freitas, R.S. Borges, C. Merlini, G.M.O. Barra, P.M. Esteves, *J. Phys. Chem. C.* **2017**, *121*, 27247–27252.
- [10] Y. Zeng, R. Zou, Y. Zhao, Adv. Mater. 2016, 28, 2855–2873.
- [11] J. Ozdemir, I. Mosleh, M. Abolhassani, L.F. Greenlee, R.R. Beitle, M.H. Beyzavi, *Front. Energy Res.* **2019**, *7*, 77.
- [12] P.J. Waller, F. Gándara, O.M. Yaghi, Acc. Chem. Res. 2015, 48, 3053–3063.
- [13] D. Ongari, A. V. Yakutovich, L. Talirz, B. Smit, ACS Cent. Sci. 2019, 5, 1663– 1675.
- [14] T.F. Willems, C.H. Rycroft, M. Kazi, J.C. Meza, M. Haranczyk, *Microporous Mesoporous Mater.* **2012**, *14*9, 134–141.