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



Felipe Lopes de Oliveira, Federal University of Rio de Janeiro





Michael Nazarkovsky, Pontifical Catholic University of Rio de Janeiro





PART 1. Introduction



Covalent Organic Frameworks (COFs)



Covalent Organic Frameworks (COFs)



Application of Covalent Organic Frameworks



Classical paradigm of materials discovery





Our computational approach



1. ONGARI, Daniele et al. Building a consistent and reproducible database for adsorption evaluation in covalent–organic frameworks. ACS central science, v. 5, n. 10, p. 1663-1675, 2019.

Some important properties

- **COF_ID** Name of the structure
- LIS Large included sphere (Pore Diameter). [Å]
- LFS Large free sphere [Å]
- LSFP large sphere free path [Å]
- $\ensuremath{\text{SSA}}$ Specific surface area $[m^2/g]$
- SSAV Specific volumetric surface area m²/cm³
- AVF Accessible volume fraction
- **Vp** Pore volume cm³/g
- N(chan) Number of accessible channels
- ChanDim Dimensionality of the channels (1D, 2D ou 3D)
- ChanSize Channel size [Å]



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



Felipe Lopes de Oliveira, Federal University of Rio de Janeiro





Michael Nazarkovsky, Pontifical Catholic University of Rio de Janeiro



PART 2. Data Pretreatment, Multivariate Analysis, Unsupervised and Supervised Predictive Modeling



Data Preparation and Preliminary Visualization

- Exlusion of 6 erroneous structures containg "0" values for channels and textural properties;
- Ordering in 10 continuous and 2 nominal inputs with a binary nominal output "type" (85% of 2D and 15% of 3D structures);



Multivariate Analysis – 2D



Hierarchical Clustering for 2D structures



Multivariate Analysis – 3D



Hierarchical Clustering for 3D structures



Portion of total variation in each column absorbed by clustering

2D

Column Summary					
Column	RSquare	.2.4.6.8			
ChanSize X, Å	0.9421				
ChanSize Z, Å	0.9420				
ChanSize Y, Å	0.9405				
LSFP, Å	0.9325				
AVF	0.9289				
Vp, cm3/g	0.9201				
SSAV, m2/cm3	0.8463				
SSA, m2/g	0.8010				

3D

Column Summary					
Column	RSquare	.2.4.6.8			
ChanSize X, Å	0.9332				
LSFP, Å	0.9322				
ChanSize Z, Å	0.9319				
ChanSize Y, Å	0.9200				
AVF	0.9128				
Vp, cm3/g	0.8600				
SSA, m2/g	0.8387				
SSAV, m2/cm3	0.8108				

2D/3D by the most massive clusters







76

37.25

204

100

128

62.75

Total

Logistic Regression

Whole Model Test

Model	-LogLikelihood	DF	χ²	p>χ²	Variable	LogWorth	p-Value
Difference	94.00774	2	188.0155	<.0001*	Vp, cm3/g	41.633	0.00000
Full	13.76301				AVE	37 202	0 00000
Reduced	107.77075				,,,,,,	57.202	0.00000

Training				Validation			
	Predi	cted			Predie	cted	
Actual	Count			Actual	Count		
type	2D	3D		type	2D	3D	
2D	99	3		2D	26	0	
3D	4	57		3D	0	15	

Metrics	Training	Validation
Entropy R ²	0.8723	0.9521
Generalized R ²	0.9332	0.9761
Mean –Log <i>p</i>	0.0844	0.0315
Misclassification Rate	0.0429	0.0000
Ν	163	41

on	Term	Estimate	Std Error	χ²	p>χ²
L	Intercept	-36.076807	11.01415	10.73	0.0011*
L	AVF	140.878677	36.651292	14.77	0.0001*
5	Vp, cm3/g	-31.147424	7.4543375	17.46	<.0001*

-36.0768067704061 + 140.878676823079AVF - 31.147423524102Vp = Lin2D

1/(1 + exp(-Lin2D)) = Probability_{2D}

1/(1 + exp(Lin2D)) = Probability_{3D}

K Nearest Neighbors (Euclidean distance. K_{max} = 162, K_{best} = 1)

Misclassification Rate

Confusion Matrix for Best K=1



Quadratic Discriminant Analysis

Score Summaries

		Number	Percent	Entropy	
Source	Count	Misclassified	Misclassified	RSquare	-2LogLikelihood
Training	163	8	4.90798	0.80767	41.4546
Validation	41	0	0.00000	0.98341	

Tra	aining		Validation		
Predicted				Predic	cted
Actual	Count		Actual	Count	
type	2D	3D	type	2D	3D
2D	102	0	2D	26	0
3D	8	53	3D	0	15

Cai	nonic	al Plot				
Canonical2	1.5 1.0 0.5 0.0 -0.5 -1.0		D D			•
		-2	0	2	4	6
				Canonical1		

Group Means						
Count	type	AVF	Vp, cm3/g			
102	2D	0.48692431	0.8707681			
61	3D	0.44836164	1.7041224			



Naive Bayes

Metrics	Training	Validation
Entropy R ²	0.6355	0.9744
Generalized R ²	0.7750	0.9874
Mean -Log p	0.2410	0.0168
Misclassification Rate	0.0736	0.0000
Ν	163	41

С	onfusio	on Ma	trix						
	Training				Validation				
	Actual	Predicted al Count			Actual	Predicted Count			
	type	2D	3D		type	2D	3D		
	2D	101	1		2D	26	0		
	3D	11	50		3D	0	15		

Independent Inputs (Overall)

Column	Main Effect	Total Effect	
Vp, cm ³ /g	0.217	0.879	
AVF	0.122	0.783	





3D

Scatterplot 3D Validation=Training



Support Vector Machine



Confusion Matrix								
Training Validation								
	Predi	cted		Predi	cted			
Actual	Count		Actual	Cou	nt			
type	2D	3D	type	2D	3D			
2D	100	2	2D	26	0			
3D	8	53	3D	0	15			



Cost = 1, Gamma = 0.5, #SV = 46

Metrics	Training	Validation
Entropy R ²	0.7849	0.9097
Generalized R ²	0.8805	0.9537
Mean -Log <i>p</i>	0.1422	0.0593
Misclassification Rate	0.0613	0.0000
N	163	41





...and PCA based Logistic Regression

Whole Model Test

Model	-LogLikelihood	DF	χ²	p>χ²
Difference	93.99795	2	187.9959	<.0001*
Full	13.77280			
Reduced	107.77075			

Metrics	Training	Validation	Confusi	on Ma	trix			
Entropy R ²	0.8722	0.9201	T	Training		Validation		
Generalized R ²	0.9331	0.9593	Actua	Predicted		Actual	Predicted	
Mean –Log p	0.0845	0.0525	type	2D	3D	type	2D	3D
Misclassification Rate	0.0061	0.0000	2D	102	0	2D	26	0
Ν	163	41	3D	1	60	3D	0	15

Term	Estimate	Std Error	χ²	p>χ²
Intercept	-0.1203949	0.6012218	0.04	0.8413
Prin1	1.80736932	0.3911632	21.35	<.0001*
Prin2	-4.004678	0.934766	18.35	<.0001*

-0.120394878783972 + 1.80736932265296Prin1 - 4.00467801810539Prin2 = Lin2D

 $1/(1 + exp(-Lin2D)) = Probability_{2D}$

1/(1 + exp(Lin2D)) = Probability_{3D}

Conclusions

Training

Model	Entropy R ²	Generalized R ²	Mean -Log p	MR	Ν	
Logistic Regression	0.8723	0.9332	0.0844	0.0429	163	∆MR
Logistic Regression by PCA	0.8722	0.9331	0.0845	0.0061	163	
Discriminant Analysis	0.8077	0.8948	0.1272	0.0491	163	d-bo-
Support Vector Machines	0.7849	0.8805	0.1422	0.0613	163	Mean I
Naive Bayes	0.6355	0.7750	0.241	0.0736	163	I R2
Validation						

Validation

Model	Entropy R ²	Generalized R ²	Mean -Log p	MR	Ν
Discriminant Analysis	0.9834	0.9919	0.0109	0.0000	41
Naive Bayes	0.9744	0.9874	0.0168	0.0000	41
Logistic Regression	0.9521	0.9761	0.0315	0.0000	41
Logistic Regression by PCA	0.9201	0.9593	0.0525	0.0000	41
Support Vector Machines	0.9097	0.9537	0.0593	0.0000	41

What's the best of the best?



Thank You!

Obrigado!

Дякую!

I'm Presenting

Acknowledgements:

Dr. David Kirmayer The Hebrew University of Jerusalem - HUJI, School of Pharmacy Jerusalem, Israel



jmp.com/discovery-america