Data analytics approaches to materials design: Critical role of a descriptor

Sergey V. Levchenko

Center for Energy Science and Technology (CEST) Skolkovo Institute of Science and Technology Moscow, Russia

Research paradigm shift



Research paradigm shift



High-throughput computational materials design

Top-down design:

target property (high activity and selectivity of a catalyst)

additional constraints (high stability, low toxicity,...)

synthesis recipe

not clear how to achieve this!

Bottom-up design:



The key issue: Complexity

$$i\frac{\partial\Psi(\boldsymbol{x}_{1},\boldsymbol{x}_{2},\ldots,\boldsymbol{x}_{n},\boldsymbol{R}_{1},\boldsymbol{R}_{2},\ldots,\boldsymbol{R}_{N},t)}{\partial t} = \hat{H}(t)\Psi(\boldsymbol{x}_{1},\boldsymbol{x}_{2},\ldots,\boldsymbol{x}_{n},\boldsymbol{R}_{1},\boldsymbol{R}_{2},\ldots,\boldsymbol{R}_{N},t)$$

1) Many-body problem (3(n + N)-dimensional)

2) Multiscale problem (tens orders of magnitude in time and space)

However, there is hope that the complexity can be treated *incrementally*

Including science in descriptors



structure descriptor: Cartesian coordinates \rightarrow changes, but properties do not change!



machine will learn symmetries, not (other) physics -- much more data will be needed for an accurate model



Simple(r) properties (bulk d-band center position and CO dissociation energy) are correlated to more complex properties (adsorption energy and reaction barrier)

The simpler quantities are called *descriptive parameters* (a *descriptor*)

J. K. Nørskov, T. Bligaard, J. Rossmeisl and C. H. Christensen, Nature Chemistry 1, 37 (2009)



A simple physical model (Newns-Anderson) motivates the *d*-band center descriptor

What if we don't know such a model, or we need a more accurate and more widely applicable model?

J. K. Nørskov, T. Bligaard, J. Rossmeisl and C. H. Christensen, Nature Chemistry 1, 37 (2009)



A simple physical model (Newns-Anderson) motivates the *d*-band center descriptor

Find descriptor from DATA!

J. K. Nørskov, T. Bligaard, J. Rossmeisl and C. H. Christensen, Nature Chemistry 1, 37 (2009)

Supervised data analysis



- 1) A descriptor d_i uniquely characterizes the material i as well as property-relevant elementary processes
- 2) The determination of the descriptor must not involve calculations as intensive as those needed for the evaluation of the property to be predicted

Target property model: Kernel ridge regression versus feature selection

Regression models: Basis set expansion in materials space

kernel ridge regressionlinear $P(d) = \sum_{i=1}^{N} c_i \exp\left(-\|d_i - d\|_2^2/2\sigma^2\right)$ P(d) = dcminimize $\|d_i - d_j\|_2^2 = \sum_{\alpha=1}^{\Omega} (d_{i,\alpha} - d_{j,\alpha})^2$

Regression: Importance of regularization



 $\min_{c} \sum_{i} (P(d_{i}, c) - P_{i})^{2} + \lambda f(c), \min_{\lambda} (\text{validation error}) \rightarrow \lambda$

Target property model: Kernel ridge regression versus feature selection

Regression models: Basis set expansion in materials space

kernel ridge regression	linear
$P(\boldsymbol{d}) = \sum_{i=1}^{N} c_i \exp\left(-\ \boldsymbol{d}_i - \boldsymbol{d}\ _2^2 / 2\sigma^2\right)$	P(d) = dc
$\sum_{i=1}^{N} (P(\boldsymbol{d}_i) - P_i)^2$ +	mize $\sum_{i=1}^N (P(oldsymbol{d}_i) - P_{oldsymbol{i}})^2$ +
$\lambda \sum_{i,j=1}^{N,N} c_i c_j \exp\left(-\ \boldsymbol{d}_i - \boldsymbol{d}_j\ _2^2/2\sigma^2\right)$	$\lambda \ \boldsymbol{c} \ _{0}$
$\ \boldsymbol{d}_{i} - \boldsymbol{d}_{j}\ _{2}^{2} = \sum_{\alpha=1}^{\Omega} (d_{i,\alpha} - d_{j,\alpha})^{2}$	

Target property model: Kernel ridge regression versus feature selection

kernel (Gaussian, Laplacian, linear $(d_i \cdot d_j)$) kernel ridge regression linear $P(\boldsymbol{d}) = \sum_{i=1}^{N} c_i \exp\left(-\|\boldsymbol{d}_i - \boldsymbol{d}\|_2^2 / 2\sigma^2\right)$ $P(\boldsymbol{d}) = \boldsymbol{d}\boldsymbol{c}$ minimize $\sum_{i=1}^{N} (P(\boldsymbol{d}_i) - P_i)^2 +$ $\sum_{i=1}^{N} (P(d_i) - P_i)^2 +$ $\lambda \sum_{i,j=1}^{N,N} c_i c_j \left(\exp\left(-\|\boldsymbol{d}_i - \boldsymbol{d}_j\|_2^2/2\sigma^2 \right) \right)$ $\lambda \| \boldsymbol{c} \|_{0}$ penalty on the number of non-zero coefficients $\|c\|_0$

penalty on similar data points

(Gaussian) kernel ridge regression example

Data: 175 linear 4-blocks periodic polymers. 7 blocks: CH₂, SiF₂, SiCl₂, GeF₂, GeCl₂, SnF₂, SnCl₂, Descriptor: 20 dimensions [# building blocks of type *i*, of *ii* pairs, of *iii* triplets]



Density Functional Theory

Pilania, Wang, ..., and Ramprasad, Scientific Reports 3, 2810 (2013). DOI: 10.1038/srep02810

- 1) A descriptor d_i uniquely characterizes the material i as well as property-relevant elementary processes
- 2) The determination of the descriptor must not involve calculations as intensive as those needed for the evaluation of the property to be predicted
- 3) The dimension Ω of the descriptor should be as low as possible (for a certain accuracy request)

Choose a physically motivated basis set!

L. M. Ghiringhelli, J. Vybiral, S. V. Levchenko, C. Draxl, and M. Scheffler, Phys. Rev. Lett. 114, 105503 (2015)

- 1) A descriptor d_i uniquely characterizes the material i as well as property-relevant elementary processes
- 2) The determination of the descriptor must not involve calculations as intensive as those needed for the evaluation of the property to be predicted
- 3) The dimension Ω of the descriptor should be as low as possible (for a certain accuracy request)

Idea: calculate many *physically motivated* quantities (features), and use these features as a basis for the physical model under compactness constraints

L. M. Ghiringhelli, J. Vybiral, S. V. Levchenko, C. Draxl, and M. Scheffler, Phys. Rev. Lett. 114, 105503 (2015)



Crystal-structure prediction was and is one of the most important, basic challenges of materials science and engineering.



Energy differences between different structures are very small.

For Si: 0.01% of the energy of a Si atom, or 0.1% of the 4 valence electrons.

Crystal-structure prediction was and is one of the most important, basic challenges of materials science and engineering.



J. A. van Vechten, Phys.
Rev. 182, 891 (1969). J. C. Phillips, Rev.
Mod. Phys. 42, 317 (1970).
J. John and A.N. Bloch, Phys. Rev. Let. 33, 1095 (1974) J. R. Chelikowsky and J. C.
Phillips, Phys. Rev. B 33, 2453 (1978)
A. Zunger, Phys. Rev. B 22, 5839 (1980).
D. G. Petifor, Solid State Commun. 51, 31 (1984). Y. Saad, D. Gao, T. Ngo, S.
Bobbit, J. R. Chelikowsky, and W.
Andreoni. Phys. Rev. B 85, 104104 (2012).

Can we predict not yet calculated structures from Z_A and Z_B ? Can we create a map: "The *ZB/W* community lives here and the *RS* community there?"



Can we predict not yet calculated structures from Z_A and Z_B ? **Can we create a** map: "The ZP/W reduction \rightarrow need a better basis



Can we predict not yet calculated structures from Z_A and Z_B ? Can we create a map: "The ZB/W community lives here and the RS community there?"



Can we predict not yet calculated structures from Z_A and Z_B ? Can we create a map: "The *ZB/W* community lives here and the *RS* community there?"



descriptor can be determined spectroscopically - properties of the solid

Can we create a map based on calculations simpler than bulk?

Primary features and feature space

ID	Description free atoms	Symbols	#
A1	Ionization Potential (IP) and Electron Affinity (EA)	IP(A) EA(A) IP(B) EA(B) [1]	4
A2	Highest occupied (H) and lowest unoccupied (L) Kohn-Sham levels	H(A) L(A) H(B) L(B)	4
A3	Radius at the max. value of s , p , and d valence radial radial probability density	$ \begin{array}{c c} r_s(\mathbf{A}) \ r_p(\mathbf{A}) \ r_d(\mathbf{A}) \\ r_s(\mathbf{B}) \ r_p(\mathbf{B}) \ r_d(\mathbf{B}) \end{array} $	6
ID	Description free dimers	Symbols	#
A4	Binding energy	$E_b(AA) E_b(BB) E_b(AB)$	3
A5	HOMO-LUMO KS gap	HL(AA) HL(BB) HL(AB)	3
A6	Equilibrium distance	$d(AA) \ d(BB) \ d(AB)$	3

How to find the best model for our target property (energy difference between different crystal structures)?

Symbolic regression: Eureqa



Uses evolutionary algorithm to find the best formula describing target property

Assumes "gene" structure of the formula \rightarrow bias

May result in an unnecessarily complex model

https://community.datarobot.com/t5/resources/introduction-to-eureqa/ta-p/2409

Primary features and feature space

ID	Description free atoms	Symbols	#
A1	Ionization Potential (IP) and Electron Affinity (EA	IP(A) EA(A) IP(B) EA(B) [1]	4
A2	Highest occupied (H) and lowest unoccupied (L)	H(A) L(A) H(B) L(B)	4
	Kohn-Sham levels		
A3	Radius at the max. value of s , p , and d	$r_s(\mathbf{A}) r_p(\mathbf{A}) r_d(\mathbf{A})$	6
	valence radial radial probability density	$r_s(\mathbf{B}) r_p(\mathbf{B}) r_d(\mathbf{B})$	
ID	Description free dimers	Symbols	#
A4	Binding energy	$E_b(AA) E_b(BB) E_b(AB)$	3
A5	HOMO-LUMO KS gap	HL(AA) HL(BB) HL(AB)	3
A6	Equilibrium distance	$d(AA) \ d(BB) \ d(AB)$	3
ID	description	prototype formula	-#
$\frac{1D}{R1}$	absolute differences and sums of A1	$ \mathbf{D}(\mathbf{A}) + \mathbf{D}(\mathbf{B}) $	12
	absolute differences and sums of A1	$ \mathbf{I} \mathbf{F}(\mathbf{A}) \pm \mathbf{I} \mathbf{F}(\mathbf{B}) $	12
B2	absolute differences and sums of A2	$ L(B) \pm H(A) $	12
$B3 \mid$	B absolute differences and sums of A3 $ r_p(A) \pm r_s(A) $		30
C3	3 squares of A3 and B3 (only sums) $r_s(A)^2, (r_p(A) + r_s(A))^2$		21
D3	23 exponentials of A3 and B3 (only sums) $\exp(r_s(A)), \exp(r_p(A) \pm r_s(A))$		21
E3	exponentials of squared $A3$ and $B3$ (only sums)	$\exp(r_s(\mathbf{A})^2), \exp(r_p(\mathbf{A}) \pm r_s(\mathbf{A})^2)$	21

We start with 23 primary features and build > 10,000 non-linear combinations

 P_j -- property value ($E_{ZB} - E_{RS}$) for material *j* (a function in materials space)

 $d_{j,l}$ -- value of feature l related to material j (e.g., $|r_s(A_j) - r_p(B_j)|$) (a basis function in materials space)

 c_l -- coefficient of the expansion of the property function in terms of basis functions:

$$P_j = \sum_l d_{j,l} c_l \qquad \text{How to find } c_l?$$

$$\sum_{j} \left(P_{j} - \sum_{l} d_{j,l} c_{l} \right)^{2} + \lambda \|\boldsymbol{c}\|_{n} \to \operatorname{argmin}(\boldsymbol{c})$$

regularization term to explore and ensure compactness of the expansion (reduce complexity)

 P_j -- property value ($E_{ZB} - E_{RS}$) for material *j* (a function in materials space)

 $d_{j,l}$ -- value of feature l related to material j (e.g., $|r_s(A_j) - r_p(B_j)|$) (a basis function in materials space)

 c_l -- coefficient of the expansion of the property function in terms of basis functions:

$$P_{j} = \sum_{l} d_{j,l}c_{l} \quad \text{How to find } c_{l}?$$

$$\sum_{j} \left(P_{j} - \sum_{l} d_{j,l}c_{l} \right)^{2} + \lambda \|c\|_{n} \rightarrow \operatorname{argmin}(c)$$

 $||c||_0$ -- number of non-zero coefficients \rightarrow NP hard! (need to try all combinations)

 P_j -- property value ($E_{ZB} - E_{RS}$) for material *j* (a function in materials space)

 $d_{j,l}$ -- value of feature l related to material j (e.g., $|r_s(A_j) - r_p(B_j)|$) (a basis function in materials space)

 c_l -- coefficient of the expansion of the property function in terms of basis functions:

$$P_{j} = \sum_{l} d_{j,l} c_{l} \quad \text{How to find } c_{l}?$$

$$\sum_{i} \left(P_{j} - \sum_{l} d_{j,l} c_{l} \right)^{2} + \lambda \| \boldsymbol{c} \|_{n} \rightarrow \operatorname{argmin}(\boldsymbol{c})$$

 $||c||_0$ -- number of non-zero coefficients \rightarrow NP hard! (need to try all combinations) $||c||_2 = \sum_l |c_l|^2$ -- ridge regression \rightarrow not most compact! $||c||_1 = \sum_l |c_l|$ -- LASSO (Least Absolute Shrinkage and Selection Operator) \rightarrow convex problem, equivalent to the NP-hard if features (columns of *d*) are uncorrelated

Compressed (compressive?) sensing





Raw: 15MB

JPEG: 150KB

Expand in a basis (wavelets) \rightarrow use LASSO to select most important basis functions \rightarrow store compressed image

 P_j -- property value ($E_{ZB} - E_{RS}$) for material *j* (a function in materials space)

 $d_{j,l}$ -- value of feature l related to material j (e.g., $|r_s(A_j) - r_p(B_j)|$) (a basis function in materials space)

 c_l -- coefficient of the expansion of the property function in terms of basis functions:



The descriptors selected with LASSO

$$\frac{\mathrm{IP}(\mathrm{B}) - \mathrm{EA}(\mathrm{B})}{r_p(\mathrm{A})^2}, \frac{|r_s(\mathrm{A}) - r_p(\mathrm{B})|}{\exp(r_s(\mathrm{A}))}, \frac{|r_p(\mathrm{B}) - r_s(\mathrm{B})|}{\exp(r_d(\mathrm{A}))}_{3\mathrm{D}}$$

$$\begin{split} \Delta E &= 0.117 \frac{\text{EA(B)} - \text{IP(B)}}{r_p(\text{A})^2} - 0.342 & \text{ID} \\ \Delta E &= 0.113 \frac{\text{EA(B)} - \text{IP(B)}}{r_p(\text{A})^2} + 1.542 \frac{|r_s(\text{A}) - r_p(\text{B})|}{\exp(r_s(\text{A}))} - 0.137 & \text{2D} \\ \Delta E &= 0.108 \frac{\text{EA(B)} - \text{IP(B)}}{r_p(\text{A})^2} + 1.790 \frac{|r_s(\text{A}) - r_p(\text{B})|}{\exp(r_s(\text{A}))} + & \text{3D} \\ &+ 3.766 \frac{|r_p(\text{B}) - r_s(\text{B})|}{\exp(r_d(\text{A}))} - 0.0267 \end{split}$$

Same features are selected for higher-dimensional descriptors, but this does not have to be the case

"The Map" -- compressed sensing -- LASSO, 2D descriptor



∆ ♦ ♦ •	= E(RS) - E(ZB) ZB, $\Delta > 0.2 \text{ eV}$ ZB, 0.1 eV $< \Delta \le 0.2 \text{ eV}$ ZB, 0.05 eV $< \Delta \le 0.1 \text{ eV}$ $- 0.05 \text{ eV} < \Delta \le 0.05 \text{ eV}$
	$-0.05 \text{ eV} < \Delta \le 0.05 \text{ eV}$ RS, $-0.1 \text{ eV} < \Delta \le -0.05 \text{ eV}$
	RS, $-0.2 \text{ eV} < \Delta \le -0.1 \text{ eV}$ RS, $\Delta \le -0.2 \text{ eV}$

$$P(j) = \boldsymbol{d}(j)\boldsymbol{c}$$

The complexity and science is in the descriptor (identified from >10,000 features).

L.M. Ghiringhelli, J. Vybiral, S.V. Levchenko, C. Draxl, and M. Scheffler, Phys. Rev. Lett. **114**, 105503 (2015).

Predictive power of the model

Hadn't we known about diamond ... we'd have predicted it!

When both carbon diamond and BN are excluded from training:

	⊿E(LDA)	∠E(predicted)
С	-2.64 eV	-1.44 eV
BN	-1.71 eV	-1.37 eV



Hadn't we known about any carbon-containing binary ... we'd have predicted carbon chemistry (from atomic features)

If all C containing binaries (C, SiC, GeC, and SnC) are excluded from training, i.e. no explicit information on C is given to the model:

	⊿E(LDA)	∠E(predicted)
С	-2.64 eV	-1.37 eV
SiC	-0.67 eV	-0.48 eV
GeC	-0.81 eV	-0.46 eV
SnC	-0.45 eV	-0.23 eV
Predictive power of the model

Mean absolute error (MAE), and maximum absolute error (MaxAE), in eV, (first two lines) and for a leave-10%-out cross validation (CV), averaged over 150 random selections of the training set (last two lines). For (Z_A^* , Z_B^*), each atom is identified by a string of three random numbers.

Descriptor	$Z_{\mathrm{A}}, Z_{\mathrm{B}}$	$Z_{\rm A}$ *, $Z_{\rm B}$ *	1 D	2D	3D	5D
MAE	1*10 ⁻⁴	3*10 ⁻³	0.12	0.08	0.07	0.05
MaxAE	8*10 ⁻⁴	0.03	0.32	0.32	0.24	0.20
MAE, CV	0.13	0.14	0.12	0.09	0.07	0.05
MaxAE, CV	0.43	0.42	0.27	0.18	0.16	0.12
	Gaussian-kernel r		L	ASSO)	

Predictive power of the model

Mean absolute error (MAE), and maximum absolute error (MaxAE), in eV, (first two lines) and for a leave-10%-out cross validation (CV), averaged over 150 random selections of the training set (last two lines). For (Z_A^* , Z_B^*), each atom is identified by a string of three random numbers.

	on -	LASSO					
MaxAE, CV	0.43	0.42	L	0.27	0.18	0.16	0.12
MAE, CV	0.13	0.14		0.12	0.09	0.07	0.05
MaxAE	8*10-4	0.03		0.32	0.32	0.24	0.20
MAE	1 * 10 ⁻⁴	3*10 ⁻³		0.12	0.08	0.07	0.05
Descriptor	$Z_{ m A}, Z_{ m B}$	$Z_{\rm A}^{*}, Z_{\rm B}^{*}$		1D	2D	3D	5D

Predictive power of the model

Mean absolute error (MAE), and maximum absolute error (MaxAE), in eV, (first two lines) and for a leave-10%-out cross validation (CV), averaged over 150 random selections of the training set (last two lines). For (Z_A^* , Z_B^*), each atom is identified by a string of three random numbers.

	γ Gaussian-kernel	ر ridge regression	L	L	ASSO	J
MaxAE, CV	0.43	0.42	0.27	0.18	0.16	0.12
MAE, CV	0.13	0.14	0.12	0.09	0.07	0.05
MAE MaxAE	1*10 ⁻⁴ 8*10 ⁻⁴	3*10 ⁻³ 0.03	0.12 0.32	0.08 0.32	0.07 0.24	0.05 0.20
Descriptor	$Z_{\rm A}, Z_{\rm B}$	$Z_{\rm A}^*, Z_{\rm B}^*$	1D	2D	3D	5D

Drawing causal inference from data



a mapping exists, even a physical intuition exists, but ΔE does not listen directly to the descriptor (intricate causality)

 $P(j) = \boldsymbol{d}(j)\boldsymbol{c}$

There are two aspects:

- 1) practical aspect -- we benefit from knowing $d \rightarrow P$ mapping for any convenient d(j) (analogy: plane waves)
- 2) physical aspect (understanding) -- we can reduce the complexity of the model and at the same time increase its applicability domain by a clever choice of d(j) (analogy: atomic orbitals and molecular-orbital picture)

We greatly benefit from d(j) providing a framework for a rational analysis

CH₄ chemical decomposition under shock-compression conditions (high *T* and *p*)

Yang, Q., Sing-Long, C. A., Reed, E. J., MRS Advances 1 (2016)

Methane at T = 3,300 K, p = 40.53 GPa: MD simulations (using a force-field description) find 2,613 different chemical reactions. Using compressed sensing it is shown that only 11% of them are relevant.

 $\min_{\widehat{k}} \|A\widehat{k} - b\|_{2}$ subject to $\widehat{k} \ge 0$, $\|\widehat{k}\|_{1} \le \lambda$ The *A* matrix has 2,613 columns, 2,395,918,510 rows



Lattice Anharmonicity and Thermal Conductivity from Compressive Sensing of First-Principles Calculations



\rightarrow predictive model for anharmonic lattice dynamics

F. Zhou, W. Nielson, Y. Xia, and Vidvuds Ozoliņš, Phys. Rev. Lett. 113, 185501 (2014)



4-body

3NN

Pair

NN

2NN

3-body

L. J. Nelson, G. L. W. Hart, F. Zhou, and V. Ozoliņš, Phys. Rev. B 87, 035125 (2013)

5-body

6-body

Vertex

distance

Enabling Feature Spaces with Billions of Elements by Sure Independence Screening

 $||c||_1 = \sum_l |c_l| - LASSO \rightarrow$ convex problem, equivalent to the NP-hard if features are uncorrelated \rightarrow not the case when many features are generated \rightarrow Sure Independence Screening plus Selection Operator (SISSO)

- 1. Systematically construct a huge feature space (10¹¹) from primary features: $\hat{R} = \{+, -, \cdot, -^{1}, ^{2}, ^{3}, \sqrt{-}, exp, log, |-|\}$ (use physically meaningful combinations!)
- 2. Select top ranked features using *Sure Independence Screening* (*SIS*)^[1] (correlation learning). Select *n* features corresponding to the *n* largest projection on the target property, i.e. largest components of the vector ($D^T y$)

 y: vector with the target property (e.g., rock saltzincblende energy differences; 82 elements)

- **D** : matrix of the feature space (e.g., 82 x 100 billion elements)
- 3. Apply a sparsifying operator (*I*₀ regularization) to the selected features to determine 1D, 2D,... descriptors
 R. Ouyang, *et al.*, Physical Review Materials 2, 083802 (2018)

SISSO: Iterative residual fitting



y: response vector P: target material property Residual: $R = P - \sum_i c_i d_i$

R. Ouyang, et al., Physical Review Materials 2, 083802 (2018)

SISSO: Performance

LASSO(+ l_0)

SISSO



SISSO: Performance



SISSO: Multitask and categorical

Multitask: Construct simultaneously SISSO models for several properties with the same descriptor

$$\min_{\boldsymbol{c}} \left(\lambda \| \boldsymbol{c}_{i}^{k} \|_{0} + \sum_{k} \frac{1}{N_{\text{samples}}^{k}} \sum_{\substack{\text{samples} \\ \text{in } k}} \left(\boldsymbol{P}^{k} - \boldsymbol{d} \boldsymbol{c}^{k} \right)^{2} \right) \to \boldsymbol{c}$$

Categorical (can be also multitask): Property - material belongs to a given class (yes/no)

$$\min_{\boldsymbol{c}} \left(\lambda \| c_i^k \|_0 + \sum_{I=1}^{N_{\text{classes}}} \sum_{J \neq I} O_{IJ}(\boldsymbol{d}, \boldsymbol{c}) \right) \to \boldsymbol{c}$$

number of data in the overlap region between domains of different classes in d-space

R. Ouyang, et al., J. Phys.: Mater. 2, 024002 (2019)

SISSO: Examples



• Perovskite phase stability (improved tolerance factor)



SISSO: Examples

Adsorption of molecules on metal surfaces

es on metal surfaces Adsorption of C, CH, CO, H, O, OH)



M. Andersen et al., ACS Catal. 9, 2752 (2019)

SISSO: Examples

• Design of topological insulators (materials for spintronics, catalysis, thermoelectricity)



G. Cao et al., arXiv:1808.04733

Data mining: Subgroup discovery under a material property 1 (y₁) Data mining: Subgroup discovery $(x_1)^2$ $(y_1)^2$ $(y_1)^2$

Subgroups are defined by selectors σ expressed as "AND" combinations of statements like "band gap < 2 eV", "atom radius > 1.4 Å", etc.

SGD algorithm: find subgroups that maximize quality function

$$f = N_{subgroup}/N_{all} \times |mean_{subgroup} - mean_{all}| \times (1 - variance_{subgroup}/variance_{all})$$

Numerical separators ("band gap < 2 eV") from k-means clustering (unsupervised learning) Search for subgroups: Monte Carlo or branch-and-bound algorithm

W. Klösgen, Advances in Knowledge Discovery and Data Mining. Palo Alto, CA: AAAI Press; 1996, 249

Data mining: Subgroup discovery



M. Boley et al., Data Min. Knowl. Disc. 31, 1391 (2017); B. Goldsmith et al., New J. Phys. 19, 013031 (2017)

Data mining: Subgroup discovery



M. Boley et al., Data Min. Knowl. Disc. 31, 1391 (2017); B. Goldsmith et al., New J. Phys. 19, 013031 (2017)

Data mining: Subgroup discovery



M. Boley et al., Data Min. Knowl. Disc. 31, 1391 (2017); B. Goldsmith et al., New J. Phys. 19, 013031 (2017)

Subgroup discovery: CO₂ activation by adsorption



Subgroup discovery: CO₂ activation by adsorption

Oxides:

dry reforming of methane: $CO_2 + CH_4 = 2H_2 + 2CO$

Me

Sabatier reaction: $CO_2 + 4H_2 = CH_4 + 2H_2O$

partial hydrogenation: $CO_2 + 3H_2 = CH_3OH + H_2O$ stable (structurally and compositionally) under increased temperatures;

. more resistant for poisoning;

 CO_2

activation is frequently observed

Subgroup discovery: CO_2 activation by adsorption O CMe

C-O bond elongation, O-C-O bending angle \rightarrow indicators of activation \rightarrow

Which surface properties lead to desired indicators?

Use subgroup discovery to find materials that optimize activation indicators

 $f = N_{subgroup}/N_{all} \times (mean_{subgroup} - mean_{all}) \times (1 - variance_{subgroup}/variance_{all})$ Maximize C-O bond length or O-C-O bending

Subgroup discovery: CO₂ activation by adsorption

19

A²⁺B⁴⁺O₃, A³⁺B³⁺O₃, A¹⁺B⁵⁺O₃, AO, BO₂, A₂O₃ (B₂O₃), A₂O, BO

,

			-			-		-		_		_	-	_			10
1 H 1.008	2											13	14	15	16	17	2 He 4.0026
3 Li 6.94	4 Be 9.0122											5 B 10.81	6 C 12.011	7 N 14.007	8 O 15.999	9 F 18.998	10 Ne 20.180
11 Na 22.990	12 Mg 24.305	3	4	5	6	7	8	9	10	11	12	13 Al 26.982	14 Si 28.085	15 P 30.974	16 S 32.06	17 Cl 35.45	18 Ar 39.948
19 K 39.098	20 Ca 40.078	21 Sc 44.956	22 Ti 47.867	23 V 50.942	24 Cr 51.996	25 Mn 54.938	26 Fe 55.845	27 Co 58.933	28 Ni 58.693	29 Cu 63.546	30 Zn 65.38	31 Ga 69.723	32 Ge 72.630	33 As 74.922	34 Se 78.97	35 Br 79.904	36 Kr 83.798
37 Rb 85.468	38 Sr 87.62	39 Y 88.906	40 Zr 91.224	41 Nb 92.906	42 Mo 95.95	43 Tc (98)	44 Ru 101.07	45 Rh 102.91	46 Pd 106.42	47 Ag 107.87	48 Cd 112.41	49 In 114.82	50 Sn 118.71	51 Sb 121.76	52 Te 127.60	53 I 126.90	54 Xe 131.29
55 Cs 132.91	56 Ba 137.33	57-71 *	72 Hf 178.49	73 Ta 180.95	74 W 183.84	75 Re 186.21	76 Os 190.23	77 Ir 192.22	78 Pt 195.08	79 Au 196.97	80 Hg 200.59	81 Tl 204.38	82 Pb 207.2	83 Bi 208.98	84 Po (209)	85 At (210)	86 Rn (222)
87 Fr (223)	88 Ra (226)	89-103 #	104 Rf (265)	105 Db (268)	106 Sg (271)	107 Bh (270)	108 Hs (277)	109 Mt (276)	110 Ds (281)	111 Rg (280)	112 Cn (285)	113 Nh (286)	114 Fl (289)	115 Mc (289)	116 Lv (293)	117 Ts (294)	118 Og (294)
	* Lanti seri	hanide es	57 La 138.91	58 Ce 140.12	59 Pr 140.91	60 Nd 144.24	61 Pm (145)	62 Sm 150.36	63 Eu 151.96	64 Gd 157.25	65 Tb 158.93	66 Dy 162.50	67 Ho 164.93	68 Er 167.26	69 Tm 168.93	70 Yb 173.05	71 Lu 174.97
	# Actir serie	iide s	89 Ac (227)	90 Th 232.04	91 Pa 231.04	92 U 238.03	93 Np (237)	94 Pu (244)	95 Am (243)	96 Cm (247)	97 Bk (247)	98 Cf (251)	99 Es (252)	100 Fm (257)	101 Md (258)	102 No (259)	103 Lr (262)
71 oxide materials																	

141 surfaces with Miller indexes ≤ 2

270 adsorption sites

Subgroup discovery: CO₂ activation by adsorption



Primary features





Subgroup discovery: Adsorbed CO2 properties

Subgroup discovery: Analysis of the OCO angle



sites delivering smaller angles (59 adsorption sites):

(energy of O 2*p* band maximum > -6.0 eV) AND (distance from O-site to first nearest cation > 1.8 Å) AND (distance from O-site to second nearest cation > 2.1 Å)



Most of the site delivering smaller OCO angles are on ionic (basic) materials

Subgroup discovery: Analysis of the C-O bond length



sites delivering larger *l*(CO) (33 sites):

(cation charge < 0.5e) AND (work function ≥ 5.2 eV) AND (distance from O site to second nearest cation ≥ 2.14 Å)

 $LaGaO_3$ – cathode material in high-temperature electrochemical CO_2 reduction;



 $KNbO_3$ – photocatalytic reduction of CO_2 into CH_4 ;

 $NaNbO_3$ – photocatalyst for CO_2 reduction with ~70% of CO selectivity;

 $NaSbO_3$ – material for CO_2 capture and storage (CCS)

Subgroup discovery: Alternative mechanisms of CO₂ activation



Longer C-O implies smaller OCO angles, but not too small \rightarrow no catalyst poisoning

Subgroup discovery: A different approach



Subgroup discovery: A different approach



large Hirshfeld charge on surface O, lower coordination for smaller angles

with adsorption energy constraint:

smaller charge on surface O, delocalized electron density, binding of O in CO₂ with surface cations

SISSO and SGD software

SISSO: https://github.com/rouyang2017/SISSO

Subgroup discovery: http://www.realkd.org/



Decision tree regression



Split criterion: $\sum (target property - \langle target property \rangle)^2 \rightarrow min within each subgroup$

Decision tree properties

- Simple to understand and interpret
- Global (important difference to subgroup discovery, which finds *locally unique* groups)
- Easy to overfit (can use LASSO-type penalty to solve this problem)
- Small change in data can lead to large change in the tree
- Relatively inaccurate

Random forest[®]

- 1) Perform tree regression or classification on several randomly selected subsets of data
- 2) In each tree, at each split choose randomly a fixed number of features, for which the best split is determined
- 3) Average predictions from the obtained trees
- **Properties:**
 - More accurate than a single tree ("each tree keeps other trees from making mistakes)
 - Interpretability of the model is lost
 - Can be use to select primary features for other approaches such as SISSO
Random forest®

Interesting application: Identify most important surface structural features that determine surface stability



Chemical Pressure-Driven Enhancement of the Hydrogen Evolving Activity of Ni₂P from Nonmetal Surface Doping Interpreted via Machine Learning

Robert B. Wexler,[†][©] John Mark P. Martirez,[‡][©] and Andrew M. Rappe^{*,†}[©]