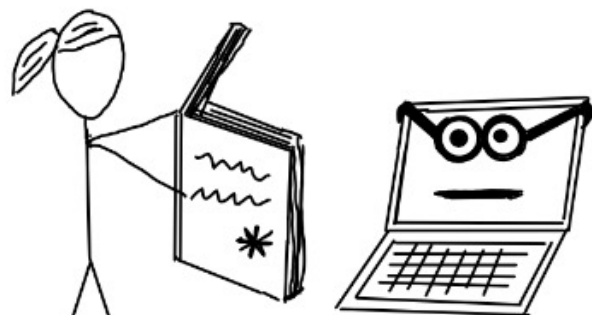


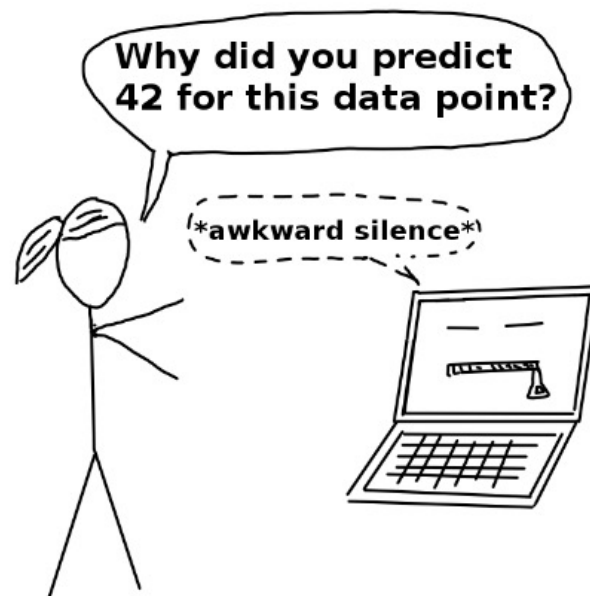
General approaches to interpretable ML

Without Machine Learning



* VERY SPECIFIC INSTRUCTIONS

With Machine Learning



Supervised learning

$$y \approx \hat{y} = f(\vec{x})$$

The equation shows y followed by an approximation symbol \approx , then \hat{y} (circled in orange), an equals sign, a function f (with a curly brace underneath), and \vec{x} (circled in orange). A small orange 'e' is written below the \vec{x} .

Why do we care about interpretability?

ACS Partner Journal

ACCOUNTS
—of materials research—

pubs.acs.org/amrcda



Article

Interpretable and Explainable Machine Learning for Materials Science and Chemistry

Felipe Oviedo,^{*,#} Juan Lavista Ferres, Tonio Buonassisi, and Keith T. Butler^{*,#}

Perspective | [Published: 17 March 2022](#)

Interpretable machine learning for knowledge generation in heterogeneous catalysis

[Jacques A. Esterhuizen](#), [Bryan R. Goldsmith](#) ✉ & [Suljo Linic](#) ✉

[Nature Catalysis](#) **5**, 175–184 (2022) | [Cite this article](#)

7787 Accesses | 21 Citations | 9 Altmetric | [Metrics](#)

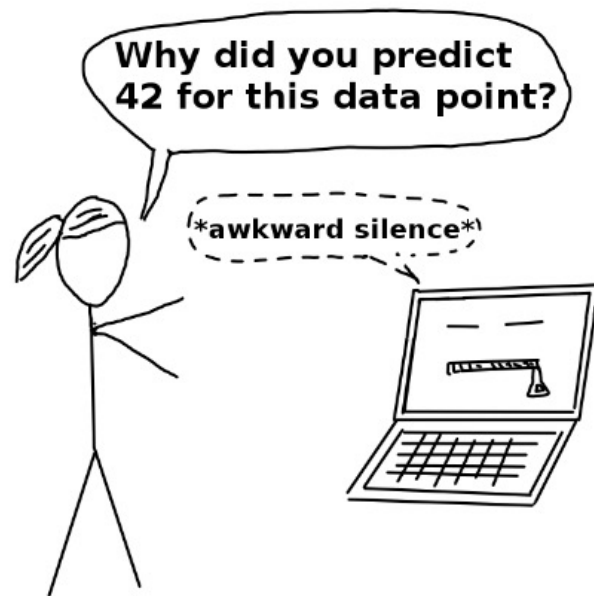
Review Article | [Published: 13 October 2020](#)

Drug discovery with explainable artificial intelligence

[José Jiménez-Luna](#), [Francesca Grisoni](#) & [Gisbert Schneider](#) ✉

[Nature Machine Intelligence](#) **2**, 573–584 (2020) | [Cite this article](#)

47k Accesses | 207 Citations | 95 Altmetric | [Metrics](#)



Why do we care about interpretability?

ACS Partner Journal

ACCOUNTS
—of materials research—

pubs.acs.org/amrcda



Article

Interpretable and Explainable Machine Learning for Materials Science and Chemistry

Felipe Oviedo,^{*,#} Juan Lavista Ferres, Tonio Buonassisi, and Keith T. Butler^{*,#}

Perspective | [Published: 17 March 2022](#)

Interpretable machine learning for knowledge generation in heterogeneous catalysis

Jacques A. Esterhuizen, Bryan R. Goldsmith & Suljo Linic

Nature Catalysis **5**, 175–184 (2022) | [Cite this article](#)

7787 Accesses | 21 Citations | 9 Altmetric | [Metrics](#)

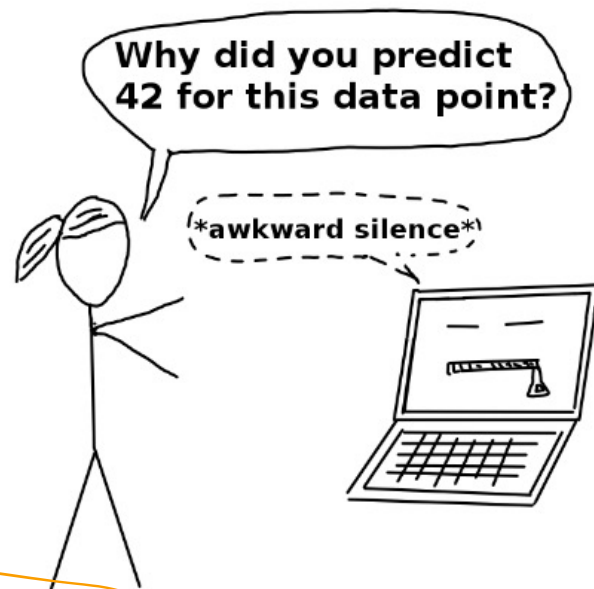
Review Article | [Published: 13 October 2020](#)

Drug discovery with explainable artificial intelligence

José Jiménez-Luna, Francesca Grisoni & Gisbert Schneider

Nature Machine Intelligence **2**, 573–584 (2020) | [Cite this article](#)

47k Accesses | 207 Citations | 95 Altmetric | [Metrics](#)



JOIN CHIME

Go to:

<https://chimein2.cla.umn.edu/join/608843>

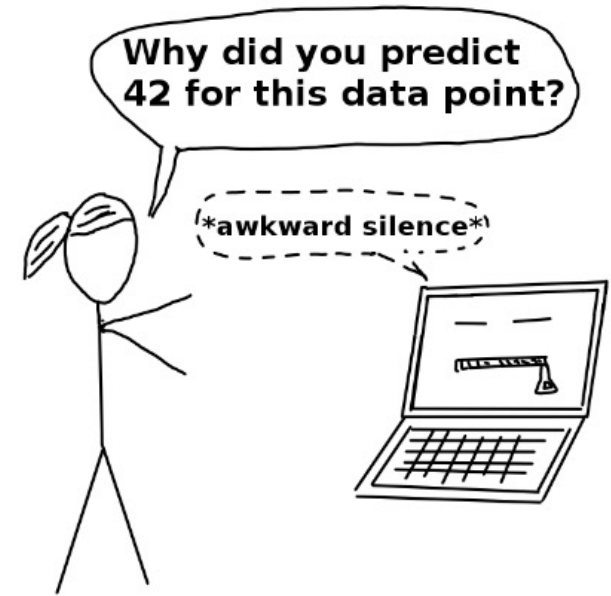
Or: visit chimein2.cla.umn.edu and enter
608-843

Why do we care about interpretability?

Reliability \rightarrow small change in x_i shouldn't lead to large changes in \hat{y}

Causality \rightarrow as x_i changes, can we anticipate $\Delta \hat{y}$

Trust \rightarrow adoption & understanding



Taxonomy of interpretability

Intrinsic \rightarrow models that are simple
(have few parameters)

Post-hoc \rightarrow analyze after training
(any black box model)

Local \rightarrow why was a prediction made (\hat{y}_i)

Global \rightarrow understand general behavior over domain (\vec{x})

Intrinsic interpretability in linear regression

$$\hat{y}(w, x) = w_0 + w_1 x_1 + \dots + w_p x_p \quad \min_w \|Xw - y\|_2^2 + \alpha \|w\|_1$$

Assumptions:

- ↳ linear relationship b/f y & x_j
- ↳ residuals are normally distributed
- ↳ features are independent

Interpretation:

- ↳ $|w_j x_j|$ = feature effect
- ↳ magnitude of each feature effect is intrinsically important



Intrinsic interpretability in linear regression

$$\hat{y}(w, x) = w_0 + w_1 x_1 + \dots + w_p x_p \quad \min_w \|Xw - y\|_2^2 + \alpha \|w\|_1$$

Advantages:

- ↳ simple
- ↳ easy to fit

Disadvantages:

- ↳ miss nonlinearities
- ↳ may be inaccurate (fit low dimension)
- ↳ confused by correlated features

Intrinsic interpretability w/ decision trees

A

```

graph TD
    A["(40, 40)"] --> B["(30, 10)"]
    A --> C["(10, 30)"]
    
```

B

```

graph TD
    B["(40, 40)"] --> D["(20, 40)"]
    B --> E["(20, 0)"]
    
```

Impurity = $-\sum_{i=0}^N p(i) \ln(p(i))$ for N classes

$y = (40, 40)$
 $x = [0, 0, 1, 1, 0, 1]$

$I_D = -2(0.5 \ln 0.5) = 0.69$

$I_A = -(\frac{3}{4} \ln \frac{3}{4} + \frac{1}{4} \ln \frac{1}{4}) = 0.56$

$I_B = -(\frac{3}{4} \ln \frac{1}{3} + \frac{2}{3} \ln \frac{2}{3}) = 0.48$

Information gain
(weighted impurity decrease)

0.13 0.21

$x \ln x + (1-x) \ln (1-x)$

JOIN CHIME

Go to:

<https://chimein2.cla.umn.edu/join/608843>

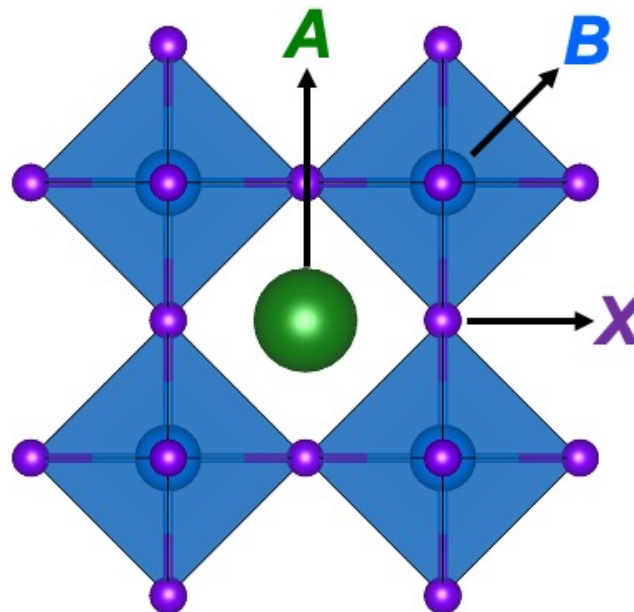
Or: visit chimein2.cla.umn.edu and enter
608-843



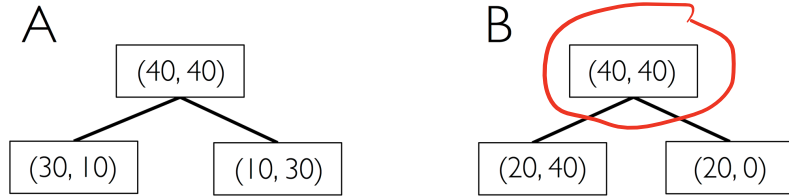
New tolerance factor to predict the stability of perovskite oxides and halides

Christopher J. Bartel^{1*}, Christopher Sutton², Bryan R. Goldsmith³, Runhai Ouyang², Charles B. Musgrave^{1,4,5}, Luca M. Ghiringhelli^{2*}, Matthias Scheffler²

$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)}$$



Intrinsic interpretability w/ tree-based methods



Advantages:

- ↳ natural (visual) representation
- ↳ allows for "what-if" explanations.

Disadvantages:

- ↳ can be too complex
- ↳ linear relationships
- ↳ unstable

From model-specific to model-agnostic methods

Model-specific

↳ linear reg ($|w_j x_j|$)

↳ decision trees ($|G|$)

Model-agnostic

↳ interpretations based on predictions (by an arbitrary model)

↳ Global → describe avg behavior

↳ Local → explains individual predictions

One global method: permutation importances

Idea: how much worse our model gets
when we randomly shuffle a feature?

1) Train a model, $\hat{y} = f(\vec{x})$

2) Compute the error, $\|y - \hat{y}\|^2$ (e_{orig})

3) For feature j in X :

a) permute x_j

b) compute e_j

4) Feature importance, $PI = \left| \frac{e_j}{e_{orig}} \right|$

$$y = 3x_1 + 2x_2^2$$



Compute on training data or validation data?

⚡ Training → how does training manifest
in dependencies between $\|y - \hat{y}\|^2$ & X_j

Validation → how essential is each feature
for generalization?

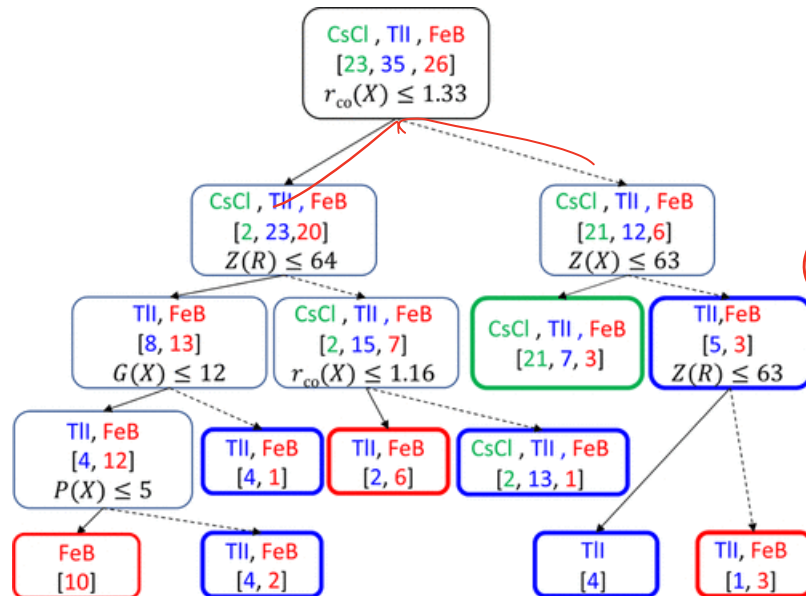
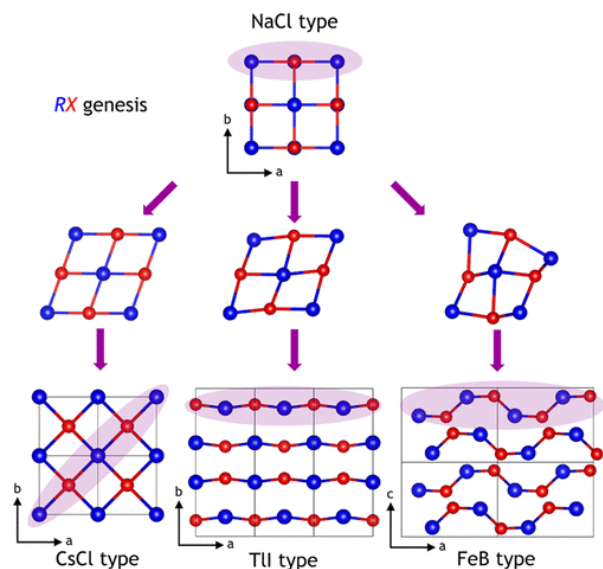
One global method: permutation importances

Revealing Hidden Patterns through Chemical Intuition and Interpretable Machine Learning: A Case Study of Binary Rare-Earth Intermetallics *RX*

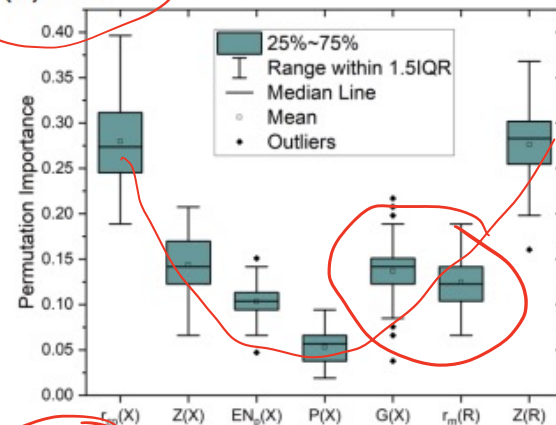
Volodymyr Gvozdetyskiy,* Balaranjan Selvaratnam, Anton O. Oliynyk, and Arthur Mar*

Cite This: *Chem. Mater.* 2023, 35, 879–890

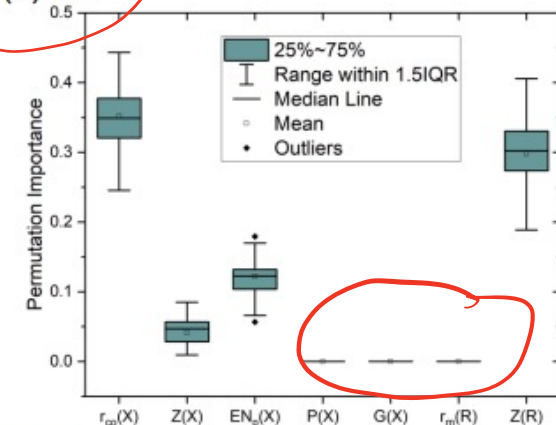
Read Online



(a) SVC



(b) DT



One global method: permutation importances

Advantages:

↳ easy to interpret

↳ does not require re-training

Disadvantages:

↳ permutations may lead to unphysical combinations of features

↳ to get robust estimates, we need many repeats

Quick survey of other global methods

Partial dependencies

↳ how does \hat{y} vary wrt a single feature, x_j

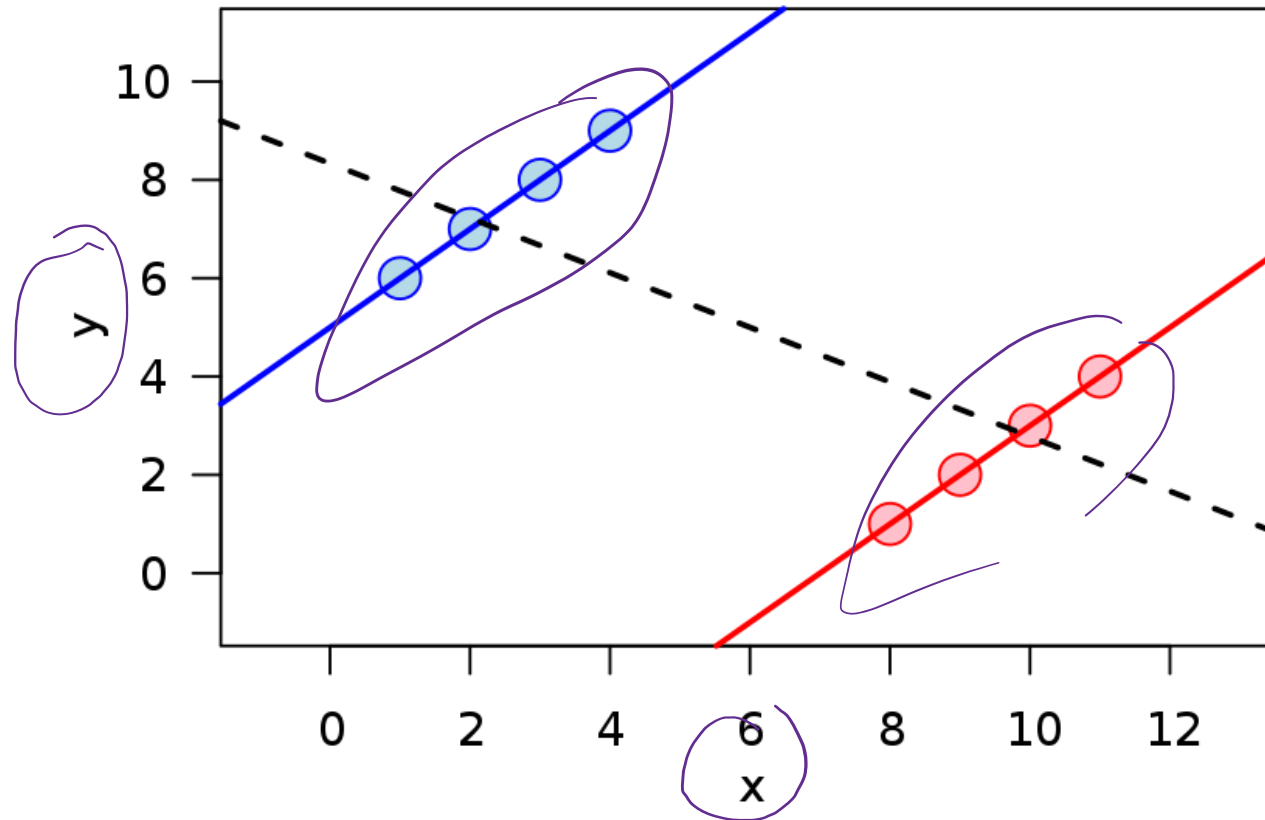
Feature interaction

↳ how does the value of one feature influence the partial dependence of another feature.

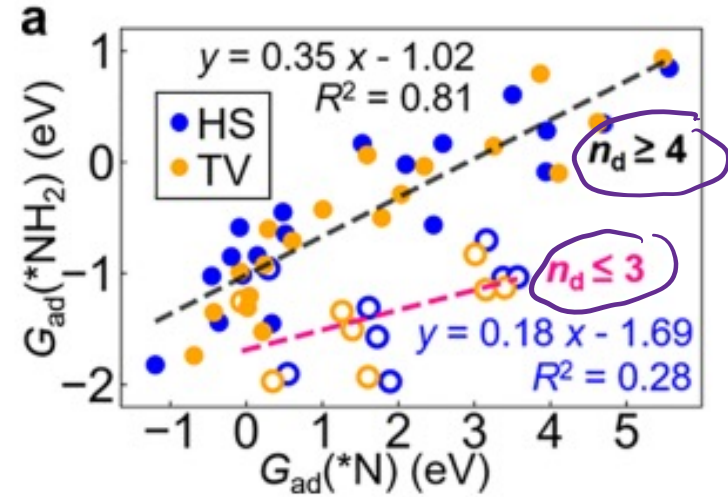
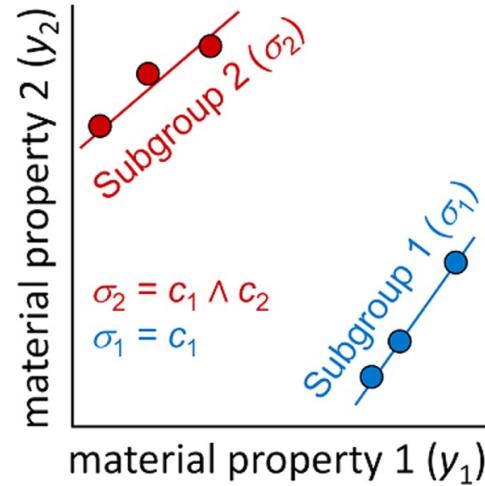
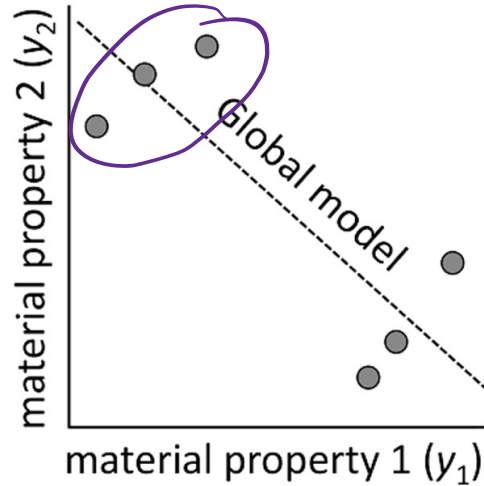
Surrogate model

↳ train an intrinsically model on my black box model's predictions, then interpret that model

From global to local methods



Our data may not be globally interpretable

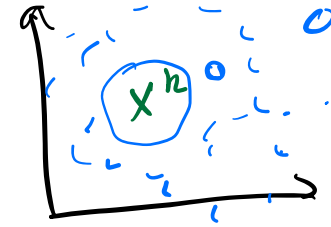


Subgroup discovery \rightarrow a systematic way to find related "sub" populations

Local version of surrogate: LIME

LIME \rightarrow Local Interpretable Model-agnostic Explanations

Given blackbox model, $\hat{y} = f(\vec{x})$



- 1) Select a point to probe, x^k
- 2) Randomly generate new data points
- 3) Predict target (\hat{y}) for new points using $f(\vec{x})$
- 4) Weight each new pt based on proximity to x^k
- 5) Train an interpretable model on new pts w/ weighted loss function
- b) Interpret

$$\propto \|y_i - \hat{y}_i\|^2$$

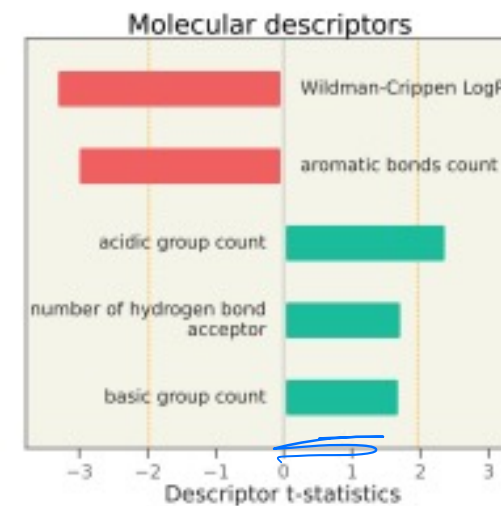
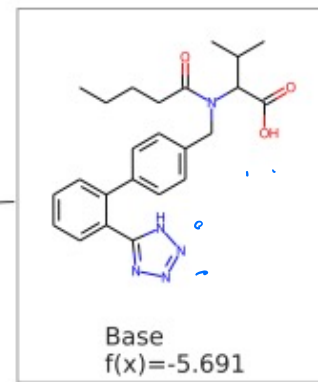
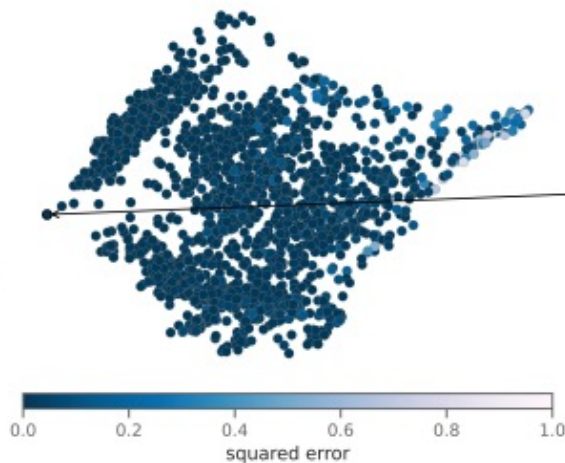
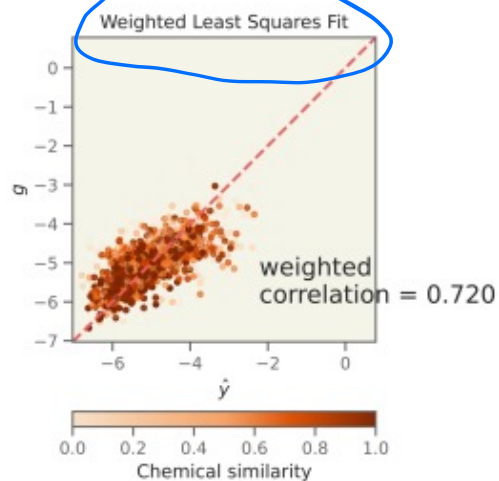
Local version of surrogate: LIME

EXPLAINING MOLECULAR PROPERTIES WITH NATURAL LANGUAGE

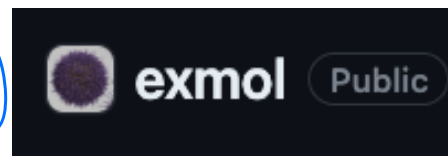
A PREPRINT

© Heta A. Gandhi
Department of Chemical Engineering
University of Rochester
Rochester, NY, 14627
hgandhi@ur.rochester.edu

© Andrew D. White*
Department of Chemical Engineering
University of Rochester
Rochester, NY, 14627
andrew.white@rochester.edu



<https://github.com/ur-whitelab/exmol>



Another local method: counterfactuals



Cite this: Chem. Sci., 2022, 13, 3697

All publication charges for this article

Model agnostic generation of counterfactual explanations for molecules†

Geemi P. Wellawatte,^a Aditi Seshadri^b and Andrew D. White^{id} ^{★b}

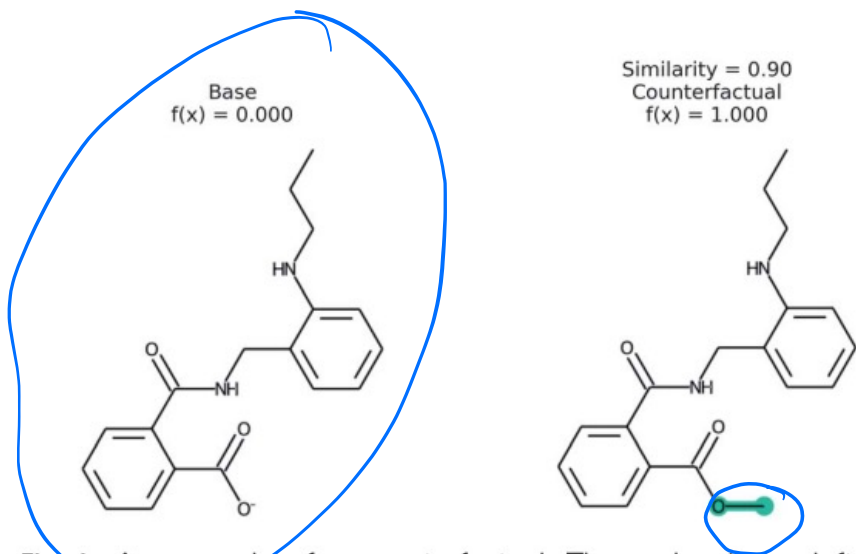
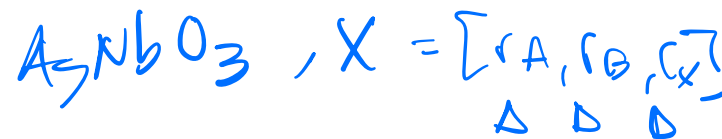
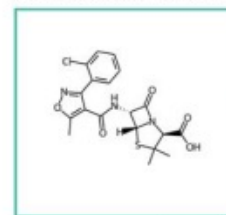
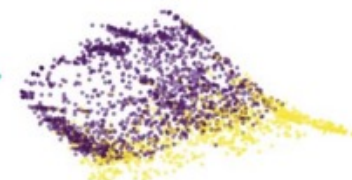


Fig. 1 An example of a counterfactual. The molecule on left was predicted to have class of 0, no activity. With the modification shown in teal, the molecule would be in class 1, active. This shows that the carboxylic acid is an explanation for lack of activity.

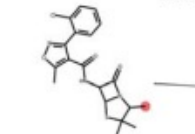
1. Molecule being predicted: base



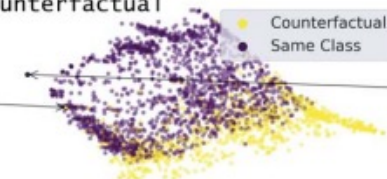
2. Expand chemical space around base



3. Identify most similar molecule with changed label: counterfactual



Counterfactual



Base

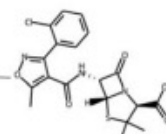


Fig. 2 Overview of MMACE. The input is a molecule to be predicted. Chemical space is expanded and clustered. Counterfactuals are selected from clusters to find succinct explanation of base molecule prediction.

Another local/global method: Shapley values and SHAP

Idea: Let features play a "game" to determine feature effects

Given: $\hat{y} = f(\vec{X})$

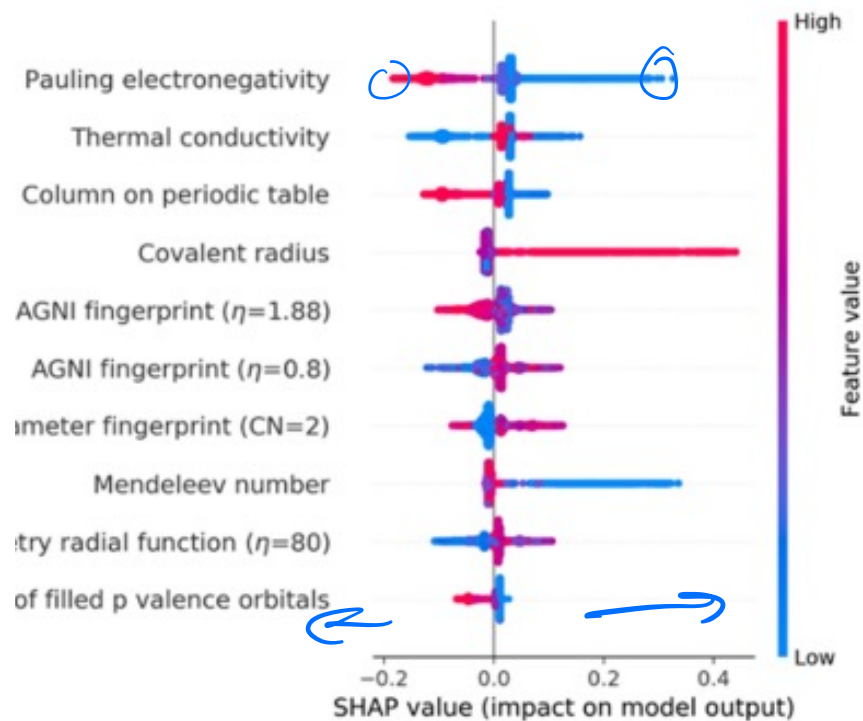
Question: How much does X^j contribute to each pts deviation from the mean

$$\Delta y_i = \langle \hat{y} \rangle - y_i$$

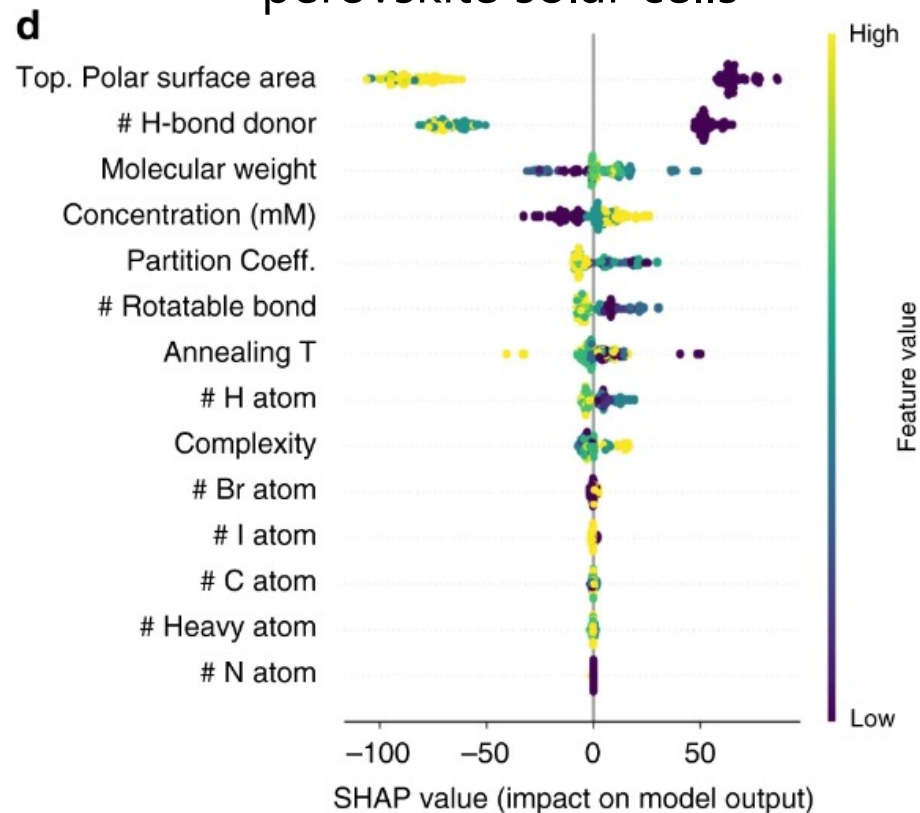
- 1) Pick a feature & data pt (X^j, y_k)
 - 2) Draw a random subset including X^j (of features)
 - 3) Compute $\hat{y} = f(\vec{X})$
 - 4) Repeat this but w/o X^j
 - 5) Repeat over all features \rightarrow subsets \rightarrow data pts
- $\rightarrow \approx$ marginal contribution to \hat{y}

Another local/global method: Shapley values and SHAP

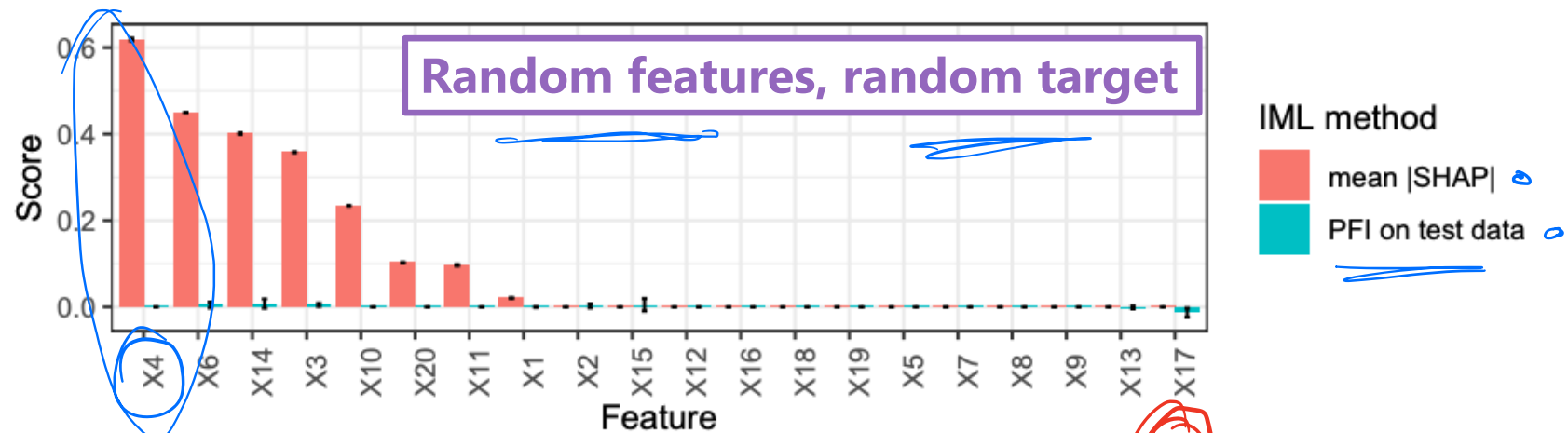
Partial charges in
nanoporous materials



Capping layers for hybrid halide
perovskite solar cells



1. Assuming one method will always work

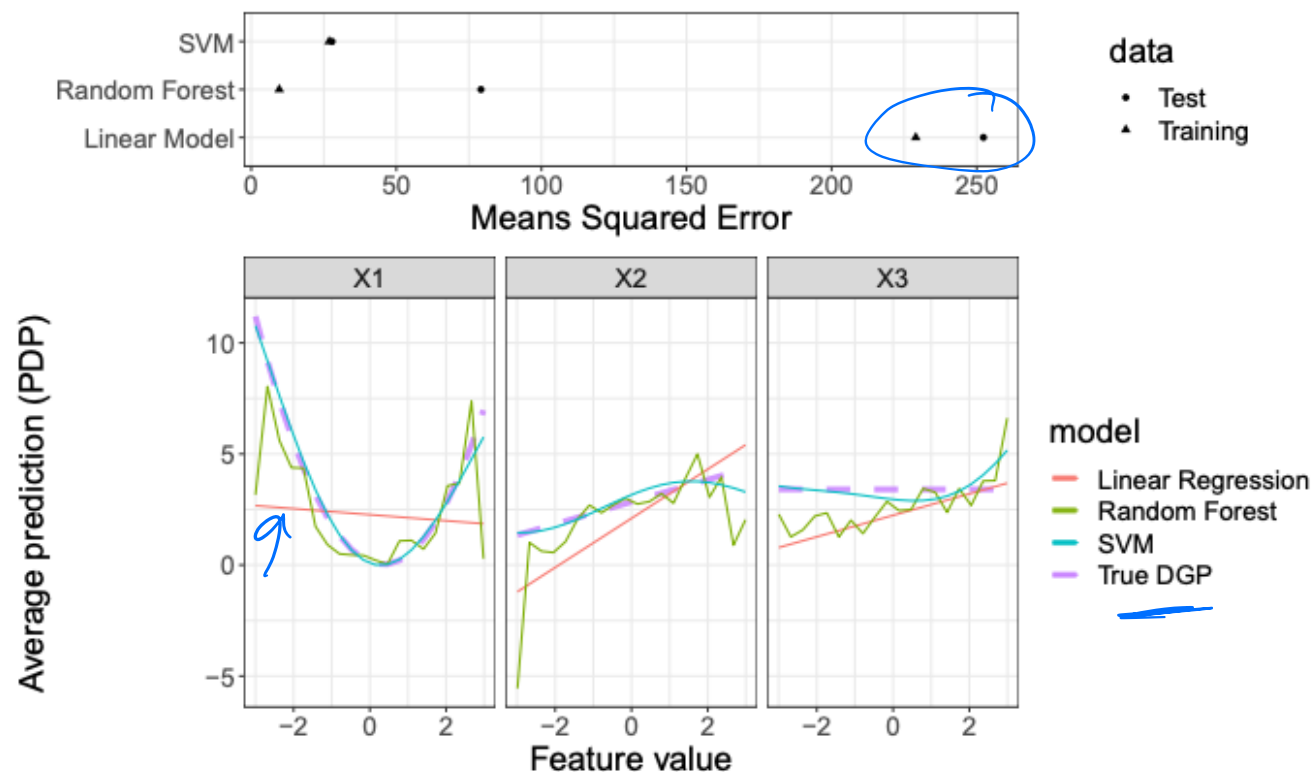


SHAP \rightarrow which x_i contribute

PFI \rightarrow which x_i control $\|y - \hat{y}\|^2$ (on test)

2. Bad model generalization

GIGO



3. Unnecessary complexity

Recall why we care about interpretability:

Reliability \rightarrow small change in x_i shouldn't lead to a large change in \hat{y}_i

Causality \rightarrow as we change x_i , can we anticipate change in \hat{y}_i

Trust \Rightarrow adoption & understanding



Intrinsic interpretability is always preferred!