

Statistical Learning for Halide Perovskite Discovery

Panayotis Manganaris¹

¹Purdue Materials Science and Engineering Mannodi Group

July 1, 2022

Outline

- 1 AI Background
- 2 Chemistry Background
- 3 Pipeline
- 4 Feature Engineering
- 5 Supervised Architectures

Artificial Intelligence

The Four Approached to AI

Thinking Humanly

- Turing test approach (The Six Fields of AI)¹
- NLP
- Knowledge Representation
- automated reasoning
- Machine Learning
- computer vision
- robotics

Thinking Rationally

- Laws of Thought
- logical positing
- proven algorithms
- correct inference
- syllogistic reason

Acting Humanly

- cognitive modeling approach
- neuromorphic algorithms

Acting Rationally

- The rational agent
- inference + reflex
- inference vs deduction

^aStuart Russell and Peter Norvig. *Artificial intelligence : a modern approach*. Upper Saddle River, New Jersey: Prentice Hall, 2010. ISBN: 9780136042594

Machine Learning I

ML Contributes to AI

- Adaptable **agent**
 - Contextual judgment of **percept** relevance
 - Autonomous utilization of **percept sequence**
- Learning
 - **function** performance improves with exposure to more percepts

Definition (Artificial Agency)

agent self-contained sensor->function->action pipeline

function Set of all possible responses for all possible percepts

percept sensory input

percept sequence history of sensory input

Machine Learning II

Supervised Training

Encourage the agent to behave "correctly"

- 1 Minimize Loss
- 2 Maximize Score

Unsupervised Training

The agent determines something principally true about its environment using mathematical/logical characterization methods.

- find eigenvectors and eigenvalues
- differentially calculate optima

Inverse Design

A Type of AI Implementation

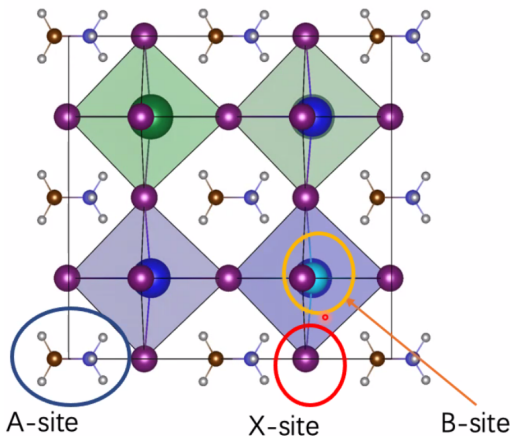
senses maps points in many dimensions

function reliably navigates it's environment searching for optima

action returns its findings to human interpreters

Perovskite Structure and Chemistry

Example² of hybrid organic-inorganic MAPbI₃



²Arun Mannodi-Kanakkithodi and Maria K. Y. Chan. “Data-Driven Design of Novel Halide Perovskite Alloys”. In: *Energy Environ. Sci.* 15 (5 2022), pp. 1930–1949.

DOI: 10.1039/D1EE02971A. URL: <http://dx.doi.org/10.1039/D1EE02971A>

Our Dataset

DFT Simulations

- 1 geometry optimization
- 2 Static band structure and optical absorption

Levels of Theory

- PBE
- HSE06
- PBE+HSE06(SOC)
- Experimental

Formula	bg _{ev}	η	LoT
MAPbCl ₃	3.03	0.002	EXP
CsPbI _{0.375} Br _{2.625}	1.68	0.153	PBE
RbSnBr _{2.625} Cl _{0.375}	1.44	NaN	HSE
CsGeCl ₃	1.05	0.176	PBE
MASr _{0.5} Pb _{0.5} Cl ₃	5.31	NaN	HSE
MABa _{0.25} Pb _{0.75} I ₃	1.99	0.015	PBE
MASnI ₃	2.57	NaN	HSE
MACa _{0.5} Pb _{0.5} Cl ₃	5.32	NaN	HSE
...

Band Gap Fidelity I

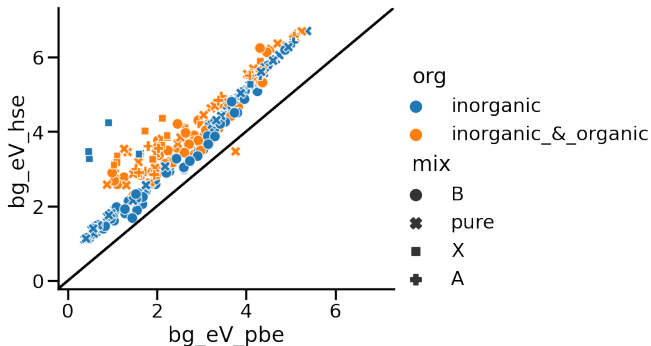


Figure 1: PBE vs HSE Band Gaps

Band Gap Fidelity II

Comparing computational with experimental³ band gaps

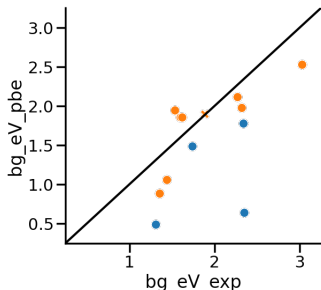


Figure 2: PBE vs Almora BG

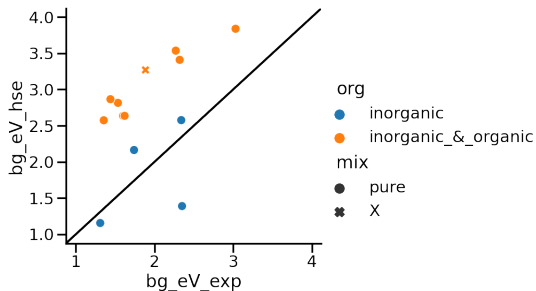


Figure 3: HSE vs Almora BG

³Osbel Almora et al. "Device Performance of Emerging Photovoltaic Materials (Version 1)". In: *Advanced Energy Materials* 11.11 (2020), p. 2002774. DOI: 10.1002/aenm.202002774. URL: <http://dx.doi.org/10.1002/aenm.202002774>

Data Pre-Processing

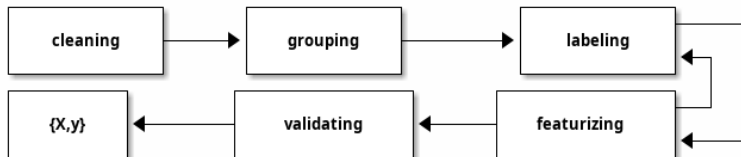


Figure 4: Data Preprocessing Workflow to Implement with Python Pandas

Machine Learning Pipeline

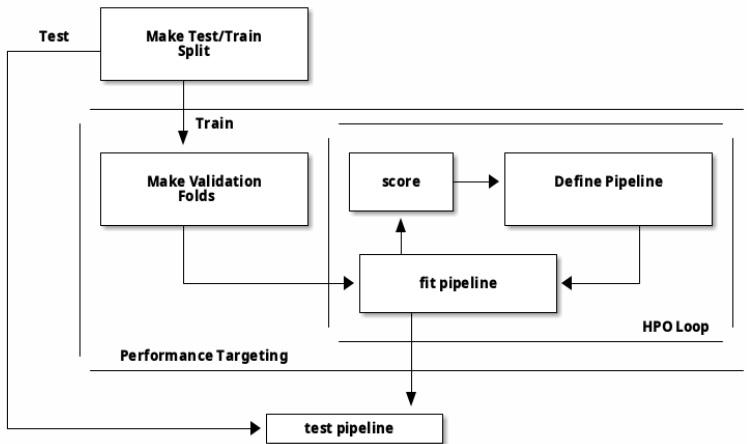


Figure 5: Machine Learning Pipeline to Implement with Python SciKit-Learn

Implementation in Jupyter Python I

```
import sys, os
sys.path.append(os.path.expanduser("~/src/cmcl"))
sys.path.append(os.path.expanduser("~/src/spyglass"))
import pandas as pd
import numpy as np
import cmcl
from spyglass.model_imaging import parityplot
from sklearn.pipeline import make_pipeline
from sklearn.compose import ColumnTransformer
from sklearn.<module> import NumPreProcessor1
from sklearn.<module> import CatPreProcessor1
from sklearn.<module> import NumPreProcessor2
from sklearn.<module> import CatPreProcessor2
from sklearn.<module> import Estimator

df = pd.read_csv('file.<data>')
df = df.groupby('Formula', as_index=False).agg(
    {'bg_eV': 'median', 'efficiency': 'median'})
```

Listing 1: Scikit-Learn mock-setup and pandas data loading + grouping

Implementation in Jupyter Python II

```
dc = df.ft.comp()
dc = dc.assign(label='label')

numeric_features = dc.select_dtypes(np.number).columns.to_list()
numeric_pipeline = make_pipeline(NumPreProcessor1(),
                                  NumPreProcessor2())
categorical_features = dc.select_dtypes('object').columns.to_list()
categorical_pipeline = make_pipeline(CatPreProcessor1(),
                                     CatPreProcessor2())
```

Listing 2: Computing Features with cmcl and Feature Engineering

```
ss = ShuffleSplit(n_splits=1, train_size=0.8, random_state=None)
train_idx, test_idx = next(ss.split(dc))
dc_tr, dc_ts = dc.iloc[train_idx], dc.iloc[test_idx]
df_tr, df_ts = df.iloc[train_idx], df.iloc[test_idx]

preprocessor = ColumnTransformer(
    transformers=[
        ("num", numeric_pipeline, numeric_features),
        ("cat", categorical_pipeline, categorical_features),
    ]
)
pipe = make_pipeline(preprocessor, Estimator())
pipe.fit(dc_tr, df_tr.<target>)
```

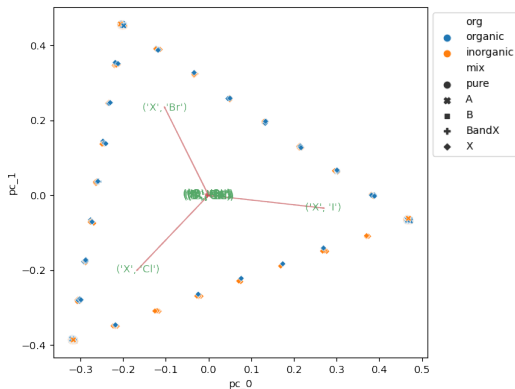
Listing 3: make Test/Train Split, Assemble, and Fit pipeline

Implementation in Jupyter Python III

```
p, data = parityplot(pipe,
                      dc_ts, df_ts.<target>.to_frame(),
                      aspect=1.0)
p.figure.show()
```

Listing 4: Evaluate Pipeline using Spyglass

PCA



$$UAU^\dagger = Q^{-1}SQ$$

Figure 6: Learn transformation matrix U to diagonalizes the matrix A . The Principal Components in Q corresponding to the largest two Singular Values in S contain the majority of the variance in the data.

tSNE

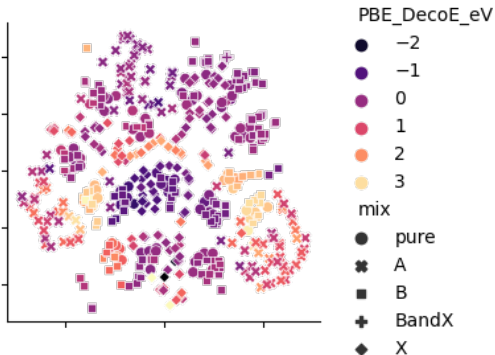


Figure 7: Learn a low-dimensional (2 or 3D) embedding space in which statistical similarity governs the proximity of high-dimensional data points

UMAP

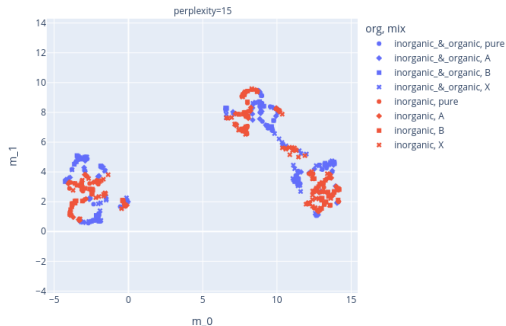


Figure 8: Learn a manifold embedding space in which nearest neighbors form clusters

Linear regressions BG Test I

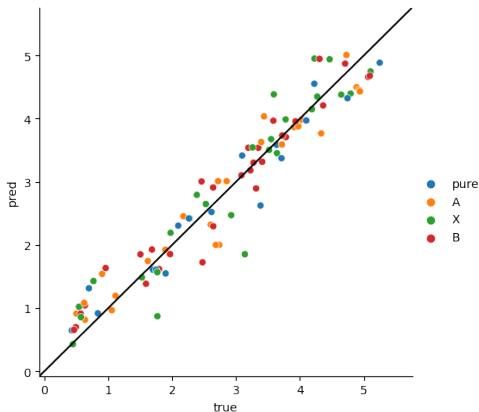


Figure 9: OLS determines \vec{w} so that $f(x) = \vec{x}^T \vec{w}$, $y_i = f(x_i) + \epsilon_i$ and all ϵ_i are as small as possible

Linear regressions BG Test II

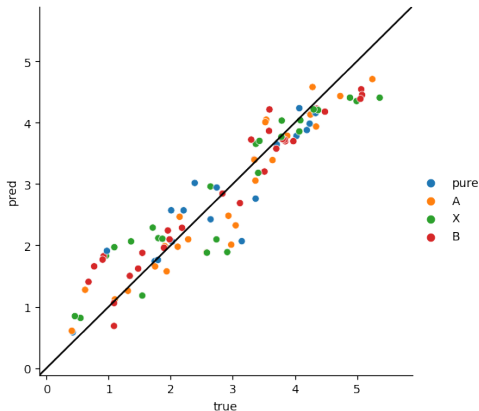


Figure 10: elasticnet determines \vec{w} as before, but also works to sparsify the model

OLS weights

site	element	
A	Cs	23.771206
A	FA	25.794831
A	K	22.774475
A	MA	25.452629
A	Rb	23.282988
B	Ba	-32.603053
B	Ca	-31.378385
B	Ge	-45.001044
B	Pb	-42.526511
B	Sn	-46.868114
B	Sr	-32.068490
X	Br	0.939374
X	Cl	1.769032
X	I	0.140658

	RSS
A	54.213044
B	95.426246
X	2.007905

elasticnet weights

site	element	
A	Cs	-0.191057
A	FA	1.589015
A	K	-1.081903
A	MA	1.214167
A	Rb	-0.530437
B	Ba	5.139688
B	Ca	6.424156
B	Ge	-5.879154
B	Pb	-3.673012
B	Sn	-7.689152
B	Sr	5.678253
X	Br	0.000000
X	Cl	0.819669
X	I	-0.786942

	RSS
A	2.342552
B	14.391222
X	1.136281

Random Forest Regression on BG I

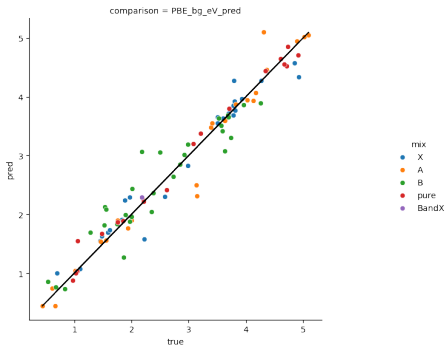


Figure 11: RFR initializes an ensemble of Decision Trees and averages their results to return its prediction. This leverages the DT's ability to strongly bias itself to the data and relies on randomness to explain variance in the underlying process

RFR Feature Importance

Gaussian Process or BG I

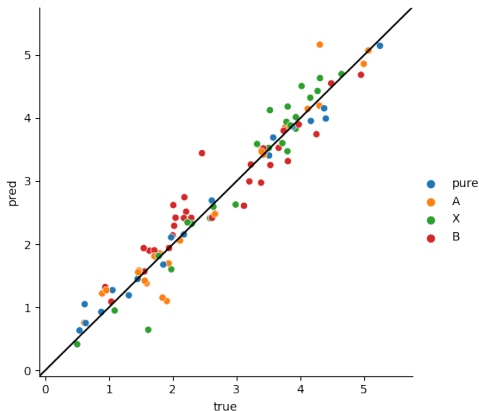


Figure 12: GPR picks functions from a distribution derived from the data covariance. The functions that satisfy the data form the fit.

Gaussian Process or BG II

Regularization with Priors

Conditional Probability $P(x|y) = \frac{P(x)P(y|x)}{P(x)}$

Conditional Odds $O(x|y) = O(x) \frac{P(x|y)}{P(x|\neg y)}$

Isolated Bayesian Prior $B = \frac{P(x|y)}{P(x|\neg y)}$

Comparing Test Scores

Table 1: Optimized Models Quantitative Performance Comparison

	OLS	elasticnet	RFR	GPR
r2	0.750403	0.750403	0.975133	
ev	0.751708	0.751708	0.975228	
maxerr	-2.984809	-2.984809	-0.873043	
rmse	-0.770263	-0.770263	-0.207796	
A _{rmse}	-0.793173	-0.793173	-0.143580	
B _{rmse}	-0.720543	-0.720543	-0.311177	
X _{rmse}	-0.703632	-0.703632	-0.130350	
Pure _{rmse}	-0.849593	-0.849593	-0.147384	



Almora, Osbel et al. “Device Performance of Emerging Photovoltaic Materials (Version 1)”. In: *Advanced Energy Materials* 11.11 (2020), p. 2002774. DOI: 10.1002/aenm.202002774. URL: <http://dx.doi.org/10.1002/aenm.202002774> (cit. on p. 10).



Mannodi-Kanakkithodi, Arun and Maria K. Y. Chan. “Data-Driven Design of Novel Halide Perovskite Alloys”. In: *Energy Environ. Sci.* 15 (5 2022), pp. 1930–1949. DOI: 10.1039/D1EE02971A. URL: <http://dx.doi.org/10.1039/D1EE02971A> (cit. on p. 7).



Russell, Stuart and Peter Norvig. *Artificial intelligence : a modern approach*. Upper Saddle River, New Jersey: Prentice Hall, 2010. ISBN: 9780136042594 (cit. on p. 3).