

# Tutorial MD00: Machine Learning in Materials Sciences — From Basic Concepts to Active Learning

**Arun Mannodi Kanakkithodi**

School of Materials Engineering  
Purdue University, West Lafayette, IN

Monday, Apr 10, 2023, 8.00 am – 5:00 pm

Email: [amannodi@purdue.edu](mailto:amannodi@purdue.edu)

# Tutorial Outline

## **8.00 am – 9.30 am, Arun Mannodi Kanakkithodi:**

Introduction to ML for materials science. Demonstration of training ML models using a DFT dataset of halide perovskite alloys.

## **9.30 am – 10.00 am: Break**

## **10.00 am – 10.30 am, Austin McDannald:**

Gaussian Process Regression: Detailed description using examples.

## **10.30 am – 11.30 am, Austin McDannald:**

Discussion of active learning, Bayesian optimization, autonomous experiments.

## **11.30 am – 12.00 pm: General Discussion**

## **12.00 pm – 1.30 pm: Lunch Break**

## **1.30 pm – 2.30 pm, Saaketh Desai:**

Overview of neural networks for prediction, convolutional neural networks for images.

## **2.30 pm – 3.00 pm: General Discussion**

## **3.00 pm – 5.00 pm (and beyond):**

Battery Informatics and ML Hackathon.

# Tutorial Instructors

- Arun Mannodi Kanakkithodi: Assistant Professor, Materials Engineering, Purdue University.  
Computational materials scientist using high-throughput DFT and ML for materials design.
- Austin McDannald: Materials Research Engineer, National Institute of Standards and Technology
- Saaketh Desai: Postdoctoral Researcher, Sandia National Laboratory
- Shijing Sun: Research Scientist, Energy and Materials, Toyota Research Institute

# MRS Bulletin Article in Press

## A Framework for Materials Informatics Education through Workshops

Arun Mannodi-Kanakkithodi<sup>1</sup>, Austin McDannald<sup>2</sup>, Shijing Sun<sup>3</sup>, Saaketh Desai<sup>4</sup>, Keith Brown<sup>5</sup> and A. Gilad Kusne<sup>2</sup>

### Abstract

The burgeoning field of materials informatics necessitates a focus on educating the next generation of materials scientists in the concepts of data science, artificial intelligence (AI), and machine learning (ML). In addition to incorporating these topics in undergraduate and graduate curricula, regular hands-on workshops present the most effective medium to initiate researchers to informatics and have them start applying the best AI/ML tools to their own research. With the help of the Materials Research Society (MRS), members of the MRS AI staging committee, and a dedicated team of instructors, we successfully conducted workshops covering the essential concepts of AI/ML as applied to materials data, at both the spring and fall meetings in 2022, with plans to make this a regular feature in future meetings. In this article, we discuss the importance of materials informatics education via the lens of these workshops, including details such as learning and implementing specific algorithms, the crucial nuts and bolts of ML, and using competitions to increase interest and participation.

# PART 1:

# Introduction to ML in Materials Science

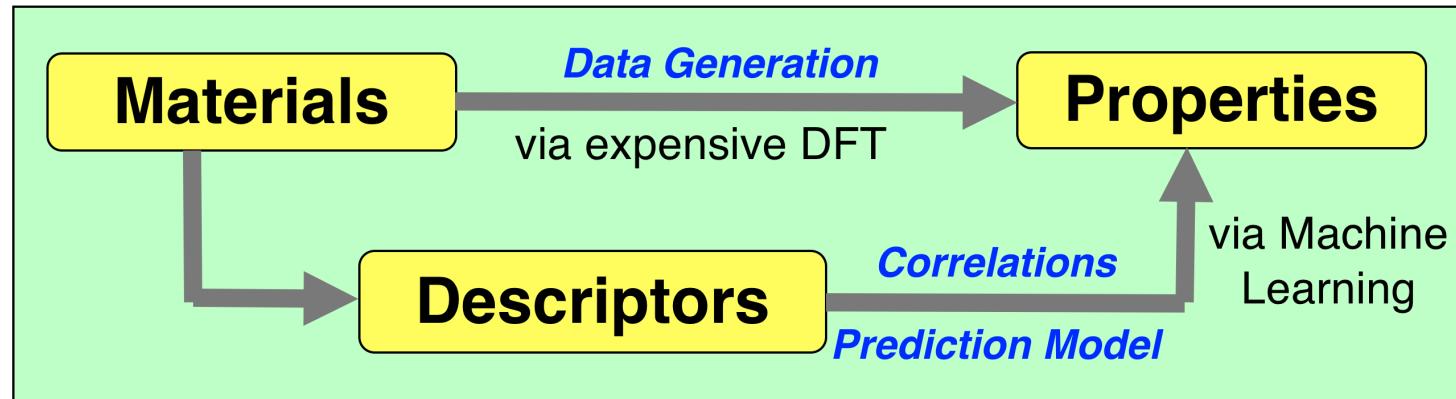


**PURDUE**  
UNIVERSITY®

School of Materials Engineering

4/10/23

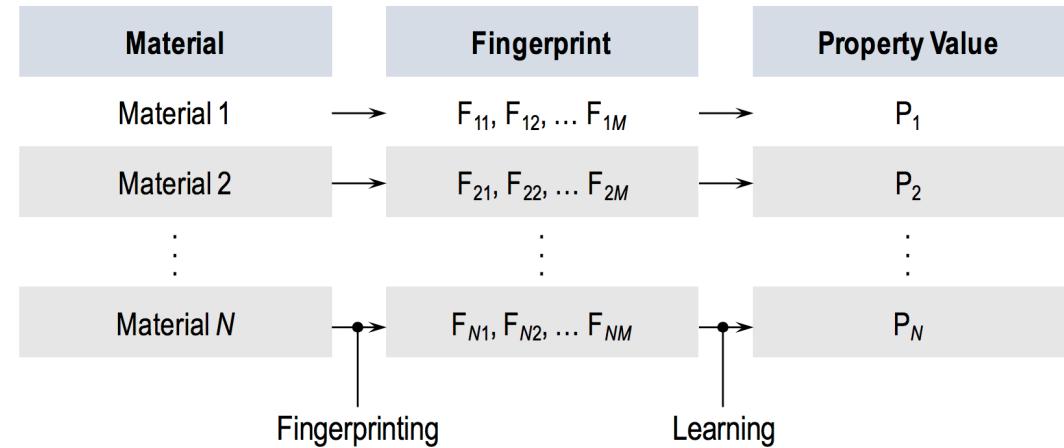
# Machine Learning in Materials Science



a Example dataset

Material	Property Value
Material 1	$P_1$
Material 2	$P_2$
⋮	⋮
Material $N$	$P_N$

c Fingerprinting, learning and prediction



b The learning problem

Material	Property Value
Material X	?

$$f(F_{i1}, F_{i2}, \dots, F_{iN}) = P_i$$

# Key Ingredient of ML: Feature Vectors / Materials Descriptors / Fingerprints

- Numerical representation of materials, input to ML.
- Definition depends on: a) application, b) domain expertise, and c) accuracy desired.
- Requirements: a) intuitive and inexpensive to calculate, b) generalizable to every material in the chemical space, and c) invariant to translation / rotation / permutation of like elements.



**PURDUE**  
UNIVERSITY®

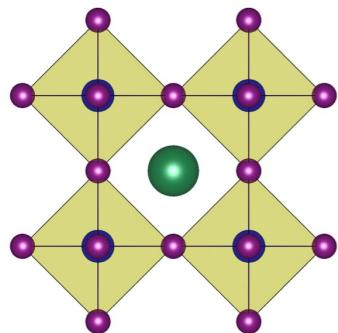
School of Materials Engineering

4/10/23

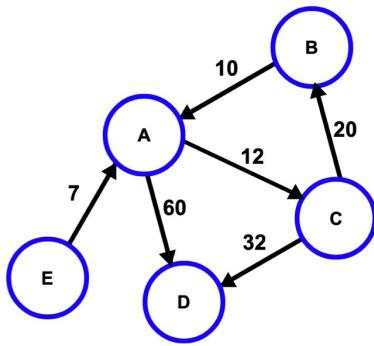
6

# Examples of Fingerprints

3D geometry:  
Atom i = (Z,x,y,z)



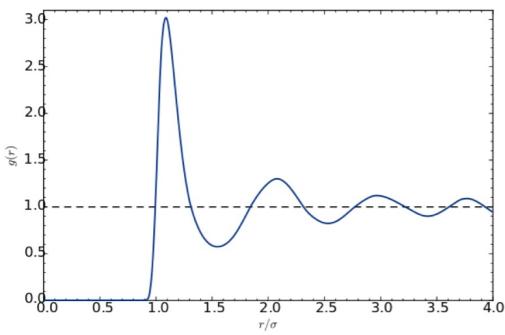
Weighted graph:  
atoms & bonds



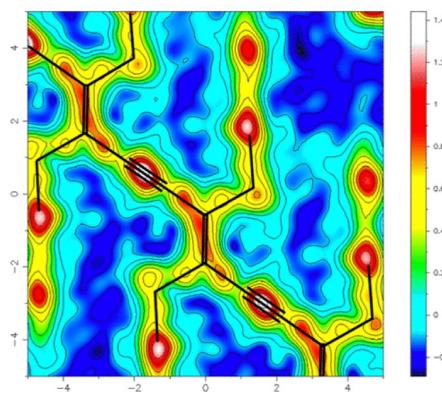
Coulomb Matrix

$$M_{IJ} = \begin{cases} 0.5 Z_I^{2.4} & \text{for } I = J \\ \frac{Z_I Z_J}{|\mathbf{R}_I - \mathbf{R}_J|} & \text{for } I \neq J \end{cases}$$

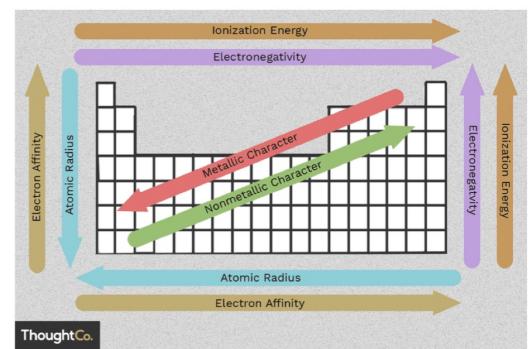
Radial Distribution Function



Electron Density Distribution



Tabulated elemental properties



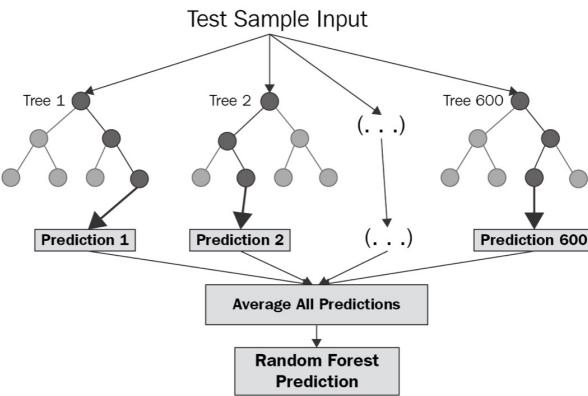
PURDUE  
UNIVERSITY®

School of Materials Engineering

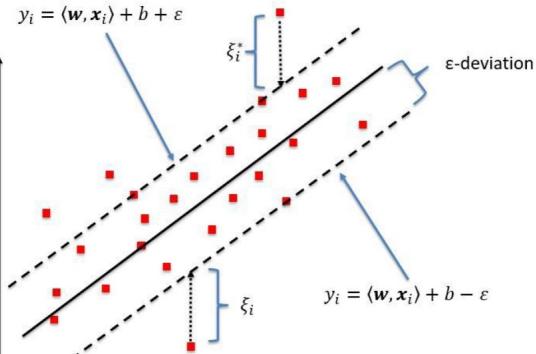
4/10/23

# Examples of ML Techniques

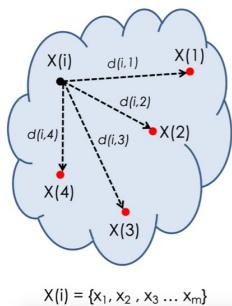
## Random Forest Regression



## Support Vector Regression



## Kernel Ridge Regression



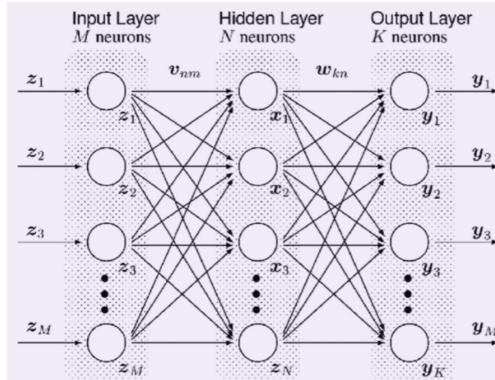
Measure of Similarity: Euclidean Distance

$$d(i,j) = \sqrt{(x_{i1} - x_{j1})^2 + \dots + (x_{im} - x_{jm})^2}$$

Property = Weighted sum of Gaussians

$$f(i) = \sum_{k=1}^N a_k \cdot \exp\left(-\frac{1}{2\sigma^2} \cdot [d(i, i_k)]^2\right)$$

## Neural Networks



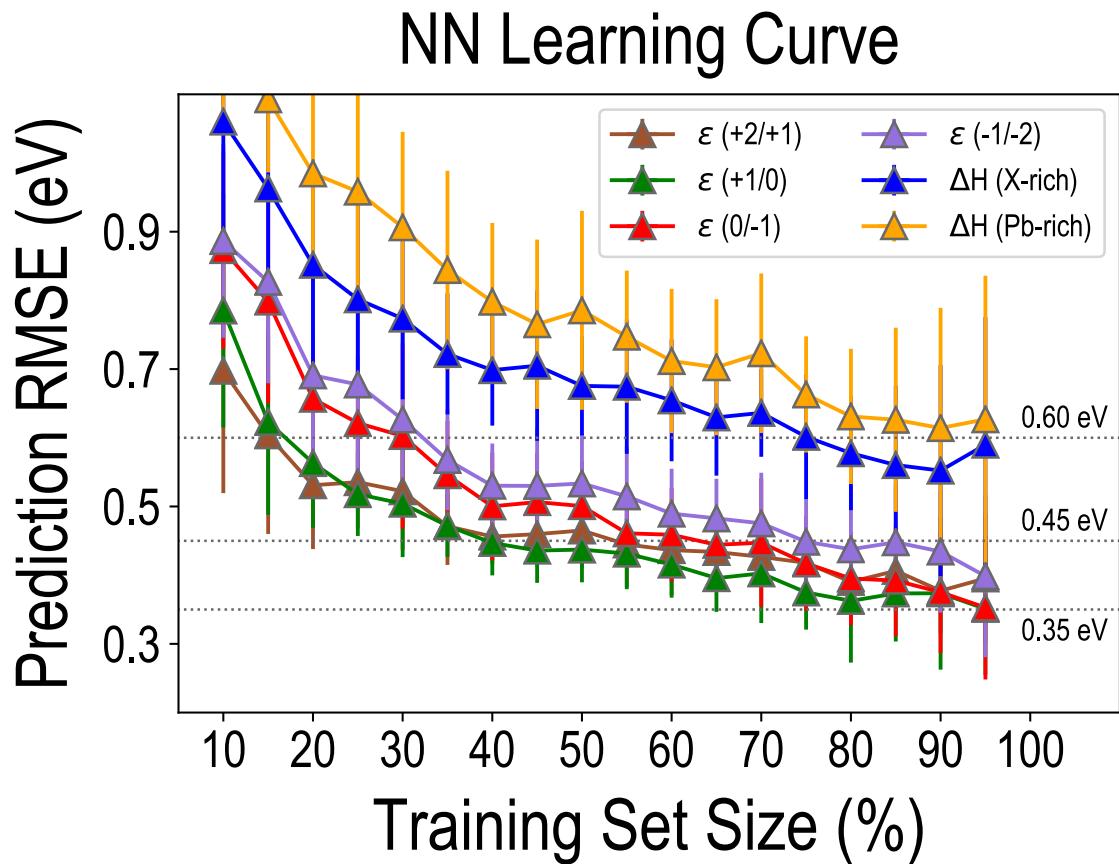
# Types of Machine Learning

- Supervised learning: From labeled training data, find the unknown function connecting known inputs to unknown outputs, based on extrapolation of patterns.
- Unsupervised learning: Find patterns in unlabeled data, leading to clustering of samples.
- Semi-supervised learning: Representations learned from a mix of unlabeled and labeled data.
- Reinforcement learning: Finding optimal or sufficiently good actions for a situation to maximize a reward.

# Nuts and Bolts of ML

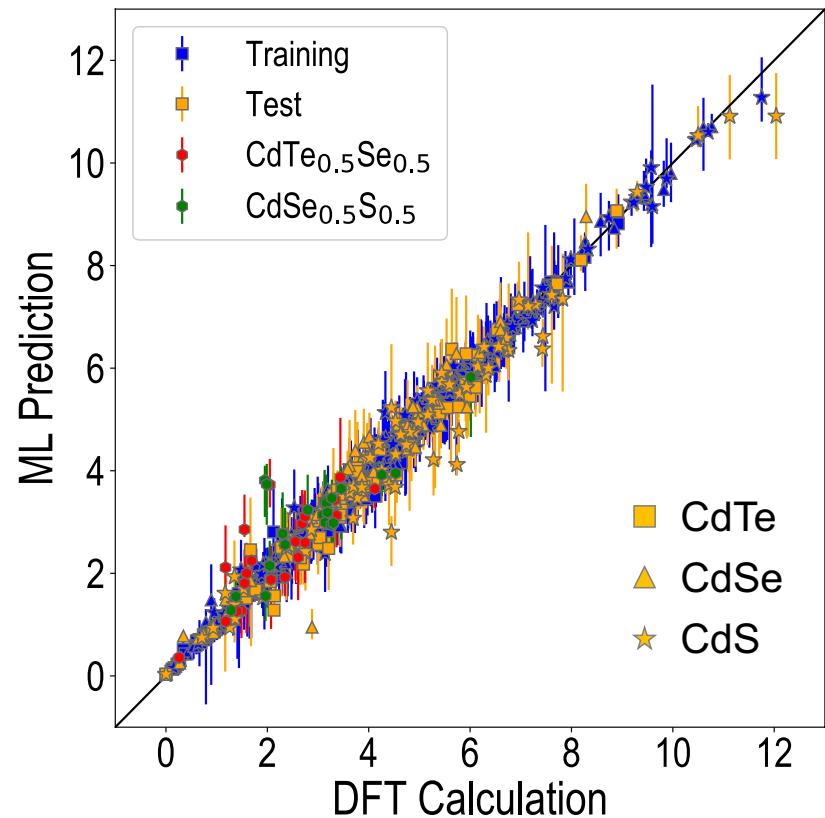
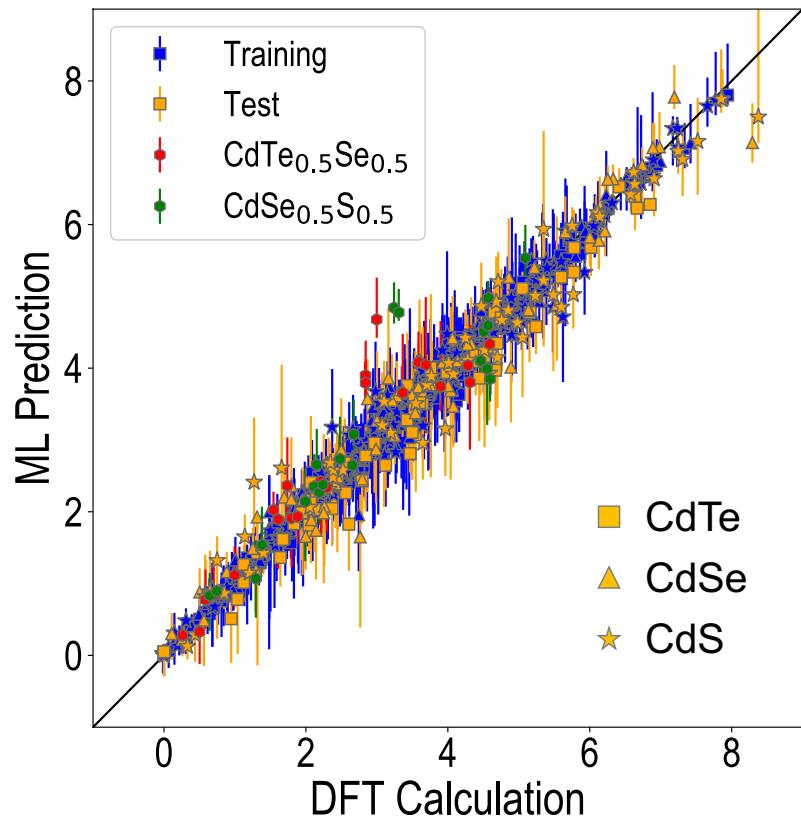
- Data: Divide into training, validation, and test sets.
- Descriptors: enumerate for all data points, perform dimensionality reduction and feature scaling.
- Cross-validation: n-fold, leave-one-out.
- Hyperparameter optimization: grid-search, Bayesian.
- Best model: optimized w.r.t. training data and descriptor size, CV, and HPO; choose error definition.
- Quality and quantity of data and descriptors: prevent underfitting. CV, HPO, regularization: prevent overfitting.

# Training Data Size: Learning Curves



Iteratively  
change training  
set size (and  
descriptor  
dimensions) until  
test error  
saturates →  
learning curve.

# Example: ML Regression



Predicting defect formation energy in semiconductors.

*npj Comput Mater.* **6**, 39 (2020)



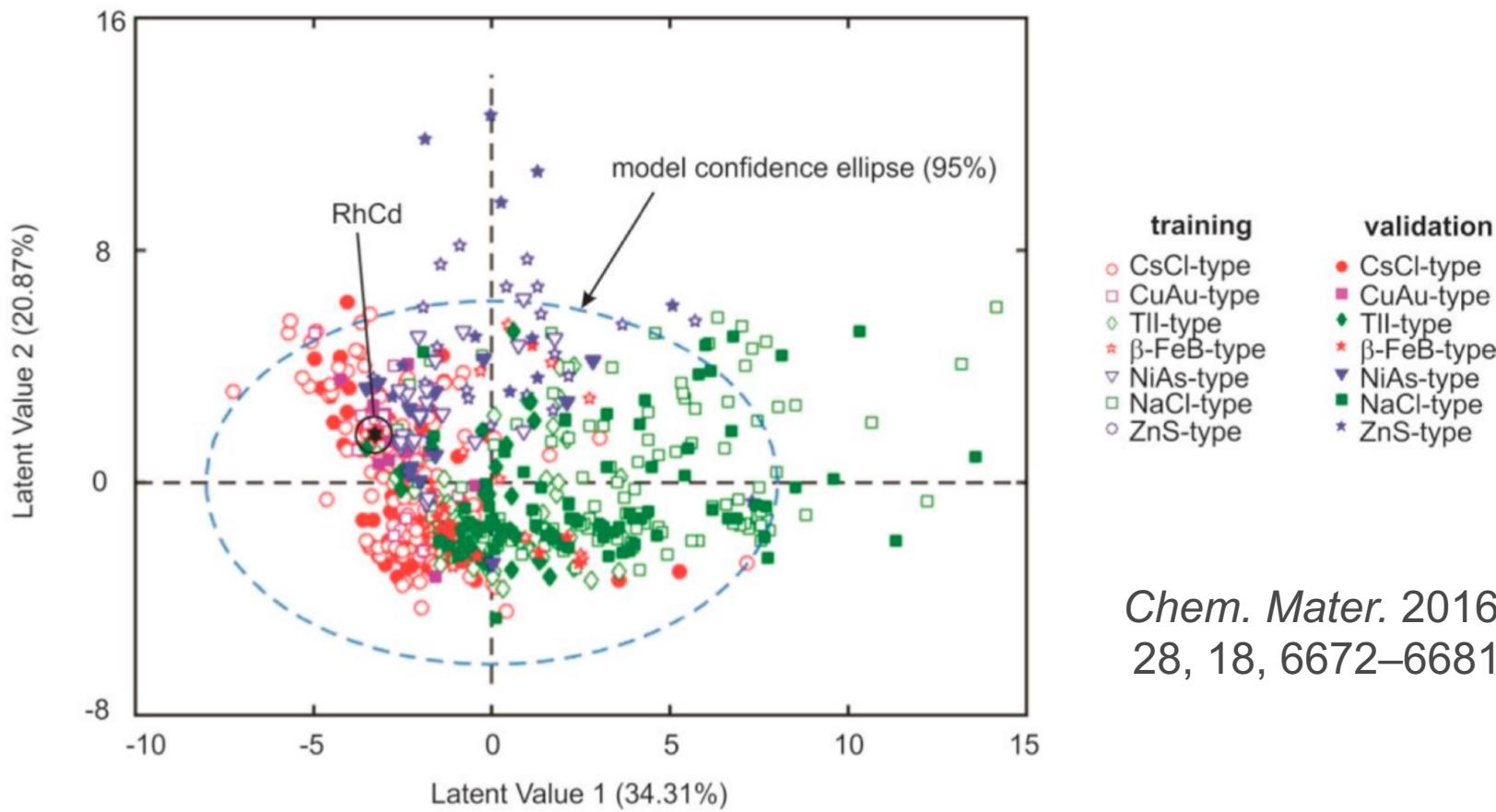
**PURDUE**  
UNIVERSITY®

School of Materials Engineering

4/10/23

12

# Example: ML Classification



Classifying AB compounds into different structure types.



PURDUE  
UNIVERSITY

School of Materials Engineering

4/10/23

13

# Linear Regression

- Linear regression model prediction form:  
 $y = \theta_0 + \theta_1x_1 + \theta_2x_2 + \dots + \theta_mx_m$ ; or,  $y = h_{\theta}(x) = \theta^T \cdot x$
- MSE cost function for a linear regression model:

$$\text{MSE}(\mathbf{X}, h_{\theta}) = \frac{1}{m} \sum_{i=1}^m (\theta^T \cdot \mathbf{x}^{(i)} - y^{(i)})^2$$

- Direct estimation of  $\theta$  using the formula below (indirect method uses GD):  
$$\hat{\theta} = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y}$$
- Easy extension to Polynomial regression.



PURDUE  
UNIVERSITY®

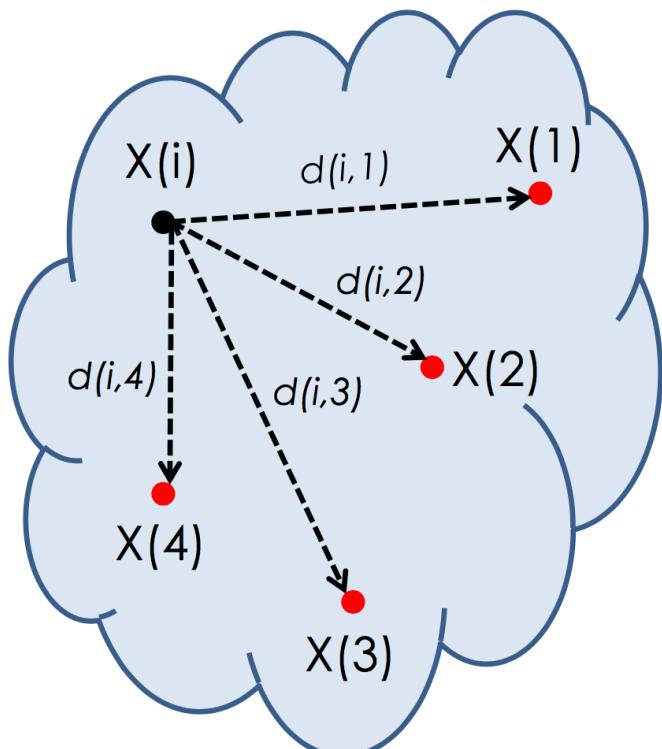
School of Materials Engineering

4/10/23

14

# Kernel Ridge Regression

*Chemical Space*



$$X(i) = \{x_1, x_2, x_3 \dots x_m\}$$

## KERNEL RIDGE REGRESSION (KRR)

Measure of Similarity: Euclidean Distance

$$d(i,j) = \sqrt{(x_{i1} - x_{j1})^2 + \dots + (x_{im} - x_{jm})^2}$$

Property = Weighted sum of Gaussians

$$f(i) = \sum_{k=1}^N a_k \cdot \exp\left(-\frac{1}{2\sigma^2} \cdot [d(i, i_k)]^2\right)$$



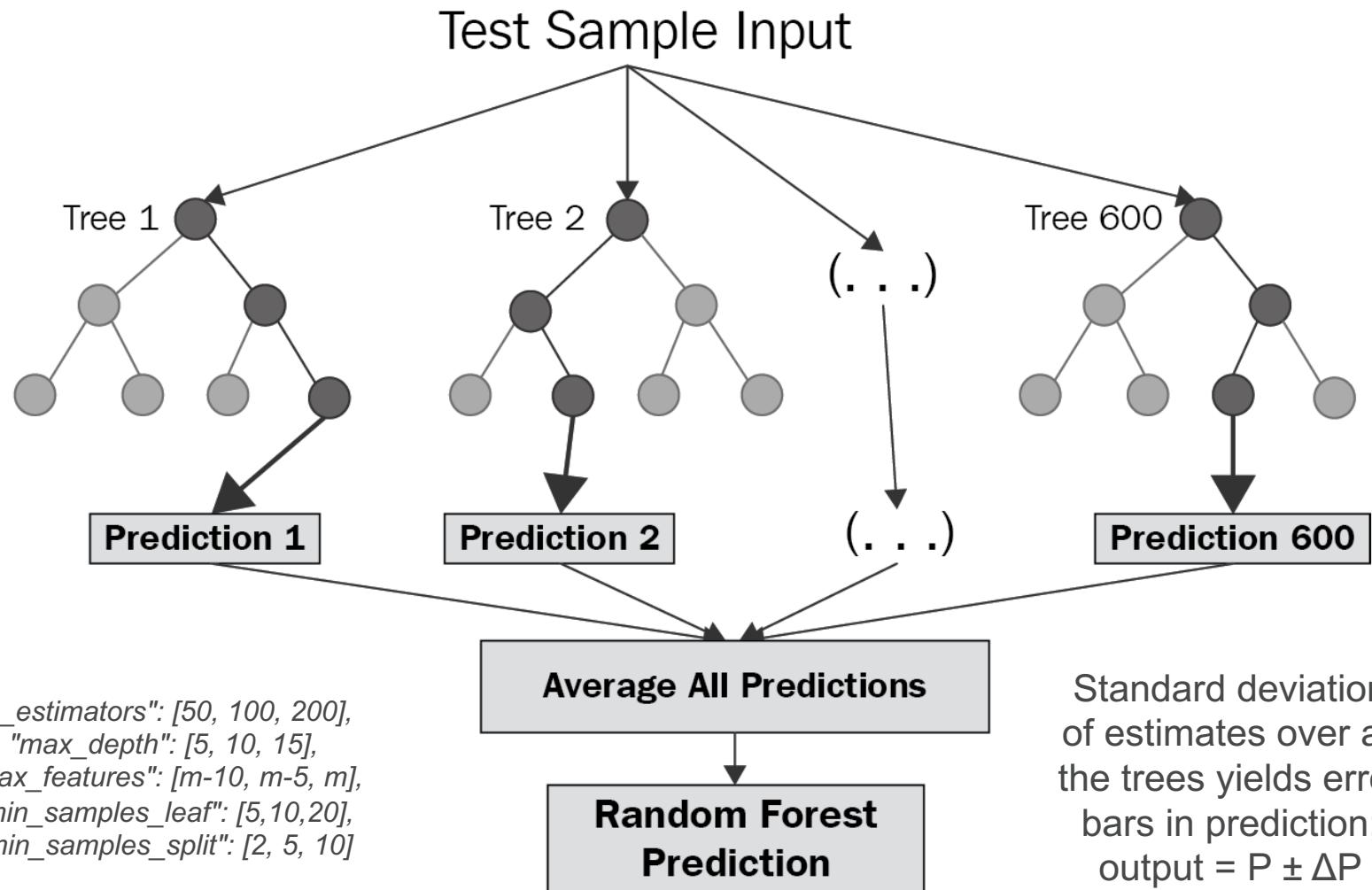
PURDUE  
UNIVERSITY®

School of Materials Engineering

4/10/23

15

# Random Forest Regression



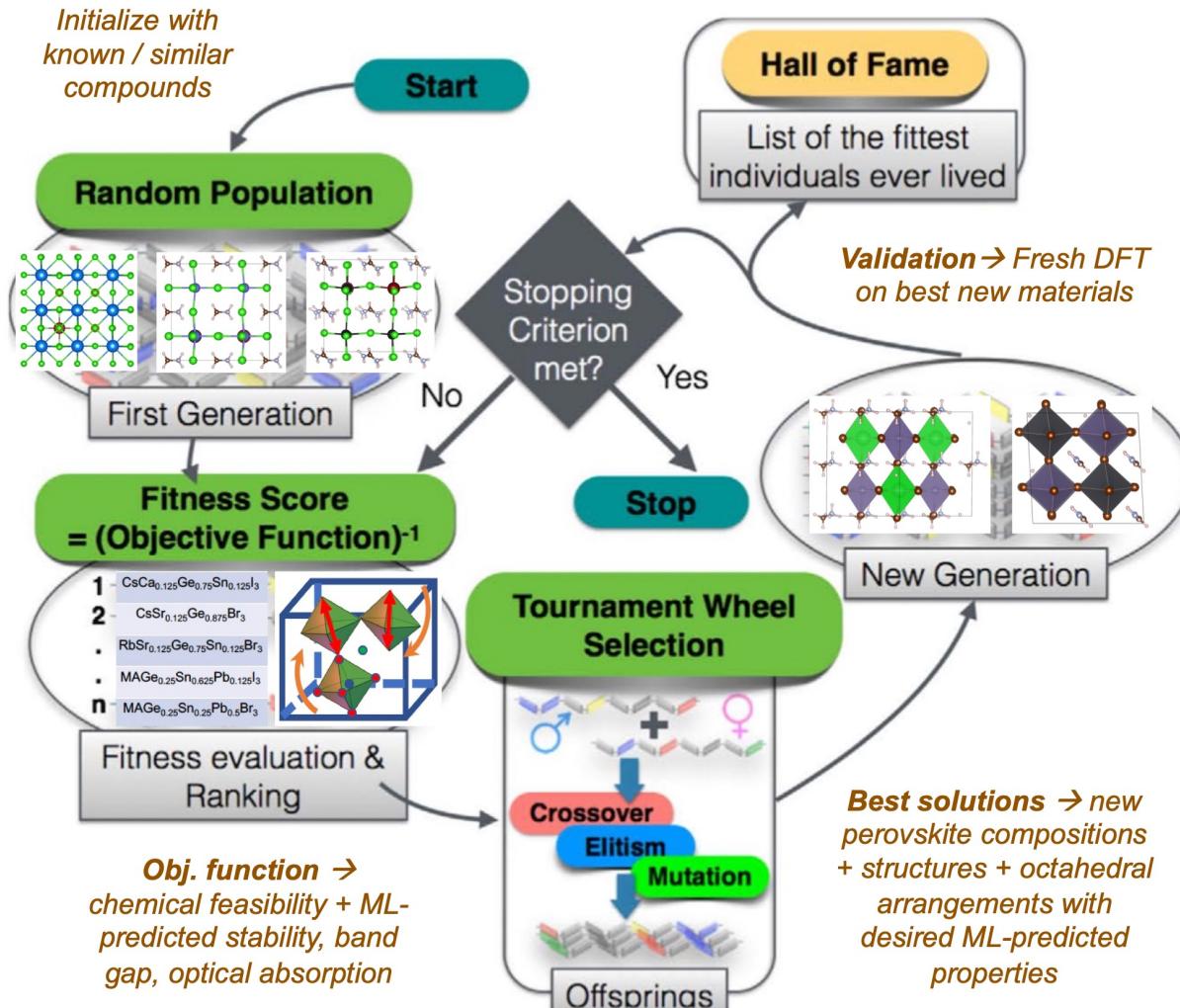
**PURDUE**  
UNIVERSITY®

School of Materials Engineering

4/10/23

16

# Genetic Algorithm



PURDUE  
UNIVERSITY®

School of Materials Engineering

4/10/23

17

# PART 1 SUMMARY

- Definition and basic principles of materials informatics.
- Key steps in ML: materials data, materials descriptors, techniques for regression and classification.
- Linear regression, RFR, KRR, GA.
- NEXT: A detailed supervised learning case study → predicting (DFT computed) properties of a chemical space of halide perovskite alloys.

# MRS Bulletin Article in Press

## A Framework for Materials Informatics Education through Workshops

Arun Mannodi-Kanakkithodi<sup>1</sup>, Austin McDannald<sup>2</sup>, Shijing Sun<sup>3</sup>, Saaketh Desai<sup>4</sup>, Keith Brown<sup>5</sup> and A. Gilad Kusne<sup>2</sup>

### Abstract

The burgeoning field of materials informatics necessitates a focus on educating the next generation of materials scientists in the concepts of data science, artificial intelligence (AI), and machine learning (ML). In addition to incorporating these topics in undergraduate and graduate curricula, regular hands-on workshops present the most effective medium to initiate researchers to informatics and have them start applying the best AI/ML tools to their own research. With the help of the Materials Research Society (MRS), members of the MRS AI staging committee, and a dedicated team of instructors, we successfully conducted workshops covering the essential concepts of AI/ML as applied to materials data, at both the spring and fall meetings in 2022, with plans to make this a regular feature in future meetings. In this article, we discuss the importance of materials informatics education via the lens of these workshops, including details such as learning and implementing specific algorithms, the crucial nuts and bolts of ML, and using competitions to increase interest and participation.

# PART 2, CASE STUDY: Data-Driven Design of Halide Perovskite Alloys

[https://github.com/mannodiarun/mrs\\_spring\\_tutorial/blob/perovs\\_dft\\_ml/DFT-ML-GA.ipynb](https://github.com/mannodiarun/mrs_spring_tutorial/blob/perovs_dft_ml/DFT-ML-GA.ipynb)

Or directly open

[https://colab.research.google.com/github/mannodiarun/mrs\\_spring\\_tutorial/blob/perovs\\_dft\\_ml/DFT-ML-GA.ipynb](https://colab.research.google.com/github/mannodiarun/mrs_spring_tutorial/blob/perovs_dft_ml/DFT-ML-GA.ipynb)

<https://tinyurl.com/257ruuu>

# PART 3: Gaussian Process Regression & Active Learning

[https://github.com/mannodiarun/mrs\\_spring\\_tutorial/tree/GP\\_and\\_AL](https://github.com/mannodiarun/mrs_spring_tutorial/tree/GP_and_AL)

GPR notebook: <https://tinyurl.com/3db88ka5>

AL notebook: <https://tinyurl.com/2zu9n7aj>

# PART 4: Neural Network Regression and Classification

[https://github.com/mannodiarun/mrs\\_spring\\_tutorial/t  
ree/nn\\_models](https://github.com/mannodiarun/mrs_spring_tutorial/tree/nn_models)

Regression notebook: <https://tinyurl.com/yf3sx74c>

Classification notebook: <https://tinyurl.com/5n7kecyr>