

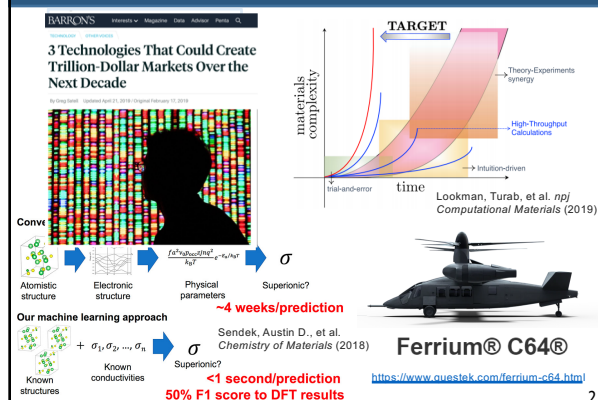
Data science and machine learning for material science

Saaketh Desai, Juan Carlos Verduzco, Ale Strachan
desai61@purdue.edu
jverduzc@purdue.edu
strachan@purdue.edu

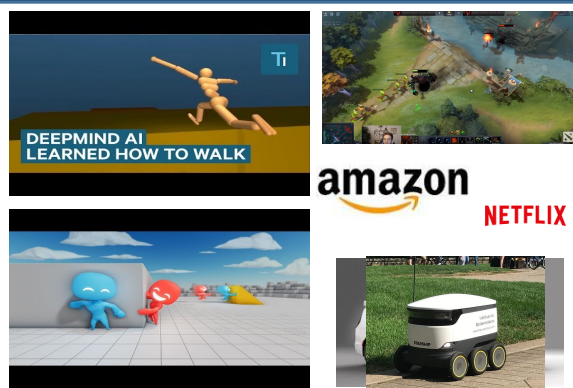
School of Materials Engineering and Birk nanotechnology center
Purdue University



Materials by design

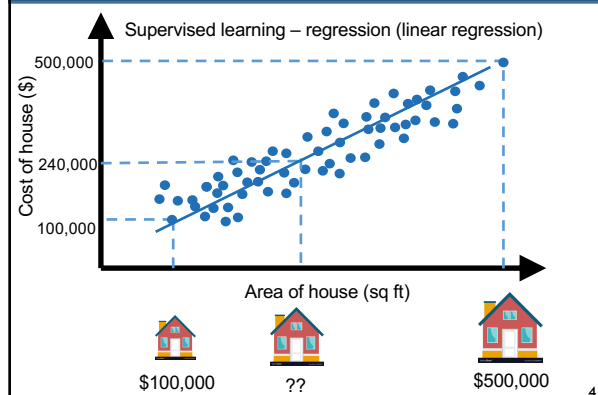


Data science, machine learning and AI



3

Machine learning and data science

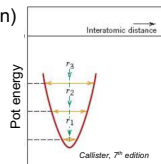


4

Linear regression to predict Young's modulus

Supervised learning – regression (linear regression)

- **Supervised learning:** Learning input-output relations based on a set of labeled 'ground truth' data
- **Regression:** Predicting a continuous relationship between inputs and outputs
- **Linear regression:** Assumed relationship to be linear



Can we predict Young's modulus based on the melting temperature?

Introduction to Machine Learning for Materials Science

The tutorials here will give you an insight into the usage of Machine Learning to approach problems related to materials science.

- **Get started:** Click on the links below to begin each tutorial.
- **Important:** To exit individual tutorials and return to this page, use File -> Close and Halt. "Terminate Session" (top right) will kill your entire Jupyter session.

Querying databases, Organizing and Plotting Data:

- Query Pymatgen and Mendeleeev for properties like Young's modulus and melting temperature
- Organize data into Pandas dataframes and python dictionaries and plot using Plotly

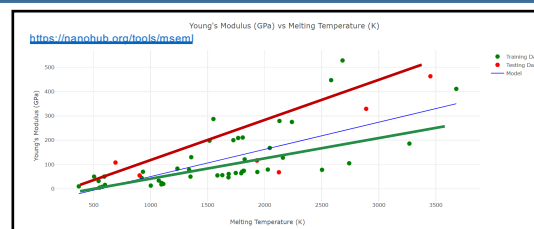
Linear Regression to predict material properties:

- Perform linear regression using the scikit learn package and predict Young's modulus
- Visualize trends in data and 'goodness of fit' of linear model

<https://nanohub.org/tools/mseml>

5

Obtaining good models – objective functions & gradient descent



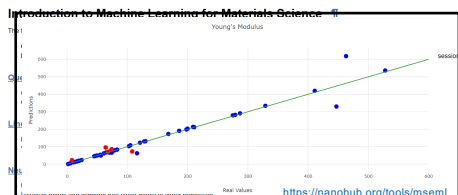
- Objective function: metric to assess model quality
- Algorithm to direct model towards objective: Gradient descent
- Objective function: minimize sum of squared distances (points should lie as close to line as possible)
- Algorithm: move in the direction of largest decrease in error

6

Non-linear models: Neural networks

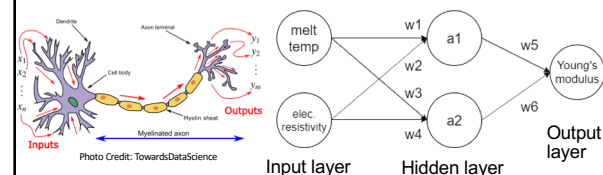
We could predict the Young's modulus given the melting temperature
Can we do better?

heat_of_formation	lattice_constant	melting_point	specific_heat	atomic_mass	atomic_radius	electrical_resistivity
284.9	4.09	1235.10	0.237	107.868200	1.60	1.630000e-08
330.9	4.05	933.50	0.900	26.981539	1.25	2.700000e-08
368.2	4.08	1337.58	0.129	196.966569	1.35	2.200000e-08
337.4	3.61	1356.60	0.385	63.546000	1.35	1.720000e-08
669.0	3.84	2683.00	0.133	192.217000	1.35	4.700000e-08



7

Neural networks – Forward propagation



$$a_1 = f_1(w_1 \text{ melt. temp} + w_2 \text{ elec. resistivity} + b_1)$$

activation weight bias

$$a_2 = f_2(w_3 \text{ melt. temp} + w_4 \text{ elec. resistivity} + b_2)$$

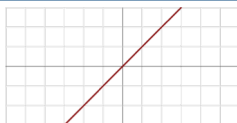
$$\text{Young's modulus} = f_3(w_5 a_1 + w_6 a_2 + b_3)$$

Training a neural network is an optimization problem where we solve for the weights and biases that result in accurate predictions

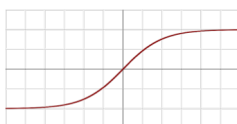
8

Neural networks – Activation functions

Linear
 $f(x) = x$



Tanh
 $f(x) = \tanh(x)$



Relu
 $f(x) = x$ if $x > 0$
 $f(x) = 0$ if $x < 0$



9

Neural networks – back propagation

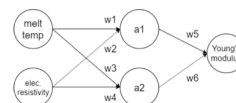
$$a_1 = f_1(w_1 \text{ melt. temp} + w_2 \text{ elec. resistivity} + b_1)$$

$$a_2 = f_2(w_3 \text{ melt. temp} + w_4 \text{ elec. resistivity} + b_2)$$

$$\text{model prediction } \hat{y} = \text{Young's modulus} = f_3(w_5 a_1 + w_6 a_2 + b_3)$$

$$\text{Objective} = \frac{1}{N_{\text{samples}}} \sum_{i=1}^{N_{\text{samples}}} (\hat{y} - y)^2$$

ground truth



We need to update weights and biases such that objective function is minimized

$$w_1 = w_1 - \alpha \frac{\partial(\text{Objective})}{\partial w_1}$$

Weight update rule: gradient descent, Adam

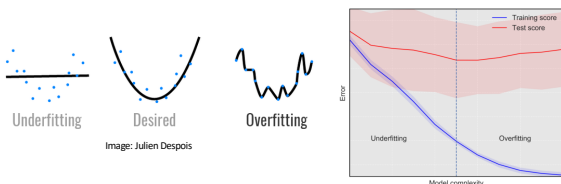
<https://neuralnetworksanddeeplearning.com/chap2.html>

10

Underfitting and overfitting

How do we judge if the model has learnt all that it could?

- Underfitting – model hasn't learnt all the trends in the training data
- Overfitting – model has "memorized" data, ignoring the underlying trend



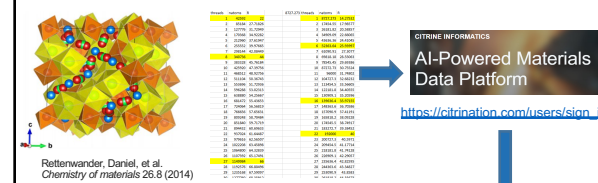
How do we train models that generalize well?

- Use a low enough learning rate
- Monitor error on validation set as a measure of model's ability to generalize

11

Sequential learning with the Citrination API

<https://nanohub.org/tools/mlday> Juan Carlos Verduzco



1. Querying a Database from Citrination

Matminer offers API tools to facilitate querying of databases like the Materials Project and Citrination. An individual Citrine Key is required for the query command CitrineDataRetrieval.

Data is stored in a Pandas DataFrame and the list of possible properties to be queried can be consulted by setting the print_properties_options parameter to True

```
In [2]: cdr = CitrineDataRetrieval('83335PqWmJnQ238QdUtt') # Citrine Key
data = cdr_get_dataframe(criteria={'data_set_id': 184822}, print_properties_options=False) # L220 Database
property_interest = 'Ionic Conductivity' # Property to be queried
display(data.head(n=10))
```

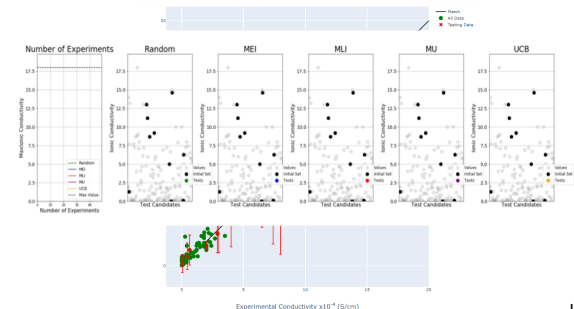
Consider making data publicly available to data platforms

12

Sequential learning with the Citrination API

<https://nanohub.org/tools/mldey> Juan Carlos Verduzco

Sequential learning accelerates identification of high ionic conductivity garnets



13

Working with small datasets – other ML techniques

Synthesis of TiO_2 nanoparticles by hydrolysis and peptization of titanium isopropoxide solution

S. Mahabadi¹, M. Aslari¹, M. Saeed Ghannam^{1,2,*}

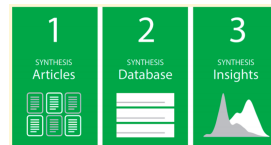
¹Department of Material Science & Eng., Sharif University of Technology, 13585-8636 Tehran, Iran

²State Key Lab of Silicon Materials, Beijing 100084, China

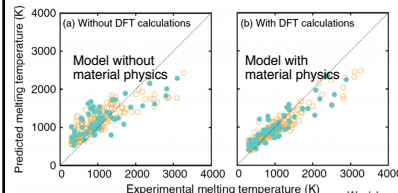
Received 12 July 2016; revised 10 December 2016; accepted 20 January 2017

The prepared precipitates were washed with ethanol and dried for several hours at 100 °C. After being washed with ethanol and dried at 100 °C in a vacuum system for 3 h, a yellow-white powder is obtained. Finally, the prepared powder was annealed at temperature ranging from 200 to 800 °C for 2 h.

ML can automate text mining and information acquisition



Kim, Edward, et al. *Chemistry of Materials* (2017)

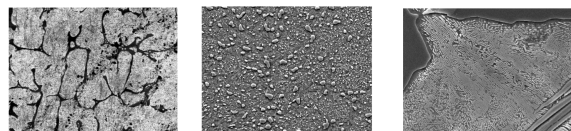


Incorporating physics into the model can significantly improve training even more smaller datasets!

Ward, Logan, and Chris Wolverton, *Current Opinion in Solid State and Materials Science* (2017)

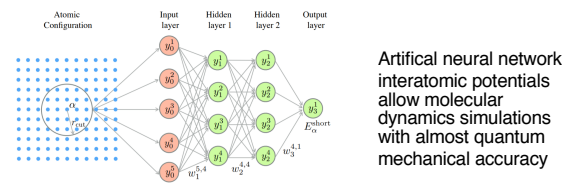
14

Other applications of neural networks



Ling, Julia, et al., *Materials Discovery* 10 (2017): 19-28

Convolutional neural networks can detect microstructural features



Artificial neural network interatomic potentials allow molecular dynamics simulations with almost quantum mechanical accuracy

Wen, Mingqian, and Elad B. Tadmor, *arXiv preprint arXiv:1909.10134* (2019)

15

Obtaining and querying data

- Materials project: <https://materialsproject.org/>
- AFLOW: <http://aflowlib.org/>
- NIST Materials Data Repository: <https://materialsdata.nist.gov/>
- Database for renewable energy materials: <https://materials.nrel.gov/>
- OQMD: <http://oqmd.org/>
- Citrination: <https://citrination.com/datasets>
- ICSD: <https://icsd.fiz-karlsruhe.de/>

Databases with Python APIs:

- Pymatgen: <https://pymatgen.org/>
- Mendelev: <https://mendelev.readthedocs.io/en/stable/data.html>
- NanoHUB: <https://nanohub.org/tools/msem/>

16

Training neural networks

- Keras: <https://materialsproject.org/>
- Scikit-learn: <https://scikit-learn.org/stable/>
- PyTorch: <https://pytorch.org/>
- Matminer: <https://hackingmaterials.lbl.gov/matminer/>
- NanoHUB: <https://nanohub.org/tools/msem/>

17