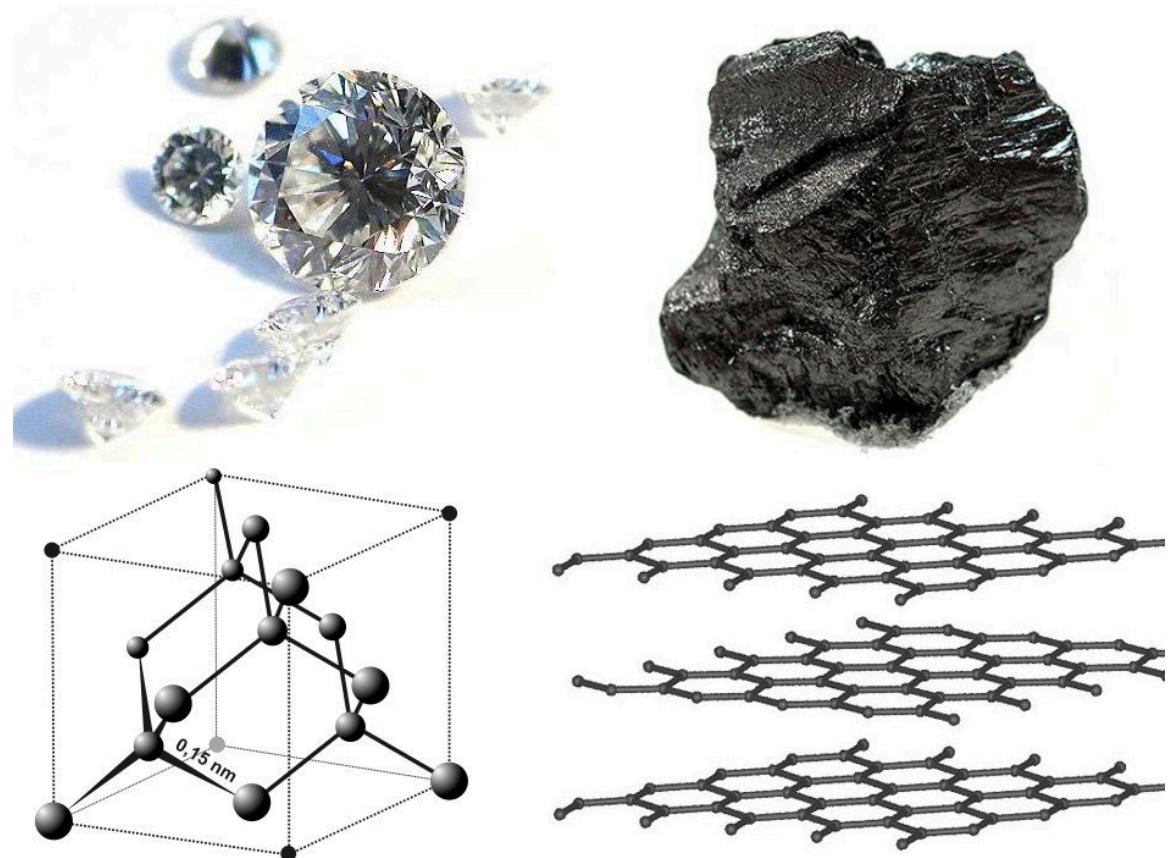


Lecture #9: Course wrap up

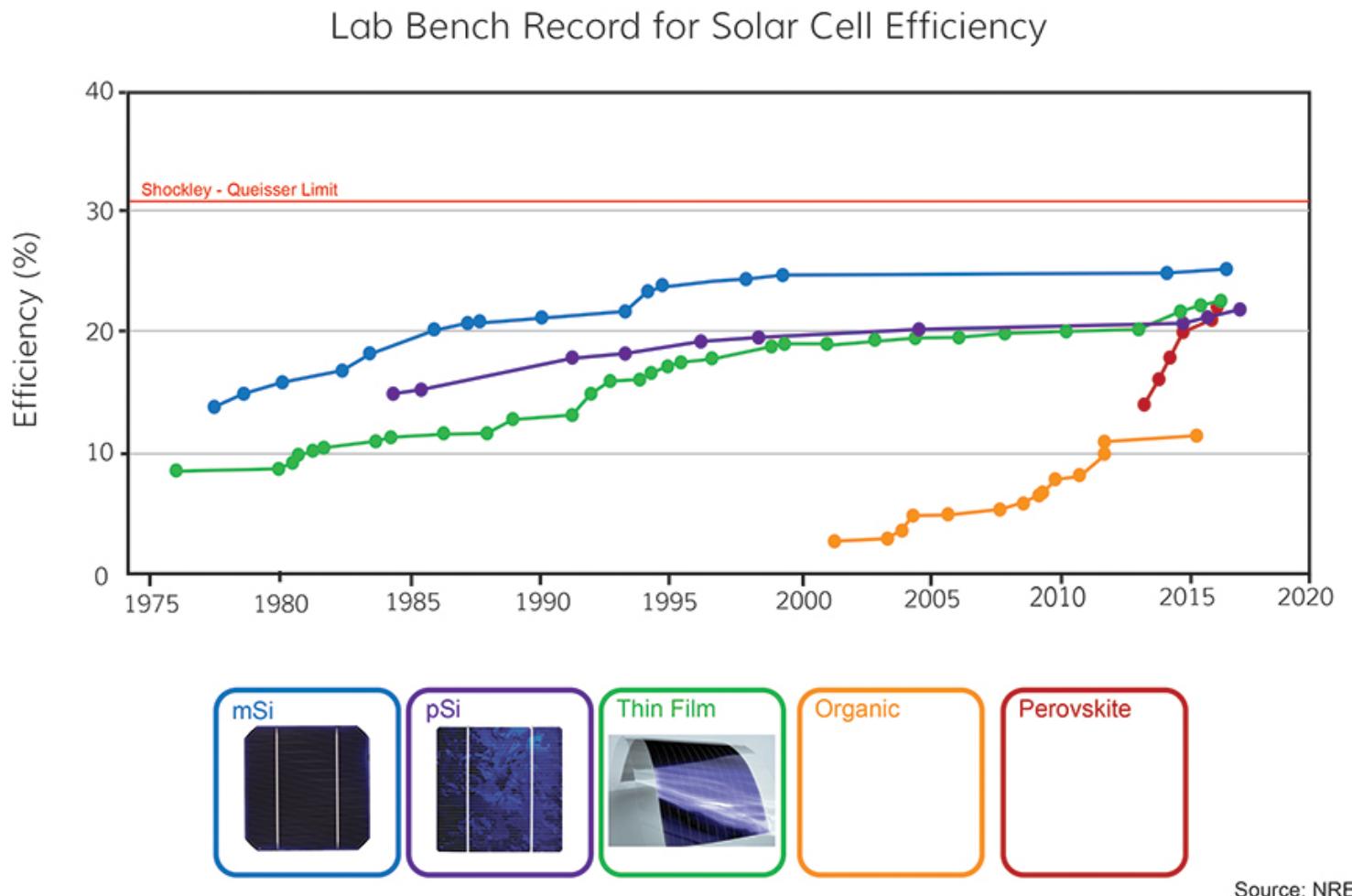
Materials science

- studies the relationship between the structure and properties of materials
- develops new materials to enable new technologies, which improve our lives (hopefully)

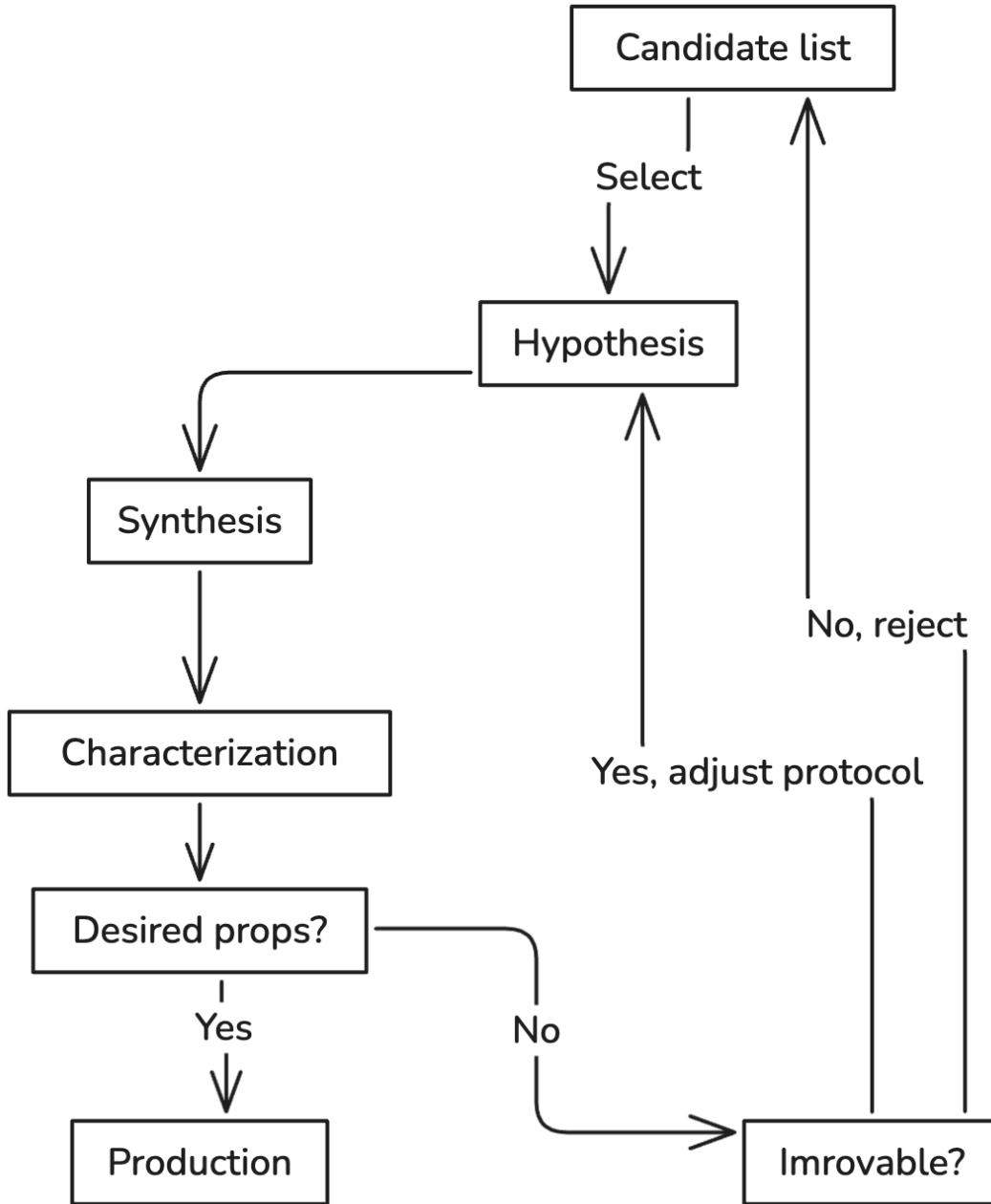


Tasks that materials science solves:

- Design
- Improve
- Understand reasons of a failure
- ...

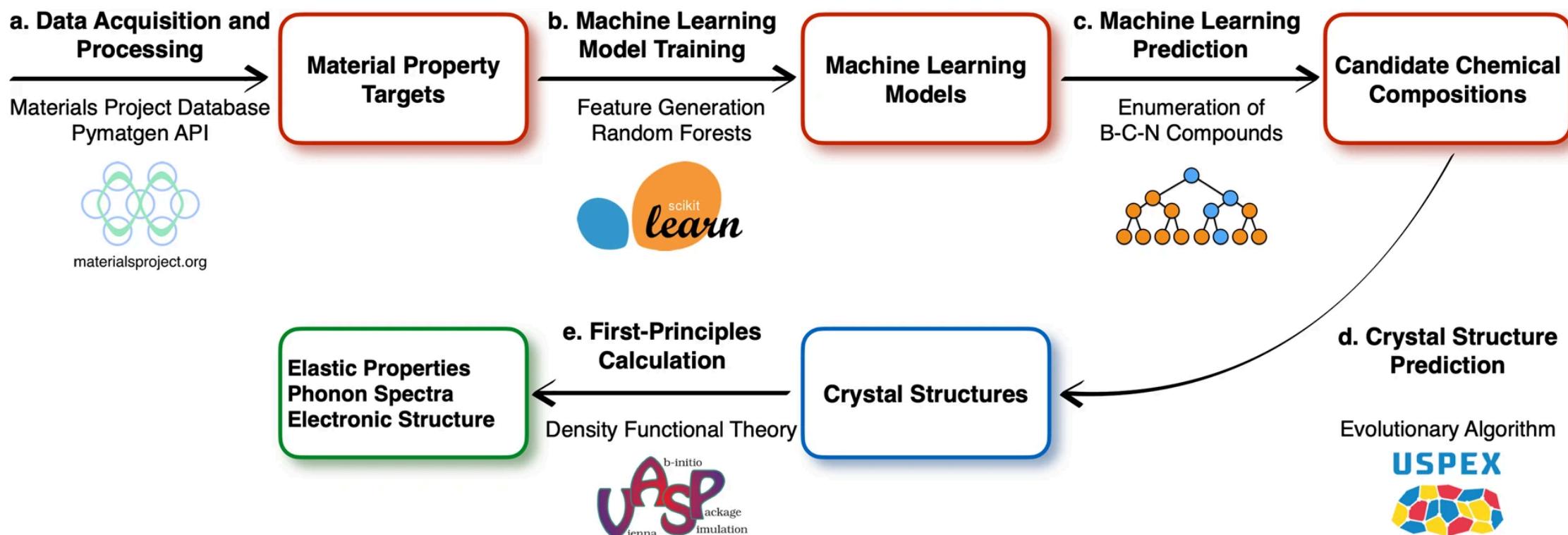


Trial-and-error



Materirals informatics

"... is a field of study that applies the principles of informatics and data science to materials science and engineering to improve the understanding, use, selection, development, and discovery of materials." ([wiki](#))



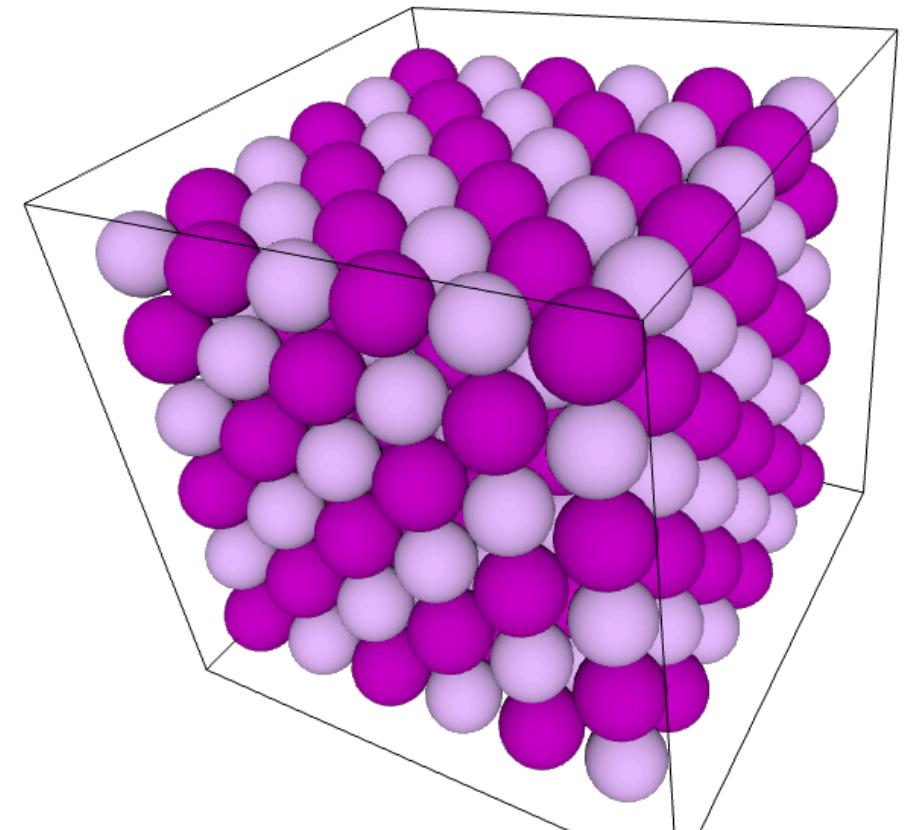
Python tools for atomistic modeling

```
from ase.io import read, write
from ase.visualize import view
from ase.build import make_supercell

atoms = read('data/LiI.cif')
sc = make_supercell(atoms,
    [[3, 0, 0],
     [0, 3, 0],
     [0, 0, 3]])
view(sc, viewer='x3d')
write('data/LiI_3x3x3.cif', sc)
```

Automation:

- saves your time
- gives you the opportunity (tools) to realize your ideas
- reduces number of mistakes (typos) made



Data

Why use someone else's data?

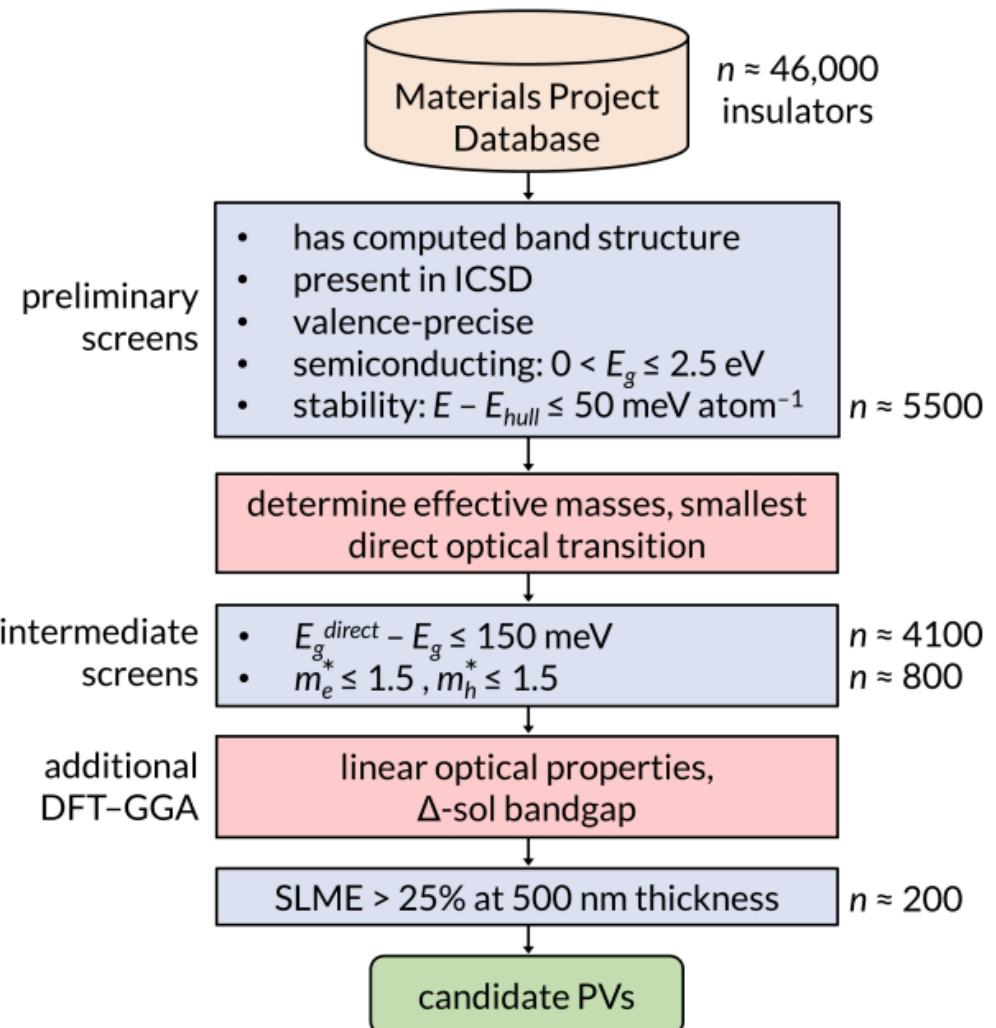
- guide to your research objective
- reference
- baseline
- insight
- explanation
- enrichment
- time



The MP's data usage example

Screening inorganic PVs

- Screened database
- Identified candidates
- Calculated properties of interest for ~800 compounds
- Shared the data



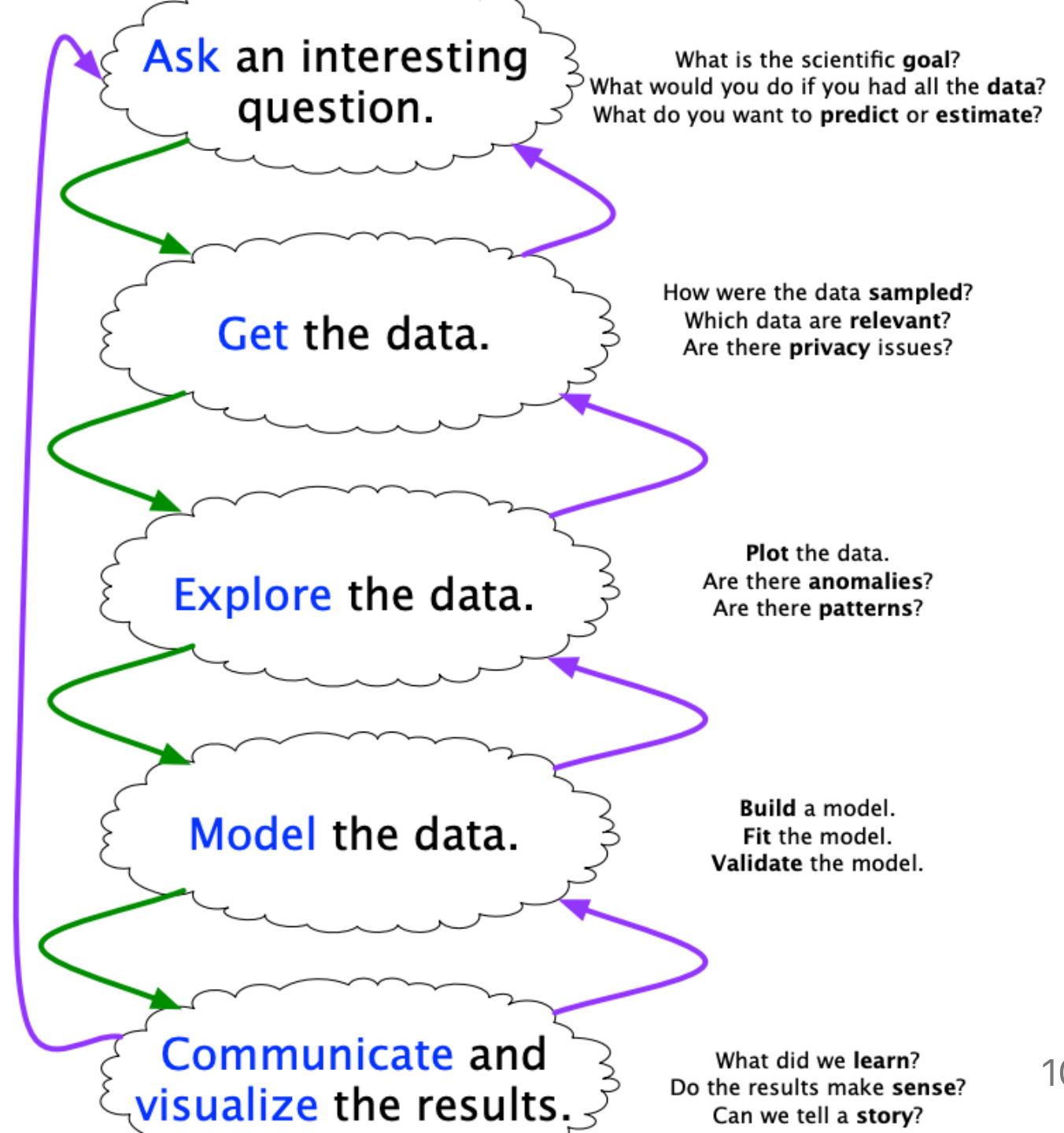
Candidate Inorganic Photovoltaic Materials from [Electronic Structure-Based Optical Absorption and Charge Transport Proxies](#)

The FAIR Guiding principles

... for scientific data management
and stewardship



The data science workflow

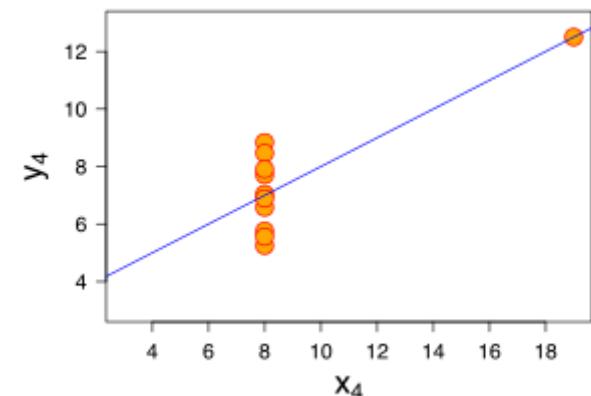
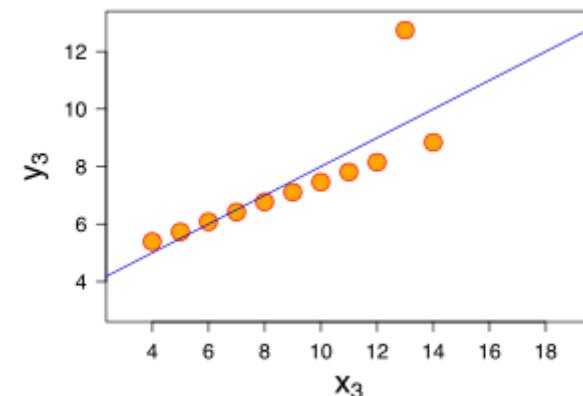
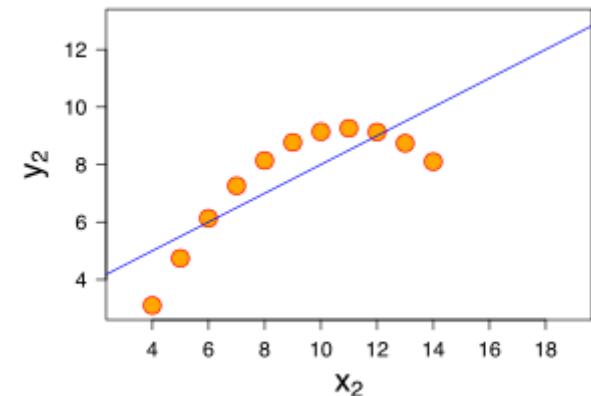
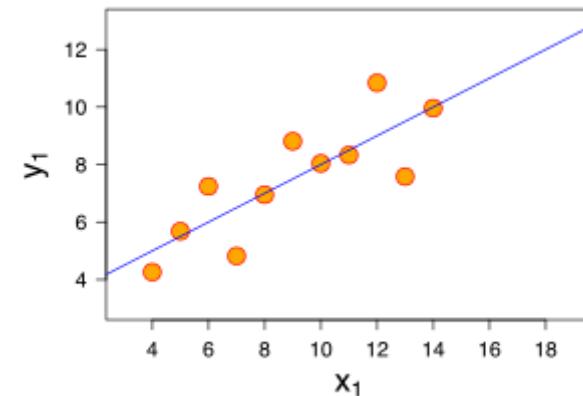


From CS 109a: Data Science, Effective Exploratory Data Analysis and Visualization by Pavlos Protopapas & Kevin Rader [slide #2](#)

Why is visual inspection of data important?

- Same descriptive statistics
- Very different distributions

https://en.wikipedia.org/wiki/Anscombe's_quartet



Exploratory data analysis pipeline

- Build data
- Clean data
- Explore global features
- Explore group features

Visulaization goals

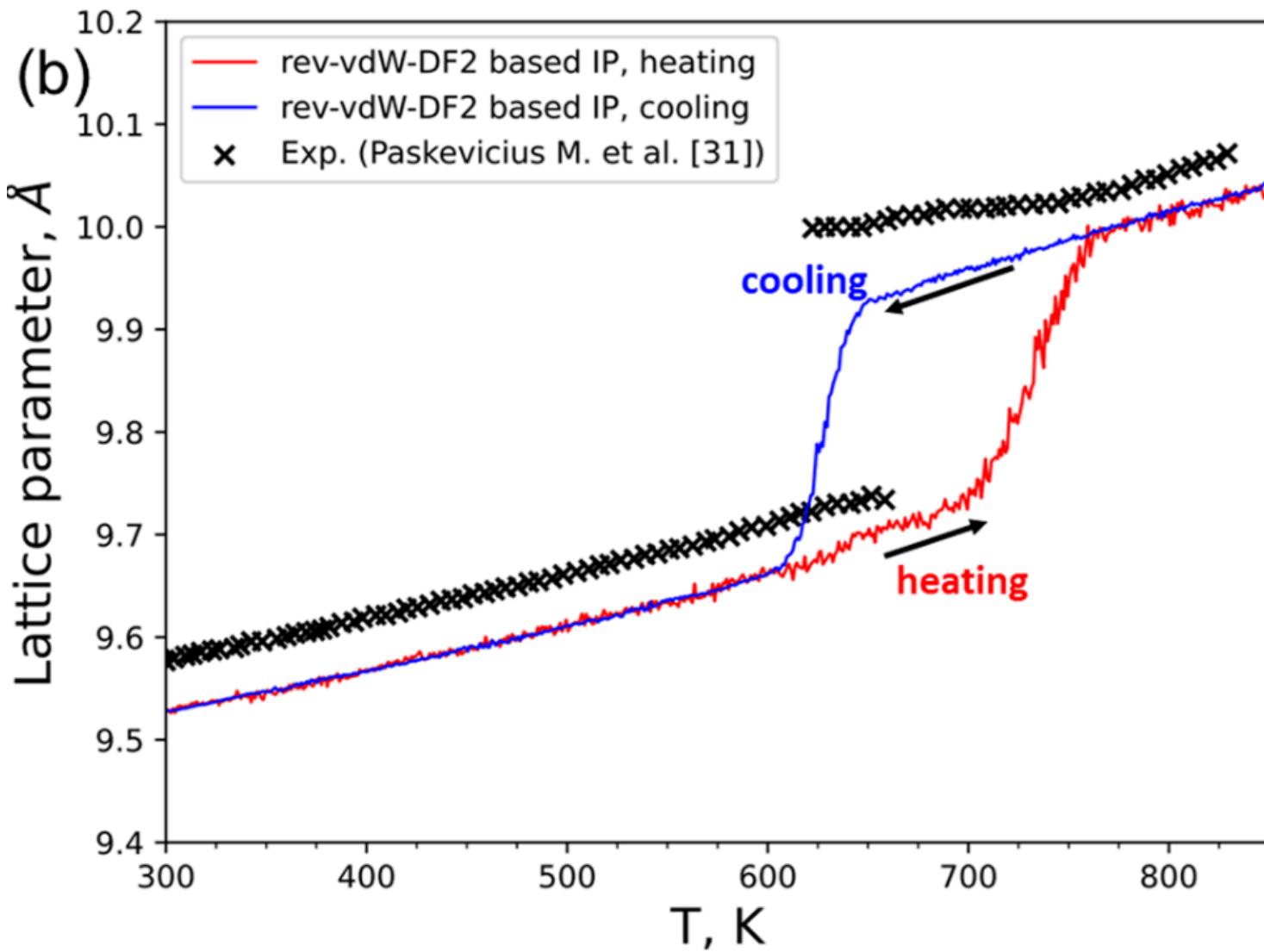
Communicate (Explanatory)

- Present data and ideas
- Explain and inform
- Provide evidence and support
- Influence and persuade

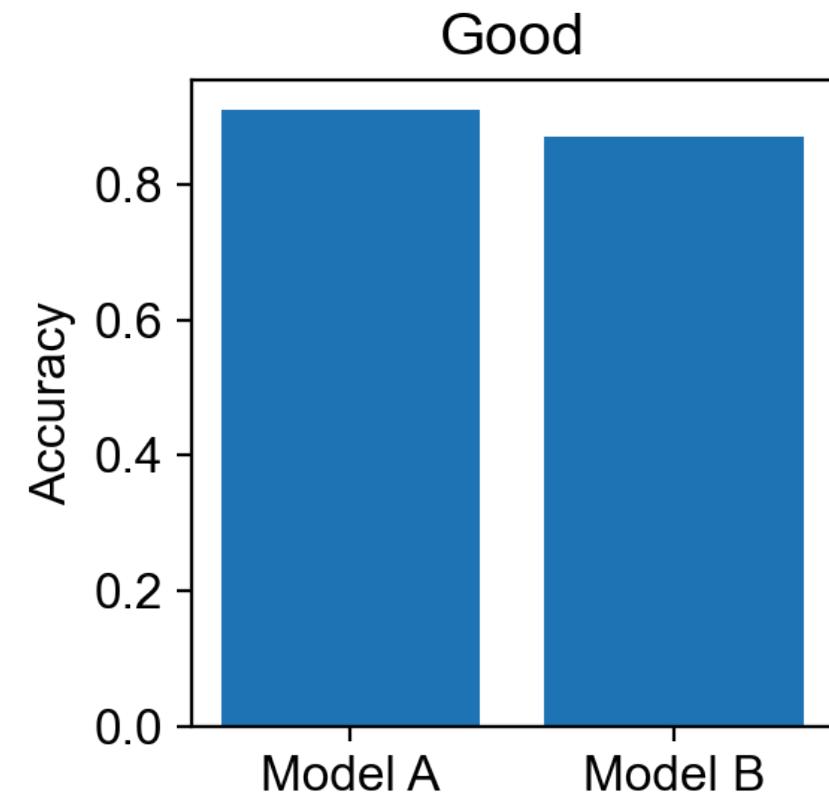
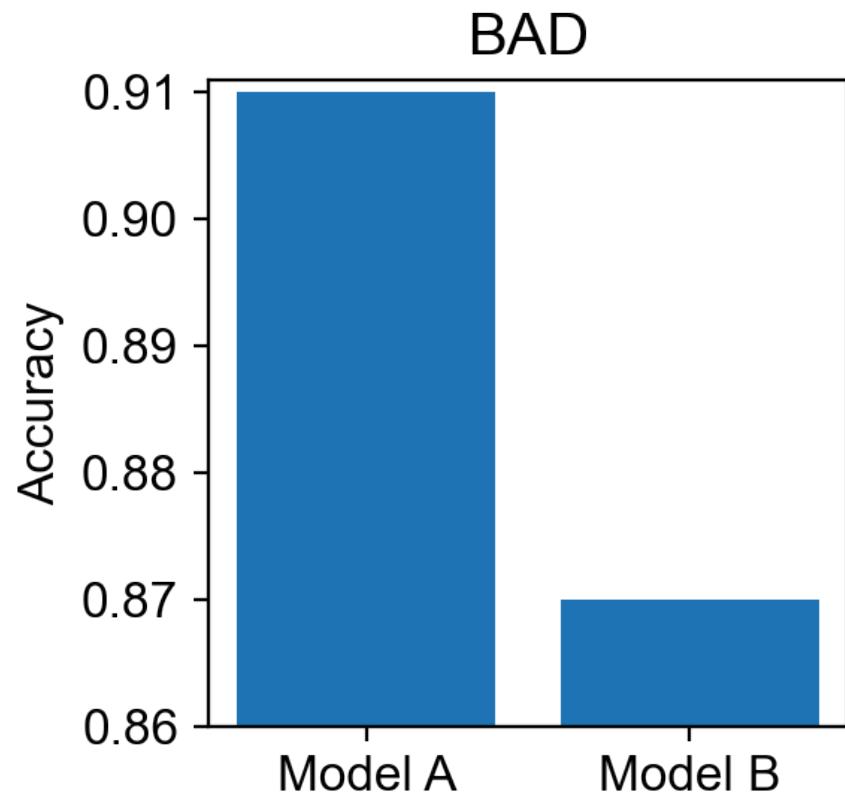
Analyze (Exploratory)

- Explore the data
- Assess a situation
- Determine how to proceed
- Decide what to do

Explore

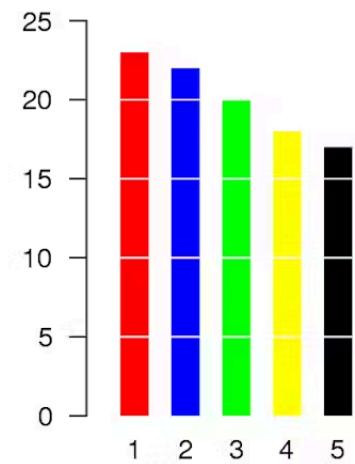
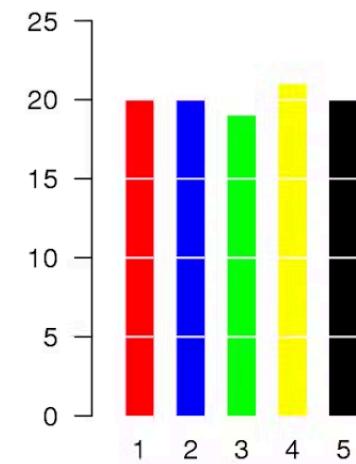
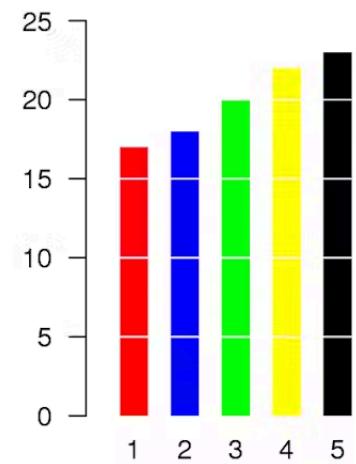
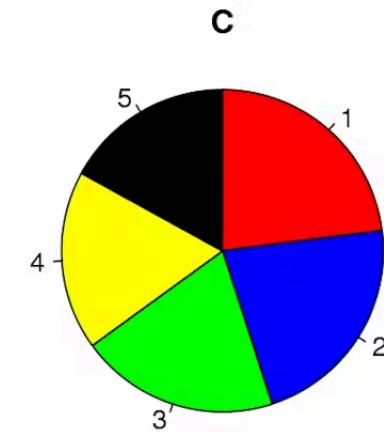
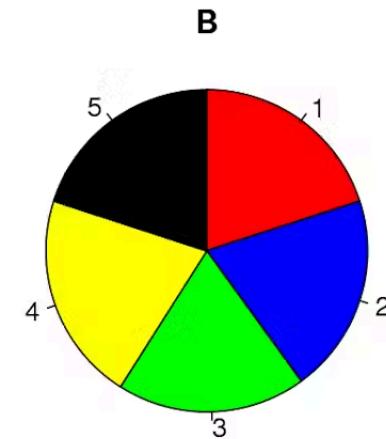
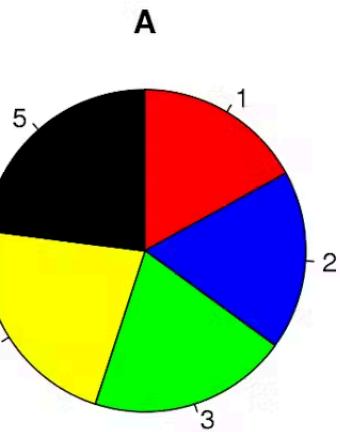


Lie factor



Don't use pie charts

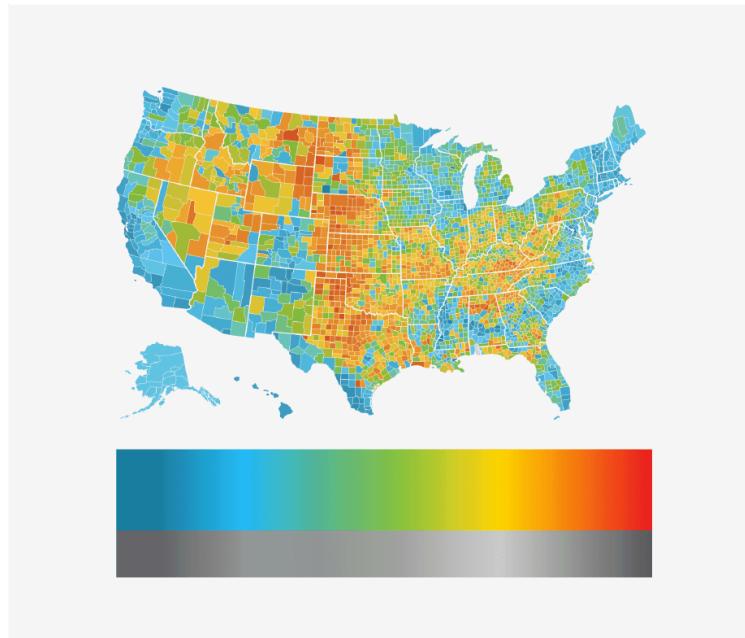
Barplots are easier to compare



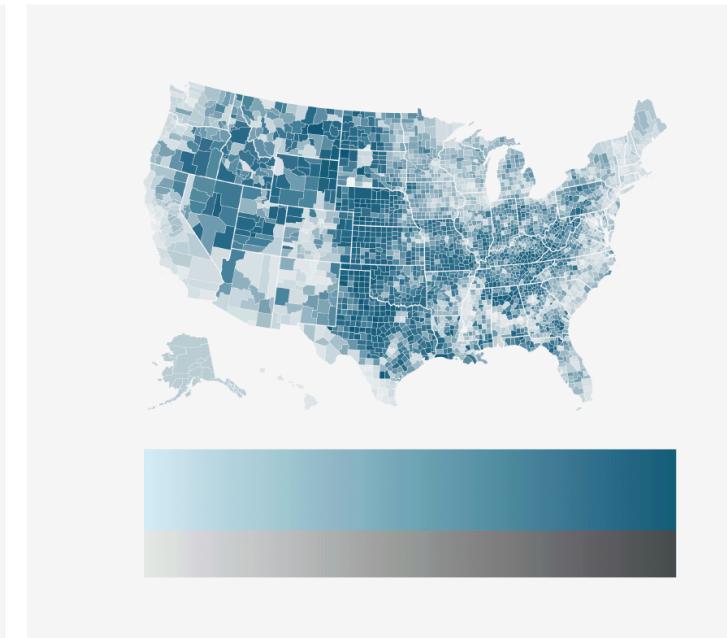
Use color

Have a look at this page:

<https://blog.datawrapper.de/colors/>



NOT IDEAL



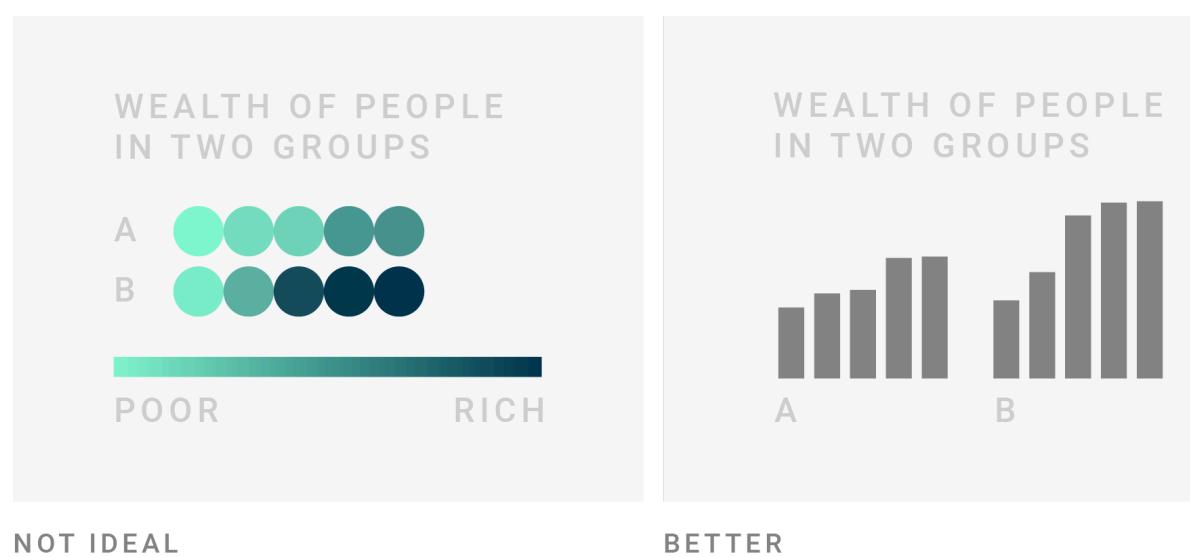
BETTER

But consider a better alternative if possible

- the simpler the better

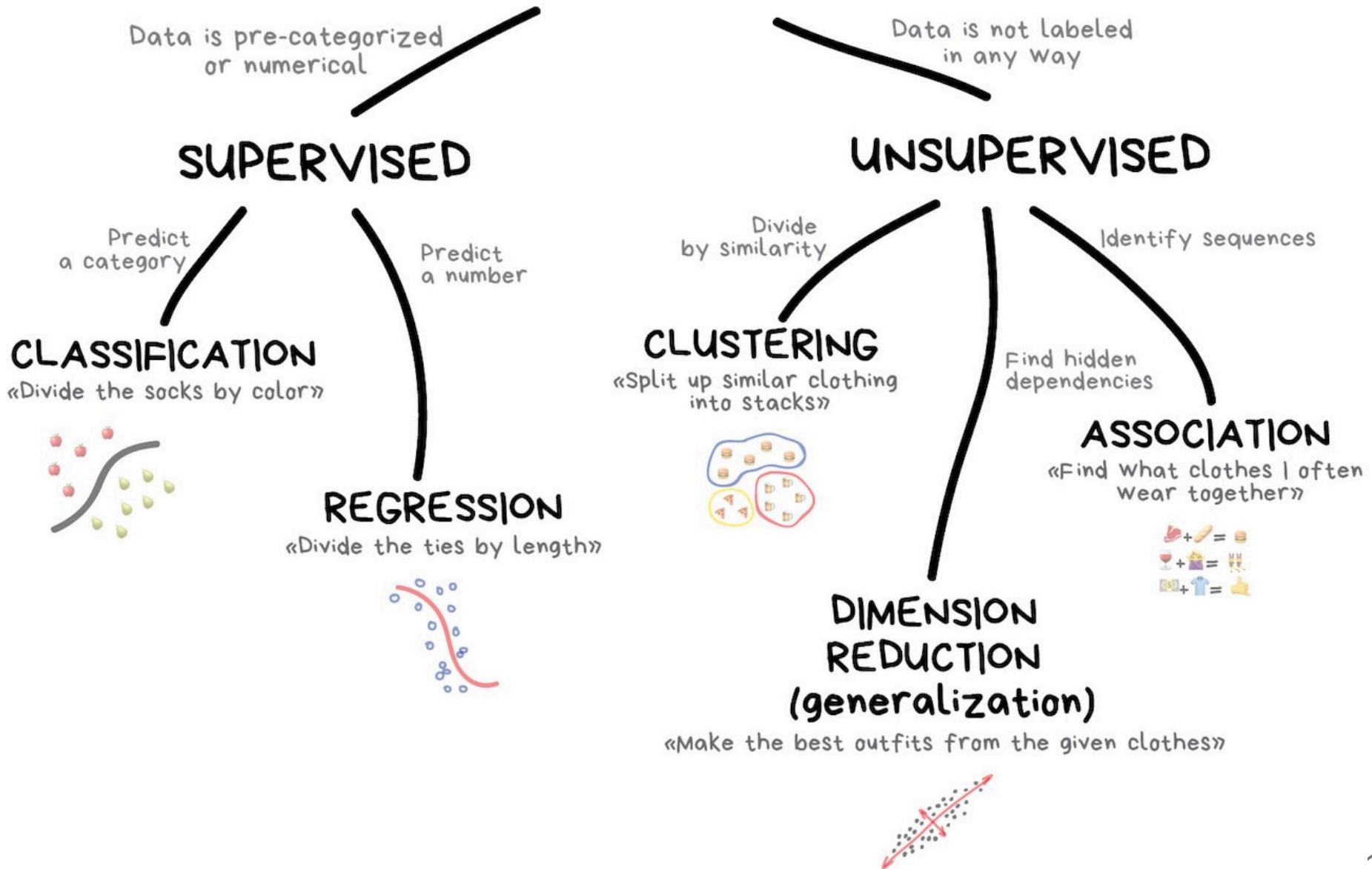
Have a look at this page:

<https://blog.datawrapper.de/colors/>



CLASSICAL MACHINE LEARNING

Classical
ML tasks



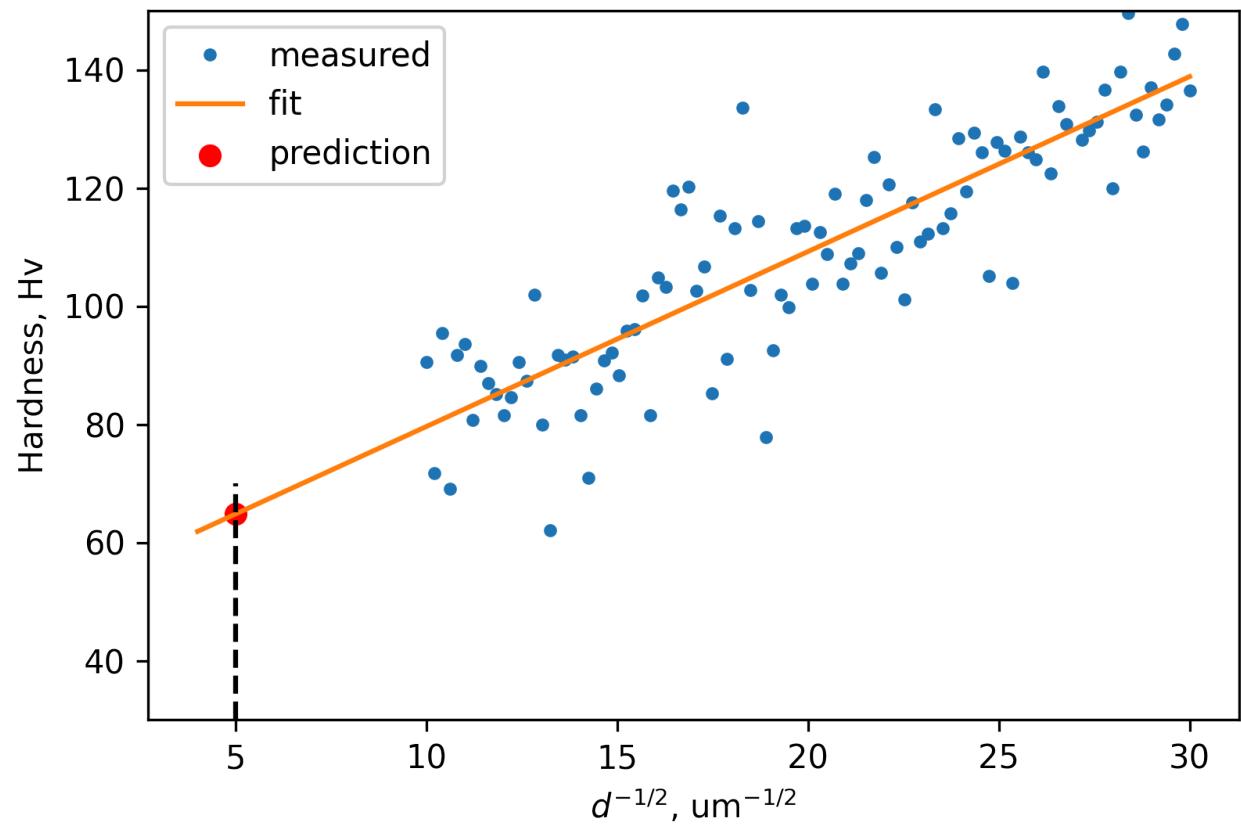
Supervised learning

"Tries to find the unknown function that connects known inputs to unknown outputs"

Model prediction:

Hardness($d^{-1/2} = 5 \text{ um}^{-1/2}$)?

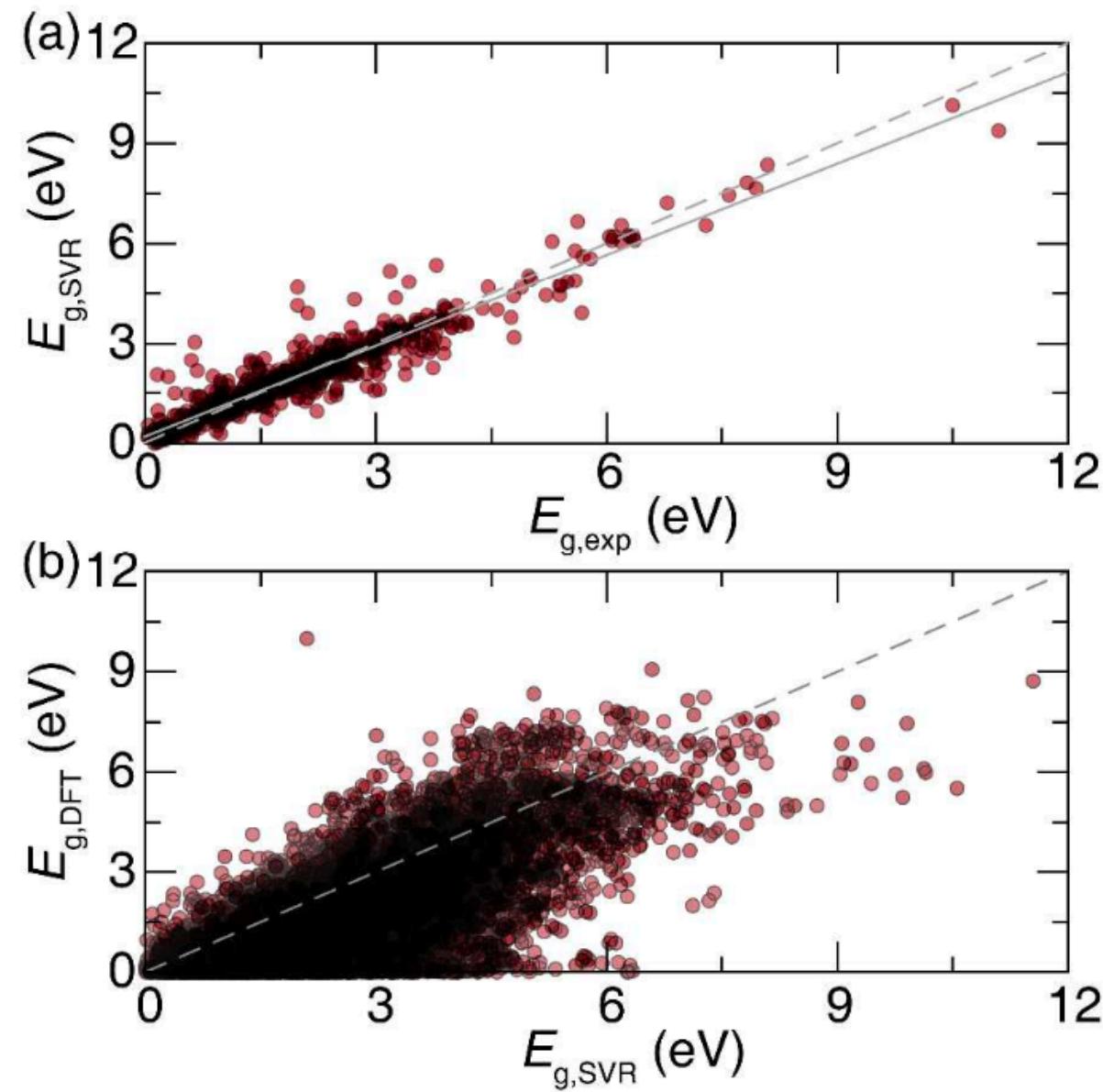
- input: grain size (known)
- output: hardness (unkown)



Surrogate model for predicting a band gap

- takes < 1 second to predict the property
- trained on existing data
- chemical composition only
- simple laptop is sufficient

SVR - Support Vector Regression

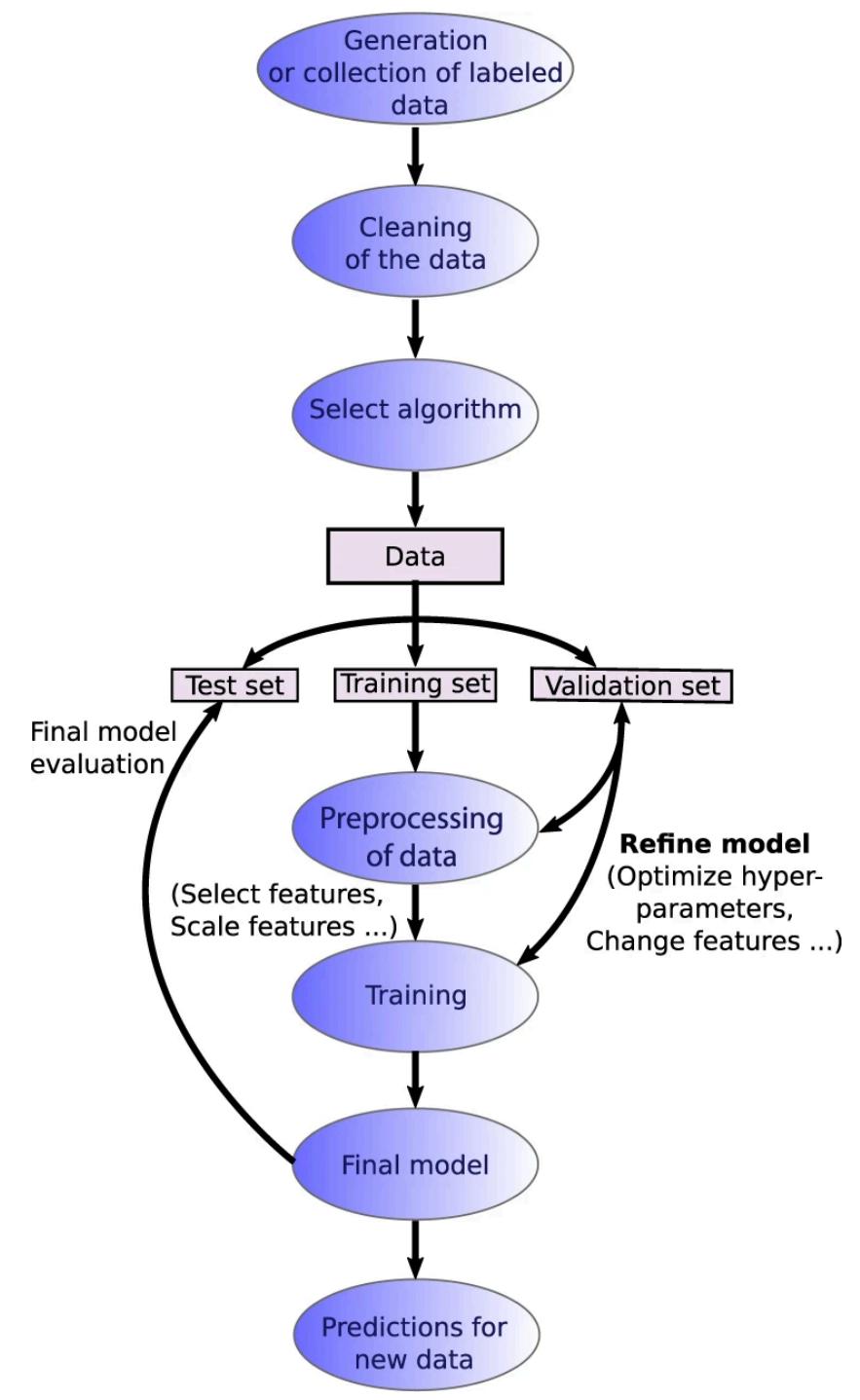


Limitations

- High quality data is required to fit the model
- Accuracy is good for screening stages (RMSE = 0.45 eV)
 - but poor for a more specific tasks

There is a trade off between speed and accuracy

Basic supervised learning workflow



Linear regression

$$J(\mathbf{w}, b) = \frac{1}{n} \sum_i^n (\mathbf{w}\mathbf{x}_i^T + b - y_i)^2$$

Optimization problem

$$J(\mathbf{w}, b) \rightarrow \min_{w,b}$$

Hyperparameters?

- During the training (learning) process we find the best **parameters** (weights) of our model
- In addition we have **hyperparameters** which are usually fixed before the actual training process begins

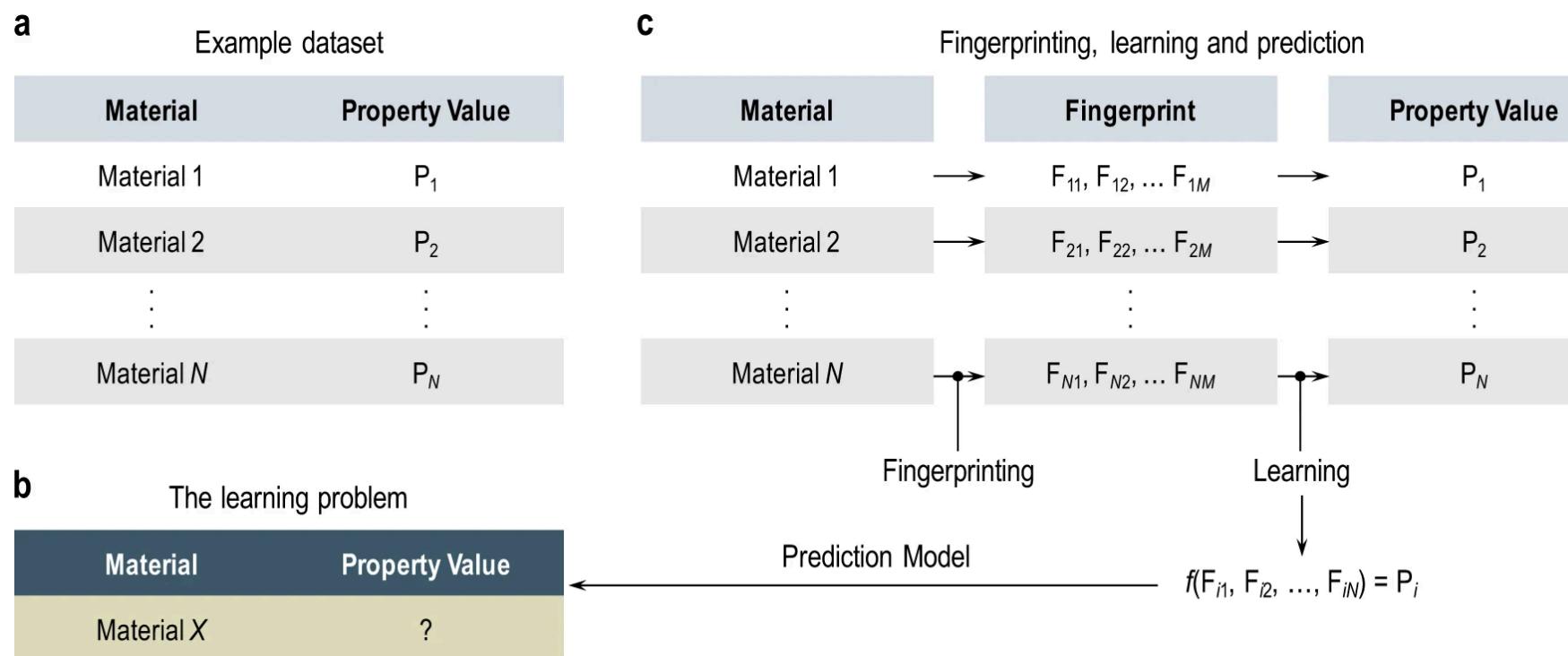
Example:

- number of trees in random forest, or depth of a decision tree
- λ in Ridge Regression
- They control the learning process itself rather than being learned from the data.

These parameters should be tuned to get better performance of the model

Features

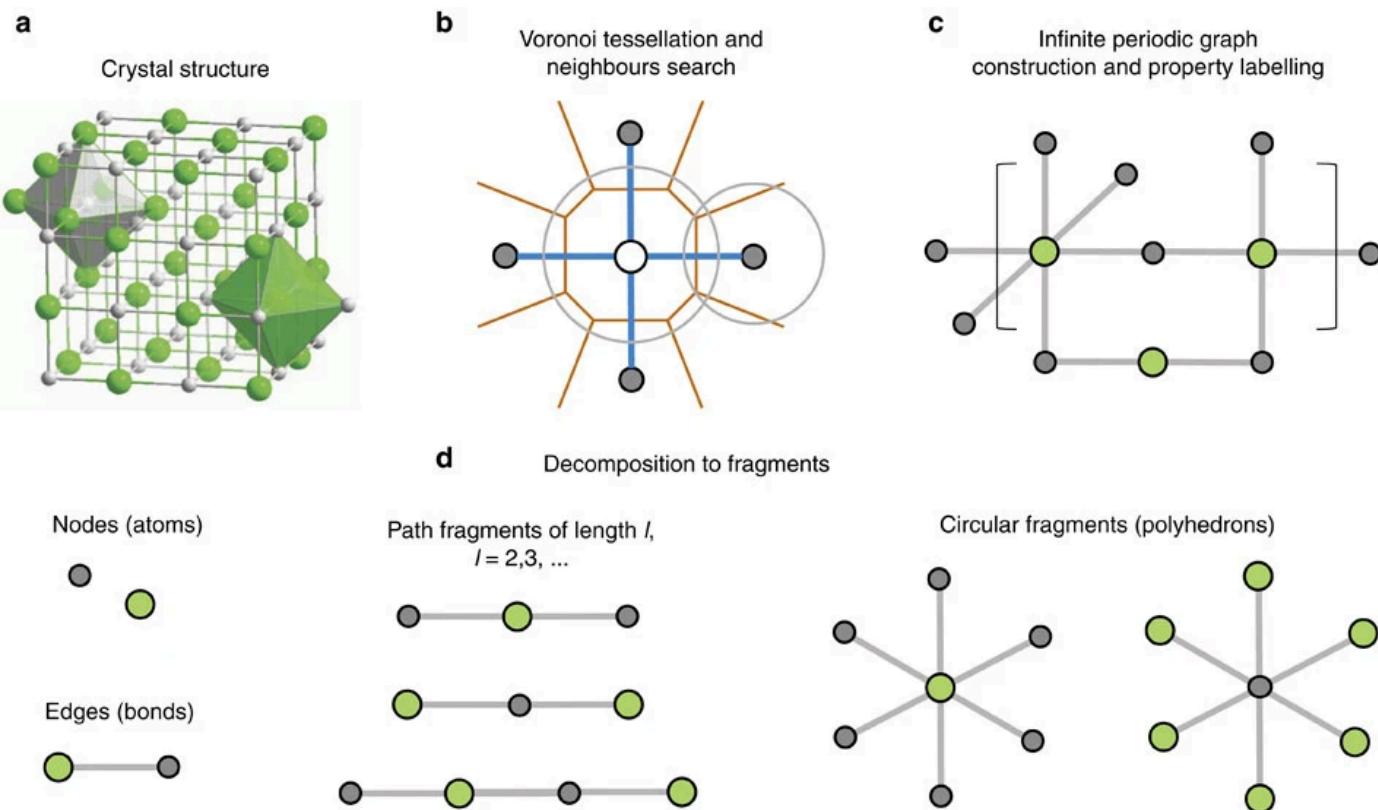
- In classical ML we use features
- to represent materials in a machine-readable format
- Think of it like a fingerprint



Hierarchy of features

Depending on the resolution we have:

- Local descriptors
 - site
- Fragment descriptors
 - bond
 - polyhedron
- Global descriptors
 - chemical family
 - structural type
 - density

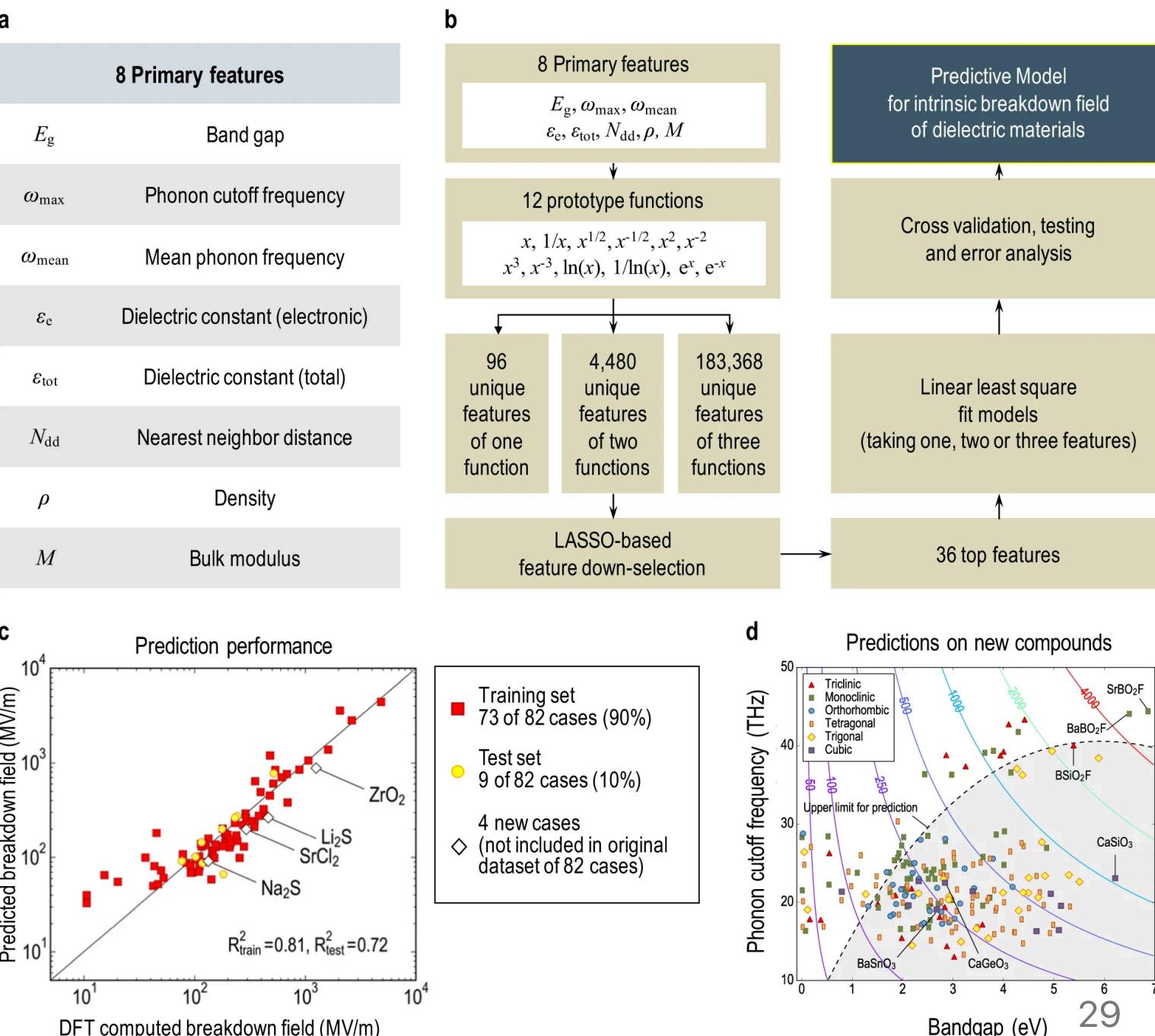


Feature engineering

- Primary descriptors are used to design more complex features
- e.g. by applying set of mathematical operators

See SISSO paper

Machine learning in materials informatics: [recent applications and prospects](#)



Permutation feature importance

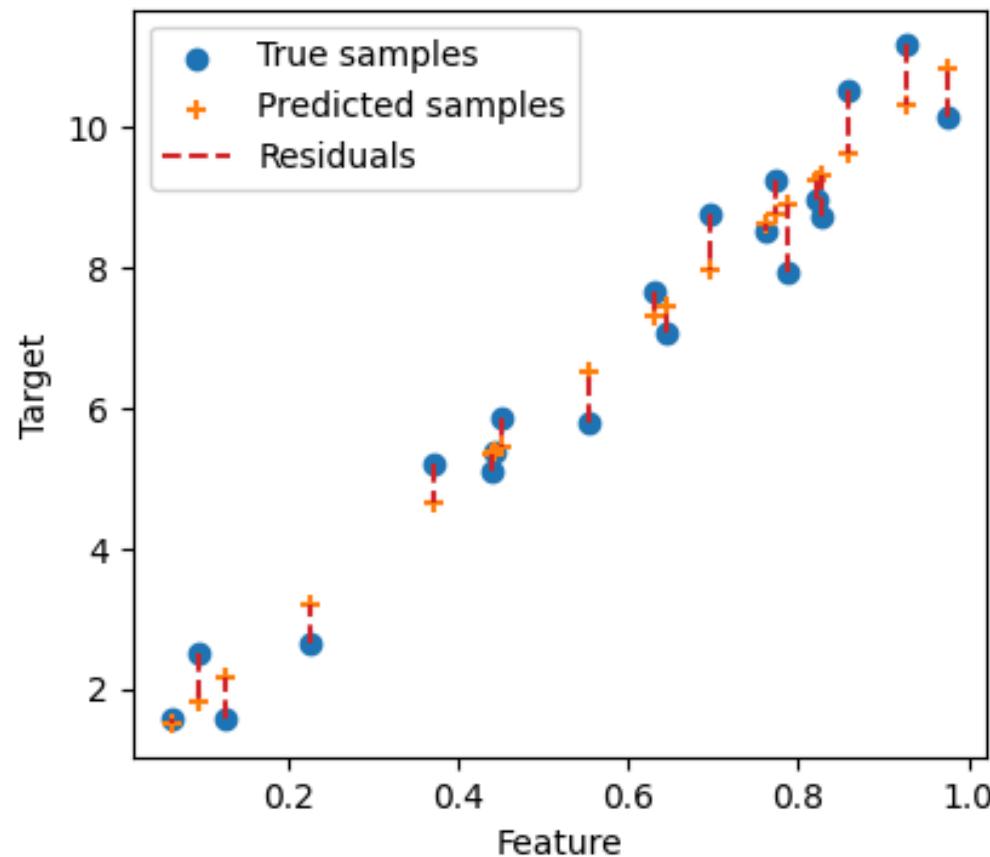
We want to know which features most influence model prediction

Given dataset $\{X, y\}$

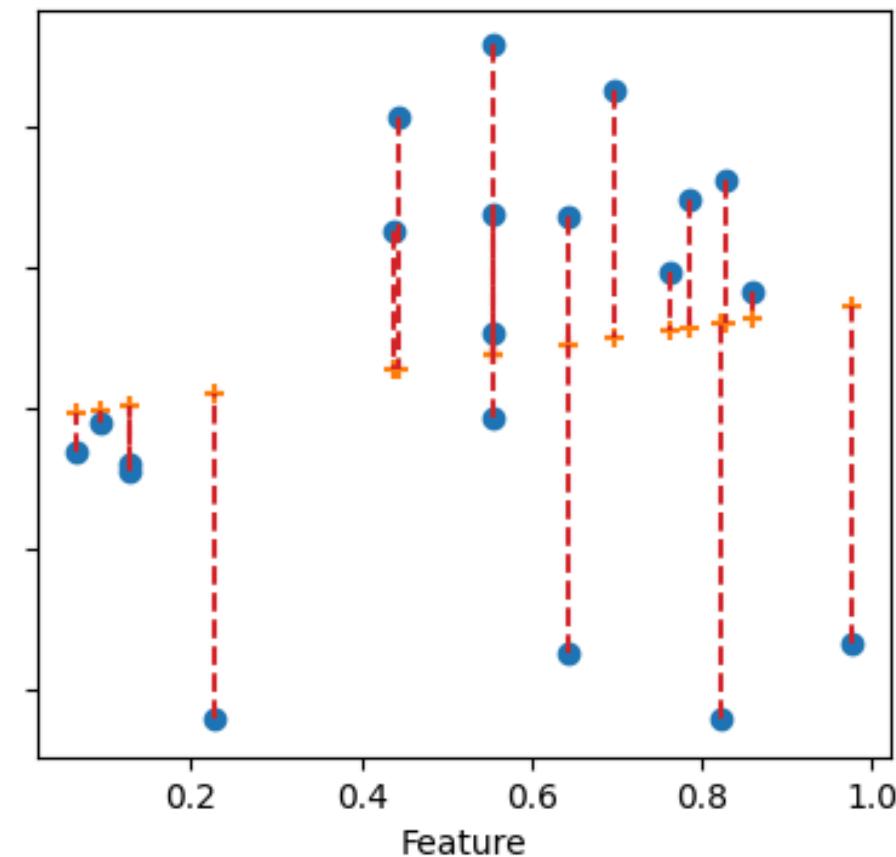
- Fit the model
- Get scores
- Randomly shuffle one of the feature vectors x_i
- Refit model
- Get scores
- The higher degradation of the model performance the more important the feature

Effect of permuting a predictive feature

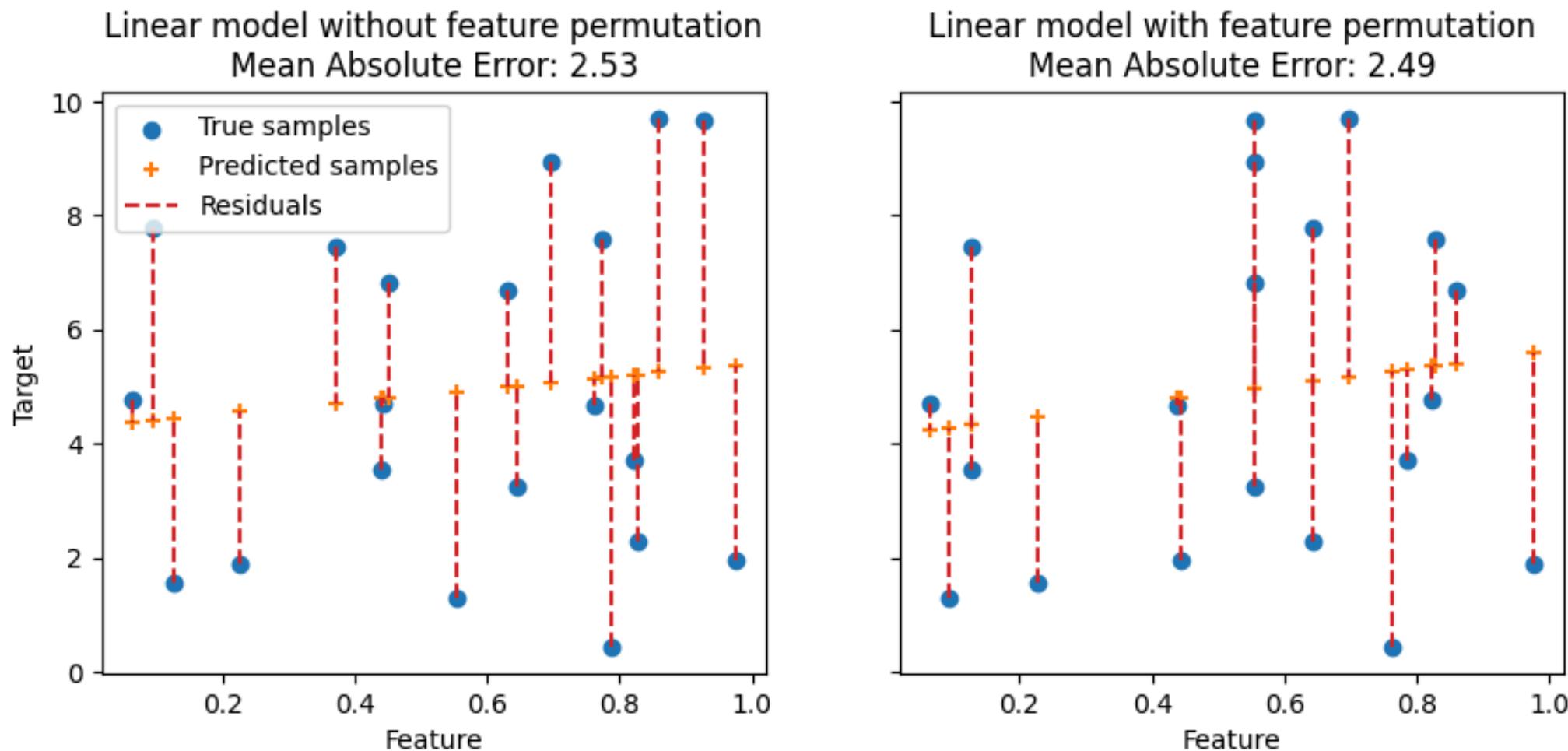
Linear model without feature permutation
Mean Absolute Error: 0.51



Linear model with feature permutation
Mean Absolute Error: 2.28



Effect of permuting a non-predictive feature



The most important properties of an ideal descriptor:

- Invariant with respect to translation of the coordinate system
- Invariant with respect to rotation of the coordinate system
- Invariant with respect to permutation of atomic indices: changing the enumeration of atoms does not affect the target
- Unique: single way to construct a descriptor and the descriptor itself corresponds to a single property
- Continuous: small changes in the atomic structure -> small changes in the descriptor
- Compact
- Computationally cheap

From [DScribe: Library of descriptors for machine learning in materials science](#) by Lauri Himanen et al. (Computer Physics Communications, 2020)

Machine learning for molecular simulations

Finite difference method

The shape of $V(R)$ is known at the given moment t

The acceleration for each particle is given by

$$a_i = \frac{F_i}{m_i} = -\frac{1}{m_i} \frac{\partial V}{\partial r_i} \text{ eq. (1)}$$

0. Initialize velocities
1. Select a small timestep δt (it should correctly describe the changes in your system)
2. Calculate forces
3. Integrate eq. (1) within $(t, t + \delta t)$ window to get velocities $v_i(t + \delta t)$
4. Update coordinates using calculated velocities for each particle
$$\mathbf{r}_i = \mathbf{r}_i + \delta t \mathbf{v}_i$$
5. Repeat 2-4 steps N times
6. Analyse the discrete trajectory obtained

Energy calculations with interatomic potentials

- the potential energy is calculated as the sum of atomic potential energies

$$E_{pot} = \sum_{i \in N_{atoms}} E_{i,atomic}$$

- forces are calculated as follows

$$\vec{F}_i = -\nabla_i E_{pot}$$

Fitting (training) the potentials

We run AIMD at elevated temperatures to sample atomic configurations

- for a given chemical system

The dataset:

Geometries: $\{R_n\}$

Energies: $\{E_n(R_n)\}$

Forces: $\{F_{i,\alpha}(R_n)\}$

Loss function

$$\mathcal{L} = \lambda_E \|\hat{E} - E\|^2 + \lambda_F \frac{1}{3N} \sum_{i=1}^N \sum_{\alpha=1}^3 \left\| -\frac{\partial \hat{E}}{\partial r_{i,\alpha}} - F_{i,\alpha} \right\|^2, \text{ where the hat symbol (e.g., } \hat{E} \text{) denotes predictions}$$

Universal interatomic potentials

- Once trained on large datasets covering a wide chemical and structural configuration space
 - For example, the Materials Project dataset
- Can be applied to the unknown systems
 - To perform downstream tasks

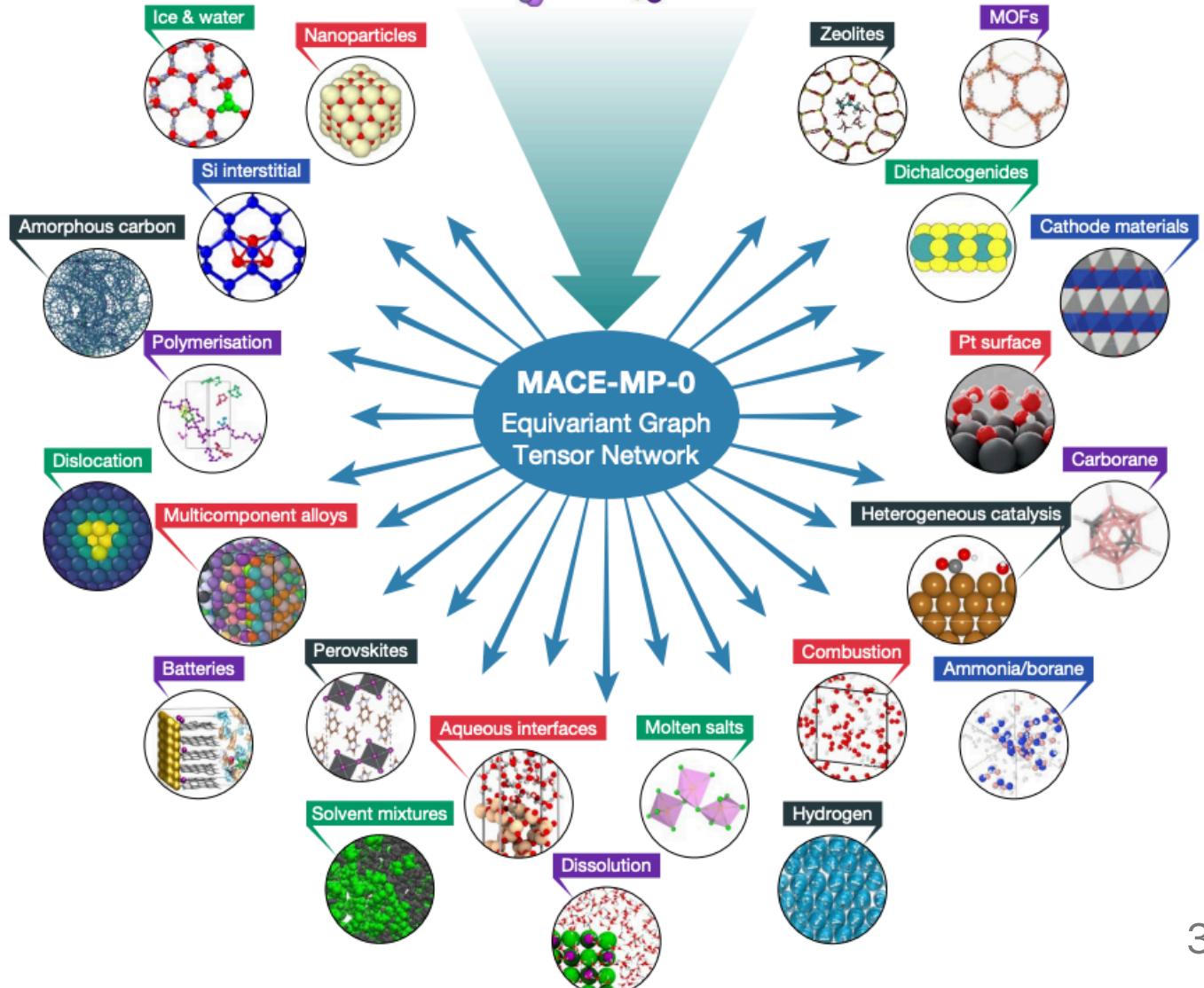
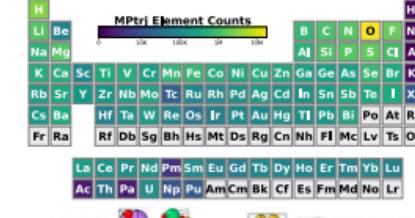
Pros:

- Works out of the box
- Can be refitted (finetuned) with a small amount of data

Cons:

- Universality is debatable. Should be carefully validated before using
- A lot of weights -> slower compared to ML potentials developed for a specific system

MATERIALS PROJECT



Example

MACE-MP-0 graph neural network

Take home message

Materials informatics

- is about effectively using data, statistics, and data science tools for accelerating the trial-and-error process in materials science
- is not always about training ML models

Python (or another programming language) these days

- is support for your development as a successful scientist
- saves your time

Data

- respect/explore/use/share data

ML

- there is a trade-off between accuracy and speed
- expertise is essential for reliable results
- don't use test data split for model training