

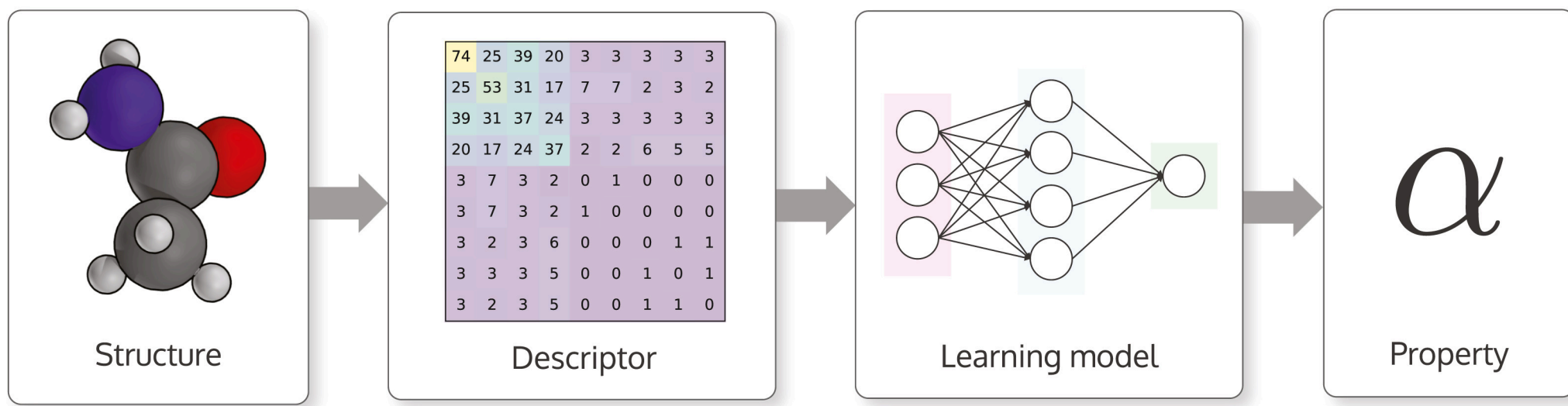
Lecture #6: Atomic structures encoding

Previously on

- Supervised machine learning
- Ridge regression
- Random forest

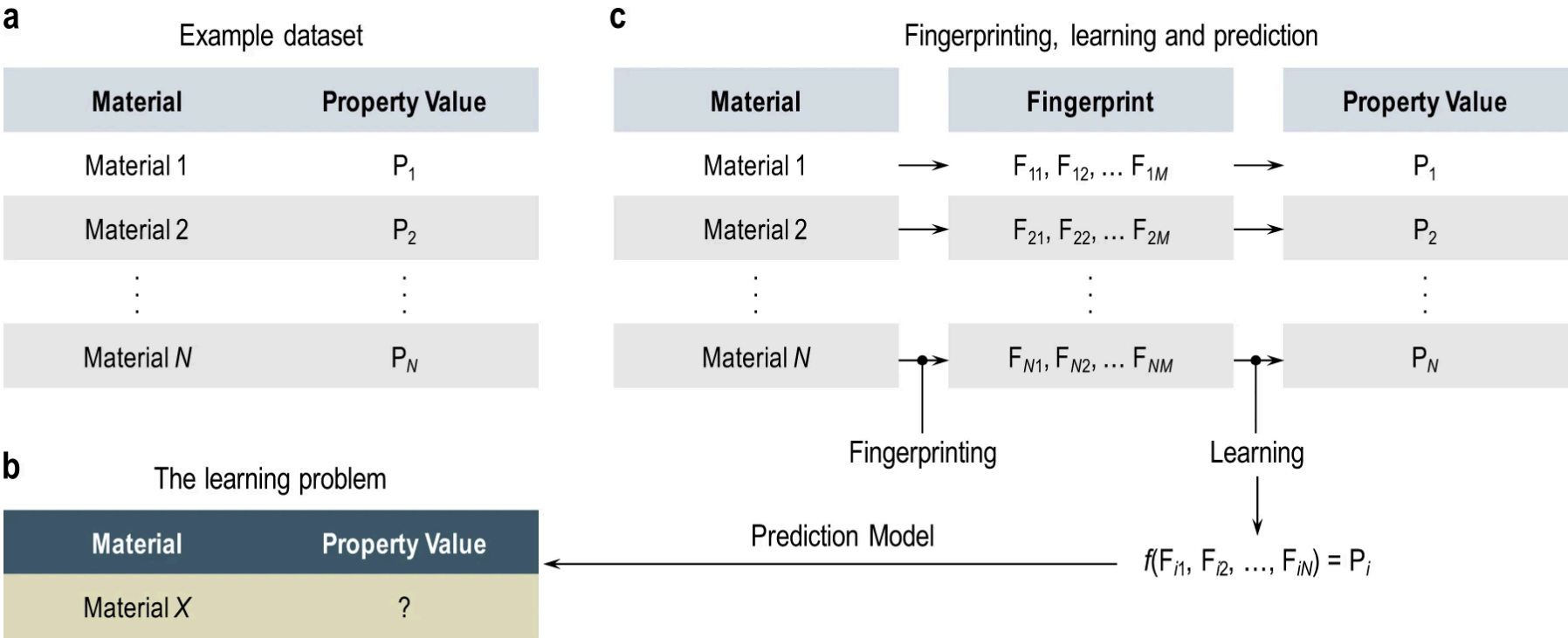
Goals/Agenda

- Hierarchy of crystal/molecular structure descriptors
- Feature importance
- Crystal structure fingerprints



Features

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



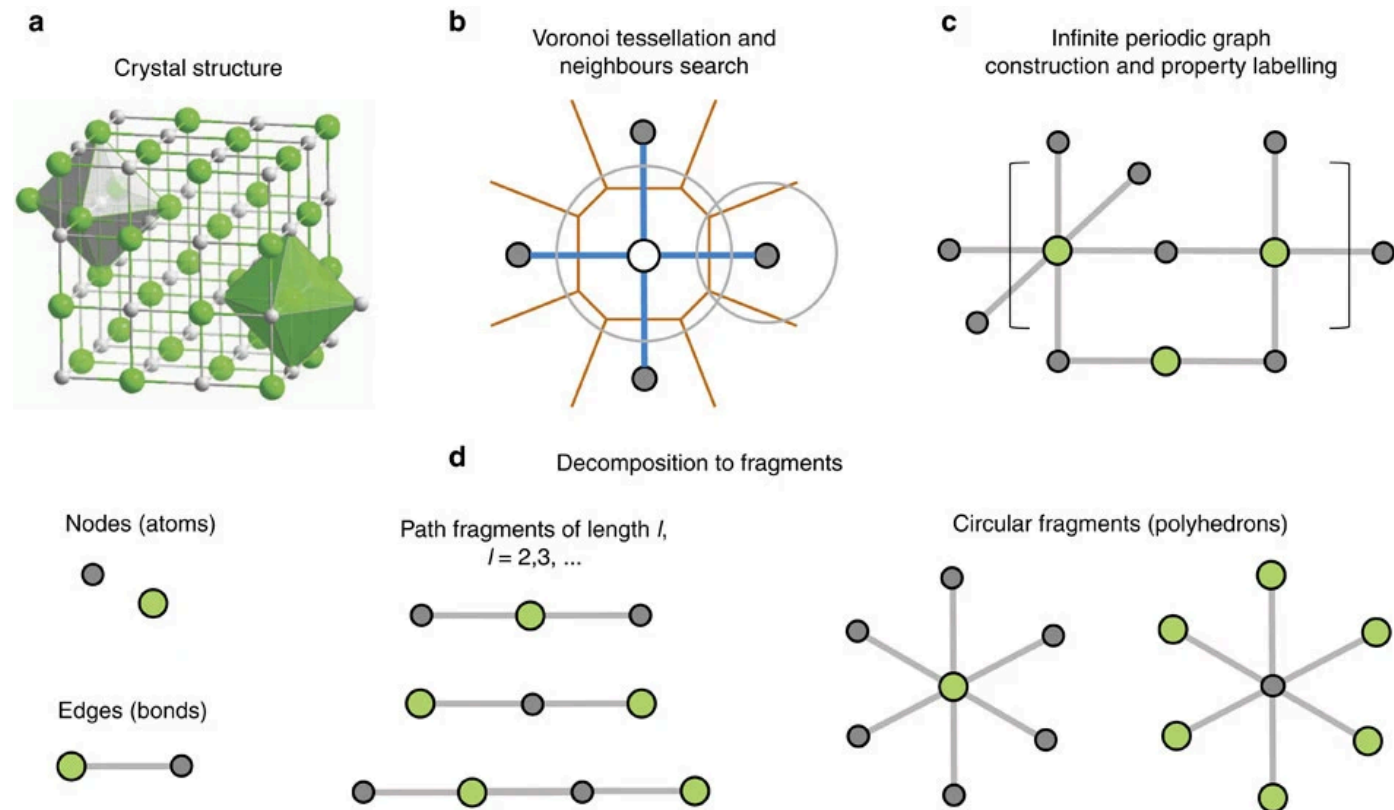
The better you find this representation for a particular problem, the better your model fits the data.

Geometrical and compositional encoding of atomic structure

Heirarchy of features

Depending on the resolution we have:

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



Geometrical (Structural) features

- Atomic packing (e.g. volume per atom, density)
- Voronoi polyhedra features
 - area, volume, face distance, solid angle
- bond distance/angle

Can be calculated for sublattices

Min/max/mean statistics

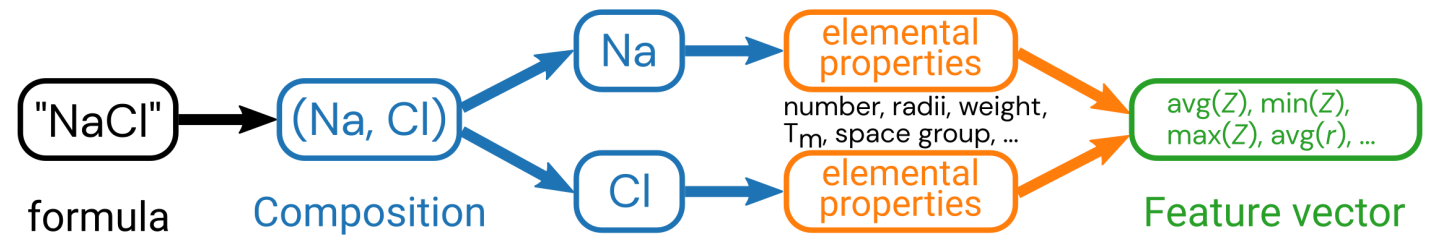
Voronoi polyhedra features can be collected by pymatgen's methodology

```
from pymatgen.analysis.local_env import VoronoiNN  
features = VoronoiNN().get_voronoi_polyhedra(structure, site_id)
```

Compositional (elemental/atomic) features

Aggregation over

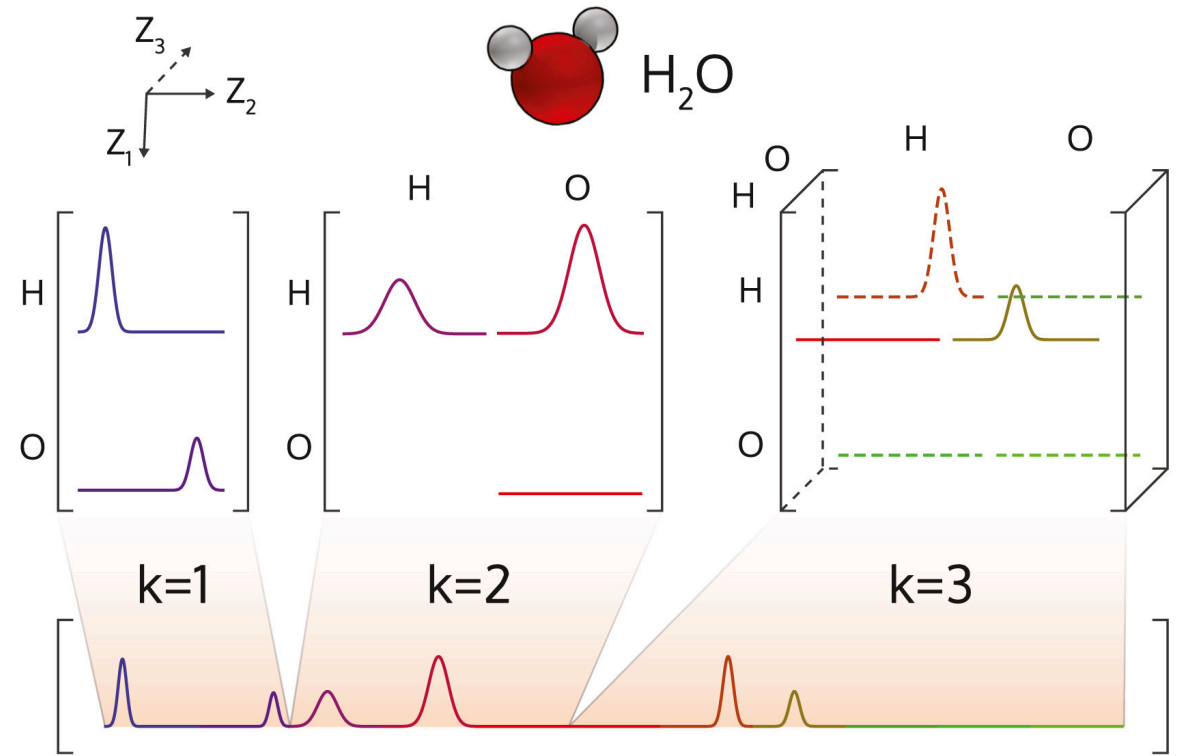
- Atomic numbers
- Valence electrons
- Covalent radii
- Atomic fraction of each element in a composition
- Stoichiometry related
- Electron affinity
- Ionic radii
- Oxidation states
- Electronegativity
- ...



Atomic structure fingerprints

(not taught in the course)

- Atom-centered Symmetry Functions
- Smooth overlap of atomic positions
- Many-body tensor representations
- Pair distribution function
- X-ray diffraction pattern

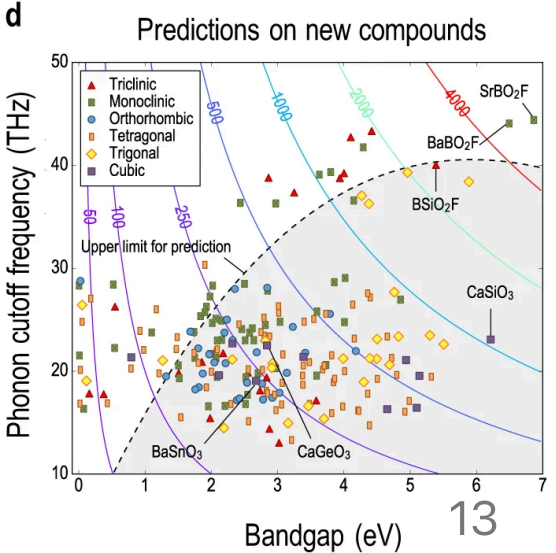
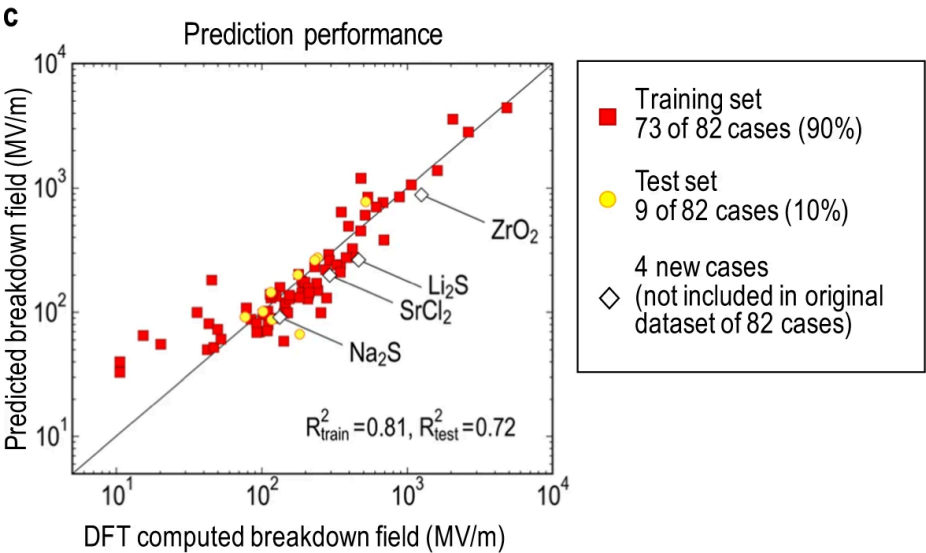
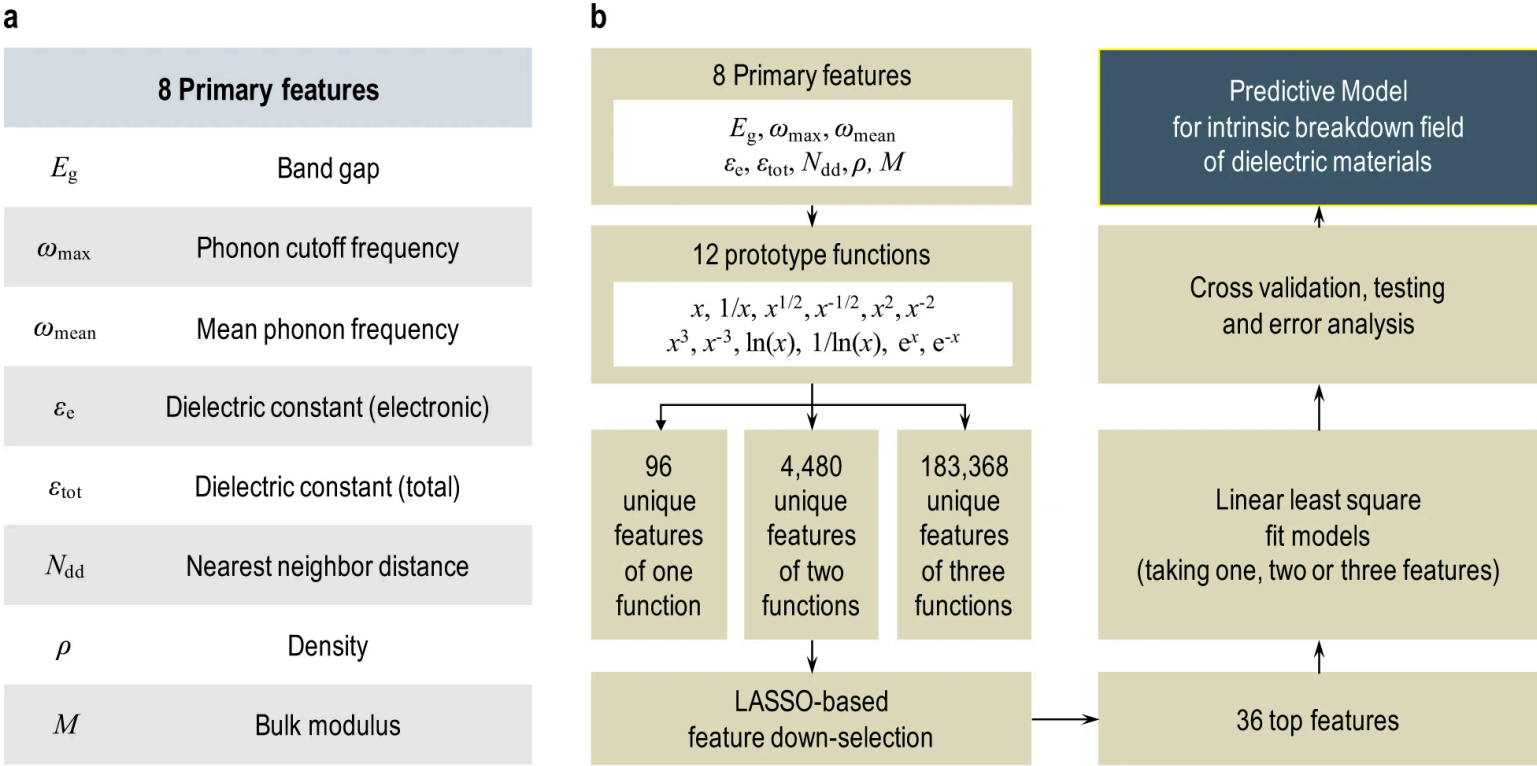


Feature engineering

- 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](#)



Pros

- Better performance compared to primary descriptors

Cons

- Increased computational complexity
- Garbage feature needs to be filtered
- Lack of interpretability

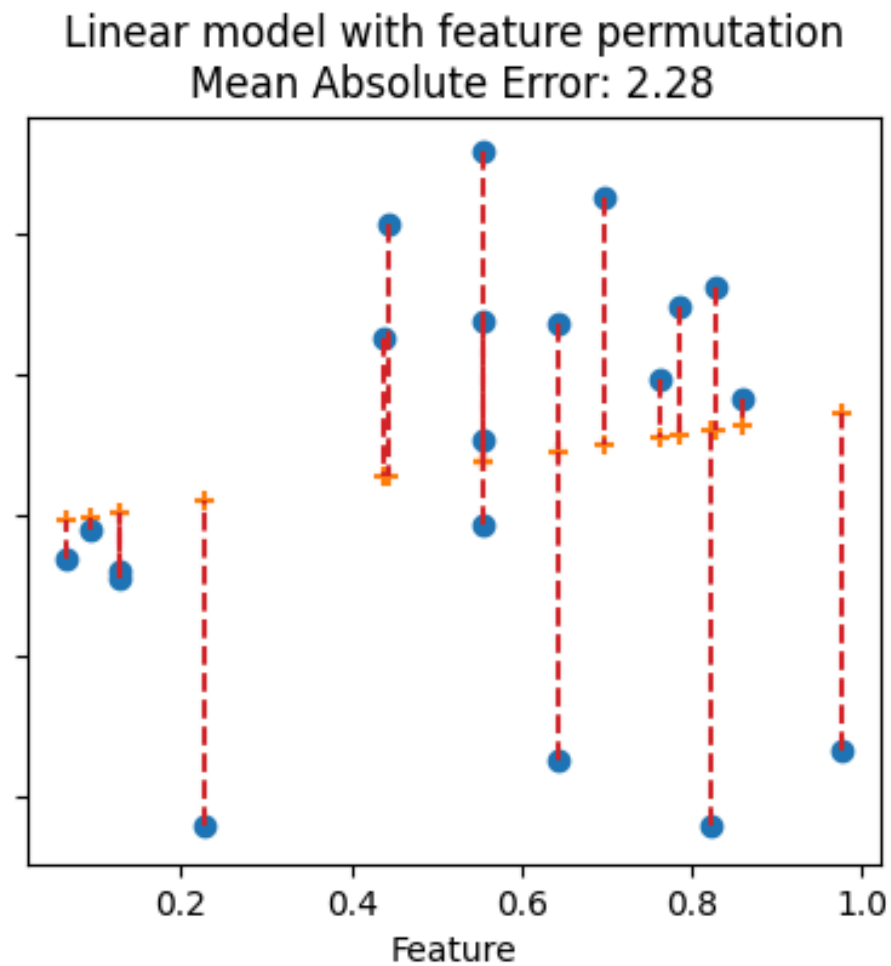
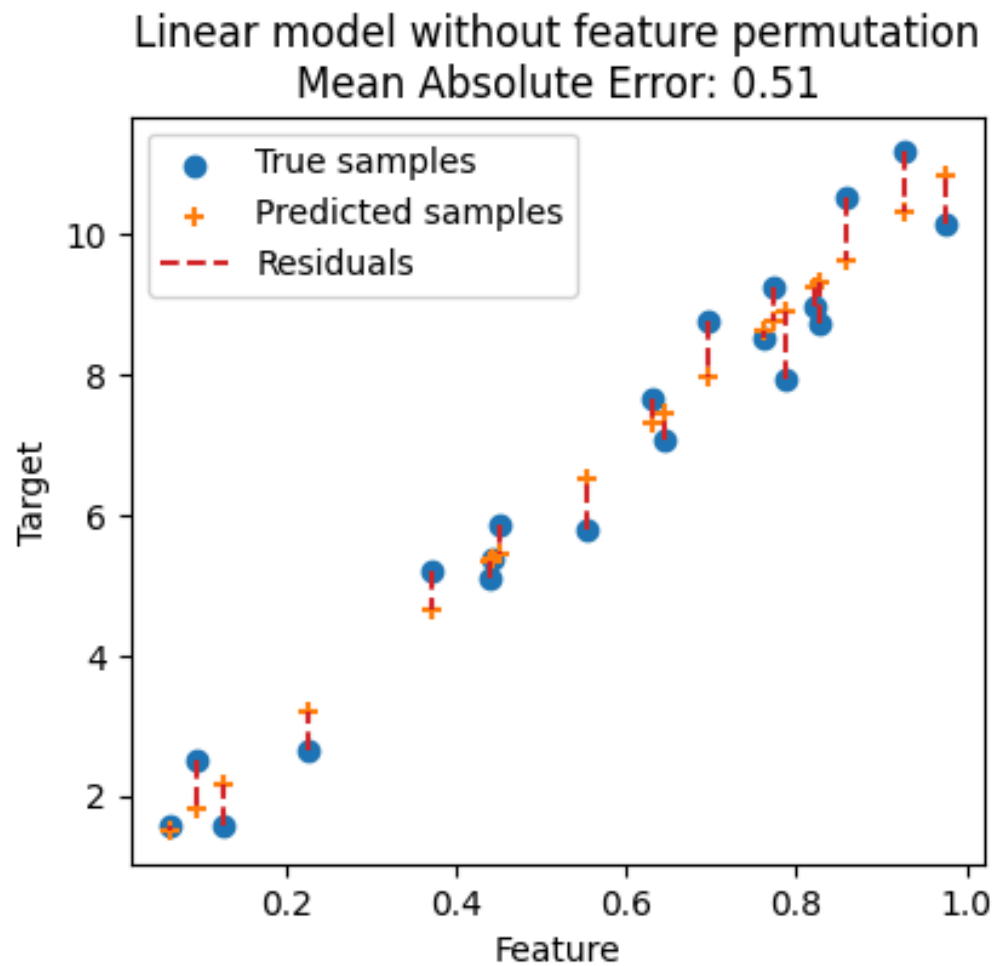
Permutation feature importance

We want to know which features most influence model prediction

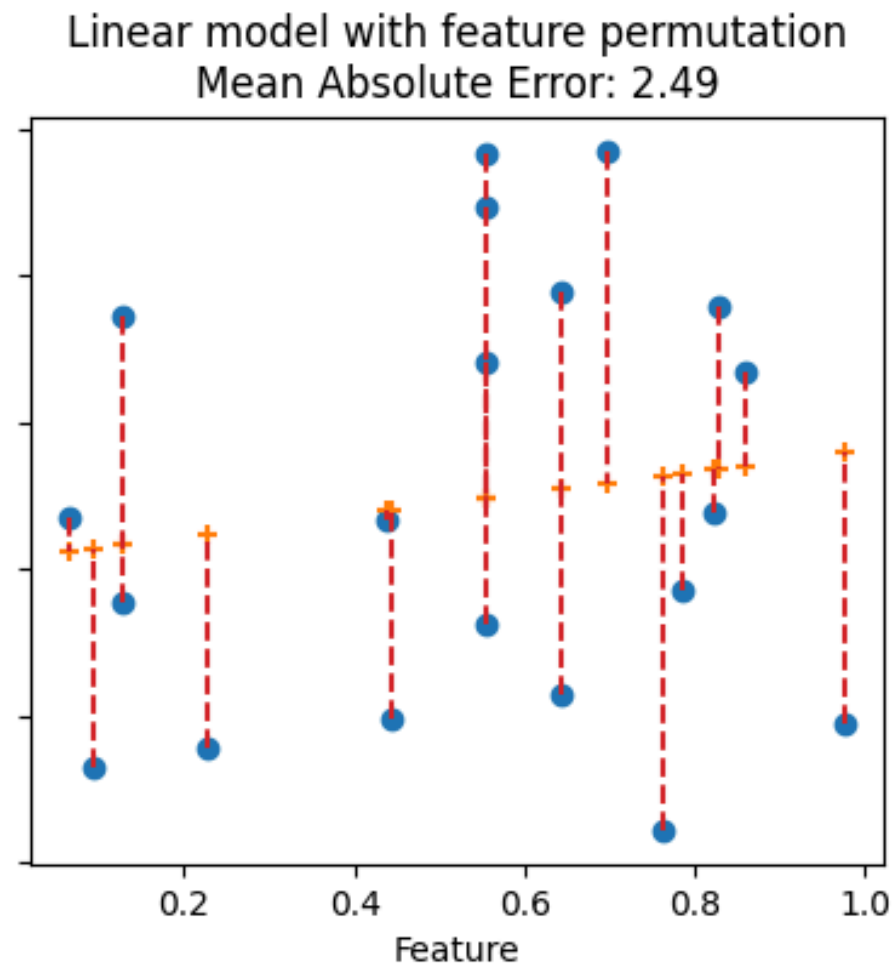
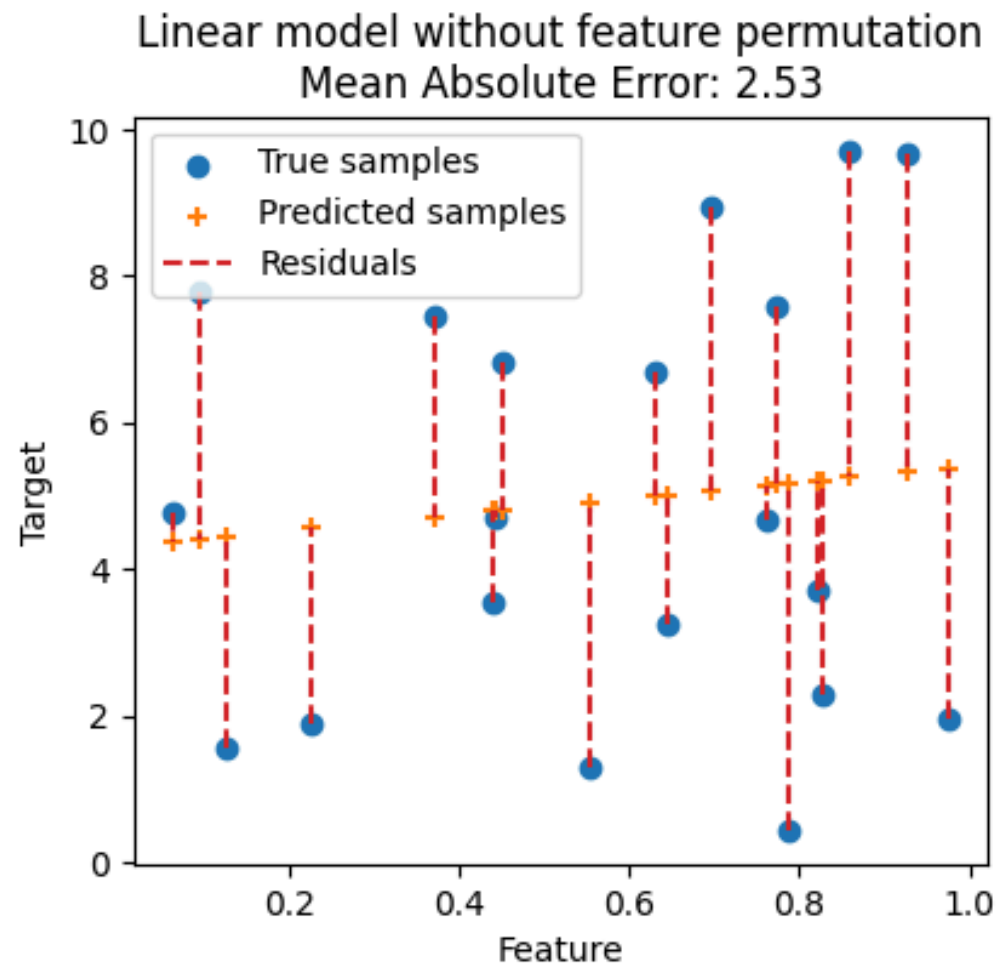
Given dataset $\{X, y\}$

- 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



Effect of permuting a non-predictive feature



Pros

- Works for any model
- Yields variance of the feature importance

Cons

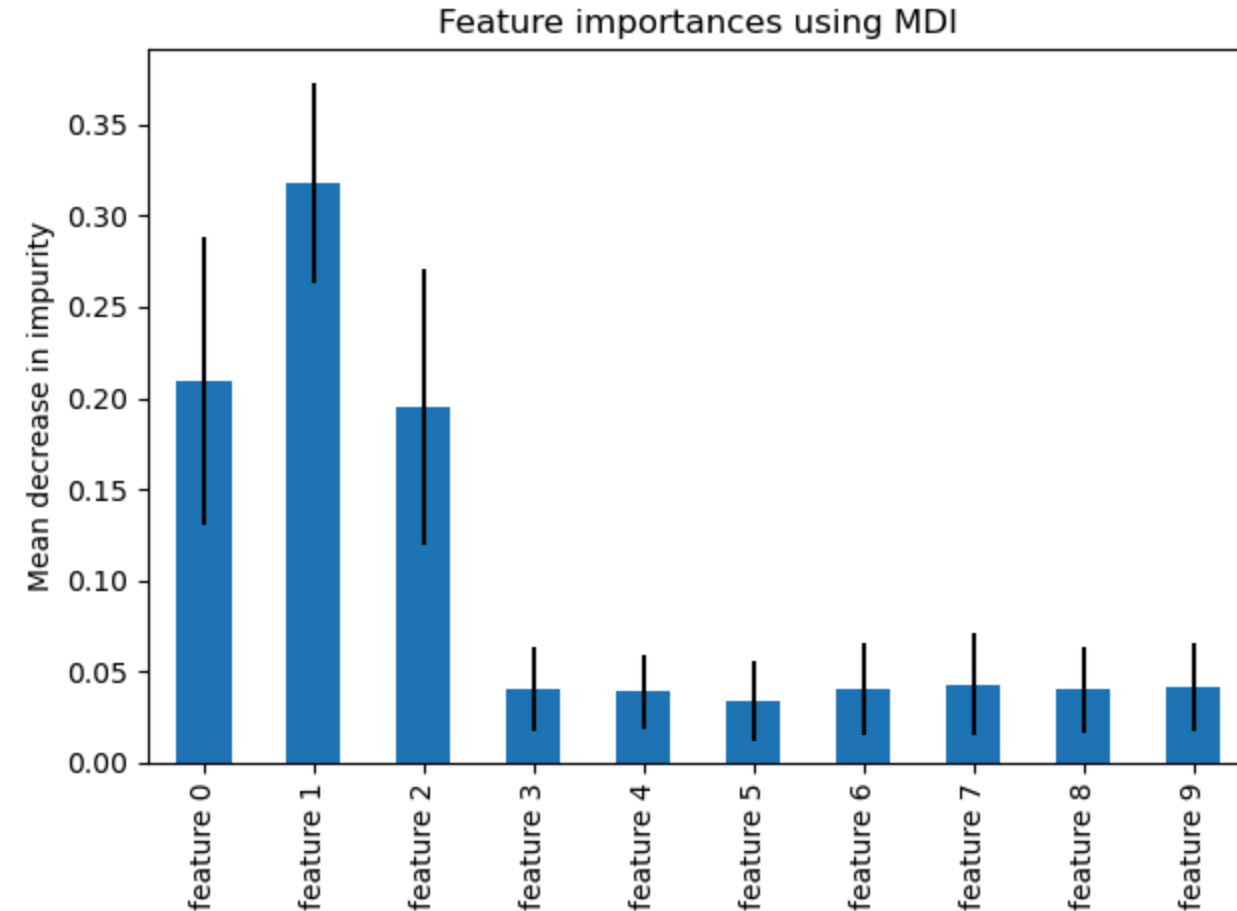
- Computational complexity
 - Requires refitting the model for each feature
- Wrong output for correlated features

Intuition behind the impurity-based feature importance

Decision tree:

- Selects the best feature for splitting at each node
- Best split is measured by gain in purity (Gini impurity or Entropy)
- The higher gain the higher importance of the feature

In the case of Random Forest the feature importances are "computed as the mean and standard deviation of accumulation of the impurity decrease within each tree."



Pros

- You have this information after fitting the model
- Yields variance of the feature importance

Cons

- Wrong output for high-cardinality features

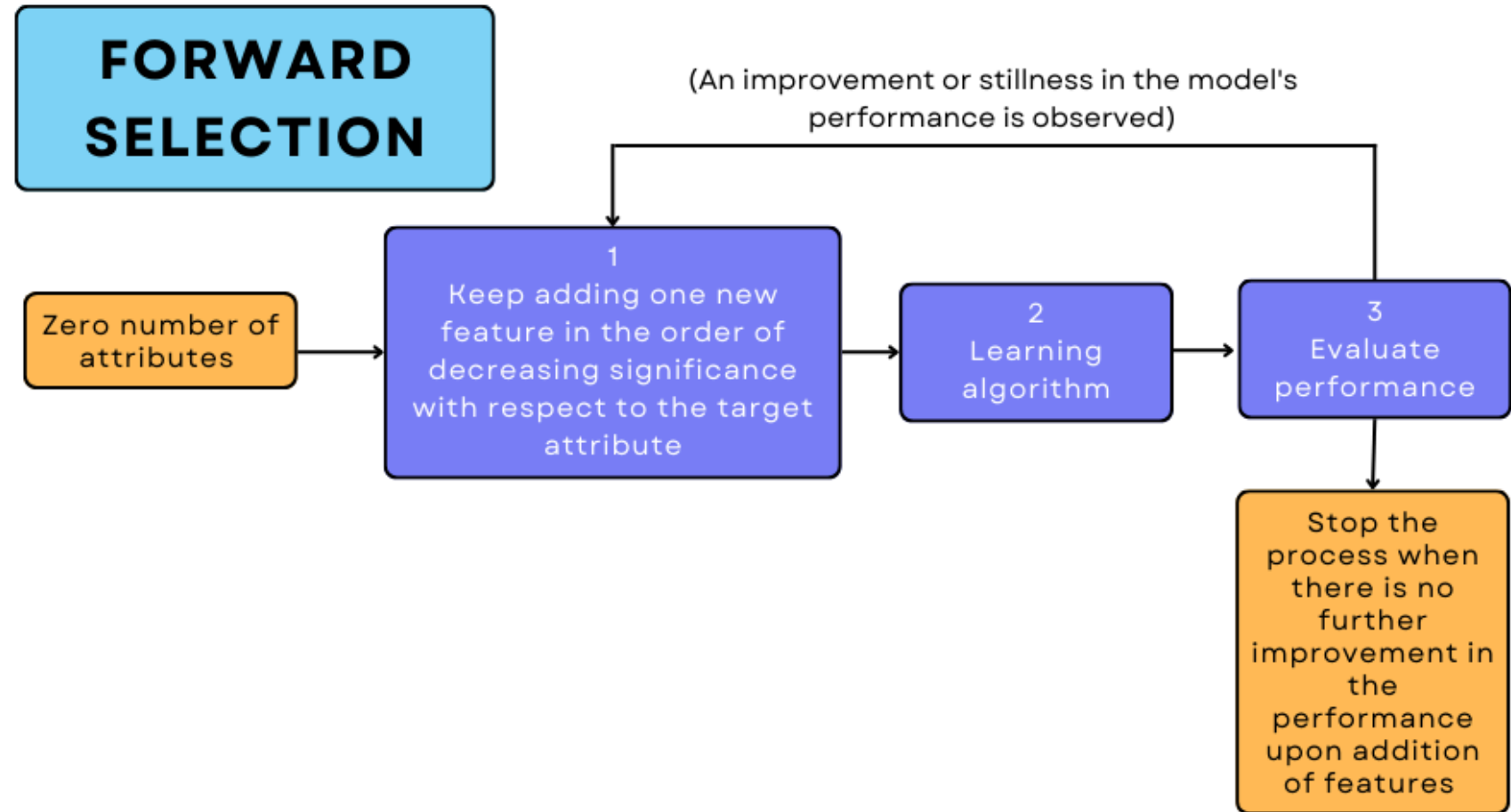
Feature selection

- Drop features with low variance
- Backward/Forward feature selection
- LASSO
- Tree-based feature selection

Why?

- Remove redundant features
- Develop robust model
- The less features the less compute is needed
- Feature collection = time = money
 - save your resources

Forward feature selection algorithm



How to collect features?

Design your own Featurizer

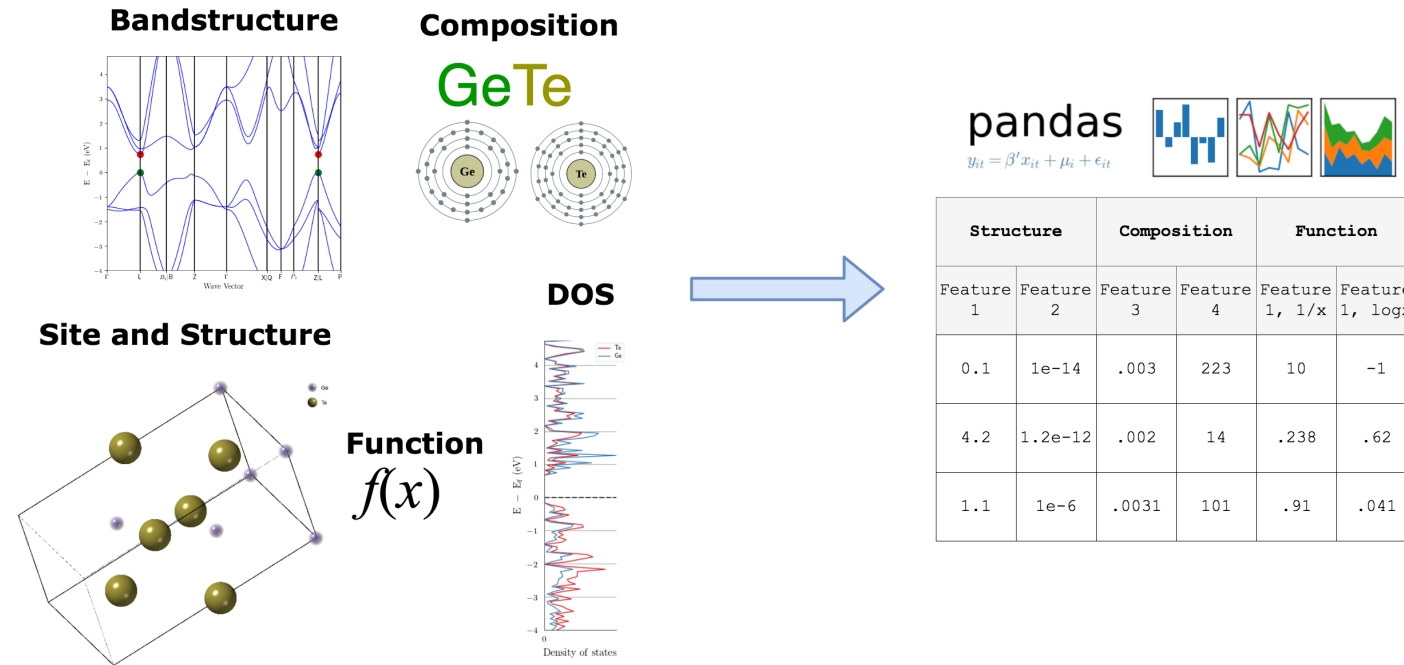
- You are an expert in your field
- Consider essential descriptors of your target
- Use pymatgen/ase/etc.

Ready-to-use Featurizers

- [Matminer](#) Python Library
- [Dscribe](#) Python library

Gross level features

- Generate gross level features non-specific to your task
- Select the best candidates



The most important properties of an ideal descriptor:

- Invariant with respect 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

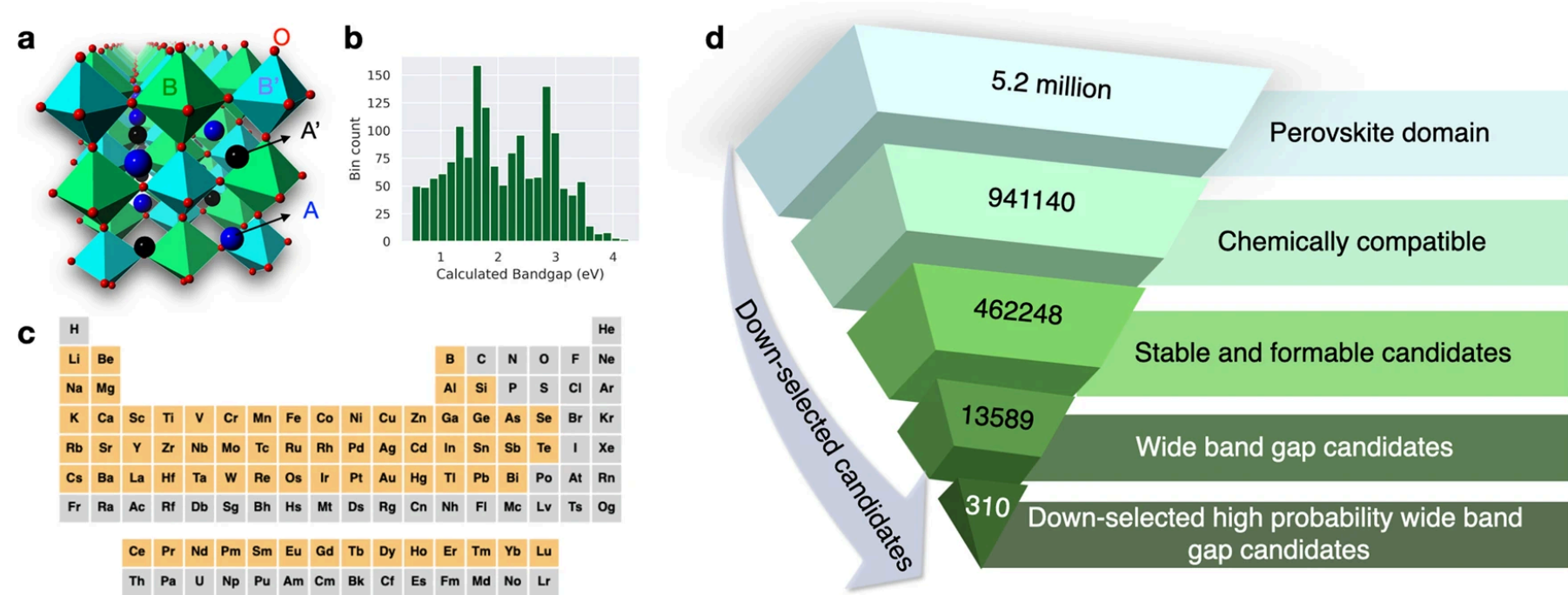
From [DScibe: Library of descriptors for machine learning in materials science](#) by Lauri Himanen et al. (Computer Physics Communications, 2020)

Which features to consider?

- What accuracy do I need?
- At which scale your property is defined?
 - Site (i.e. defect formation energy)
 - Bond (i.e. migration barrier)
 - Structure (i.e. hardness)
- Size of the chemical space?
 - Several elements or the whole periodic table
- Diversity of crystal structures?
 - Fixed structural type? Not fixed?
- Do we know (expect) any relationship?
- How many samples do I have?
 - ~100? ~100,000?

You can find the right feature design by answering these questions

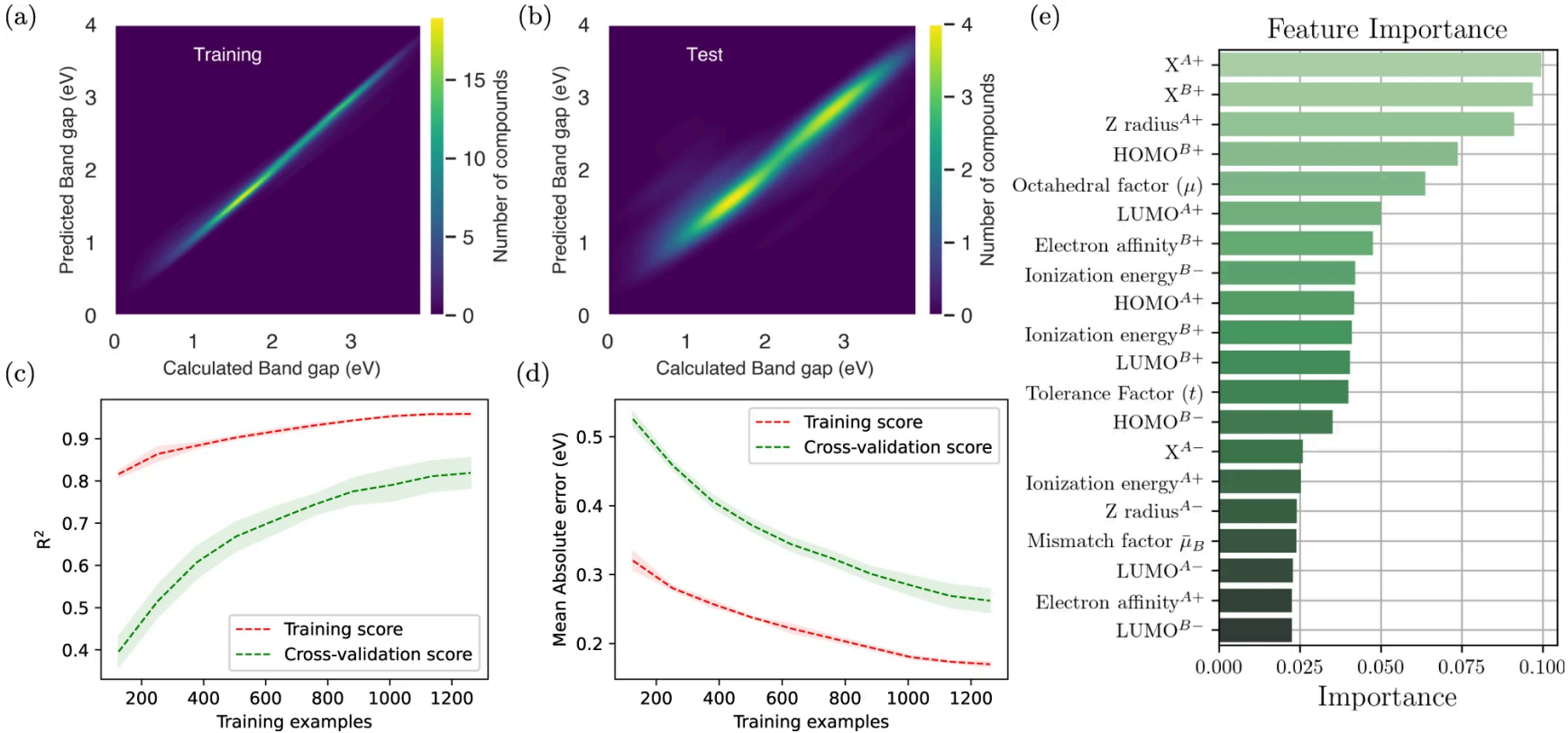
Example



Considered features

Abbreviation	Feature
<i>Elemental</i>	
HOMO	Highest Occupied Molecular Orbital (eV)
LUMO	Lowest Unoccupied Molecular Orbital (eV)
IE	Ionization energy (kJ/mol)
X	Pauling Electronegativity
Z radius	Zunger's Pseudopotential radius (a.u.)
EA	Electron affinity (kJ/mol)
<i>Geometric</i>	
t	Tolerance factor
μ	Octahedral factor
$\bar{\mu}$	Mismatch factor

Features for E_g prediction



Band gap predictions of double perovskite oxides using machine learning

Thank you for your attention!