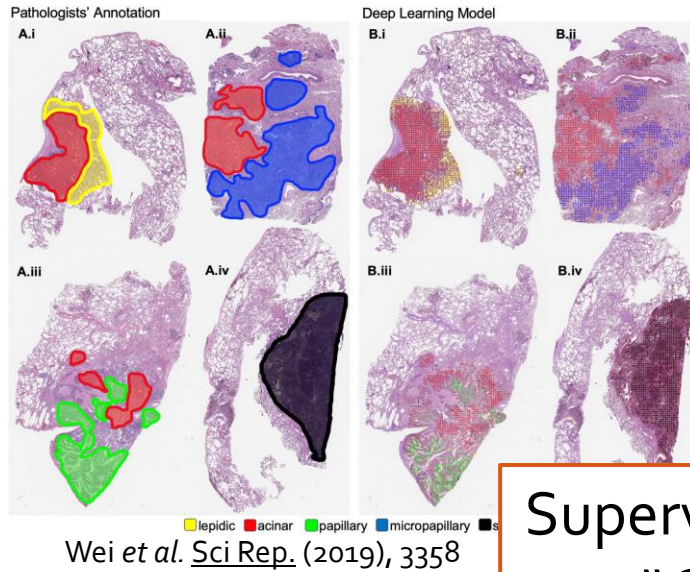


REPRESENTATIONS: HOW TO PERFORM SUPERVISED LEARNING ON MATERIALS DATA

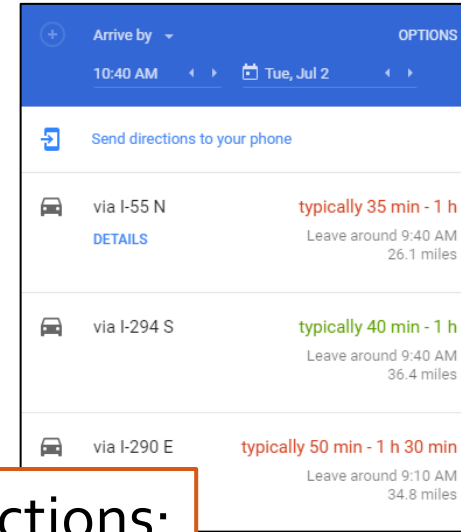
Logan Ward
Asst. Computational Scientist
Argonne National Laboratory

15 January 2022

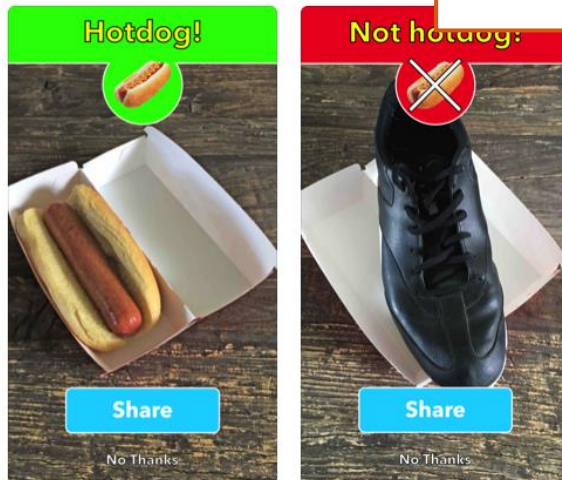
Supervised Learning: The ML you've definitely used



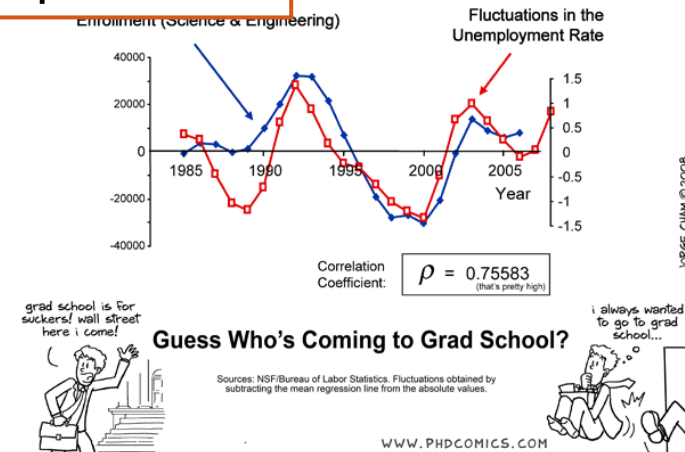
Supervised Learning models functions:
"Given inputs, predict output"



Google Maps

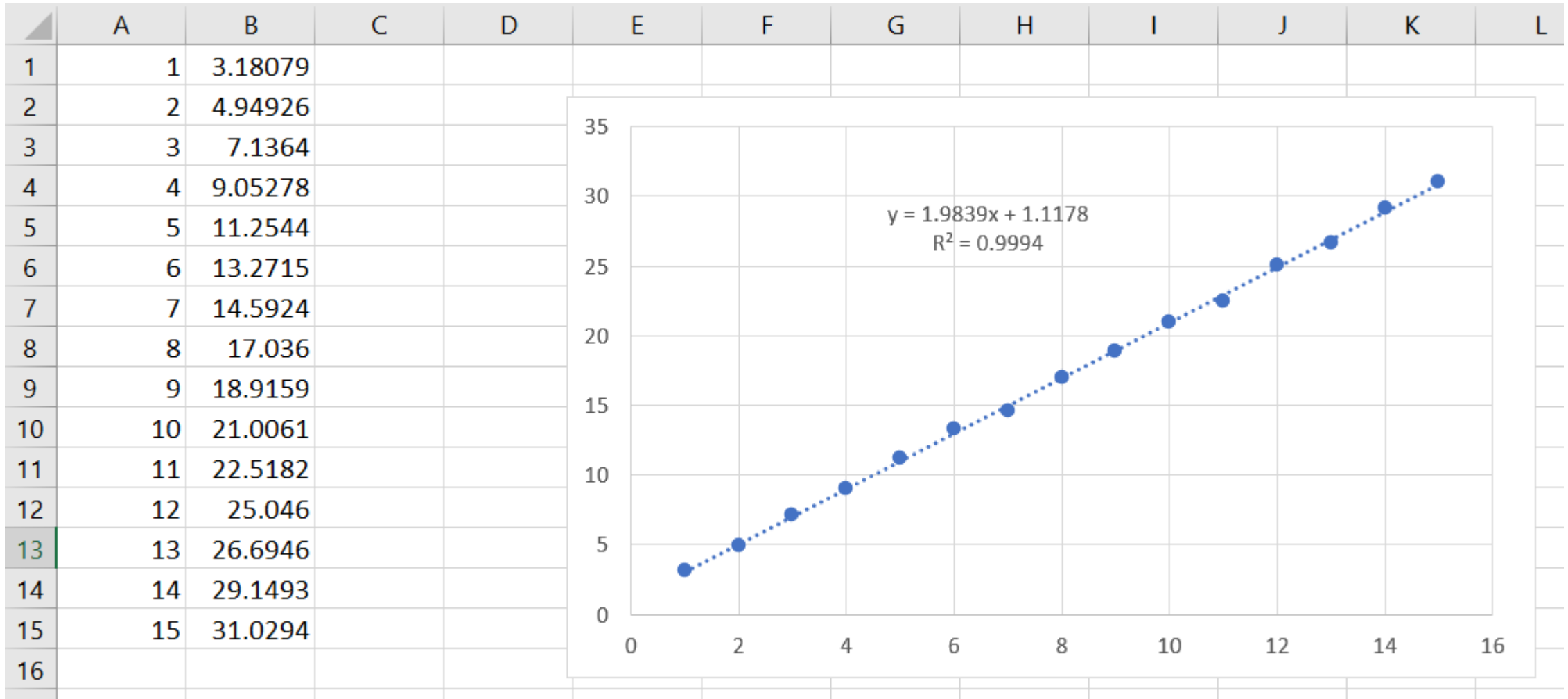


<https://apps.apple.com/us/app/not-hotdog/id1212457521>



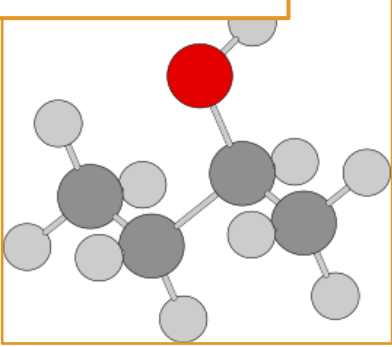
<http://phdcomics.com/comics/archive.php?comid=1078>

You've already done machine learning



Problem: Not all data are tensors

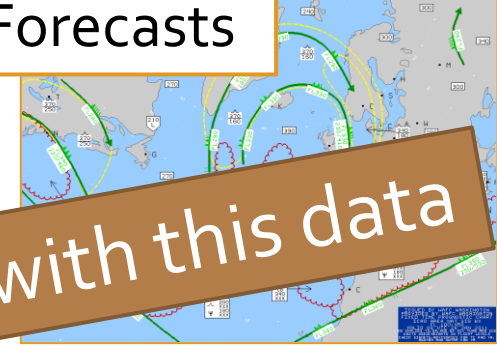
Molecules



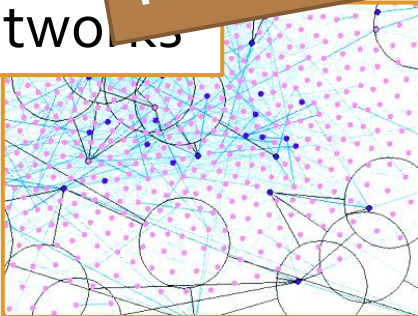
Images



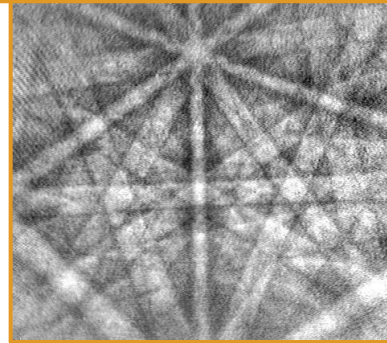
Weather
Forecasts



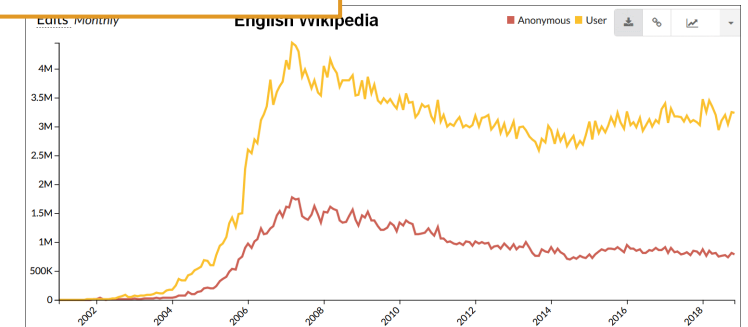
Social
Networks



EBSD Patterns

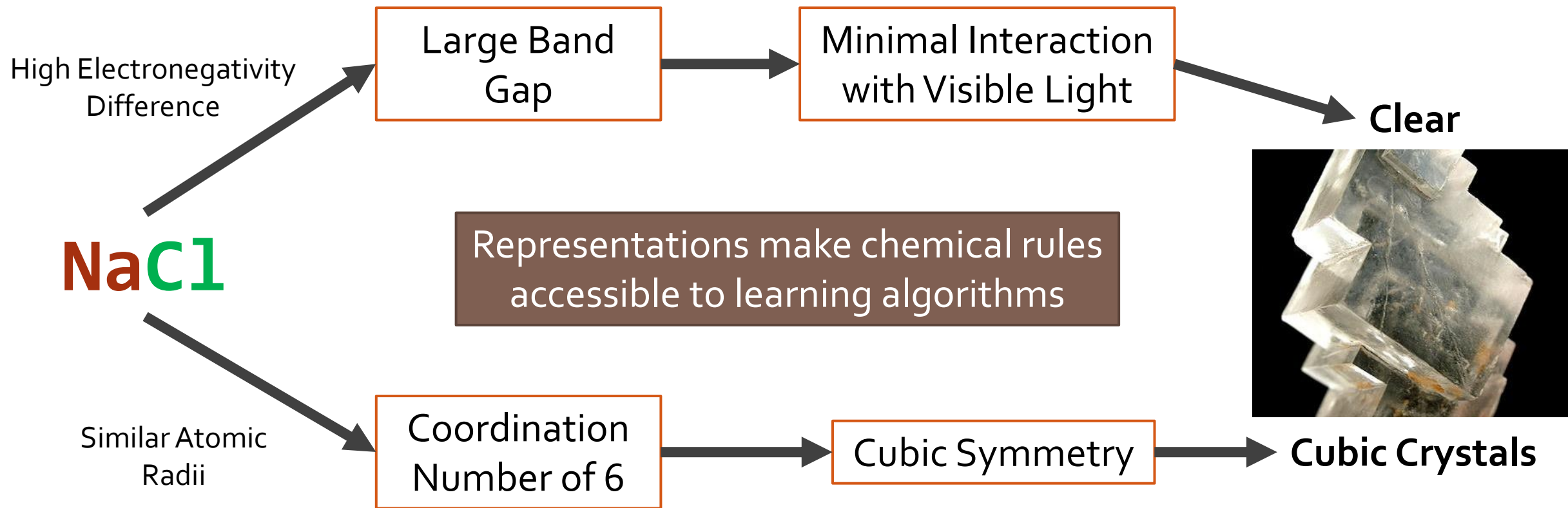


Timeseries



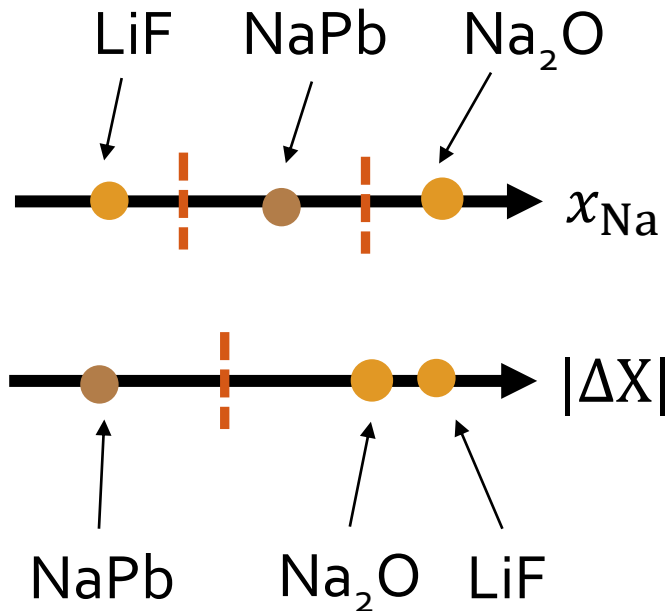
It's still possible to use machine learning with this data

What do NaCl crystals look like?



How to translate chemistry to a computer?

Representation: *Set of quantitative attributes that describe a material, molecule, ...*



Why do we need them?

- Machine learning tools take tensors

What makes a good one?

- Easy to learn generalizable rules

Rules for representations: Faber et al. (2015)

Complete: Different materials should have different representations

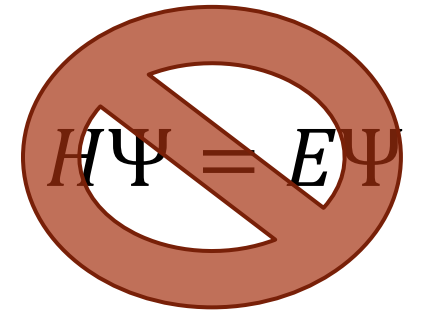
$$A \neq B$$

$$x_A \neq x_B$$

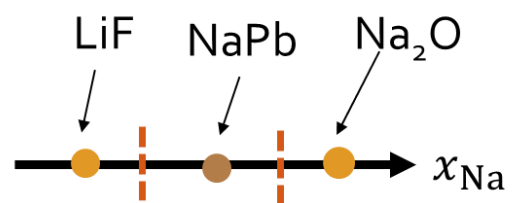
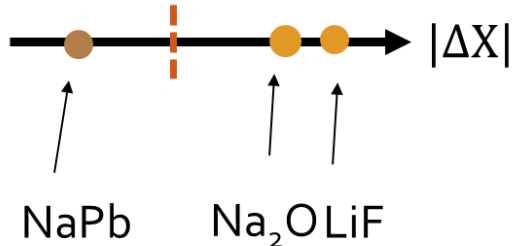
Compact: Representations should not include redundant features

Descriptive: Similar materials should have similar representations

Simple: Representations should be fast to compute

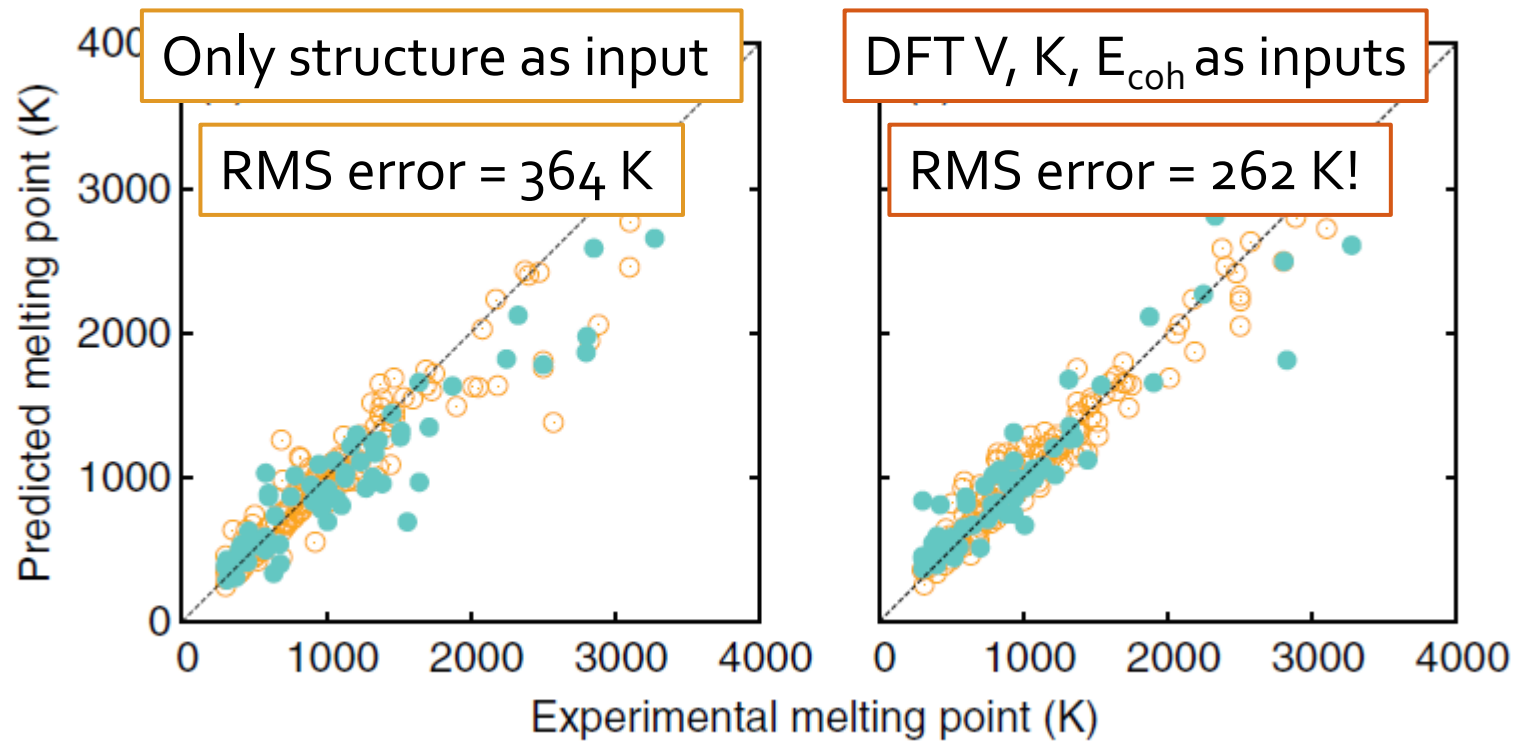


Example: Back to “Is it Ionic?”

	Fraction of Na	Difference in EN
	Complete	✓ Yes, no overlapping points
	Compact	✓ Yes, single feature!
	Descriptive	✗ No, rules are complicated
	Simple	✓ Yes, no compute
		✓ Yes, minimal compute

Representations standards are not clear cut

Seko et al. used DFT-computed properties as inputs to an ML model



Ok for 10^2 compounds. Not OK for 10^5 compounds

Key Concept: There is not and will not be a “one representation for all uses.”

Types of representations: Discriminative vs Descriptive

Discriminative

Make features that capture intuition

Element properties ($|\Delta X|$)
Interocular distances

Learn from little data
(potentially) Interpretable models

Introduces biases into the model
Need to know physics to solve problem

Concept

Examples

Advantages

Challenges

Descriptive

Make features that distinguish examples

Atomic fractions (x_{Na})
Pixel values in images

Maximum expressivity

Requires more data
(learn model *and features*)

“Every feature” vs “Hand Selected”

Use Every Features You Can Think Of

- Use when physics poorly understood
- Can find unexpected linkages

Advantages

- Danger of overfitting
 - Find patterns not due to physics
 - More complex models
- Requires algorithms with feature selection

Disadvantages

Selecting by Hand

- Avoid fitting to bias in data
- Little need for feature selection

- Requires understanding physics
- Limited to rules you expect to find

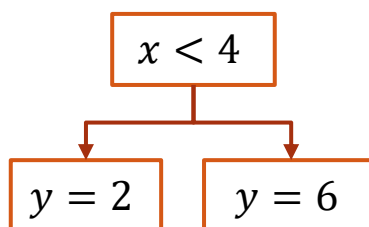
Different representations for different algorithms

Kernel Ridge Regression

$$f(x) = K(x, x'_i)$$

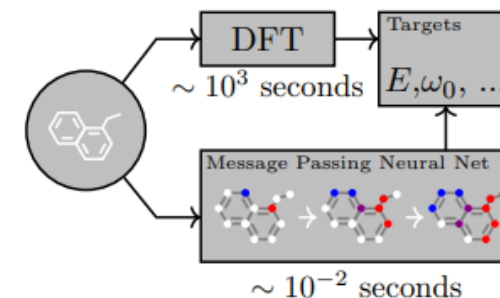
- Model requires a “kernel” and not “features”
- Only needs a single feature
- Kernels must be relevant to target property

Random Forest Regression



- Can use with 10^2 - 10^3 features
- Features work best if they compose simple rules
- Automatic feature selection
- Each feature must have same number of features

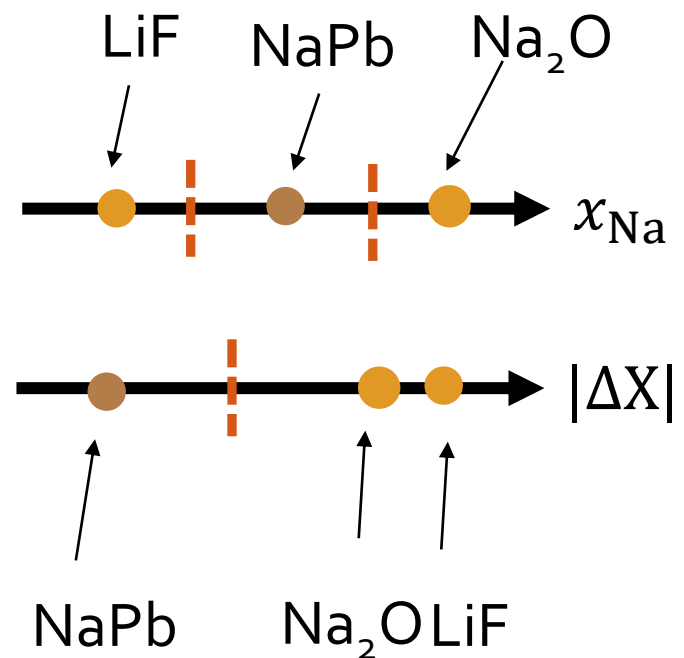
Neural Networks



- Automatically learn representations from data
 - Variable number of features, input shapes, ...
- if you can design the architecture

Take-away Points

- Representations are the key for encoding chemistry/physics into machine learning
- Representations should be:
 - Complete
 - Compact
 - Descriptive
 - Simple
- There is no best representation for “all problems”



Learning goals for the following lectures

- Explain why representations for kernel ridge regression and neural networks are different
- Identify appropriate representations for inorganic and organic materials
- Design studies for testing different approaches that expose problems in models