



ALCF AI FOR SCIENCE TRAINING SERIES



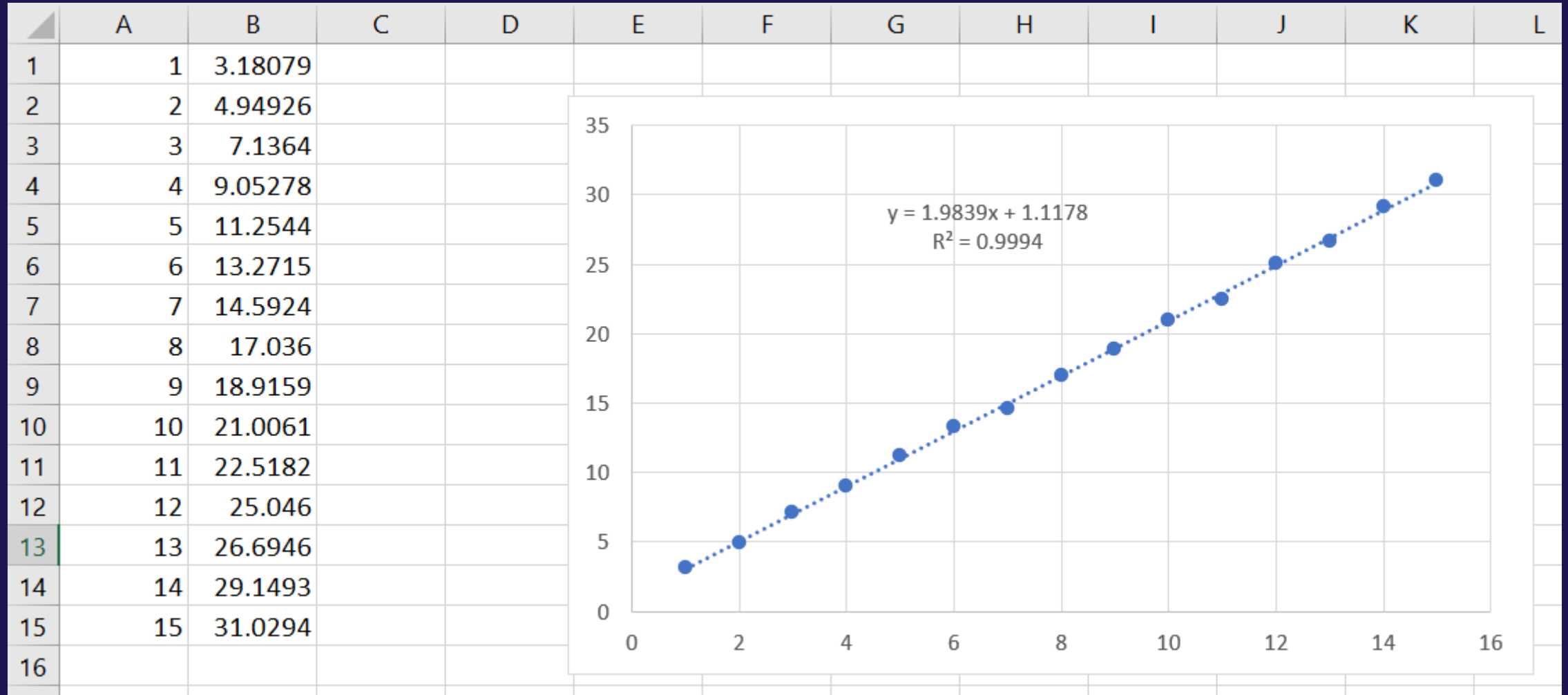
# Introduction To Machine Learning

Adapting scientific data, picking the right approach, and testing it well

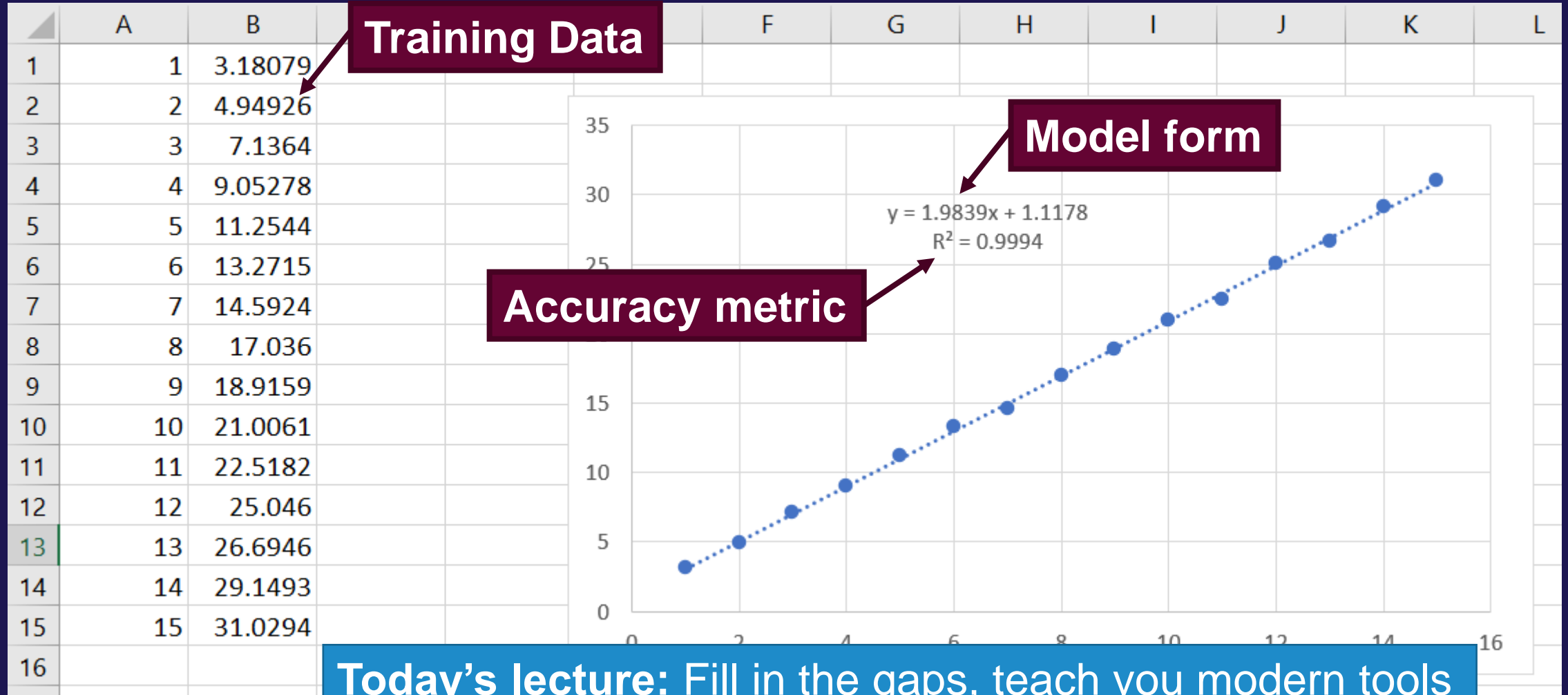
Logan Ward

Data Science and Learning Division, Argonne National Laboratory

# You've already done machine learning



# ... and already have the main ideas



# Is it different than conventional modeling? *Somewhat*

## “Semi-Empirical” Models

- Compose **equations**
- Fit unknown terms to data
- Evaluate to solve problem

Tersoff Potential  
12 parameters

Understand model by  
understanding *physics*

## Machine Learning

- Compose **model forms**
- Fit unknown terms to data
- Evaluate to solve problem

Gaussian Approximation Potential  
300 parameters

Understand model by  
understanding *algorithms*

**Many Similarities**

**Difference is  
in Scale**

**Bottom Line**

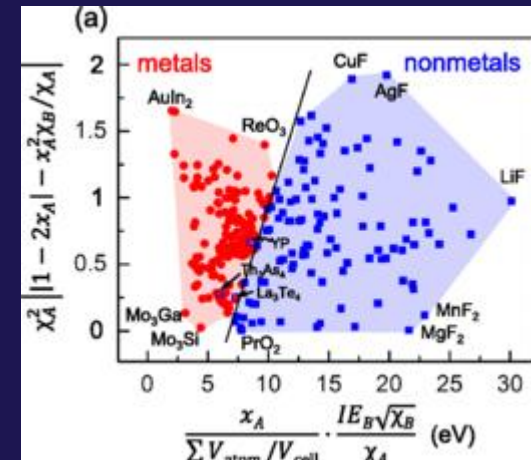
# Is it different than conventional empiricism? *Mostly*



$$M_d30 = 497 - 462(x_C + x_N) - 8.1x_{Mn}$$

Engineering is rife with empirical models,  
fit to data, used to understand physics

Modern AI techniques delegate more  
creativity to the algorithm,  
but more interpretation to the human



Ouyang et al. PRM. (2018), 083802

# Ok, why do we have an “AI4Sci” tutorial series?

(And 2 hours on “machine learning”)

**Key reason:** AI requires different skills

1. Understanding algorithms is key for success
2. Complex models must be used with care
3. We still do not understand all uses of AI
4. Tools are not commonly taught or used elsewhere

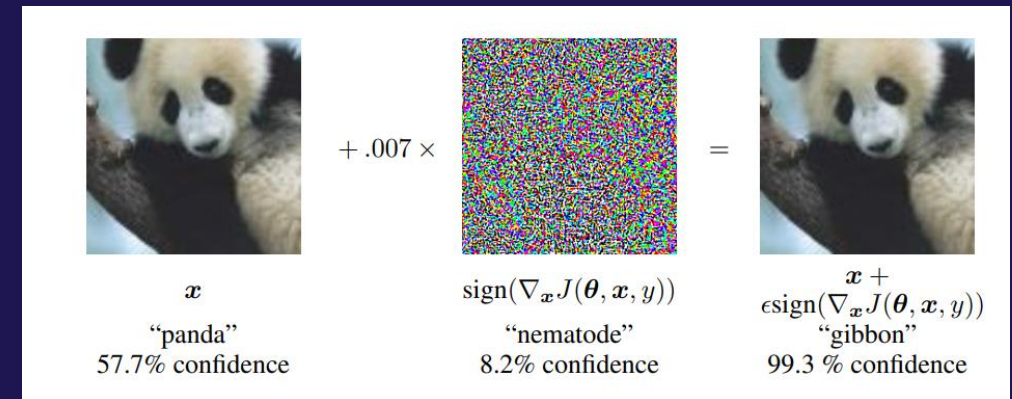


Figure: Goodfellow et al. ICLR (2015)

**Today's Goal:** Introduce “machine learning”

1. What are the key algorithms?
2. How do I use them for “science?”





# PART 1: THE ALGORITHMS!

7

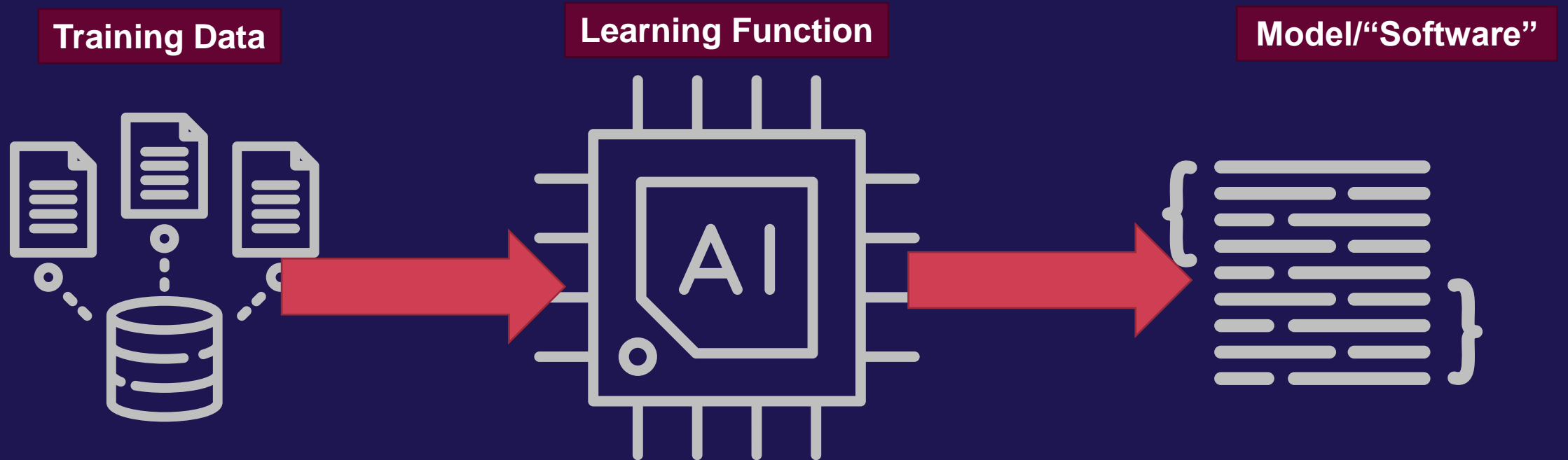


Argonne National Laboratory is a  
U.S. Department of Energy laboratory  
managed by UChicago Argonne, LLC.



# Anatomy of a machine learning algorithms

*All have roughly the same parts*



- Entry/Data Point/Instance  
Individual piece of data
- Labels  
measured about a record

- Objective Function  
*What marks "better"?*
- Optimization Routine  
*How do I make it better?*

- Architecture  
*What form does the model take?*
- Weights  
Values of adjustable parameters

Let's go through a few key examples



# Anatomy of a machine learning algorithms

## Training Data

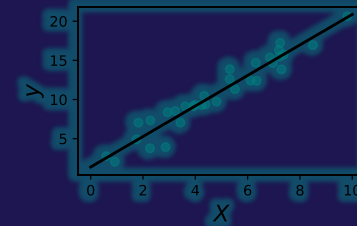
## Learning Function

## Model/“Software”

### Supervised Learning

Entries:  $X$   
Labels:  $y$

Find  $f$  where  
 $y \approx f(X)$



Linear models  
Lookup tables  
Decision Trees

### Clustering

Entries:  $X$

Discover similar subsets  
with  $X$

k Means  
Agglomerative Clustering

### Dimensionality Reduction

Entries:  $X$

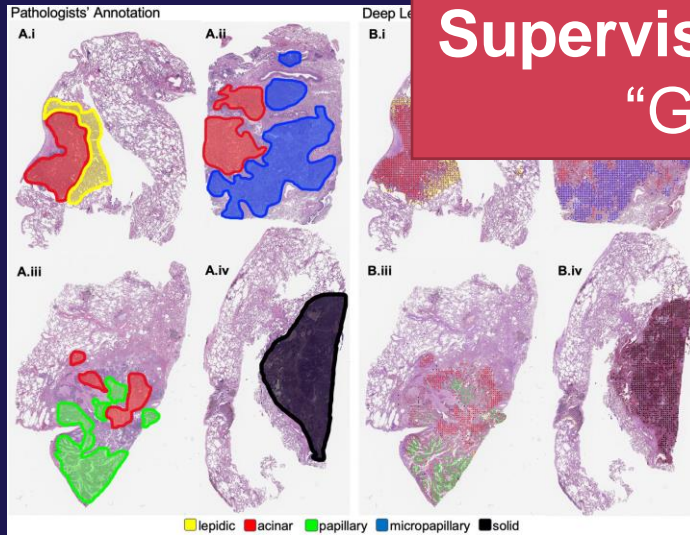
Learn  $h$  where  
 $X \approx h(X)$

Principle Component Analysis  
Autoencoders

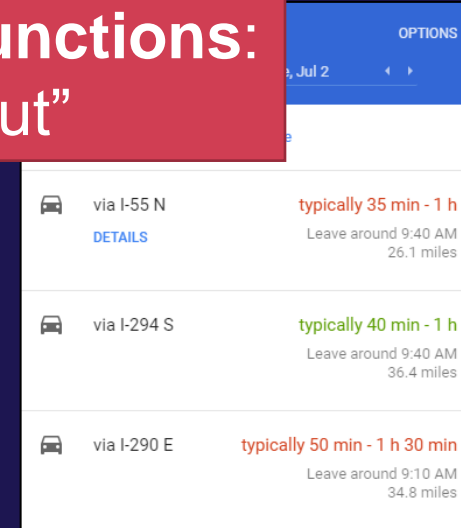
There are many types “learning” you can do. This list is not complete!

# Supervised Learning: The ML you've done before

Supervised Learning models functions:  
“Given inputs, predict output”

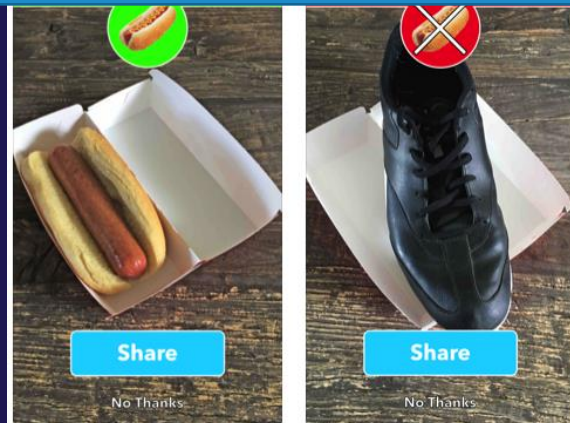


Wei et al. *Sci Rep.* (2019), 3358



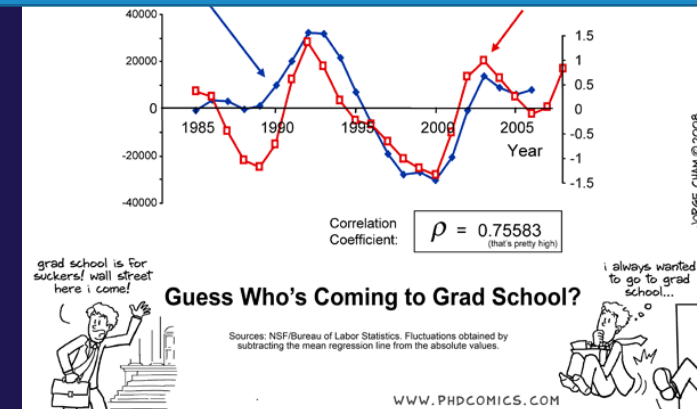
Google Maps

Classification: Outputs are *discrete*



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

Regression: Outputs are *continuous*



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

# An old friend: Simple Linear Regression

... but let's make it sound modern

## Model Architecture

$$f(x; m, b) = mx + b$$

**Training Data:** Inputs ( $x_i$ ) and outputs ( $y_i$ )

**Goal:** Determine  $m$  and  $b$  that minimize

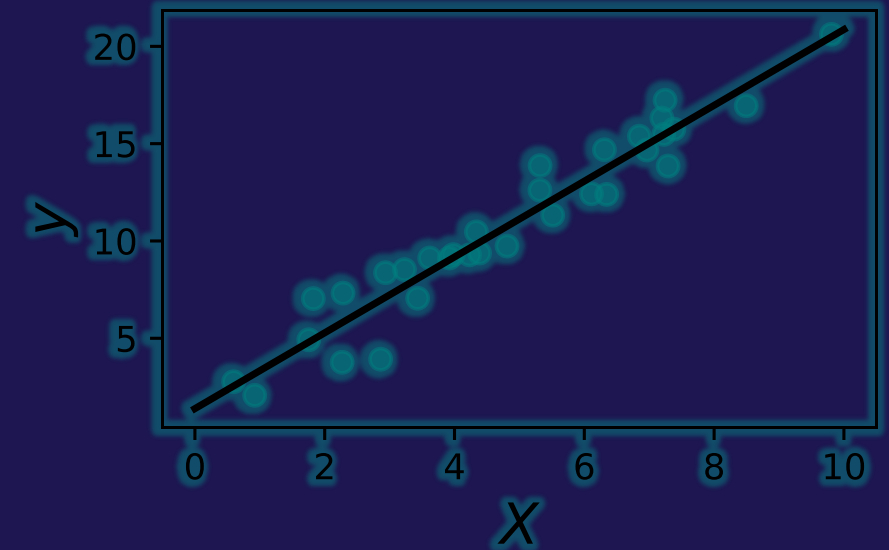
## Loss Function

$$\sum_i (f(x_i; m, b) - y_i)^2$$

by computing

$$m = \text{Cov}[x, y] / \text{Var}[x]$$
$$b = \bar{y} - m\bar{x}$$

## Optimizer



# Simple Logistic Regression

A version of Linear Regression suitable for classification

## Model Architecture

$$f(x; m, b) = \frac{1}{1 + e^{-(mx+b)}}$$

**Training Data:** Inputs ( $x_i$ ) and outputs ( $y_i$ )

**Goal:** Determine  $m$  and  $b$  that minimize

## Loss Function

$$L(m, b) = \sum_i y_i \ln(f(x_i)) + (1 - y_i) \ln(1 - f(x_i))$$

by computing

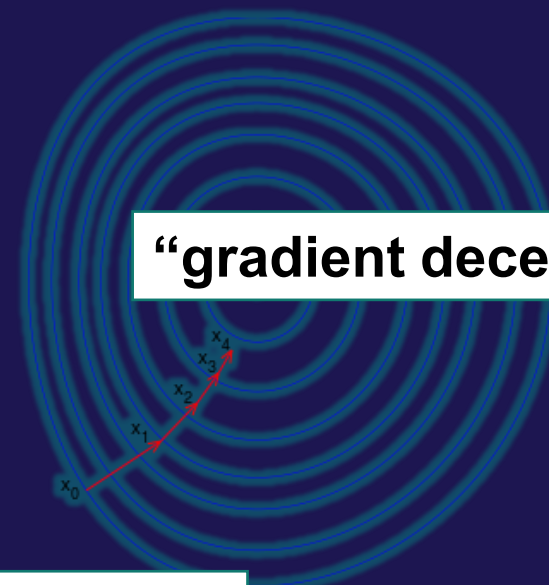
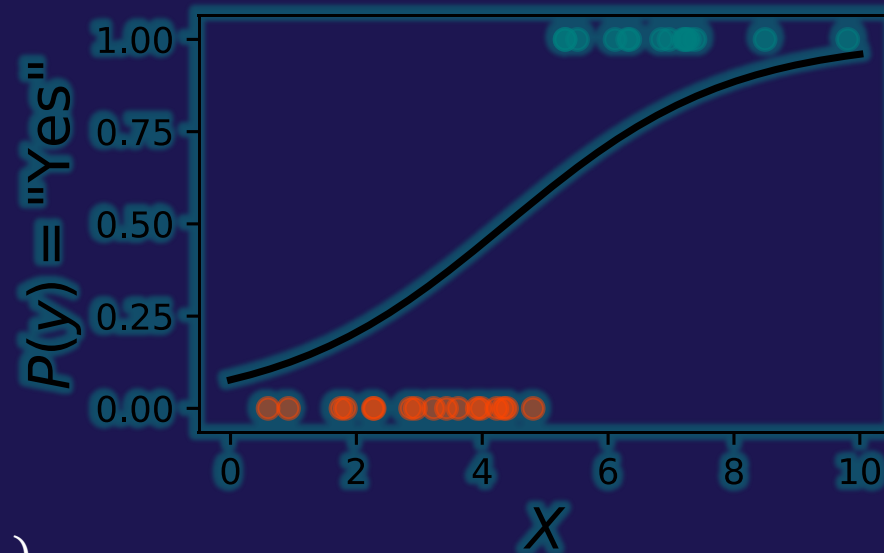
$$x_0 = (1, 0) \\ x_{n+1} = x_n + \gamma \nabla L(m, b)$$

“log loss”\*

## Optimizer

“gradient decent”\*

**Architecture + Loss + Optimizer = ML Algorithm  
For Regression *and* Classification**

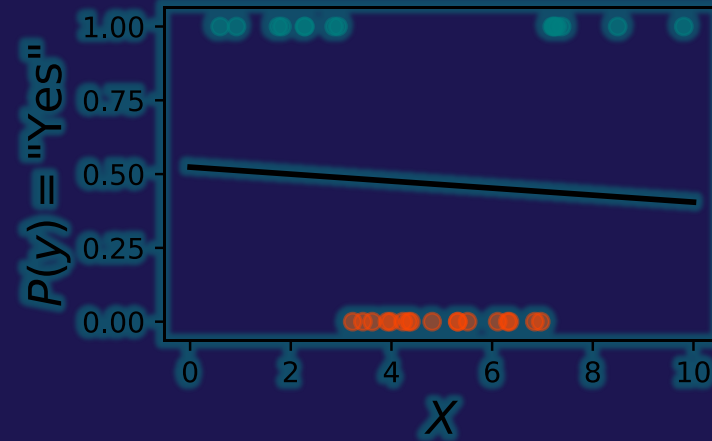
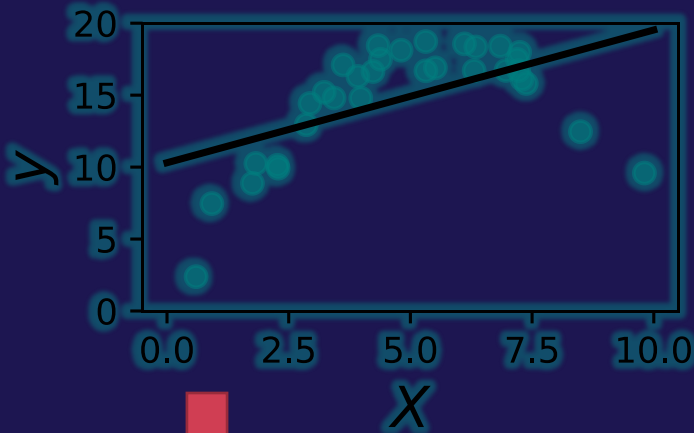


# Linear Models Are Not Sufficient

Otherwise, this would be a very short lecture

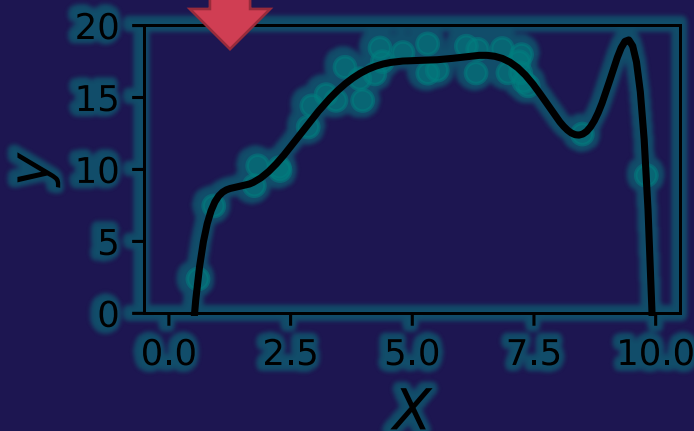
**Why Not?** Model complexity is limited

**“underfit”\***



*... and adding complexity comes with risks*

**“overfit”\***



**Key Questions for Supervised ML:**

1. How to add more complexity?
2. How to know when “overfit”?

# There is a ZOO of ML Algorithms

Variations of Bayes' Theorem

Linear regression, etc.

*k*-Nearest neighbors

*Many* kinds of decision trees

Trees vary in...

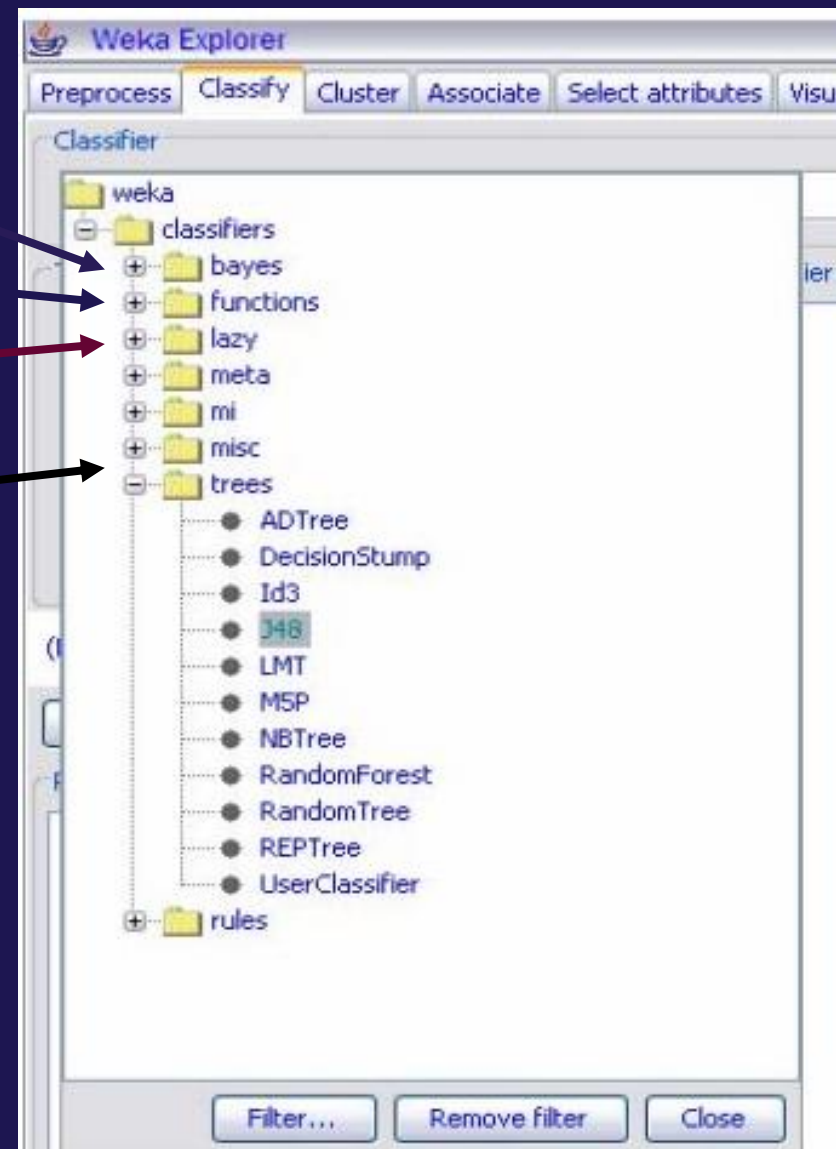
- whether to process inputs
- ways picking splits
- what are on the “leaves”
- how to prune after training

$$x_0 < 4$$

$$y = 2$$

$$y = 6$$

There is no “algorithm to rule them all”





# Linear Models are quite important

**Key Concept:** “Penalized” regression methods

$$\|y - \alpha X\|_2 + \lambda \|\alpha\|_p \leftarrow \text{“Regularizer”}$$

**Subtlety:** More than one type of “regularization” loss

where  $\|\alpha\|_p$  is the “ $L^p$  norm”:  $\|\alpha\|_p = (|\alpha_0|^p + \dots + |\alpha_n|^p)^{\frac{1}{p}}$

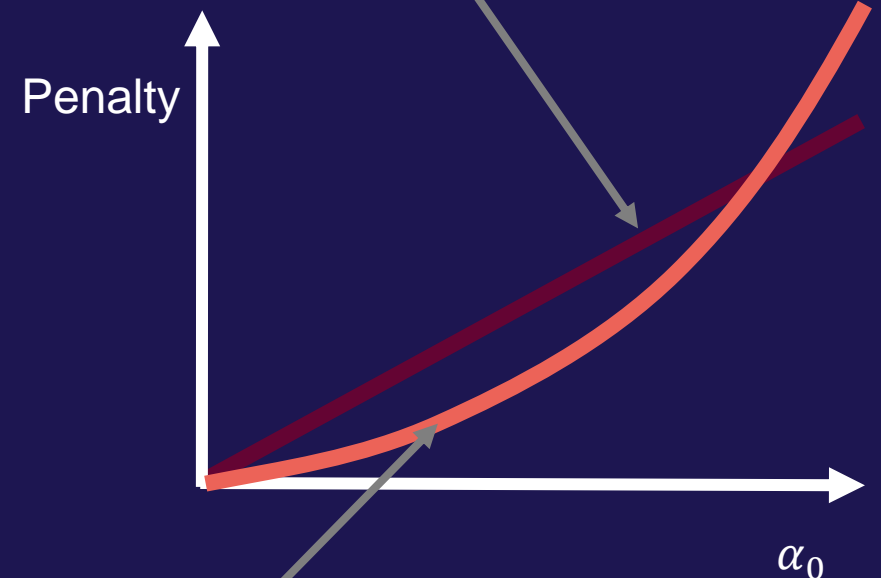
**Advantages:**

1. Fast to train!  $L_2$  has an analytic solution
2. Interpretable, especially for  $L_1$

**Disadvantage:** Limited complexity

**$L^1$ -norm:** “sparse linear regression”

- Moderate penalties for small parameters
- Sparse models (many  $\alpha = 0$ )



**$L^2$ -norm:** “ridge regression”

- Small penalties for small parameters
- Many parameters (No  $\alpha = 0$ )

# Key algorithm: Decision Trees

**Model form:** Series of rules

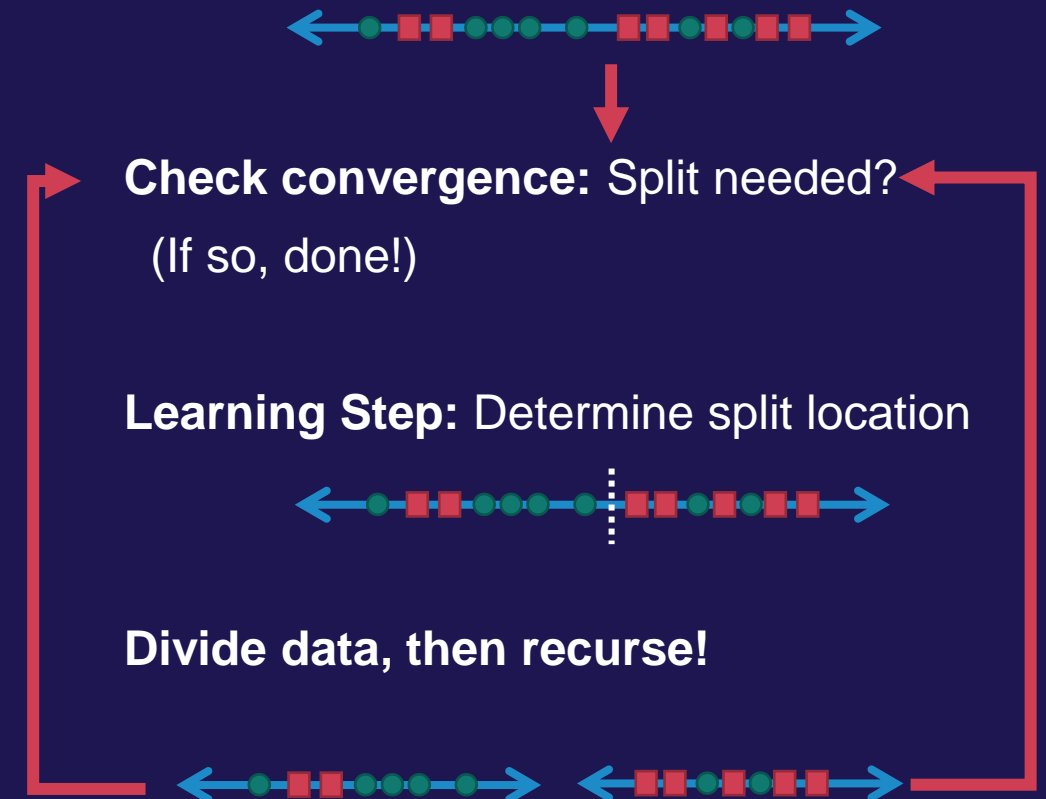
**Learning objective:** Find rules that maximize performance

**Learning algorithm:** Recursive rule selection

*(Many variations of these components!)*

**Why might I choose decision trees:**

- Fast to train ( $O(N \log N)$ )
- Mostly interpretable by humans

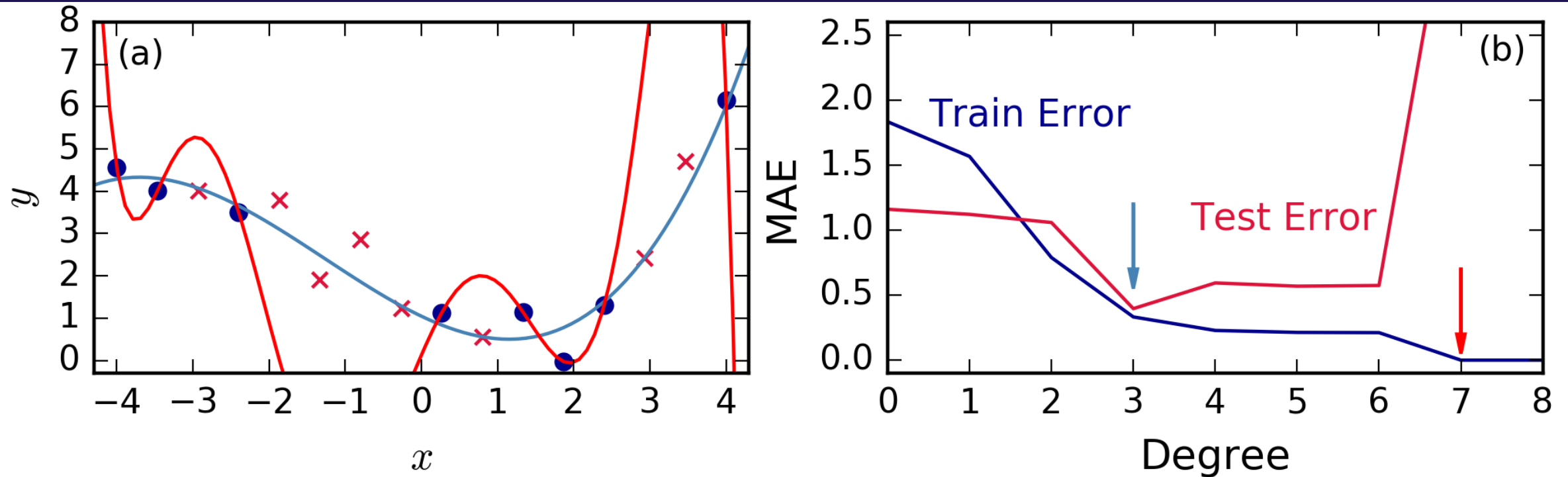


# MODEL SELECTION: THE KEY TASK IN SUPERVISED ML

17

# A Reprise: Overfitting and Complexity

Training accuracy vs “generalizability”

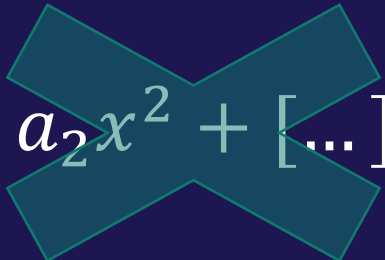


# How do We Adjust COMPLEXITY?

We'll just talk linear models for now, but these are general ideas

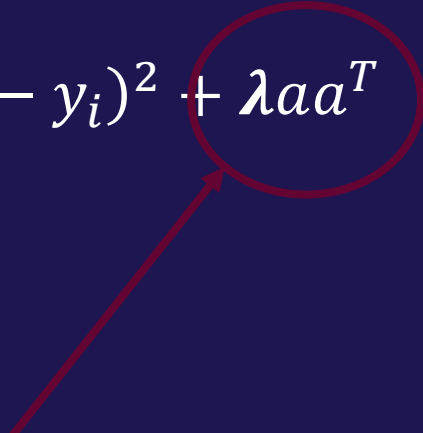
## LIMIT COMPLEXITY

*Reduce degrees of freedom*

$$f(x) = a_0 + a_1x + a_2x^2 + [\dots]$$


## PENALIZE COMPLEXITY

*Add “complexity” to loss function*

$$\sum_i (f(x_i) - y_i)^2 + \lambda a a^T$$


Larger coefficients = bigger penalty

# Two MAJOR Concepts

Don't leave this room without them!

## CROSS VALIDATION

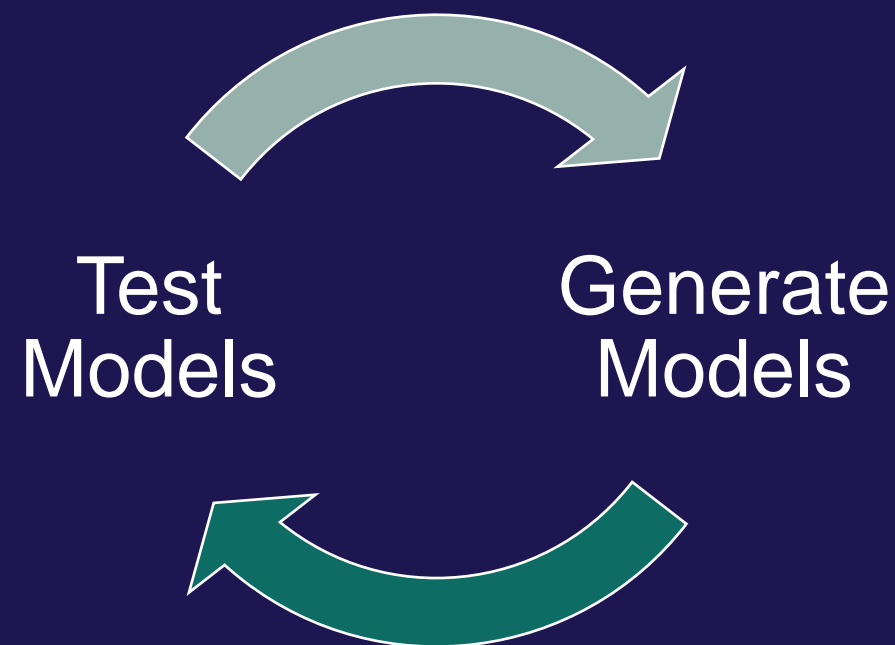
*Given available data,  
test whether model is predictive*

### Basic Techniques:

- *Shuffle Split*
  1. Pick 10% as test set
  2. Train on remaining 90%
  3. Test model on test set
- *Leave-one-out*
  1. Pick one entry
  2. Train on remaining entries
  3. Test model on held-out entry
  4. Repeat using each entry

## HYPERPARAMETER OPTIMIZATION

*Adjust model settings  
to maximize CV performance*





# Summary for Part 1

**What is machine learning?**

Algorithms that generate software from desired outcomes

**What makes kinds of learning different?**

Desired inputs and outputs, learning algorithms

**Example:** “Supervised Learning” finds  $y = f(x)$

**Why are there so many learning algorithms?**

Each have their own advantages and drawbacks

**Example:** Kernel methods train time is  $O(N^3)$

**Why must I optimize hyperparameters?**

Adjust been “under” and “over” fitting

**Example:** Number of polynomial terms

# PRACTICAL EXERCISE 1: BUILDING MODELS WITH SKLEARN!

22

# PART 2: USING IT FOR SCIENCE (MOSTLY MATERIALS SCIENCE, BECAUSE I'M A MATERIALS ENGINEER)

23



# Not all data are vectors

And that's OK!

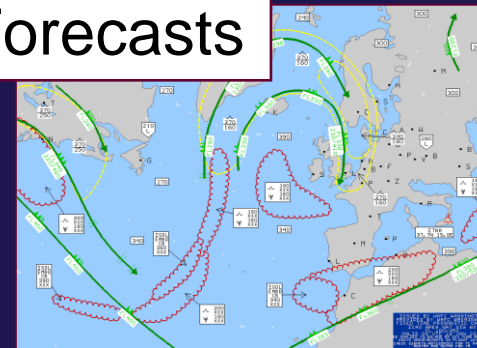
Molecules



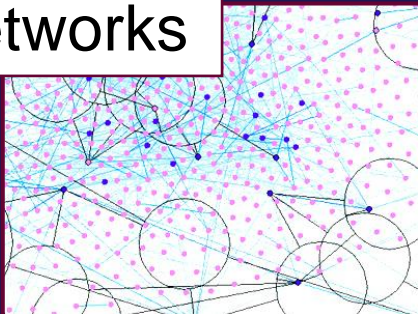
Images



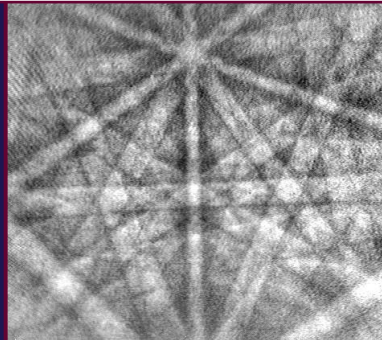
Weather  
Forecasts



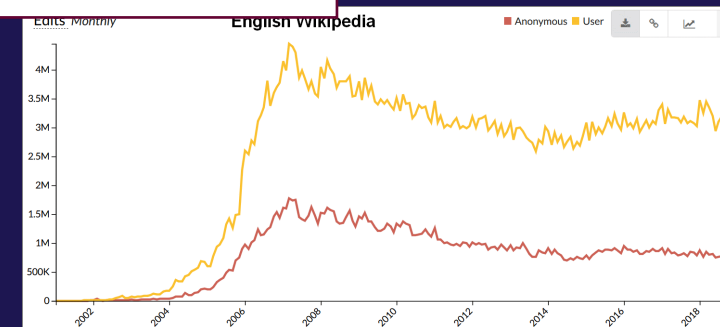
Social  
Networks



EBSD Patterns



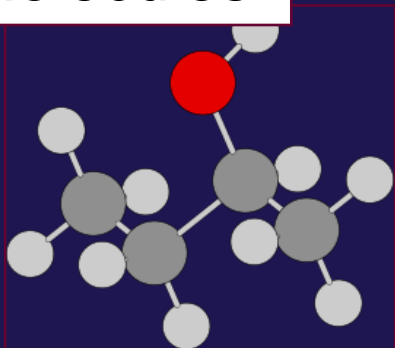
Timeseries



# Not all data are vectors

And that's OK!

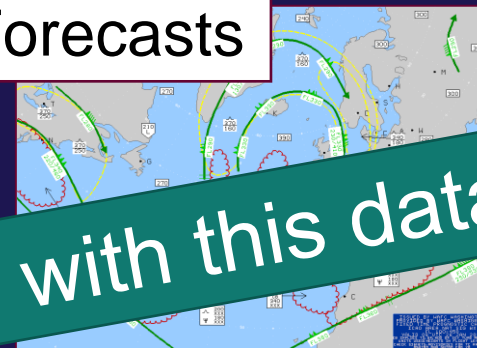
Molecules



Images

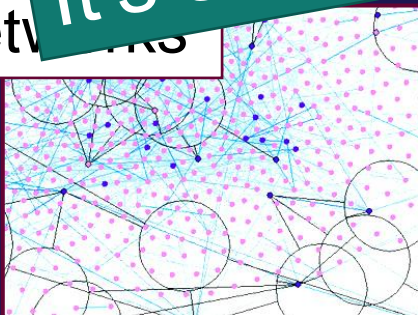


Weather  
Forecasts

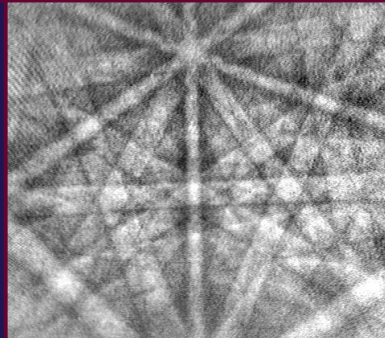


It's still possible to use machine learning with this data

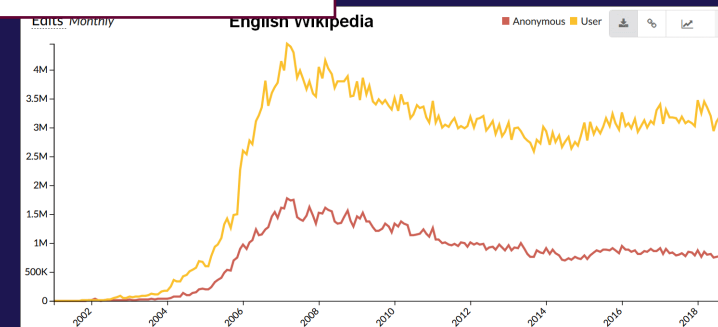
Social  
Networks



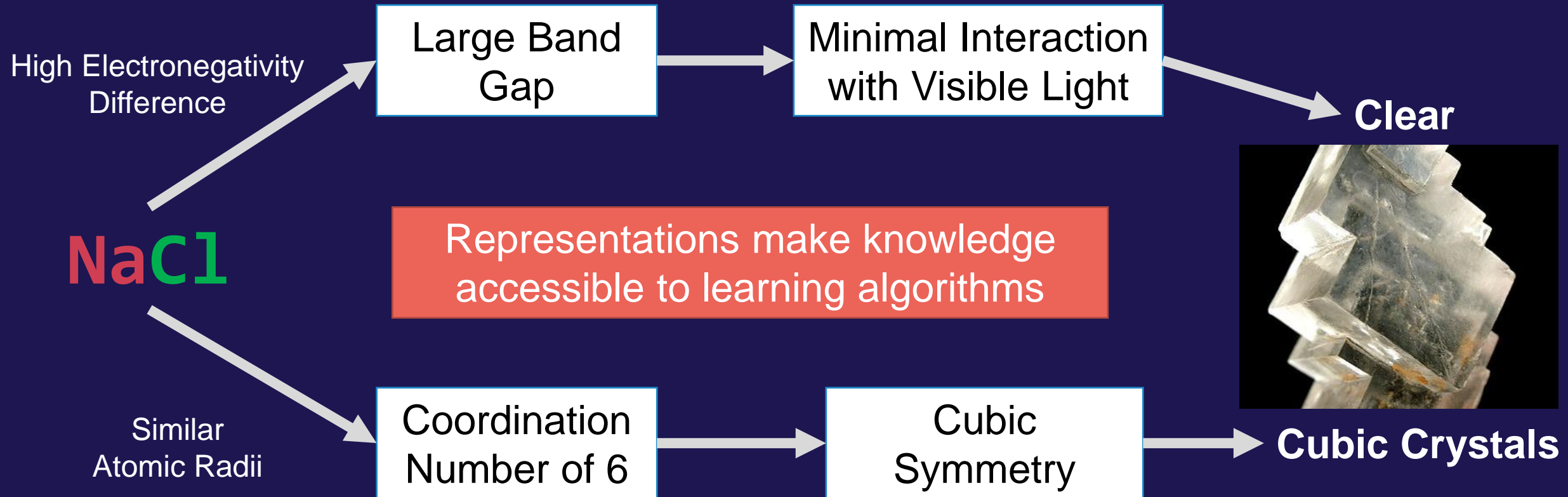
EBSD Patterns



Timeseries



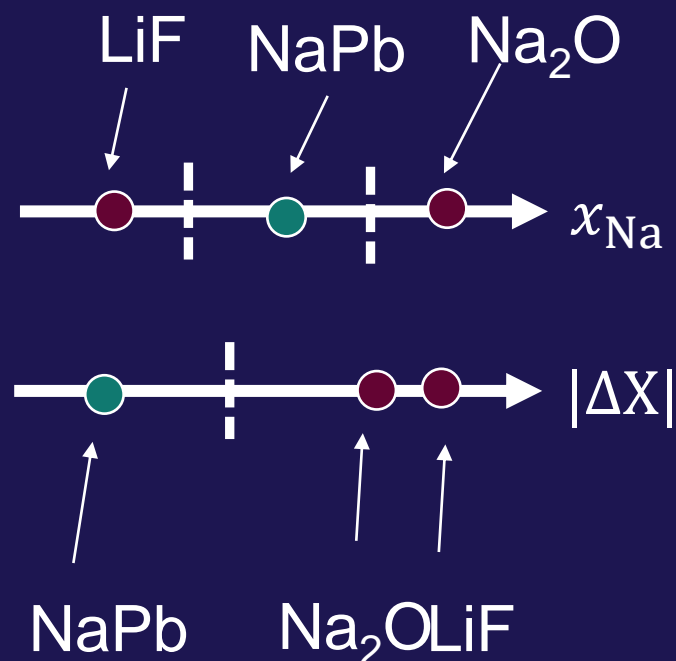
# What do NaCl crystals look like?





# How to translate chemistry/physics/... to a computer?

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



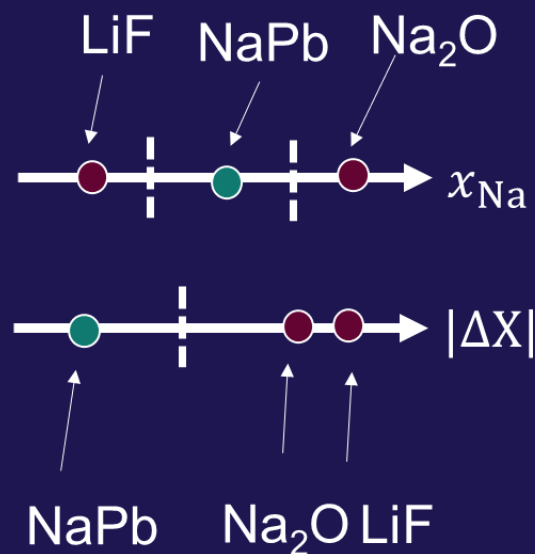
**Why do we need them?**

- Machine learning tools take tensors

**What makes a good one?**

- Easy to learn generalizable rules

# What makes a good representation?



**Complete**

✓ Yes, no overlapping points

✓ Yes, no overlapping

**Compact**

✓ Yes, single feature!

✓ Yes, single feature!

**Descriptive**

✗ No, rules are complicated

✓ Yes, matches textbook

**Simple**

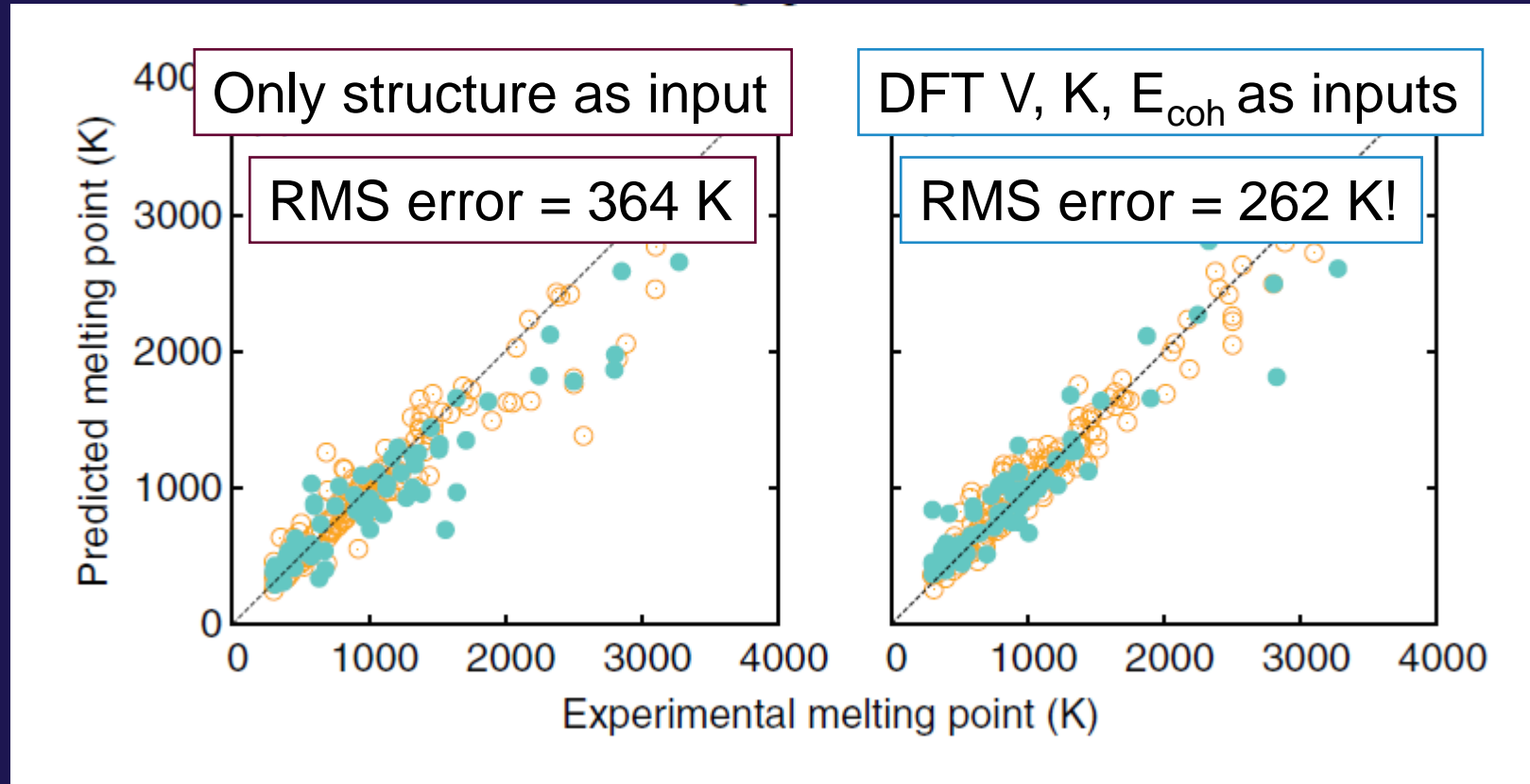
✓ Yes, no compute

✓ Yes, minimal compute

Ref: Faber et al. Int. J. Quantum Chem. (2015)

# 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

Baking basis into the model

Typical for  
“conventional” ML

## Concept

## Examples

## Advantages

## Challenges

## Descriptive

Make features that distinguish examples

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

Maximum expressivity

Requires more data  
(learn features *and* model)

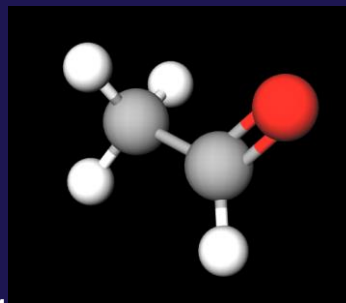
Typical for  
“deep learning”

# Molecular Descriptors

## Discriminative

## Line-notations for structure:

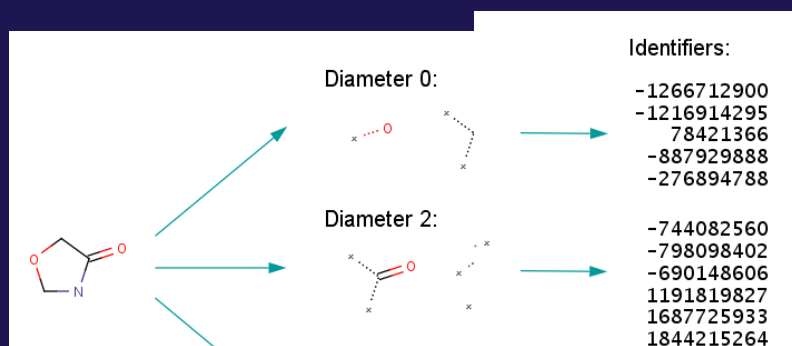
- SMILES (ex: "CC=O")
- InChI (ex: InChI=1S/C2H4O/c1-2-3/h2H,1H3)



## Descriptive

## Extremely Well-Studied

## *Fingerprints:*



Identifier list representation:

~~-1266712900~~   ~~-1216914295~~   ~~78421366~~   ~~-887929888~~   ~~-276894788~~   ~~-744082560~~   ~~-798098402~~   ~~-690148606~~   ~~1191819827~~  
~~1687725933~~   ~~1844215264~~   ~~-252457408~~   ~~132019747~~   ~~-2036474688~~   ~~-1979958858~~   ~~-1104704513~~

Fixed-length binary representation:

A horizontal sequence of 100 bits: 010000000001000001100001000110000000000101000000000000000000000000000000010010100100000000001000000000000. Above the sequence, teal arrows indicate dependencies: from bit 1 to 2, 3 to 4, 5 to 6, 7 to 8, 9 to 10, 11 to 12, 13 to 14, 15 to 16, 17 to 18, 19 to 20, 21 to 22, 23 to 24, 25 to 26, 27 to 28, 29 to 30, 31 to 32, 33 to 34, 35 to 36, 37 to 38, 39 to 40, 41 to 42, 43 to 44, 45 to 46, 47 to 48, 49 to 50, 51 to 52, 53 to 54, 55 to 56, 57 to 58, 59 to 60, 61 to 62, 63 to 64, 65 to 66, 67 to 68, 69 to 70, 71 to 72, 73 to 74, 75 to 76, 77 to 78, 79 to 80, 81 to 82, 83 to 84, 85 to 86, 87 to 88, 89 to 90, 91 to 92, 93 to 94, 95 to 96, 97 to 98, 99 to 100.

### Bit collisions



## Handbook of Molecular Descriptors

Author(s): Prof. Dr. Roberto Todeschini, Dr. Viviana Consonni

First published: 22 September 2000

Print ISBN: 9783527299133 | Online ISBN: 9783527613106 | DOI: 10.1002/9783527613106

Copyright © 2000 WILEY-VCH Verlag GmbH

Book Series: Methods and Principles in Medicinal Chemistry

 Free Access

## Bibliography (Pages: 524-667)

[First Page](#) | [PDF](#) | [References](#) | [Request permissions](#)

## Many types of descriptors:

1. Constitutional (ex: "how many Ns?")
2. Structural (ex: Solvent-Accessible Surface Area)
3. Quantum-chemical (ex: partial charges)

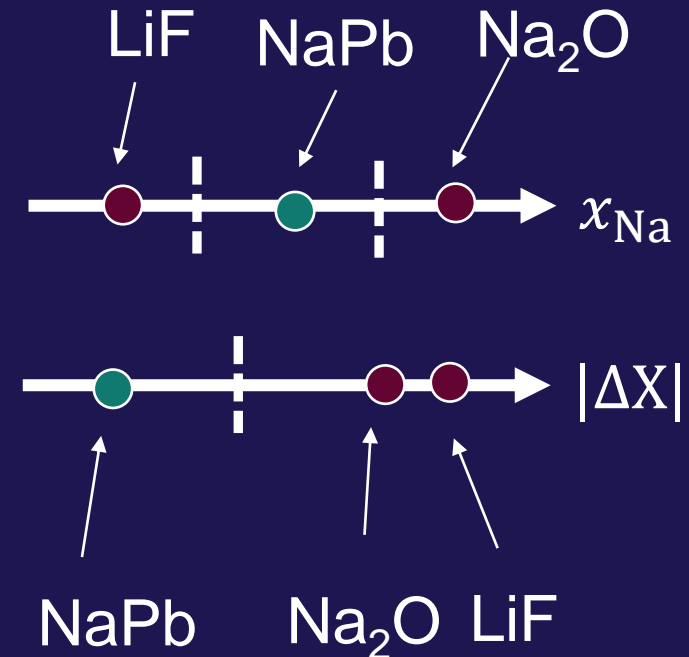
# Conclusion: How to translate science to ML?

**Key concept:** Express your data in a way that captures knowledge

Many ways to achieve “Guiding Principles:”

- **Complete:** Separate different entries
- **Compact:** Minimal complexity
- **Descriptive:** Maximum relevance
- **Simple:** Fast to gather

Work here can yield significant benefits,  
and is a great way to engage domain experts





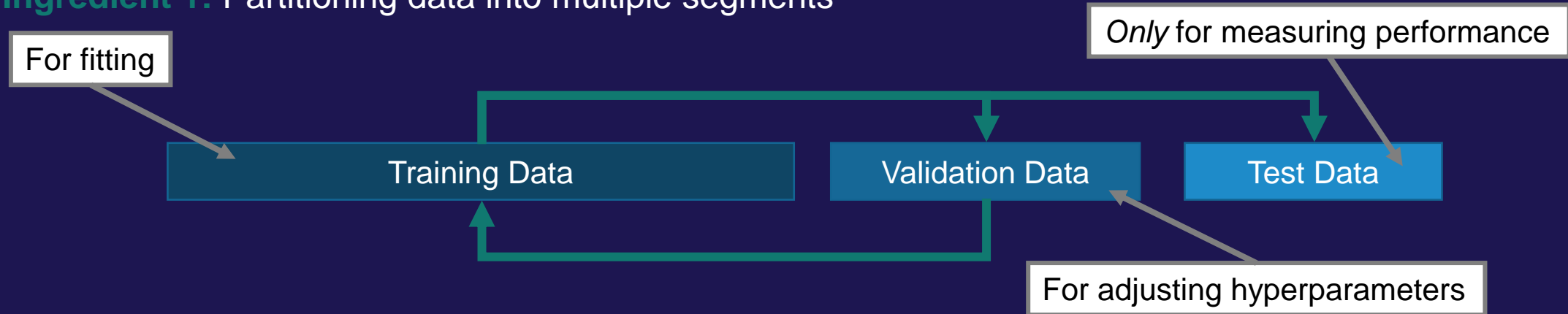
# PRACTICAL EXERCISE 2: MACHINE LEARNING WITH CHEMISTRY DATA

# PART 2B: VALIDATE LIKE A PRO

34

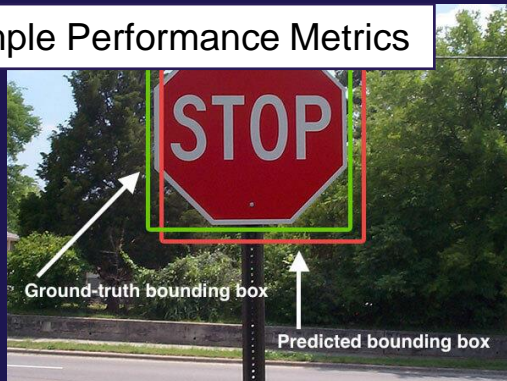
# General approach to testing: “cross validation”

**Ingredient 1:** Partitioning data into multiple segments

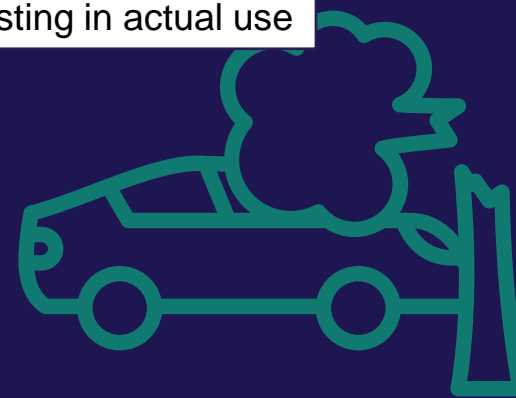


**Ingredient 2:** Choose how you measure performance

✓ Simple Performance Metrics



🏎️ Testing in actual use



**My advice:** Think carefully about how to mimic your application

Image: [pyimagesearch](https://www.pyimagesearch.com/)



# A standard approach: k-fold cross validation

**Step 1:** Split dataset randomly into  $k$  equally-sized chunks



**Step 2:** Hold out 1 chunk as test data, use remainder as train and validation data



**Step 3:** Measure performance on remaining chunk

**Step 4:** Repeat from Step 2 with a different choice from the  $k$  chunks, until all chunks used



**Step 5:** Report performance as average over all iterations

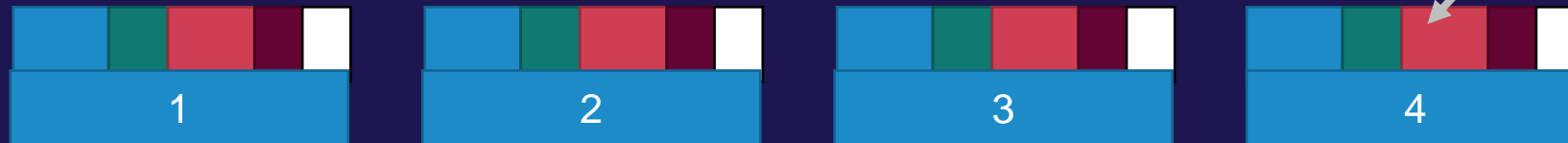
# A wrinkle: Data has hidden structure

**Common case:** Some data are more alike than others



**Question:** How would the model perform on new types of data?

*K-fold validation does not answer this question!*



Always predicting within same group!

Simulate predicting outside of groups.

*Alternative: Use “hidden” groups as training and test sets*



**Aside:** I avoid using the term “extrapolation.” Ask me later...



# Many examples of “data with hidden structure”

## Time Series

### ✗ Predicting past *and* future

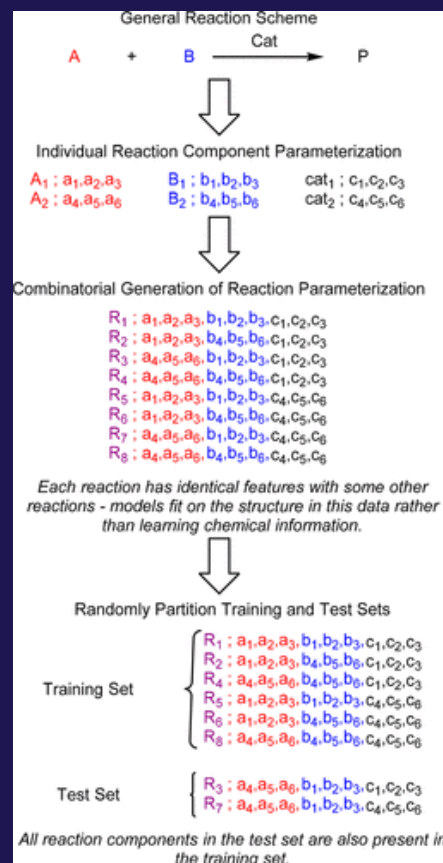
Split 1:	Test set	Training set			
Split 2:	Training set	Test set	Training set		
Split 3:	Training set	Test set	Training set		
Split 4:	Training set	Test set	Training set		
Split 5:	Training set	Test set	Training set		
	Time 1	Time 2	Time 3	Time 4	Time 5

### ✓ Predict only the future

Split 1:	Training set	Test set			
Split 2:	Training set	Test set			
Split 3:	Training set	Test set			
Split 4:	Training set	Test set			
	Time 1	Time 2	Time 3	Time 4	Time 5

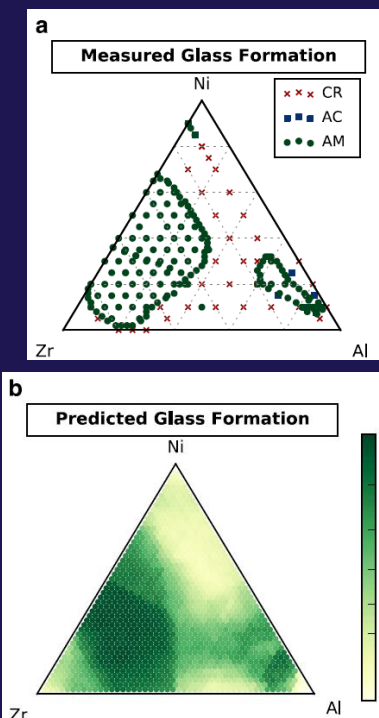
Split 1:	Training set	Test set			
Split 2:	Training set	Test set			
Split 3:	Training set	Test set			
Split 4:	Training set	Test set			
	Time 1	Time 2	Time 3	Time 4	Time 5

## Combinatorial Chemistry



## Metallic Alloys

### ✓ Excluding an alloy system



Source: [CrossValidated](#)

Ref: [Zahrt et al. ACS Combi. \(2020\)](#)

Ref: [Ward et al. npj Comp Mat. \(2016\)](#)

# There are many ways to measure performance

		Predicted condition	
		Positive (PP)	Negative (PN)
Actual condition	Total population = P + N		
	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection

Is “No missed diagnoses”  
your goal?

Or is it “No false convictions”?

Table from [Wikipedia](#)

# These only scratch the service for classification

Sources: [15][16][17][18][19][20][21][22] view · talk · edit

		Predicted condition			
Total population = P + N		Positive (PP)	Negative (PN)	Informedness, bookmaker informedness (BM) = TPR + TNR - 1	Prevalence threshold (PT) = $\frac{\sqrt{TPR \times FPR} - FPR}{TPR - FPR}$
Actual condition	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power = $\frac{TP}{P} = 1 - FNR$	False negative rate (FNR), miss rate = $\frac{FN}{P} = 1 - TPR$
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection	False positive rate (FPR), probability of false alarm, fall-out = $\frac{FP}{N} = 1 - TNR$	True negative rate (TNR), specificity (SPC), selectivity = $\frac{TN}{N} = 1 - FPR$
Prevalence = $\frac{P}{P + N}$	Positive predictive value (PPV), precision = $\frac{TP}{PP} = 1 - FDR$	False omission rate (FOR) = $\frac{FN}{PN} = 1 - NPV$	Positive likelihood ratio (LR+) = $\frac{TPR}{FPR}$	Negative likelihood ratio (LR-) = $\frac{FNR}{TNR}$	
Accuracy (ACC) = $\frac{TP + TN}{P + N}$	False discovery rate (FDR) = $\frac{FP}{PP} = 1 - PPV$	Negative predictive value (NPV) = $\frac{TN}{PN} = 1 - FOR$	Markedness (MK), deltaP ( $\Delta p$ ) = PPV + NPV - 1	Diagnostic odds ratio (DOR) = $\frac{LR+}{LR-}$	
Balanced accuracy (BA) = $\frac{TPR + TNR}{2}$	F <sub>1</sub> score = $\frac{2PPV \times TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$	Fowlkes–Mallows index (FM) = $\sqrt{PPV \times TPR}$	Matthews correlation coefficient (MCC) = $\frac{\sqrt{TPR \times TNR \times PPV \times NPV} - \sqrt{FNR \times FPR \times FOR \times FDR}}{\sqrt{FNR \times FPR \times FOR \times FDR}}$	Threat score (TS), critical success index (CSI), Jaccard index = $\frac{TP}{TP + FN + FP}$	

You should think about what success for your model looks like!

# Better regression validation: An example

**Great Example:** Predicting Density of States (phDOS) of materials ([Kong et al., 2021](#))

Best model for predicting  $C_v$  (a derived property)...

ML model	Setting		phDOS prediction				Calculated $C_v$		Calculated $\bar{\omega}$	
	Scaling	Loss	$R^2$	MAE	MSE	WD	MAE	MSE	MAE	MSE
Mat2Spec	SumNorm	WD	0.57	0.085	0.026	<b>21</b>	<b>1.32</b>	<b>10</b>	<b>10.6</b>	<b>348</b>
E3NN	SumNorm	KL	0.48	0.105	0.036	50	4.88	77	41.1	3718
GATGNN	SumNorm	KL	-1.05	0.177	0.057	215	22.4	756	205	51609
Mat2Spec	SumNorm	KL	0.62	<b>0.078</b>	<b>0.023</b>	24	1.96	11	17.1	625

... is not the one that predicts phDOS the best.

Choices of which performance metric to use will affect your model choice

# Study your problem, it will give you better models

**Ingredient 1:** Cross-validation that simulates how you want to use the model

✗ Predicting past *and* future

Split 1:	Test set	Training set			
Split 2:	Training set	Test set	Training set		
Split 3:	Training set		Test set	Training set	
Split 4:	Training set			Test set	Training set
Split 5:	Training set				Test set
	Time 1	Time 2	Time 3	Time 4	Time 5

✓ Predict only the future

Split 1:	Training set	Test set			
Split 2:		Training set	Test set		
Split 3:			Training set	Test set	
Split 4:				Training set	Test set
	Time 1	Time 2	Time 3	Time 4	Time 5

Source: [CrossValidated](#)

**Ingredient 2:** Evaluate the model how you actually measure performance

ML model	Setting		phDOS prediction				Calculated $C_v$		Calculated $\bar{\omega}$	
	Scaling	Loss	$R^2$	MAE	MSE	WD	MAE	MSE	MAE	MSE
Mat2Spec	SumNorm	WD	0.57	0.085	0.026	<b>21</b>	<b>1.32</b>	<b>10</b>	<b>10.6</b>	<b>348</b>
E3NN	SumNorm	KL	0.48	0.105	0.036	50	4.88	77	41.1	3718
GATGNN	SumNorm	KL	-1.05	0.177	0.057	215	22.4	756	205	51609
Mat2Spec	SumNorm	KL	0.62	<b>0.078</b>	<b>0.023</b>	24	1.96	11	17.1	625

Source: [Kong et al., 2021](#)

**My advice:** Think carefully about how to mimic your application



# WRAP UP: WHAT DID YOU LEARN?



Argonne National Laboratory is a  
U.S. Department of Energy laboratory  
managed by UChicago Argonne, LLC.





# Ok, why do we have an “AI4Sci” tutorial series?

(And 2 hours on “machine learning”)

**Key reason:** AI requires different skills

1. Understanding algorithms is key for success
2. Complex models must be used with care
3. We still do not understand all uses of AI
4. Tools are not commonly taught or used elsewhere

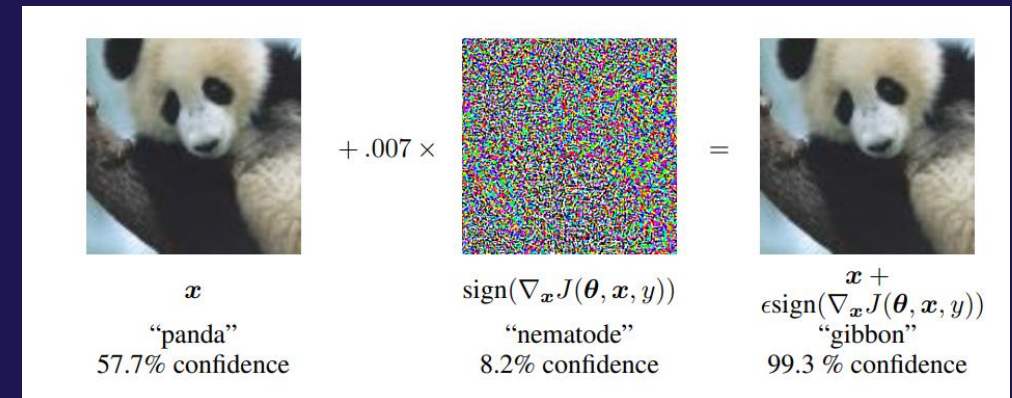


Figure: Goodfellow et al. ICLR (2015)

**Today's Goal:** Introduce “machine learning”

1. What are the key algorithms?
2. How do I use them for “science?”

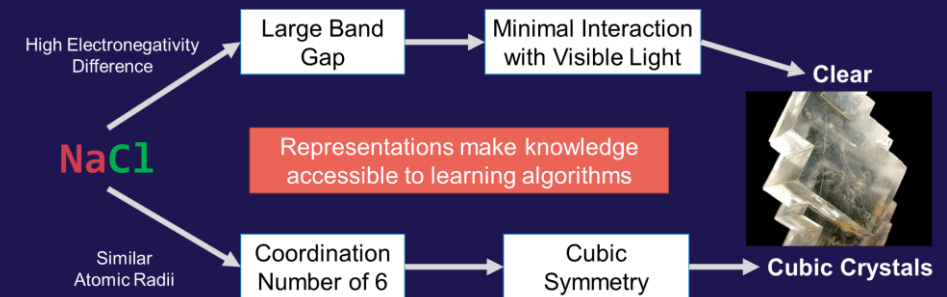


# Key takeaways!

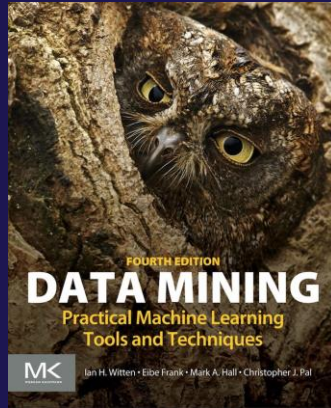
1. A view into the depths of how many ML algorithms there are, and knowledge of how to use a few *classic* ones

2. Understanding that data representation is one key to success

3. ... and that knowing how to validate is the other!

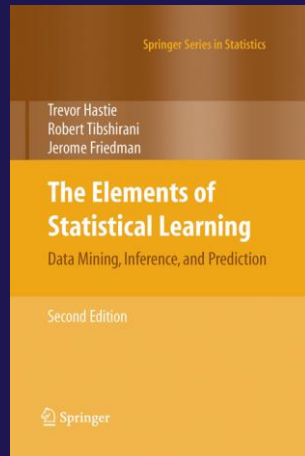


# Where to next? I like reading books...



I learned with a previous version of this book,  
and liked the level of mathematics vs application.  
- It teaches Weka, which is Java based

Recommended by Bethany Lusch (ALCF).  
Similar focus on application, but teaches scikit-learn



Deeper dive into the mathematics,  
very thorough coverage of key algorithms

