

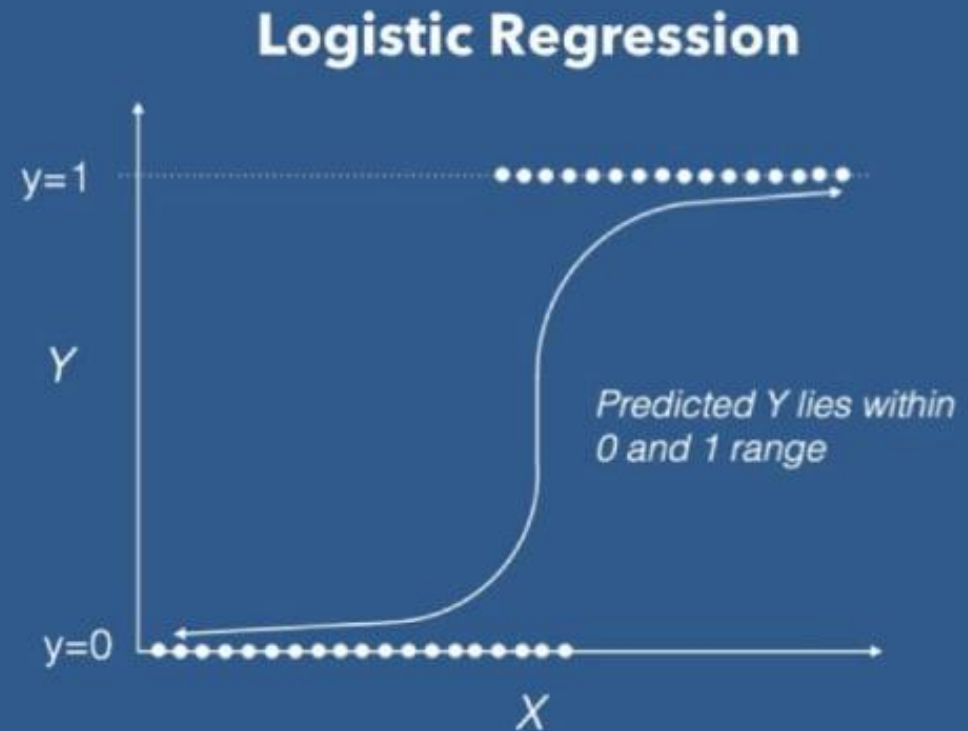
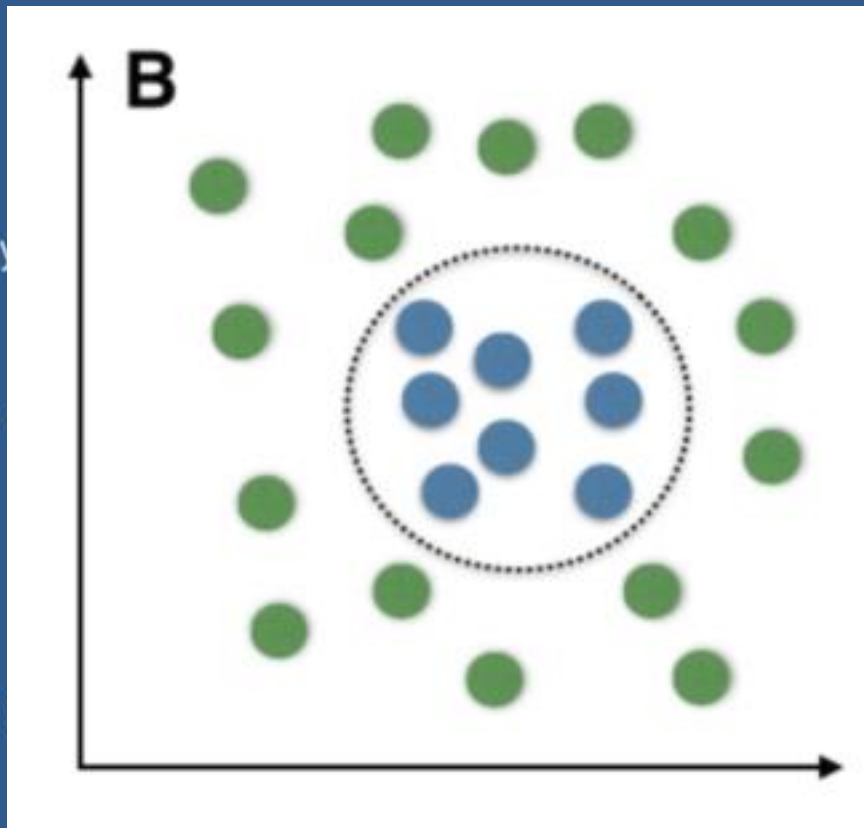
Summary of key points (so far)

- 1) General supervised learning:**
- 2) Success depends on:**
- 3) Validation is essential:**
- 4) Leverage your domain knowledge:**

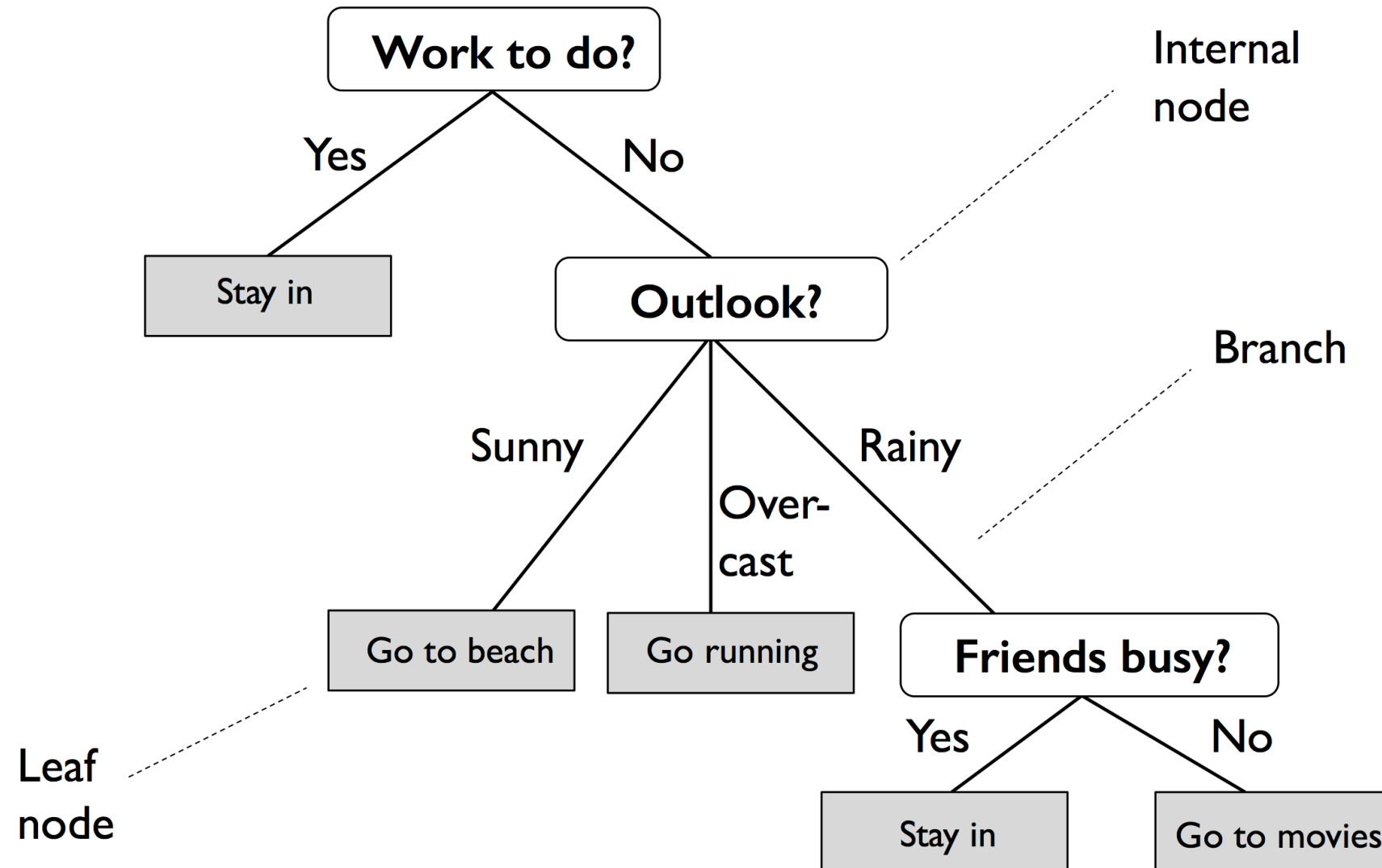
What about nonlinear classification?

For binary classification, y is no longer continuous, but binomial:

$$y = [1, 1, 1, -1, -1, 1, -1, -1, \dots]$$

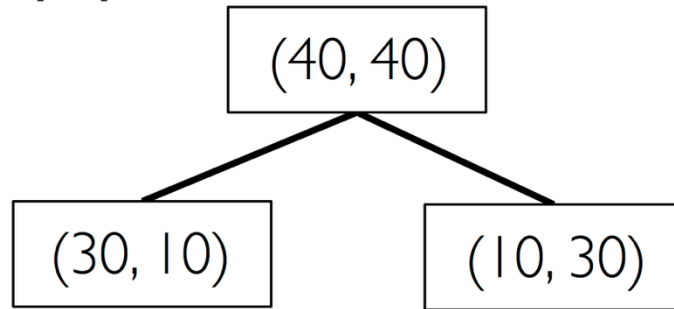


Brief intro to decision trees

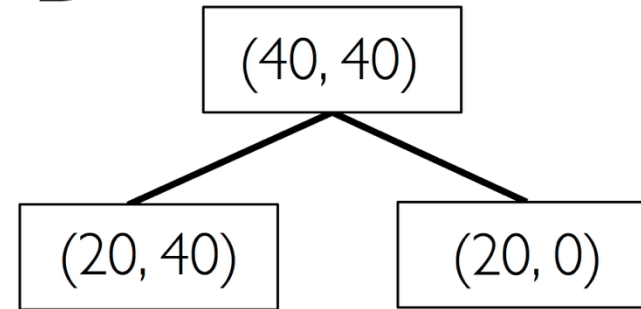


Which split is better?

A



B



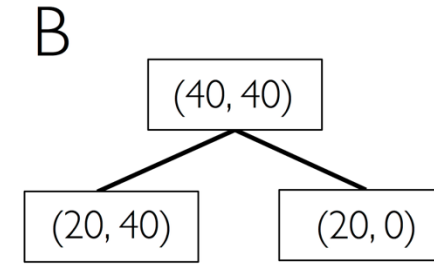
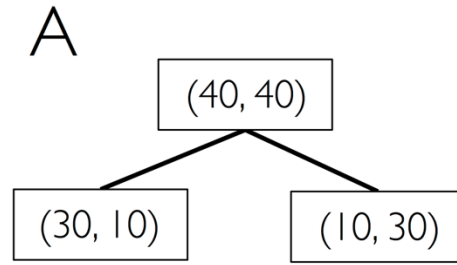
Brief intro to decision trees

Determine splits by **maximizing information gain (IG)**
minimizing weighted impurity, I

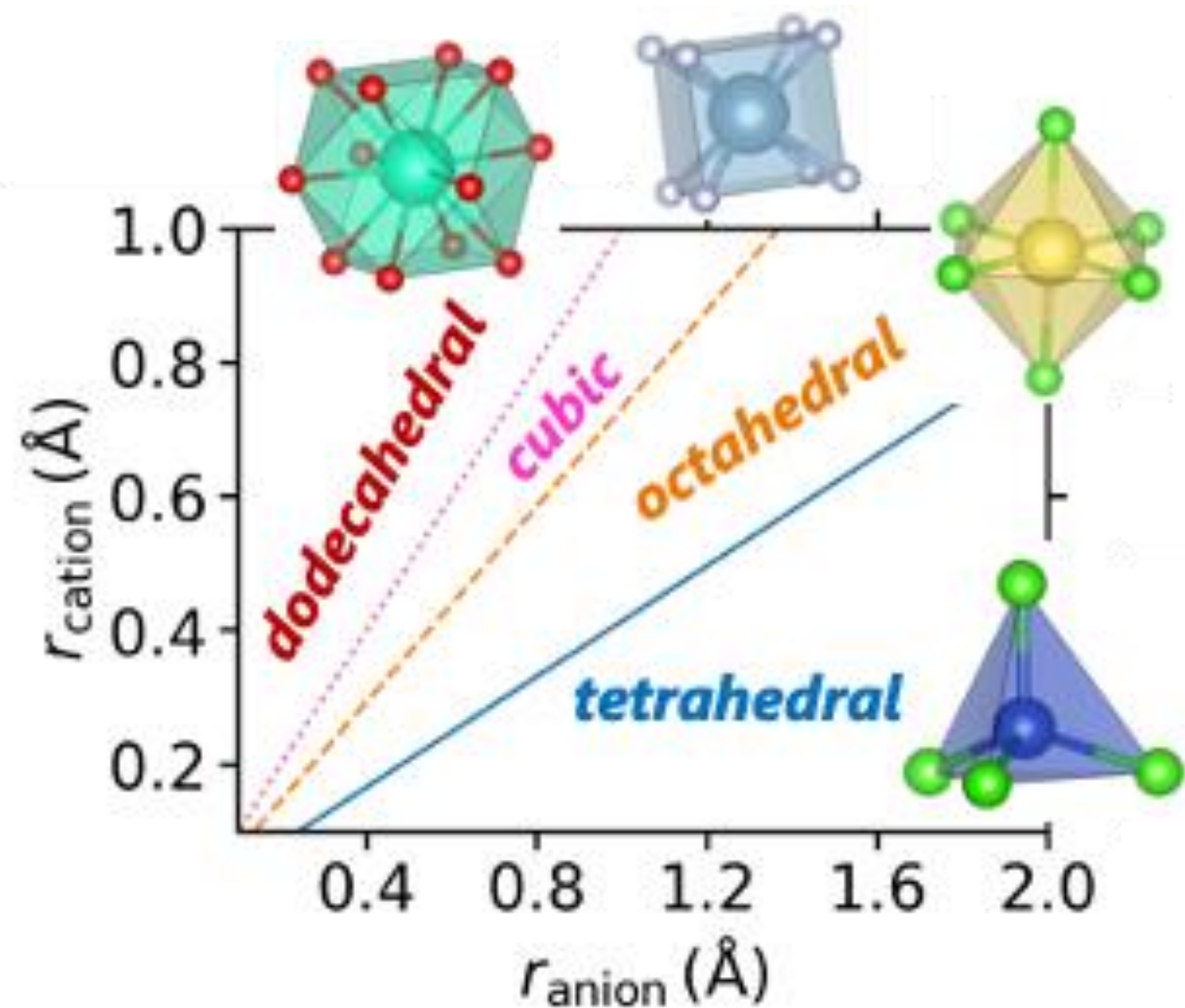
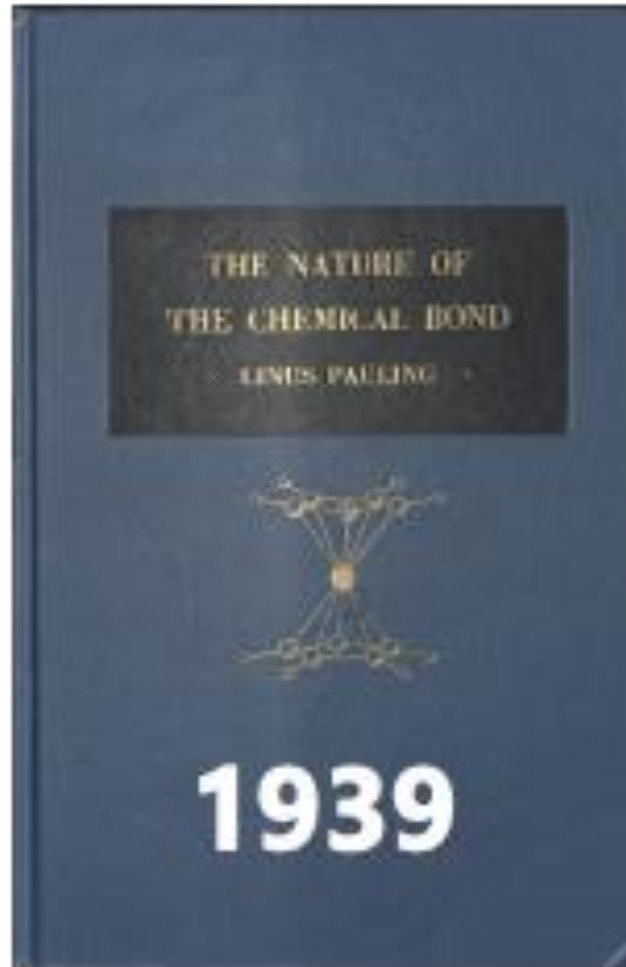
Brief intro to decision trees

Determine splits by **maximizing information gain (IG)**
minimizing weighted impurity, I

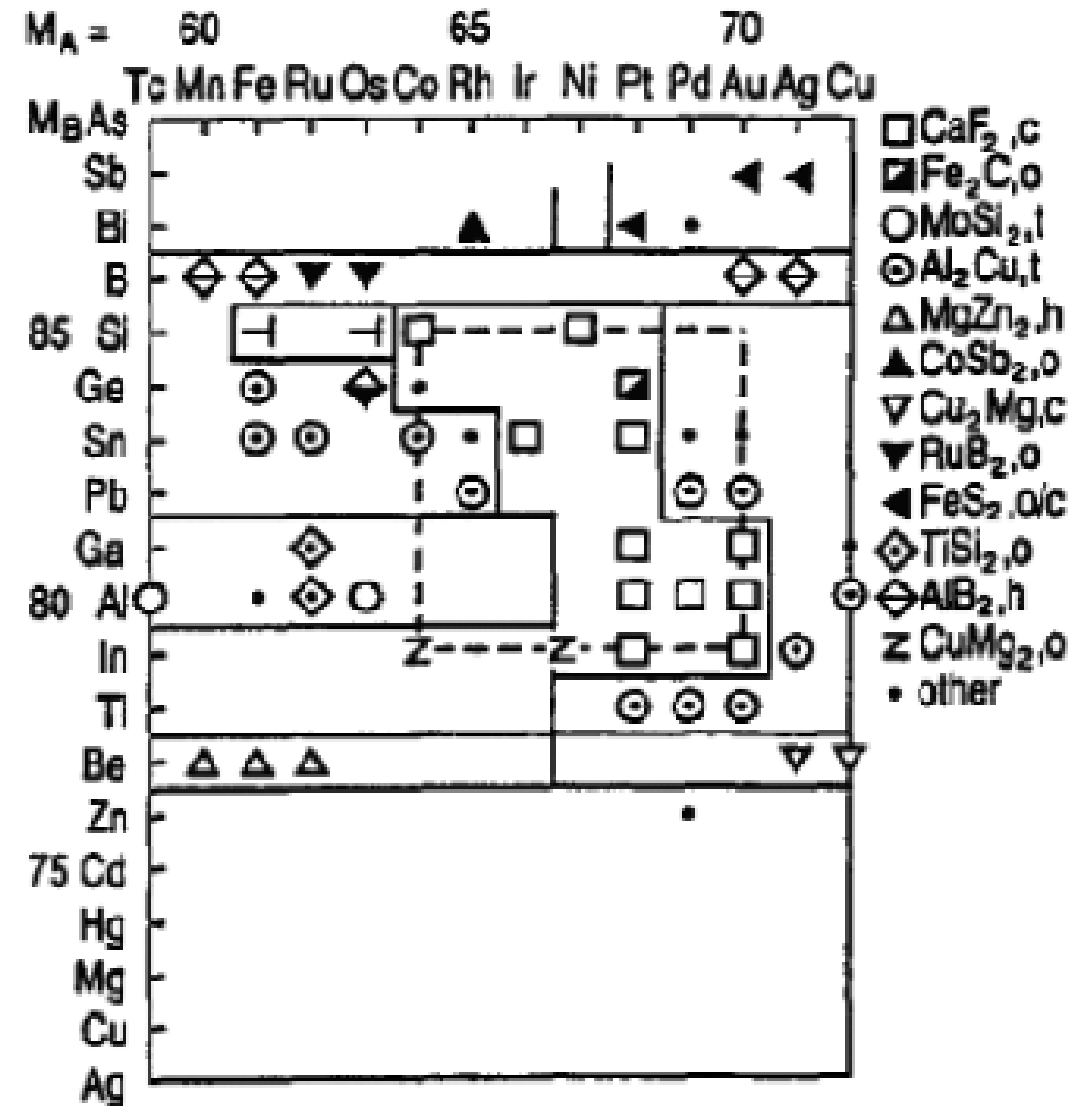
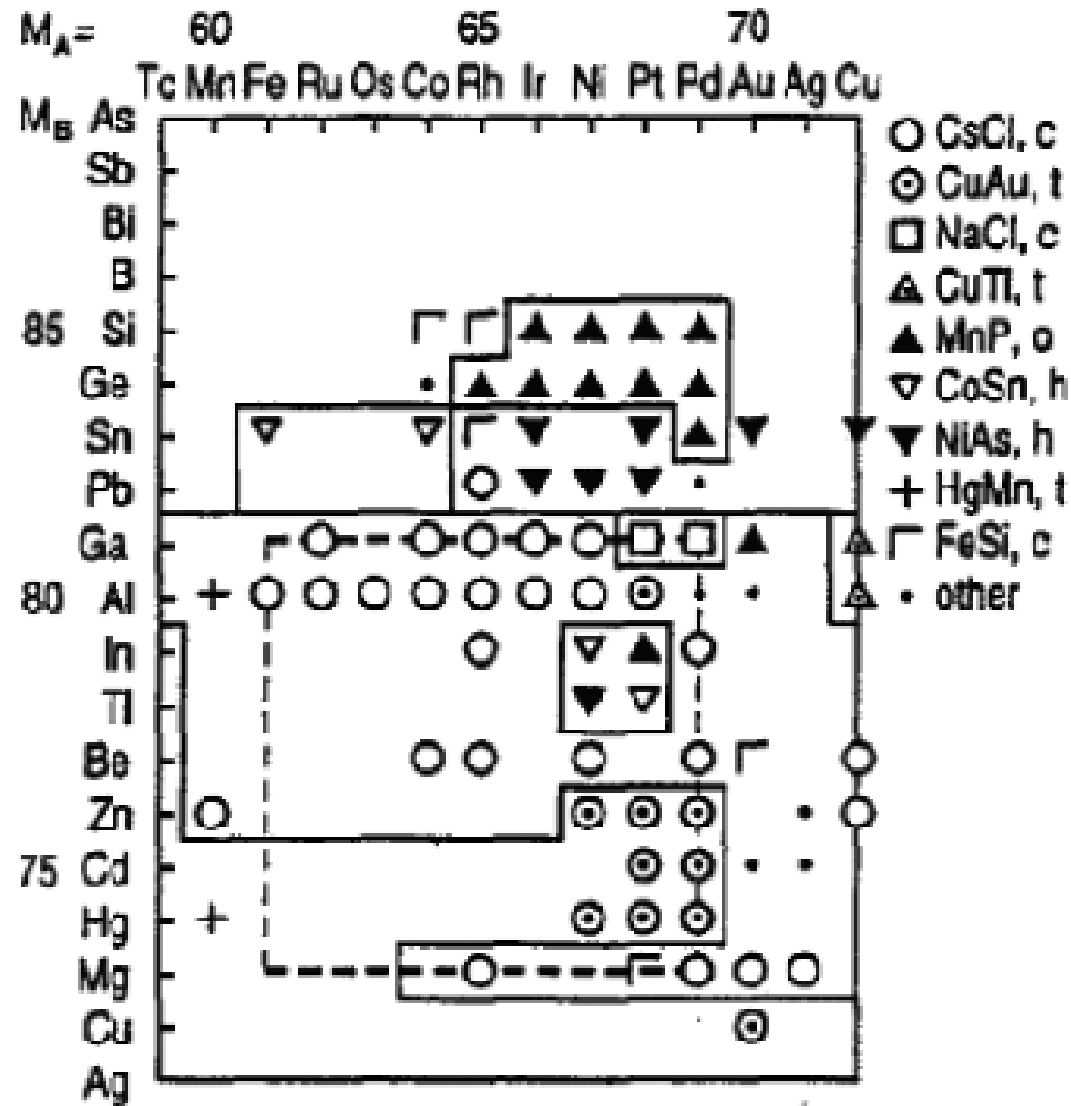
Brief intro to decision trees



Classifying crystal structures



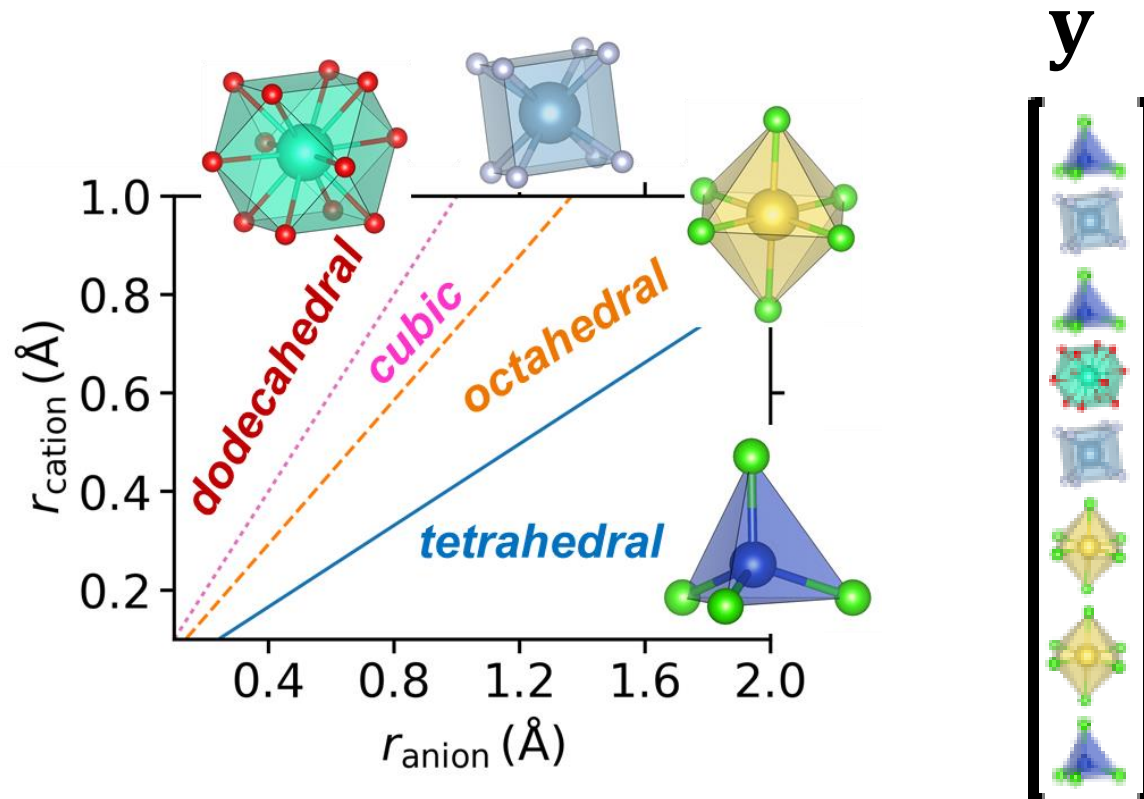
Classifying crystal structures



Finding simple models w/ supervised ML

y – target property (observable)

y – data you find or generate



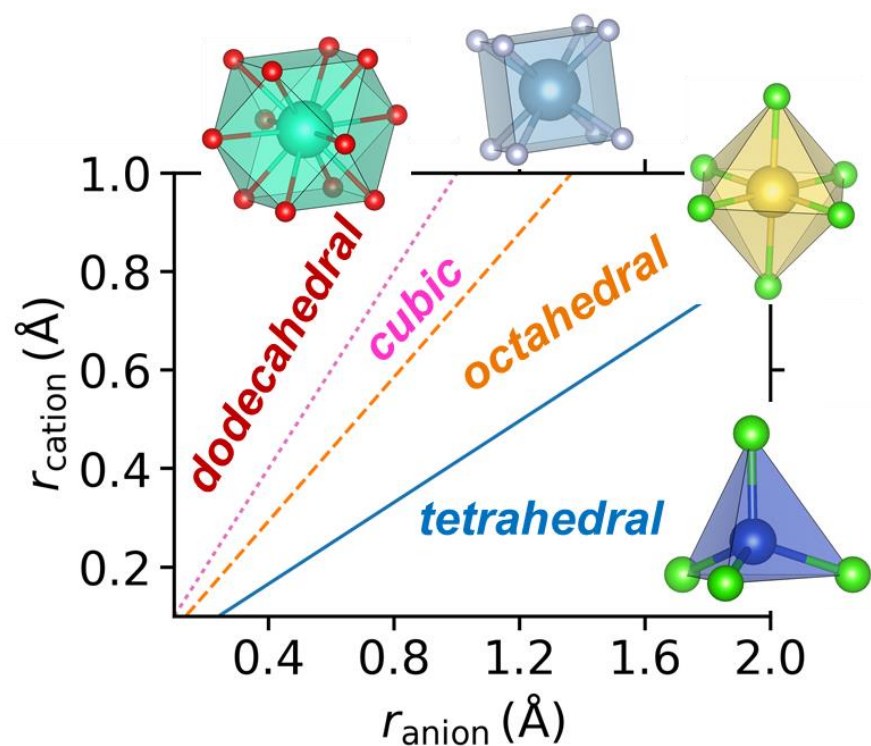
Finding simple models w/ supervised ML




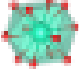




y – target property (observable)

X – feature space (representation)

y – data you find or generate

X – stuff you hope relates to **y**



y	X	
	r_{cation}	r_{anion}
	0.5	1.7
	0.7	1.1
	0.3	1.2

		
		
		
		

Finding simple models w/ supervised ML

y – target property (observable)

X – feature space (representation)

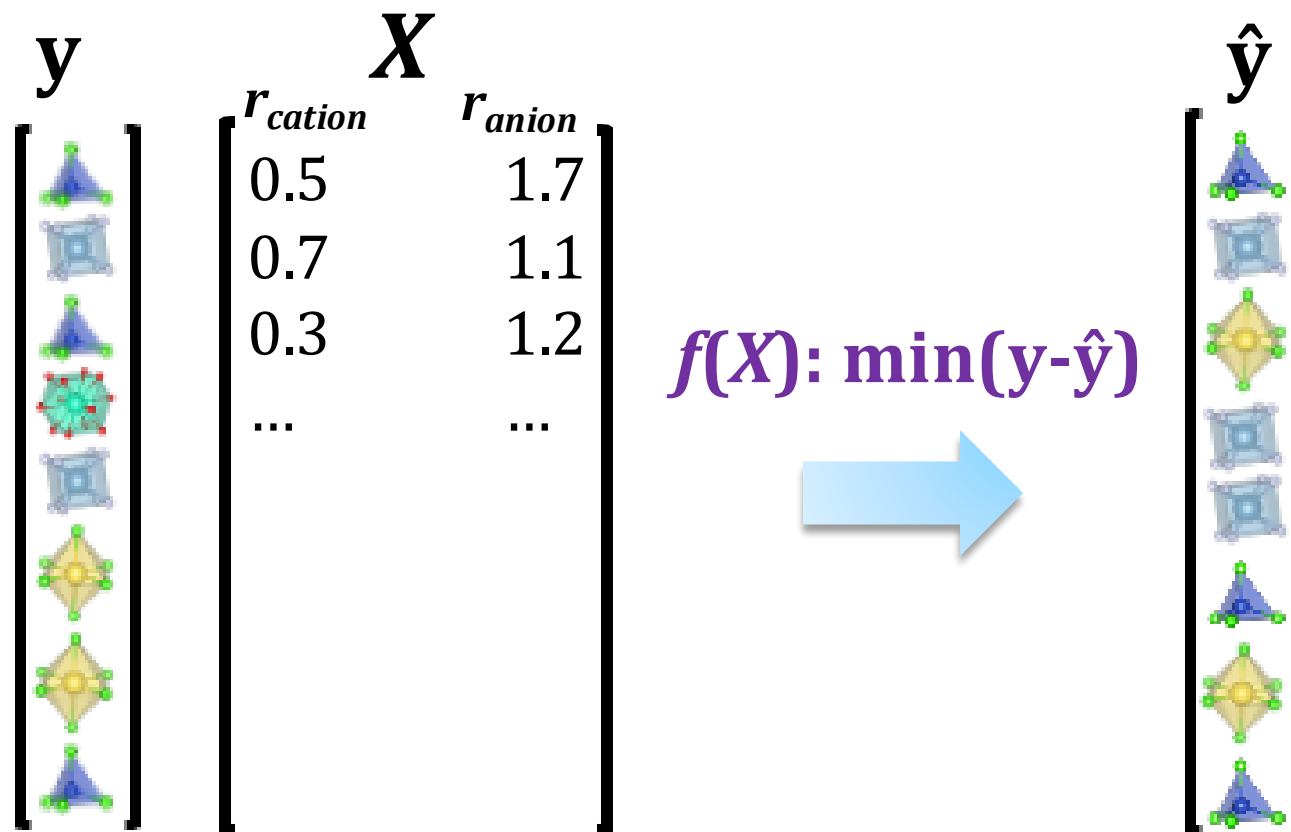
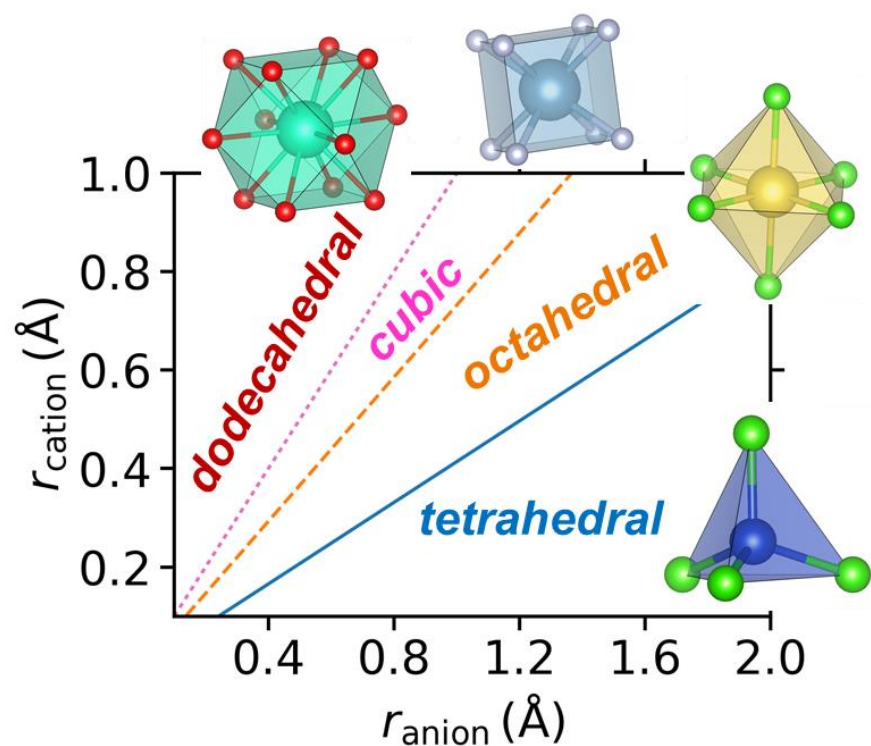
$f(X)$ – model (descriptor)

\hat{y} – prediction (model output)

y – data you find or generate

X – stuff you hope relates to y

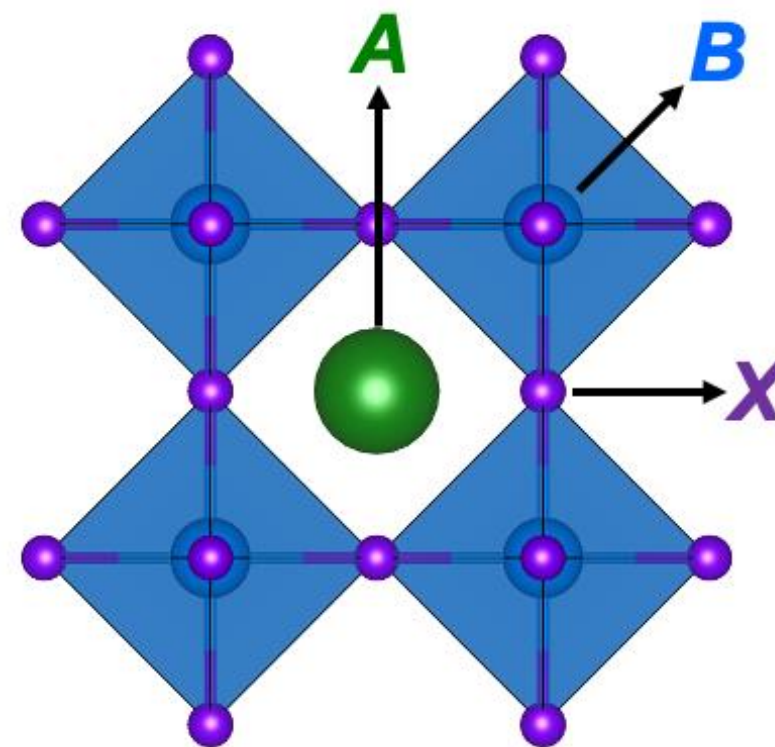
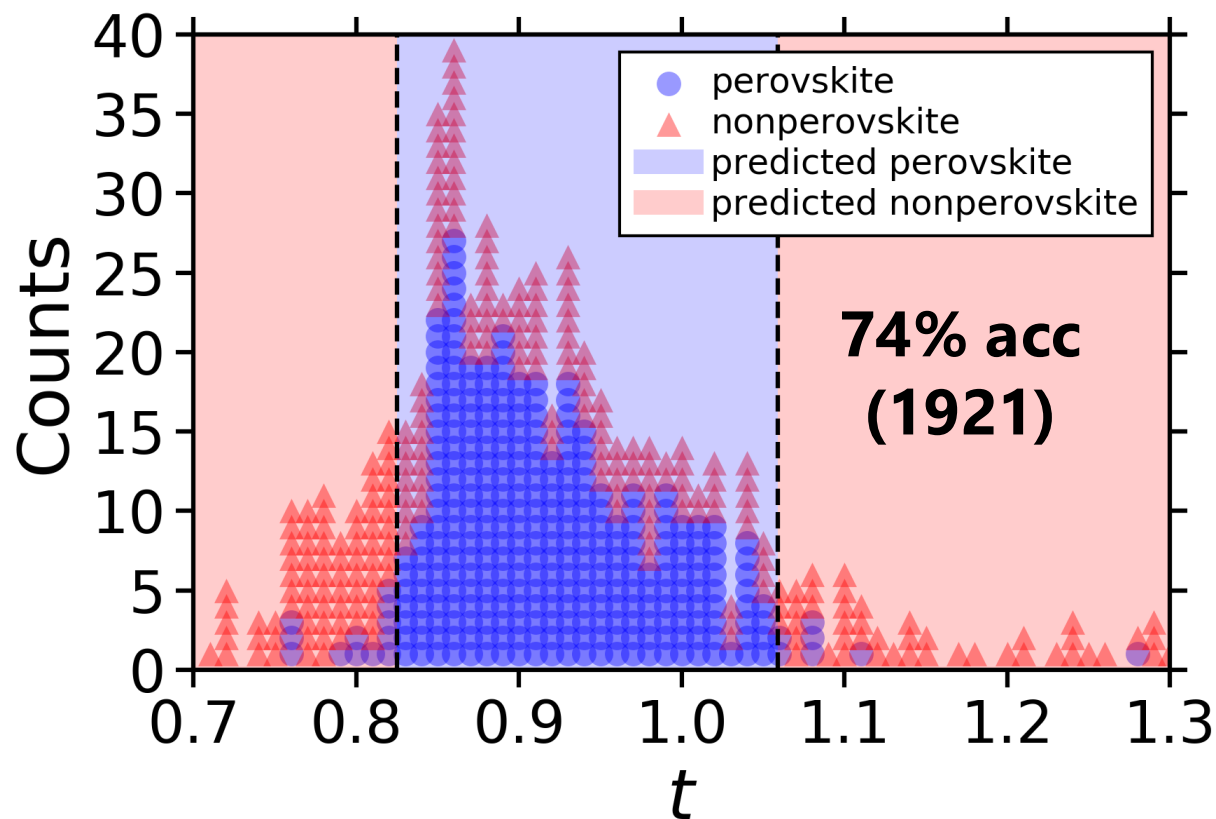
f – the learned mapping of X to y



Goldschmidt's tolerance factor for perovskite stability

For 576 experimentally characterized ABX_3 compounds

$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)}$$



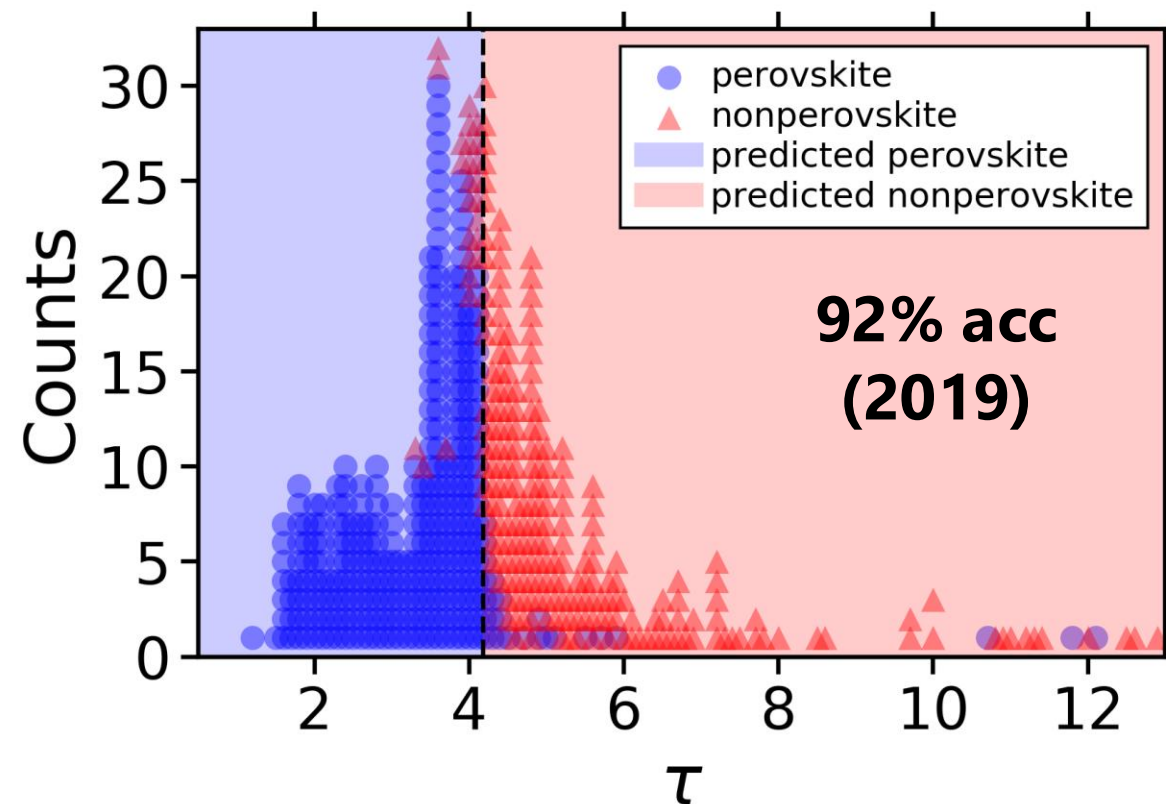
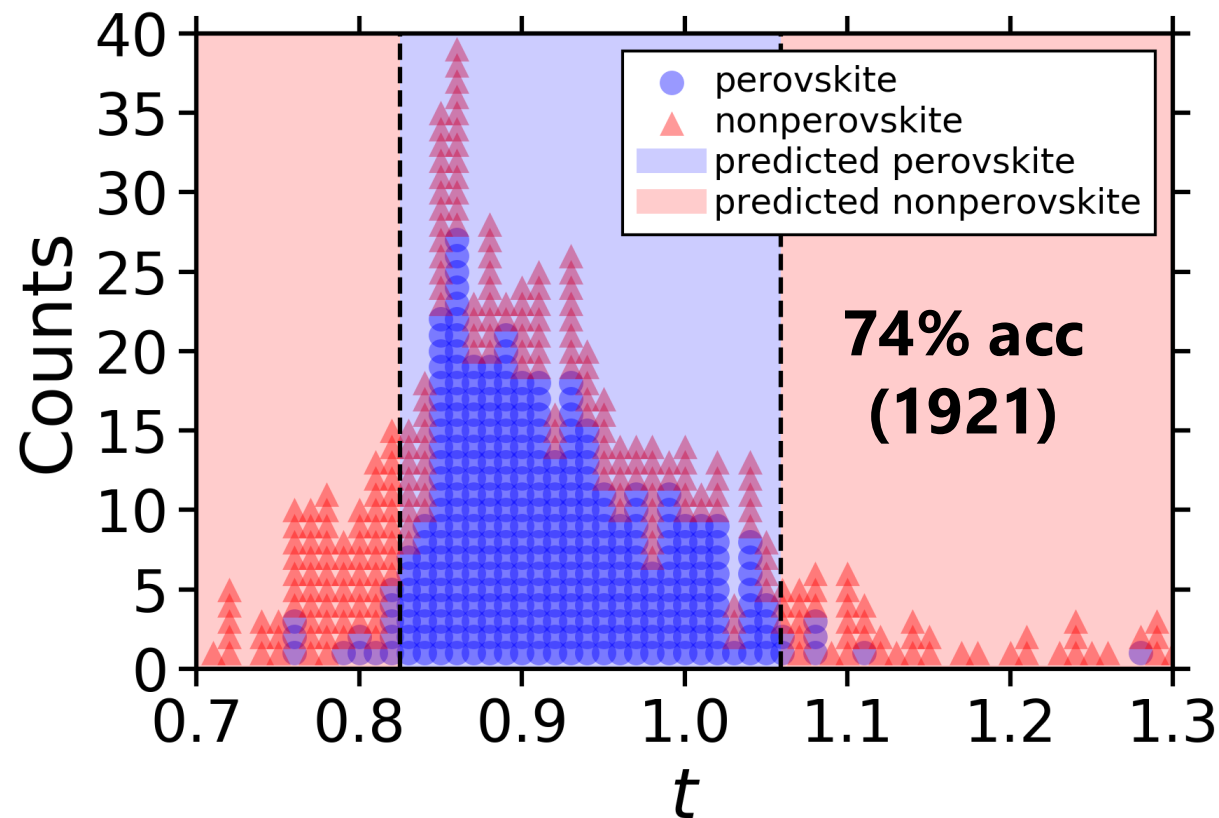
New tolerance factor!

For 576 experimentally characterized ABX_3 compounds

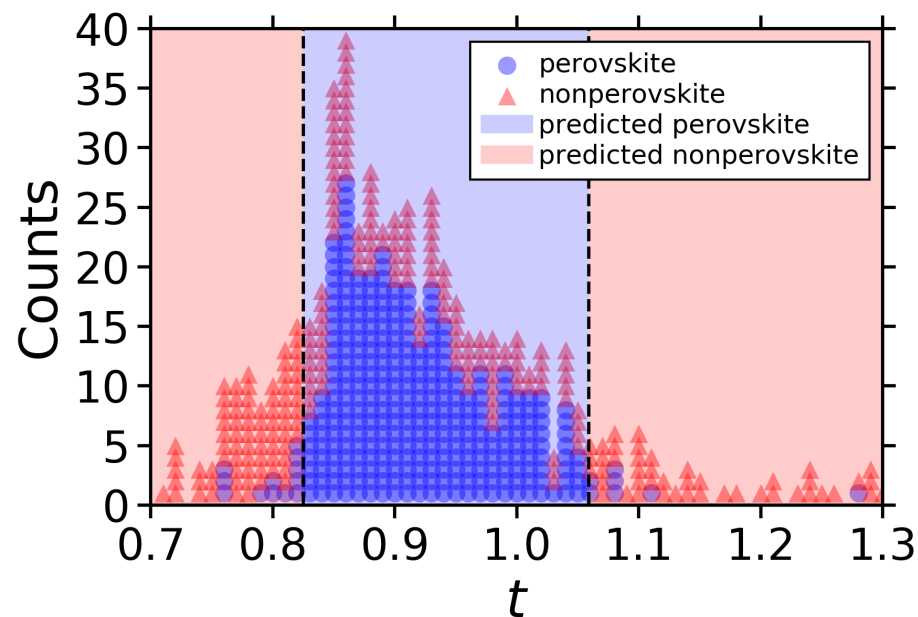
$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)}$$

SISSO

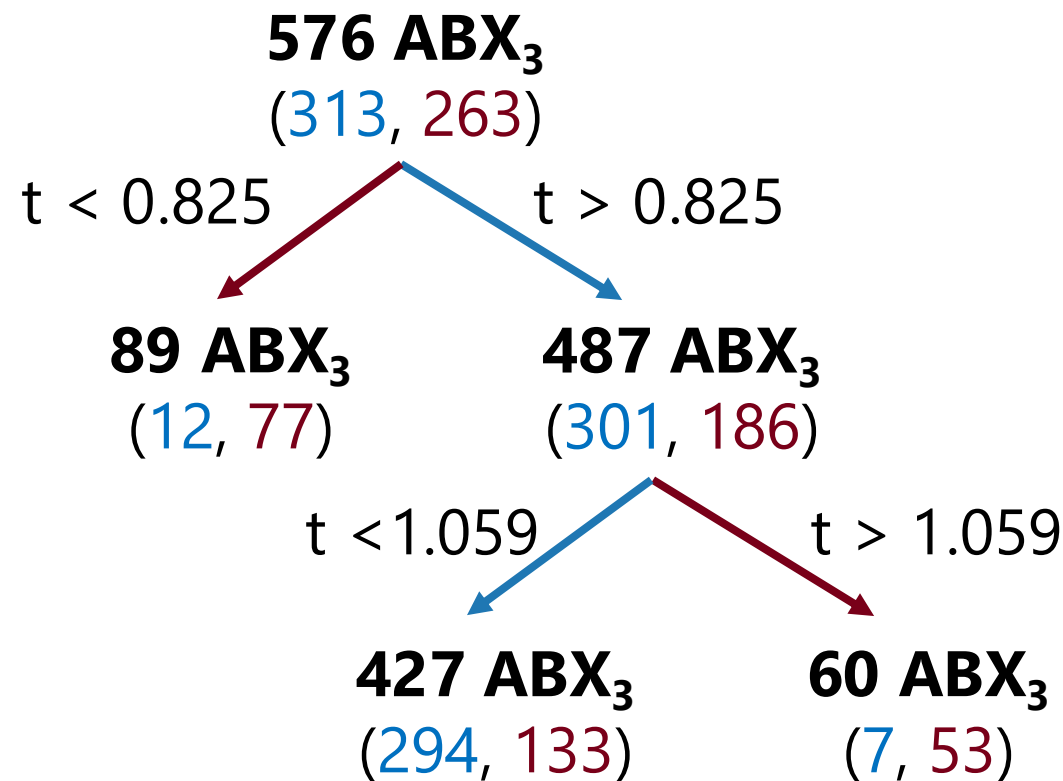
$$\tau = \frac{r_X}{r_B} - n_A \left(n_A - \frac{r_A/r_B}{\ln r_A/r_B} \right)$$



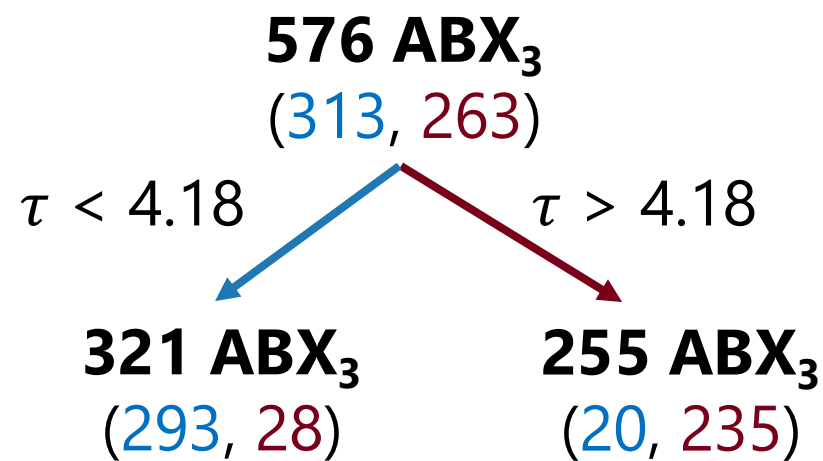
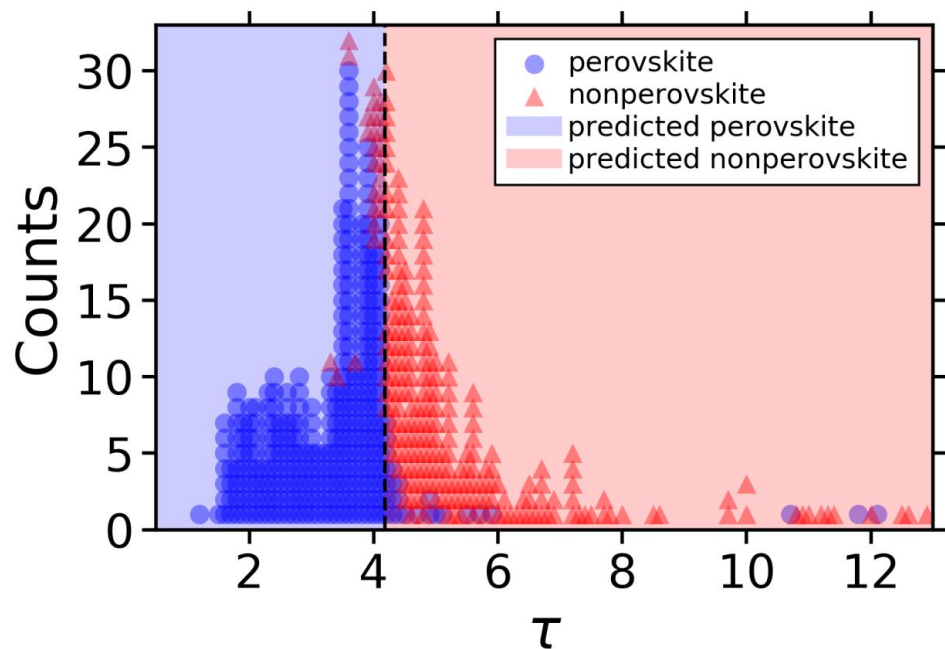
Decision trees w/ Goldschmidt's t



$$\mathbf{y} = \begin{bmatrix} 1 \\ -1 \\ 1 \\ \dots \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)} \end{bmatrix}$$



Decision trees w/ τ



$$\begin{array}{c} \mathbf{y} \\ \left[\begin{array}{c} 1 \\ -1 \\ 1 \\ \dots \end{array} \right] \end{array} \quad \begin{array}{c} \mathbf{x} \\ \left[\tau = \frac{r_X}{r_B} - n_A \left(n_A - \frac{r_A/r_B}{\ln r_A/r_B} \right) \right] \end{array}$$