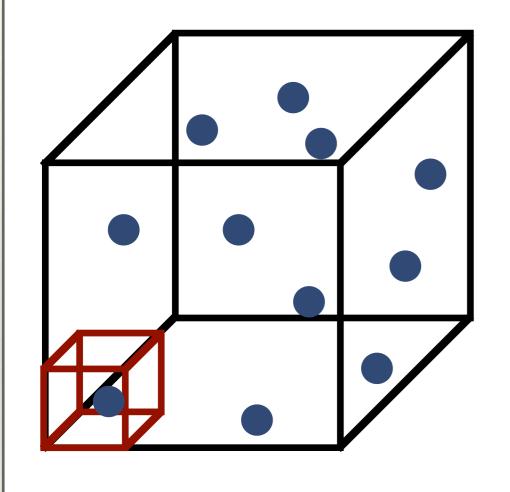
Feature Selection

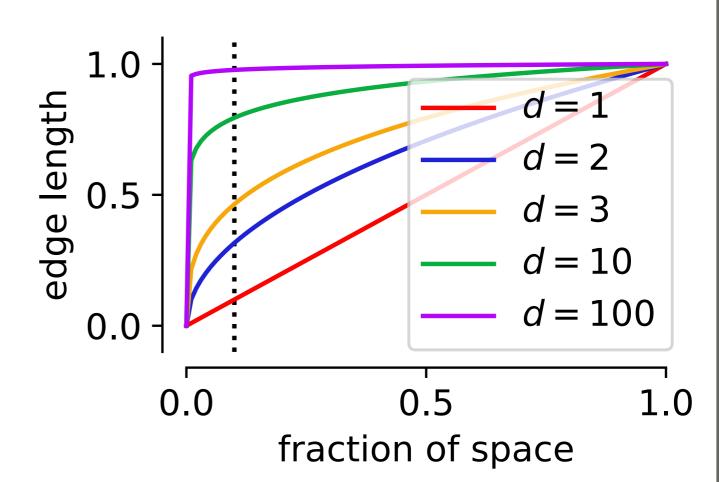
Key points:

- Why do we do this?
- Some techniques that we can do for that

The Curse of Dimensionality - No Locality in High Dimensions

 $edgelength = fraction\ of\ space^{1/dimensionality}$





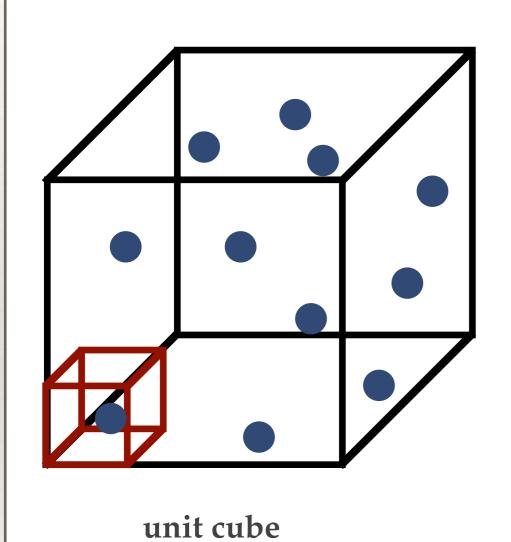
unit cube

Methods that are based on similarity (kNN/KRR) might fail in high dimensional spaces!

We want to **reduce the dimensionality** of our feature matrix!

The Curse of Dimensionality - No Locality in High Dimensions

 $edgelength = fraction\ of\ space^{1/dimensionality}$



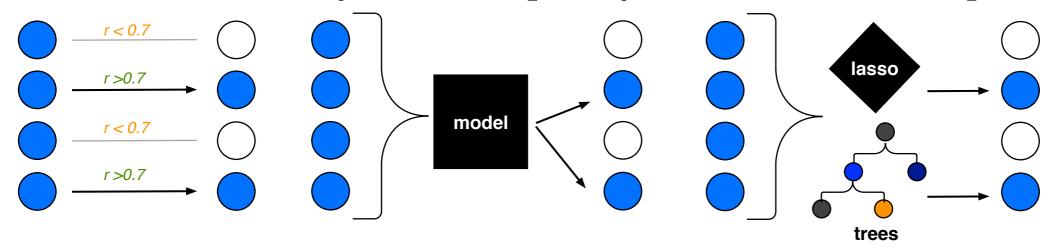


Methods that are based on similarity (kNN/KRR) might fail in high dimensional spaces!

We want to **reduce the dimensionality** of our feature matrix!

Feature Projection and Feature Selection

Reduce dimensionality of feature space by feature selection (compression)

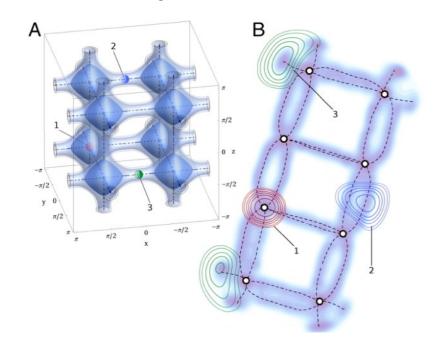


a Univariate filters.

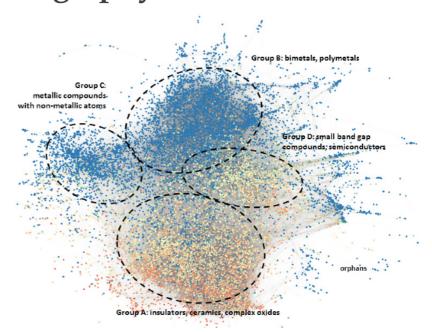
b Wrapper methods.

c Shrinkage or direct.

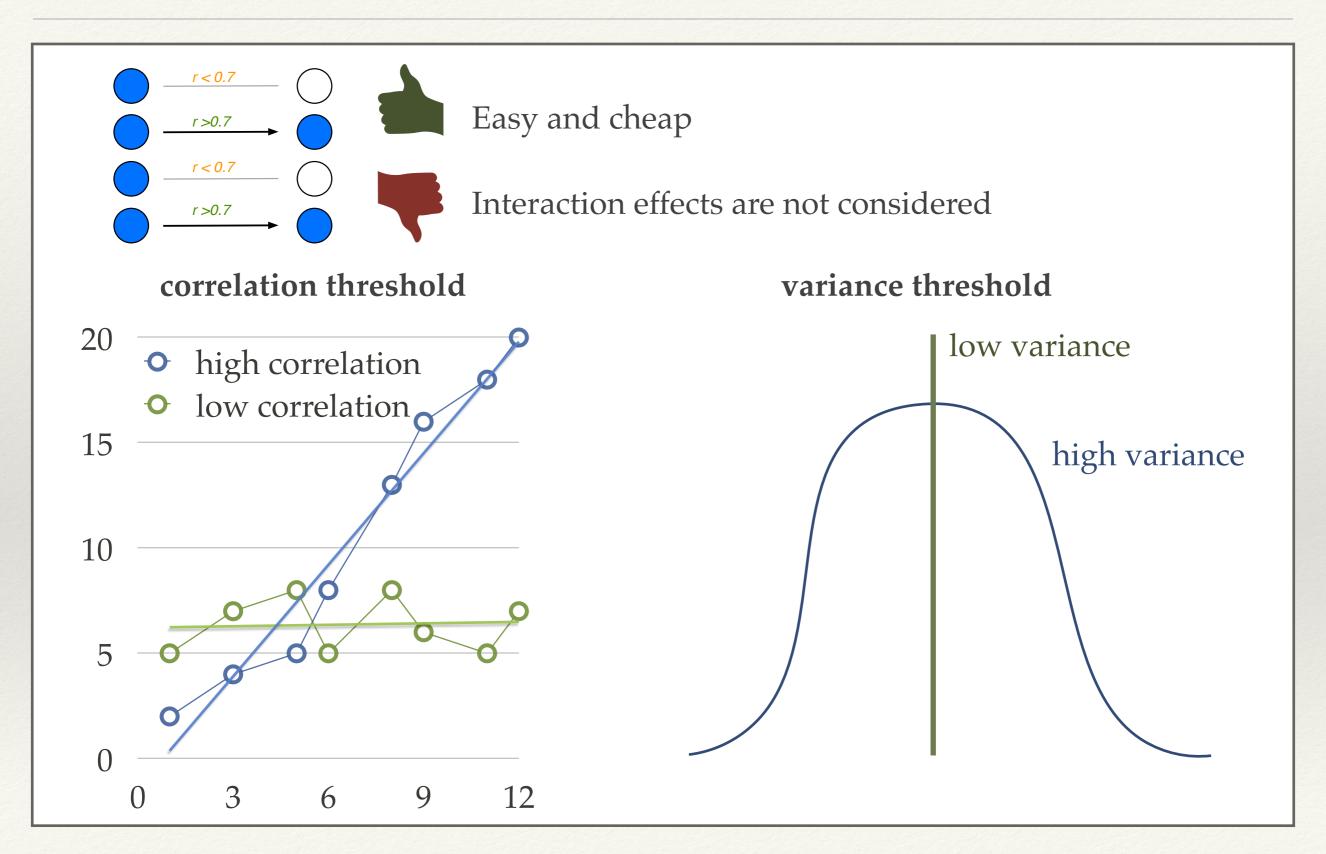
Reduce size of feature space by dimensionality reduction (feature projection)



Visualize data and Materials Cartography



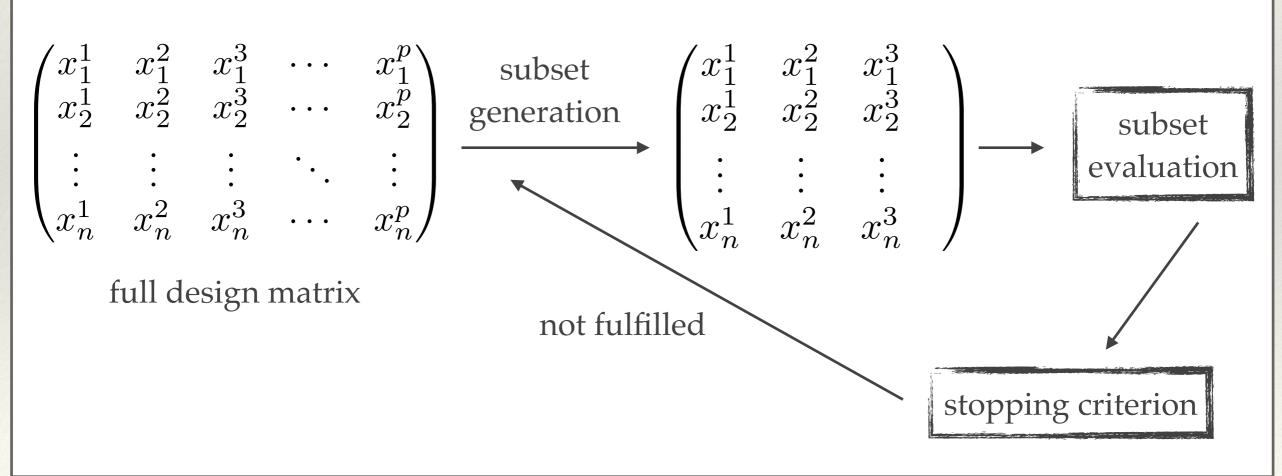
Feature Selection: Filter Methods



Feature Selection: Wrapper Methods



For example recursive feature addition or elimination



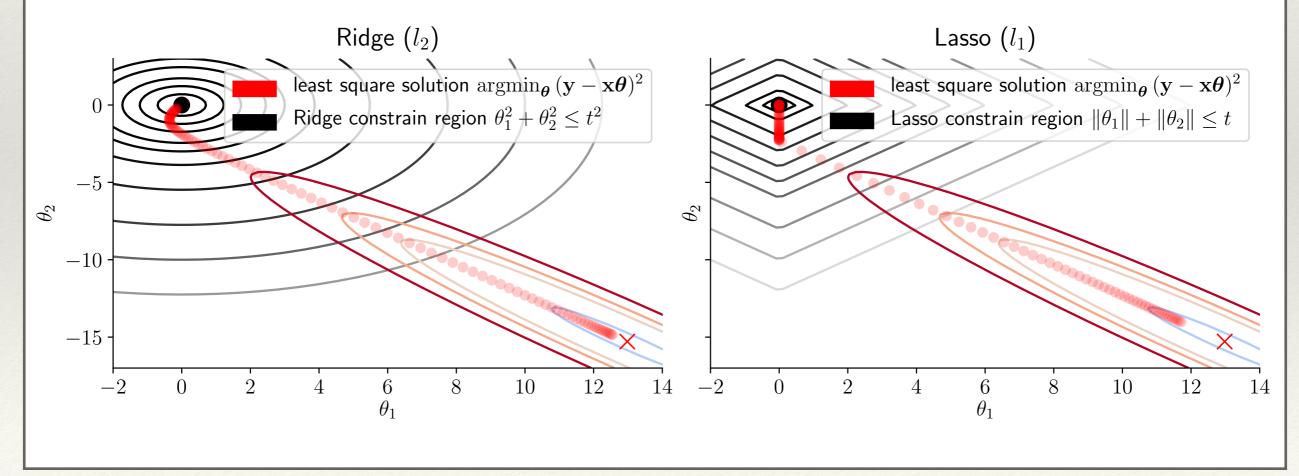
Feature Selection: Just Relax Best Subset Selection

The basic problem: Best subset selection

$$\min_{\beta \in \mathbb{R}^p} \|y - X\beta\|^2 \text{ subject to } \|\beta\|_0 \le k \qquad \|\beta\|_0 = \sum_{j=1}^p 1\{\beta_j \ne 0\}$$

But this is our hard problem (NP hard) ...

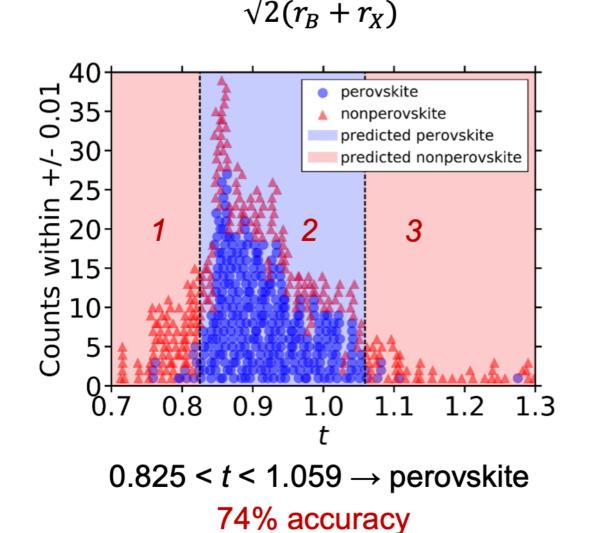
- ... hence we relax the constraint to have a problem that is convex
- ... the Lasso gives use sparsity as the most feasible approximation to l_0

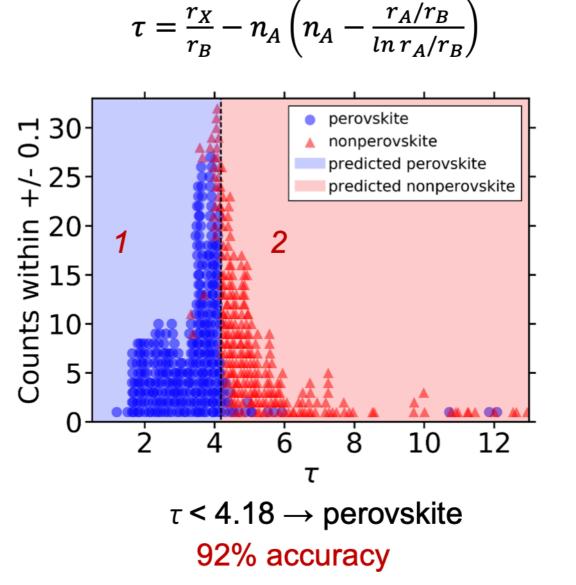


Hastie, T.; Tibshirani, R.; Wainwright, M. *Statistical Learning with Sparsity: The Lasso and Generalizations*; Monographs on statistics and applied probability; CRC Press, Taylor & Francis Group: Boca Raton, 2015.

Lasso in Practice: Finding New Tolerance Factors For Perovskites (Developing Causal Models)

primary features $\xrightarrow{\hat{R}}$ billions of feature candidates $\xrightarrow{\text{(SI)} + \text{Lasso}}$ subset of features primary features = $\{r_A, r_B, n_a, n_b, ...\}$ $\hat{R} = \{+, -, \cdot, \exp, \lg, -1, 2, 3, \sqrt{,} || \cdot || \}$ $t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_Y)}$ $\tau = \frac{r_X}{r_B} - n_A \left(n_A - \frac{r_A/r_B}{\ln r_A/r_B}\right)$





Feature Selection

Key points:

- Why do we do this?
 - Curse of dimensionality
- Some techniques that we can use for that
 - Filter, Wrapper, Direct
 - Difference between selection and projection