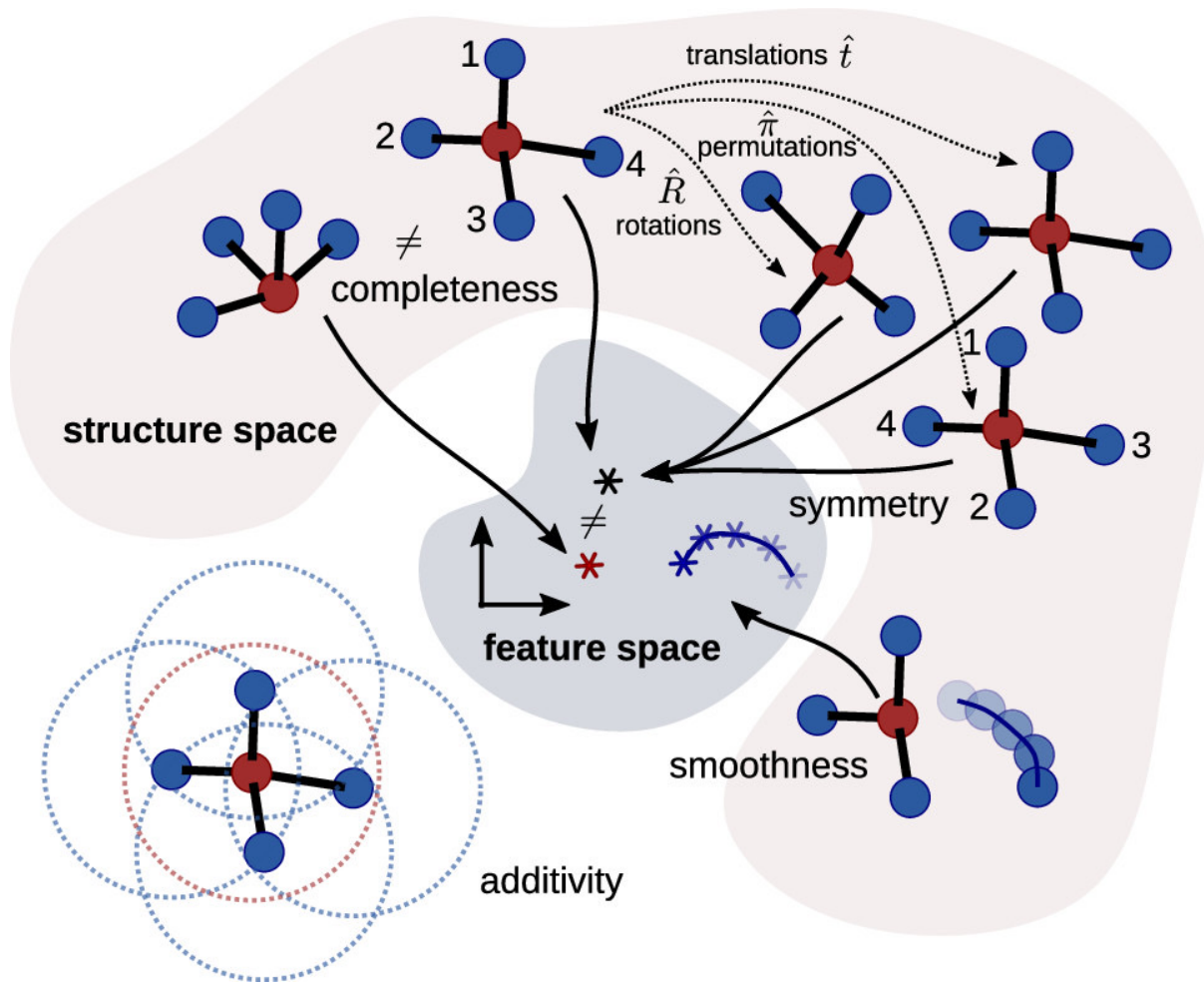


# Representations





Expected error given sufficient data

$$\exp(-c/h_{X,\Omega})$$

Positive constant  $c$

Function of the „fill distance“

$$h_{X,\Omega} \equiv \sup_{y \in \Omega} \min_{x_j \in X} \|y - x_j\|_2$$

Training data  $X$   
Domain  $\Omega$

Shorter fill distance: steeper learning curve, i.e. more data efficient model

# Categorical: One-Hot

37

Categorical data vs regression

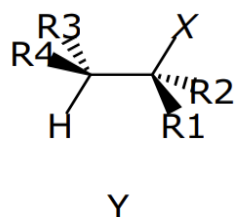
Solution: „binary“, „dummy“, „one-hot“ encoding

Encode  $n$  categories as vector of length  $n-1$  with one category (arbitrary) being the null vector.

Example:

- A: (0, 0)
- B: (1, 0)
- C: (0, 1)

Chemistry example: (A B B A | C D)  $\rightarrow$  (0 0 0 0 | 1 0 0 0 | 1 0 0 0 | 0,0,0,0 | 0,1,0,0,1)



	A	B	C	D	E
Rk	H	NO <sub>2</sub>	CN	CH <sub>3</sub>	NH <sub>2</sub>
X	F	Cl	Br		
Y	H	F	Cl	Br	

# Graph-based

38

Often: molecules = well-defined bonds

## Adjacency matrix

- 1 if atoms  $i$  and  $j$  are bonded
- 0 otherwise

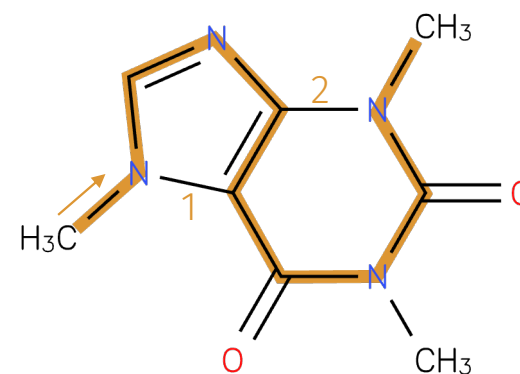
## Bond order matrix

- Bond order if atoms  $i$  and  $j$  are bonded
- 0 otherwise

## SMILES

- Bonds: Nothing (1), = (2), # (3), \$ (4)
- Partial charges
- Fragments: .
- Rings: labels
- Branches: parentheses

O=C=O  
[Na+]  
[Na+].[Cl-]  
C1CCCCC1

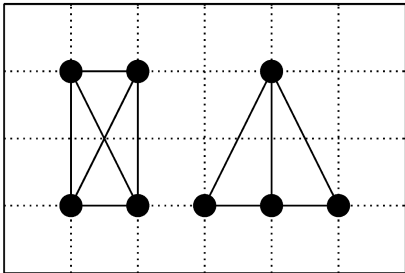

$$\begin{matrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{matrix}$$
$$\begin{matrix} 0 & 2 & 0 \\ 2 & 0 & 2 \\ 0 & 2 & 0 \end{matrix}$$


CN1C=NC2=C1C(=O)N(C)C(=O)N2C

# Distance-based

39

Many-body descriptions, e.g. all pairwise distances



	s	l	m
s		m	l
l	m		s
m	l	s	

	l	l	m
l		m	s
l	m		s
m	s	s	

Even three- and four-body interactions not unique [1].

Chemistry: **Coulomb Matrix** [2]

- Diagonal:
- Off-diagonal:
- Problem: sorting, uniqueness

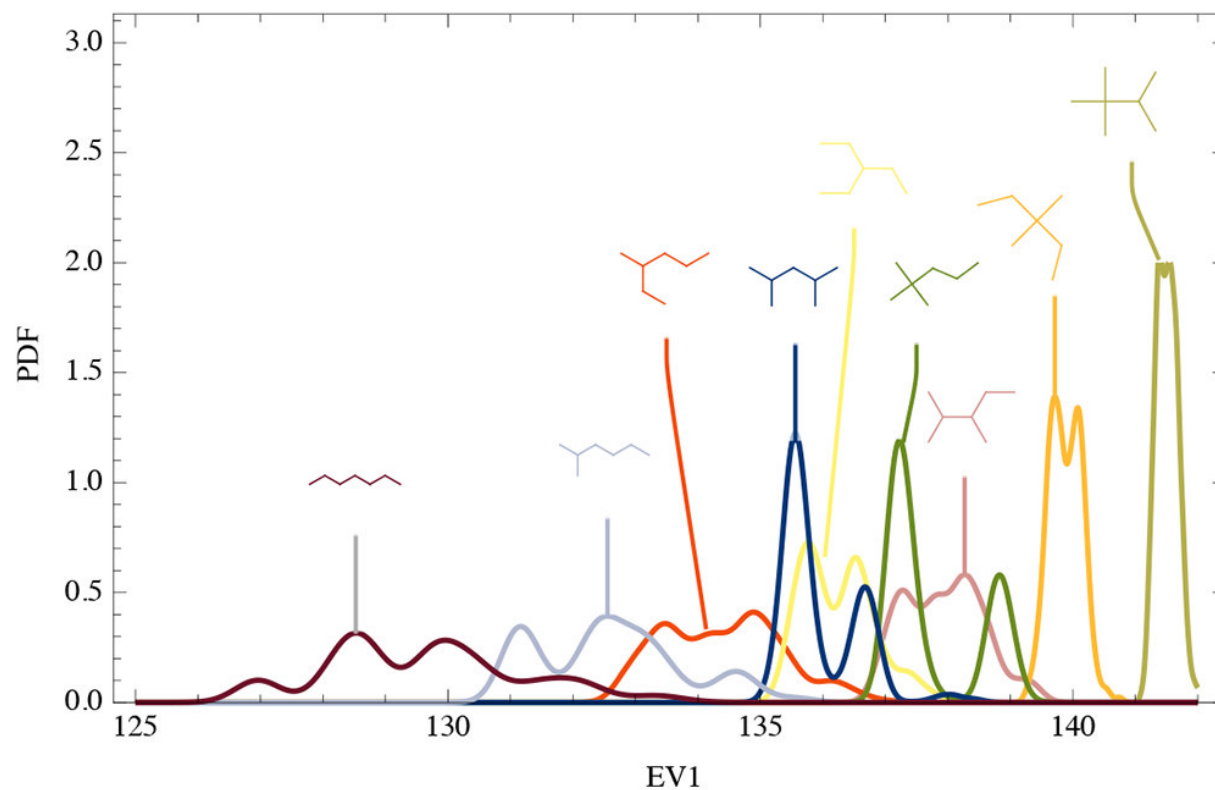
$$0.5 Z_i Z_j / \|\mathbf{R}_i - \mathbf{R}_j\|$$

[1] Pozdnyakov et al., *Phys Rev Lett* 2020. [2] M. Rupp et al., *Phys Rev Lett* 2012.



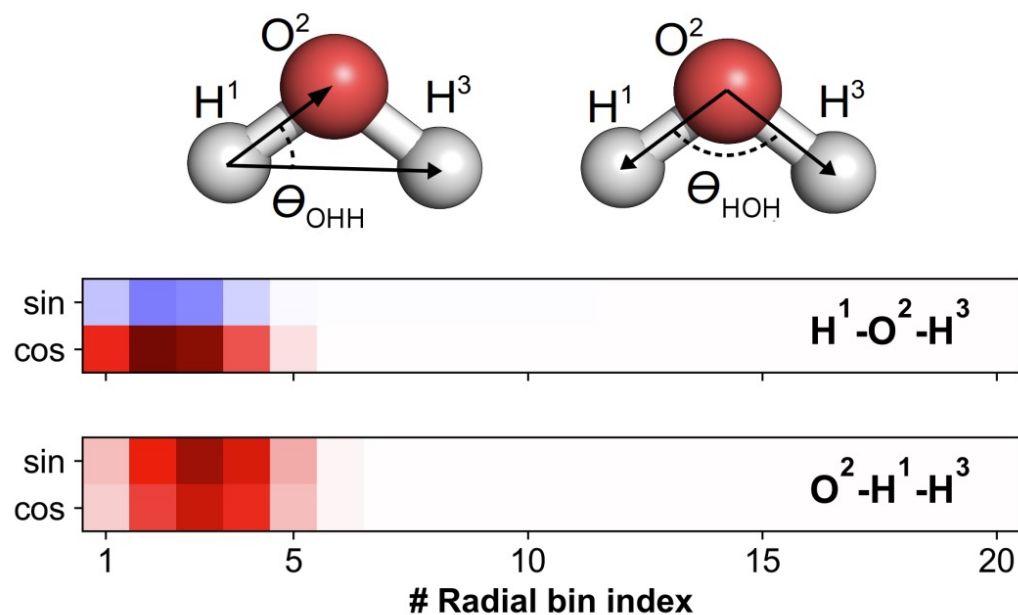
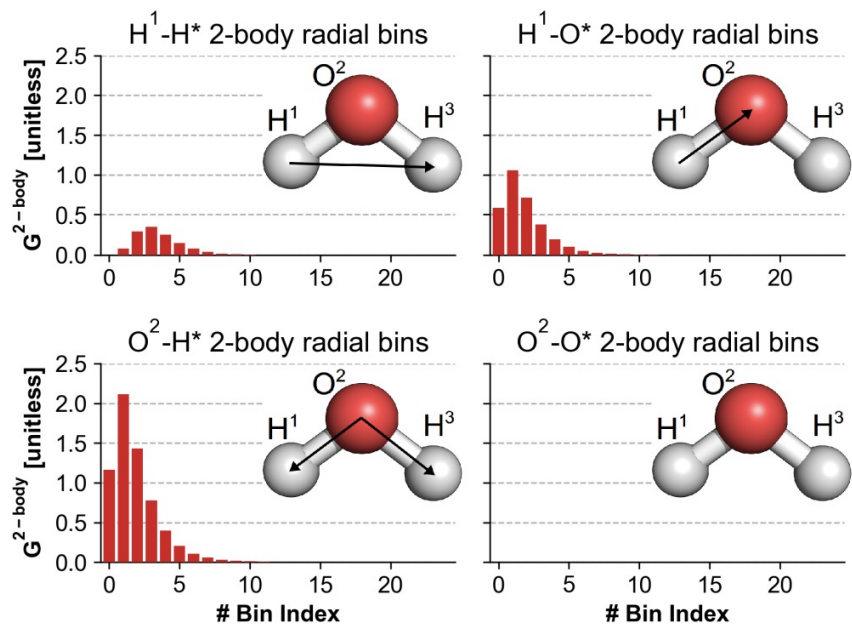
Condense a matrix into its eigenvalues

- No sorting issue, as permutationally invariant
- Smaller:  $N$  instead of  $N^2$
- Lossy



J. Schrier, *J Chem Inf Model*, 2020.

Smear positions into densities (FCHL / SOAP)



# Summary Representations

43

- Helpful properties of a representation
  - Unique
  - Smooth in representation space
  - Changes in geometry smooth in representation
  - Complete
  - Invariant under transformations (translation, rotation, permutation)
  - Compact
  - Additive
  - Invertible
  - Fast
  - Simple
  - Containing many-body effects