

Machine Learning for Quantum Mechanics

Matthias Rupp

Fritz Haber Institute of the Max Planck Society, Berlin, Germany

2017 KITS Workshop on
Machine Learning and Many-Body Physics
Tutorials, June 28–30, Beijing, China



Outline

1. Rationale

quantum mechanics, machine learning

2. Kernel learning

kernel trick, kernels, regression

3. Model building

overfitting, validation, hyperparameters

4. Applications

examples, representation

Rationale

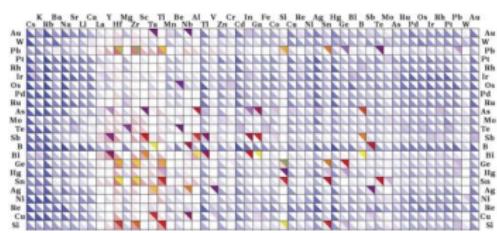
Rationale

“The underlying physical laws necessary for [. . .] chemistry are thus completely known, and the difficulty is only that the exact application of these laws leads to equations much too complicated to be soluble. It therefore becomes desirable that approximate practical methods of applying quantum mechanics should be developed.”

Paul A.M. Dirac

Challenges in quantum mechanical simulations

High-throughput screening



Large systems

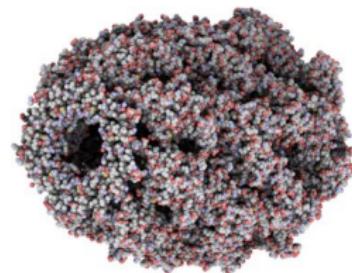
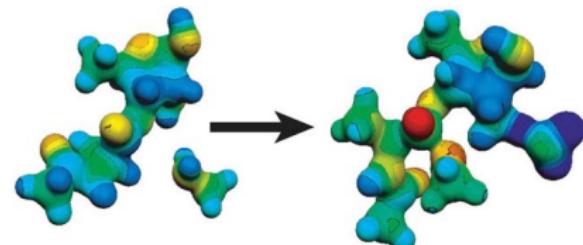


Image: Tarini et al, *IEEE Trans Visual Comput Graph* 2006

Long simulations



Quantum effects



Liwo et al, *Proc Natl Acad Sci USA* 102: 2362, 2005

Image: Hiller et al, *Nature* 476: 236, 2011

Approximations

Hierarchy of numerical approximations to Schrödinger's equation:

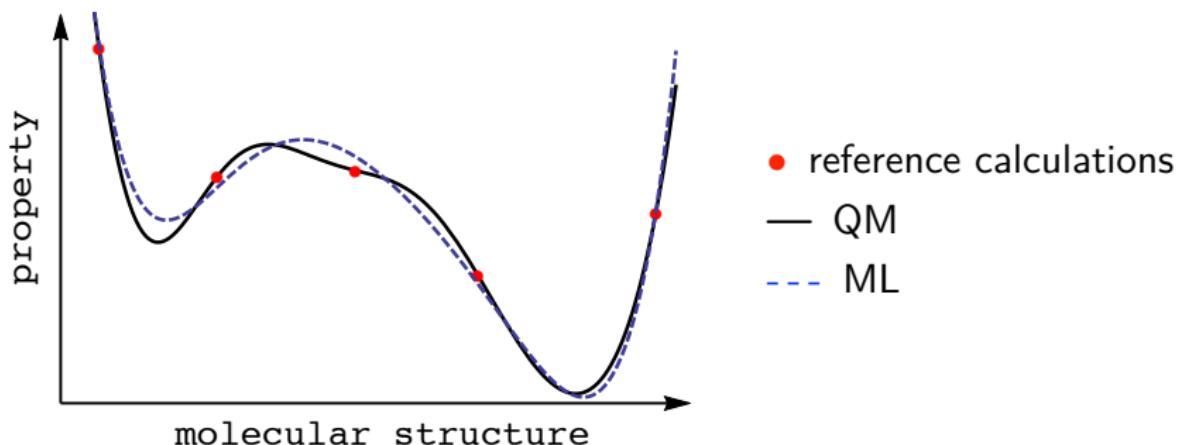
Abrv.	Method	Runtime
FCI	Full Configuration Interaction (CISDTQ)	$O(N^{10})$
CC	Coupled Cluster (CCSD(T))	$O(N^7)$
FCI	Full Configuration Interaction (CISD)	$O(N^6)$
MP2	Møller-Plesset second order perturbation theory	$O(N^5)$
HF	Hartree-Fock	$O(N^4)$
DFT	Density Functional Theory (Kohn-Sham)	$O(N^{3-4})$
TB	Tight Binding	$O(N^3)$
MM	Molecular Mechanics	$O(N^2)$

N = system size

Is it possible to be both accurate and fast?

The key idea

- Exploit redundancy in related QM calculations
- **Interpolate** between QM calculations using ML
- Smoothness assumption (regularization)



Machine learning

Machine learning (ML) studies algorithms whose performance **improves with data** (“learning from experience”).

Mitchell, McGraw Hill, 1997



- widely applied, many problem types and algorithms
- systematic identification of regularity in data for prediction & analysis
- interpolation in high-dimensional spaces
- inductive, data-driven; empirical in a principled way
- connections to statistics, mathematics, computer science, physics, . . .
example: information theory

Problem types

Unsupervised learning: Data do not have labels

Given $\{x_i\}_{i=1}^n$, find structure

- dimensionality reduction

Burges, now Publishers, 2010

Supervised learning: Data have labels

Given $\{(x_i, y_i)\}_{i=1}^n$, predict \tilde{y} for new \tilde{x}

- novelty detection
- classification
- regression
- structured output learning

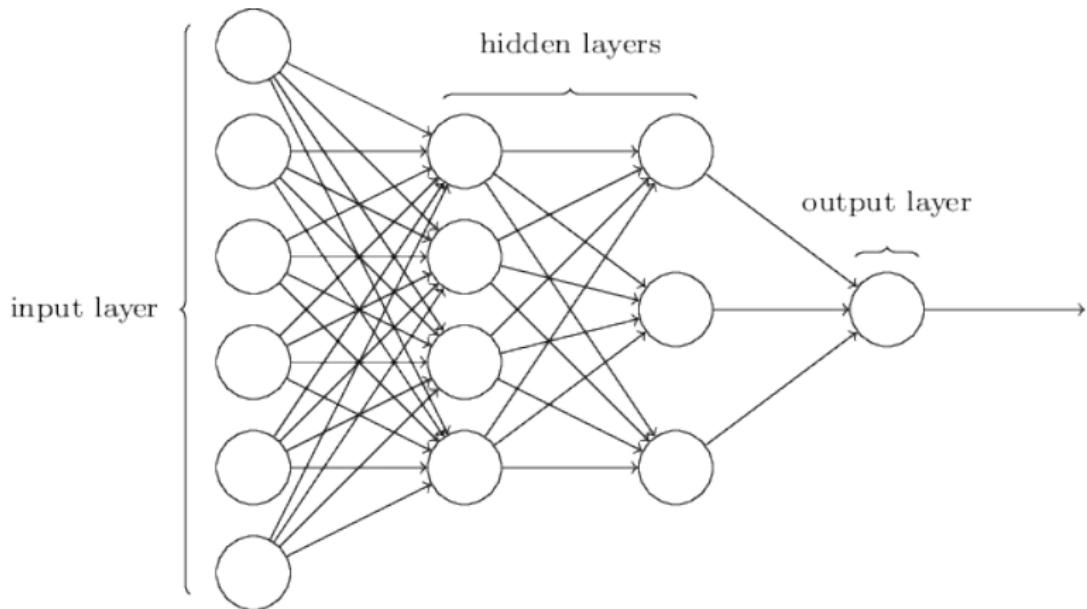
Semi-supervised learning: Some data have labels

Given $\{(x_i, y_i)\}_{i=1}^n$ and $\{x_i\}_{i=1}^m$, $m \gg n$, predict \tilde{y} for new \tilde{x}

Active learning: Algorithm chooses data to label

Choose n data $\{x_i\}_{i=1}^n$ to predict \tilde{y} for new \tilde{x}

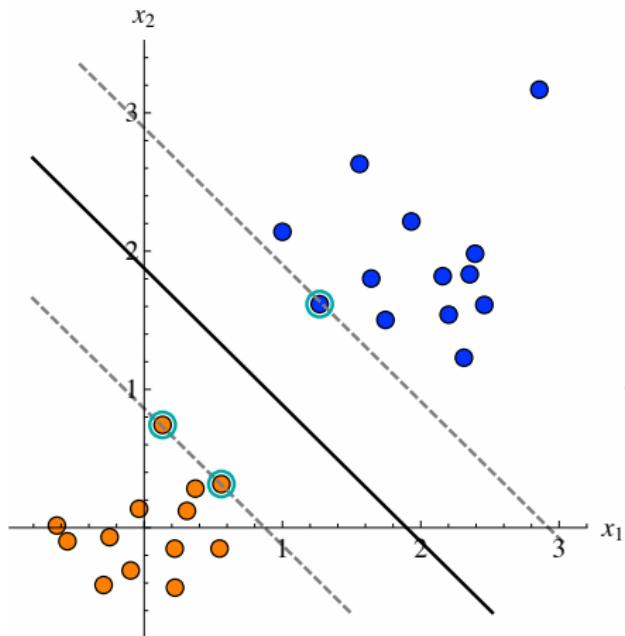
Artificial neural networks



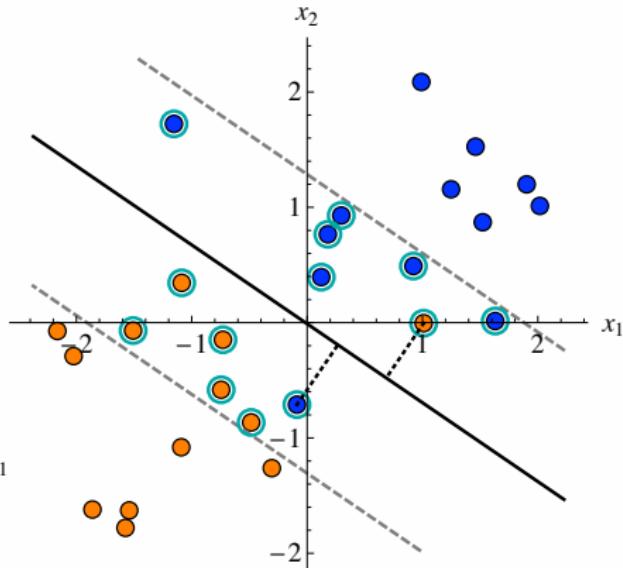
$$f(x_{i,j}) = h\left(\sum_{k=1}^{n_i} w_{i-1,k} f(x_{i-1,k})\right)$$

- parametric model
- universal function approximator
- training via non-convex optimization

Support vector machines



linear separable problem

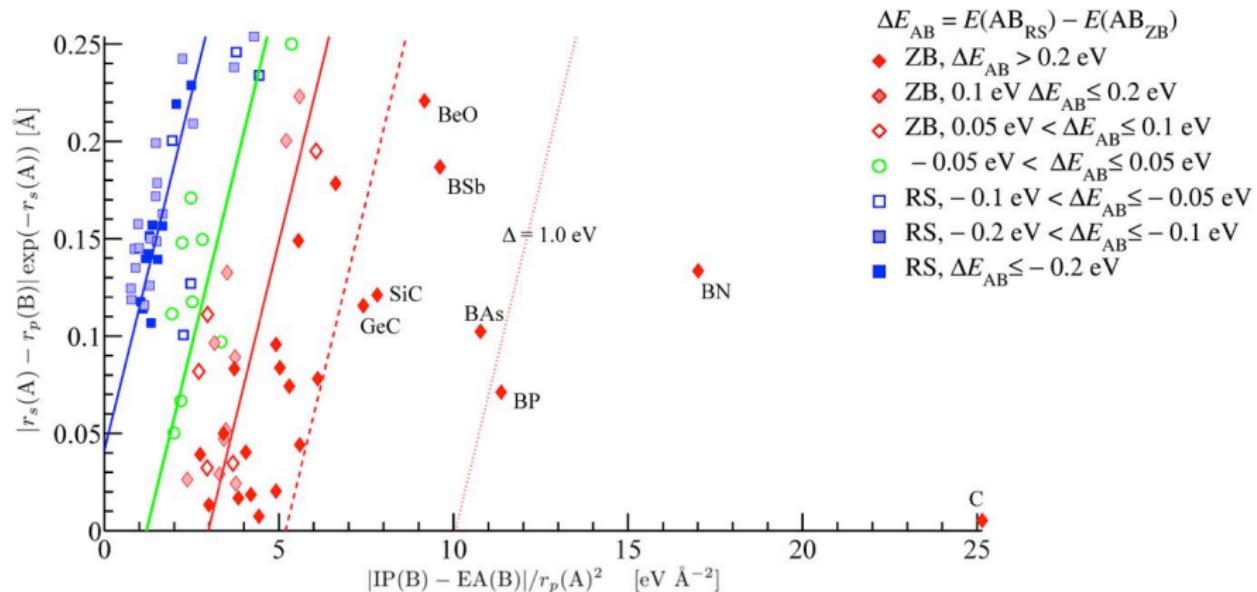


linear inseparable problem

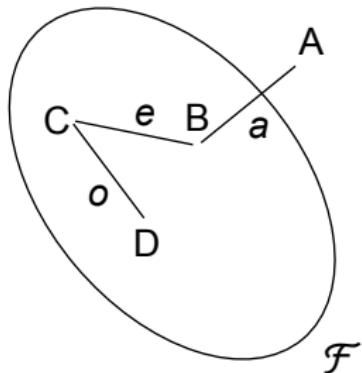
maximal margin plane bisects (reduced) convex hull closest points

Symbolic regression

- stochastic search in the space of analytic functions
- fast, interpretable models



Learning theory



prediction error =
approximation error a
+ estimation error e
+ optimization error o

\mathcal{F} = model class, A = true model, B = best model in class, C = best identifiable model (data), D = best identifiable model (optimization)

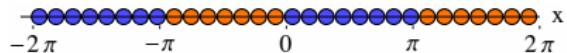
Changes in size of $\mathcal{F} \Leftrightarrow a$ vs. $e \Leftrightarrow$ **bias-variance trade-off**

Kernel learning

The kernel trick

Idea:

- **Transform** samples into higher-dimensional space
- **Implicitly** compute inner products there
- Rewrite linear algorithm to use only inner products

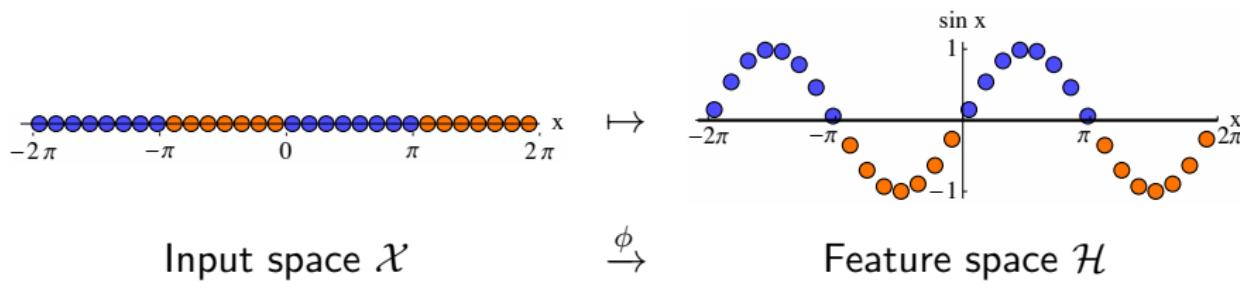


Input space \mathcal{X}

The kernel trick

Idea:

- **Transform** samples into higher-dimensional space
- **Implicitly** compute inner products there
- Rewrite linear algorithm to use only inner products



$$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}, \quad k(x, z) = \langle \phi(x), \phi(z) \rangle$$

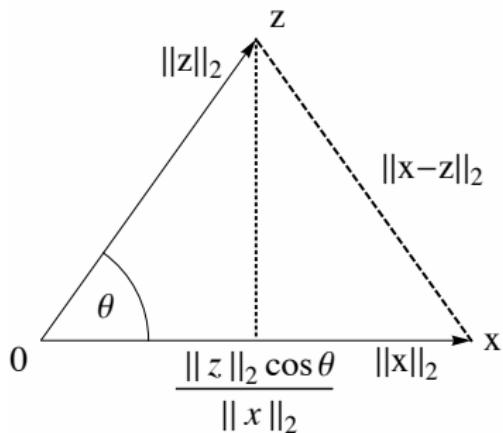
Kernel functions

Kernels correspond to **inner products**.

If $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is symmetric positive semi-definite, then $k(x, z) = \langle \phi(x), \phi(z) \rangle$ for some $\phi : \mathcal{X} \rightarrow \mathcal{H}$.

Inner products encode information about lengths and angles:

$$\|x - z\|^2 = \langle x, x \rangle - 2 \langle x, z \rangle + \langle z, z \rangle, \quad \cos \theta = \frac{\langle x, z \rangle}{\|x\| \|z\|}.$$



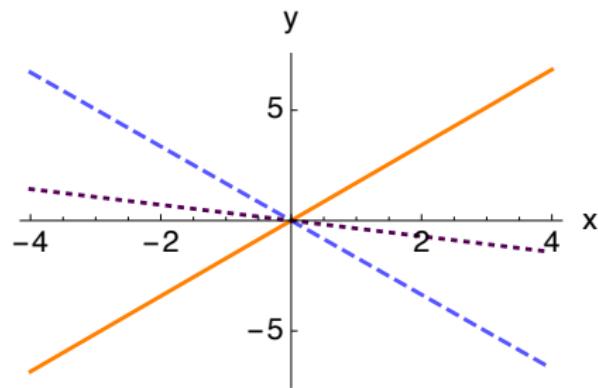
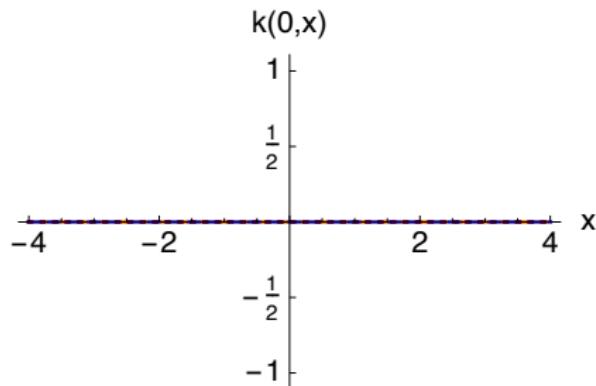
- well characterized function class
- closure properties
- access data only by $K_{ij} = k(x_i, x_j)$
- \mathcal{X} can be any non-empty set

Example: quadratic kernel

→ blackboard

Examples of kernel functions

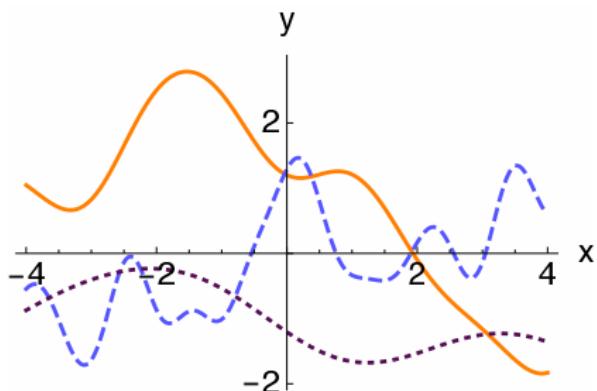
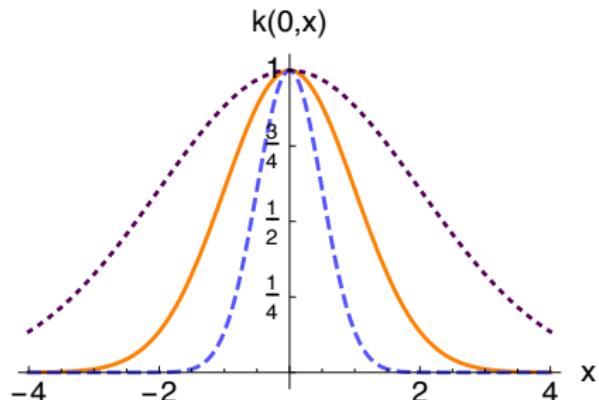
Linear kernel $k(\mathbf{x}, \mathbf{z}) = \langle \mathbf{x}, \mathbf{z} \rangle$



- recovers original linear model

Examples of kernel functions

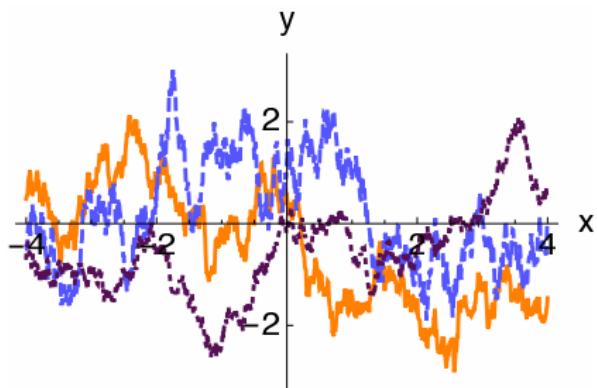
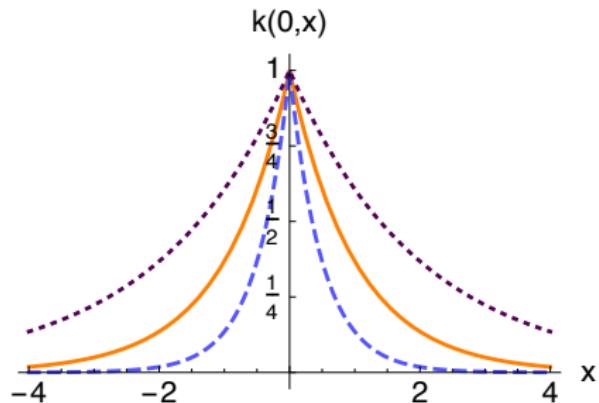
$$\text{Gaussian kernel } k(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2}\right)$$



- length scale σ
- infinite dimensional feature space
- universal local approximator

Examples of kernel functions

$$\text{Laplacian kernel } k(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|_1}{\sigma}\right)$$



- length scale σ

From linear regression to kernel ridge regression

- **linear regression → blackboard**
problem, model form, optimization problem, solution
- **ridge regression → blackboard**
correlated inputs, overfitting, “ridge” penalization, meaning
- **kernel ridge regression → blackboard**
kernel trick, solution

Comparison of linear and kernel ridge regression

Ridge regression

Minimizing

$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2 + \lambda \|\beta\|^2$$

yields

$$\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

for models

$$f(\mathbf{x}) = \sum_{i=1}^d \beta_i \mathbf{x}_i$$

Kernel ridge regression

Minimizing

$$\min_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

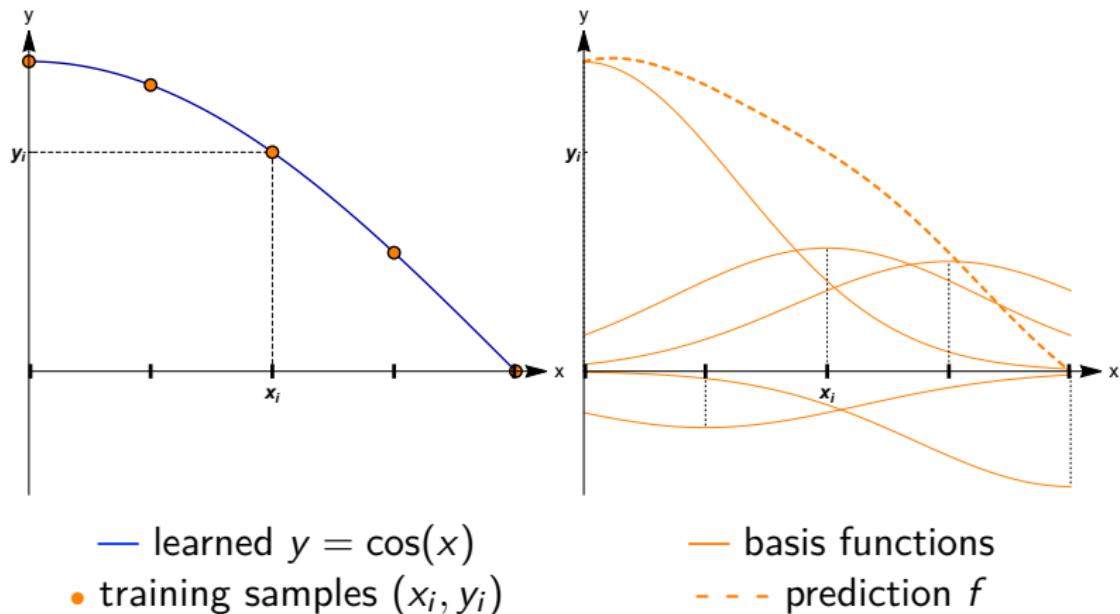
yields

$$\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}.$$

for models

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

The basis function picture



Representer theorem

Kernel models have form

$$f(\mathbf{z}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{z})$$

due to the **representer theorem**:

Any function minimizing a regularized risk functional

$$\ell\left((\mathbf{x}_i, y_i, f(\mathbf{x}_i))_{i=1}^n\right) + g(\|f\|)$$

admits to above representation.

Intuition:

- model lives in space spanned by training data
- weighted sum of basis functions

Centering in kernel feature space

Centering \mathbf{X} and \mathbf{y} is equivalent to having a bias term b .

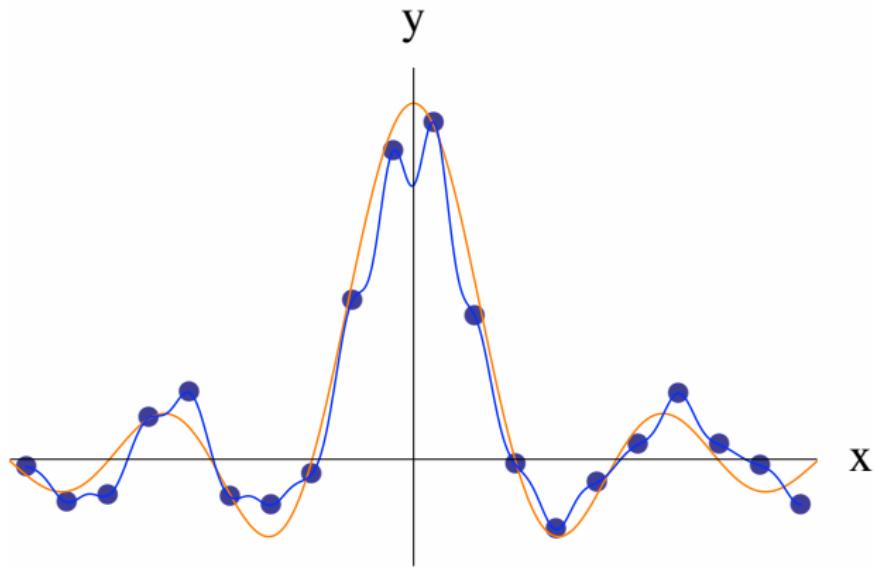
For kernel models, center in kernel feature space:

$$\begin{aligned}\tilde{k}(\mathbf{x}, \mathbf{z}) &= \left\langle \phi(\mathbf{x}) - \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i), \phi(\mathbf{z}) - \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \right\rangle \\ &\Rightarrow \tilde{\mathbf{K}} = \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \right) \mathbf{K} \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \right)\end{aligned}$$

Some kernels like Gaussian and Laplacian kernels do not need centering
Poggio *et al.*, Tech. Rep., 2001

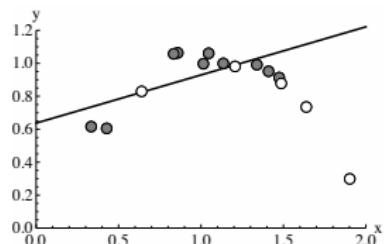
Model building

How regularization helps against overfitting



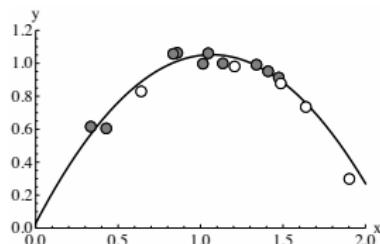
Effect of regularization

Underfitting



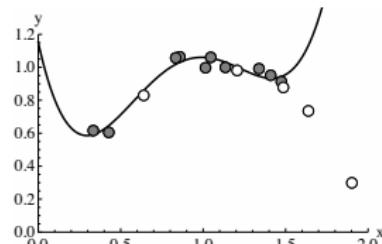
0.123 / 0.443

Fitting

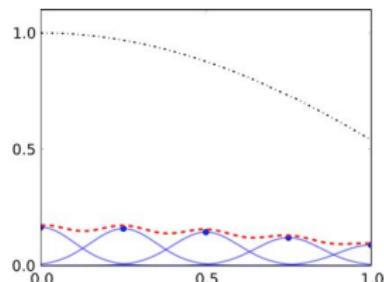


0.044 / 0.068

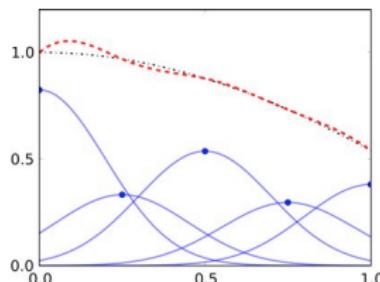
Overfitting



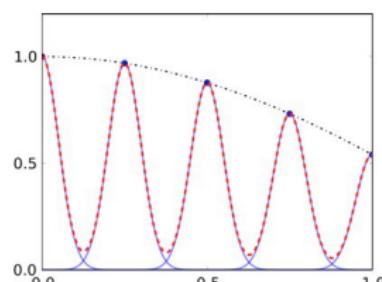
0.036 / 0.939



λ too large

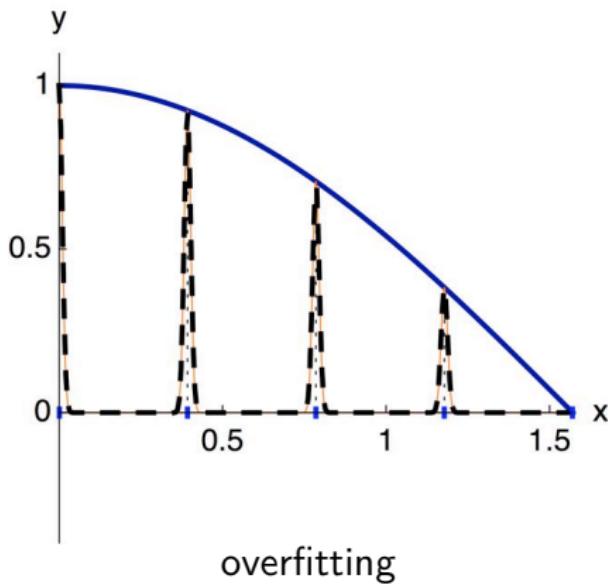
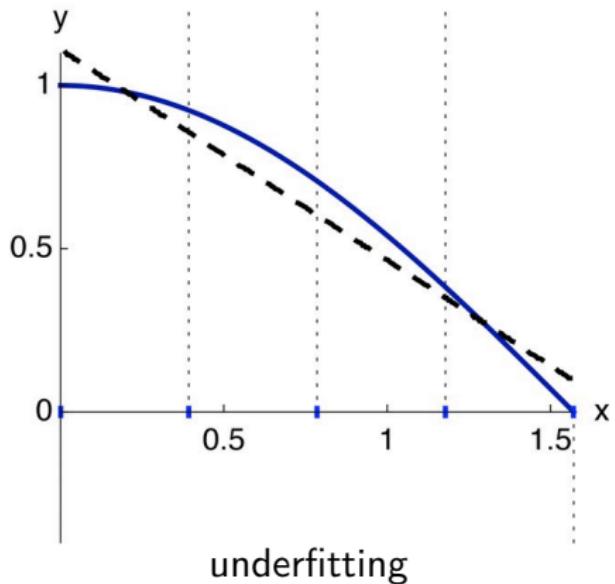


λ right



λ too small

Overfitting and underfitting in the limit



Validation

Why?

- assess model performance
- optimize free parameters (hyperparameters)

Which statistics?

- root mean squared error (RMSE)
- mean absolute error (MAE)
- maximum error
- squared correlation coefficient (R^2)

What else can we learn from validation?

- distribution of errors, not only summary statistics
- convergence of error with number of samples

Validation

Golden rule:

Never use training data for validation

Violation of this rule leads to overfitting
by measuring flexibility in fitting instead of generalization ability
rote learner example

If there is sufficient data:

- divide data into two subsets, training and validation
- build model on training subset
- estimate error of trained model on validation subset

Sometimes an external validation set is used in addition.

Statistical validation

If too few data, statistical re-sampling methods can be used, such as cross-validation, bagging, bootstrapping, jackknifing

k-fold cross-validation:

- divide data into k evenly sized subsets
- for $i = 1, \dots, k$,
build model on union of subsets $\{1, \dots, k\} \setminus \{i\}$
and validate on subset i

All model building steps must be repeated for data splits:

- all pre-processing such as feature selection and centering
- optimization of hyperparameters

Hyperparameters: physically motivated choices

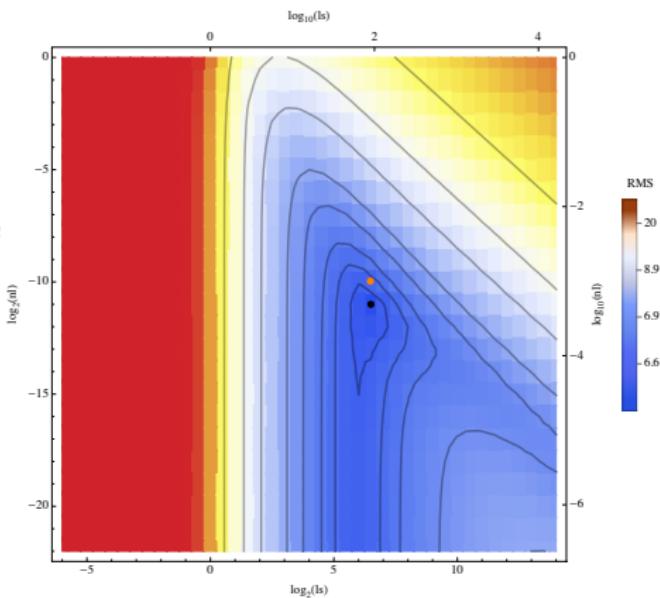
Length scale σ :

$$\sigma \approx \|x - z\|_1$$

median nearest neighbor distance

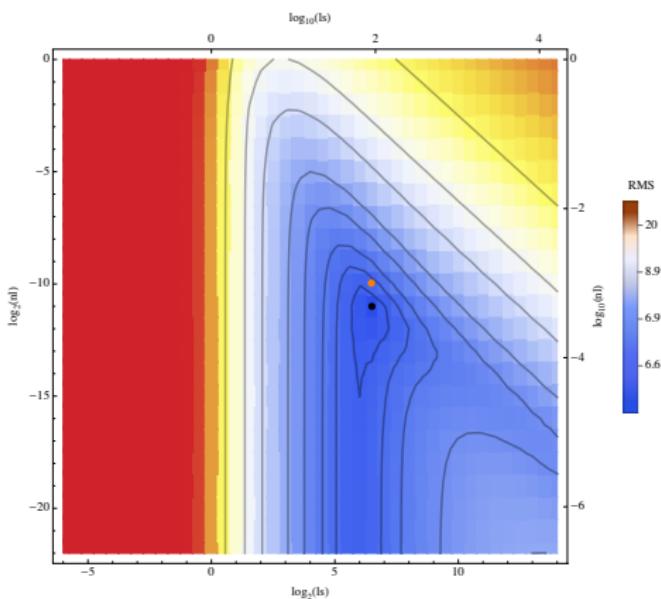
Regularization strength λ :

- $\hat{=}$ noise variance (Bayesian)
- $\hat{=}$ leeway around y_i for fitting
- \Rightarrow target accuracy



Hyperparameters: statistically motivated choices

- data-driven method for choosing hyperparameters
- optimize using grid search or gradient descent
- use statistical validation to estimate error
- for validation and hyperparameter optimization, use nested data splits



Nested data splits

- **never use data from training in validation**
- for performance assessment **and** hyperparameter optimization, use nested cross-validation or nested hold-out sets
- beware of overfitting

Example 1: plain overfitting

- ✗ train on all data, predict all data
- ✓ split data, train, predict

Example 2: centering

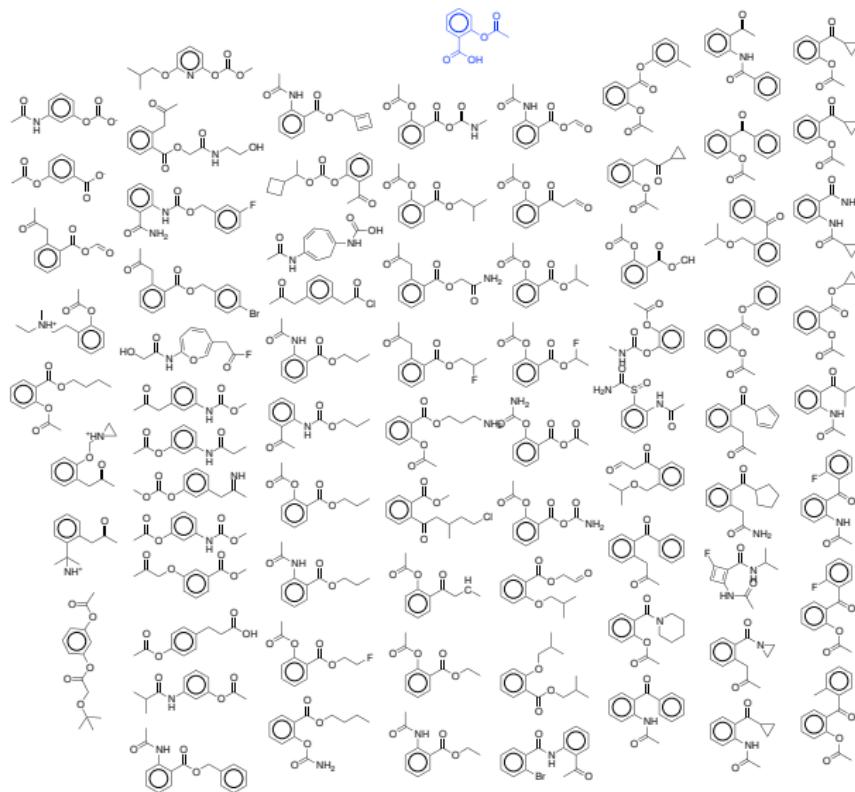
- ✗ center data, split data, train & predict
- ✓ split data, center training set, train, center test set, predict

Example 3: cross-validation with feature selection

- ✗ feature selection, cross-validation
- ✓ feature selection for each split of cross-validation

Applications

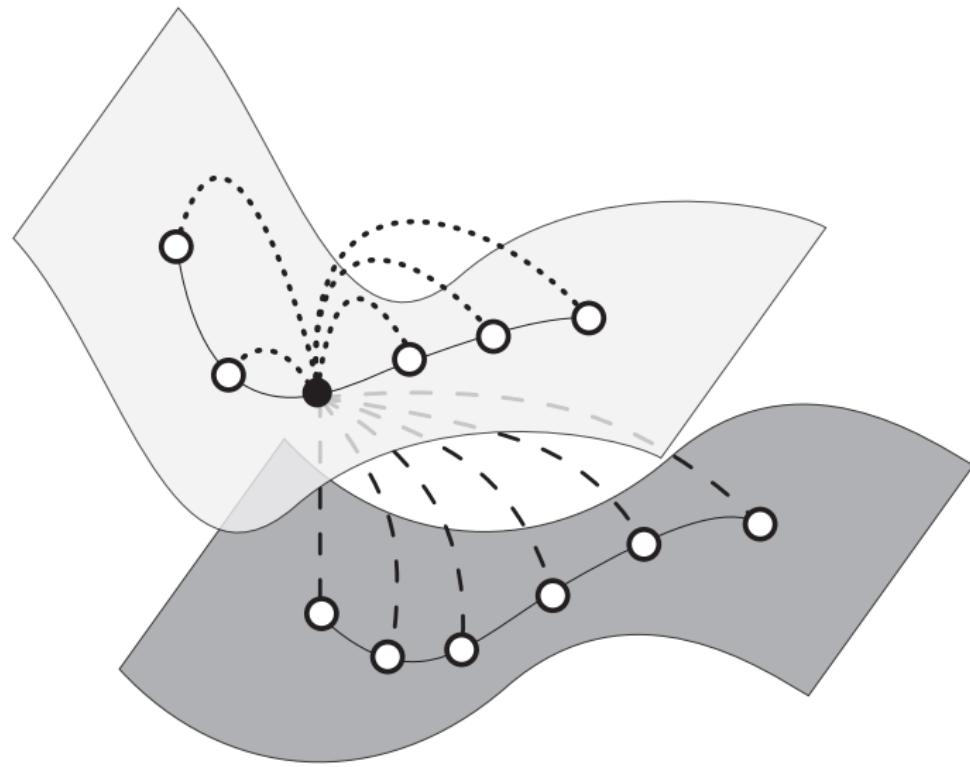
The combinatorial nature of chemical/materials space



- molecule space:
graph theory
- materials space:
group theory
- combinatorial
explosion

aspirin derivatives

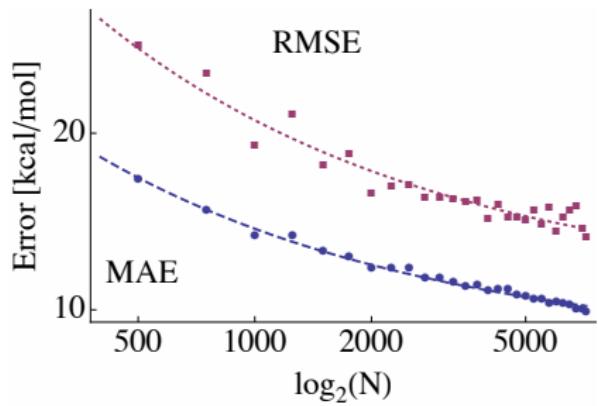
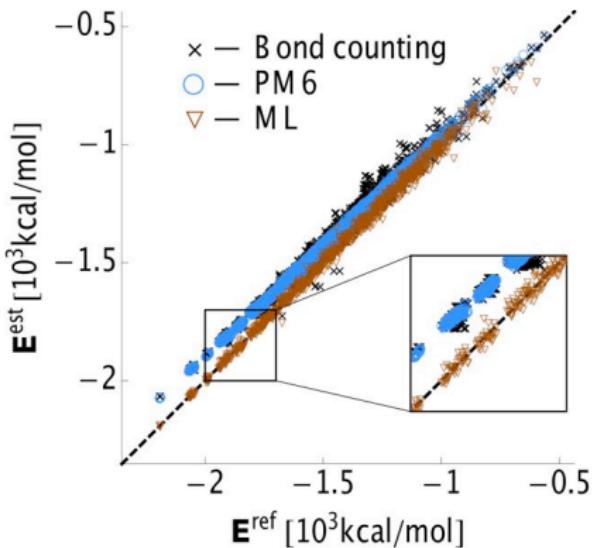
Learning across chemical space



Chang, von Lilienfeld: CHIMIA 68, 602, 2014
von Lilienfeld, *Int. J. Quant. Chem.* 113, 1676, 2013.

Predicting atomization energies

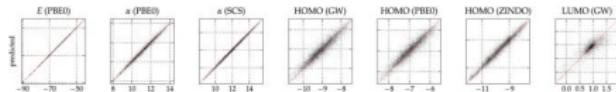
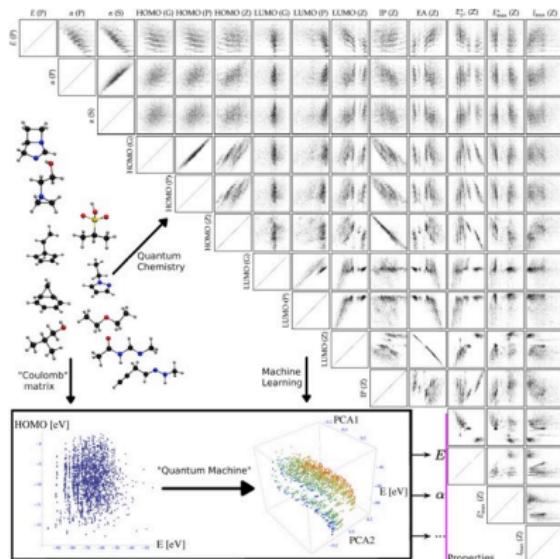
- 7 165 small organic molecules (H,C,N,O,S; 1–7 non-H atoms)
- DFT PBE0 atomization energies
- kernel ridge regression, Gaussian kernel $k(\mathbf{M}, \mathbf{M}') = \exp\left(-\frac{d^2(\mathbf{M}, \mathbf{M}')}{2\sigma^2}\right)$



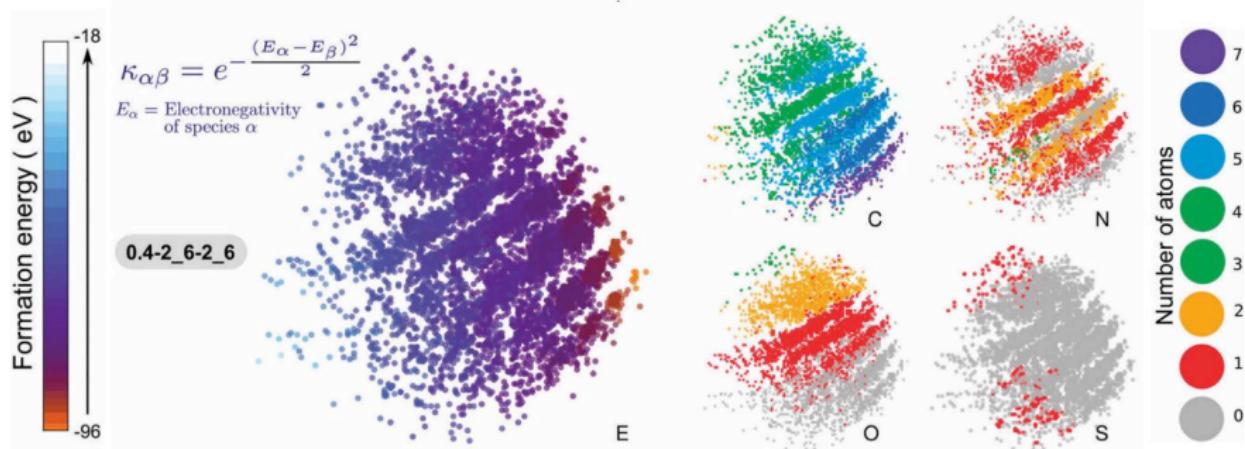
Extension to other properties

Learning the map from molecular structure to molecular properties

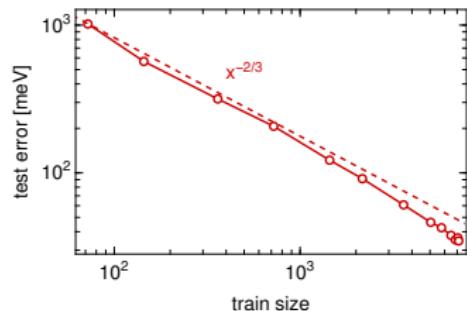
- various properties
- various levels of theory
- small organic molecules
- Coulomb matrix representations
- kernel learning,
deep neural networks
- for 5k training molecules, errors
are comparable to the reference



Molecular energies

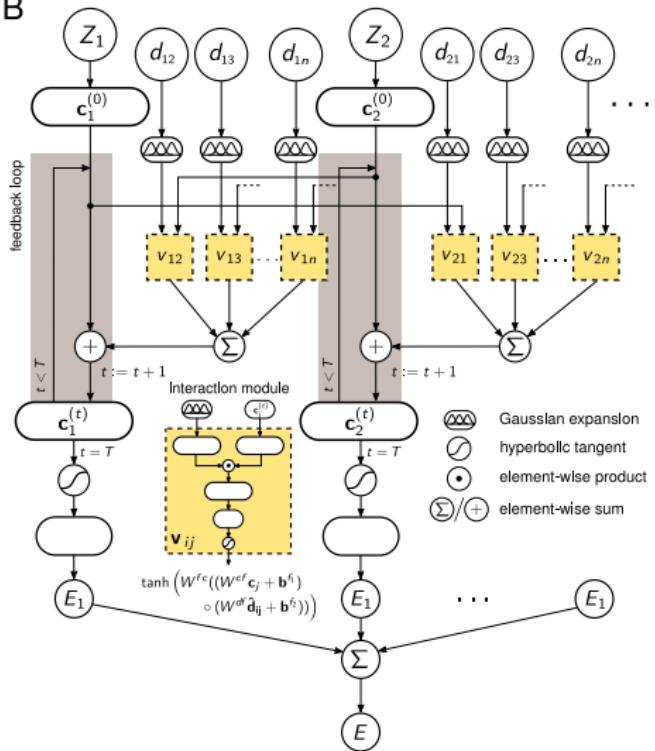


- Gaussian process regression
- regularized entropy match kernel (Sinkhorn distance) with smooth overlap of atomic positions representation
- MAE = 0.6 kcal/mol, RMSE = 0.9 kcal/mol

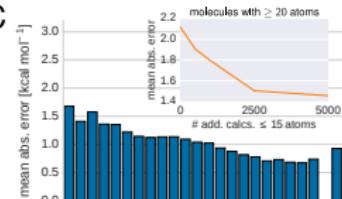


Deep tensor neural networks

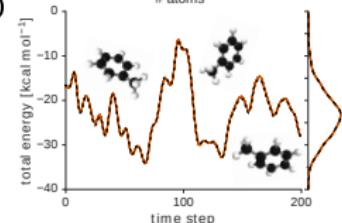
B



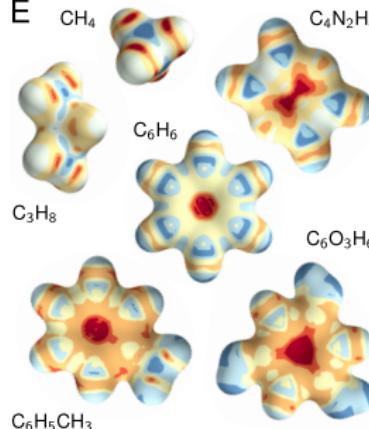
C



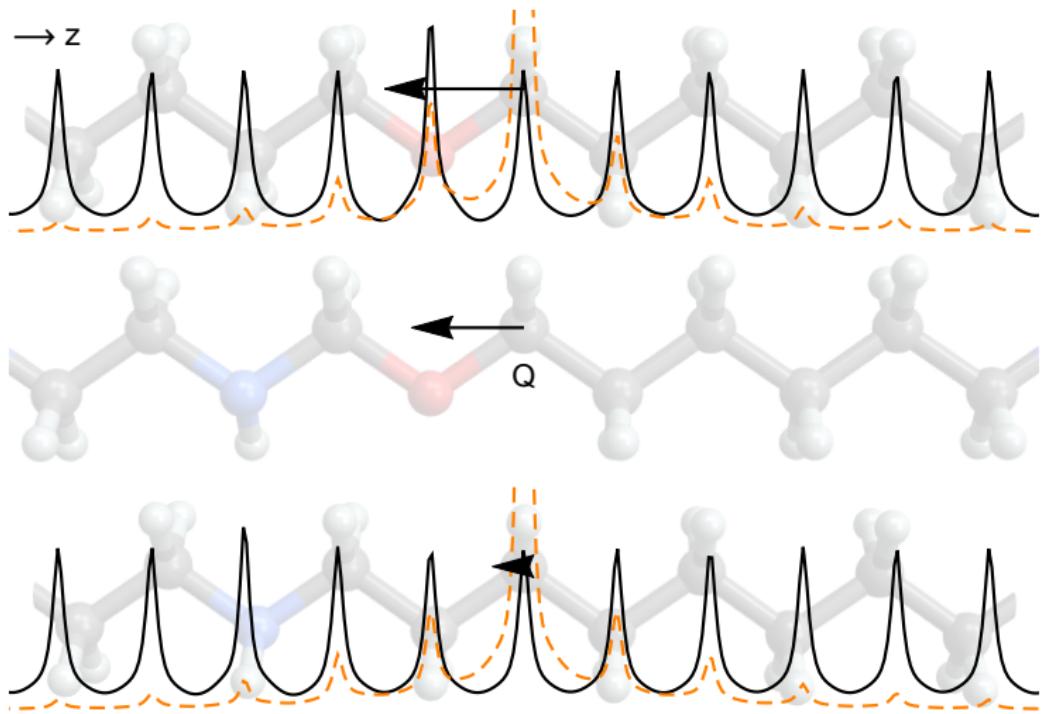
D



E

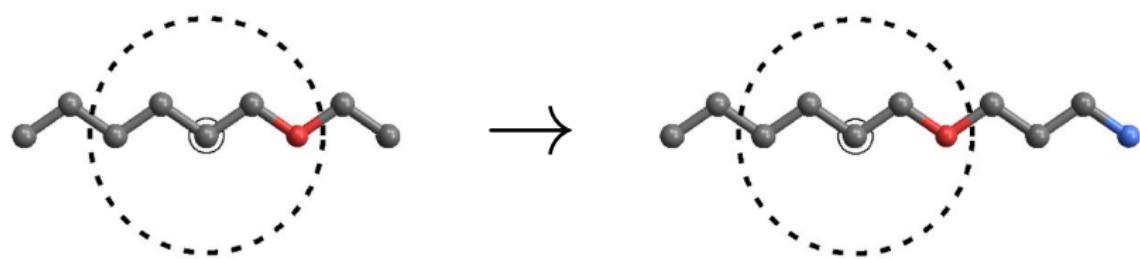


Local properties



Local properties

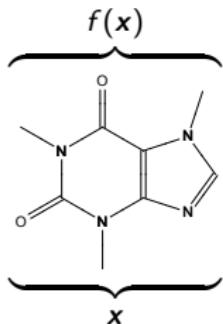
Local interpolation is global extrapolation.



- **linear scaling** of computational effort with system size
- size consistent in the limit
- requires **partitioning** for global properties

Local properties

Molecular model

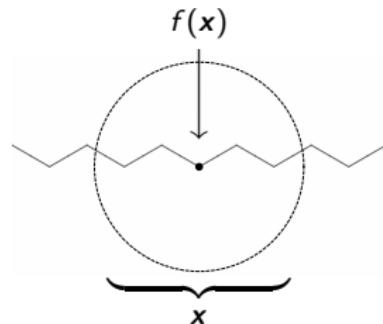


$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

n = number of molecules

\mathbf{x} = representation of molecule

Atomic model

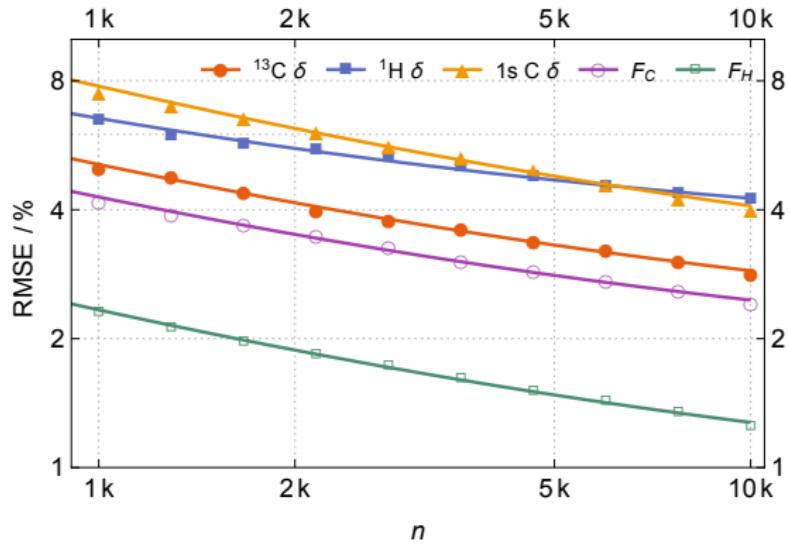


$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

n = number of atoms

\mathbf{x} = representation of atom

Local properties



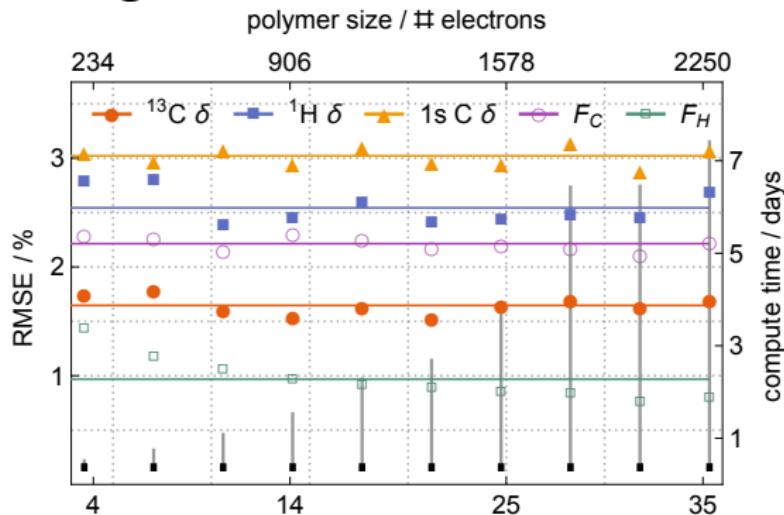
Property	Ref.	RMSE	maxAE	R^2
$^{13}\text{C} \delta / \text{ppm}$	2.4	5.8 ± 0.3	36 ± 8	0.988 ± 0.001
$^1\text{H} \delta / \text{ppm}$	0.11	0.42 ± 0.02	3.2 ± 1.1	0.954 ± 0.005
$1\text{s C} \delta / \text{mE}_h$	7.5	6.5 ± 0.3	34 ± 17	0.971 ± 0.002
$F_C / \text{mE}_h / a_0$	1	4.7 ± 0.15	29 ± 5.5	0.983 ± 0.002
$F_H / \text{mE}_h / a_0$	1	1.1 ± 0.03	7.4 ± 2.6	0.996 ± 0.003

Local properties

Dataset

- linear polyethylene ($\text{CH}_2\text{CH}_2)_n$, doped with N and O
- varying length in multiples of basic unit (29 non-H atoms)
- DFT / PBE0 / def2TZVP using Gaussian 09

Scaling



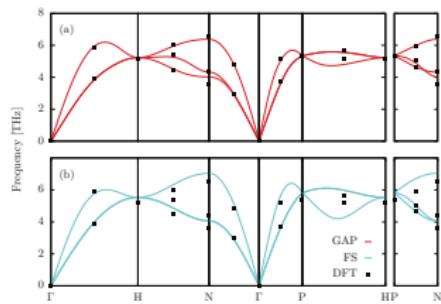
- training on shortest polymers only
- prediction of polymers of increasing size (up to x10)

Molecular dynamics—adsorption on surfaces

- since early 1990s > 35 studies on molecules
- many studies using artificial neural networks for potential energy surface interpolation

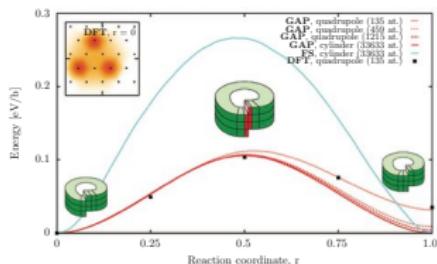
Year	Ref.	System	Reference method
1995	Blank <i>et al</i> [77]	CO @ Ni(111)	empirical PES
1995	Blank <i>et al</i> [77]	H ₂ @ Si(100)-(2×1)	DFT (LDA)
2004	Lorenz <i>et al</i> [223]	H ₂ @ K(2×2)/Pd(100)	DFT (PW91)
2005	Behler <i>et al</i> [123, 224, 225]	O ₂ @ Al(111)	DFT (RPBE)
2006	Lorenz <i>et al</i> [226]	H ₂ @ Pd(100)	empirical PES
2006	Lorenz <i>et al</i> [226]	H ₂ @ (2×2)S/Pd(100)	empirical PES
2007	Ludwig and Vlachos [227]	H ₂ @ Pt(111)	empirical PES
2007	Ludwig and Vlachos [227]	H ₂ @ Pt(111)	DFT (PW91)
2008	Behler <i>et al</i> [225]	O ₂ @ Al(111)	DFT (PBE)
2008	Latino <i>et al</i> [228]	ethanol @ Au(111)	DFT (B3LYP)
2008	Carbogno <i>et al</i> [229]	O ₂ @ Al(111)	DFT (RPBE)
2009	Manzhos <i>et al</i> [97]	N ₂ O @ Cu(100)	DFT
2009	Carbogno <i>et al</i> [230]	O ₂ @ Al(111)	DFT (RPBE)
2010	Latino <i>et al</i> [231]	ethanol @ Au(111)	DFT (B3LYP)
2010	Manzhos and Yamashita [232]	N ₂ O @ Cu(100)	DFT
2012	Goikoetxea <i>et al</i> [233]	O ₂ @ Ag(111)	DFT (PBE)
2013	Liu <i>et al</i> [234]	HCl @ Au(111)	DFT (PW91)

Molecular dynamics—tungsten

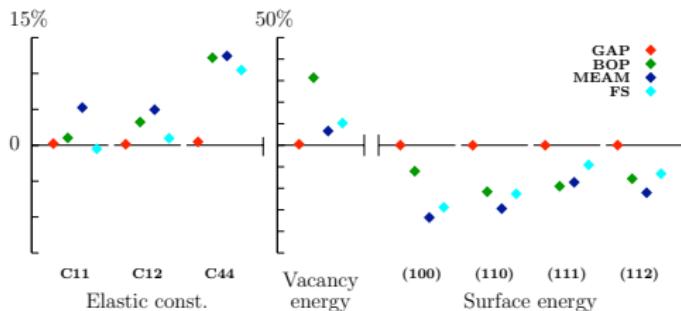


Phonon spectrum

- tungsten in bcc crystal phase
- Gaussian approximation potential
- DFT (PBE, plane waves, pseudopotentials) reference
- screw dislocation



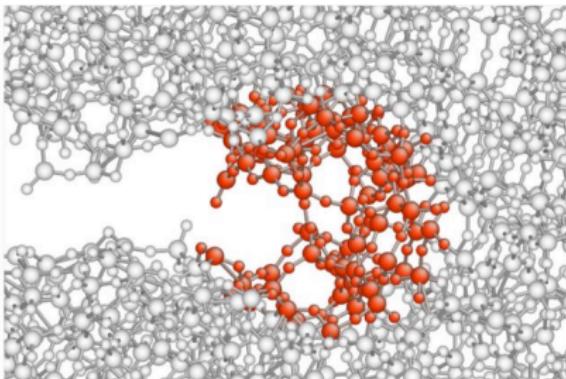
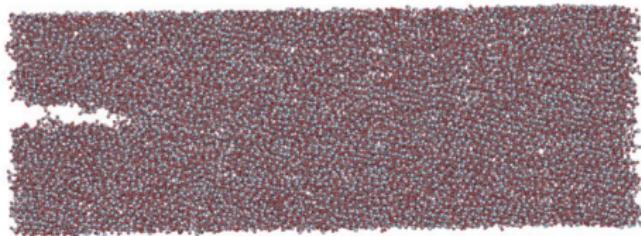
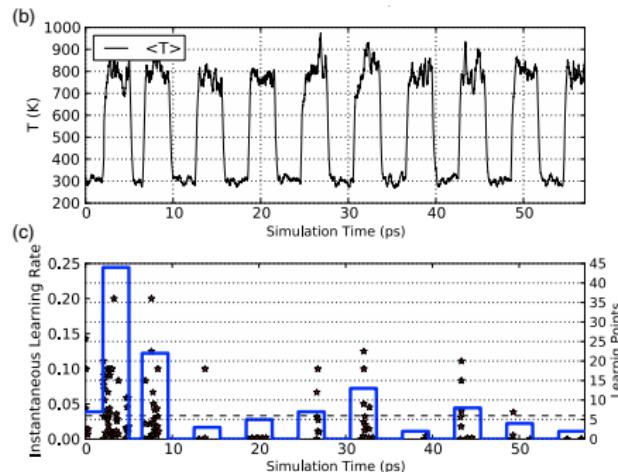
Peierls barrier



Errors on properties

Molecular dynamics—crack propagation

- crack propagation in silicon
- learning on the fly (model is updated when leaving domain)
- form of active learning
- k -step predictor/corrector

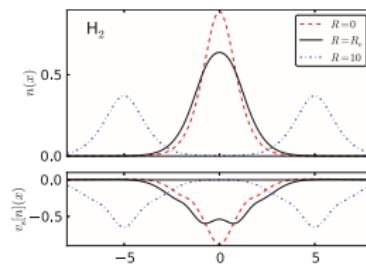


Caccin et al, *Int J Quant Chem* 115: 1129, 2015.
Li et al, *Phys Rev Lett* 114: 096405, 2015.

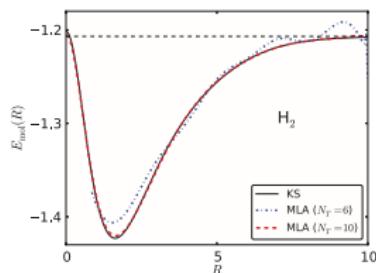
Density functional theory

Learning the map from electron density to kinetic energy

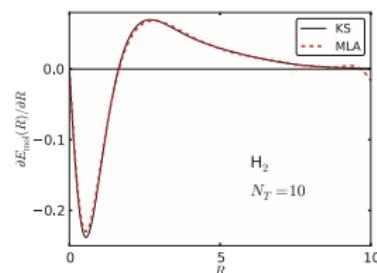
- orbital-free DFT
- 1D toy system
- DFT/LDA as reference
- error decays to zero
- self-consistent densities
- bond breaking and formation



H_2 potential



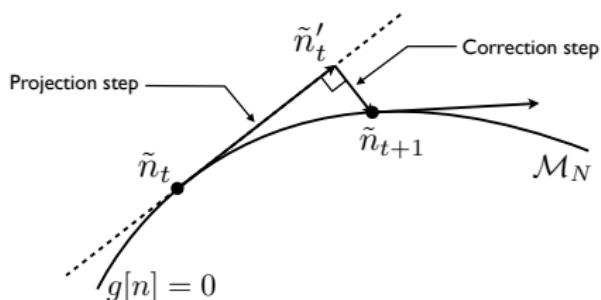
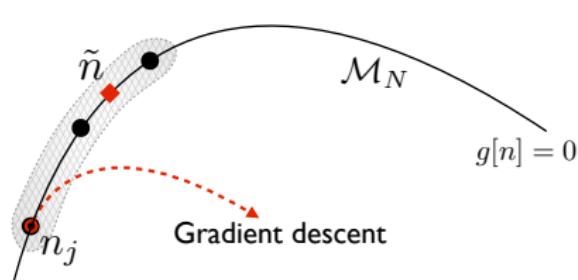
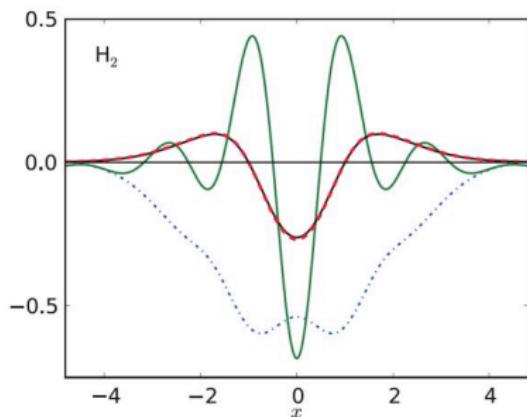
H_2 binding curve



H_2 forces

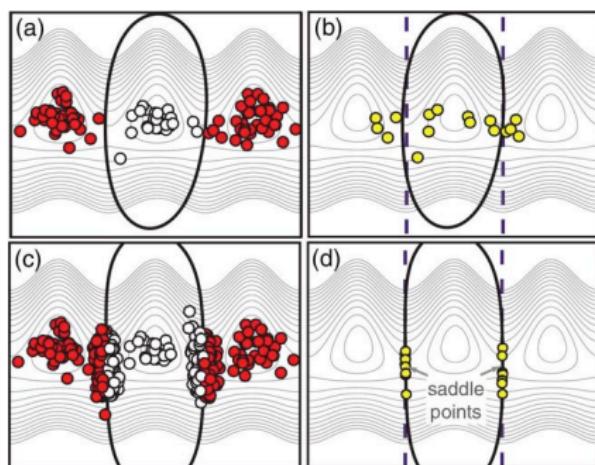
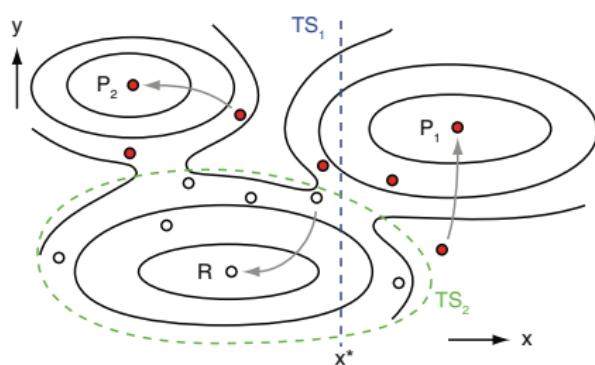
Electron densities—projected gradients

- kinetic energy of electron densities
- Gaussian process regression
- orbital-free DFT, 1D toy system
- error decays to zero
- projected gradients for self-consistent densities (“non-linear gradient denoising”)



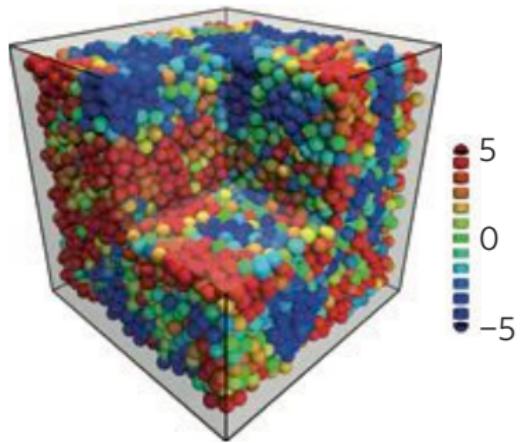
Transition state theory

- characterization of dividing surfaces
- support vector machine for classification
- alternate between learning and sampling
- no prior information required
- iteratively refined by biased sampling along dividing surface

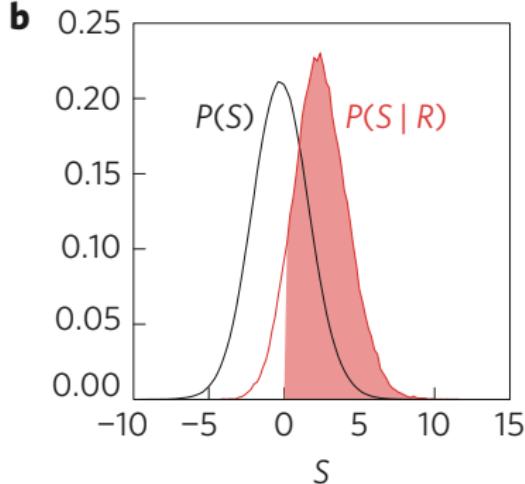


Relaxation in glassy liquids

a

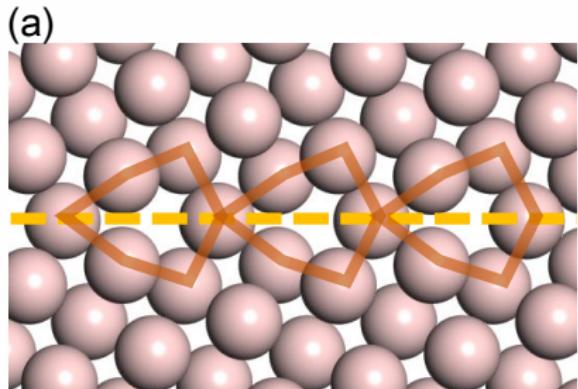


b



- identify subtle structural changes (“softness”) in glassy dynamics
- softness correlates to probability of rearrangement in near future

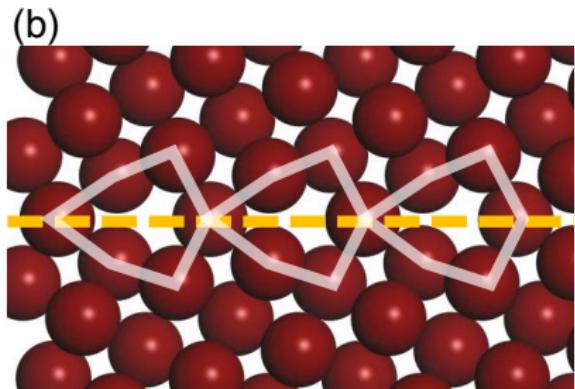
Stable interface search



Exhaustive calculations

GB energy=0.96J/m²

Number of energy calculations
=16,983



Bayesian optimization

GB energy=0.96J/m²

Number of energy calculations
=69

Summary

- machine learning finds regularity in data for analysis or prediction, improving with more data
- the kernel trick for implicit transformation to high-dimensional spaces
- for validation, avoid over-fitting by following the golden rule
- examples of predicting computational outcomes

Tutorial

Matthias Rupp:

Machine Learning for Quantum Mechanics in a Nutshell

International Journal of Quantum Chemistry 15(16): 1058–1073, 2015

<http://doi.org/10.1002/qua.24954>

Links

<http://mrupp.info> (Publications)

<http://qmml.org> (Datasets)