



EFFICIENT MACHINE LEARNING FROM SPARSE DATA IN HIGH- D FEATURE SPACES

Sergei Manzhos

Ihara-Manzhos Laboratory

<https://sites.google.com/site/sergeimanzhos/>

<http://www.chemeng.titech.ac.jp/~iharalab/>

School of Materials and Chemical Technology



ML IS A VAST AND BRANCHED FIELD...



- We will focus on regression problems
- And methods *based on* neural networks and kernel methods

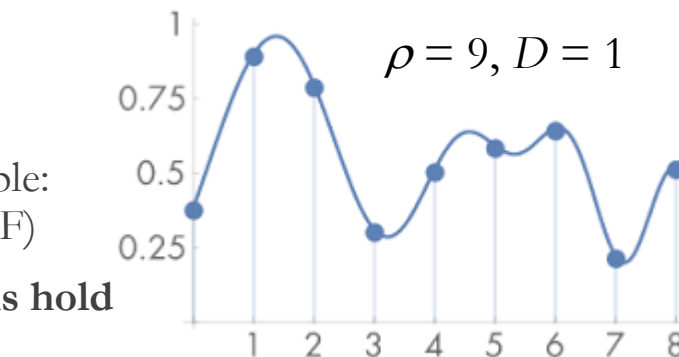
OFF-THE-SHELF ML METHODS WORK WELL IN MODERATELY HIGH-D WITH MODERATELY SPARSE DATA

- Curses of dimensionality

- Exponentially difficult to fill space: $N_{\text{data}} = \rho^D$
- Data are **always** sparse in multidimensional spaces (example: 10^6 points in 20D ~ 2 points/DOF, $10^7 \sim 2.2$ points/DOF)
- Exponential growth of required no. of data and terms hold for direct product type representations (think a regular grid / Fourier expansion)

- ML methods avoid them

- What is local in low-d may be non-local in high-D:
1d: 68% of quadrature within σ ; 6d: 10% of quadrature



$$f(\mathbf{x}) = \prod_{i=1}^D (2\pi\sigma_i^2)^{-\frac{1}{2}} \exp\left(-\frac{(x_i - \bar{x}_i)^2}{2\sigma_i^2}\right)$$

Consequential for
kernel regressions

- Possible blessings of dimensionality

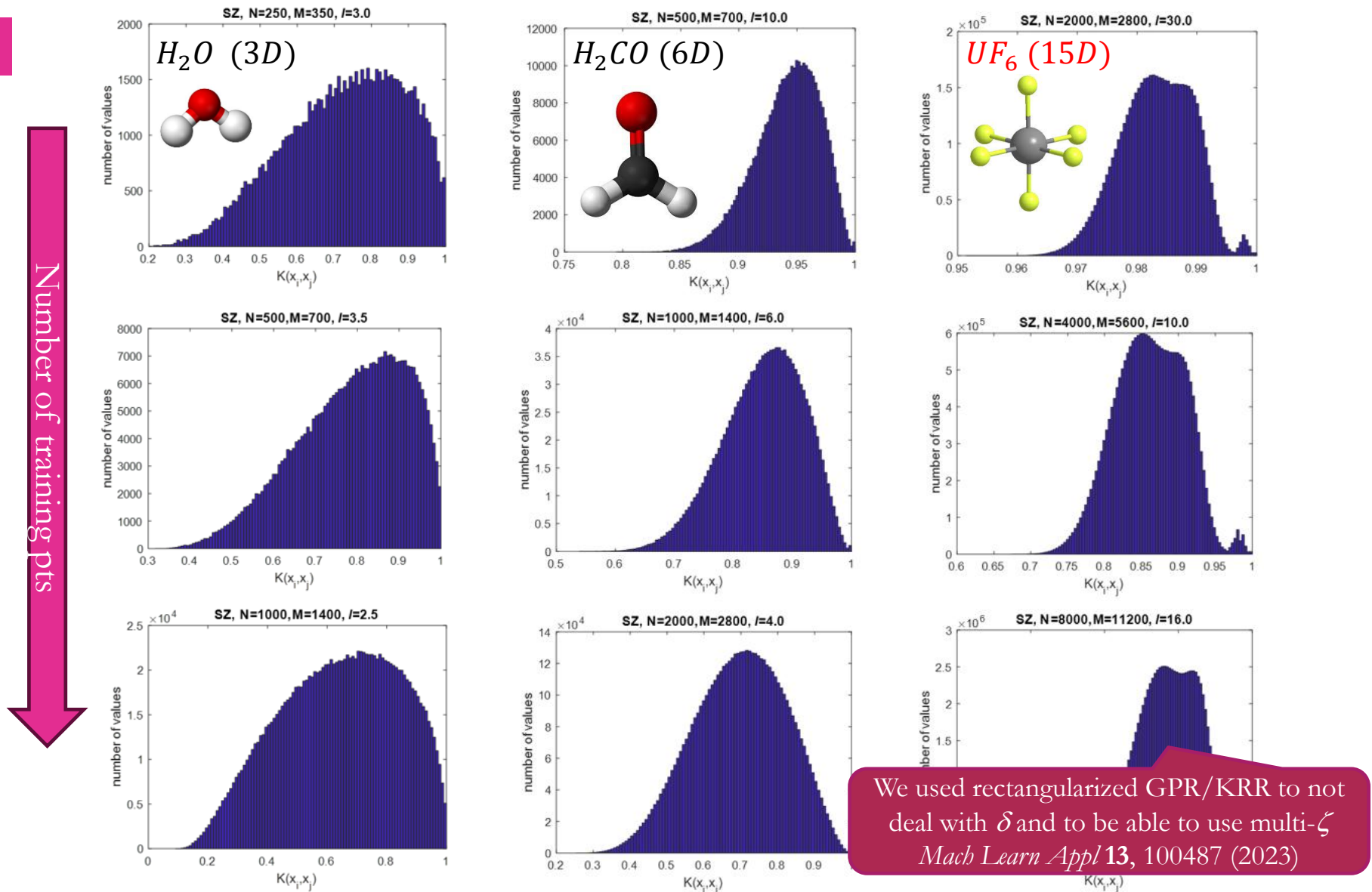
- Concentration of metrics:

$$\lim_{D \rightarrow \infty} E\left(\frac{\text{dist}_{\max}(D) - \text{dist}_{\min}(D)}{\text{dist}_{\min}(D)}\right) \rightarrow 0$$

- Some ML methods work better in high than in low dimensionality (e.g. NN)

- But in very high-D or with very low density of data specialized ML methods are desired

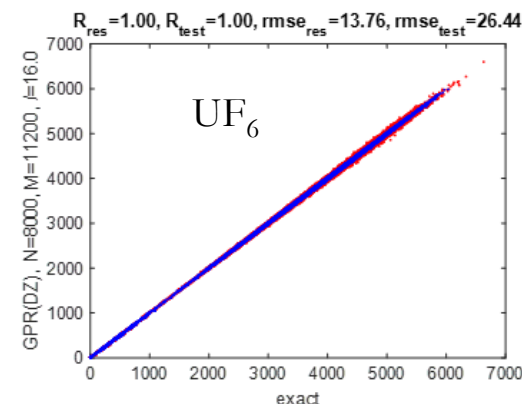
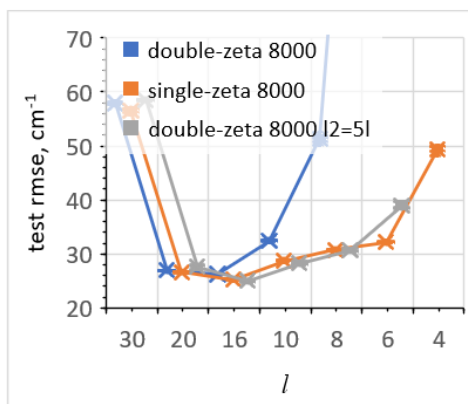
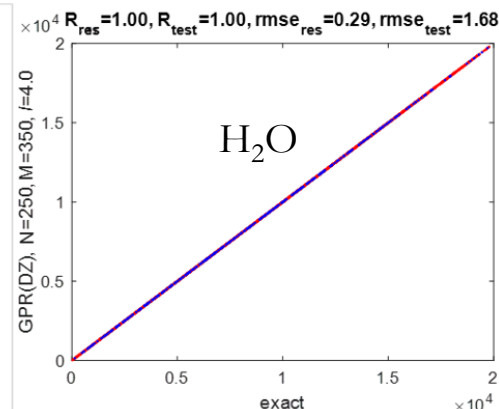
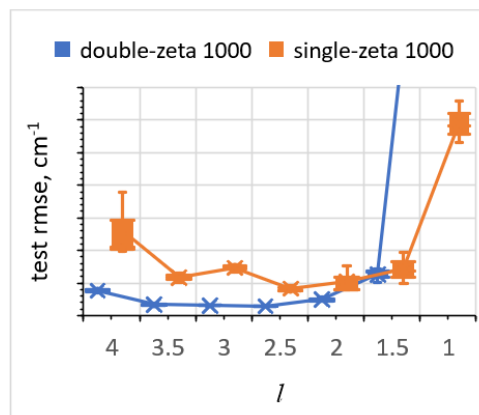
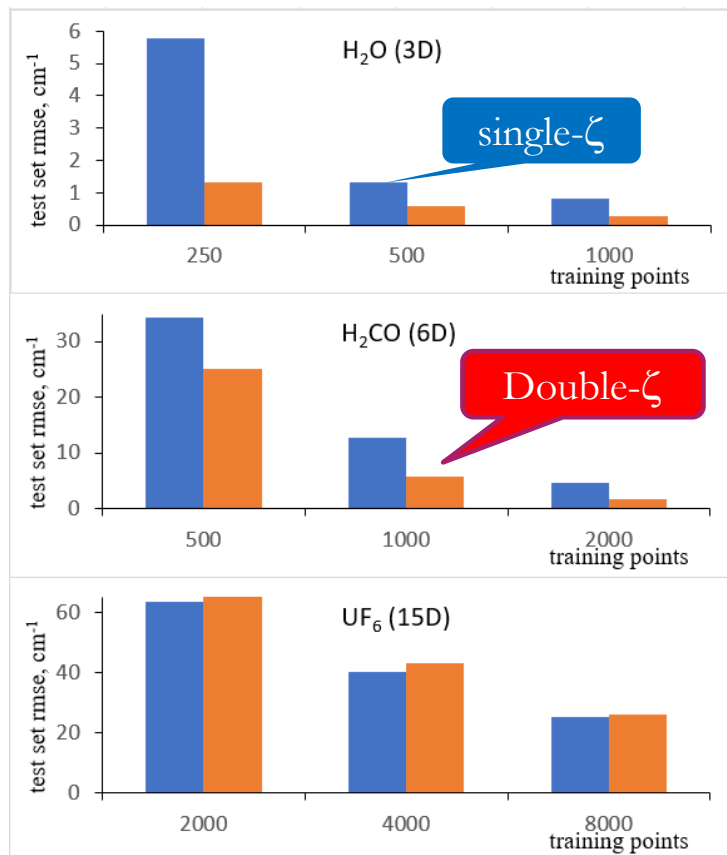
HOW HIGH-D MATERN KERNELS DIE: SPARSE DATA



■ Matern type kernels lose the property of locality...

J. Chem. Phys. **158**, 044111 (2023)

ADVANTAGE OF MULTI- ζ () DIES WHEN LOCALITY IS GONE



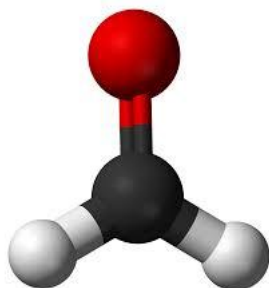
Error in, cm ⁻¹	PES	Δ ZPE	Lowest 50 transitions		Lowest 100 transitions	
			mae	rmse	mae	Rmse
Single-zeta kernel	4.60	0.077	0.138	0.166	0.139	0.165
Double-zeta kernel	1.67	-0.026	0.013	0.016	0.014	0.019

Double- ζ kernel PES representation useful for properties when it is meaningful (H₂CO, 2000 kernel centers)

HOW HIGH-D KERNELS DIE: SPARSE DATA

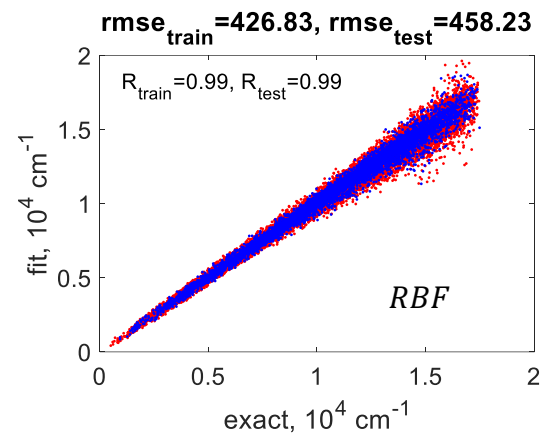
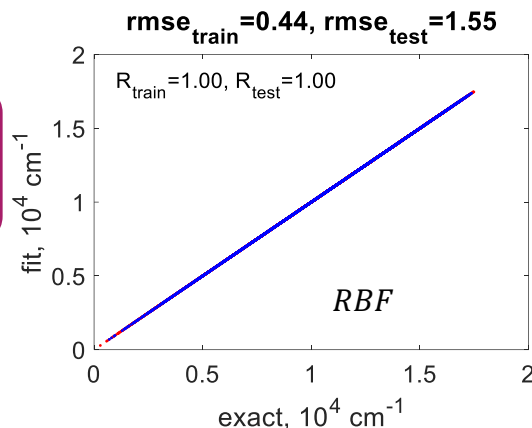
5,600 train pts
20,000 test pts

$l_{opt} \sim 5$

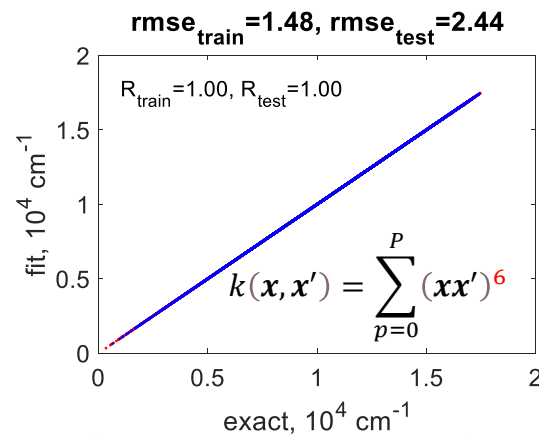
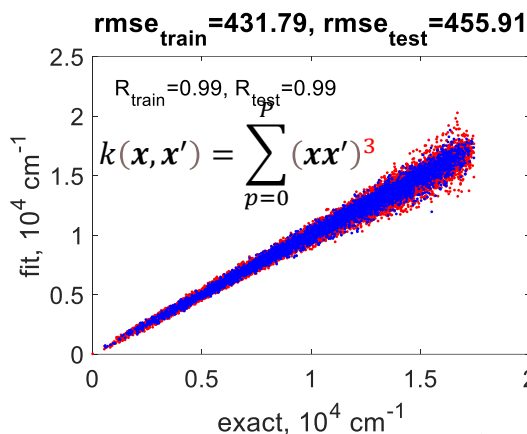


Product kernel =
polynomial
regression

$$k(\mathbf{x}, \mathbf{x}') = \chi + (\mathbf{x}\mathbf{x}')^2:$$



$l=50$



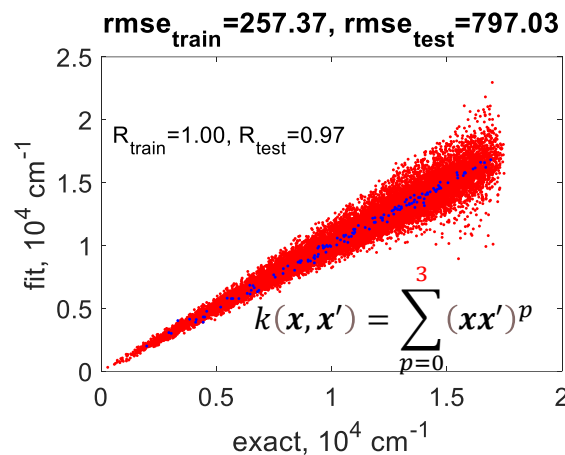
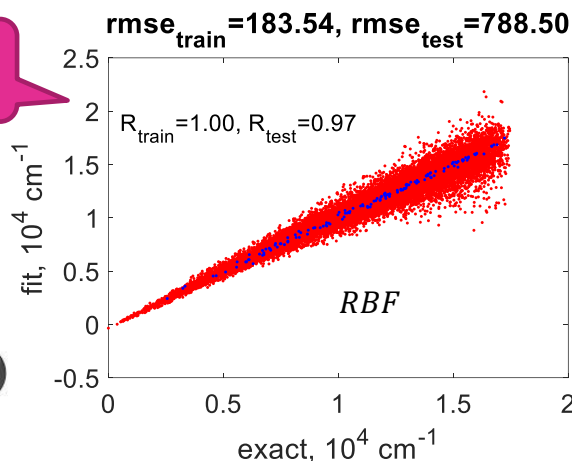
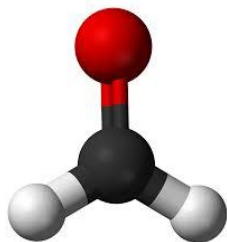
$$f(\mathbf{x}) = \left(\chi \sum_{n=1}^N c_n \right) + \sum_{d=1}^D \left(\sum_{n=1}^N c_n (x_d^{(n)})^2 \right) x_d^2 + \sum_{d \neq d'}^D \left(\sum_{n=1}^N c_n x_d^{(n)} x_{d'}^{(n)} \right) x_d x_{d'}$$

$$= b_0 + \sum_{d=1}^D b_d^{(2)} x_d^2 + \sum_{d \neq d'}^D b_{dd'} x_d x_{d'} \quad \text{J. Chem. Phys., 160, 012201(2024)}$$

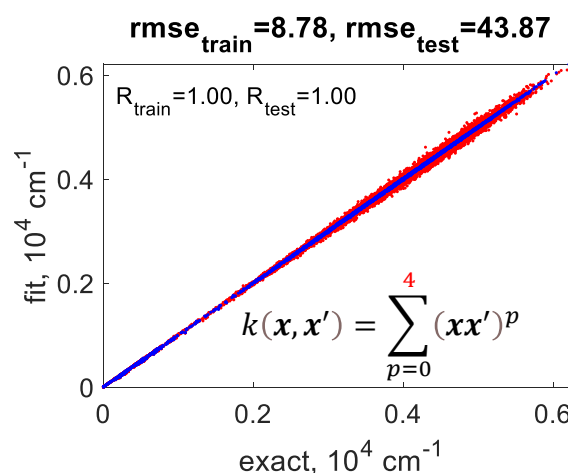
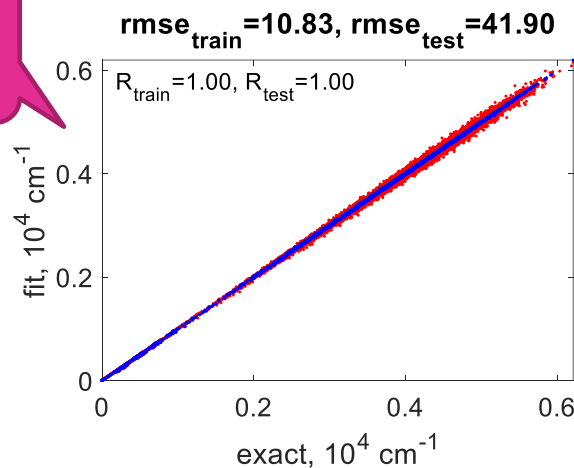
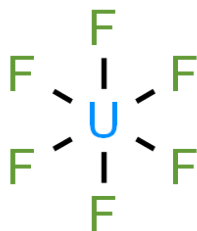
- ...and may degenerate into a low-order polynomial basis

HOW HIGH-D MATERN KERNELS DIE: SPARSE DATA

150 train pts
 $l_{opt} \sim 15$



15D
5600 train pts
 $l_{opt} \sim 30$



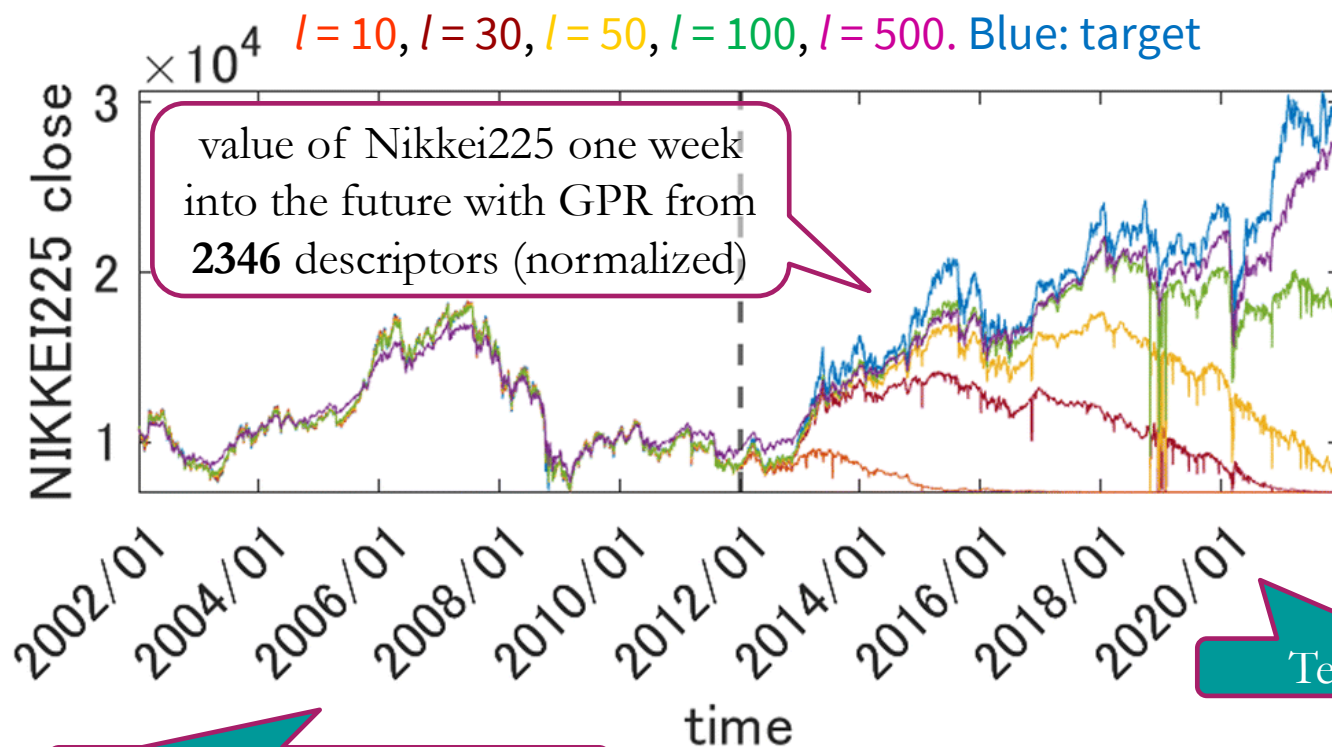
■ Higher D → lower data density → higher l_{opt}

J. Chem. Phys., **160**, 012201(2024)

FAILURE OR COLLAPSE OF GPR WITH MATERN KERNELS IN VERY HIGH-D

EXAMPLE FROM QUANTITATIVE FINANCE

$$y(t + \Delta t) = f(x_1(t), x_2(t), \dots, x_D(t)) \quad \exp\left(-\frac{|\mathbf{x} - \mathbf{x}'|^2}{l^2}\right) = \prod_d^D \exp\left(-\left(\frac{x_d - x'_d}{l}\right)^2\right)$$



A product of very many functions each of value < 1

Test: 2012-2020

Training: 2002-2011

GPR collapses unless l is so large it loses any advantage over simple linear regression

HIGH-DIMENSIONAL MODEL REPRESENTATION

(HDMR)

$$x \in R^D; f(x) \approx f_0 + \sum_{i=1}^D f_i(x_i) + \sum_{i=1, j=i+1}^D f_{ij}(x_i, x_j) + \dots$$

1st order additive

$$+ \sum_{i_1, i_2, \dots, i_d}^D f_{i_1, i_2, \dots, i_d}(x_{i_1}, x_{i_2}, \dots, x_{i_d}) + \dots$$

In some form dates back 100 years (ANOVA)

HDMR is formalized in a series of papers by Rabitz's group at Princeton
J Math Chem A **1999**, 25, 197 etc

$C_d(D)$ component functions at each d

Particular cases: physics – MBE, ...

... computational spectroscopy: N-mode approximation

$$E = \sum_{I=1}^{N_b} E_I + \sum_I^{N_b} \sum_{J>I}^{N_b} \Delta E_{IJ} + \sum_I^{N_b} \sum_{J>I}^{N_b} \sum_{K>J}^{N_b} \Delta E_{IJK} + \dots$$

Can consider only highest- d terms:

$$f(x) \approx \sum_{i_1, i_2, \dots, i_d}^D f_{i_1, i_2, \dots, i_d}(x_{i_1}, x_{i_2}, \dots, x_{i_d})$$

Comput Phys Commun **2009**, 180, 2002

J Phys Chem A **2020**, 124, 7598

Comput Phys Commun **2022**, 271, 108220

- Formalization of **representation with lower-dimensional functions (many-body, N-mode)**
 - Easier to build: component functions are lower-dimensional, could be built with fewer data w/out overfitting (e.g. *J Math Chem* **61** (2023) 7)
 - Fitting method can work in its comfort zone
 - Easier to use (integrals etc)
 - Provides elements of insight (see e.g. *Comput Phys Commun* **2022**, 271, 108220)

HDMR-ML: A POWERFUL TOOL TO WORK WITH SPARSE DATA

Review: *Artificial Intelligence Chemistry* **1**, 100008 (2023)

$$f(\mathbf{x}) \approx f_0 + \sum_{i=1}^D f_i(x_i) + \sum_{i=1, j=i+1}^D f_{ij}(x_i, x_j) + \cdots + \sum_{i_1, i_2, \dots, i_d}^D f_{i_1, i_2, \dots, i_d}(x_{i_1}, x_{i_2}, \dots, x_{i_d})$$

$$f_{k_1 k_2 \dots k_d}(x_{k_1}, x_{k_2}, \dots, x_{k_d}) = f(\mathbf{x}) - \sum_{\substack{\{i_1 i_2 \dots i_d\} \in \{1 2 \dots D\} \\ \{i_1 i_2 \dots i_d\} \neq \{k_1 k_2 \dots k_d\}}}^D \text{NN}_{i_1 i_2 \dots i_d}(x_{i_1}, x_{i_2}, \dots, x_{i_d})$$

Cycle through all functions
until convergence

Comput Phys Commun **180** (2009) 2002
Chem Rev **121** (2021) 10187 & refs therein

$$f_{k_1 k_2 \dots k_d}(x_{k_1}, x_{k_2}, \dots, x_{k_d}) = f(\mathbf{x}) - \sum_{\substack{\{i_1 i_2 \dots i_d\} \in \{1 2 \dots D\} \\ \{i_1 i_2 \dots i_d\} \neq \{k_1 k_2 \dots k_d\}}}^D \text{GPR}_{i_1 i_2 \dots i_d}(x_{i_1}, x_{i_2}, \dots, x_{i_d})$$

J Phys Chem A **124** (2020) 7598

Comput Phys Commun **271** (2022) 108220 (with example of KED)

- Component functions are easier to build from fewer data: HDMR is helpful when data are sparse, *which they always are in high-D* (full D-dimensional coupling terms may not be recoverable)
- Using ML for f_{i_1, i_2, \dots, i_d} simplifies greatly their calculation (avoids debilitating $D - d$ dimensional integrals which formally express f_{i_1, i_2, \dots, i_d})

HDMR REPRESENTATION ACHIEVED VIA GPR KERNEL DESIGN

$$k(\mathbf{x}, \mathbf{x}') = \sum_{\{i_1 i_2 \dots i_d\} \in \{12 \dots D\}} A_{i_1 i_2 \dots i_d} k_{i_1 i_2 \dots i_d}(\mathbf{x}_{i_1 i_2 \dots i_d}, \mathbf{x}'_{i_1 i_2 \dots i_d})$$

Amplitudes of individual terms are unimportant and can be omitted

$$f(\mathbf{x}) = \sum_{\{i_1 i_2 \dots i_d\} \in \{12 \dots D\}} A_{i_1 i_2 \dots i_d} \sum_{n=1}^M k_{i_1 i_2 \dots i_d}(\mathbf{x}_{i_1 i_2 \dots i_d}, \mathbf{x}_{i_1 i_2 \dots i_d}^{(n)}) c_n$$

$$\equiv \sum_{\{i_1 i_2 \dots i_d\} \in \{12 \dots D\}} B_{i_1 i_2 \dots i_d} \tilde{f}_{i_1 i_2 \dots i_d}(\mathbf{x}_{i_1 i_2 \dots i_d})$$

$$f_{k_1 k_2 \dots k_d}(\mathbf{x}_{k_1}, \mathbf{x}_{k_2}, \dots, \mathbf{x}_{k_d}) = \mathbf{K}_{i_1 i_2 \dots i_d}^* \mathbf{c}$$

In the simplest case $k(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^D k_i(x_i, x'_i)$

$$f(\mathbf{x}) \approx \sum_{i=1}^D f_i(x_i) = \left(\sum_{i=1}^D \mathbf{K}_i^* \right) \mathbf{c}'$$

$$\mathbf{c}' = \left(\sum_{i=1}^D \mathbf{K}_i \right)^{-1} \mathbf{t}$$

\mathbf{K}_i^* is a row vector with elements $k_i(x_i, x_i^{(n)})$

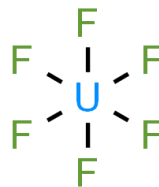
Duvenaud et al. *Adv Neural Inf Process Syst* (2011) 226–234
 Manzhos et al. *Machin Learn Sci Technol* **2** (2022) 01LT02
 Rasmussen GPR books also talks about additive kernels

If the kernel is in HDMR form the final function representation also is

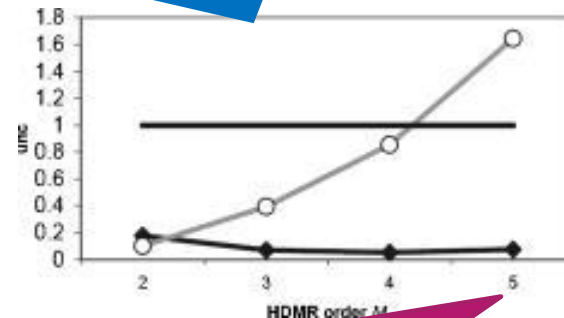
USE OF HDMR TO GENERATE SYNTHETIC DATA TO HELP TUNE HYPERPARAMETERS

$$f(\mathbf{x}) \approx f_0 + \sum_{i=1}^D f_i(x_i) + \sum_{\substack{i=1, \\ j=i+1}}^D f_{ij}(x_i, x_j) + \sum_{i_1, i_2, \dots, i_d}^D f_{i_1, i_2, \dots, i_d}(x_{i_1}, x_{i_2}, \dots, x_{i_d}) + \dots$$

Low-order terms can be reliably determined from few data

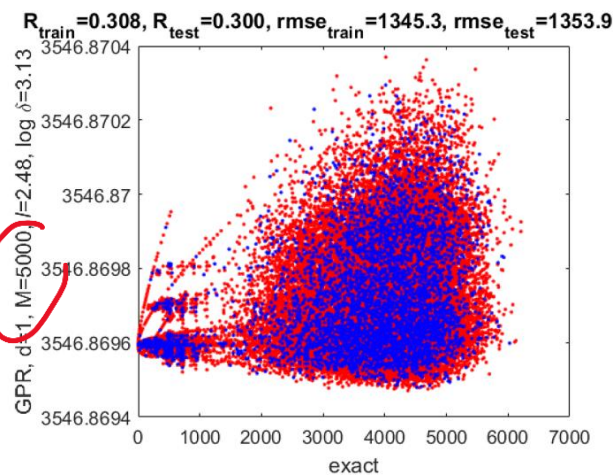
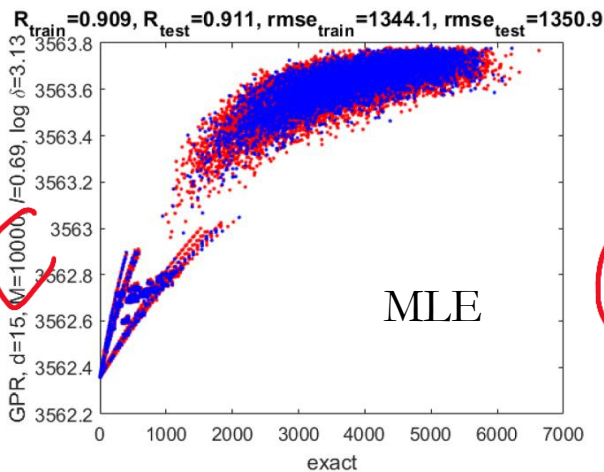
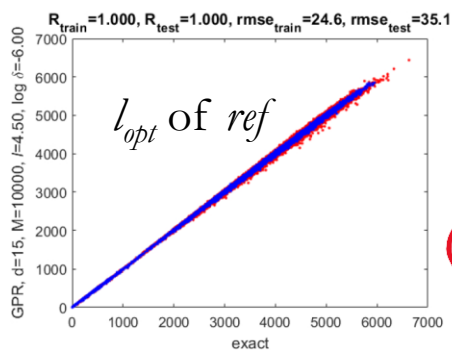
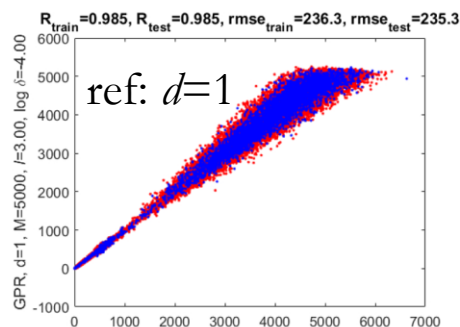
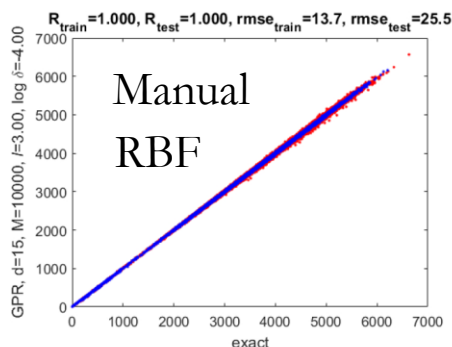


Relative uncertainties of NN parameters in an HDMR-NN fit of H₂O₂ PES (*J. Chem. Phys.* **125** (2006) 084109)



Higher order terms may not be recoverable

- Achievable quality of fit is limited by the density of sampling



WITH SPARSE DATA ONLY LOW-ORDER TERMS ARE RECOVERABLE : EXAMPLE OF UF6 INTERATOMIC POTENTIAL

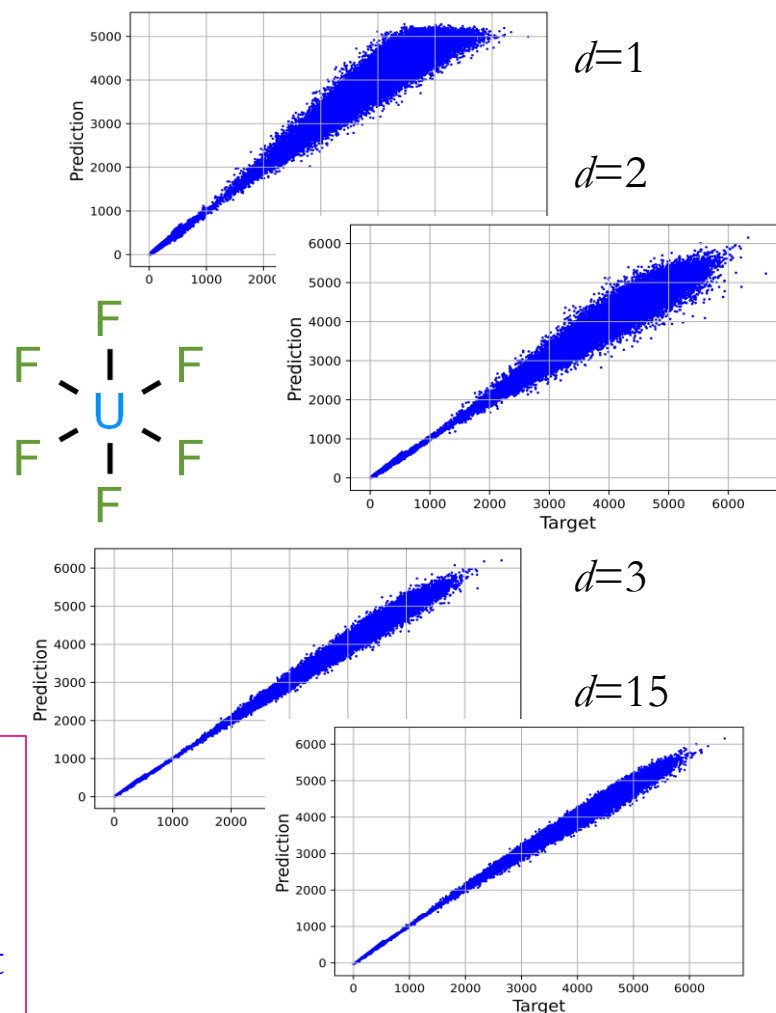
$$V^{UF_6}(\mathbf{x}) \approx \sum_{i=1}^N f_i^{\text{GPR}}(x_1^{(i)}, x_2^{(i)}, \dots, x_d^{(i)})$$

\mathbf{x} : 15(= D) modes of vibration

Test set of
50,000 pts

rmse, cm ⁻¹	N	5,000	3,000	2,000
Full-D ($d = D$)	1	42.2	75.4	106.7
$d = 1$	15	234.6	236.4	237.3
$d = 2$	105	168.1	178.6	190.3
$d = 3$	455	65.6	78.0	97.4

- With 2000 points in 15D space
 - 1.66 data per degree of freedom
 - rmse(3d-HDMR-GPR) < RMSE(full-D GPR)
- A finite dataset in D-dimensional space is not a D-dimensional object



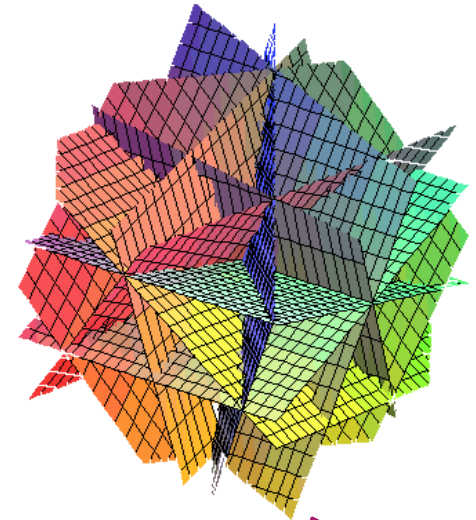
HOW TO DEAL WITH OVERFITTING AND NON-LINEAR OPTIMIZATION? - BEYOND TRADITIONAL NN

traditional NN: \mathbf{w} , b define hyperplanes that are optimized for given data, $\sigma(\mathbf{x})$, and N_{neurons}

$$\sigma(\mathbf{w}_n \mathbf{x} + b_n) \quad f(\mathbf{x}) = \sum_{n=0}^N c_n \sigma_{\mathbf{n}}(\mathbf{w}_n \mathbf{x} + b_n)$$

We will instead define \mathbf{w} , b with rule and then **optimize** $\sigma(\mathbf{x})$ for given data and N_{neurons}

- \mathbf{w} define planes' orientations – sample them
- b define distance from origin – we will not need to deal with them



“ridge NN”

- Avoid non-linear optimization
- No particular advantage using the same activation function for all neurons if parameters are not optimized
 - Can optimize neuron activation functions for a particular problem
 - This makes sense if we can build cheaply... which we can

NN AS AN ADDITIVE MODEL IN REDUNDANT COORDINATES

The only extra cost vs standard GPR is summation in the kernel (parallelizable)

$$k(\mathbf{y}, \mathbf{y}') = \sum_{i=1}^N k_i(y_i, y'_i) \quad k_i(y_i, y'_i) = \exp\left(-\frac{(y_i - y'_i)^2}{2l^2}\right)$$

$$f(\mathbf{x}) = \sum_{i=1}^N c_i \sigma_i(\mathbf{w}_i \mathbf{x} + b_i) \equiv \sum_{m=1}^M b_m k(\mathbf{y}, \mathbf{y}^{(m)}) = \sum_{i=1}^N \sum_{m=1}^M b_m k_i(y_i, y_i^{(m)}) = \sum_{i=1}^N f_i(y_i)$$

NN
in \mathbf{x}

$$\mathbf{y}_n = \mathbf{x}^T \mathbf{s}_n$$

1d-GPR
in \mathbf{y}

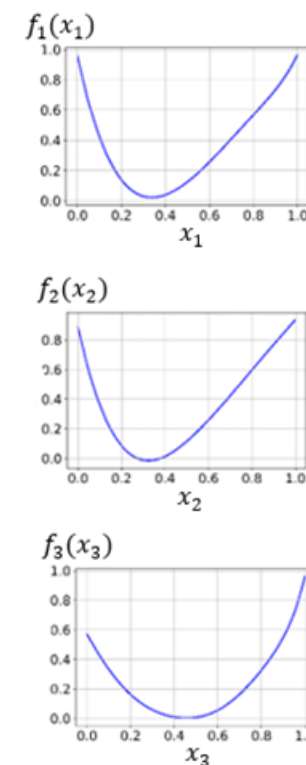
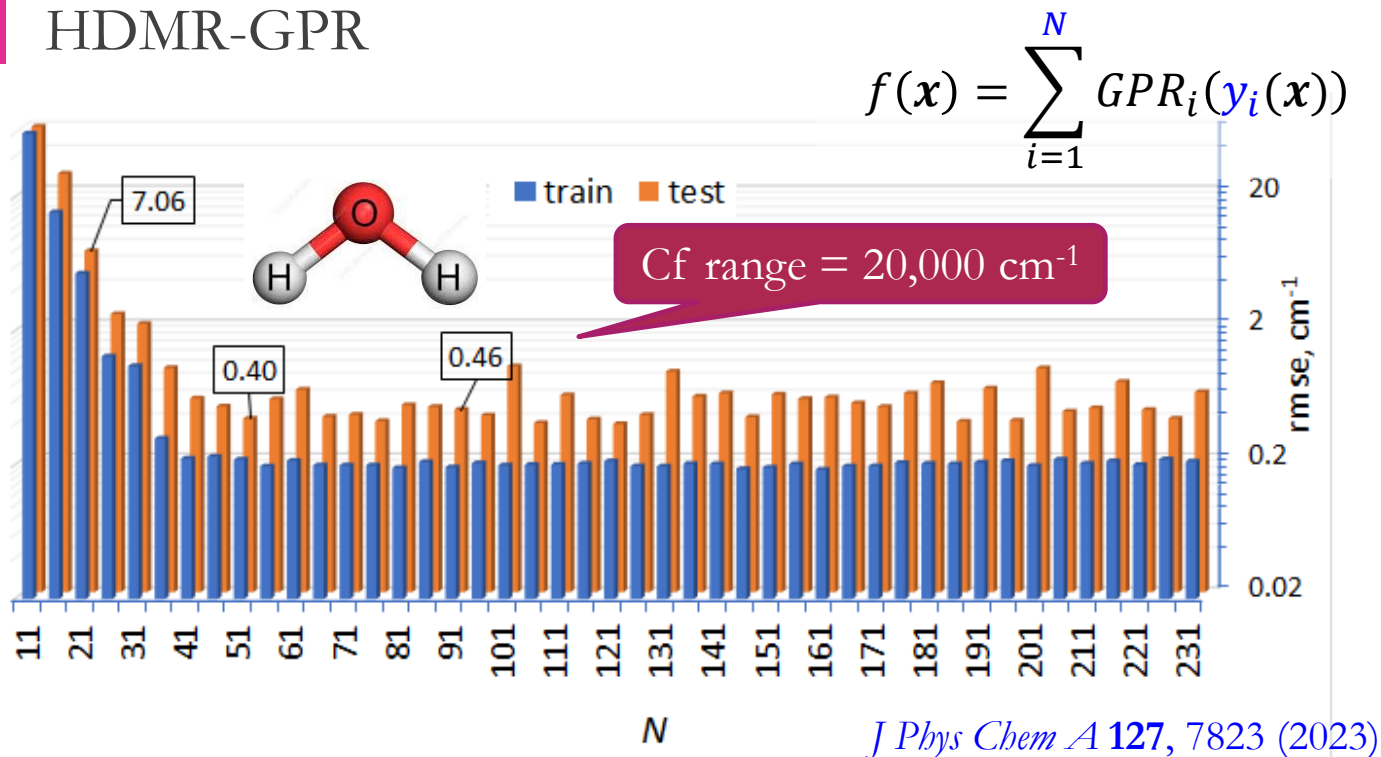
Elements of Sobol sequence
 \mathbf{w}_i define “directions” of \mathbf{y}_i -
distribute them pseudorandomly

- A first order additive model in redundant coordinates \mathbf{y}
- f_i can be obtained from 1st order HDMR-GPR in one go with an additive kernel... for given \mathbf{W}
- Can reduce N by using different (optimal) neuron activation functions for different neurons

J Phys Chem A **127**, 7823 (2023)

- Avoid nonlinear parameter optimization with rule-based NN weights
 - With no nonlinear optimization, no advantage for all neuron activation functions to be the same
- Without nonlinear parameter optimization... what happens to overfitting as the no. of neurons is increased beyond optimum?

NN WITH OPTIMAL NN NEURONS BUILT WITH 1ST ORDER HDMR-GPR



- H₂O PES, training on 1000 samples, large test set of 9000 samples
- **no overfitting** as no. of neurons grows due to absence of nonlinear optimization
- Combines expressive power of NN with robustness of linear regression

- See Huixin Liu (刘辉鑫)'s poster for use of GPR-NN for nuclear mass prediction
- See *J. Phys. Chem. Lett.* **15**, 6974 (2024) for use in studies of neuromorphic computing

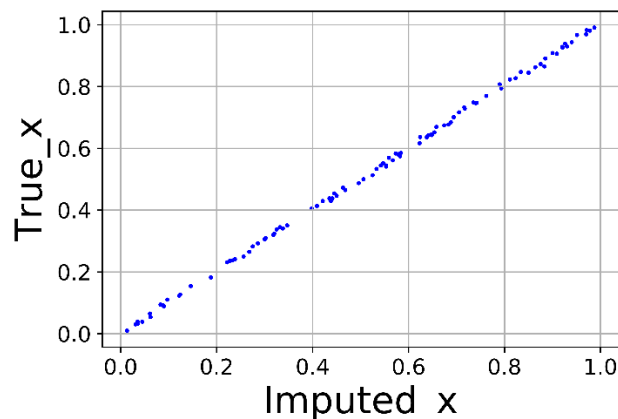
OTHER ADVANTAGES OF HMDR-GPR: HELPS IMPUTE MISSING DATA...

$$f(\mathbf{x}) \approx f_0 + \sum_{i=1}^D f_i(x_i)$$

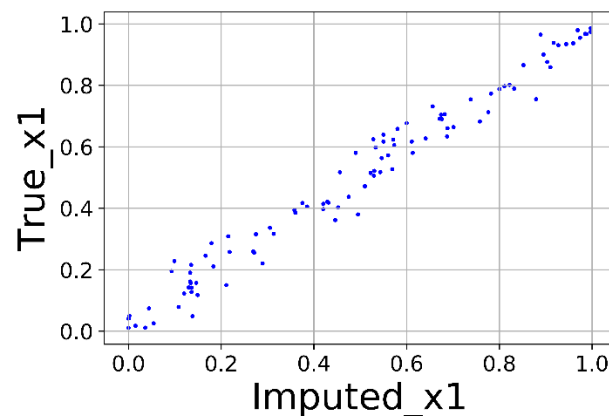
$$f_i(x_i^m) = f(\mathbf{x}_m) - f_0 - \sum_{\substack{j=1, \\ j \neq i}}^D f_j(x_j^m)$$

$$x_i^m = f_i^{-1}(y_i^m)$$

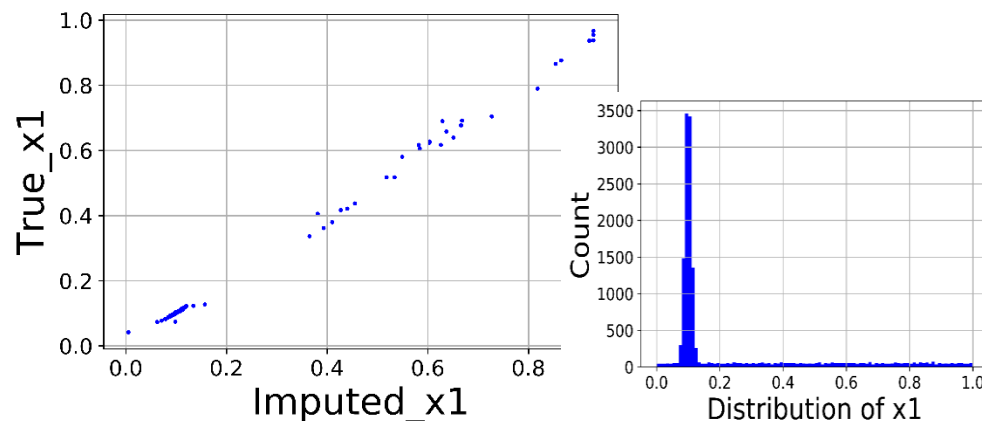
$$f(x, y, z) = x + y + z$$



$$f(x, y, z) = x + y + z + \text{noise}$$



$$f(x, y, z) = x + y + z \text{ with bad distribution}$$



Comput Phys Commun, **271**, 108220 (2022)

Code: <https://github.com/owen-ren0003/rshdmrgpr>

... and helps find optimal hyperparameters
when data are sparse: *J Math Chem* **61**, 7 (2023)

ML-TO-FORMULA: ML-GUIDED KEF CONSTRUCTION

Kinetic energy in Orbital-free DFT

$$E_{kin} = \int \tau(\mathbf{r}) d\mathbf{r}$$

$$\mathbf{x} = (\overline{\tau_{TF}}, \overline{\tau_{TF}p}, \overline{\tau_{TF}p^2}, \overline{\tau_{TF}pq}, \overline{\tau_{TF}q^2}, \overline{\rho V_{eff}})$$

$$\bar{\tau}(\mathbf{x}) = \sum_{i=1}^N GPR_i(\mathbf{y}_i(\mathbf{x}))$$

$$E_{kin} = V\bar{\tau}$$

Scaled gradient expansion (4th order)

$$\tau_{GE4}(\mathbf{r})$$

$$= \tau_{TF}(\mathbf{r}) \left(1 + \frac{5}{27} p(\mathbf{r}) + \frac{20}{9} q(\mathbf{r}) + \frac{8}{81} q(\mathbf{r})^2 - \frac{1}{9} p(\mathbf{r}) q(\mathbf{r}) + \frac{8}{243} p(\mathbf{r})^2 \right)$$

$$p = \frac{|\nabla \rho|^2}{4(3\pi^2)^{2/3} \rho^{8/3}}, \quad q = \frac{\Delta \rho}{4(3\pi^2)^{2/3} \rho^{5/3}}$$

“insightful” ML

See *Digital Discovery* **3**, 1967 (2024),
J. Mater. Inform. **5**, 38 (2025) for uses
in materials informatics

Additive model,
feature importance

Coupling
terms

$$\bar{\tau}(\mathbf{x}) = \sum_{i=1}^D f(\mathbf{x}_i) + \sum_{i=D+1}^N f_i(\mathbf{y}_i(\mathbf{x}))$$

Functional form

ML-guided analytic KEF

$$\bar{\tau} = \sum_{n=1}^6 \sum_{p=1}^{P_n} a_{np} (x'_n)^p$$

Mach. Learn. Sci. Technol. **6**, 035002 (2025)

GPR-NN WITH OPTIMIZED REDUNDANT COORDINATES COULD OBVIATE DEEP NN IN SOME APPLICATIONS

Cf. kernel regression $f(\mathbf{x}) = \sum_{i=1}^M c_i k(\xi | \mathbf{x}, \mathbf{x}^{(i)})$

$$f(\mathbf{x}) = \sum_{k_n=0}^{N_n} w_{l,k_n}^{(n)} \sigma_{n,k_n} \left(\sum_{k_{n-1}=0}^{N_{n-1}} w_{k_n,k_{n-1}}^{(n-1)} \sigma_{n-1,k_{n-1}} \left(\dots \sum_{k_1=0}^{N_1} w_{k_2,k_1}^{(2)} \sigma_{1,k_1} \left(\sum_{i=0}^d w_{k_1,i}^{(1)} x_i \right) \right) \right)$$

$$\equiv \sum_{i=0}^{N_n} c_i \theta_i(\mathbf{W}|\mathbf{x}) \quad \theta_i(\mathbf{x}) = \theta_i(\mathbf{x}_0) + \nabla \theta_i(\mathbf{x}_0) \Delta \mathbf{x} + \frac{1}{2} (\Delta \mathbf{x})^T \mathbf{H} \Delta \mathbf{x} + \dots$$

$$f_{GPRNN}(\mathbf{x}) = \sum_{i=1}^N c_i \sigma_i(\mathbf{w}_i \mathbf{x} + b_i) \equiv \sum_{m=1}^M b_m k(\mathbf{y}, \mathbf{y}^{(m)}) = \sum_{i=1}^N \sum_{m=1}^M b_m k_i(y_i, y_i^{(m)}) = \sum_{i=1}^N f_i(\mathbf{y}_i)$$

$$\mathbf{y} = \mathbf{W}\mathbf{x}$$

Approaches expressive
power of deep NN in
diadic approximation

$$f_n(y_n(\mathbf{x})) = f_n(y_n(\mathbf{x}_0)) + \left(\mathbf{w}_n \frac{df_n}{dy_n} \Big|_{y_{n0}} \right) \Delta \mathbf{x} + \frac{1}{2} (\Delta \mathbf{x})^T \left(\mathbf{w}_n \mathbf{w}_n^T \frac{d^2 f_n}{dy_n^2} \Big|_{y_{n0}} \right) \Delta \mathbf{x} + \dots$$

EASE OF BUILDING AN ORDERS-OF-COUPLING (HDMR) REPRESENTATION

0.5	0.5	0.5	0	0
0.75	0.25	0.75	0	0
0.25	0.75	0.25	0	0
0.375	0.375	0.625	0	0
0	0.5	0.5	0.5	0
0	0.75	0.25	0.75	0
0	0.25	0.75	0.25	0
0	0.375	0.375	0.625	0
0	0	0.5	0.5	0.5
0	0	0.75	0.25	0.75
0	0	0.25	0.75	0.25
0	0	0.375	0.375	0.625
0	0.5	0	0.5	0.5
0	0.75	0	0.25	0.75
0	0.25	0	0.75	0.25
0	0.375	0	0.375	0.625
0.5	0	0	0.5	0.5
0.75	0	0	0.25	0.75
0.25	0	0	0.75	0.25
0.375	0	0	0.375	0.625

...

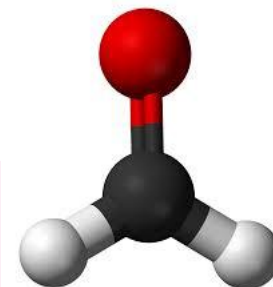
$$f(\mathbf{x}) = \sum_{i=1}^N c_i \sigma_i(\mathbf{w}_i \mathbf{x} + b_i) = \sum_{i=1}^N f_i(y_i)$$

$$y_n = \mathbf{x}^T \mathbf{s}_n$$

D-dimensional vector with

- d elements of d -dimensional Sobol sequence
- $D - d$ elements are zero

Order of coupling d	No. of coupling terms	Test rmse, cm^{-1} (this work)	Test rmse, cm^{-1} (Ref. [44])
1	6	1302.7	1315.6
2	15	385.5	410.0
3	20	20.4	29.1
4	15	14.6	16.1
5	6	16.5	14.4
6	1	19.6	23.8



When data density is low, full-D terms may not be recoverable

CONCLUSIONS

- ML methods are useful in high-dimensional problems as they avoid direct product representations, but when data density is very low or dimensionality is high, off-the-shelf methods may fail
 - Many non-linear parameters in NNs that also scale with D
 - Data are always sparse in high- D and high-order coupling terms (full- D function) may not be recoverable - Cannot be palliated by just adding more data
- GPR/KRR are attractive as they combine high expressive power of nonlinear kernel and robustness of linear regression
 - Can collapse in very high D
 - Loss of kernel resolution / property of locality of Matern kernels, loss of advantage of multi- ζ bases
 - Loss of advantage vs low-order polynomial regression
- One can stabilize the solution by using orders of coupling representations
 - Low-order (low- d) terms are easier to build (fit) and to use
 - It is convenient to use machine-learned component functions: HDMR-NN, HDMR-GPR
 - Low-order HDMR helps impute missing data and optimize hyperparameters
- It is possible to combine advantages of NN and GPR:
 - GPRNN: NN in original feature space, additive GPR in redundant coordinates
 - No nonlinear optimization
 - The shapes of all neuron activation functions are gotten in one go from 1d-HDMR-GPR
 - Optimal for given data, weights, and NN size
 - Provides an easy way to build orders-of coupling representations
 - With optimized redundant coordinates conceptually approaches expressive power of deep NN
- We provided examples from computational chemistry but **these approaches are general**

THANKS

Select collaborators, group members

Manabu Ihara
Johann Lüder
Tucker Carrington
Dmitry Voytsekhovsky
Owen Ren
Pavlo Golub
Eita Sasaki
Shunsaku Tsuda
Hyojae Lee
Heejoo Yoon
Mohamed Boussaidi



Ministry of Education
SINGAPORE

Funders

NSERC
Mitacs
Ministry of Education of Singapore
JST
Dlab
Digital Research Alliance of Canada



Digital Research
Alliance of Canada

Alliance de recherche
numérique du Canada



APPENDIX

RECAP OF SOME MAJOR ML METHODS: NEURAL NETWORKS

$$f_k(\mathbf{x}) = \sum_{n=0}^N c_{nk} \sigma_n(\mathbf{w}_n \mathbf{x} + b_n)$$

Training set

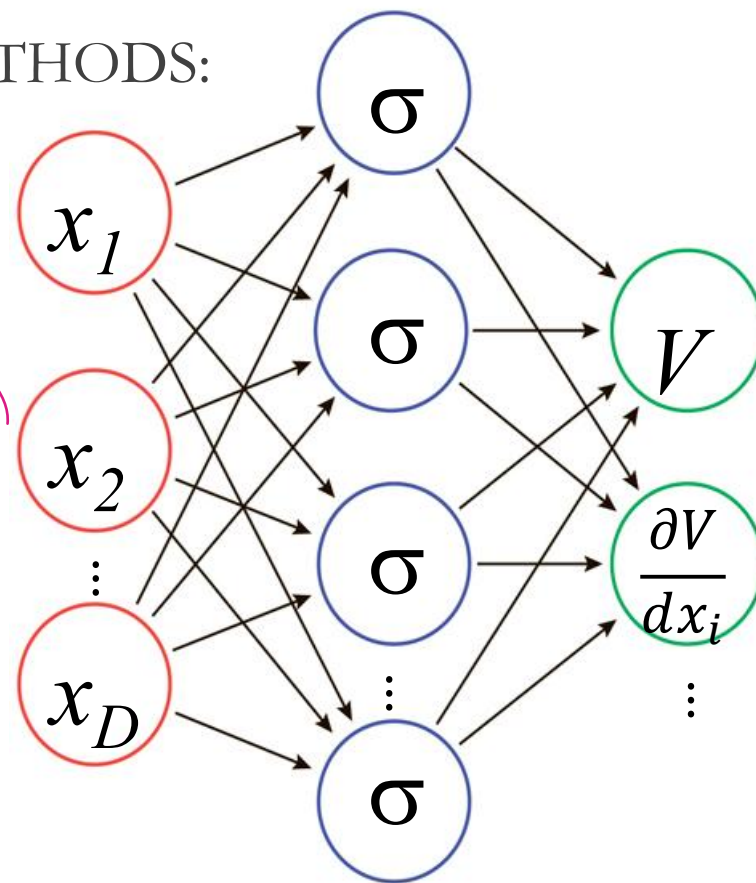
$$\{f^{(j)}, \mathbf{x}^{(j)}\}, j = 1, \dots, M; \mathbf{x} \in \mathbb{R}^D$$

weights

biases

Often $\sigma(x) = (e^x - e^{-x}) / (e^x + e^{-x})$
but can be **any smooth nonlinear** function

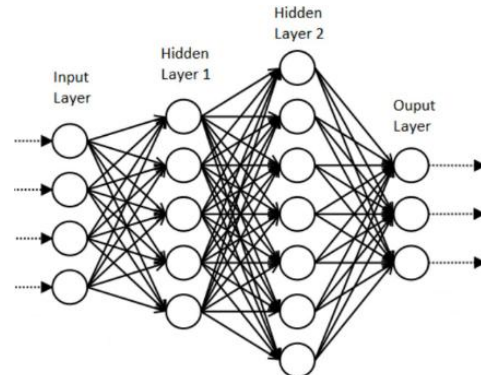
Key works by *Kolmogorov, Specher, Hornik, Gorban*



- An analogy with a biologic neural network is often made
- As a mathematical object used for non-linear regression, a view of a representation in a **non-direct product basis** is useful
- A single hidden layer NN is a universal approximator
 - Think carefully if you actually need a “deep” NN as it comes at a price of a large no. of non-linear parameters

PROS AND CONS OF FFNN

Int J Quant Chem **2015**, *115*, 1012; *Chem Rev* **2021**, *121*, 10187



“Deep” NN:

Non-direct product basis

$$f_l = \sigma_{out} \left(\sum_{k_n=0}^{N_n} w_{l,k_n}^{(n)} \sigma_{n,k_n} \left(\sum_{k_{n-1}=0}^{N_{n-1}} w_{k_n,k_{n-1}}^{(n-1)} \sigma_{n-1,k_{n-1}} \left(\dots \sum_{k_1=0}^{N_1} w_{k_2,k_1}^{(2)} \sigma_{1,k_1} \left(\sum_{i=0}^d w_{k_1,i}^{(1)} x_i \right) \right) \right) \right)$$

$$\sigma = \exp(\quad): f(\mathbf{x}) = \sum_{i=0}^M w_i^{(2)} \prod_{k=0}^d e^{w_{ik}^{(1)} x_k}$$

$$\int f(\mathbf{x}) d\mathbf{x} = \sum \prod \int dx_i f_i(x_i) dx_i$$

J Chem Phys **2006**, *125*, 194105

$$f(\mathbf{x}) = \mu_1^{(2)} + \sum_{k_1=1}^{n_1} w_{1,k_1}^{(2)} \prod_{k_0=1}^J \operatorname{erf} \left(\mu_{k_1 k_0}^{(1)} + w_{k_1 1 k_0}^{(1)} x_{k_0} \right)$$

Refs in color are ours

Neural Computation **2001**, *14*, 241

Refs in black are literature

- A single hidden layer NN is a universal approximator
 - Think carefully if you actually need a deep NN as it comes at a price of **a large no. of non-linear parameters**
- Easy to achieve sum-of-products
 - Important for multidimensional integration

$$\forall \delta > 0, \exists N < \infty: \left| t(\mathbf{x}) - \sum_{n=0}^N c_n \sigma(\mathbf{w}_n \mathbf{x} + b_n) \right| < \delta$$

Universal approximator theorems do not consider the data aspect

GPR/KRR:

LINEAR REGRESSION WITH A NONLINEAR BASIS

what is the expected value of f at \mathbf{x} given the set $\{f^{(n)}, \mathbf{x}^{(n)}\}$?

$$f(\mathbf{x}) = \mathbf{K}^* (\mathbf{K}^{-1} \mathbf{f}) = \sum_{n=1}^M b_n k(\mathbf{x}, \mathbf{x}^{(n)})$$

$$\text{var}(f(\mathbf{x})) = \mathbf{K}^{**} - \mathbf{K}^* \mathbf{K}^{-1} \mathbf{K}^{*T}$$

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) + \delta & k(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) & \dots & k(\mathbf{x}^{(1)}, \mathbf{x}^{(M)}) \\ k(\mathbf{x}^{(2)}, \mathbf{x}^{(1)}) & k(\mathbf{x}^{(2)}, \mathbf{x}^{(2)}) + \delta & \dots & k(\mathbf{x}^{(2)}, \mathbf{x}^{(M)}) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}^{(M)}, \mathbf{x}^{(1)}) & k(\mathbf{x}^{(M)}, \mathbf{x}^{(2)}) & \dots & k(\mathbf{x}^{(M)}, \mathbf{x}^{(M)}) + \delta \end{pmatrix}$$

$$\mathbf{K}^* = (k(\mathbf{x}, \mathbf{x}^{(1)}) \quad k(\mathbf{x}, \mathbf{x}^{(2)}) \quad \dots \quad k(\mathbf{x}, \mathbf{x}^{(M)})) \quad \mathbf{K}^{**} = k(\mathbf{x}, \mathbf{x})$$

Matern family of functions:

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\sqrt{2\nu} \frac{|\mathbf{x} - \mathbf{x}'|}{l} \right)^\nu K_\nu \left(\sqrt{2\nu} \frac{|\mathbf{x} - \mathbf{x}'|}{l} \right)$$

particular cases: exp, Gaussian

https://en.wikipedia.org/wiki/Matern_covariance_function

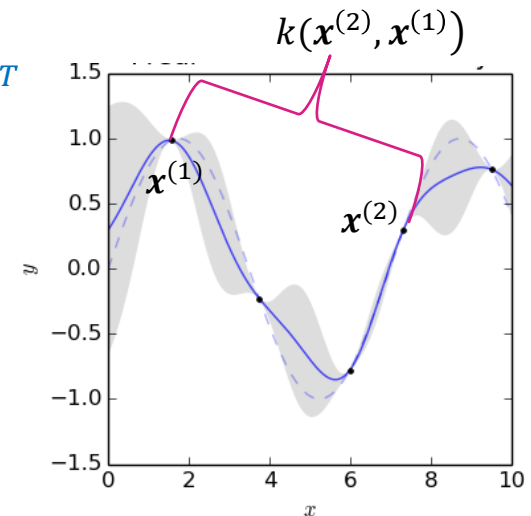
hyperparameters

■ Robust (a type of linear regression) but

- Hard to apply to large datasets
- Defining eqs do not scale with D but can collapse in very high-D

(*Phys Chem Chem Phys* **25** (2023) 1546)

Matrix \mathbf{K} can describe how correlated is each pair of data points



Optimized l_i informs on relevance of the i -th variable (ARD, automated relevance determination)

$$\nu = \infty: (\text{RBF}): \sigma^2 \exp \left(-\frac{|\mathbf{x} - \mathbf{x}'|^2}{2l^2} \right)$$

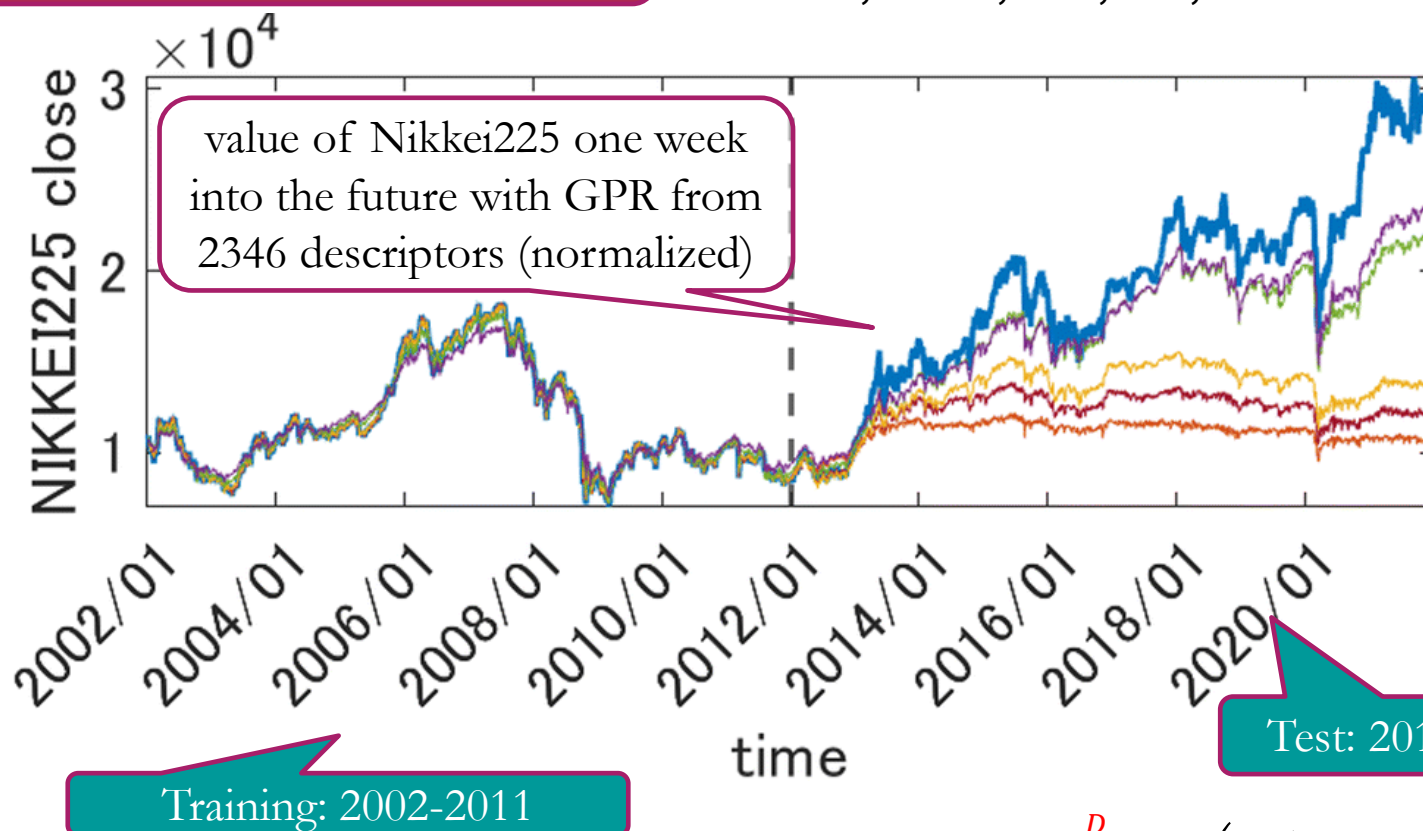
$$\nu = \frac{1}{2}: \sigma^2 \exp \left(-\frac{|\mathbf{x} - \mathbf{x}'|}{l} \right)$$

Can be palliated with additive GPR (HDMR)

HDMR AVOIDS COLLAPSE OF GPR IN VERY HIGH-D

kernel length parameters small enough to preserve the nonlinear nature of the model

$l = 0.1, l = 0.5, l = 1, l = 5, l = 10$. Blue: target



Phys. Chem. Chem. Phys. **25**, 1546 (2023)

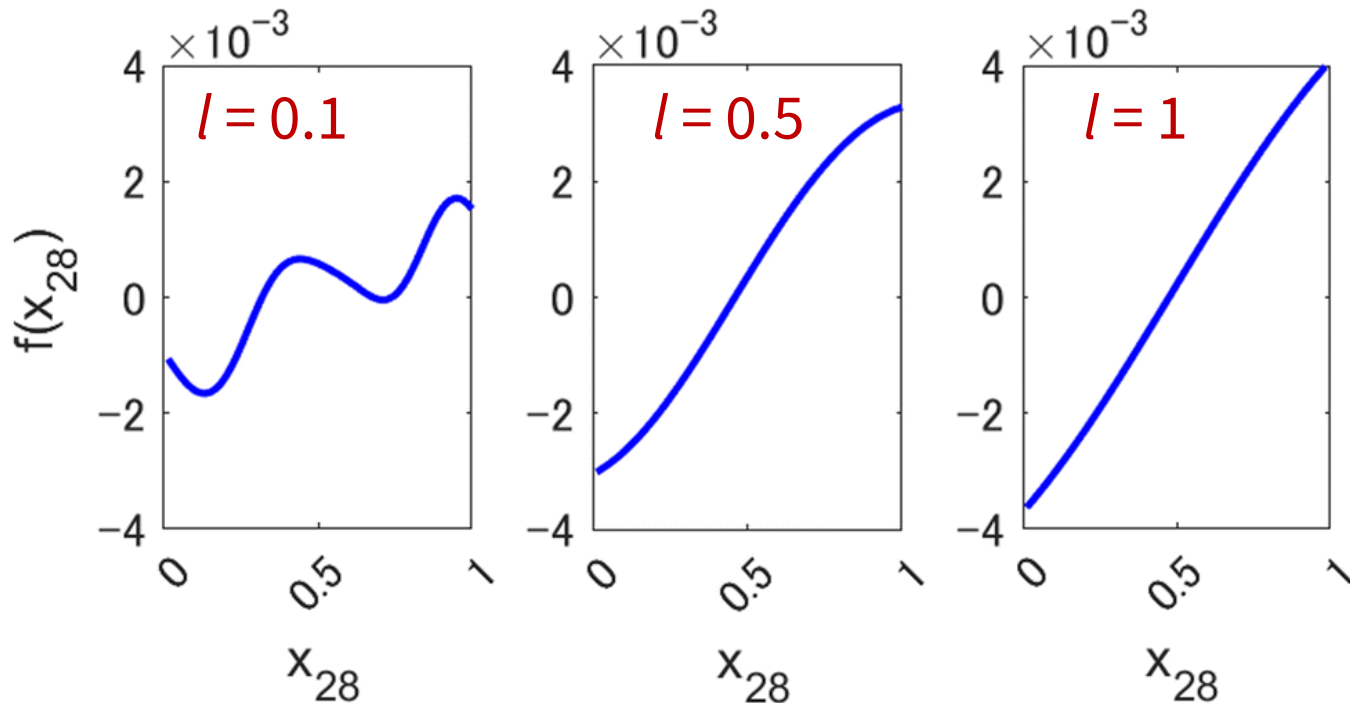
$$k(\mathbf{x}, \mathbf{x}') = \sum_{d=1}^D \exp \left(- \left(\frac{x_d - x'_d}{l} \right)^2 \right)$$

HDMR ABLE TO RECOVER THE VALUE OF NONLINEARITY

$f_i(x_i)$ in the example of N225 forecasting

$$f(\mathbf{x}^*) = f_1(x_1^*) + \cdots + f_D(x_D^*)$$

compare to linear regression using a basis: $f(\mathbf{x}) = \sum_{i=1}^D c_i \theta_i(x_i)$

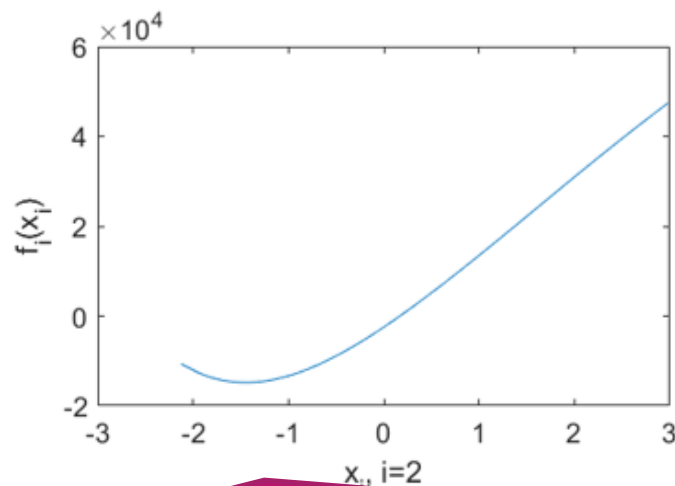
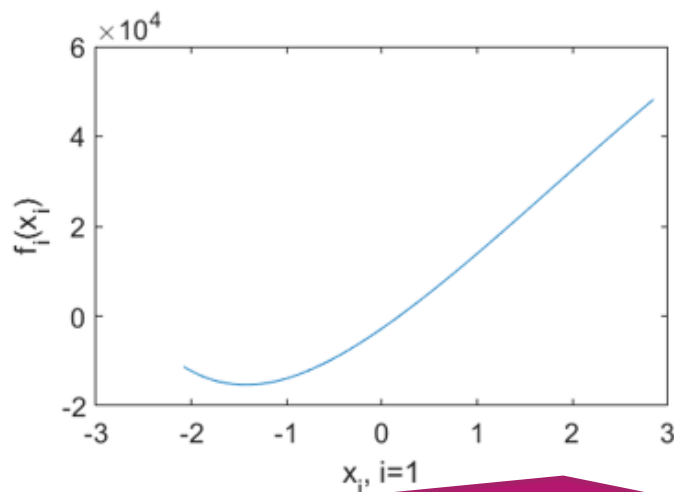
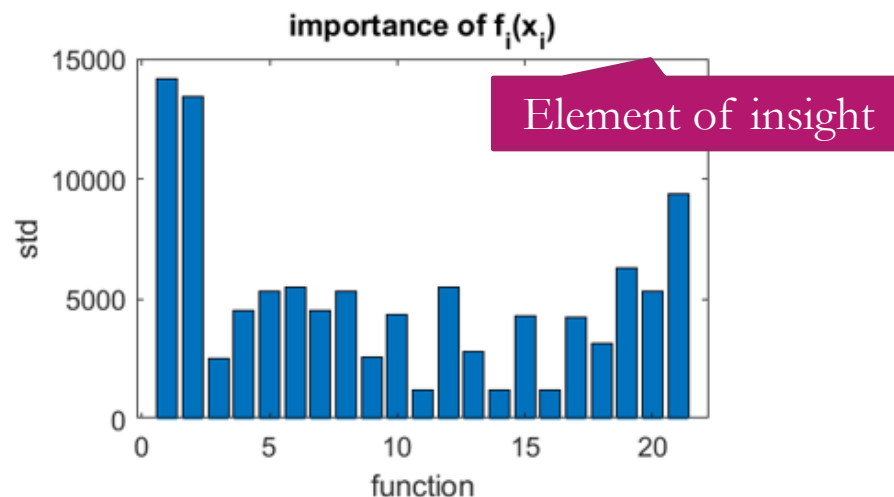
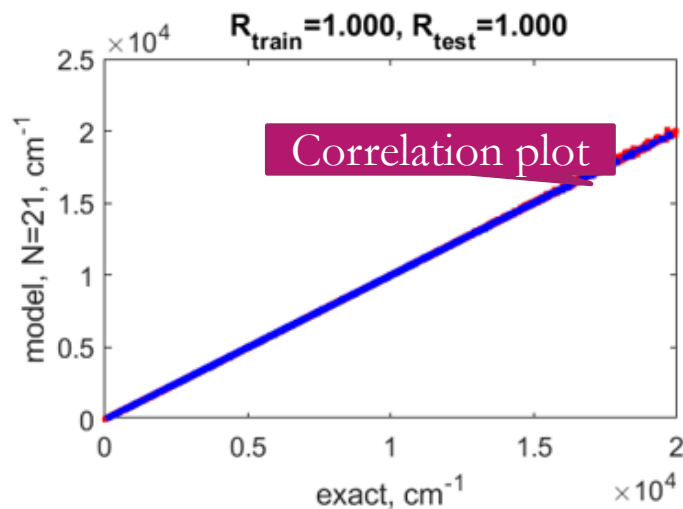


With
 $D = 2346$
even 2nd order
terms ($f_{ij}(x_i, x_j)$)
are hard

$$\sum_{i_1, i_2, \dots, i_d}^D f_{i_1, i_2, \dots, i_d}(x_{i_1}, x_{i_2}, \dots, x_{i_d})$$

... preserving a higher expressive power of a nonlinear model

NN WITH OPTIMAL NN NEURONS BUILT WITH 1ST ORDER HDMR-GPR



Two neuron activation functions for terms with the largest magnitude