



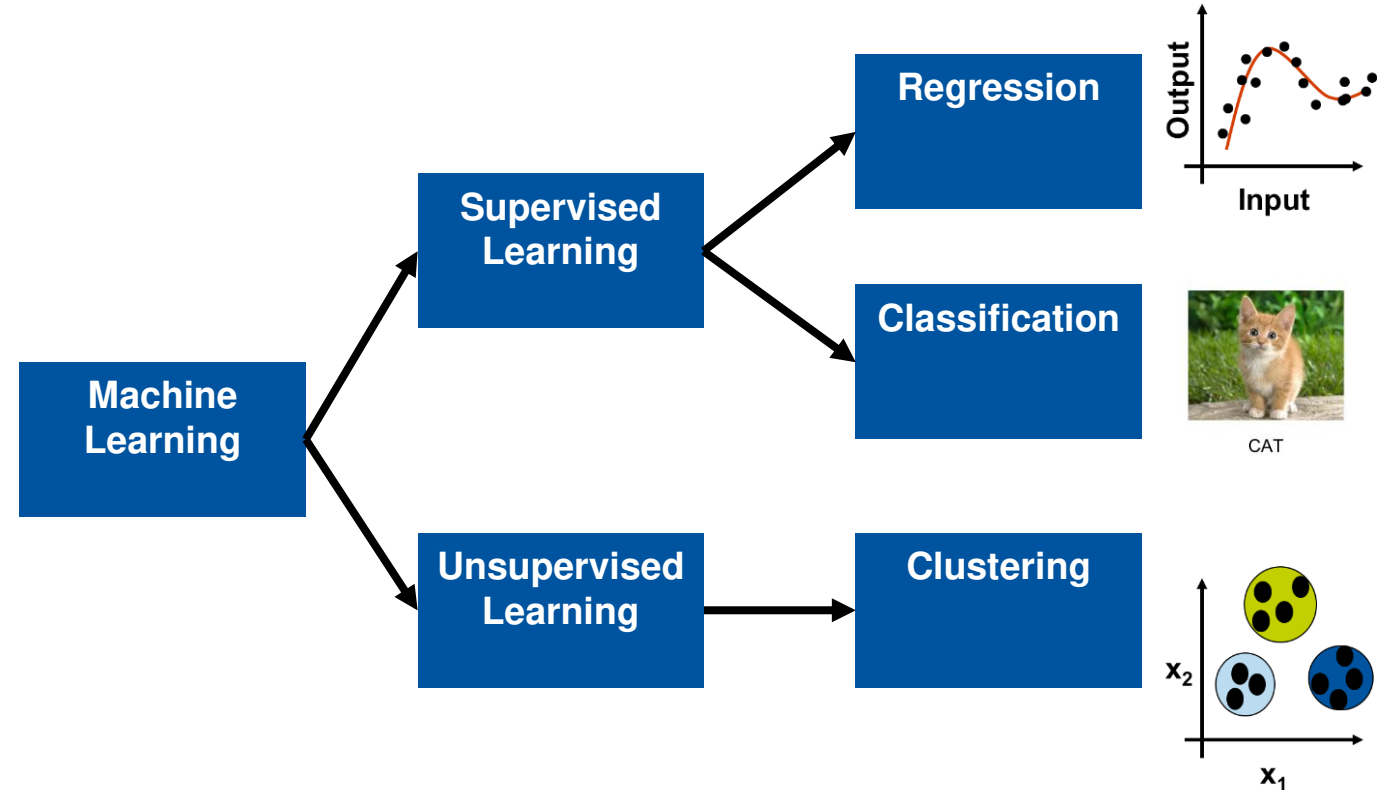
Applied Numerical Optimization

Prof. Alexander Mitsos, Ph.D.

Machine learning and hybrid modeling

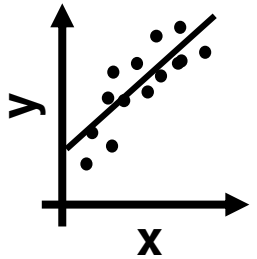
Machine Learning¹

- Programming computers using example data or past experience
- Learning general models from example data
- Key ingredient: **data**
 - formats vary widely
 - labeled?
 - structured?
 - relevant?

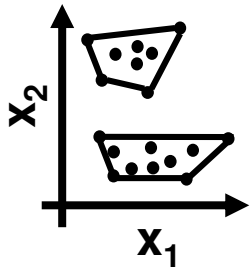


Data Driven Model Structures

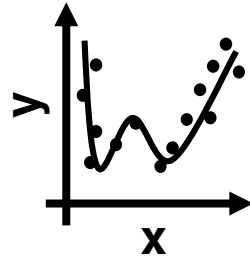
Linear



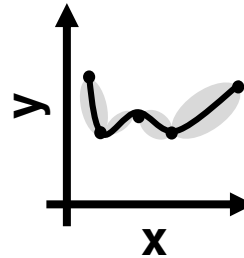
Convex region
linear surrogate



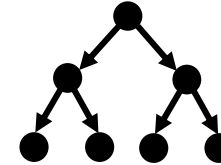
Basis
functions



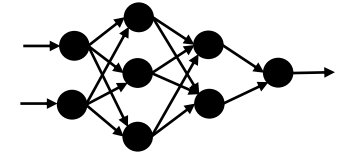
Gaussian
processes



Gradient-boosted
trees



Artificial neural
networks

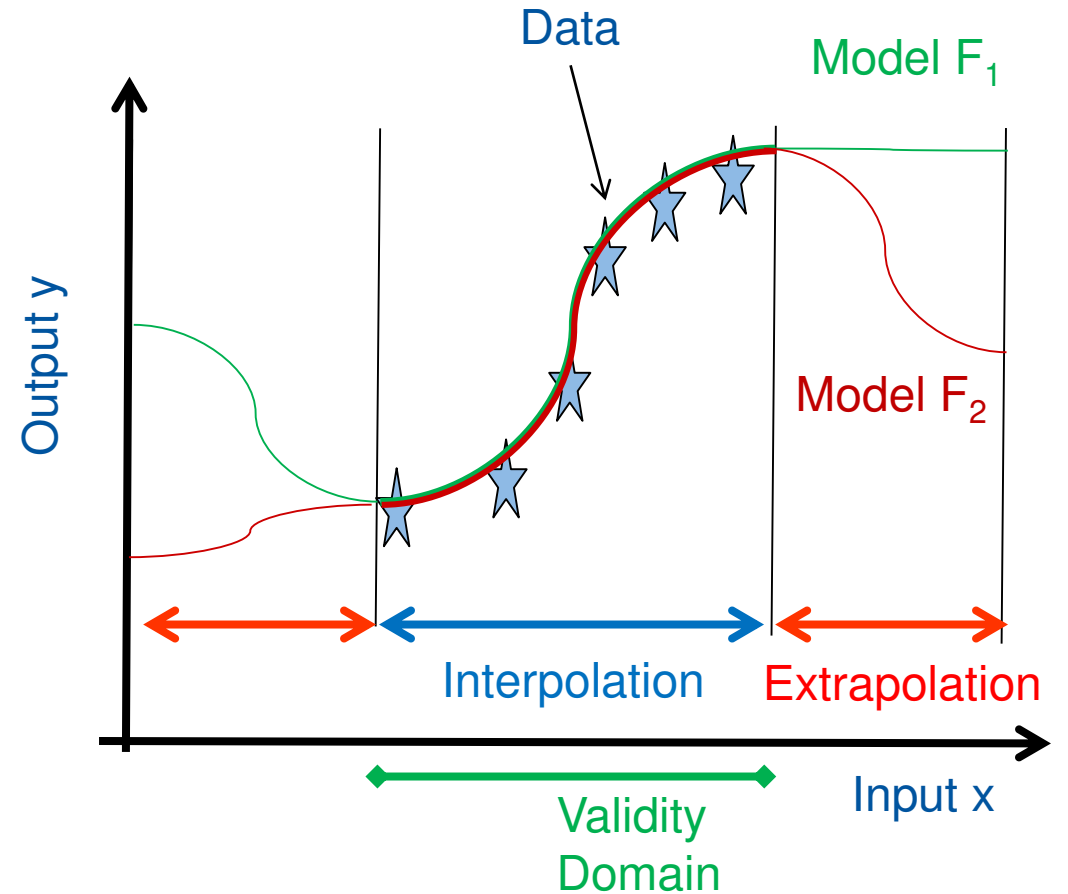


Surrogate model complexity

Interpolation and Extrapolation

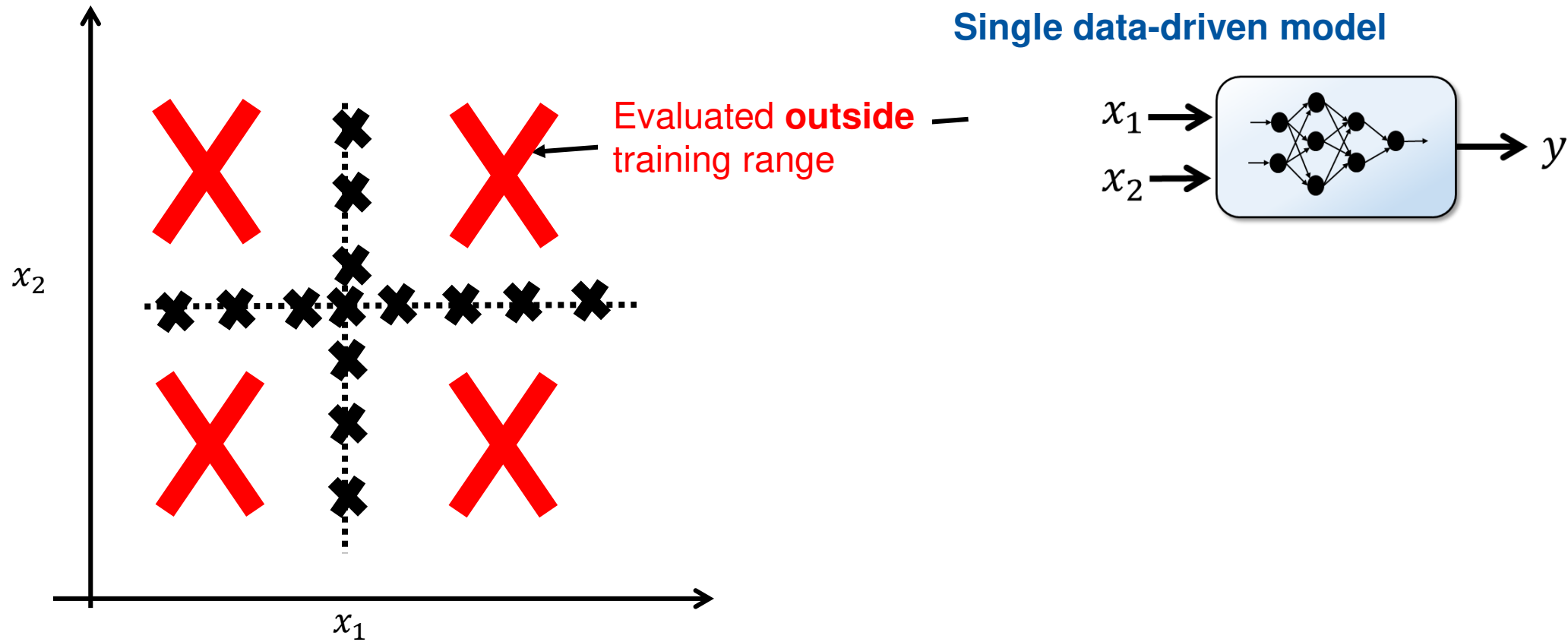
- **Interpolation**: prediction within the convex hull of training data
- **Extrapolation**: prediction outside the convex hull of training data

→ Data-driven models cannot extrapolate

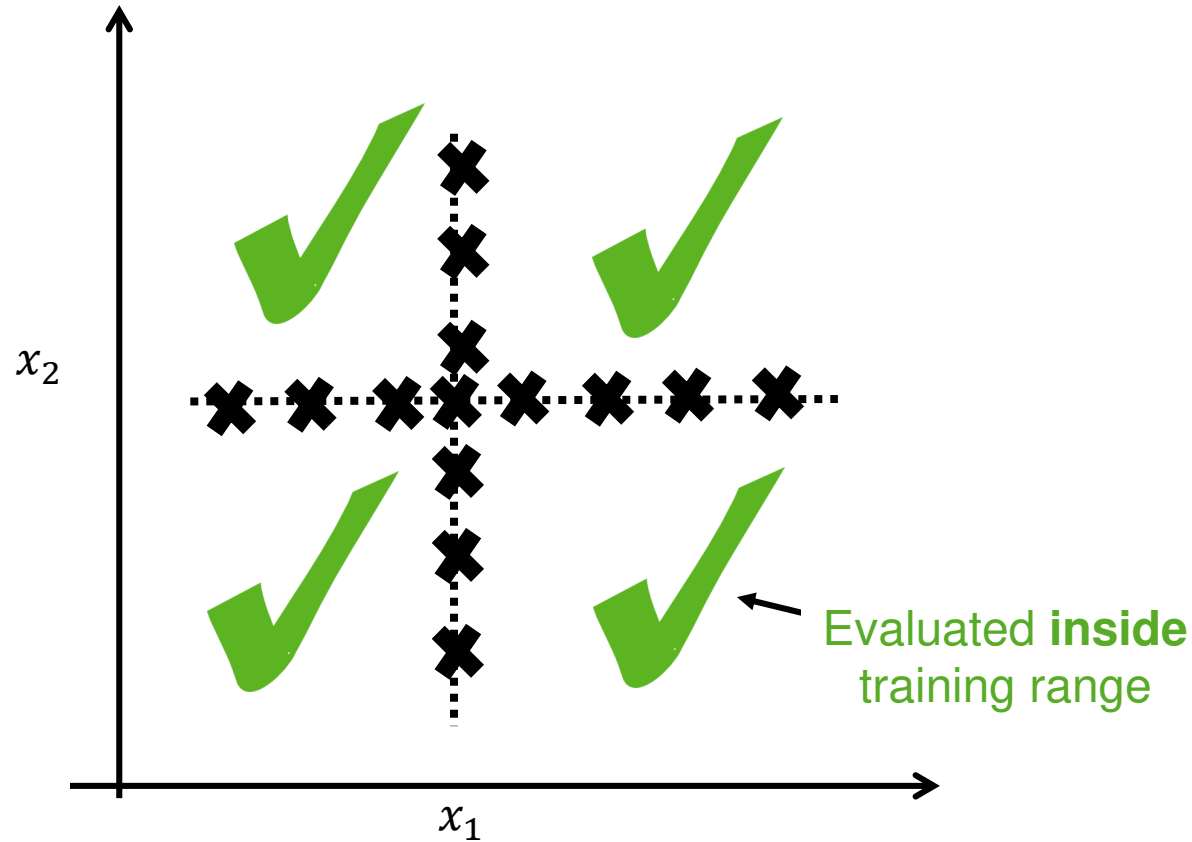


→ Staying in validity domain is essential for black-box models

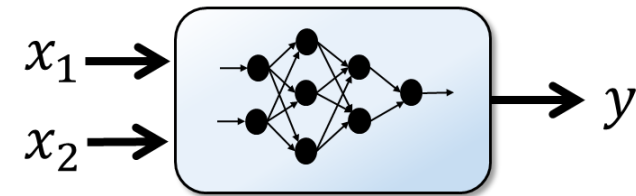
Hybrid Modelling – Combining Mechanistic and Data-driven Models



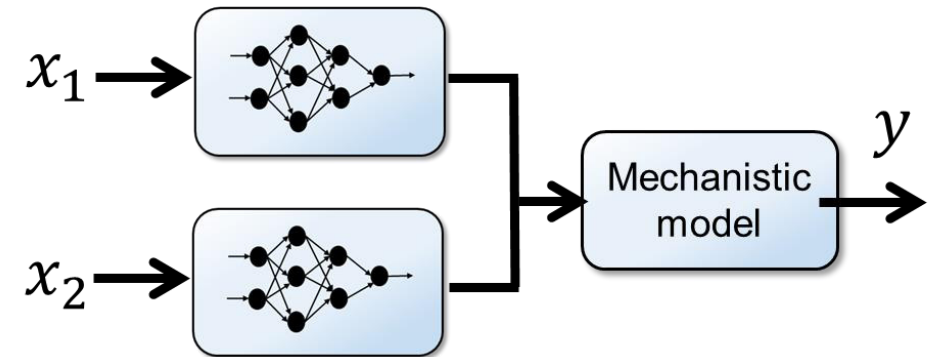
Hybrid modeling – Combining Mechanistic and Data-Driven Models



Single data-driven model

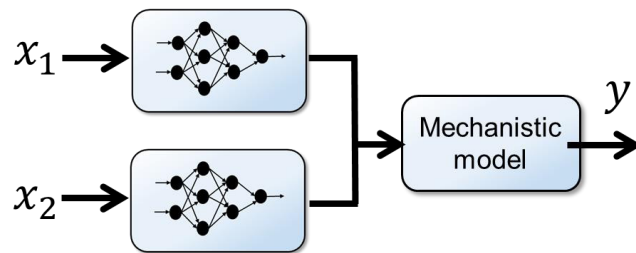


Hybrid model structure



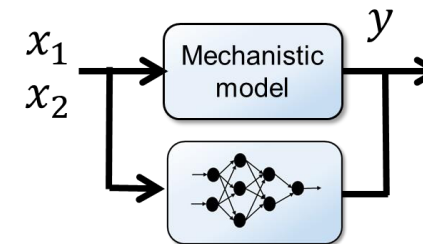
Training/Fitting of Hybrid Models

- Different training approaches possible, e.g.,
 - train hybrid model
 - use different data for each data-driven part
 - train separately mechanistic and data-driven model
- Different topologies possible, e.g.,:
 - serial and structured approach



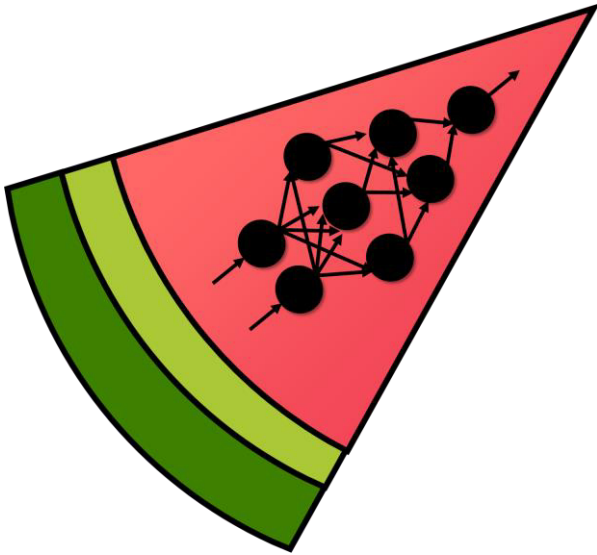
models relations within
the mechanistic model






parallel approach



compensates errors
of mechanistic model

MeLOn – Machine Learning Models for Optimization



- MeLOn provides machine learning models with training scripts and examples
- The models can be used within our open-source global solver **MAiNGO** or exported to GAMS
- Interfaces to  **Keras**  **TensorFlow**  **scikit-learn**  **MATLAB**  **python**
- Includes artificial neural networks, Gaussian processes,...

MeLOn

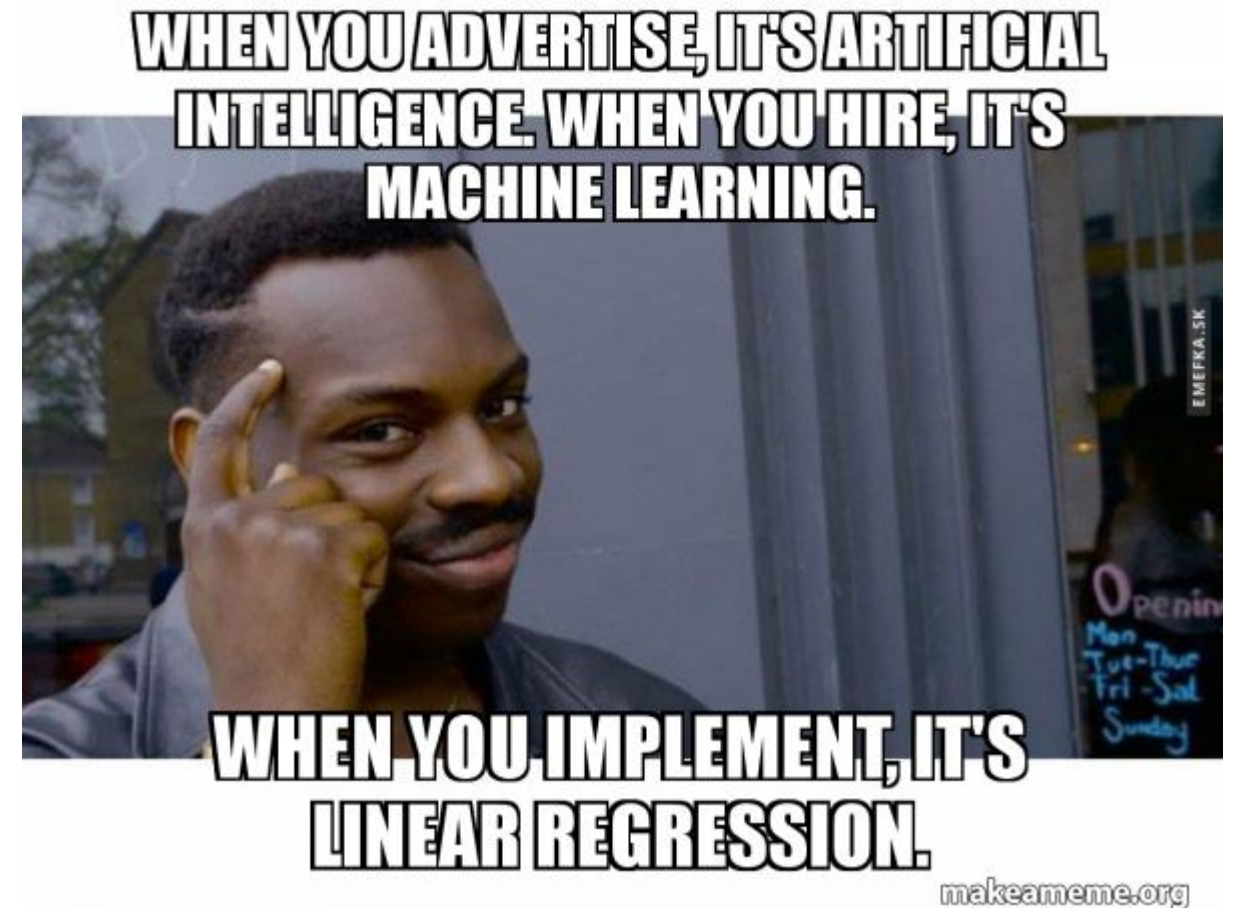
Open-source

<https://git.rwth-aachen.de/avt.svt/public/MeLOn>



Check Yourself

- What is machine learning?
- Can data-driven model extrapolate from their training data set?
- Why is it beneficial to use hybrid models?



<https://media.makeameme.org/created/when-you-advertise-f81897f53a.jpg>

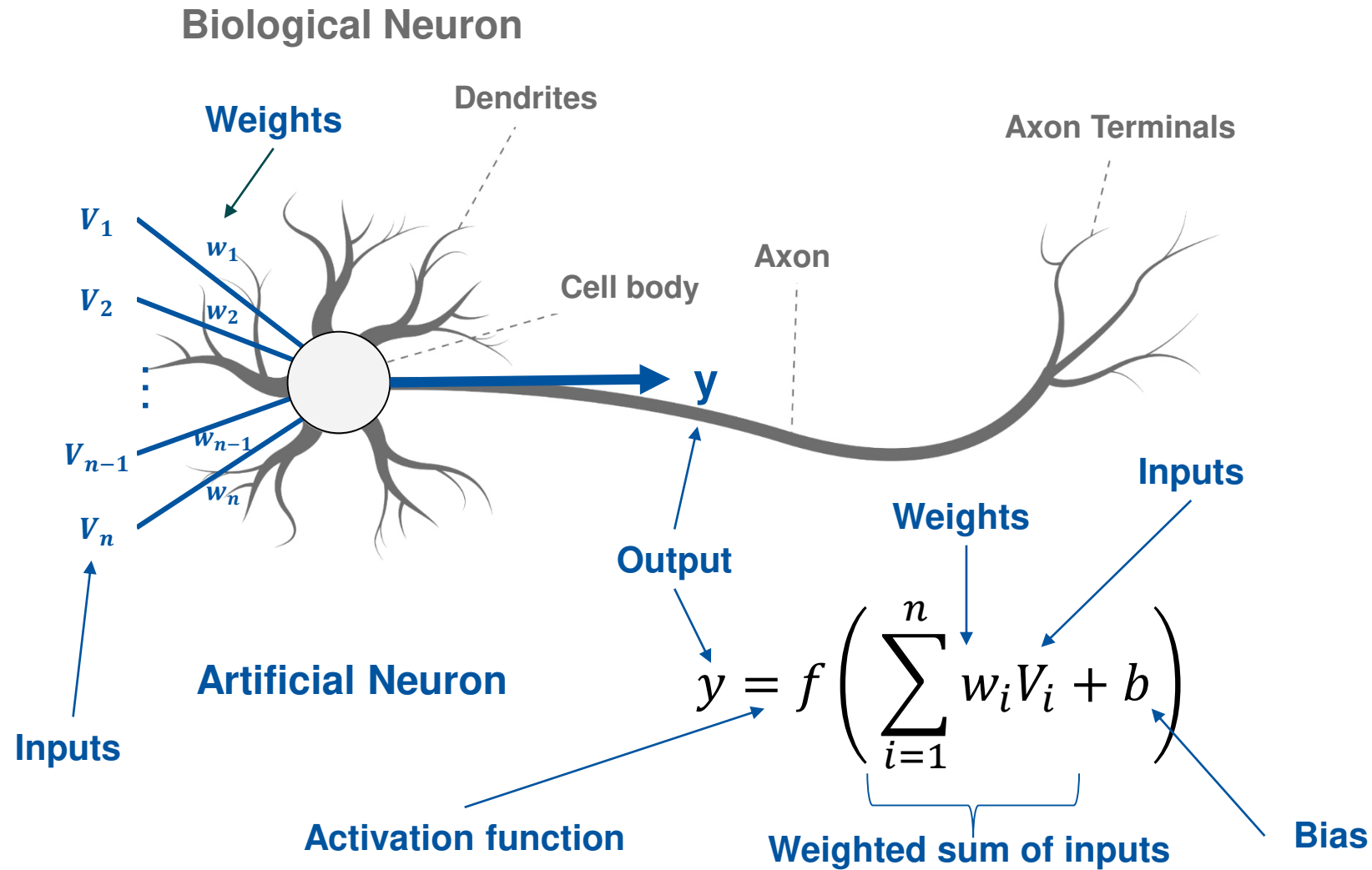


Applied Numerical Optimization

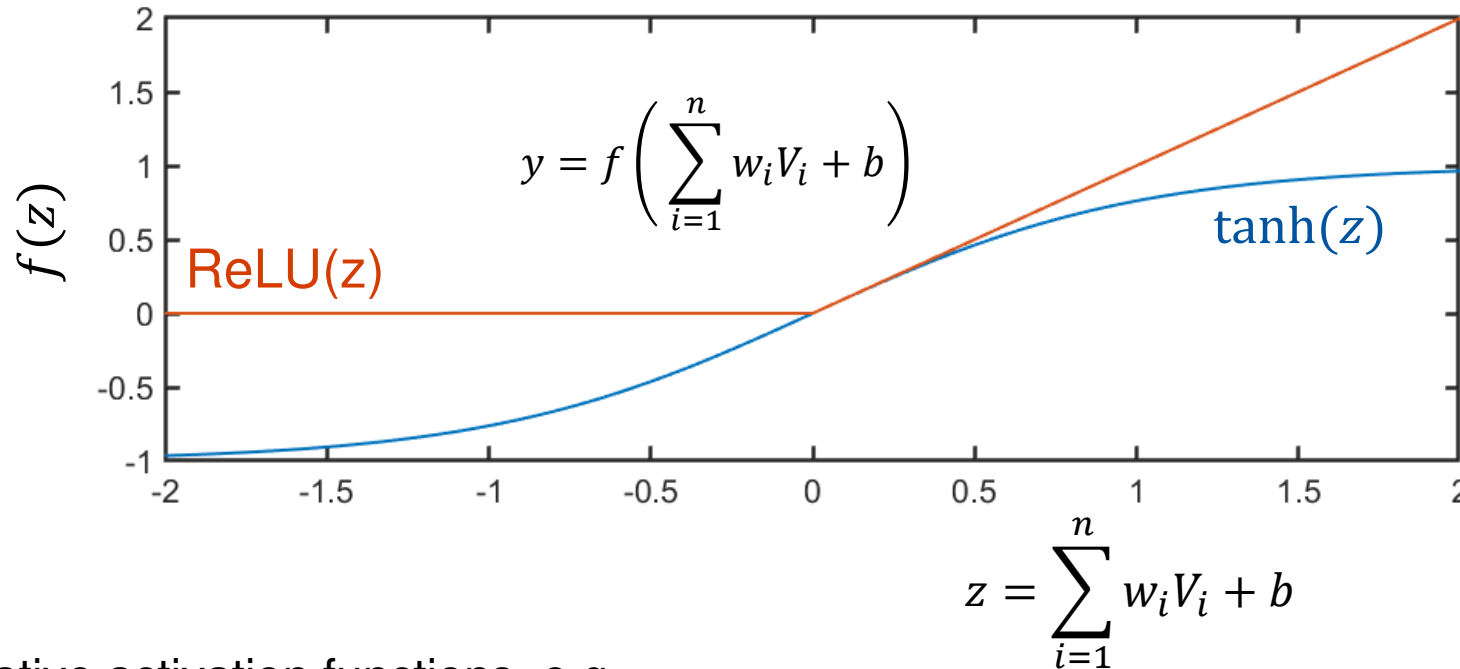
Prof. Alexander Mitsos, Ph.D.

Artificial Neural Networks

Biological and Artificial Neuron



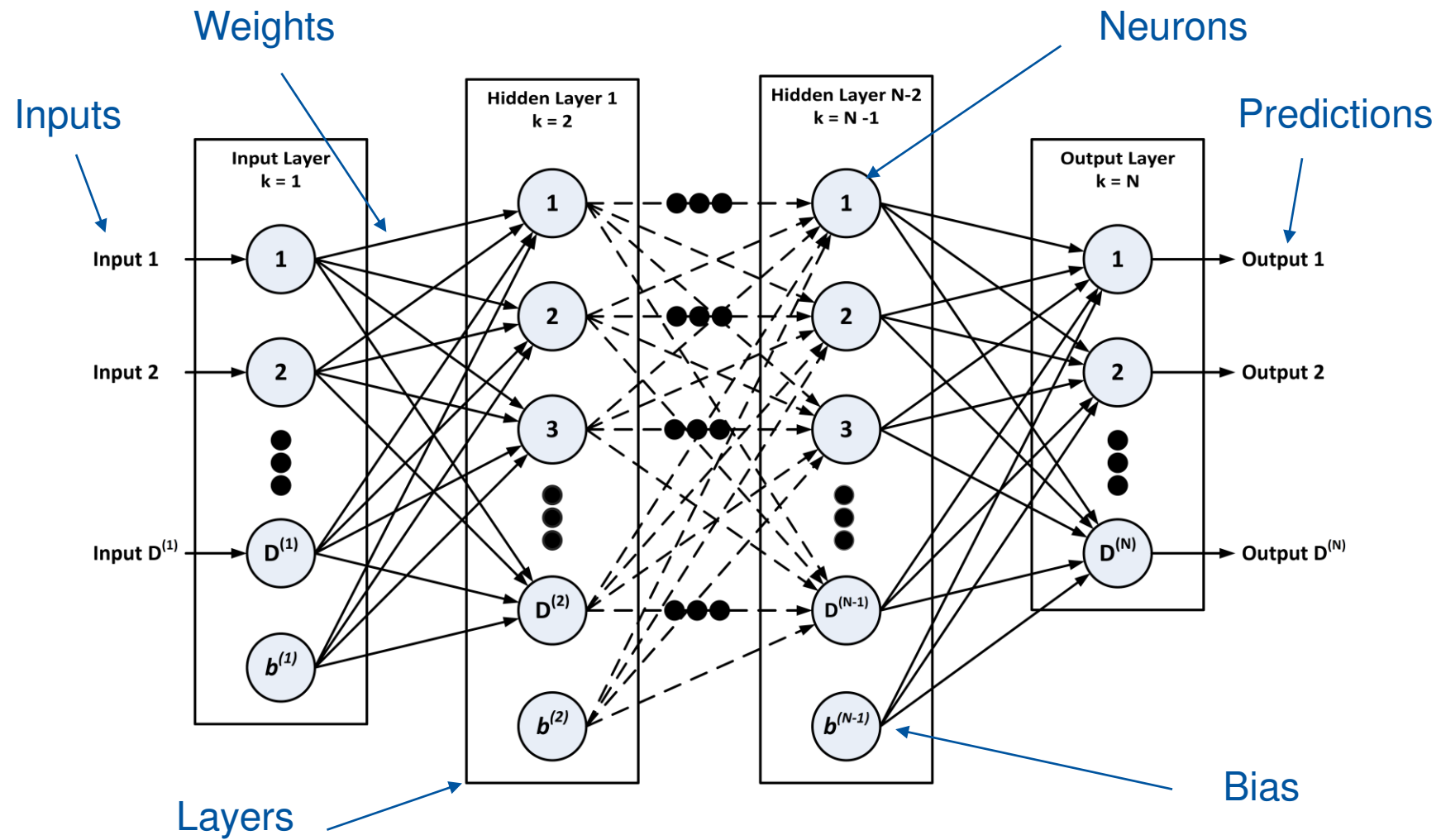
Artificial Neural Networks: Activation Functions



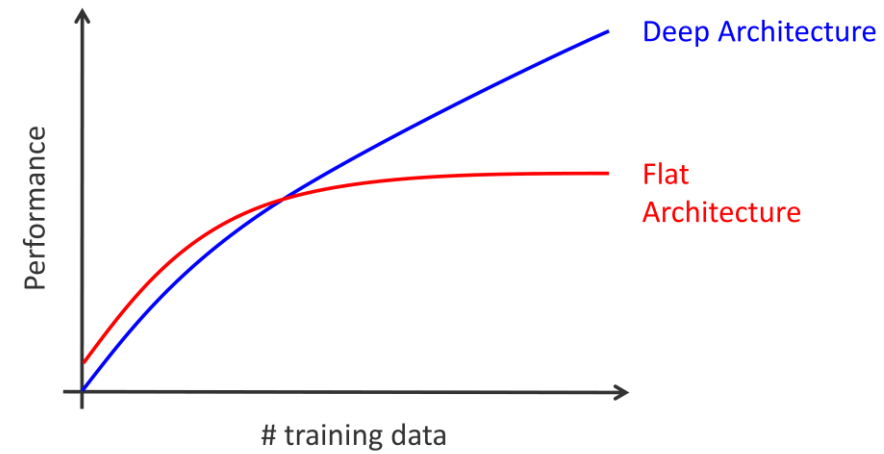
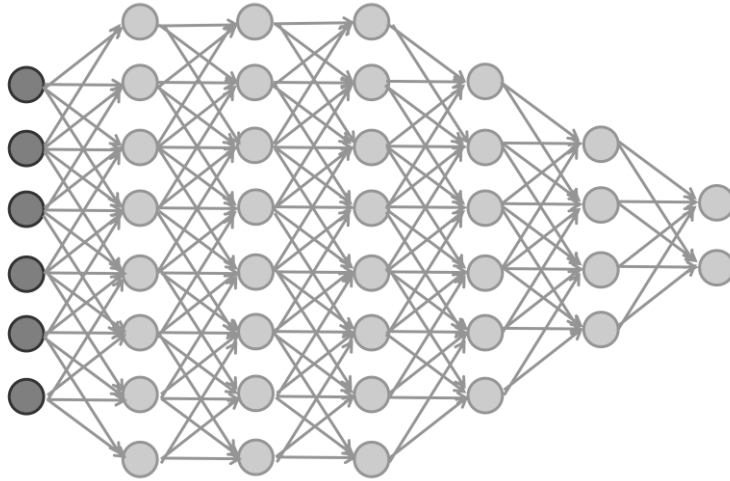
Several alternative activation functions, e.g.,

- Hyperbolic tangent function $f(z) = \tanh(z)$
 - Saturates quickly
- Rectified Linear Unit (ReLU) $f(z) = \max(0, z)$
 - Used in deep ANNs

Feed-Forward ANNs



Deep Learning



- Deep networks show better performance for large data sets, e.g., image recognition
 - AlexNet (2012): 8 (learning) layers with 60 mio parameters
 - Inception-v3 (2015): 47 (learning) layers with 25 mio parameters
- Success factors: e.g., GPU computing, ReLU activation function, dropout regularization

Check Yourself

- What are ANNs?
- What is deep learning?

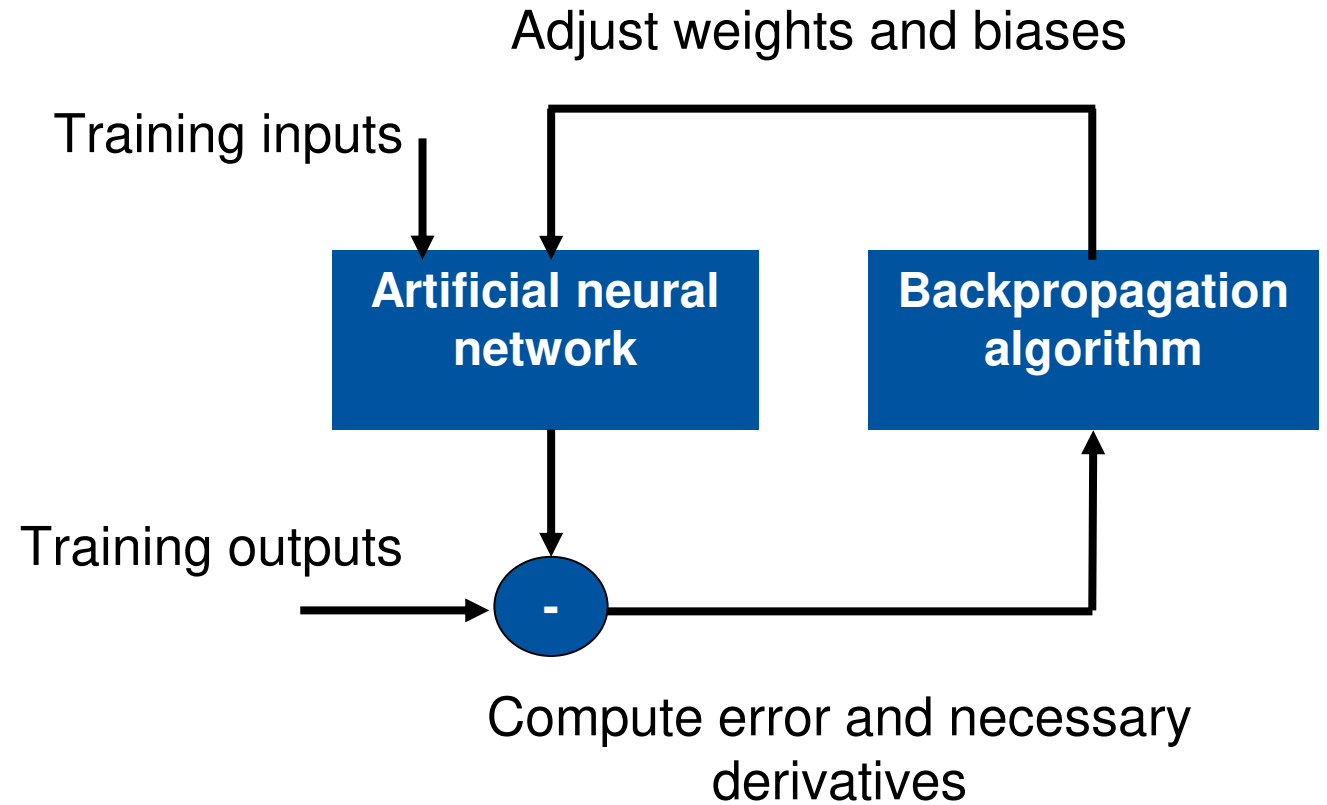


Applied Numerical Optimization

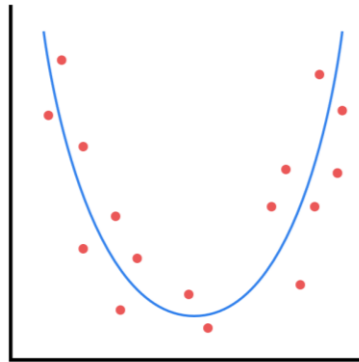
Prof. Alexander Mitsos, Ph.D.

Training of data-driven models: ANN

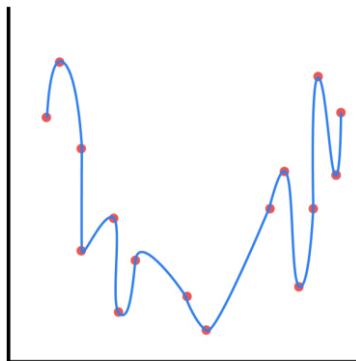
Training of ANNs



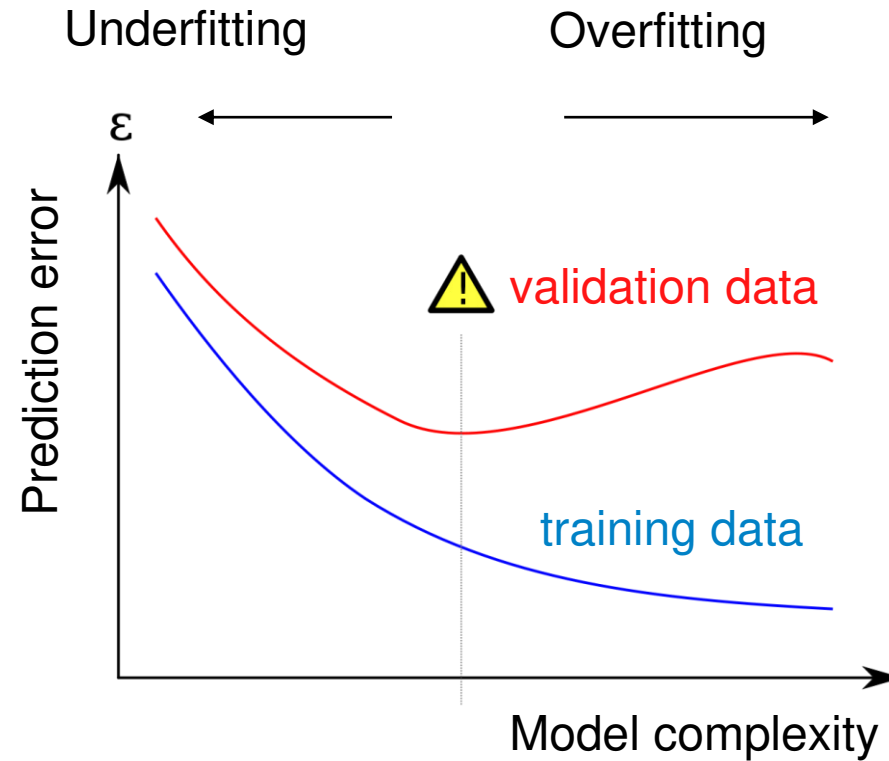
Overfitting & Underfitting



A well-fit model

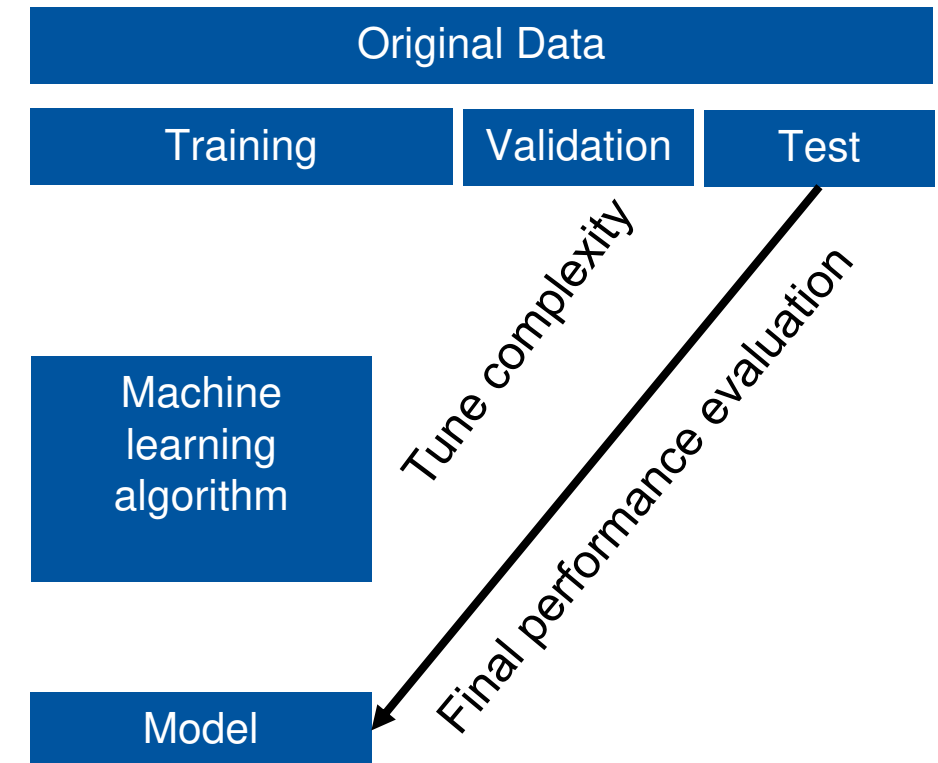


An overfit model



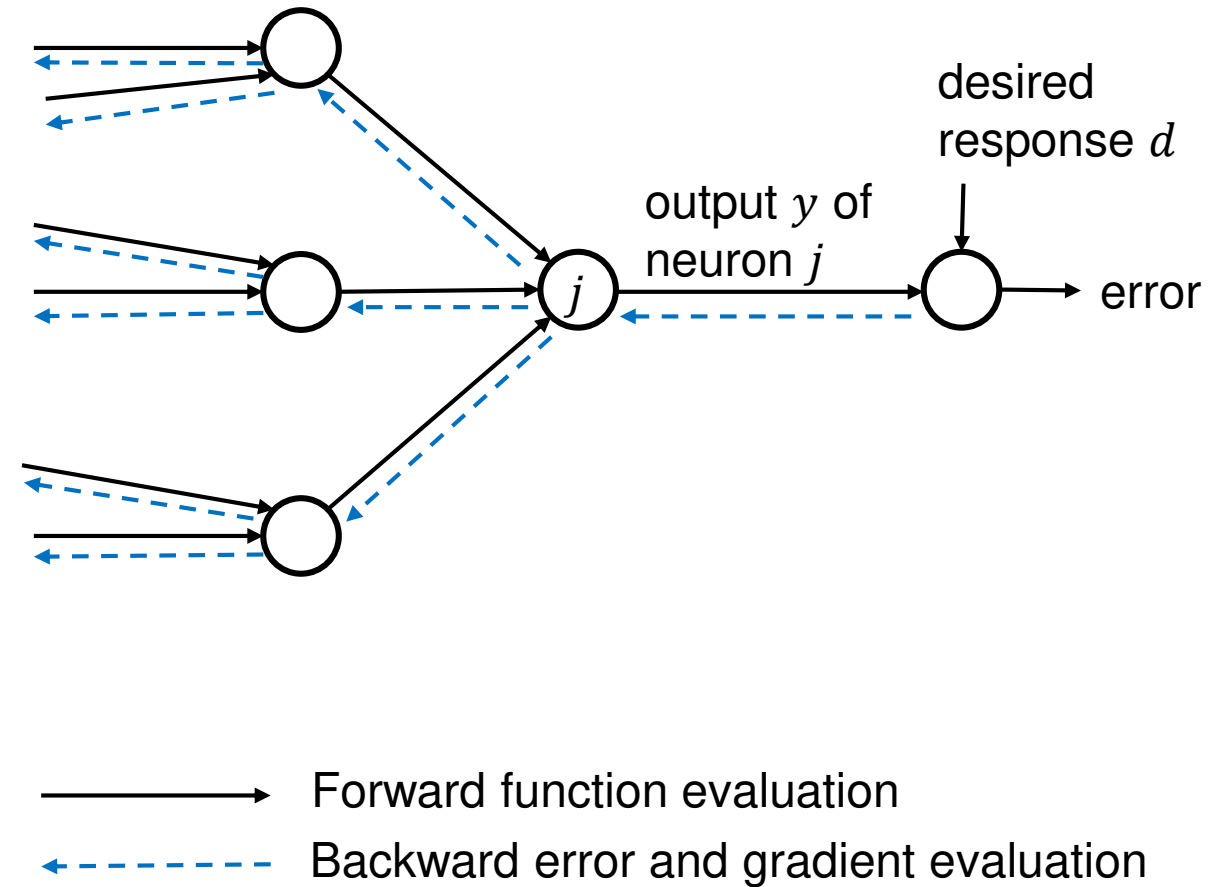
Training, Test and Validation Data

- Training: minimize error between model and data
 - Adjust weights and biases of neurons
 - Use of gradient-based updates beneficial
- Divide data randomly into 3 sets: training, validation and test
 - Predefined ratio, e.g., 70%:15%:15%
 - Fit weights and bias with the **training set**
 - Choose network complexity with the **validation set**
 - Evaluate performance with the **test set**
- Different language in ML vs. optimization
 - Training \leftrightarrow regression
 - loss function \leftrightarrow objective function
 - weights, biases \leftrightarrow decision variables (or parameters)
 - backpropagation algorithm \leftrightarrow gradient based method



Backpropagation Algorithm for Training ANNs

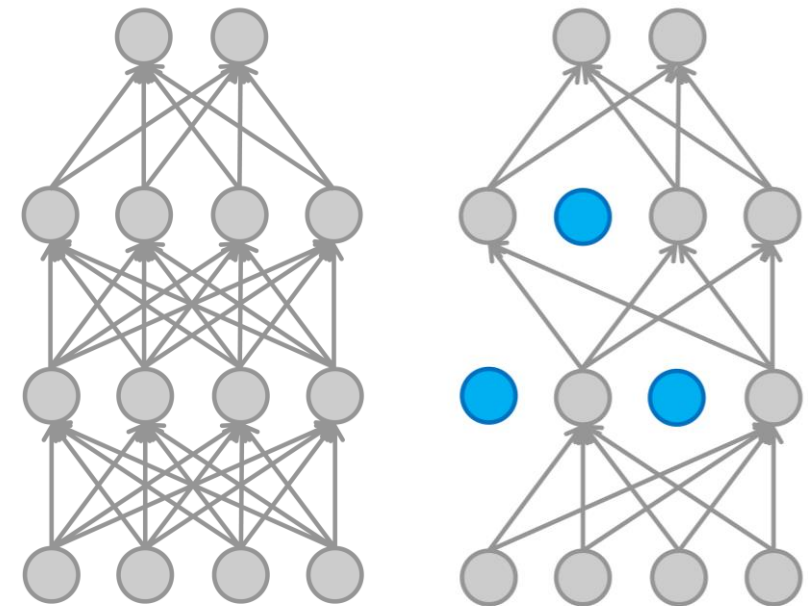
- Standard method for training
- Two passes of the ANN
 1. Evaluate functions in forward pass
 2. Evaluate gradients in backward pass
- Need chain rule: weights of neuron affect subsequent neurons



Regularization in Training

- Large neural networks include many parameters
→ risk of overfitting
- Regularization mitigates risk of overfitting, e.g.,
 - **Weight decay**: penalize large weights
 - **Dropout**: repeat training iterations with randomly dropped units

Dropout



Check Yourself

- How does the backpropagation algorithm work in principle?
- How is the available data split? What is done with each of the chunks?
- What regularization is used to avoid overfitting?

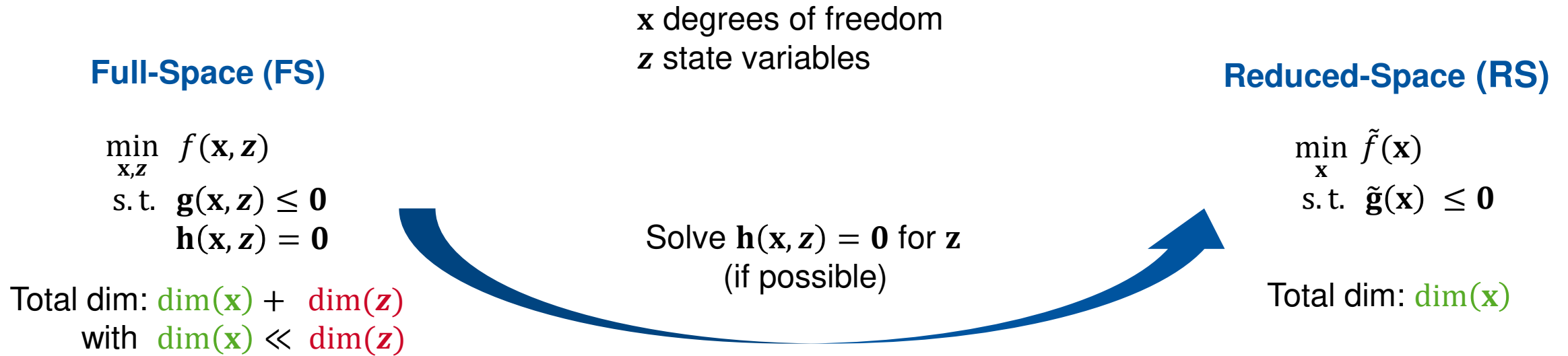


Applied Numerical Optimization

Prof. Alexander Mitsos, Ph.D.

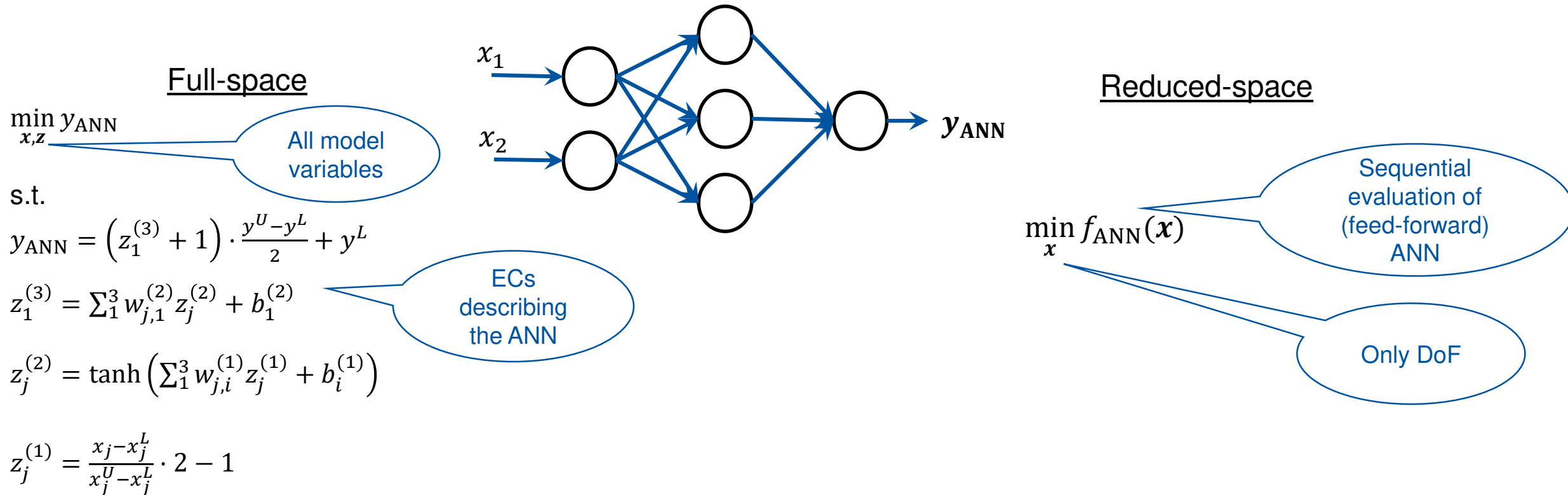
Deterministic global optimization with ANN embedded

Recap: Full-Space vs. Reduced-Space¹⁻⁵



[1] Epperly & Pistikopoulos, JOGO, 11(3), 287-311 (1997) [2] Byrne & Bogle, Ind. Eng. Chem. Res, 39(11), 4296-4301 (2000). [3] Mitsos, Chachuat & Barton, SIOPT, 20(2), 573-601 (2009) [4] Bongartz, & Mitsos, JOGO, 69(4), 761-796 (2017) [5] Bongartz and Mitsos, JOGO, 69(4), 761-796 (2018) [6] Schweidtmann, A. M., & Mitsos, A..JOTA, 180(3), 925-948 (2019).

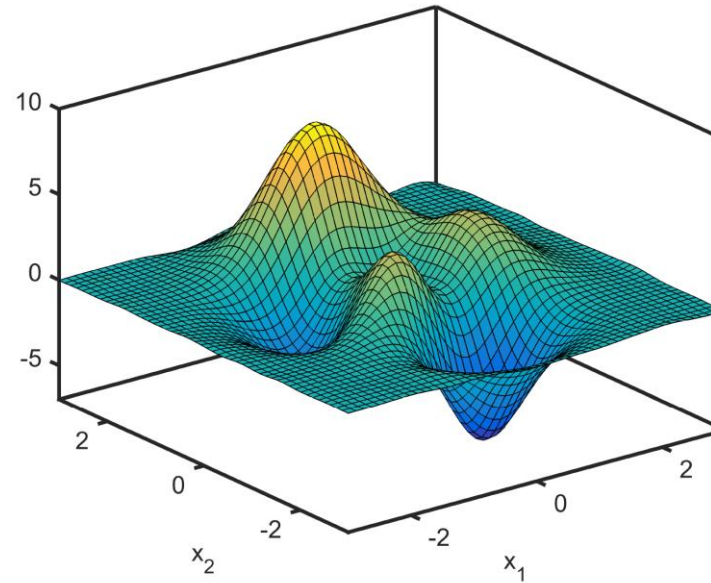
Reduced Space Formulation with ANNs-Embedded



- ☺ Suitable for standard solvers
- ☺ Simple functions
- ☹ Large-scale NLPs

- ☺ Small scale: $\dim(\mathbf{z}) \gg \dim(\mathbf{x})$
- ☹ Relaxations of f_{ANN} required

Example: The Peaks Function is learned by an ANN and Subsequently Optimized



- Train a feed-forward ANN
 - 500 Data points from peaks function
 - 1 hidden layer with 47 neurons
- Optimize the predictions of the ANN

Full space

101 optimization variables
99 equality constraints

Reduced-space

2 optimization variables
0 equality constraints

MAiNGO

161.35 seconds**

0.50 seconds**

Speedup by factor about 300

21.08 to $>10^5$ seconds*

G A M S BARON

*using different algebraic formulations of tanh activation function, ** using the implemented envelope of tanh

Example: Hybrid Model for Chemical Process Optimization

1. Decomposition

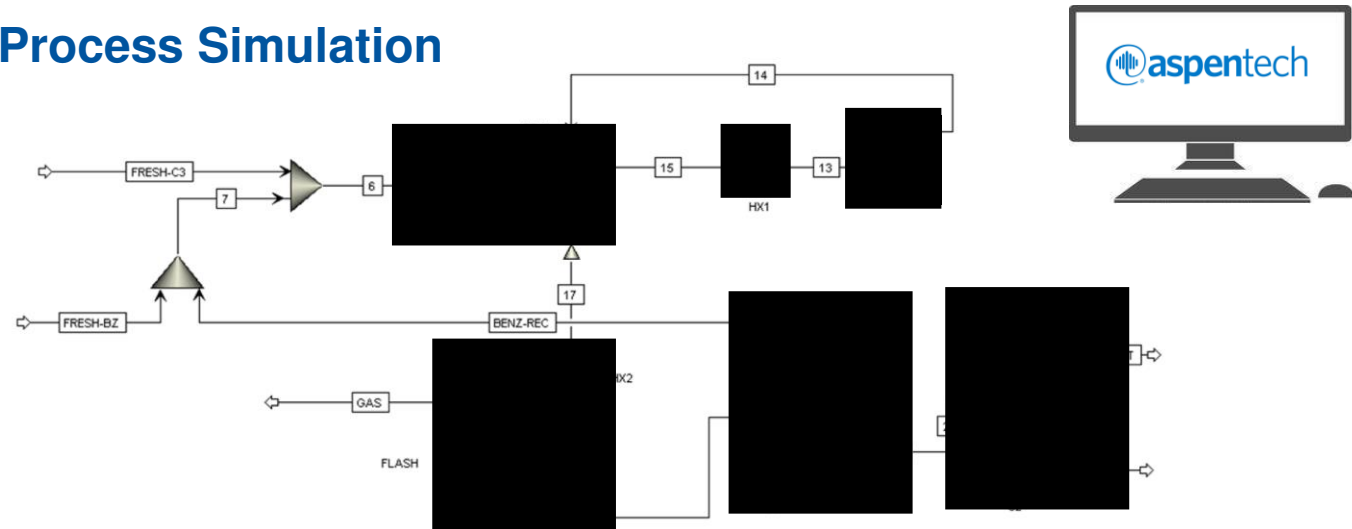
2. Data generation

3. Training

4. Hybrid model

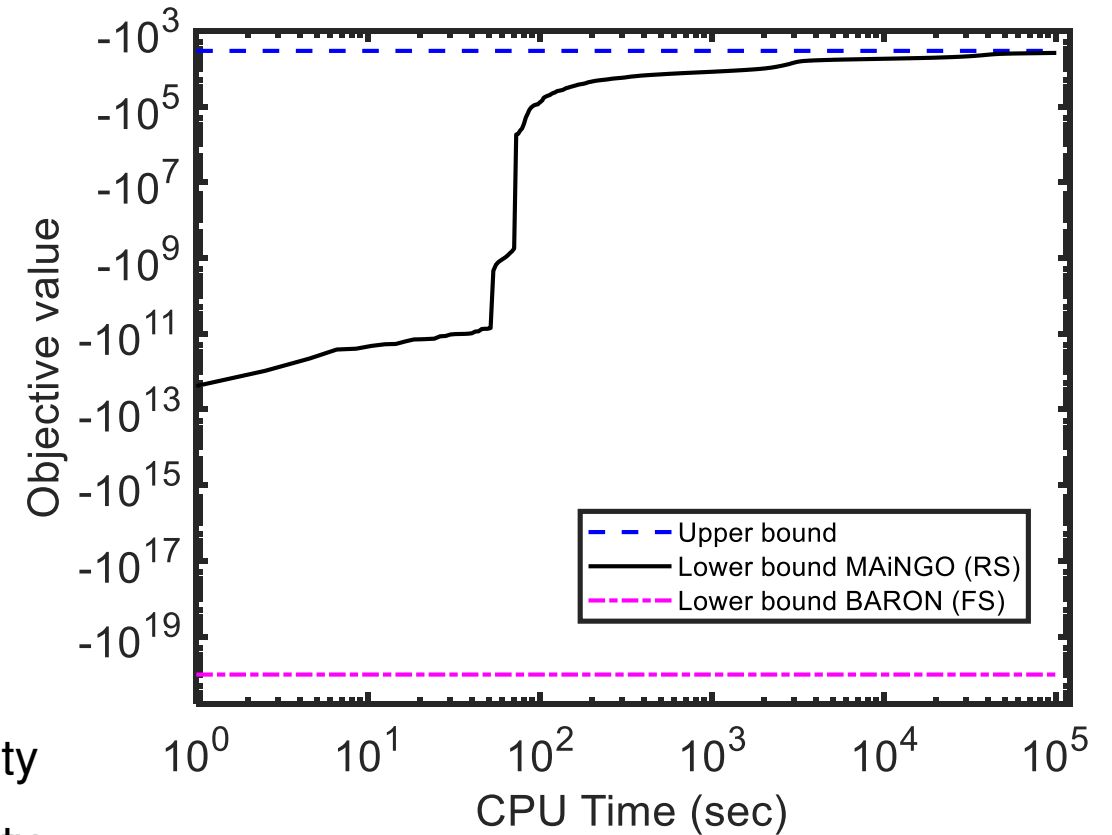
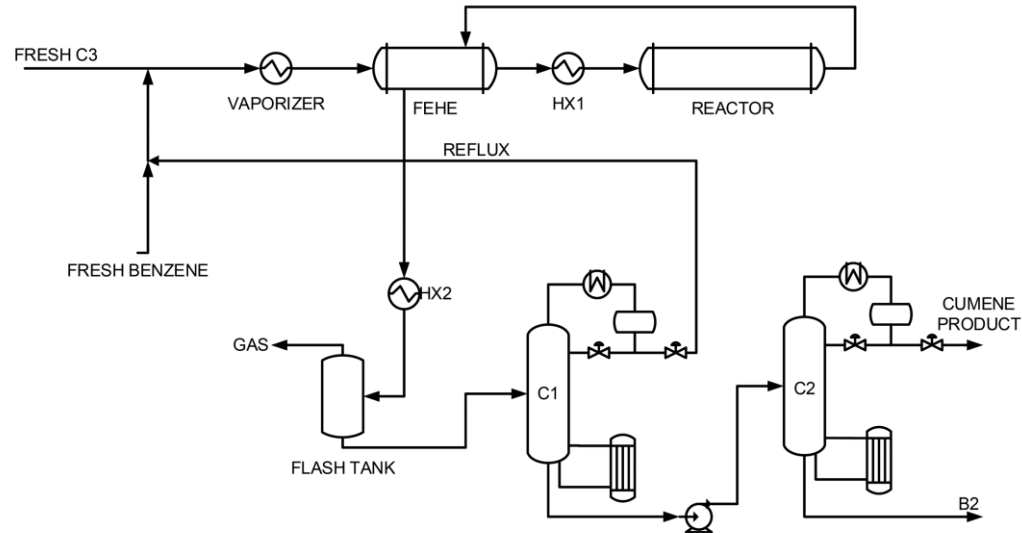
5. Optimization

Process Simulation



- Not single black box
- Hybrid process flowsheet including 14 ANNs
- ANNs connected via balance equations

Example: Hybrid Model for Chemical Process Optimization – Results



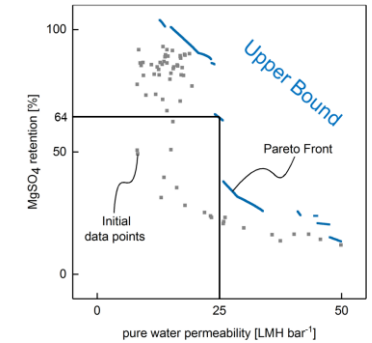
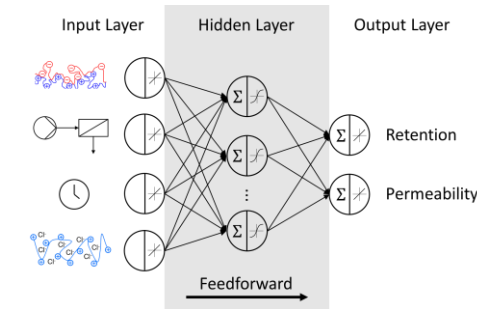
Problem size:

Full-space (FS): 794 variables, 789 equality, 1 inequality

Reduced-space (RS): 5 variables, 0 equality, 1 inequality

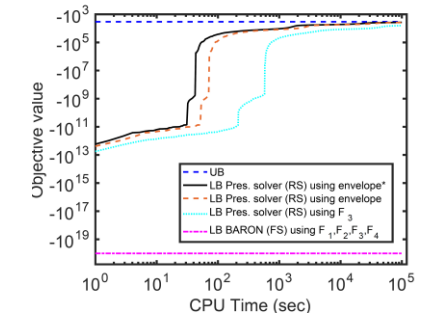
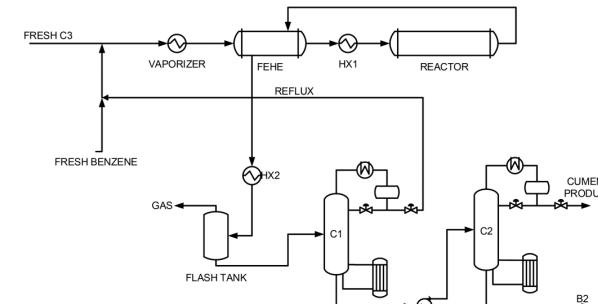
Rational design of ion separation membranes¹

- Learn membrane properties as a function of fabrication parameters from experimental data
- Multi-objective optimization to identify Pareto optimal membranes (retention and permeability)



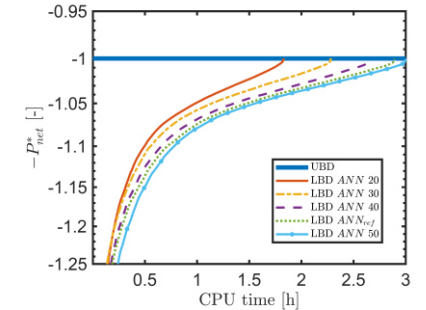
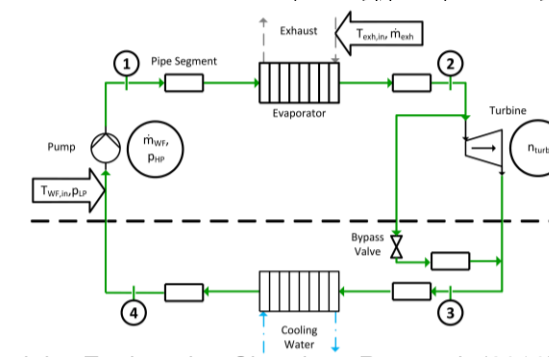
Cumene process optimization^{2,3}

- Unit operations simulated in AspenPlus
- Individual learning of unit operation
- Aggregation to hybrid process model



Learning single species thermodynamic properties⁴

- Organic Rankine cycle process optimization
- Learn thermodynamic properties from simulated data
- Faster optimization compared to Helmholtz EOS



[1] Rall et al., Journal of Membrane Science (2019) [2] Schweidtmann & Mitsos, JOTA (2019) [3] Schultz & Jorge, Industrial & Engineering Chemistry Research (2016)

[4] Schweidtmann, Huster, & Mitsos, Computers & Chemical Engineering (2019) [5] Bongartz et al., MAiNGO technical report (2018)

Check Yourself

- What is reduced space formulation for optimization with ANN embedded?
- Why is it beneficial to use reduced space formulation/ MAiNGO for global optimization with ANNs embedded?



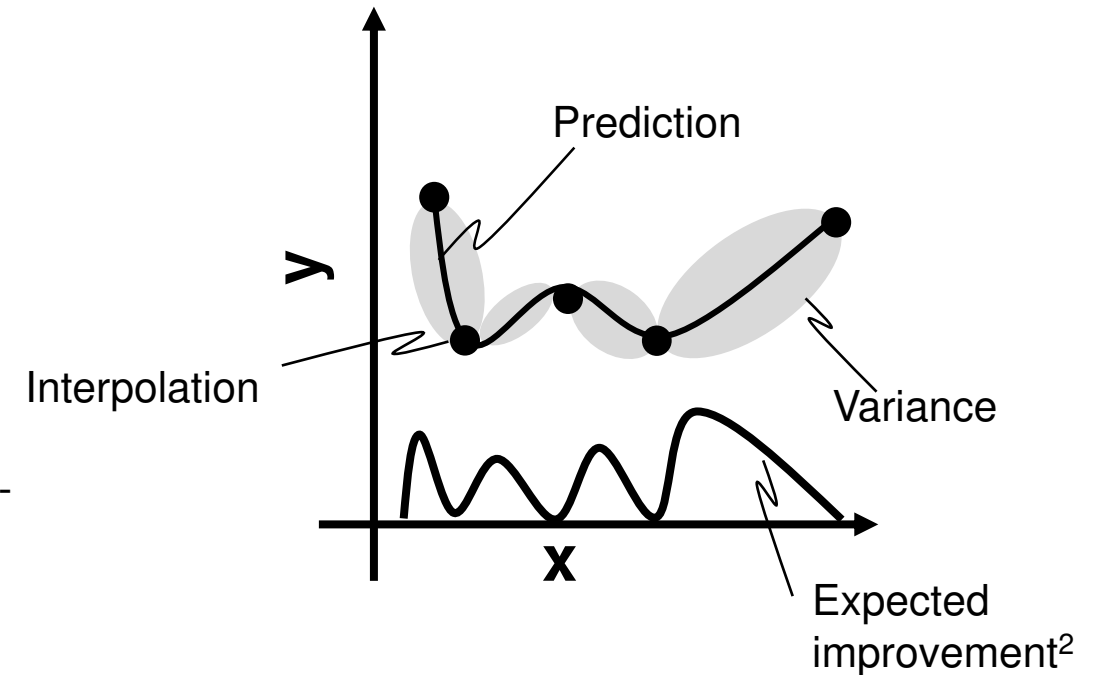
Applied Numerical Optimization

Prof. Alexander Mitsos, Ph.D.

Gaussian Processes (Kriging) and Bayesian Optimization

Motivation

- Most data-driven models provide only a prediction
 - No information about uncertainty of predictions
- Gaussian processes (GPs) a.k.a. Kriging¹
 - Generalization of multivariate Gaussian distributions
 - Nonlinear interpolation method
 - Provide prediction and variance
 - Well-suited for sparse data
 - Can be used in Bayesian optimization²
- Many applications for model-free optimization of expensive-to-evaluate functions
 - Process simulations³, NMPC for bioprocesses⁴, hardware-in-the-loop⁵, hyperparameter optimization in ML⁶



[1] Krige, D. G. (1951). Journal of the Southern African Institute of Mining and Metallurgy, 52(6), 119-139. [2] Jones, Schonlau, & Welch, *Journal of Global Optimization* (1998). 13(4). 455-492. [3] Helmdach, D., Yaseneva, P., Heer, P. K., Schweidtmann, A. M., & Lapkin, A. A. (2017). ChemSusChem, 10(18), 3632-3643. [4] Bradford, E., Schweidtmann, A. M., Zhang, D., Jing, K., & del Rio-Chanona, E. A. (2018). Computers & Chemical Engineering, 118, 143-158. [5] Schweidtmann, A. M., Clayton, A. D., Holmes, N., Bradford, E., Bourne, R. A., & Lapkin, A. A. (2018). Chemical Engineering Journal, 352, 277-282. [6] Snoek, J., Larochelle, H., & Adams, R. P. (2012). In Advances in neural information processing systems (pp. 2951-2959).

Optimization With Gaussian Processes Embedded is Difficult

Full-Space

$$\min_{\mathbf{x}_*, \bar{\mathbf{y}}_*, \mathbf{K}_*, k, r \dots} \bar{\mathbf{y}}_*$$

$$\text{s.t. } \bar{\mathbf{y}}_* = \mathbf{K}_* \cdot \mathbf{K}^{-1} \cdot \mathbf{y}$$

Training
targets

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

Covariance
matrix ($N \times N$)

$$\mathbf{K}_* = [k(\mathbf{x}_*, \mathbf{x}_1) \quad k(\mathbf{x}_*, \mathbf{x}_2) \quad \dots \quad k(\mathbf{x}_*, \mathbf{x}_N)]$$

Covariance
function

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 e^{-\frac{r_{ij}}{l}}, k \in [k^l, k^u]$$

$$r_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\| \quad r_{ij} \in [r_{ij}^l, r_{ij}^u] \quad i, j \in [1, N]$$

$$\mathbf{x}_* \in [\mathbf{x}^l, \mathbf{x}^u] \subset \mathbb{R}^D$$

Euclidian
distance

Total dim:

Variables: $\mathcal{O}(N + D)$

usually

Equalities: $\mathcal{O}(N + D)$

$N \gg D$

D : Dimension of training points

N : Number of training points

→ Gaussian processes lead to large-scale nonlinear problems in the full-space formulation

Gaussian Processes in a Reduced-Space

Full-Space

$$\min_{\mathbf{x}_*, \bar{\mathbf{y}}_*, \mathbf{K}_*, k, r \dots} \bar{\mathbf{y}}_*$$

$$\text{s.t. } \bar{\mathbf{y}}_* = \mathbf{K}_* \cdot \mathbf{K}^{-1} \cdot \mathbf{y}$$

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

$$\mathbf{K}_* = [k(\mathbf{x}_*, \mathbf{x}_1) \quad k(\mathbf{x}_*, \mathbf{x}_2) \quad \dots \quad k(\mathbf{x}_*, \mathbf{x}_N)]$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 e\left(-\frac{r_{ij}}{l}\right), k \in [k^l, k^u]$$

$$r_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\| \quad r_{ij} \in [r_{ij}^l, r_{ij}^u] \quad i, j \in [1, N]$$

$$\mathbf{x}_* \in [\mathbf{x}^l, \mathbf{x}^u] \subset \mathbb{R}^D$$

Total dim: # Variables: $\mathcal{O}(N + D)$; # Equalities: $\mathcal{O}(N + D)$

- Solve equalities for \mathbf{y}_{out}
- Propagate McCormick relaxations and (sub) gradient through code

Reduced-space

$$\min_{\mathbf{x}_*} \mathcal{GP}(\mathbf{x}_*)$$

$$\text{s.t. } \mathbf{x}_* \in [\mathbf{x}^l, \mathbf{x}^u] \subset \mathbb{R}^D$$

Sequential evaluation of Gaussian process

Only DoF

Total dim:
Variables: D
Equalities: 0

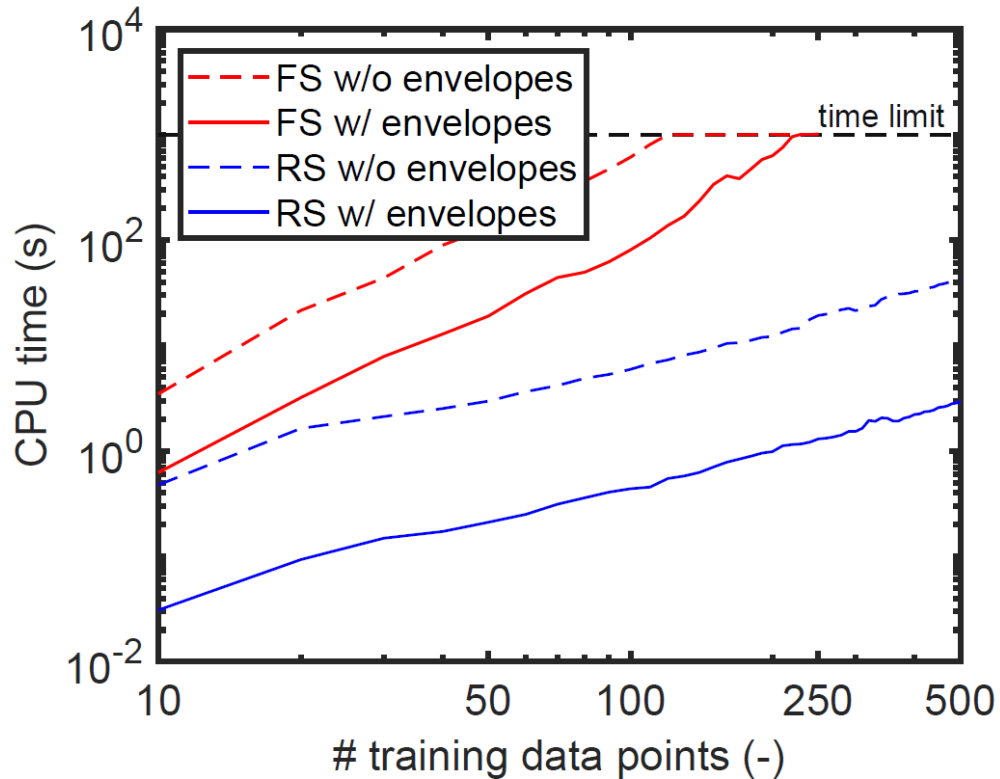
→ The reduced-space formulation reduces problem size drastically

D : Dimension of training points

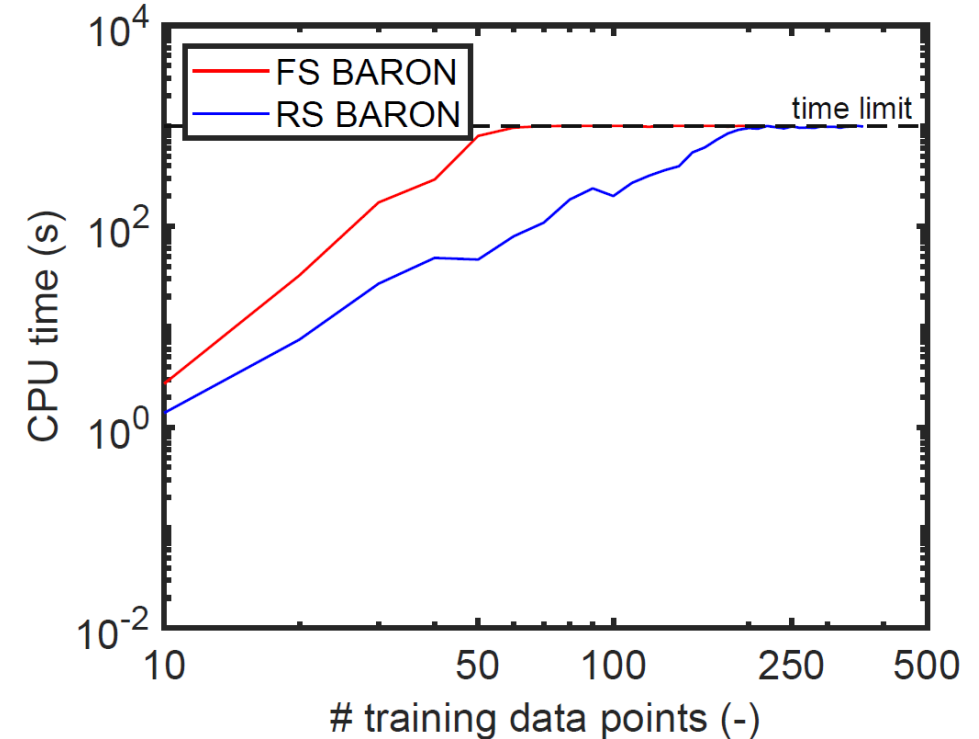
N : Number of training points

Scaling of Computational Performance with Training Set Size on Peaks Function

MAiNGO



BARON



→ The reduced-space formulation outperforms the full-space formulation for Gaussian processes

[1] Schweidtmann, A. M., Bongartz, D., Grothe, D., Kerkenhoff, T., Lin, X., Najman, J., & Mitsos, A. (2020). Global Optimization of Gaussian processes. arXiv preprint arXiv:2005.10902.

Mixed-integer Bayesian Optimization of Membrane Synthesis¹

Mixed-Integer optimization problem:

max Expected Improvement_{Na₂SO₄ Retention}

s.t. $x_{LB,i} \leq x_i \leq x_{UB,i} \quad i = 1, \dots, 3$

2 continuous variables:

sodium chloride concentration: $c_{NaCl} \in [0, 0.5] \text{gL}^{-1}$

deposited polyelectrolyte mass: $m_{PE} \in [0, 5] \text{gm}^{-2}$

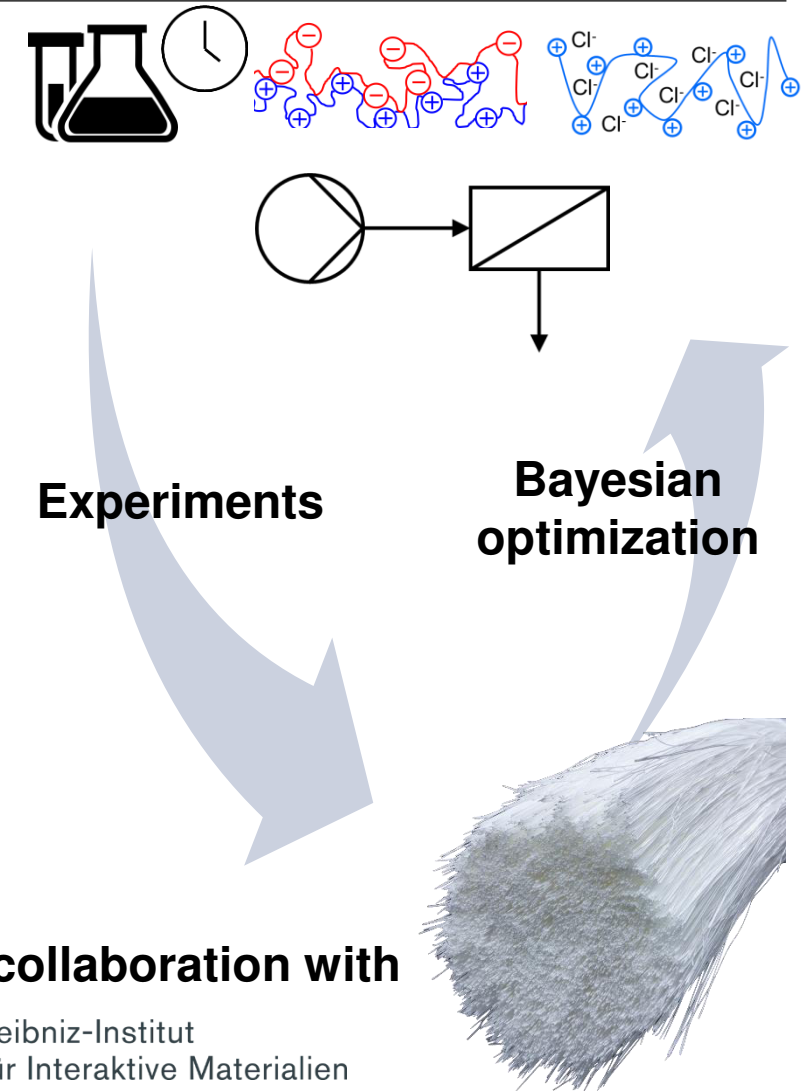
1 integer variable:

number of layers: $N_{Layers} \in \{1, 2, 3 \dots 10\}$

Global solution found within 7 minutes

- $c_{NaCl} = 0.362 \text{gL}^{-1}$, $m_{PE} = 0 \text{gm}^{-2}$, $N_{Layers} = 4$
- Expected retention 85.32 with standard deviation $\sigma = 14.8$

[1] Schweidtmann, Bongartz, Grothe, Kerkenhoff, Lin, Najman, Mitsos (2020). <https://arxiv.org/abs/2005.10902>



Check Yourself

- What are GP?
- What is Bayesian optimization?
- Why is it beneficial to use MAiNGO for global optimization with ANNs/GP embedded?