

# Supervised machine learning with tensor-based kernel machines

Kim Batselier, Afra Kilic, Eva Memmel,  
Clara Menzen, Albert Saiapin, Frederiek Wesel



## Setting the stage...

### Parametric supervised learning problem:

Given a set of measurements  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$  and a parameterization  $\mathbf{w}$  of  $\hat{y}_n(\cdot)$  in

$$y_n = \hat{y}_n(\mathbf{x}_n | \mathbf{w}) + e_n,$$

determine  $\mathbf{w}$ .

## Kernel machines

$$\hat{y}_n(\mathbf{x}_n|\mathbf{w}) = \varphi(\mathbf{x}_n)^T \mathbf{w}$$

$\varphi(\cdot)$ : collection of basis functions

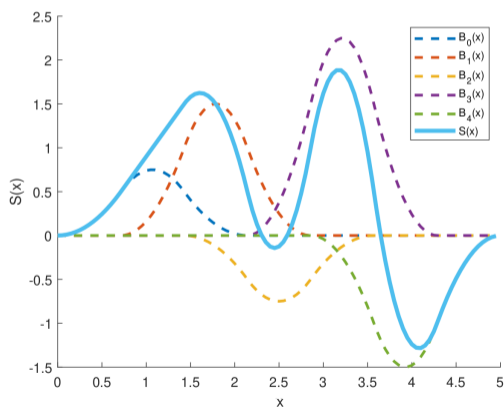
$$\text{OE: } \hat{y}_n = (u_n \quad \cdots \quad \hat{y}_{n-1} \quad \cdots) \mathbf{w}$$

$$\text{ARMAX/ Box Jenkins: } \hat{y}_n = (u_n \quad \cdots \quad y_{n-1} \quad \cdots \quad \hat{y}_{n-1} \quad \cdots) \mathbf{w}$$

⇒ Learn  $\mathbf{w}$  from data = solving multivariate polynomial system

## Pick your favorite basis function

$$\hat{y}(x) = w_0 B_0(x) + w_1 B_1(x) + w_2 B_2(x) + w_3 B_3(x) + \dots$$



## A simple univariate example

$$\begin{aligned}\hat{y}(x) &= w_0 + w_1 x + w_2 x^2 \\ &= \underbrace{(1 \quad x \quad x^2)}_{\varphi(x)^T} \underbrace{\begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}}_{\mathbf{w}}\end{aligned}$$

## Another simple univariate example

Ground truth:

$$y = f(x) + e = 3.2 x^2 - 9.6 \cos(5.1 x^3 - e^{-3x}) + e$$

Assume  $x \in [-1, 1]$ . How many basis functions  $I$  do we need?

$$f(x) = 3.2 x^2 - 9.6 \cos(5.1 x^3 - e^{-3x}) \approx \sum_{i=1}^I \varphi_i(x) w_i$$

The total number of basis functions  $I$  we need depends on  $\varphi_i(x)$  and  $f(x)$ .

## Multivariate basis functions as products of univariate basis functions

$$\underbrace{\varphi(\mathbf{x})}_{I^D} = \underbrace{\varphi(x_1)}_I \otimes \underbrace{\varphi(x_2)}_I \otimes \cdots \otimes \underbrace{\varphi(x_D)}_I$$



## Bivariate monomial basis functions from a Kronecker product

$$\varphi(\mathbf{x}) = \varphi(x_1) \otimes \varphi(x_2)$$

$$\begin{pmatrix} 1 \\ x_1 \\ x_1^2 \\ x_2 \\ x_1x_2 \\ x_1^2x_2 \\ x_2^2 \\ x_1x_2^2 \\ x_1^2x_2^2 \end{pmatrix} = \begin{pmatrix} 1 \\ x_1 \\ x_1^2 \end{pmatrix} \otimes \begin{pmatrix} 1 \\ x_2 \\ x_2^2 \end{pmatrix}$$

# Learning problem

$$\mathbf{y} = \Phi \mathbf{w} + \mathbf{e}$$

$$\min_{\mathbf{w}} L(\mathbf{y}, \Phi, \mathbf{w})$$

Possible loss functions  $L(\mathbf{y}, \Phi, \mathbf{w})$ :

- logistic (logistic regression)
- hinge (support vector machines / classification)
- Vapnik  $\epsilon$ -insensitive (support vector machines / regression)
- squared (least-squares support vector machines, Gaussian Processes)

## Squared loss

$$\min_{\mathbf{w}} \|\mathbf{y} - \Phi \mathbf{w}\|_2^2$$

$$\mathbf{w} = \underbrace{(\Phi^T \Phi)^{-1}}_{I^D \times I^D} \Phi^T \mathbf{y}$$

### Assumptions:

- $N \geq I^D$
- Persistency of excitation:  $\text{rank}(\Phi) = I^D$

What if  $\text{rank}(\Phi) < I^D$ ?

Kernel ridge regression

$$\min_{\mathbf{w}} \|\mathbf{y} - \Phi \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

$$\mathbf{w} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y}$$

## A simple example

- $D = 18$  variables  $x_1, \dots, x_{18}$
- $I = 20$  basis functions
- $w$  has  $20^{18} = 2.6 \cdot 10^{23}$  entries  $\approx 10^{10}$  ChatGPTs

Tensor networks allow you to **efficiently** solve

$$(\Phi^T \Phi + P) w = \Phi^T y$$

**without** explicitly making / inverting  $(\Phi^T \Phi + P)$ .

# A tensor is a high-dimensional generalization of vectors and matrices



vector

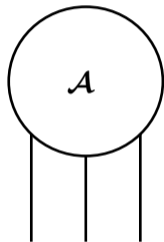
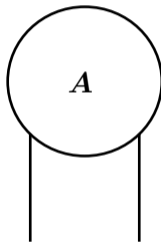
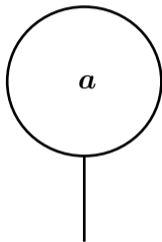
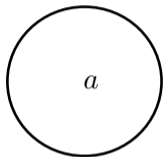


matrix



tensor

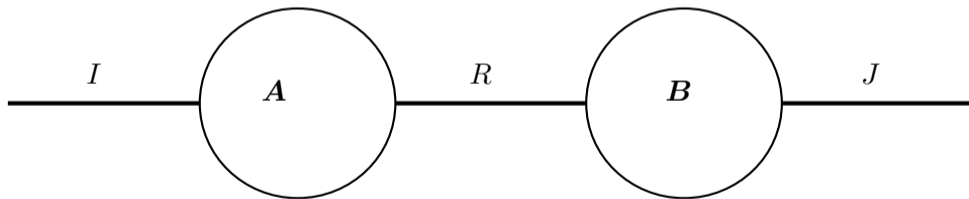
# Tensor diagrams



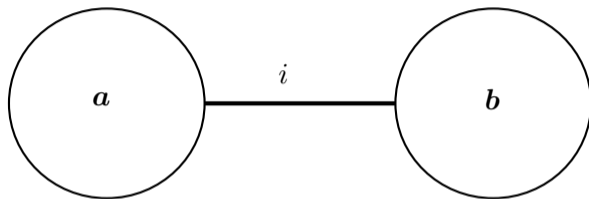


Index summation is visualized as a common edge

$$C(i, j) = \sum_{r=1}^R A(i, r) B(r, j)$$



## Brain teaser



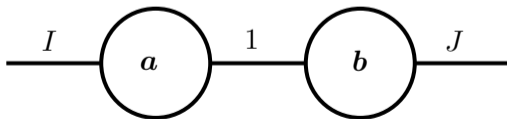
inner product: 
$$\sum_{i=1}^I \mathbf{a}(i) \mathbf{b}(i) = \mathbf{a}^T \mathbf{b}$$

An important operation: Outer product of vector  $\mathbf{a}$  with vector  $\mathbf{b}$

$$\mathbf{C} = \mathbf{a} \mathbf{b}^T = \mathbf{a} \circ \mathbf{b}$$

$$C(i, j) = \mathbf{a}(i) \mathbf{b}(j)$$

$$C(i, j) = \sum_{r=1}^1 \mathbf{a}(i, r) \mathbf{b}(r, j)$$

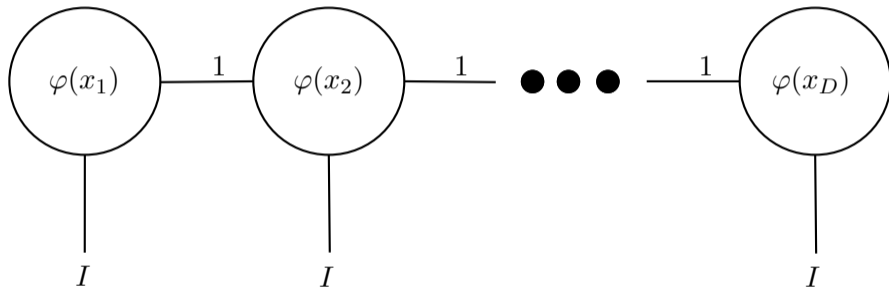


## Key idea

Replacing  $\Phi, P, w$  by networks of much smaller tensors avoids explicit construction/inversion of  $\Phi, P, w$  and enables **efficient** learning.

## Back to our kernel machines

$$\phi = \varphi(x_1) \circ \varphi(x_2) \circ \dots \circ \varphi(x_D)$$



$$I^D \rightarrow D I$$

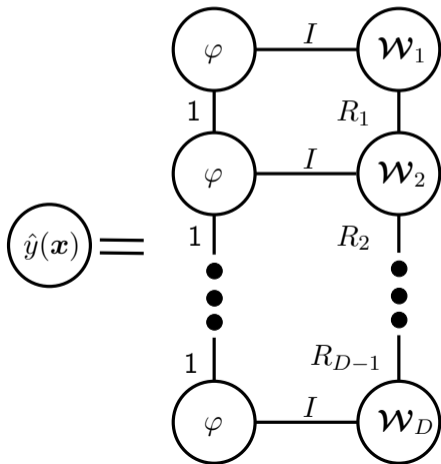
## Revisit bivariate monomial basis functions as a tensor product

$$\phi = \varphi(x_1) \circ \varphi(x_2)$$

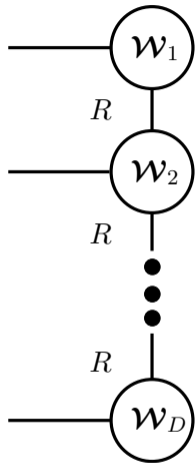
$$\begin{pmatrix} 1 & x_2 & x_2^2 \\ x_1 & x_1x_2 & x_1x_2^2 \\ x_1^2 & x_1^2x_2 & x_1^2x_2^2 \end{pmatrix} = \begin{pmatrix} 1 \\ x_1 \\ x_1^2 \end{pmatrix} \circ \begin{pmatrix} 1 \\ x_2 \\ x_2^2 \end{pmatrix}$$

# Tensor-based kernel machine

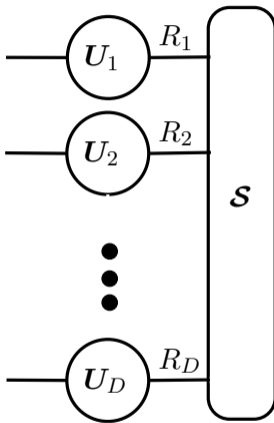
$$y = \langle \phi, \mathbf{w} \rangle + e$$



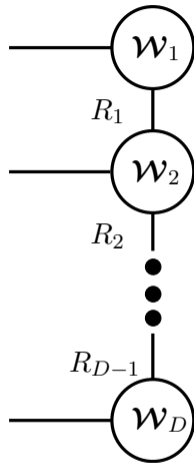
## Possible tensor network topologies...



Polyadic  
 $DIR$



Tucker  
 $R^D + DIR$



Tensor Train  
 $DIR^2$



## From linear least-squares to nonlinear least-squares

$$\min_{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_D} \|\mathbf{y} - \langle \phi, \text{TN}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_D) \rangle\|_2^2 + R(\text{TN}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_D))$$

Solving a polynomial system of equations!

## Exploiting tensor network structure enables efficient learning

$\text{TN}(\mathcal{W}_1, \mathcal{W}_2, \dots, \mathcal{W}_D)$  is a linear map in  $\mathcal{W}_1$

$$\text{TN}(\mathcal{W}_1, \mathcal{W}_2, \dots, \mathcal{W}_D) \rightarrow \text{TN}_1(\mathcal{W}_1) = \mathbf{T}_1 \mathbf{w}_1$$

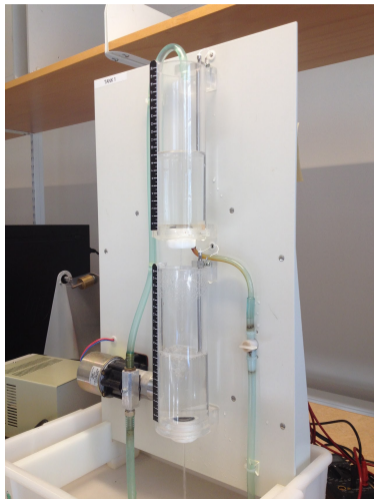
$$\min_{\mathbf{w}_1} \|\mathbf{y} - \Phi \mathbf{T}_1 \mathbf{w}_1\|_2^2 + \mathcal{R}(\mathbf{T}_1 \mathbf{w}_1)$$

Computational cost polyadic network:  $O(N(IR)^2 + (IR)^3)$

## Example: Cascaded tanks

- $u_n$  : pump voltage
- $y_n$  : water level lower tank
- Task: learn

$$\hat{y}_n(u_{n-1}, u_{n-2}, \dots, y_{n-1}, y_{n-2}, \dots)$$



## Example: SUSY binary classification problem

$I = 20, D = 18, N = 5,000,000$

| Technique            | AUC                                 |                                     | 1-Accuracy (%)   | Time (s)       |
|----------------------|-------------------------------------|-------------------------------------|------------------|----------------|
|                      | Low-level                           | Complete                            | Complete         | Complete       |
| BDT                  | $0.850 \pm 0.003$                   | $0.863 \pm 0.003$                   | NA               | NA             |
| NN                   | $0.867 \pm 0.002$                   | $0.875 \pm 0.001$                   | NA               | NA             |
| Dropout NN           | $0.856 \pm 0.001$                   | $0.873 \pm 0.001$                   | NA               | NA             |
| DNN                  | $0.872 \pm 0.001$                   | $0.876 \pm 0.001$                   | NA               | NA             |
| Dropout DNN          | <b><math>0.876 \pm 0.001</math></b> | <b><math>0.879 \pm 0.001</math></b> | NA               | NA             |
| VISH                 | $0.859 \pm 0.001$                   | NA                                  | NA               | NA             |
| CPD-SIP ( $R = 5$ )  | $0.860 \pm 0.001$                   | $0.871 \pm 0.001$                   | $20.14 \pm 0.04$ | $1416 \pm 12$  |
| CPD-SIP ( $R = 10$ ) | $0.869 \pm 0.001$                   | $0.872 \pm 0.002$                   | $19.97 \pm 0.10$ | $3567 \pm 1$   |
| CPD-SIP ( $R = 15$ ) | $0.871 \pm 0.001$                   | $0.873 \pm 0.001$                   | $19.86 \pm 0.10$ | $5628 \pm 153$ |
| CPD-SIP ( $R = 20$ ) | $0.872 \pm 0.001$                   | $0.875 \pm 0.001$                   | $19.74 \pm 0.04$ | $9437 \pm 721$ |

# References

- Batselier K., Low-rank tensor decompositions for nonlinear system identification, in IEEE Control Systems Magazine, vol. 42, no. 1, pp. 54-74, Feb. 2022
- Wesel F., Batselier K., Large-Scale Learning with Fourier Features and Tensor Decompositions, in Proc. of the 2021 Conference on Neural Information Processing Systems (NeurIPS), December 2021
- Karagoz R., Batselier K., Nonlinear system identification with regularized Tensor Network B-splines, Automatica, vol. 122, 2020
- Wesel F., Batselier K., Tensor-based Kernel Machines with Structured Inducing Points for Large and High-Dimensional Data, in Proc. of The 26th International Conference on Artificial Intelligence and Statistics (AISTATS), April 2023
- Wesel F., Batselier K., Quantized Fourier and Polynomial Features for more Expressive Tensor Network Models, in Proc. of The 27th International Conference on Artificial Intelligence and Statistics (AISTATS), April 2024.
- Schuurmans J., Batselier K., Kooij J., How informative is the Approximation Error from Tensor Decomposition for Neural Network Compression?, in Proc. of The Eleventh International Conference on Learning Representations (ICLR), May 2023