

EE613
Machine Learning for Engineers

LINEAR REGRESSION II

Sylvain Calinon
Robot Learning & Interaction Group
Idiap Research Institute
Nov. 14, 2019

EE613 - List of courses

19.09.2019 (JMO) Introduction

26.09.2019 (JMO) Generative I

03.10.2019 (JMO) Generative II

10.10.2019 (JMO) Generative III

17.10.2019 (JMO) Generative IV

24.10.2019 (JMO) Decision-trees

31.10.2019 (SC) Linear regression I

07.11.2019 (JMO) Kernel SVM

14.11.2019 (SC) Linear regression II

21.11.2019 (FF) MLP

28.11.2019 (FF) Feature-selection and boosting

05.12.2019 (SC) HMM and subspace clustering

12.12.2019 (SC) Nonlinear regression I

19.12.2019 (SC) Nonlinear regression II

Outline

Linear Regression II (Nov 14)

- Logistic regression
- Tensor-variate regression

HMM: preliminaries (Nov 14)

- Expectation-maximization (EM)
- Covariance structures in HMM

Logistic regression

Python notebook:

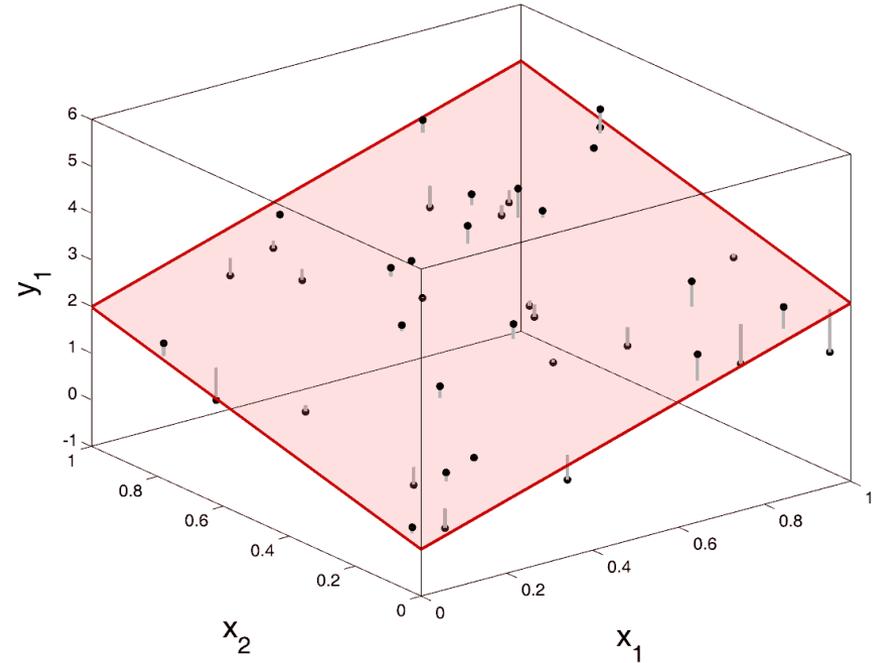
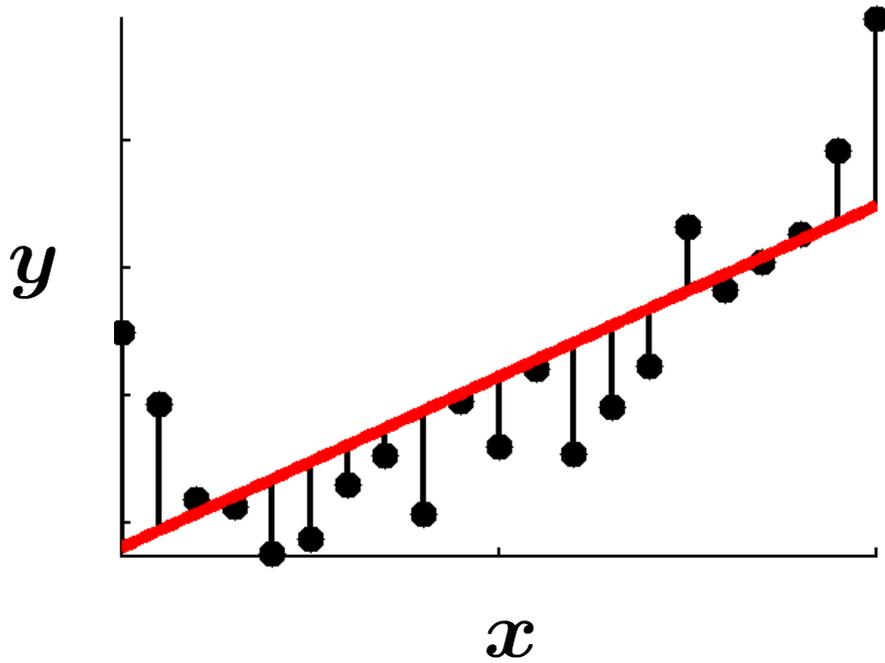
demo_LS_IRLS_logRegr.ipynb

Matlab code:

demo_LS_IRLS_logRegr01.m

Last course: Linear regression

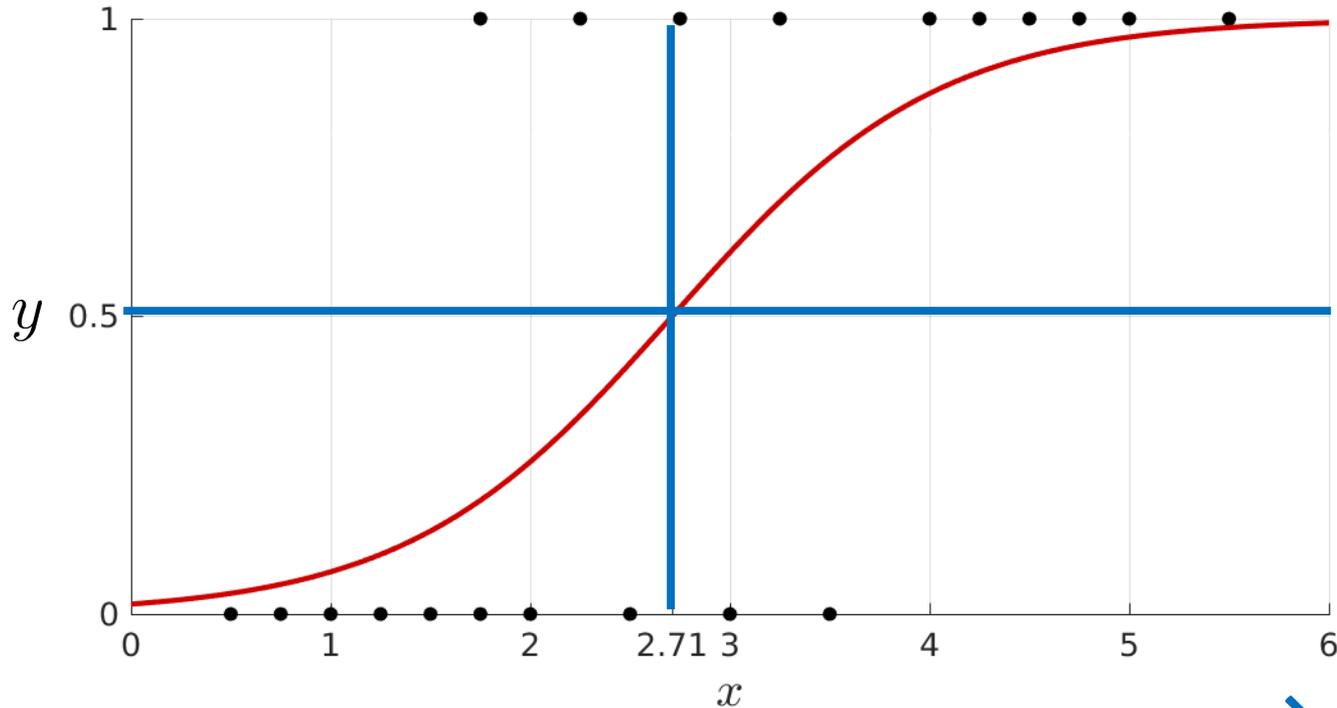
$$\hat{a} = X^\dagger y$$



→ Fitting a line/plane model

Logistic regression

Pass/fail in function of the time spent to study at an exam:

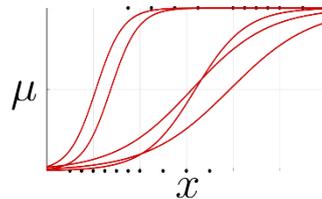


→ Classification

Logistic function:

$$\mu(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{a}^T \mathbf{x}}}$$

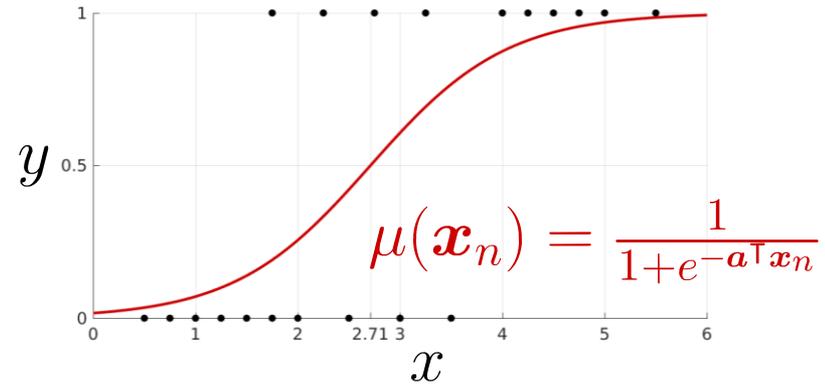
$$\mu(x) = \frac{1}{1 + e^{-(a_1 + a_2 x)}}$$



Logistic regression

Likelihood:

$$\mathcal{L} = \prod_n \mu(\mathbf{x}_n)^{y_n} (1 - \mu(\mathbf{x}_n))^{(1-y_n)}$$



Cost function as negative log-likelihood:

$$c = - \sum_n y_n \log(\mu(\mathbf{x}_n)) + (1 - y_n) \log(1 - \mu(\mathbf{x}_n))$$

$$\frac{\partial c}{\partial \mathbf{a}} = - \sum_n y_n \mu^{-1} \mu (1 - \mu) \mathbf{x}_n - (1 - y_n) (1 - \mu)^{-1} \mu (1 - \mu) \mathbf{x}_n$$

$$= - \sum_n y_n (1 - \mu) \mathbf{x}_n - (1 - y_n) \mu \mathbf{x}_n$$

$$= \sum_n (\mu - y_n) \mathbf{x}_n$$

$$\mu(t) = \frac{1}{1+e^{-t}}$$
$$\frac{\partial \mu}{\partial t} = \mu(1 - \mu)$$

Logistic regression

$$\frac{\partial c}{\partial \mathbf{a}} = \sum_n (\mu - y_n) \mathbf{x}_n$$

It can for example be solved by a Newton-Raphson iterative optimization scheme $\mathbf{a} \leftarrow \mathbf{a} - \mathbf{H}^{-1} \mathbf{g}$,

with gradient $\mathbf{g} = \sum_n (\mu(\mathbf{x}_n) - y_n) \mathbf{x}_n = \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y})$ and Hessian $\mathbf{H} = \mathbf{X}^\top \mathbf{W} \mathbf{X}$, with diagonal matrix $\mathbf{W} = \text{diag}(\boldsymbol{\mu} * (\mathbf{1} - \boldsymbol{\mu}))$.

We then obtain

$$\mathbf{a} \leftarrow \mathbf{a} - \mathbf{H}^{-1} \mathbf{g}$$

$$\leftarrow \mathbf{A} - (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y})$$

$$\leftarrow (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{W} \mathbf{X} \mathbf{a} - \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y}))$$

$$\leftarrow (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{W} \mathbf{X} \mathbf{a} + \mathbf{y} - \boldsymbol{\mu})$$

$$\leftarrow (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{z},$$

Hadamard
(elementwise)
product

with *working response* $\mathbf{z} = \mathbf{X} \mathbf{a} + \mathbf{W}^{-1} (\mathbf{y} - \boldsymbol{\mu})$.

→ IRLS procedure

Tensor-variate regression

Python notebook:
`demo_tensorRegr.ipynb`

Matlab code:
`demo_tensorRegr01.m`

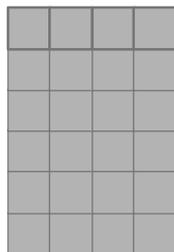
Tensors



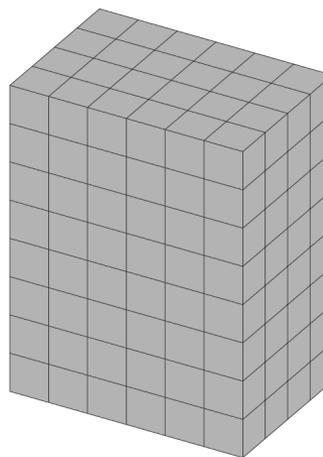
TensorFlow



1st-order
tensors



2nd-order
tensors



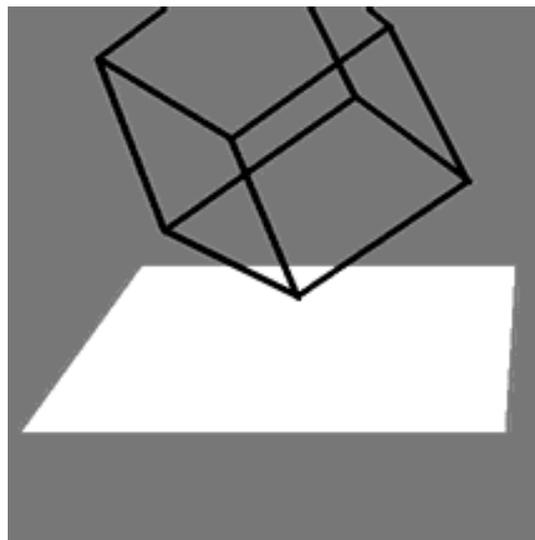
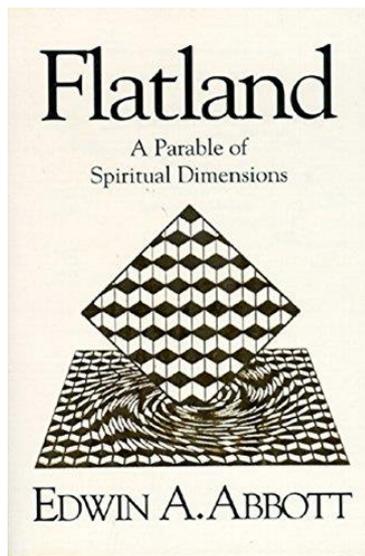
3rd-order
tensors

...

Examples of data organized as tensors:

- Recommender systems (e.g., age, M/F, city, income)
- Images (e.g., x, y, rgb channels)
- Videos (e.g., x, y, rgb channels, time)
- Robot motions (left/right arm, xyz, time)

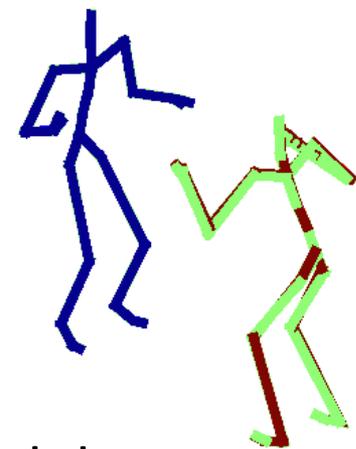
Tensor methods - Motivation



agent joint coordinate
sample time step

$$\mathcal{X} \in \mathbb{R}^{10 \times 2 \times 31 \times 3 \times 100}$$

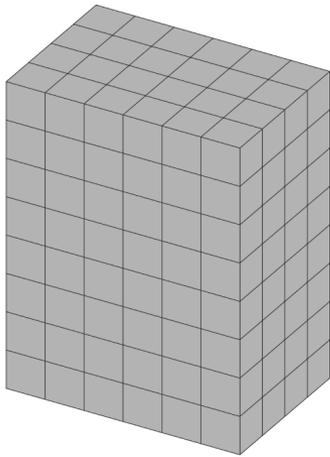
$$\mathbf{X} \in \mathbb{R}^{10 \times 18600}$$



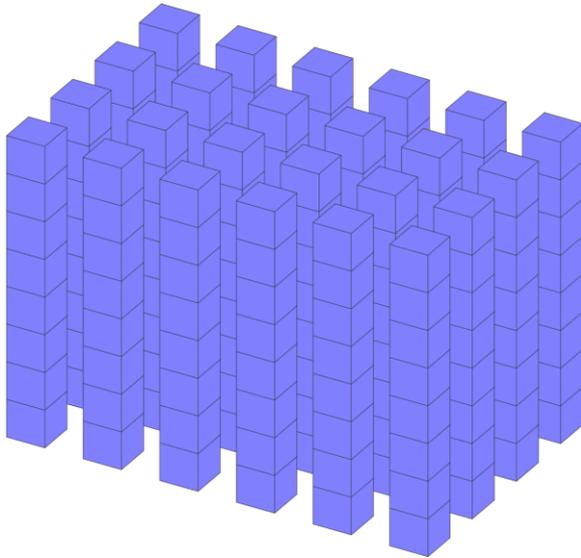
Tensor factorization keeps the structure of the original data
→ Multiway analysis of the data

Tensor indexing - Fibers

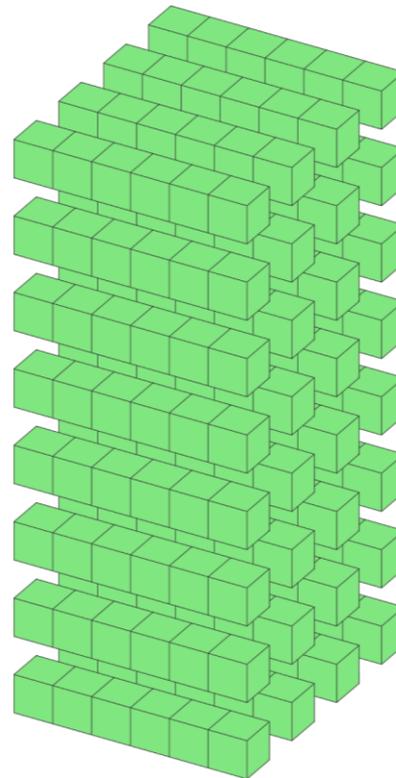
\mathcal{X} tensor
 X matrix
 \mathbf{x} vector
 x scalar



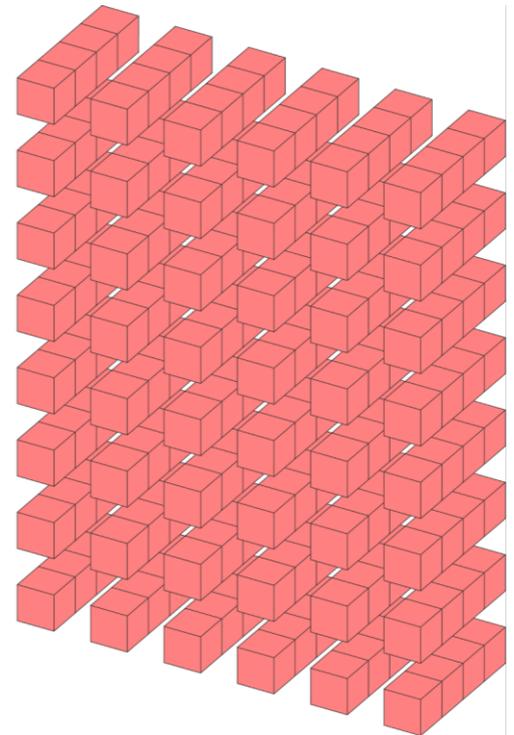
$$\mathcal{X} \in \mathbb{R}^{8 \times 6 \times 4}$$



$\mathbf{x}_{:,j,k}$ (column)

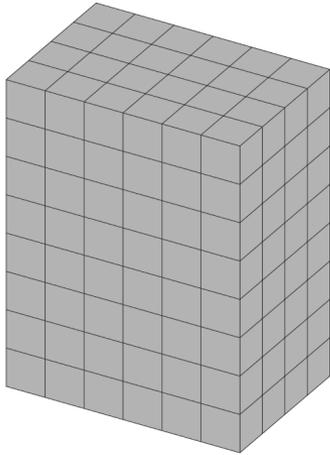


$\mathbf{x}_{i,:,k}$ (row)



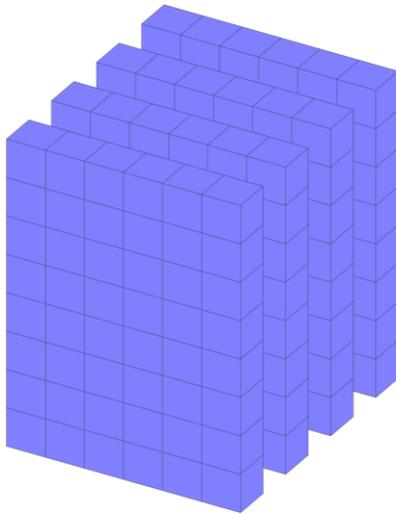
$\mathbf{x}_{i,j,:}$ (tube)

Tensor indexing - Slices

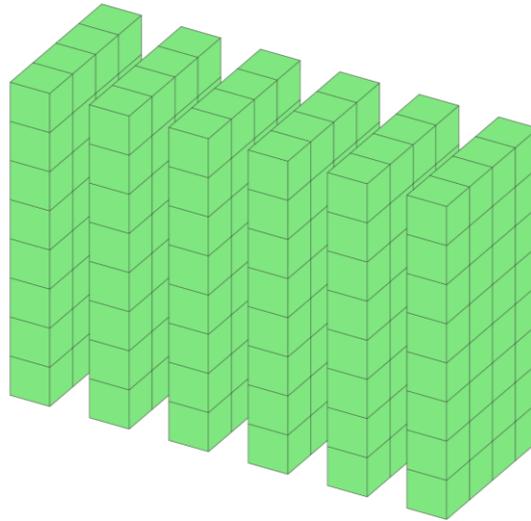


$$\mathcal{X} \in \mathbb{R}^{8 \times 6 \times 4}$$

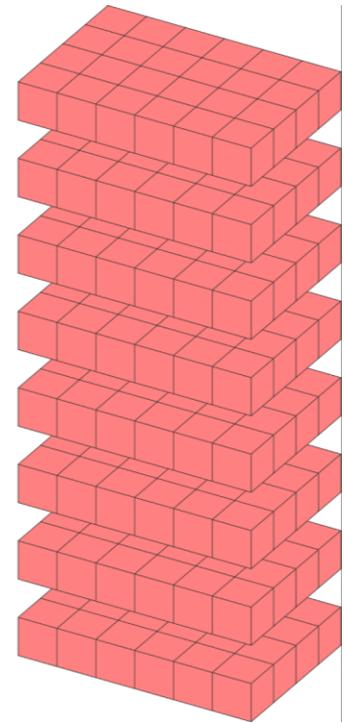
- \mathcal{X} tensor
- \mathbf{X} matrix
- \mathbf{x} vector
- x scalar



$$\mathbf{X}_{:,:,k} \text{ (frontal)}$$



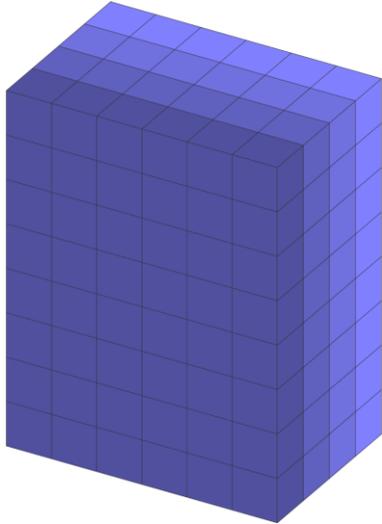
$$\mathbf{X}_{:,j,:} \text{ (lateral)}$$



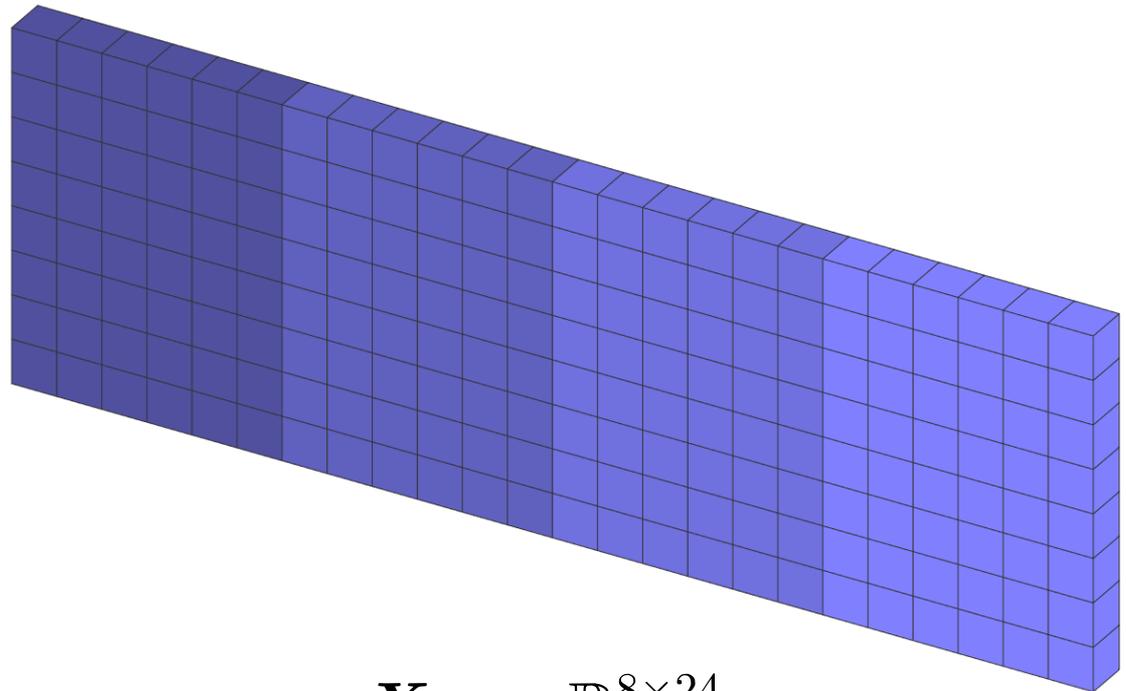
$$\mathbf{X}_{i,:,:} \text{ (horizontal)}$$

Tensor matricization / unfolding

A matrix $\mathbf{X}_{(n)} \in \mathbb{R}^{I_n \times (I_1 \cdots I_{n-1} I_{n+1} \cdots I_N)}$ results from the mode- n matricization (unfolding) of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$, which consists of turning the mode- n fibers of \mathcal{X} into the columns of a matrix $\mathbf{X}_{(n)}$.



$$\mathcal{X} \in \mathbb{R}^{8 \times 6 \times 4}$$



$$\mathbf{X}_{(1)} \in \mathbb{R}^{8 \times 24}$$

(mode-1 unfolding)

Products (Hadamard, Kronecker, Khatri-Rao)

Hadamard $A * B =$
(elementwise)

$$\begin{bmatrix} a_{1,1}b_{1,1} & a_{1,2}b_{1,2} & \cdots & a_{1,J}b_{1,J} \\ a_{2,1}b_{2,1} & a_{2,2}b_{2,2} & \cdots & a_{2,J}b_{2,J} \\ \vdots & \vdots & \ddots & \vdots \\ a_{I,1}b_{I,1} & a_{I,2}b_{I,2} & \cdots & a_{I,J}b_{I,J} \end{bmatrix}$$

$$\begin{aligned} A &\in \mathbb{R}^{I \times J} \\ B &\in \mathbb{R}^{I \times J} \\ A * B &\in \mathbb{R}^{I \times J} \end{aligned}$$

Kronecker $A \otimes B =$

$$\begin{bmatrix} a_{1,1}B & a_{1,2}B & \cdots & a_{1,J}B \\ a_{2,1}B & a_{2,2}B & \cdots & a_{2,J}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{I,1}B & a_{I,2}B & \cdots & a_{I,J}B \end{bmatrix}$$

$$\begin{aligned} A &\in \mathbb{R}^{I \times J} \\ B &\in \mathbb{R}^{K \times L} \\ A \otimes B &\in \mathbb{R}^{IK \times JL} \end{aligned}$$

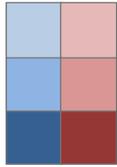
Khatri-Rao $A \odot B =$

$$\begin{bmatrix} a_{1,1}\mathbf{b}_1 & a_{1,2}\mathbf{b}_2 & \cdots & a_{1,K}\mathbf{b}_K \\ a_{2,1}\mathbf{b}_1 & a_{2,2}\mathbf{b}_2 & \cdots & a_{2,K}\mathbf{b}_K \\ \vdots & \vdots & \ddots & \vdots \\ a_{I,1}\mathbf{b}_1 & a_{I,2}\mathbf{b}_2 & \cdots & a_{I,K}\mathbf{b}_K \end{bmatrix}$$

$$\begin{aligned} A &\in \mathbb{R}^{I \times K} \\ B &\in \mathbb{R}^{J \times K} \\ A \odot B &\in \mathbb{R}^{IJ \times K} \end{aligned}$$

Hadamard (elementwise) product

Example



A

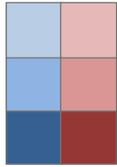


B

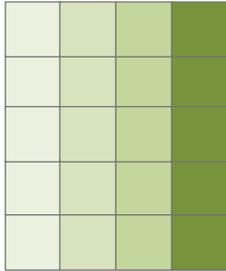
$$A * B = \begin{bmatrix} \text{light blue} & \text{light green} & \text{light red} & \text{light purple} \\ \text{light blue} & \text{light green} & \text{light red} & \text{light purple} \\ \text{dark blue} & \text{dark green} & \text{dark red} & \text{dark purple} \end{bmatrix}$$

$$A \in \mathbb{R}^{3 \times 2}$$
$$B \in \mathbb{R}^{3 \times 2}$$
$$A * B \in \mathbb{R}^{3 \times 2}$$

Kronecker product - Example

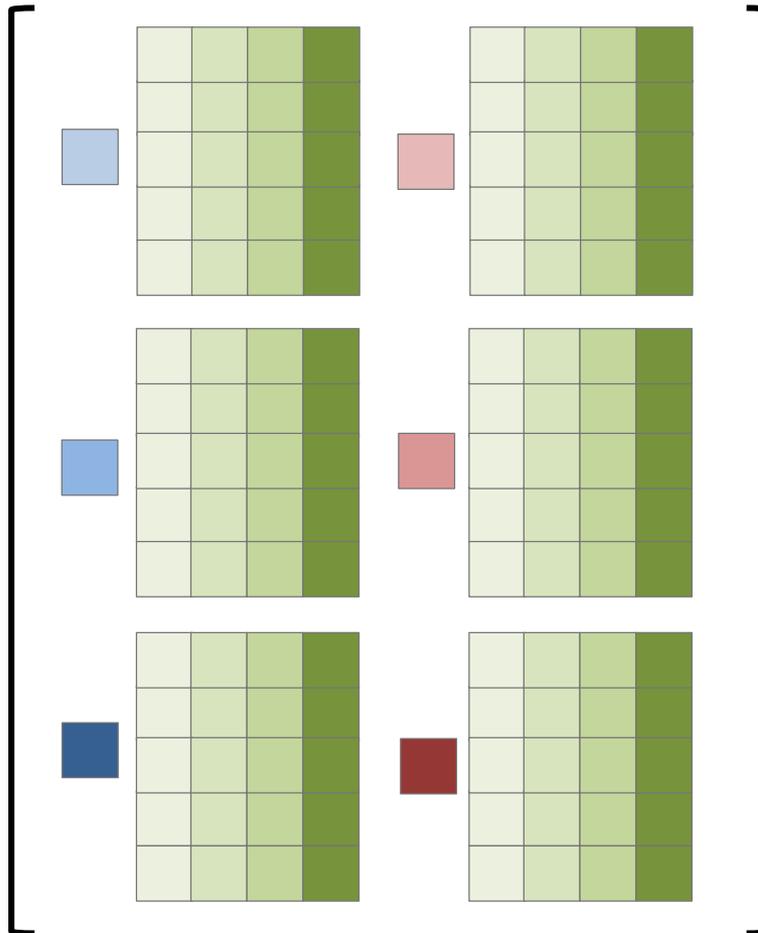


A



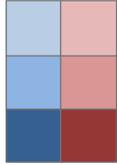
B

$A \otimes B =$

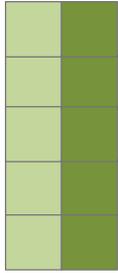


$$\begin{aligned} A &\in \mathbb{R}^{3 \times 2} \\ B &\in \mathbb{R}^{5 \times 4} \\ A \otimes B &\in \mathbb{R}^{15 \times 8} \end{aligned}$$

Khatri-Rao product - Example

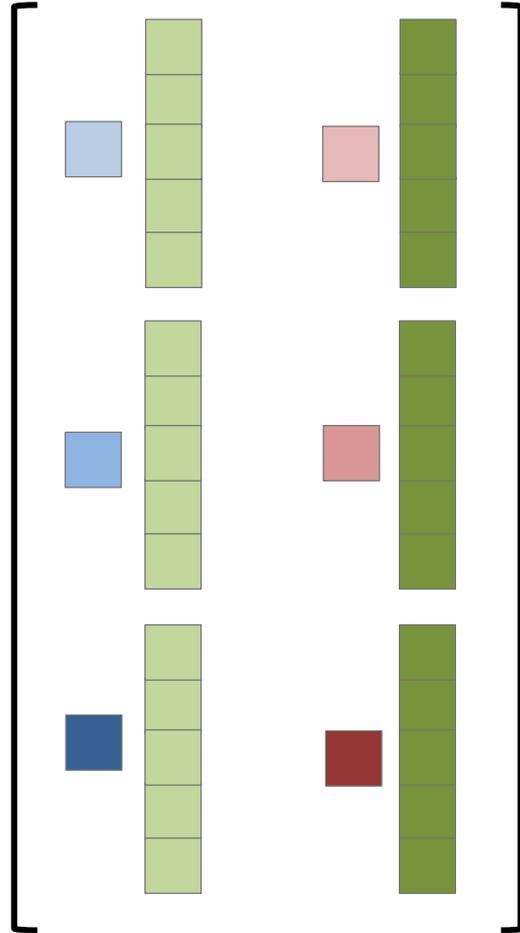


A



B

$A \odot B =$

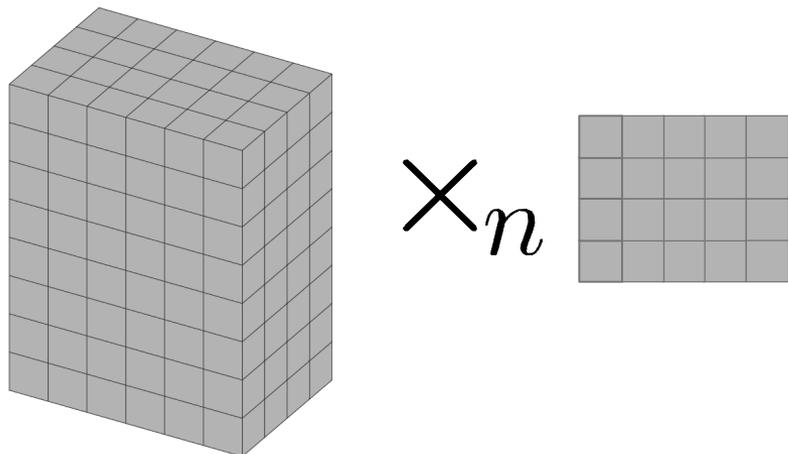


$$A \in \mathbb{R}^{3 \times 2}$$

$$B \in \mathbb{R}^{5 \times 2}$$

$$A \odot B \in \mathbb{R}^{15 \times 2}$$

Mode- n product



$$\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$$

$$M \in \mathbb{R}^{J \times I_n}$$

$$\mathcal{Y} \in \mathbb{R}^{I_1 \times \cdots \times I_{n-1} \times J \times I_{n+1} \times \cdots \times I_N}$$

$$\mathcal{Y} = \mathcal{X} \times_n M$$

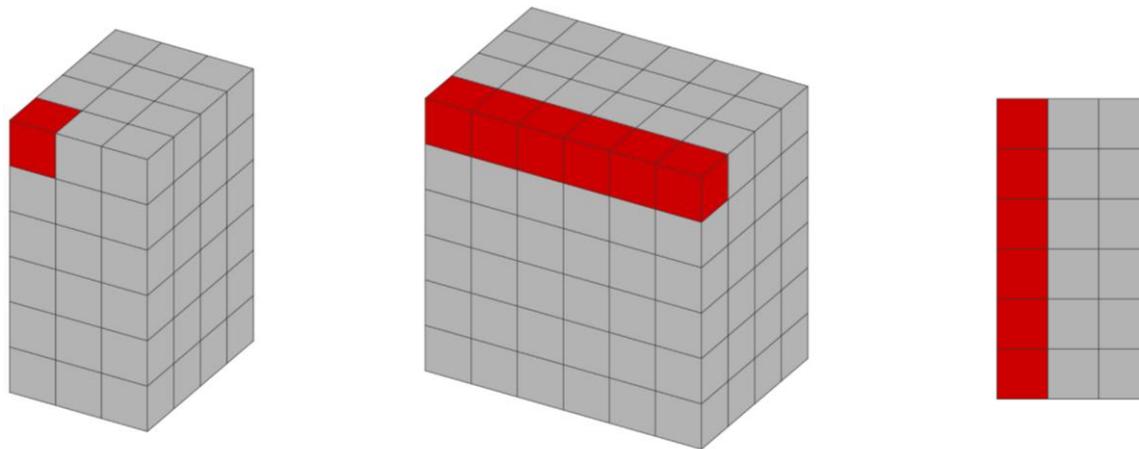
$$Y_{(n)} = M X_{(n)} \quad (\text{matricized form})$$

$$y_{i_1, \dots, i_{n-1}, j, i_{n+1}, \dots, i_N} = \sum_{i_n=1}^{I_n} x_{i_1, \dots, i_N} m_{j, i_n} \quad (\text{elementwise})$$

Intuitively, the operation corresponds to multiplying each mode- n fiber of \mathcal{X} by the matrix M .

Mode-n product: Example

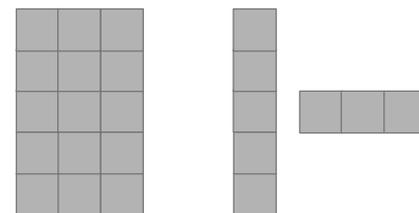
$$\begin{aligned}\mathcal{X} &\in \mathbb{R}^{8 \times 6 \times 4} \\ M &\in \mathbb{R}^{6 \times 3} \\ \mathcal{Y} &\in \mathbb{R}^{8 \times 3 \times 4}\end{aligned}$$



$$\mathcal{Y} = \mathcal{X} \times_2 M$$

Outer product and inner product

The **outer product** of two vectors $\mathbf{a} \in \mathbb{R}^I$ and $\mathbf{b} \in \mathbb{R}^J$ results in a matrix $\mathbf{X} \in \mathbb{R}^{I \times J}$ denoted by $\mathbf{X} = \mathbf{a} \circ \mathbf{b} = \mathbf{a}\mathbf{b}^\top$.



$$\begin{aligned} \mathbf{X} &= \mathbf{a} \mathbf{b}^\top \\ &= \mathbf{a} \circ \mathbf{b} \end{aligned}$$

(outer product)

The **outer product** of three (or more) vectors $\mathbf{a} \in \mathbb{R}^I$, $\mathbf{b} \in \mathbb{R}^J$ and $\mathbf{c} \in \mathbb{R}^K$ results in a tensor $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ denoted by $\mathcal{X} = \mathbf{a} \circ \mathbf{b} \circ \mathbf{c}$ with elements $x_{i,j,k} = a_i b_j c_k$.

The **inner product** of two vectors $\mathbf{a} \in \mathbb{R}^I$ and $\mathbf{b} \in \mathbb{R}^I$ results in a scalar $x = \langle \mathbf{a}, \mathbf{b} \rangle = \mathbf{a}^\top \mathbf{b} = \sum_{i=1}^I a_i b_i$.



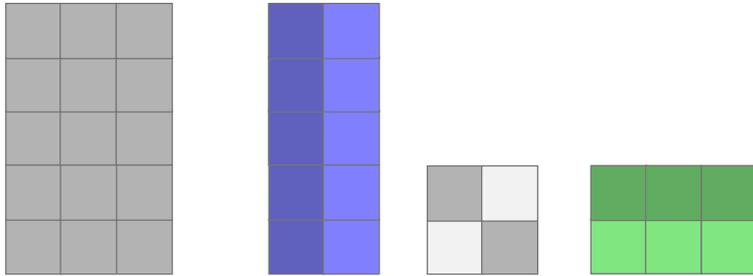
$$\begin{aligned} x &= \mathbf{a}^\top \mathbf{b} \\ &= \langle \mathbf{a}, \mathbf{b} \rangle \end{aligned}$$

(inner product)

The formulation can be extended to tensors \mathcal{A} and \mathcal{B} of the same size. We have

$$\langle \mathcal{A}, \mathcal{B} \rangle = \langle \mathbf{A}_{(n)}, \mathbf{B}_{(n)} \rangle = \langle \text{vec}(\mathcal{A}), \text{vec}(\mathcal{B}) \rangle.$$

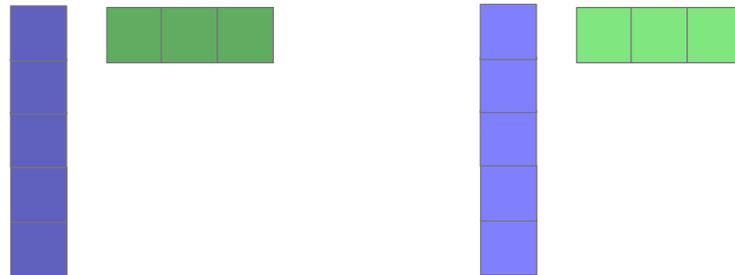
Singular value decomposition (SVD)



$$X = U \Sigma V^T$$

$$= \sigma_1^2 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2^2 \mathbf{u}_2 \mathbf{v}_2^T$$

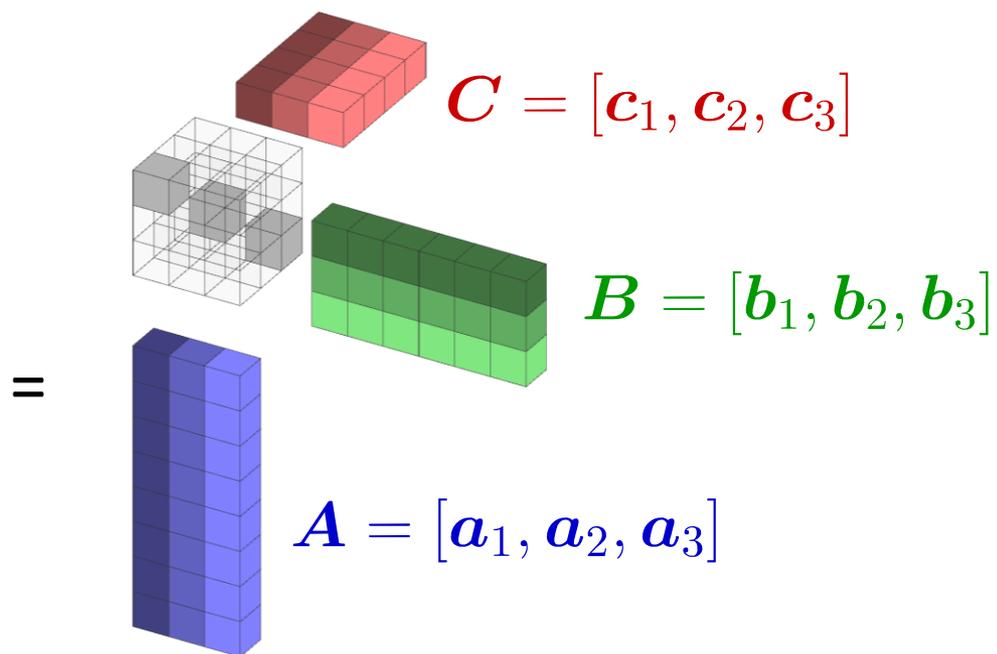
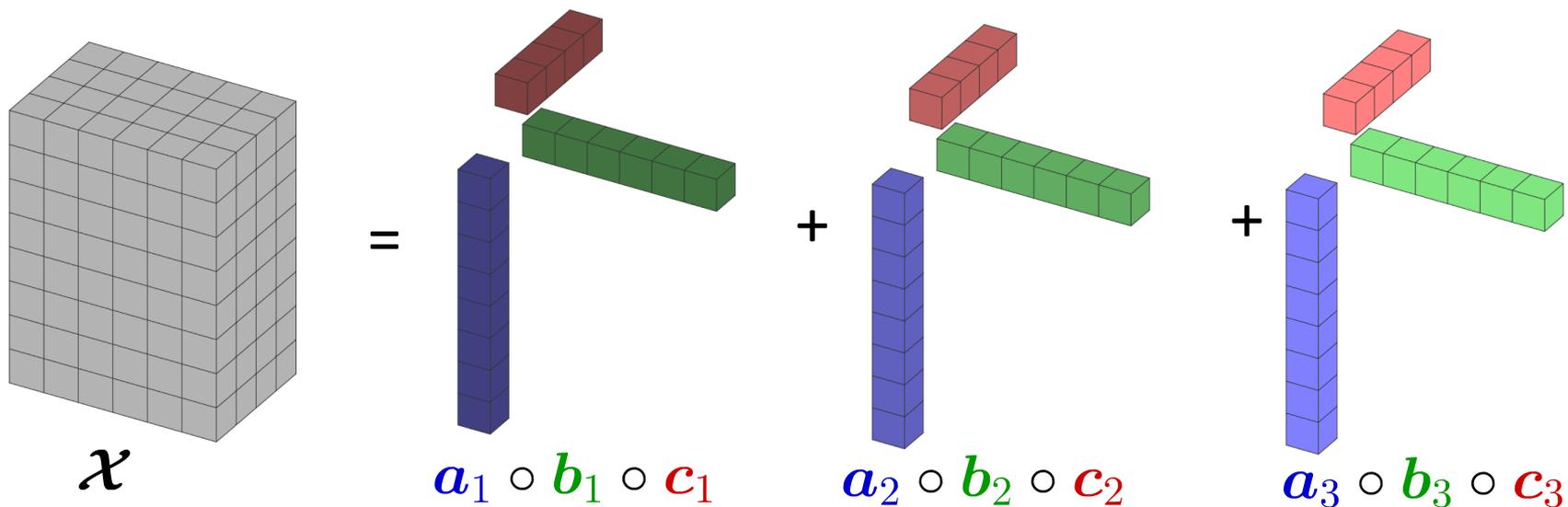
$$= \tilde{\mathbf{u}}_1 \tilde{\mathbf{v}}_1^T + \tilde{\mathbf{u}}_2 \tilde{\mathbf{v}}_2^T$$



$$\tilde{\mathbf{u}}_i = \sigma_i \mathbf{u}_i$$

$$\tilde{\mathbf{v}}_i = \sigma_i \mathbf{v}_i$$

CP decomposition



CP decomposition

$$\begin{aligned}\mathcal{X} &= \sum_{r=1}^R \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r \\ &= [\mathbf{A}, \mathbf{B}, \mathbf{C}]\end{aligned}$$

$$\begin{aligned}\mathbf{A} &\in \mathbb{R}^{I \times R} \\ \mathbf{B} &\in \mathbb{R}^{J \times R} \\ \mathbf{C} &\in \mathbb{R}^{K \times R}\end{aligned}$$

Matricized form:

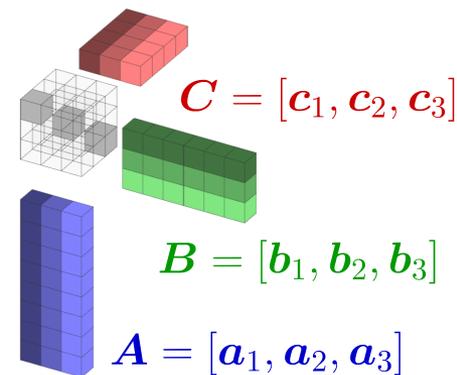
$$\begin{aligned}\mathbf{X}_{(1)} &= \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top \\ \mathbf{X}_{(2)} &= \mathbf{B}(\mathbf{C} \odot \mathbf{A})^\top \\ \mathbf{X}_{(3)} &= \mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top\end{aligned}$$

Vectorized form: $\text{vec}(\mathcal{X}) = (\mathbf{C} \odot \mathbf{B} \odot \mathbf{A}) \mathbf{1}_R$

Elementwise: $x_{i,j,k} = \sum_{r=1}^R a_{i,r} b_{j,r} c_{k,r}$

$\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_R]$ is called a factor matrix.

The *tensor rank* R corresponds to the smallest number of components required in the CP decomposition.



Parameters estimation: Alternating least squares

The CP decomposition can be solved by alternating least squares (ALS), by repeating

$$\mathbf{A} \leftarrow \arg \min_{\mathbf{A}} \left\| \mathbf{X}_{(1)} - \mathbf{A}(\mathbf{C} \odot \mathbf{B})^{\top} \right\|_{\text{F}}^2$$

$$\mathbf{B} \leftarrow \arg \min_{\mathbf{B}} \left\| \mathbf{X}_{(2)} - \mathbf{B}(\mathbf{C} \odot \mathbf{A})^{\top} \right\|_{\text{F}}^2$$

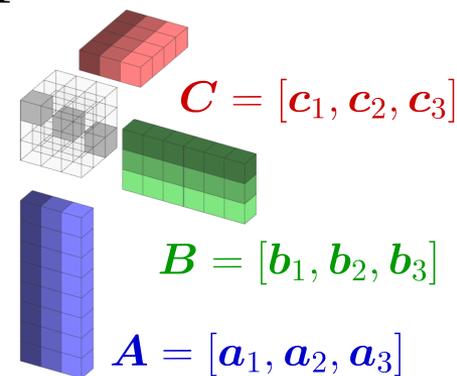
$$\mathbf{C} \leftarrow \arg \min_{\mathbf{C}} \left\| \mathbf{X}_{(3)} - \mathbf{C}(\mathbf{B} \odot \mathbf{A})^{\top} \right\|_{\text{F}}^2$$

until convergence, yielding the update rules

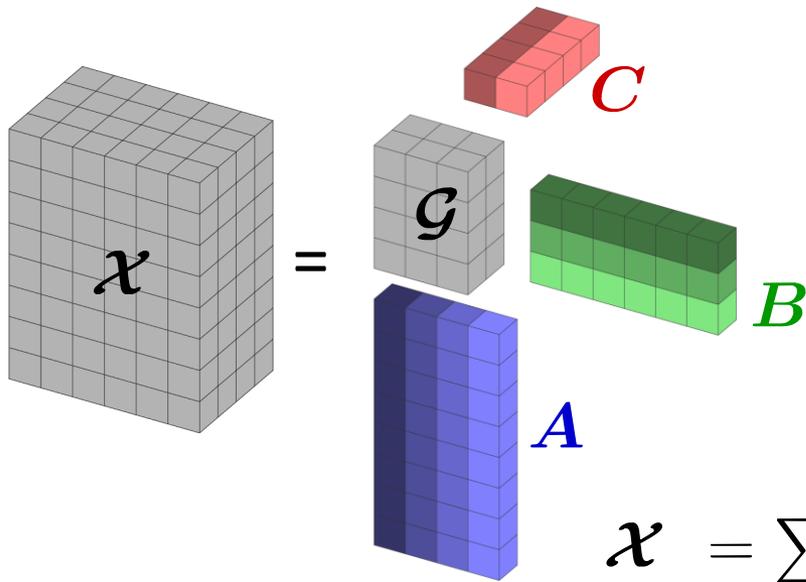
$$\mathbf{A} \leftarrow \mathbf{X}_{(1)} \left((\mathbf{C} \odot \mathbf{B})^{\top} \right)^{\dagger}$$

$$\mathbf{B} \leftarrow \mathbf{X}_{(2)} \left((\mathbf{C} \odot \mathbf{A})^{\top} \right)^{\dagger}$$

$$\mathbf{C} \leftarrow \mathbf{X}_{(3)} \left((\mathbf{B} \odot \mathbf{A})^{\top} \right)^{\dagger}$$



Tucker decomposition



Core tensor

$$\begin{aligned} \mathcal{G} &\in \mathbb{R}^{P \times Q \times R} \\ \mathbf{A} &\in \mathbb{R}^{I \times P} \\ \mathbf{B} &\in \mathbb{R}^{J \times Q} \\ \mathbf{C} &\in \mathbb{R}^{K \times R} \end{aligned}$$

$$\begin{aligned} \mathcal{X} &= \sum_{p=1}^P \sum_{q=1}^Q \sum_{r=1}^R g_{p,q,r} \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r \\ &= \mathcal{G} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C} \\ &= [\mathcal{G}; \mathbf{A}, \mathbf{B}, \mathbf{C}] \end{aligned}$$

Matricized form:

$$\begin{aligned} \mathbf{X}_{(1)} &= \mathbf{A} \mathbf{G}_{(1)} (\mathbf{C} \otimes \mathbf{B})^\top \\ \mathbf{X}_{(2)} &= \mathbf{B} \mathbf{G}_{(2)} (\mathbf{C} \otimes \mathbf{A})^\top \\ \mathbf{X}_{(3)} &= \mathbf{C} \mathbf{G}_{(3)} (\mathbf{B} \otimes \mathbf{A})^\top \end{aligned}$$

Elementwise:

$$x_{i,j,k} = \sum_{p=1}^P \sum_{q=1}^Q \sum_{r=1}^R g_{p,q,r} a_{i,p} b_{j,q} c_{k,r}$$

Parameters estimation:

Higher-order orthogonal iteration (HOOI)

$$\min_{\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C}} \left\| \mathcal{X} - \llbracket \mathcal{G}; \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket \right\|_{\text{F}}^2 \text{ s.t. } \mathbf{A}^\top \mathbf{A} = \mathbf{I}_P, \mathbf{B}^\top \mathbf{B} = \mathbf{I}_Q, \mathbf{C}^\top \mathbf{C} = \mathbf{I}_R$$

which can be solved by repeating

$$\mathcal{Y}^A \leftarrow \mathcal{X} \times_2 \mathbf{B}^\top \times_3 \mathbf{C}^\top$$

$$\mathcal{Y}^B \leftarrow \mathcal{X} \times_1 \mathbf{A}^\top \times_3 \mathbf{C}^\top$$

$$\mathcal{Y}^C \leftarrow \mathcal{X} \times_1 \mathbf{A}^\top \times_2 \mathbf{B}^\top$$

$$\mathbf{A} \leftarrow P \text{ leading singular vectors of } \mathbf{Y}_{(1)}^A$$

$$\mathbf{B} \leftarrow Q \text{ leading singular vectors of } \mathbf{Y}_{(2)}^B$$

$$\mathbf{C} \leftarrow R \text{ leading singular vectors of } \mathbf{Y}_{(3)}^C$$

until convergence, with \mathcal{G} finally evaluated as

$$\mathcal{G} \leftarrow \mathcal{X} \times_1 \mathbf{A}^\top \times_2 \mathbf{B}^\top \times_3 \mathbf{C}^\top$$

In contrast to CP, the Tucker decomposition is generally not unique

→ \mathbf{A} , \mathbf{B} and \mathbf{C} constrained to be orthogonal matrices

Parameters estimation:

Higher-order orthogonal iteration (HOOI)

The problem can be recast as a series of maximization subproblems

$$\begin{aligned} \mathbf{A} &\leftarrow \arg \max_{\mathbf{A}} \left\| \mathbf{A}^\top \mathbf{X}_{(1)}(\mathbf{C} \otimes \mathbf{B}) \right\|_{\text{F}}^2 & \text{s.t.} \quad \mathbf{A}^\top \mathbf{A} = \mathbf{I}_P \\ \mathbf{B} &\leftarrow \arg \max_{\mathbf{B}} \left\| \mathbf{B}^\top \mathbf{X}_{(2)}(\mathbf{C} \otimes \mathbf{A}) \right\|_{\text{F}}^2 & \text{s.t.} \quad \mathbf{B}^\top \mathbf{B} = \mathbf{I}_Q \\ \mathbf{C} &\leftarrow \arg \max_{\mathbf{C}} \left\| \mathbf{C}^\top \mathbf{X}_{(3)}(\mathbf{B} \otimes \mathbf{A}) \right\|_{\text{F}}^2 & \text{s.t.} \quad \mathbf{C}^\top \mathbf{C} = \mathbf{I}_R \end{aligned}$$

which can be solved by repeating

$$\begin{aligned} \mathbf{A} &\leftarrow P \text{ leading singular vectors of } \mathbf{X}_{(1)}(\mathbf{C} \otimes \mathbf{B}) \\ \mathbf{B} &\leftarrow Q \text{ leading singular vectors of } \mathbf{X}_{(2)}(\mathbf{C} \otimes \mathbf{A}) \\ \mathbf{C} &\leftarrow R \text{ leading singular vectors of } \mathbf{X}_{(3)}(\mathbf{B} \otimes \mathbf{A}) \end{aligned}$$

until convergence, with \mathcal{G} finally evaluated as

$$\mathcal{G} \leftarrow \mathcal{X} \times_1 \mathbf{A}^\top \times_2 \mathbf{B}^\top \times_3 \mathbf{C}^\top$$

Tensor-variate linear regression

y predicted output
 \mathbf{w} vector of weights
 b bias
 ϵ Gaussian noise

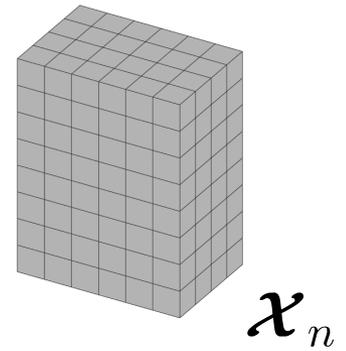
For vector-variate \mathbf{x} :
$$y = \mathbf{x}^\top \mathbf{w} + b + \epsilon$$
$$= \langle \mathbf{x}, \mathbf{w} \rangle + b + \epsilon$$

For matrix-variate \mathbf{X} :
$$y = \mathbf{w}^{(1)\top} \mathbf{X} \mathbf{w}^{(2)} + b + \epsilon$$
$$= \langle \mathbf{X}, \mathbf{w}^{(1)} \circ \mathbf{w}^{(2)} \rangle + b + \epsilon$$

For tensor-variate \mathcal{X} :
$$y = \langle \mathcal{X}, \mathbf{w}^{(1)} \circ \dots \circ \mathbf{w}^{(M)} \rangle + b + \epsilon$$
$$= \langle \mathcal{X}, \mathcal{W} \rangle + b + \epsilon$$

\Rightarrow for \mathcal{W} of rank R :
$$y = \langle \mathcal{X}, \sum_{r=1}^R \mathbf{w}_r^{(1)} \circ \dots \circ \mathbf{w}_r^{(M)} \rangle + b + \epsilon$$
$$= \langle \mathcal{X}, \mathcal{W} \rangle + b + \epsilon$$

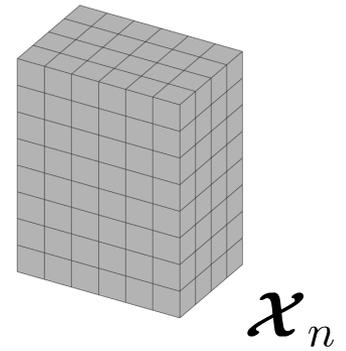
Tensor-variate linear regression: Parameters estimation



$$\begin{aligned} y_n &= \left\langle \mathcal{X}_n, \underbrace{\sum_{r=1}^R \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r}_{\phi_{1,n}} \right\rangle + b \\ &= \left\langle \mathbf{X}_{(1),n}, \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top \right\rangle \\ &= \left\langle \mathbf{X}_{(1),n}(\mathbf{C} \odot \mathbf{B}), \mathbf{A} \right\rangle \\ &= \left\langle \text{vec}(\mathbf{X}_{(1),n}(\mathbf{C} \odot \mathbf{B})), \text{vec}(\mathbf{A}) \right\rangle \\ &= \underbrace{\text{vec}(\mathbf{X}_{(1),n}(\mathbf{C} \odot \mathbf{B}))^\top}_{\phi_{1,n}} \text{vec}(\mathbf{A}) \end{aligned}$$

$$\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^\top \mathbf{v}$$

Tensor-variate linear regression: Parameters estimation



$$y_n = \left\langle \mathcal{X}_n, \sum_{r=1}^R \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r \right\rangle + b$$

$$= \underbrace{\text{vec}\left(\mathbf{X}_{(1),n}(\mathbf{C} \odot \mathbf{B})\right)}_{\phi_{1,n}}{}^\top \text{vec}(\mathbf{A})$$

$$= \underbrace{\text{vec}\left(\mathbf{X}_{(2),n}(\mathbf{C} \odot \mathbf{A})\right)}_{\phi_{2,n}}{}^\top \text{vec}(\mathbf{B})$$

$$= \underbrace{\text{vec}\left(\mathbf{X}_{(3),n}(\mathbf{B} \odot \mathbf{A})\right)}_{\phi_{3,n}}{}^\top \text{vec}(\mathbf{C})$$

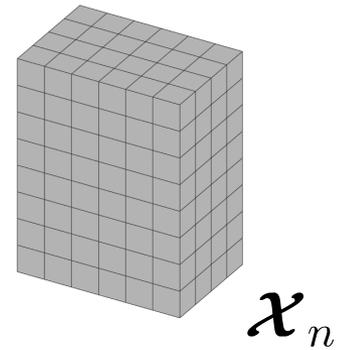
$$\mathbf{y} - \mathbf{1}b = \Phi_1 \text{vec}(\mathbf{A})$$

$$\mathbf{y} - \mathbf{1}b = \Phi_2 \text{vec}(\mathbf{B})$$

$$\mathbf{y} - \mathbf{1}b = \Phi_3 \text{vec}(\mathbf{C})$$

$$\left(\begin{array}{l} \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \\ \Phi_i = \begin{bmatrix} \Phi_{i,1} \\ \Phi_{i,2} \\ \vdots \\ \Phi_{i,N} \end{bmatrix} \end{array} \right)$$

Tensor-variate linear regression: Parameters estimation



Alternating least squares (ALS)

update rules:

$$\text{vec}(\mathbf{A}) \leftarrow \Phi_1^\dagger (\mathbf{y} - \mathbf{1}b)$$

$$\text{vec}(\mathbf{B}) \leftarrow \Phi_2^\dagger (\mathbf{y} - \mathbf{1}b)$$

$$\text{vec}(\mathbf{C}) \leftarrow \Phi_3^\dagger (\mathbf{y} - \mathbf{1}b)$$

$$b \leftarrow \frac{1}{N} \sum_{n=1}^N \left(y_n - \langle \mathbf{X}_{(1),n}, \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top \rangle \right)$$

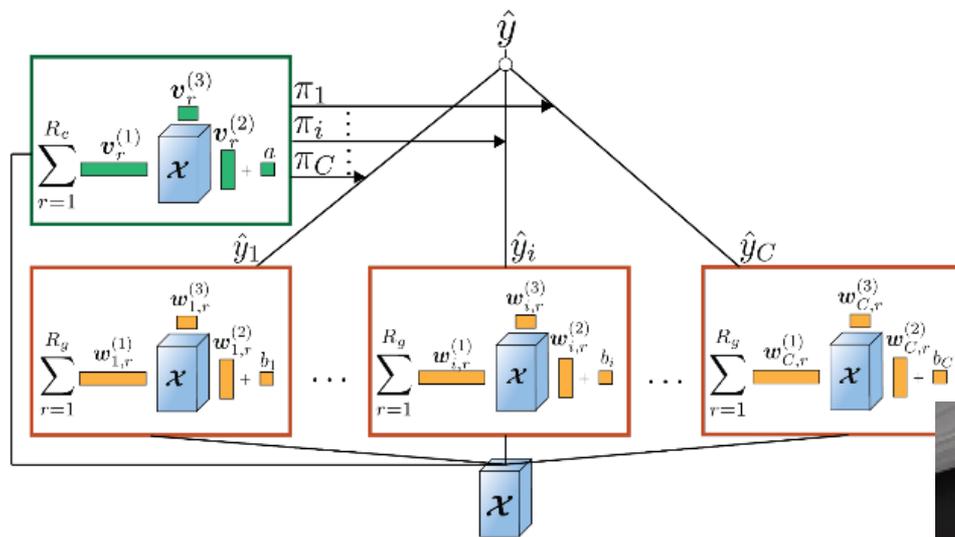
$$\mathbf{y} - \mathbf{1}b = \Phi_1 \text{vec}(\mathbf{A})$$

$$\mathbf{y} - \mathbf{1}b = \Phi_2 \text{vec}(\mathbf{B})$$

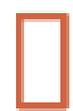
$$\mathbf{y} - \mathbf{1}b = \Phi_3 \text{vec}(\mathbf{C})$$

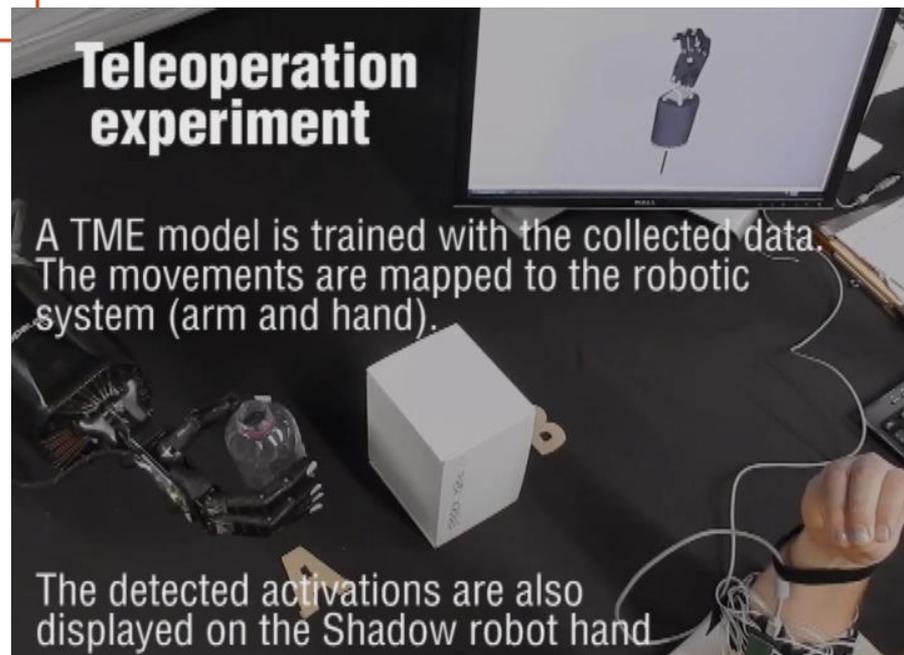
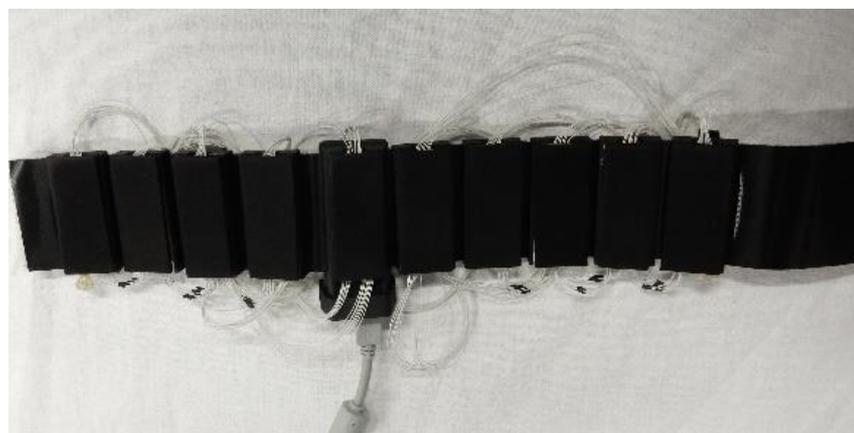
Example of application:

Tensor-variate mixture of experts



 logistic regression (gates)

 ridge regression (experts)



References

Logistic regression

Walker, SH, Duncan, DB (1967) Estimation of the probability of an event as a function of several independent variables. *Biometrika* 54 (1/2): 167–178.

Tensor-variate regression

Kolda T, Bader B (2009) Tensor decompositions and applications. *SIAM Review* 51(3):455-500

Comon P (2014) Tensors: A brief introduction. *IEEE Signal Processing Magazine* 31(3):44-53

Rabanser S, Shchur O, Günnemann S (2017) Introduction to tensor decompositions and their applications in machine learning. arXiv:171110781 pp 1-13

Sorber L, Van Barel M, De Lathauwer L (2015) Structured data fusion. *IEEE Journal of Selected Topics in Signal Processing* 9(4):586-600

Tensor methods - Softwares

<http://tensorly.org> (Python)

<https://www.tensorlab.net> (Matlab)