

---

# Statistical Learning and Data Science

**Xinwei Deng**

xdeng@vt.edu

Department of Statistics

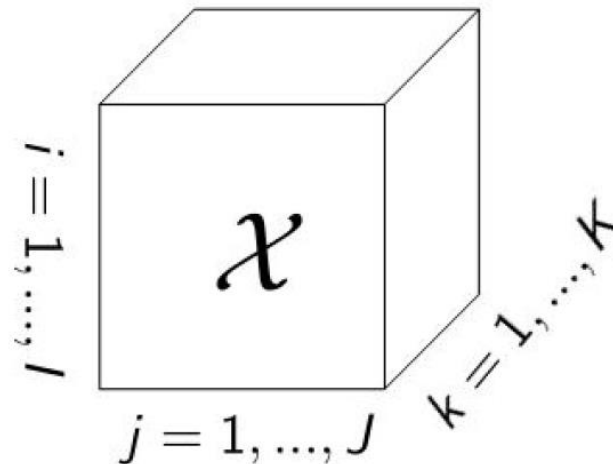
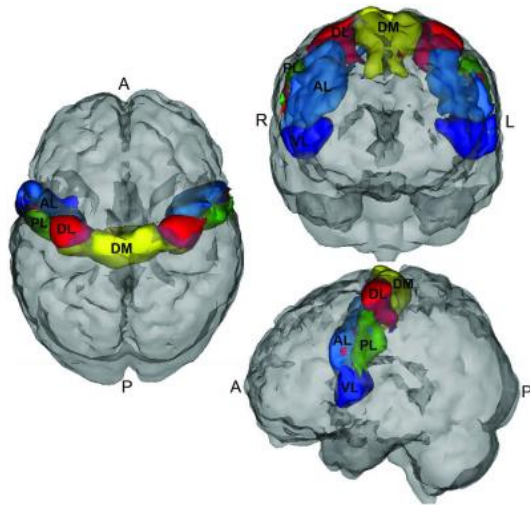
# Tensor Regression and Beyond

---

- Introduction to Tensor Data.
- Tensor and Its Decomposition.
- Tensor Regression
- Modeling and Analysis of Tensor Data
- Future Direction

# Tensor: A Generalization of Matrix

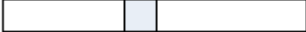
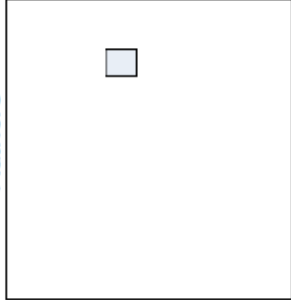
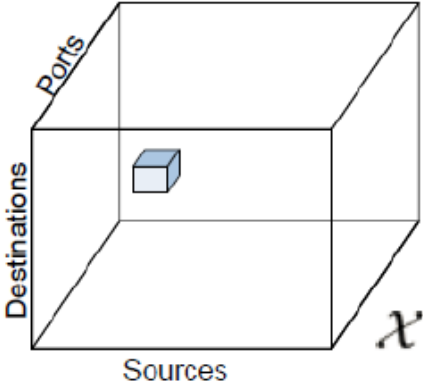
- Tensor: a *multi-way* array.
- Example: fMRI data in Neuroscience: 3D brain images.



# What is Tensor?

A tensor is formally denoted as  $\mathcal{X} \in \mathbb{R}^{l_1 \times l_2 \times \dots \times l_N}$

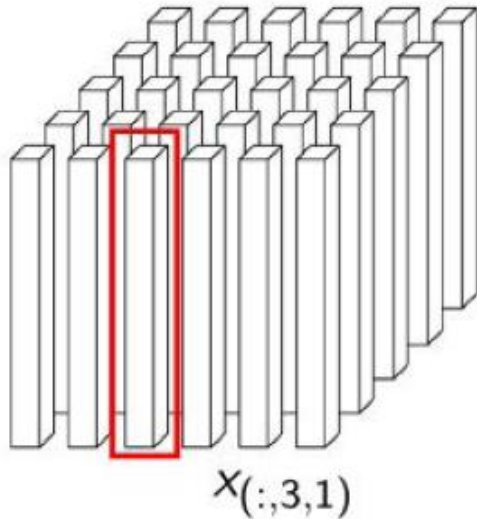
- generalization of vector and matrix
- represented as multi-dimensional array

Order	1st	2 <sup>nd</sup>	3 <sup>rd</sup>
Correspondence	Vector	Matrix	3D array
Example	 <p>Sensors</p>	 <p>Keywords</p> <p>Authors</p>	 <p>Ports</p> <p>Destinations</p> <p>Sources</p> <p><math>\mathcal{X}</math></p>

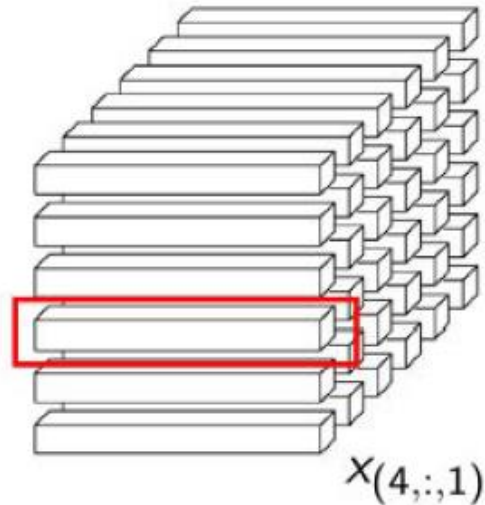
# Definition of Fiber

- Fibers** are created when fixing all but one index:

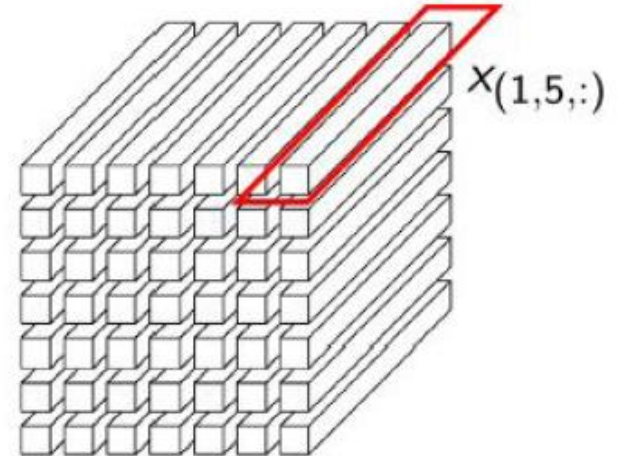
Column(Mode 1)Fibers



Column(Mode 2)Fibers

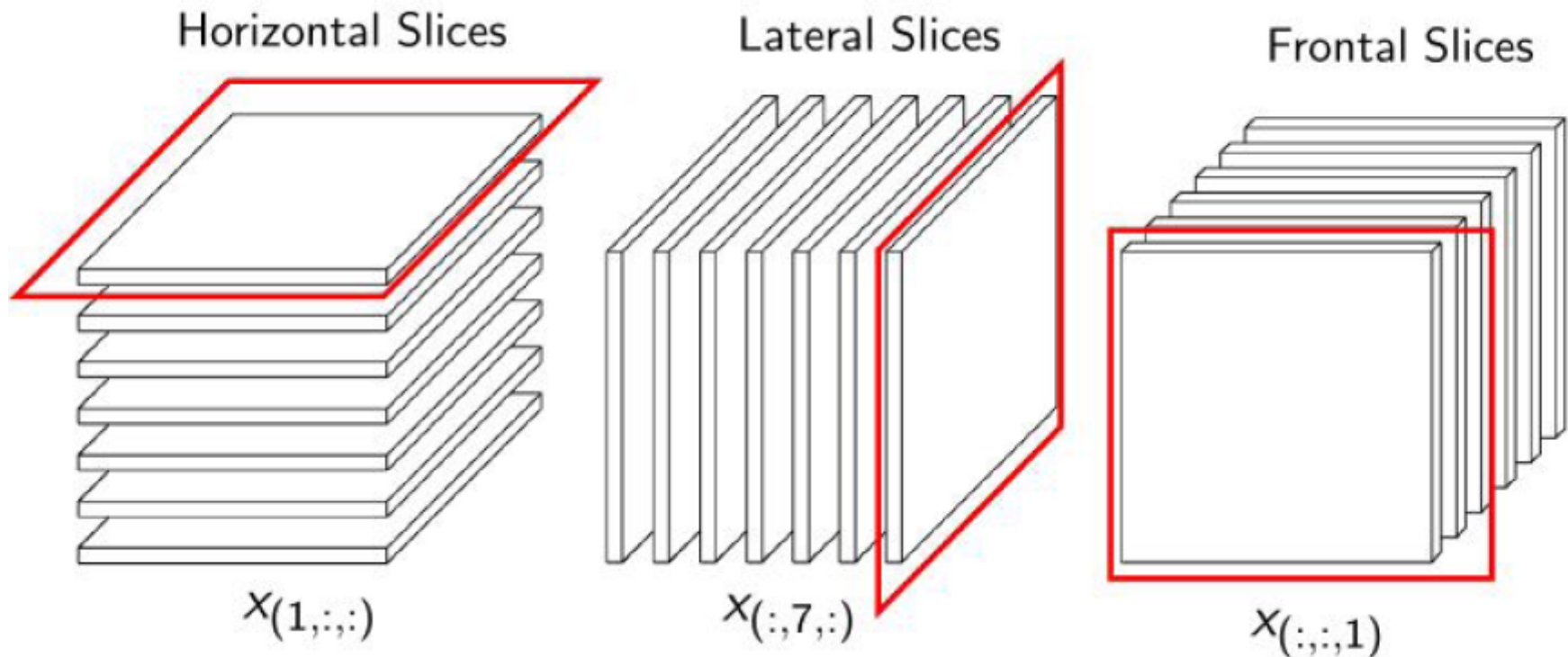


Column(Mode 3)Fibers



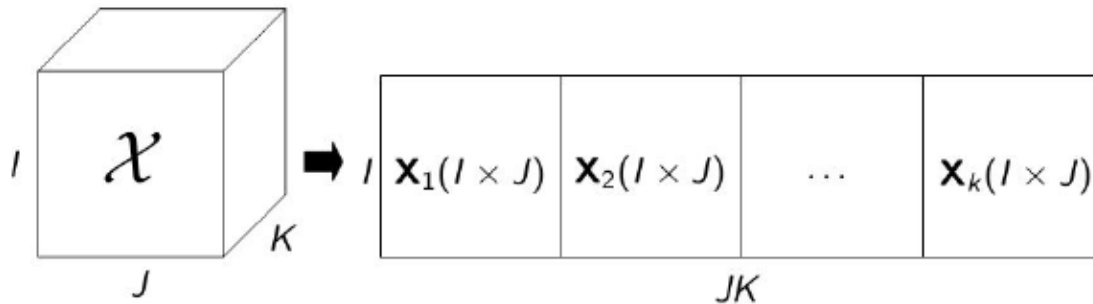
## Definition of Slice

- **Slices** (or slabs) are created when fixing all but two indices.

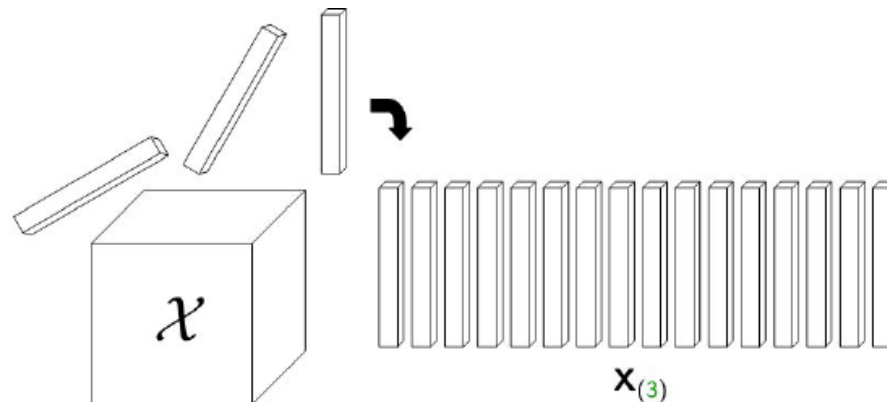


# Matricization based on Slices

- **Vectorization** is to reorder a tensor into a vector.
- **Matricization** is to reorder a tensor into a matrix.



- One can also think to rearrange the fibers into the columns of a matrix.



# The h-Mode Multiplication

Let  $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ ,  $\mathbf{B} \in \mathbb{R}^{M \times J}$ , the 2-mode product of  $\mathcal{X}$  with  $\mathbf{B}$  is defined by

$$\mathcal{Y} = \mathcal{X} \times_2 \mathbf{B} \in \mathbb{R}^{I \times M \times K}$$

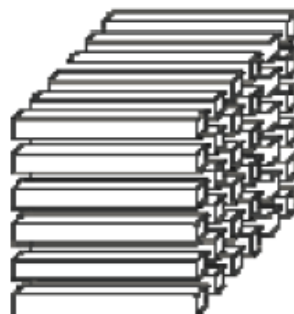
Elementwise

$$y_{imk} = \sum_j x_{ijk} b_{mj}$$

In matrix form

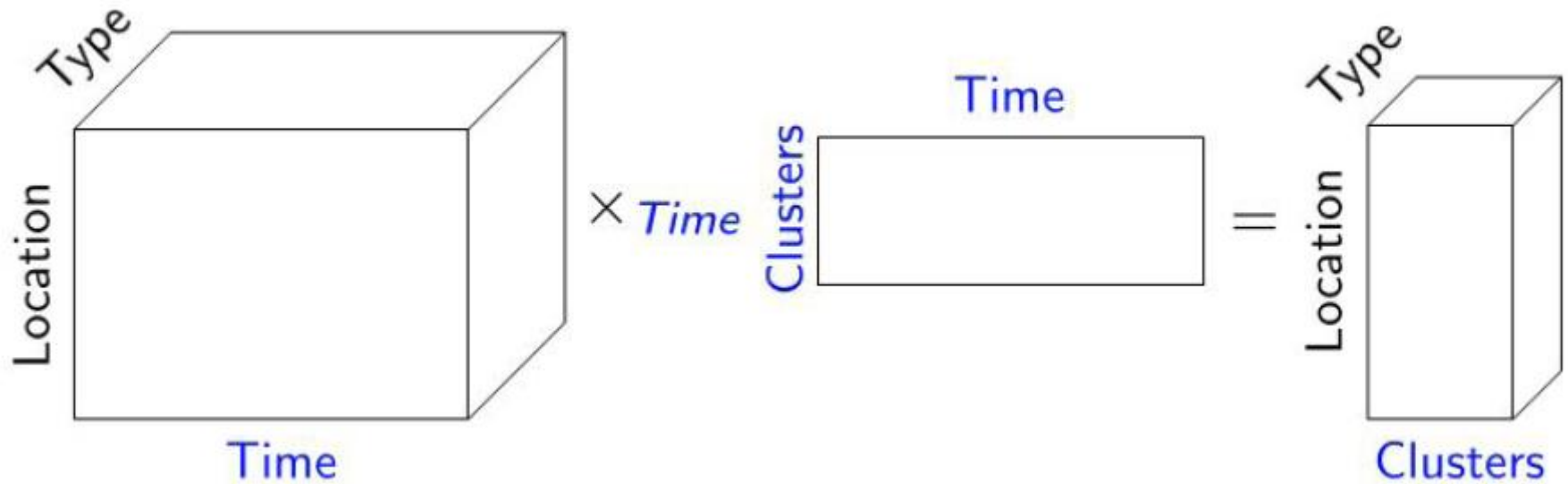
$$\mathbf{Y}_{(2)} = \mathbf{B} \mathbf{X}_{(2)}$$

Multiply each  
row (mode-2)  
fiber by  $\mathbf{B}$





# Illustration of h-Mode Multiplication



# Matrix Decomposition

- Rank decomposition of a matrix

$$M = AB^T \quad \text{with} \quad M \in \mathbb{R}^{n \times m}, A \in \mathbb{R}^{n \times r}, B^T \in \mathbb{R}^{r \times m}$$

where  $r$  represents the rank of the decomposition.

- The rotation problem

$$M \approx A R R^1 B^T$$

- The optimization does not have a unique solution.

$$\min_{\hat{M}} ||M - \hat{M}|| \quad \text{with} \quad \hat{M} = AB^T$$

- Need certain conditions, such as orthogonality, to make matrix decompositions unique.

# Tensor and Its Decomposition

- Tensor is a multi-way array: An **order- $k$**  tensor can be expressed as

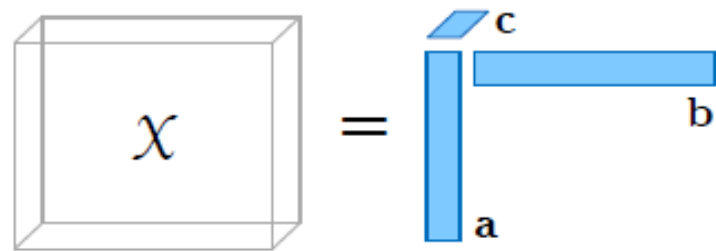
$$\mathbf{X} = (x_{i_1, i_2, \dots, i_k}) \in R^{I_1 \times I_2 \times \dots \times I_k}.$$

- Inner product vs outer product

$$\langle a, b \rangle = a^T b \quad \text{vs} \quad a \odot b = ab^T$$

- Rank-one Tensor:

$$\begin{aligned} \mathbf{Z} &= \mathbf{a}^{(1)} \odot \mathbf{a}^{(2)} \odot \dots \odot \mathbf{a}^{(k)} \\ &= \left( a_{i_1}^{(1)} a_{i_2}^{(2)} \dots a_{i_k}^{(k)} \right) \end{aligned}$$



- Tensor rank: The rank of tensor  $\mathbf{X}$  is defined as the minimum number of rank-one tensors needed to produce  $\mathbf{X}$  as their sum.

## Tensor and Its Decomposition (Con't)

---

- Thus a tensor  $\mathcal{X}$  of rank-L can be expressed as

$$\begin{aligned}\mathcal{X} &= \sum_{r=1}^L \lambda_r \mathbf{a}_r^{(1)} \odot \mathbf{a}_r^{(2)} \odot \cdots \odot \mathbf{a}_r^{(k)} \\ &= \llbracket \boldsymbol{\lambda}; \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(k)} \rrbracket\end{aligned}$$

where  $\mathbf{A}^{(j)} = [\mathbf{a}_1^{(j)}, \dots, \mathbf{a}_L^{(j)}]$  is called the *factor matrix*.

- Tensor Decomposition: generalizing the idea of SVD.
  - Canonical Polyadic Decomposition (CPD)
  - Tucker Decomposition (TD)

## Norm for Tensors

- An easy norm: *Hilbert–Schmidt norm*, defined as

$$\|A\| = \sqrt{\langle A, A \rangle} = \left( \sum_{i_1, \dots, i_d=1}^{n_1, \dots, n_d} |a_{i_1 \dots i_d}|^2 \right)^{\frac{1}{2}}.$$

which can be viewed as an extension of Frobenius norm.

- A *spectral norm*, is defined as

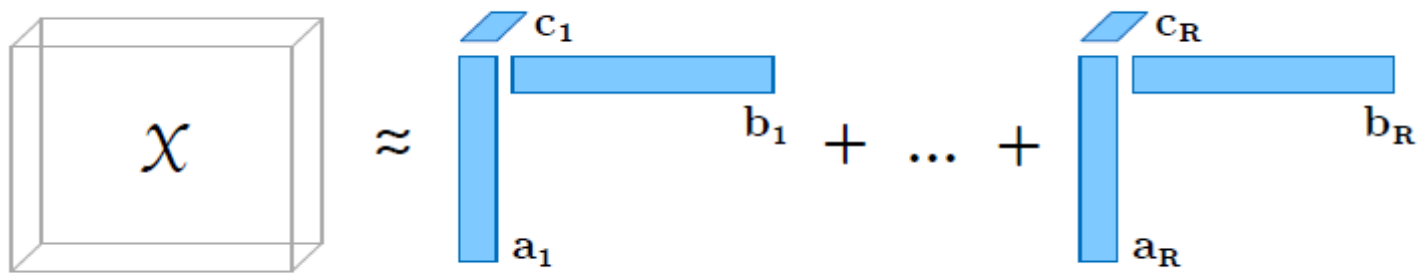
$$\begin{aligned} \|A\|_{\sigma, \mathbb{F}} &:= \sup \left\{ \frac{|\langle A, x_1 \odot \dots \odot x_d \rangle|}{\|x_1\| \cdots \|x_d\|} : 0 \neq x_k \in \mathbb{F}^{n_k} \right\} \\ &= \sup \left\{ |\langle A, u_1 \odot \dots \odot u_d \rangle| : \|u_k\| = 1 \right\}. \end{aligned}$$

- A singular-value norm: *nuclear norm*, is defined as

$$\begin{aligned} \|A\|_{*, \mathbb{F}} &= \inf \left\{ \sum_{i=1}^r |\lambda_i| : A = \sum_{i=1}^r \lambda_i \mathbf{u}_i^{(1)} \odot \dots \odot \mathbf{u}_i^{(d)}, \|\mathbf{u}_i^{(j)}\| = 1, \quad r \in \mathbb{N} \right\} \\ &= \inf \left\{ \sum_{i=1}^r \|\mathbf{u}_1\| \cdots \|\mathbf{u}_d\| : A = \sum_{i=1}^r \mathbf{u}_i^{(1)} \odot \dots \odot \mathbf{u}_i^{(d)} \right\} \end{aligned}$$

# Canonical Polyadic Decomposition (CPD)

- The CPD is rank decomposition, which is to express a tensor as the sum of a finite number of rank-one tensors.

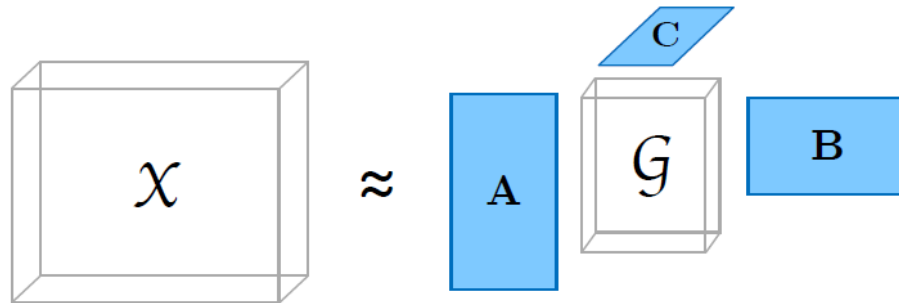


- An order-3 tensor CPD can be formulated as

$$\min_{\hat{\mathcal{X}}} \|\mathcal{X} - \hat{\mathcal{X}}\| \quad \text{where} \quad \hat{\mathcal{X}} = \sum_{r=1}^R a_r \odot b_r \odot c_r = \llbracket A, B, C \rrbracket$$

# Tucker Decomposition (TD)

- Tucker decomposition
  - Decomposes a tensor into a **core tensor** and multiple matrices which correspond to different core scalings along each mode.
  - Tucker decomposition can be seen as a higher-order PCA.

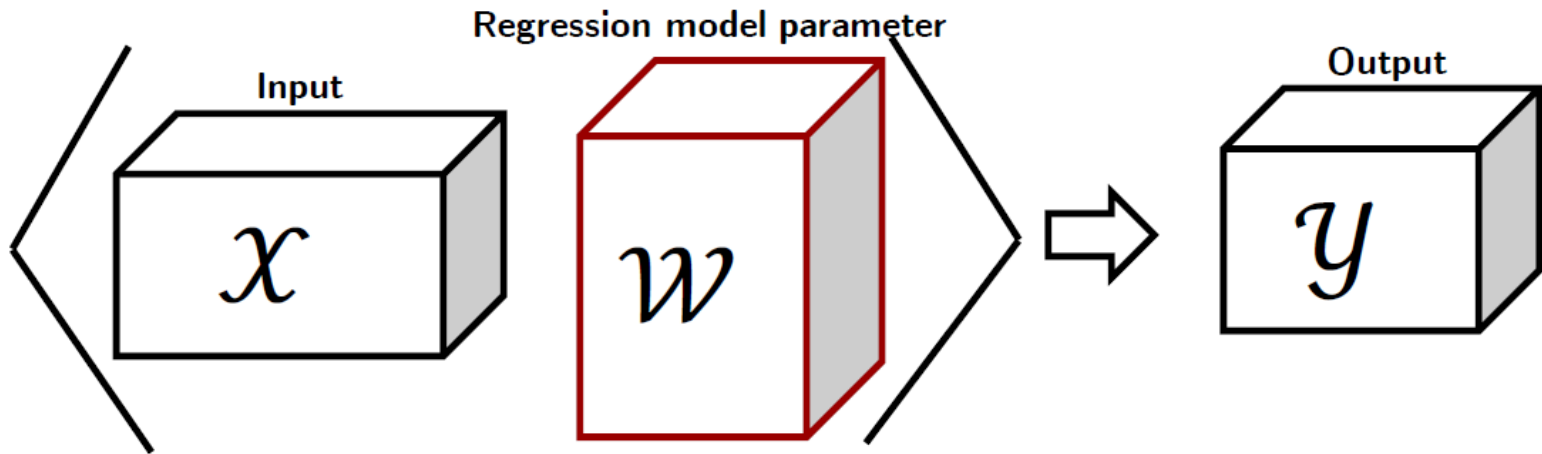


- A order-3 tensor decomposition can be formulated as

$$\begin{aligned} \min_{\hat{\mathcal{X}}} ||\mathcal{X} - \hat{\mathcal{X}}|| \quad \text{with} \quad \hat{\mathcal{X}} &= \sum_{p=1}^P \sum_{q=1}^Q \sum_{r=1}^R g_{pqr} \mathbf{a}_p \odot \mathbf{b}_q \odot \mathbf{c}_r \\ &= \mathcal{G} \times_1 A \times_2 B \times_3 C \\ &= \llbracket \mathcal{G}; A, B, C \rrbracket \end{aligned}$$

# Tensor Regression

- Tensor Regression: large-scale supervised learning from multi-way data
- Goal: learn a regression model with multi-linear parameters





# Low-Rank Representation

- Low-rank structures can capture multi-linear correlations.



A 4x4 matrix representing user ratings for movies. The rows represent users (indicated by avatars on the left) and the columns represent movies (indicated by movie covers on top). The ratings are as follows:

	5	2	4	3
	4	0	1	1
	3	?	?	1
	5	?	?	?

**Collaborative Filtering**

- Tucker decomposition: high-order SVD

$$\begin{array}{c} I \\ \text{---} \\ \text{---} \mathcal{W} \text{---} \\ \text{---} \\ J \quad K \end{array} \approx_{R_1} \begin{array}{c} \text{---} \\ \text{---} \mathcal{S} \text{---} \\ \text{---} \\ R_2 \quad R_3 \end{array} \times_1 \begin{array}{c} I \\ \text{---} \\ \text{---} U_1 \text{---} \\ \text{---} \\ R_1 \end{array} \times_2 \begin{array}{c} J \\ \text{---} \\ \text{---} U_2 \text{---} \\ \text{---} \\ R_2 \end{array} \times_3 \begin{array}{c} \text{---} \\ \text{---} U_3 \text{---} \\ \text{---} \\ K \quad R_3 \end{array}$$

# Low-Rank Tensor Regression

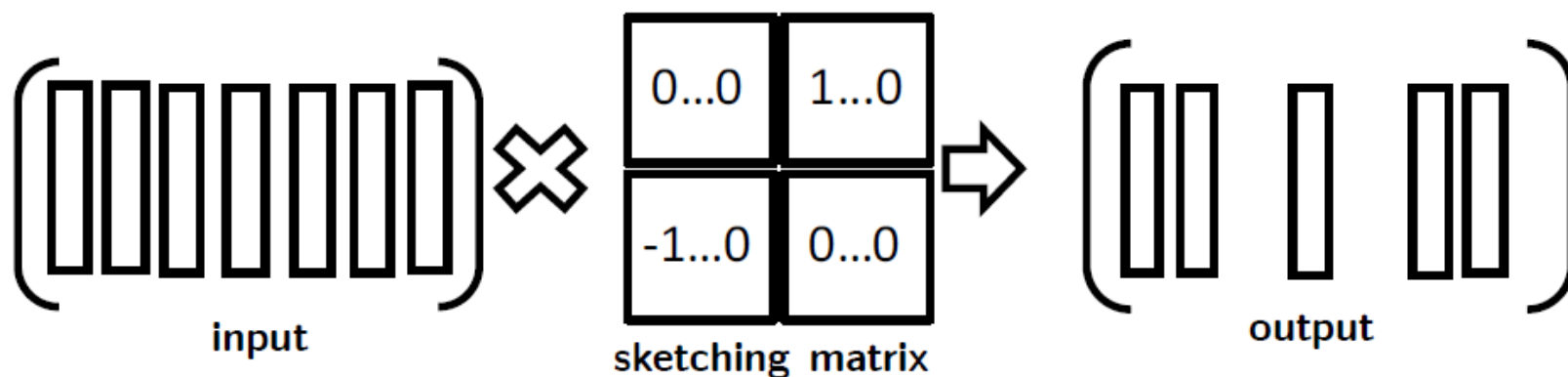
---

- Predictor tensor  $\mathcal{X}$ ; response tensor  $\mathcal{Y}$
- Regression model  $\langle \mathcal{X}, \mathcal{W} \rangle$ : e.g.  $\sum_{m=1}^M \mathcal{X}_{::,m} \mathcal{W}_{::,m}$
- Loss function  $\mathcal{L}(\hat{\mathcal{Y}}; \mathcal{Y})$ : e.g.  $\|\hat{\mathcal{Y}} - \mathcal{Y}\|_F^2$
- Goal: Learn a parameter tensor  $\mathcal{W}$  with low-rank constraint

$$\begin{aligned} \mathcal{W}^* = \operatorname{argmin}_{\mathcal{W}} & \hat{\mathcal{L}}(f(\mathcal{X}, \mathcal{W}); \mathcal{Y}) \\ \text{s.t.} \quad & \operatorname{rank}(\mathcal{W}) \leq R \end{aligned}$$

# Subsampled Tensor Projected Gradient (TPG)

- Data  $\Rightarrow$  Random sketching [Woodruff 2014]
- Model  $\Rightarrow$  Iterative hard thresholding [Thomas and Davies 2009]



- Projected gradient descent:  $\mathcal{W}^{k+1} = P_R (\mathcal{W}^k - \eta \nabla \mathcal{W}^k)$ 
  1. Gradient descent step
  2. Low-rank projection step

# Brief on Random Sketching

---

- Let us consider a least squares estimation problem,

$$\beta_{OLS} = \arg \min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|_2^2$$

where  $Y$  is response vector and  $X$  is  $n$ -by- $p$  regression matrix.

- Suppose that there is a  $r$ -by- $n$  **sketching matrix  $S$** . Then the sketched problem for parameter estimation

$$\beta_S \in \arg \min_{\beta \in \mathbb{R}^p} \|SY - SX\beta\|_2^2.$$

- It is shown that in Drineas et al. (2011, 2012), for any arbitrary  $(X; Y)$ ,

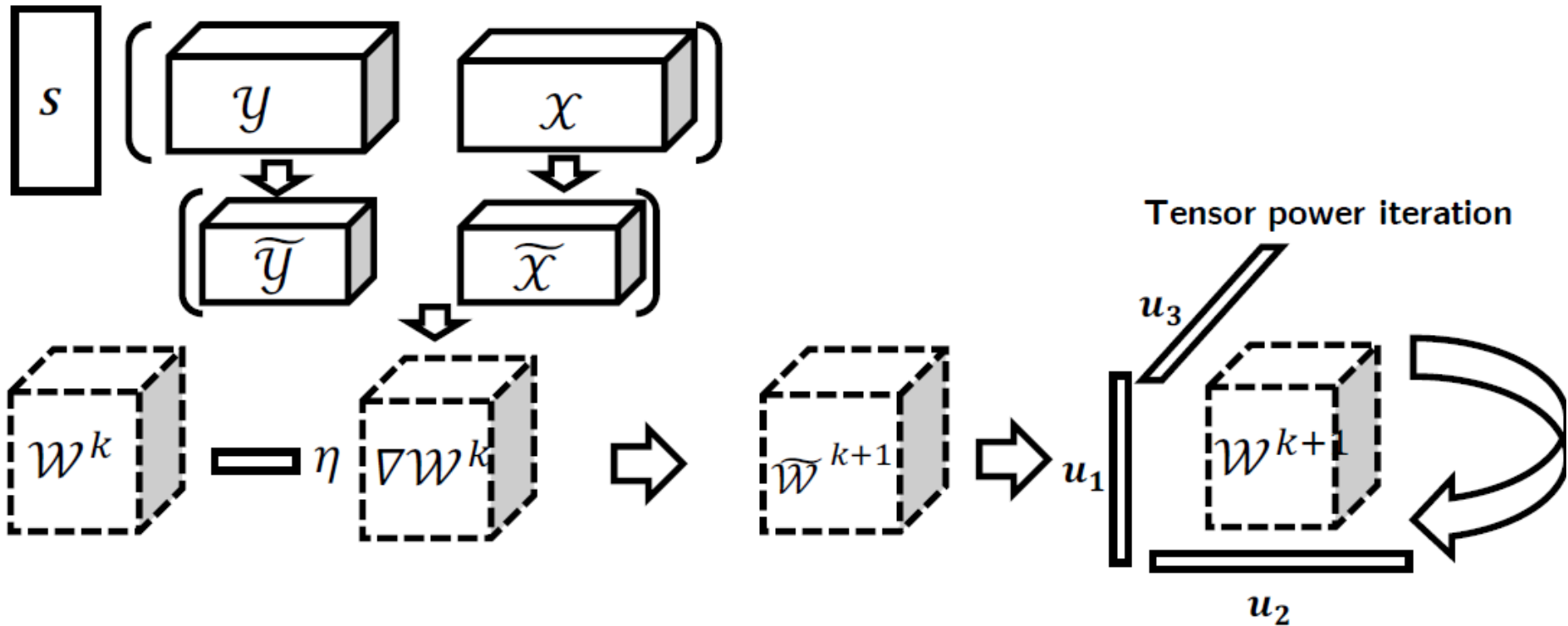
$$\|Y - X\beta_S\|_2^2 \leq (1 + \kappa) \|Y - X\beta_{OLS}\|_2^2,$$

with high probability for some pre-specified error parameter  $\kappa \in (0,1)$ .

# Subsampled Tensor Projected Gradient (TPG)

- Random sketching as data subsampling
- Iterative hard thresholding as dimensional reduction

Sketching  
matrix



## Example: GLM

---

The standard linear regression model  $\mathbf{x} \in \mathbb{R}^p$ ,  $y = \beta^T \mathbf{x} + \alpha + \varepsilon$ ,  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$  can be written

$$\mu = \beta^T \mathbf{x} + \alpha \quad y \sim \mathcal{N}(\mu, \sigma^2)$$

where  $\mu = \mathbb{E}(Y|\mathbf{x})$

A **generalized linear regression model (GLM)** extends this to

$$g(\mu) = \beta^T \mathbf{x} + \alpha \quad y \sim \mathcal{EF}(\mu, \phi)$$

- $\mathcal{EF}(\mu, \phi)$  is any exponential family distribution (e.g. Normal, Poisson, Binomial)
- $g(\cdot)$  is any smooth monotonic link function
- $\beta^T \mathbf{x} + \alpha (= \eta)$  is the linear predictor

## Example: GLM with Matrix Predictor

---

In classical **GLM**  $Y$  belongs to an exponential family with **PMF**

$$p(y|\theta, \phi) = \exp\left\{\frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi)\right\}$$

The **GLM** relates  $\mathbf{x} \in \mathbb{R}^p$  to the mean  $\mu = \mathbb{E}(Y|\mathbf{x})$  by

$$g(\mu) = \eta = \alpha + \beta^T \mathbf{x}$$

The **GLM** for the matrix predictor  $\mathbf{X}$  given by

$$g(\mu) = \eta = \alpha + \gamma^T \mathbf{z} + \beta_1^T \mathbf{X} \beta_2$$

## Example: GLM with Tensor Predictor

---

The **GLM** with the systematic part for tensor predictor given by

$$g(\mu) = \eta = \alpha + \gamma^T \mathbf{z} + \langle \mathcal{B}, \mathcal{X} \rangle$$

- $D$ -dimensional tensor predictor  $\mathcal{X} \in \mathbb{R}^{p_1 \times \cdots \times p_D}$
- $D$ -dimensional coefficient tensor  $\mathcal{B} \in \mathbb{R}^{p_1 \times \cdots \times p_D}$
- $\mathcal{B}$  has  $\prod_{d=1}^D p_d$  parameters, which is **ultrahigh dimensional** and **far exceeds sample size**



## Example: GLM with Tensor Predictor (Con't)

- Univariate outcome  $Y$  belongs to exponential family
- Tensor covariate  $\mathcal{X} \in \mathbb{R}^{p_1 \times \dots \times p_D}$
- Assume coefficient tensor  $\mathcal{B}$  has a rank- $R$  decomposition  $[\mathbf{B}_1, \dots, \mathbf{B}_D]$  where  $\mathbf{B}_d \in \mathbb{R}^{p_d \times R}$

**Generalized linear CP tensor regression model** (Zhou et al. 2013) with the systematic part given by

$$g(\mu) = \eta = \alpha + \gamma^T \mathbf{z} + \left\langle \sum_{r=1}^R \beta_1^{(r)} \circ \dots \circ \beta_D^{(r)}, \mathcal{X} \right\rangle$$

$$= \alpha + \gamma^T \mathbf{z} + \left\langle (\mathbf{B}_D \odot \dots \odot \mathbf{B}_1) \mathbf{1}_R, \text{vec}(\mathcal{X}) \right\rangle$$

- **substantial reduction in dimensionality** to the scale of  $R \times \sum_{d=1}^D p_d$

# Parameter Estimation via Regularization

---

Maximize a regularized log-likelihood function

$$\ell(\alpha, \gamma, \mathbf{B}_1, \dots, \mathbf{B}_D) - \sum_{d=1}^D \sum_{r=1}^R \sum_{i=1}^{p_d} P_{\lambda}(|\beta_{di}^{(r)}|, \rho)$$

- scalar penalty function  $P_{\lambda}(|\beta|, \rho)$
- **power family**  $P_{\lambda}(|x|, \rho) = \rho|\beta|^{\lambda}$ ,  $\lambda \in (0, 2]$
- in particular lasso ( $\lambda = 1$ )

## Comments

---

***Thank You!***