

Tensor Regression

CAIT

Skolkovo Institute of Science and Technology (SKOLTECH), Moscow, Russia

March 17, 2022

Problem (Feature extraction for 2-D samples)

- A set of K data matrices $\mathbf{X}^{(k)} \in \mathbb{R}^{I_1 \times I_2}$, ($k = 1, \dots, K$) that belong to C different classes.
- Feature extraction for all samples by simultaneous matrix factorizations:

$$\mathbf{X}^{(k)} \approx \mathbf{A}^{(1)} \mathbf{F}^{(k)} \mathbf{A}^{(2)T}$$

$\mathbf{A}^{(1)} \in \mathbb{R}^{I_1 \times R_1}$ and $\mathbf{A}^{(2)} \in \mathbb{R}^{I_2 \times R_2}$,
are two basis matrices

$\mathbf{F}^{(k)} \in \mathbb{R}^{R_1 \times R_2}$ are features,
 $R_n \leq I_n$.

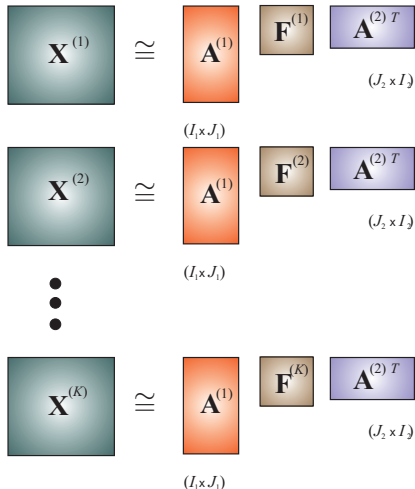


Figure: Simultaneous matrix factorizations

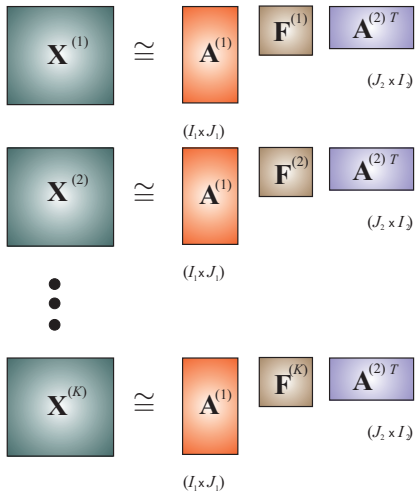


Figure: Simultaneous approximations of 2-D samples

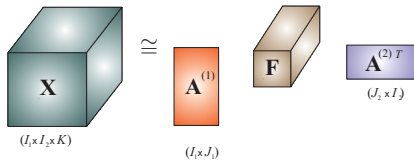


Figure: Equivalent 3D tensor decomposition: TUCKER-2 model.

- If $\mathbf{F}^{(k)}$ are diagonal matrices, the model corresponds to CPD
- If $\mathbf{A}^{(n)}$ are orthogonal and $\mathbf{F}^{(k)}$ are dense, the model becomes HOSVD or multi-way PCA De Lathauwer et al. (2001)
- If $\mathbf{A}^{(n)}$ are nonnegative, the model becomes Tri Nonnegative Matrix Factorization of data $\mathbf{X}^{(k)}$.
- if $\mathbf{X}^{(k)}$ are positive-definite covariance or cumulant matrices, model becomes closely related to Joint Diagonalization, ICA.

Problem (N-way feature extraction)

Find N common factors

$\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R_n}, (n = 1, 2, \dots, N)$

from K simultaneous

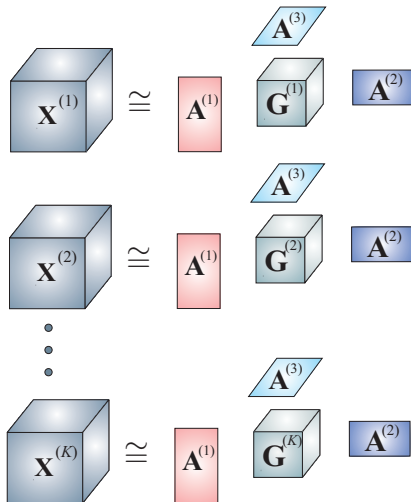
decompositions of K samples

$\mathcal{X}^{(k)} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$

$\mathcal{X}^{(k)} \approx \mathcal{G}^{(k)} \times_1 \mathbf{A}^{(1)} \dots \times_N \mathbf{A}^{(N)},$

where $\mathcal{G}^{(k)} \in \mathbb{R}^{R_1 \times R_2 \times \dots \times R_N}$

consist of features of $\mathcal{X}^{(k)}$.



Problem (Global TUCKER decomposition)

The N common bases of K samples $\mathcal{X}^{(k)} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ are factor matrices, $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R_n}$, for $n = 1, \dots, N$, in the TUCKER- N decomposition of the concatenation tensor along the mode- $(N+1)$, that is

$$\mathcal{X} \approx \mathcal{G} \times_1 \mathbf{A}^{(1)} \times_2 \mathbf{A}^{(2)} \dots \times_N \mathbf{A}^{(N)},$$

$$\mathcal{X} = \text{cat}(N+1, \mathcal{X}^{(1)}, \mathcal{X}^{(2)}, \dots, \mathcal{X}^{(K)})$$

\mathcal{G} represents extracted features for the training samples.

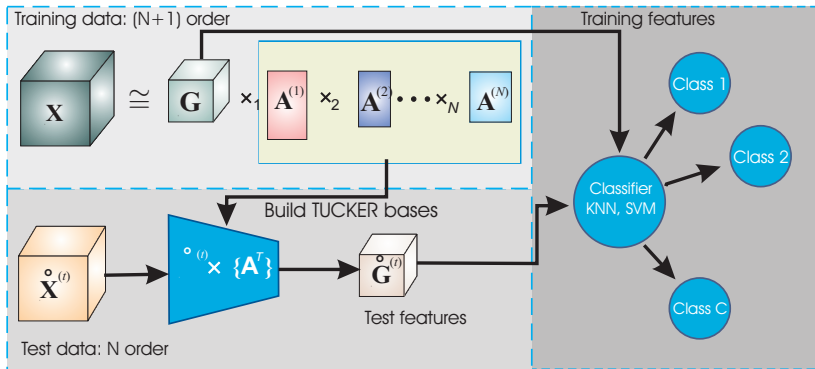


Figure: Conceptual diagram illustrating a classification procedure based on the TUCKER decomposition of the concatenated tensor of all sampling training data. Reduced features are obtained by projecting the data tensor onto the feature subspace spanned by basis factors (bases).

Algorithm 1: HOOI for Orthogonal Bases

Input: \mathcal{X} : concatenation tensor of all training samples

$$I_1 \times I_2 \times \cdots \times I_N \times K,$$

Output: N orthogonal factors $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times R_n}$ and a core tensor $\mathcal{G} \in \mathbb{R}^{R_1 \times R_2 \times \cdots \times R_N \times K}$.

begin

 HOSVD or random initialization for all factors $\mathbf{U}^{(n)}$

repeat

for $n = 1$ to N **do**

$\mathcal{W}^{(-n)} = \mathcal{X} \times_{-(n, N+1)} \{\mathbf{U}^{(T)}\}$

$[\mathbf{U}^{(n)}, \Sigma^{(n)}, \mathbf{V}^{(n)}] = \text{svds}(\mathcal{W}_{(n)}^{(-n)}, R_n, \text{'LM'})$

end

until a stopping criterion is met

$\mathcal{G} = \mathcal{W}^{(-N)} \times_N \mathbf{U}^{(N)T}$

end

Discriminant Analysis

Suppose we have a set of training samples $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ belonging to C classes. (L)DA maximizes the between-class scatter \mathbf{S}_b while minimizing the within-class scatter \mathbf{S}_w .

Problem (Multilinear Discriminant Analysis)

$$\arg \max_{\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(N)}} \frac{\sum_{c=1}^C K_c \left\| \bar{\mathcal{G}}^{(c)} - \bar{\bar{\mathcal{G}}} \right\|_F^2}{\sum_k^K \left\| \mathcal{G}^{(k)} - \bar{\mathcal{G}}^{(c_k)} \right\|_F^2}$$

$\bar{\mathcal{G}}^{(c)} = \left(\sum_{k \in \mathcal{I}_c} \mathcal{G}^{(k)} \right) / K_c$: mean tensor of the c -th class

$\bar{\bar{\mathcal{G}}} = \left(\sum_{k=1}^K \mathcal{G}^{(k)} \right) / K$: mean tensor of the whole training features

c_k class to which $\mathcal{X}^{(k)}$ belongs

Within-class scatter matrix

$$\begin{aligned}\sum_{k=1}^K \left\| \mathbf{g}^{(k)} - \bar{\mathbf{g}}^{(c_k)} \right\|_F^2 &= \sum_{k=1}^K \left\| \left(\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(c_k)} \right) \times_{-n} \{ \mathbf{U}^T \} \times_n \mathbf{U}^{(n)T} \right\|_F^2 \\ &= \sum_{k=1}^K \text{tr} \left[\mathbf{U}^{(n)T} \langle \tilde{\mathbf{z}}^{(k)-n}, \tilde{\mathbf{z}}^{(k)-n} \rangle_{-n} \mathbf{U}^{(n)} \right] \\ &= \text{tr} \left[\mathbf{U}^{(n)T} \mathbf{S}_w^{-n} \mathbf{U}^{(n)} \right],\end{aligned}$$

The within-class scatter matrix \mathbf{S}_w^{-n} is defined as

$$\begin{aligned}\mathbf{S}_w^{-n} &= \langle \tilde{\mathbf{z}}^{-n}, \tilde{\mathbf{z}}^{-n} \rangle_{-n} \\ \tilde{\mathbf{x}}^{(k)} &= \mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(c_k)}, \\ \tilde{\mathbf{z}}^{-n} &= \tilde{\mathbf{x}} \times_{-(n,N+1)} \{ \mathbf{U}^T \},\end{aligned}$$

Between-class scatter matrix

$$\begin{aligned}\sum_{c=1}^C K_c \left\| \bar{\mathbf{g}}^{(c)} - \bar{\bar{\mathbf{g}}} \right\|_F^2 &= \sum_{c=1}^C K_c \left\| \left(\bar{\mathbf{x}}^{(c)} - \bar{\bar{\mathbf{x}}} \right) \times_{-n} \{ \mathbf{U}^T \} \times_n \mathbf{U}^{(n)T} \right\|_F^2 \\ &= \sum_{c=1}^C \text{tr} \left[\mathbf{U}^{(n)T} \langle \check{\mathbf{z}}^{(c)-n}, \check{\mathbf{z}}^{(c)-n} \rangle_{-n} \mathbf{U}^{(n)} \right] \\ &= \text{tr} \left[\mathbf{U}^{(n)T} \mathbf{S}_b^{-n} \mathbf{U}^{(n)} \right],\end{aligned}$$

The between-class scatter matrix \mathbf{S}_b^{-n} is defined as

$$\begin{aligned}\mathbf{S}_b^{-n} &= \langle \check{\mathbf{z}}^{-n}, \check{\mathbf{z}}^{-n} \rangle_{-n} \\ \check{\mathbf{x}}^{(c)} &= \sqrt{K_c} \left(\bar{\mathbf{x}}^{(c)} - \bar{\bar{\mathbf{x}}} \right) \\ \check{\mathbf{z}}^{-n} &= \check{\mathbf{x}} \times_{-(n, N+1)} \{ \mathbf{U}^T \}\end{aligned}$$

Multilinear Discriminant Analysis

Problem (Discriminant Bases)

$$\max_{\mathbf{U}^{(n)}} \frac{\text{tr}[\mathbf{U}^{(n)T} \mathbf{S}_b^{-n} \mathbf{U}^{(n)}]}{\text{tr}[\mathbf{U}^{(n)T} \mathbf{S}_w^{-n} \mathbf{U}^{(n)}]} = \max_{\mathbf{U}^{(n)}} \frac{\text{tr}[\mathbf{S}_b]}{\text{tr}[\mathbf{S}_w]}$$

find R_n leading left eigenvectors of matrices

- $\mathbf{S}_w^{-1} \mathbf{S}_b$

- $\mathbf{S}_b - \varphi \mathbf{S}_w$

$$\arg \max_{\mathbf{U}^{(n)}} \frac{\text{tr}[\mathbf{S}_b]}{\text{tr}[\mathbf{S}_w]} = \arg \max_{\mathbf{U}^{(n)}} \frac{\text{tr}[\mathbf{S}_b + \mathbf{S}_w]}{\text{tr}[\mathbf{S}_w]}$$

$$\text{tr}[\mathbf{S}_t] = \text{tr}[\mathbf{S}_b + \mathbf{S}_w] = \sum_{k=1}^K \|\mathbf{g}^{(k)} - \bar{\mathbf{g}}\|_F^2$$

Multilinear Discriminant Analysis

Problem (Discriminant Bases)

$$\max_{\mathbf{U}^{(n)}} \frac{\text{tr}[\mathbf{U}^{(n)T} \mathbf{S}_b^{-n} \mathbf{U}^{(n)}]}{\text{tr}[\mathbf{U}^{(n)T} \mathbf{S}_w^{-n} \mathbf{U}^{(n)}]} = \max_{\mathbf{U}^{(n)}} \frac{\text{tr}[\mathbf{S}_b]}{\text{tr}[\mathbf{S}_w]}$$

find R_n leading left eigenvectors of matrices

- $\mathbf{S}_w^{-1} \mathbf{S}_b$

- $\mathbf{S}_b - \varphi \mathbf{S}_w$

$$\arg \max_{\mathbf{U}^{(n)}} \frac{\text{tr}[\mathbf{S}_b]}{\text{tr}[\mathbf{S}_w]} = \arg \max_{\mathbf{U}^{(n)}} \frac{\text{tr}[\mathbf{S}_b + \mathbf{S}_w]}{\text{tr}[\mathbf{S}_w]}$$

$$\text{tr}[\mathbf{S}_t] = \text{tr}[\mathbf{S}_b + \mathbf{S}_w] = \sum_{k=1}^K \left\| \mathbf{g}^{(k)} - \bar{\mathbf{g}} \right\|_F^2$$

Regularized Multilinear Discriminant Analysis

Problem (Regularization)

$$\varphi = \max_{\mathbf{U}^{(n)}} \frac{\text{tr}[\mathbf{S}_t]}{\text{tr}[\alpha \mathbf{S}_w + (1 - \alpha) \mathbf{I}]}, \quad (1)$$

where $0 \leq \alpha \leq 1$.

The basis factors $\mathbf{U}^{(n)}$ are R_n leading eigenvectors of matrices

$$\bullet (\alpha \mathbf{S}_w + (1 - \alpha) \mathbf{I})^{-1} \mathbf{S}_b \qquad \bullet \mathbf{S}_b - \varphi (\alpha \mathbf{S}_w + (1 - \alpha) \mathbf{I})$$

- For $\alpha = 1$, the optimization problem (1) simplifies MDA.
- For $\alpha = 0$, the optimization problem (1) simplifies the TUCKER decomposition of the training samples after centering.

Algorithm 2: HODA Algorithm for Feature Extraction

input : \mathcal{X} : Concatenated tensor of K training samples $I_1 \times I_2 \times \dots \times I_N \times K$

output: $\mathbf{U}^{(n)}$: N orthogonal basis factors $I_n \times R_n$ ($n = 1, 2, \dots, N$)

begin

Initialize $\mathbf{U}^{(n)}$

Calculate $\tilde{\mathcal{X}}$, and $\check{\mathcal{X}}$

repeat

for $n = 1$ to N **do**

$$\tilde{\mathcal{Z}} = \tilde{\mathcal{X}} \times_{-(n, N+1)} \{\mathbf{U}^T\}, \quad \mathbf{S}_w^{-n} = \langle \tilde{\mathcal{Z}}, \tilde{\mathcal{Z}} \rangle_{-n}$$

$$\check{\mathcal{Z}} = \check{\mathcal{X}} \times_{-(n, N+1)} \{\mathbf{U}^T\}, \quad \mathbf{S}_b^{-n} = \langle \check{\mathcal{Z}}, \check{\mathcal{Z}} \rangle_{-n}$$

$$\varphi = \frac{\text{trace}(\mathbf{u}^{(n)T} \mathbf{S}_b^{-n} \mathbf{u}^{(n)})}{\text{trace}(\mathbf{u}^{(n)T} \mathbf{S}_w^{-n} \mathbf{u}^{(n)})}$$

$$[\mathbf{U}^{(n)}, \Lambda] = \text{eigs}(\mathbf{S}_b^{-n} - \varphi \mathbf{S}_w^{-n}, R_n, \text{'LM'})$$

end

until a criterion is met

end

Algorithm 2: HODA Algorithm for Feature Extraction

input : \mathcal{X} : Concatenated tensor of K training samples $I_1 \times I_2 \times \cdots \times I_N \times K$

output: $\mathbf{U}^{(n)}$: N orthogonal basis factors $I_n \times R_n$ ($n = 1, 2, \dots, N$)

begin

 Initialize $\mathbf{U}^{(n)}$

 Calculate $\tilde{\mathcal{X}}$, and $\check{\mathcal{X}}$

repeat

for $n = 1$ to N **do**

$\tilde{\mathcal{Z}} = \tilde{\mathcal{X}} \times_{-(n, N+1)} \{\mathbf{U}^T\}$, $\mathbf{S}_w^{-n} = \langle \tilde{\mathcal{Z}}, \tilde{\mathcal{Z}} \rangle_{-n}$

$\check{\mathcal{Z}} = \check{\mathcal{X}} \times_{-(n, N+1)} \{\mathbf{U}^T\}$, $\mathbf{S}_b^{-n} = \langle \check{\mathcal{Z}}, \check{\mathcal{Z}} \rangle_{-n}$

$[\mathbf{U}^{(n)}, \Lambda] = \text{eigs}(\mathbf{S}_b^{-n}, \mathbf{S}_w^{-n}, R_n, \text{'LM'})$

end

until *a criterion is met*

end

- Ranking the features in a descending order of their Fisher scores

$$\varphi(i) = \frac{\sum_{c=1}^C K_c (\bar{g}_i^{(c)} - \bar{\bar{g}}_i)^2}{\sum_{k=1}^K (g_i^{(k)} - \bar{g}_i^{(c_k)})^2}, \quad (i = 1, 2, \dots, L), \quad (2)$$

where $\mathbf{g}^{(k)} = \text{vec}(\mathcal{G}^{(k)})$.

- Choose significant features

Graph Regularizer I

- Denote by $\tilde{\mathbf{x}}_k = \mathbf{x}_k - \bar{\mathbf{x}}$ the centered samples,
 $\tilde{\mathbf{X}} = [\dots, \tilde{\mathbf{x}}_k, \dots]$
- The Total scatter matrix

$$\begin{aligned}\mathbf{S}_t &= \mathbf{S}_w + \mathbf{S}_b \\ &= \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}\end{aligned}$$

- The between-scatter matrix

$$\mathbf{S}_b = \tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{X}}$$

where $\mathbf{W} = \text{diag}(\frac{1}{K_1} \mathbf{1}_{K_1 \times K_1}, \dots, \frac{1}{K_C} \mathbf{1}_{K_C \times K_C})$

- Graph Regularizer is given in form of

$$\frac{1}{2} \sum_{i \neq j} w_{i,j} \|\mathcal{G}_i - \mathcal{G}_j\|_F^2 \quad (3)$$

where $\mathbf{W} = [w_{i,j}]$ is an adjacency matrix, defines the similarity between i -th and j -th samples.

- Similarity can be calculated as cosine similarity

$$w_{i,j} = \frac{\mathbf{x}_i^T \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|} \quad (4)$$

or

$$w_{i,j} = \begin{cases} \exp(-\frac{d_{i,j}^2}{2\sigma^2}), & \text{if } (i, j) \text{ in the same class} \\ 0 & \end{cases} \quad (5)$$

Graph Regularizer II

where $d_{i,j}^2$ measures Euclidean distance between two samples, σ controls the sensitivity of the graph weights to the distances

- Laplacian matrix $\mathbf{L} = \mathbf{D} - \mathbf{W}$ where \mathbf{L} is Laplacian matrix of the sample similarity based graph and \mathbf{D} is the degree matrix. The degree matrix \mathbf{D} is a diagonal matrix with its i -th diagonal entry computed as $\sum_{j=1} w_{ij}$

Local Discriminant Analysis

Denote by k_w and k_b the number of within and between nearest neighbours

$$\mathbf{W}_w(i, j) = \begin{cases} \frac{1}{k_w} & j \in I_{w, k_w}(i) \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

where $I_{w, k_w}(i)$ denotes index set of top- k_w samples closest to this to the i -th sample and in the same class with this sample.

$$\mathbf{W}_b(i, j) = \begin{cases} \frac{1}{k_w + k_b} - \frac{1}{k_w} & j \in I_{b, k_b}(i) \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

where $I_{b, k_b}(i)$ denotes index set of top- k_b samples closest to this to the i -th sample but not in the same class with this sample.

De Lathauwer, L., de Moor, B., and Vandewalle, J. (2001). A multilinear singular value decomposition. *SIAM Journal of Matrix Analysis and Applications*, 21:1253–1278.