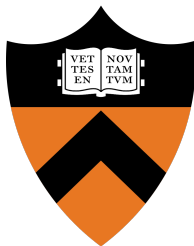


Tensor decomposition and completion



Yuxin Chen

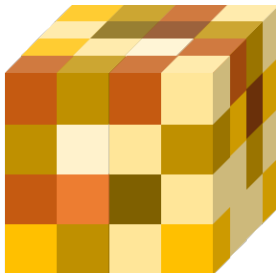
Princeton University, Fall 2018

Outline

- Tensor decomposition
- Latent variable models & tensor decomposition
- Tensor power method
- Tensor completion

Tensor decomposition

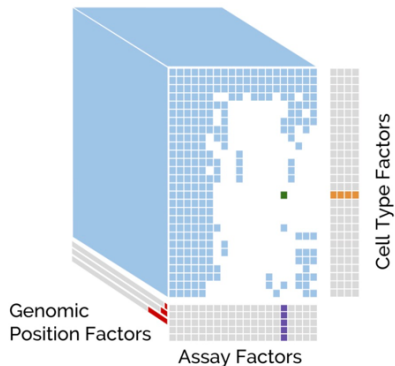
Tensor



An order- d tensor $\mathbf{T} = [T_{i_1, \dots, i_d}]_{1 \leq i_1, \dots, i_d \leq n}$ is a d -way array

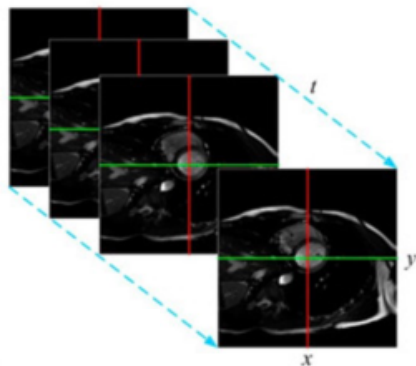
- a matrix is a tensor of order 2

Ubiquity of high-dimensional tensor data



computational genomics

— *fig. credit: Schreiber et al. 19*



dynamic MRI

— *fig. credit: Liu et al. 17*

Basics

- **Rank-1 tensor:** $T = x \otimes x \otimes x$ denotes a tensor such that

$$T_{i_1, \dots, i_d} = x_{i_1} x_{i_2} \cdots x_{i_d}$$

- the inner product of two tensors T and A :

$$\langle T, A \rangle := \sum_{i_1, \dots, i_d} T_{i_1, \dots, i_d} A_{i_1, \dots, i_d}$$

- the Frobenius norm of a tensor T :

$$\|T\|_F := \sqrt{\sum_{i_1, \dots, i_d} T_{i_1, \dots, i_d}^2}$$

- the operator norm of an order- d tensor T :

$$\|T\| = \max_{\{\mathbf{u}_i\}: \|\mathbf{u}_i\|_2=1} \langle T, \mathbf{u}_1 \otimes \cdots \otimes \mathbf{u}_d \rangle$$

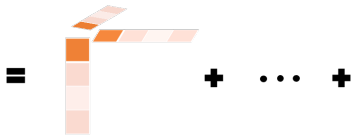
Tensor decomposition

Suppose we observe an order- d tensor

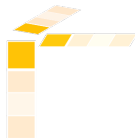
$$\mathbf{T} = \sum_{i=1}^r \lambda_i \mathbf{u}_i \otimes \mathbf{u}_i \otimes \cdots \otimes \mathbf{u}_i$$



true tensor



rank-1 tensor



rank-1 tensor

Question: can we recover $\{\mathbf{u}_i\}$ and $\{\lambda_i\}$ given \mathbf{T} ?

- if $d = 2$ (matrix case), it is often not recoverable; what if $d \geq 3$?
- this question arises in a number of latent-variable models

Latent variable models and tensor decomposition

Notation

- probability simplex

$$\Delta_n := \{\mathbf{z} \in \mathbb{R}^n \mid z_i \geq 0, \forall i; \mathbf{1}^\top \mathbf{z} = 1\}$$

- any vector $\mathbf{w} \in \Delta_n$ represents a distribution (or probability mass function) over n objects

A simple topic model

Consider a collection of documents

- r : the number of distinct topics
- n : the number of distinct words in vocabulary

A simple topic model

Consider a collection of documents

- each time, draw 3 words as follows
 - pick $\underbrace{\text{a topic } h}_{\text{latent variable}}$ according to distribution $[w_1, \dots, w_r] \in \Delta_r$ s.t.
$$\mathbb{P}\{h = j\} = w_j, \quad 1 \leq j \leq r$$
 - given topic h , draw 3 independent words from this topic according to the distribution

$$\underbrace{\mu_h}_{\text{determined only by the topic}} \in \Delta_n$$

Goal: recover $\{\mu_i\}$ and $\{w_i\}$ from the collected samples

Moment method for the topic model

Denote the 3 words we draw as $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)} \in \mathbb{R}^n$:

$$\mathbf{x}^{(i)} = \mathbf{e}_j \quad \text{if the } i\text{th word is } j$$

It is straightforward to check

$$\mathbf{M}_2 := \mathbb{E}[\mathbf{x}^{(1)} \otimes \mathbf{x}^{(2)}] = \sum_{i=1}^r w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i$$

$$\mathbf{M}_3 := \mathbb{E}[\mathbf{x}^{(1)} \otimes \mathbf{x}^{(2)} \otimes \mathbf{x}^{(3)}] = \sum_{i=1}^r w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i$$

- $\mathbf{M}_2, \mathbf{M}_3$ can be reliably estimated when we have many samples
- recovering $\{\boldsymbol{\mu}_i\}$ and $\{w_i\}$ from $\mathbf{M}_2, \mathbf{M}_3$
 \iff tensor decomposition

Latent Dirichlet allocation (LDA)

More complicated topic models: **mixed membership models**, where each data might belong to multiple latent classes simultaneously

This means: the latent variable h is no longer an indicator of topics, but rather, a topic mixture $h \in \Delta_r$

Latent Dirichlet allocation (LDA)

- n : the number of distinct words in the vocabulary
- r : the number of distinct topics
- topic i has word distribution $\boldsymbol{\mu}_i \in \Delta_n$ ($1 \leq i \leq n$)
- each time, draw 3 words as follows
 - draw $\underbrace{\text{topic mixture } \mathbf{h}}_{\text{latent variables}} \in \Delta_r$ according to Dirichlet distribution

$$p_{\boldsymbol{\alpha}}(\mathbf{h}) = \frac{\Gamma(\alpha_0)}{\prod_{i=1}^r \Gamma(\alpha_i)} \prod_{i=1}^r h_i^{\alpha_i - 1}$$

- draw $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)} \in \mathbb{R}^n$ independently according to the *mixed distribution* $\sum_{i=1}^r h_i \boldsymbol{\mu}_i$

Moment method for latent Dirichlet allocation

$$\mathbf{M}_1 := \mathbb{E}[\mathbf{x}^{(1)}]$$

$$\mathbf{M}_2 := \mathbb{E}[\mathbf{x}^{(1)} \otimes \mathbf{x}^{(2)}] - \frac{\alpha_0}{\alpha_0 + 1} \mathbf{M}_1 \otimes \mathbf{M}_1 = \sum_{i=1}^r \frac{\alpha_i}{(\alpha_0 + 1)\alpha_0} \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i$$

$$\begin{aligned} \mathbf{M}_3 &:= \mathbb{E}[\mathbf{x}^{(1)} \otimes \mathbf{x}^{(2)} \otimes \mathbf{x}^{(3)}] - \frac{\alpha_0}{\alpha_0 + 2} \\ &\quad \cdot \left(\mathbb{E}[\mathbf{x}^{(1)} \otimes \mathbf{x}^{(2)} \otimes \mathbf{M}_1] + \mathbb{E}[\mathbf{x}^{(1)} \otimes \mathbf{M}_1 \otimes \mathbf{x}^{(2)}] + \mathbb{E}[\mathbf{M}_1 \otimes \mathbf{x}^{(1)} \otimes \mathbf{x}^{(2)}] \right) \\ &\quad + \frac{2\alpha_0^2}{(\alpha_0 + 2)(\alpha_0 + 1)} \mathbf{M}_1 \otimes \mathbf{M}_1 \otimes \mathbf{M}_1 \\ &= \sum_{i=1}^r \frac{2\alpha_i}{(\alpha_0 + 2)(\alpha_0 + 1)\alpha_0} \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \end{aligned}$$

- estimate \mathbf{M}_1 , \mathbf{M}_2 , \mathbf{M}_3 from samples (assuming α_0 is known)
- recover $\{\boldsymbol{\mu}_i\}$ and $\{\alpha_i\}_{i \geq 1}$ from \mathbf{M}_2 , \mathbf{M}_3 (tensor decomposition)

Gaussian mixture model

- r Gaussian distributions $\mathcal{N}(\boldsymbol{\mu}_i, \sigma^2 \mathbf{I}_n)$ ($1 \leq i \leq r$)
- a sample $\mathbf{x} \in \mathbb{R}^n$ is drawn as follows
 - the latent indicator variable h is generated according to distribution $[w_1, \dots, w_r] \in \Delta_r$ s.t.

$$\mathbb{P}(h = i) = w_i, \quad 1 \leq i \leq r$$

- generate \mathbf{x} from $\mathcal{N}(\boldsymbol{\mu}_h, \sigma^2 \mathbf{I}_n)$

Moment method for Gaussian mixture model

$$\mathbf{M}_2 := \mathbb{E}[\mathbf{x} \otimes \mathbf{x}] - \sigma^2 \mathbf{I} = \sum_{i=1}^r w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i$$

$$\begin{aligned} \mathbf{M}_3 &:= \mathbb{E}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] \\ &\quad - \sigma^2 \sum_{i=1}^n (\mathbb{E}[\mathbf{x}] \otimes \mathbf{e}_i \otimes \mathbf{e}_i + \mathbf{e}_i \otimes \mathbb{E}[\mathbf{x}] \otimes \mathbf{e}_i + \mathbf{e}_i \otimes \mathbf{e}_i \otimes \mathbb{E}[\mathbf{x}]) \\ &= \sum_{i=1}^r w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \end{aligned}$$

- \mathbf{M}_2 , \mathbf{M}_3 and $\mathbb{E}[\mathbf{x}]$ can all be reliably estimated when there are many samples
- recover $\{\boldsymbol{\mu}_i\}$ and $\{w_i\}$ from \mathbf{M}_2 , \mathbf{M}_3 (tensor decomposition)

Tensor power method

Main task

Given

$$M_2 = \sum_{i=1}^r \lambda_i \mathbf{u}_i \otimes \mathbf{u}_i$$
$$M_3 = \sum_{i=1}^r \lambda_i \mathbf{u}_i \otimes \mathbf{u}_i \otimes \mathbf{u}_i$$

where $\lambda_i > 0$

Question: can we recover $\{\lambda_i\}$ and $\{\mathbf{u}_i\}$ from M_2 and M_3 ?

An easier case: orthogonal decomposition

Given

$$M_2 = \sum_{i=1}^r \lambda_i \mathbf{u}_i \otimes \mathbf{u}_i$$
$$M_3 = \sum_{i=1}^r \lambda_i \mathbf{u}_i \otimes \mathbf{u}_i \otimes \mathbf{u}_i$$

where $\lambda_i > 0$, $r \leq n$, and $\{\mathbf{u}_i\}$ are orthonormal

Question: can we recover $\{\lambda_i\}$ and $\{\mathbf{u}_i\}$ from M_2 and M_3 ?

Tensor power method

Define

$$\mathbf{T}(\mathbf{I}, \mathbf{x}, \dots, \mathbf{x}) := \sum_{i=1}^r \lambda_i (\mathbf{u}_i^\top \mathbf{x})^{d-1} \mathbf{u}_i$$

- if $d = 2$ (matrix case): $\mathbf{T}(\mathbf{I}, \mathbf{x}) = \mathbf{T}\mathbf{x}$

Algorithm 5.1 Tensor power method

- 1: **initialize** $\mathbf{x}_0 \leftarrow$ random unit vector
 - 2: **for** $t = 1, 2, \dots$ **do**
 - 3: $\mathbf{x}_t = \mathbf{T}(\mathbf{I}, \mathbf{x}_{t-1}, \dots, \mathbf{x}_{t-1})$ (power iteration)
 - 4: $\mathbf{x}_t \leftarrow \frac{1}{\|\mathbf{x}_t\|_2} \mathbf{x}_t$ (re-normalization)
-

Convergence analysis

Theorem 5.1 (Convergence of tensor power method)

Suppose $\{\mathbf{u}_i\}$ are orthonormal, $\lambda_i > 0$ ($1 \leq i \leq r$), $r \leq n$, and $d = 3$. Then for any $1 \leq i \leq r$,

$$1 - \frac{(\mathbf{u}_i^\top \mathbf{x}_t)^2}{\|\mathbf{x}_t\|_2^2} \leq \lambda_i^2 \sum_{j:j \neq i} \lambda_j^{-2} \left(\frac{\lambda_j \mathbf{u}_j^\top \mathbf{x}_0}{\lambda_i \mathbf{u}_i^\top \mathbf{x}_0} \right)^{2^{t+1}}$$

- tensor power method converges quadratically to some \mathbf{u}_i
- it converges to a point \mathbf{u}_i associated with the largest $\lambda_i \mathbf{u}_i^\top \mathbf{x}_0$
 - both the eigenvalue and the initial point matter!

Proof of Theorem 5.1

Note that removing “re-normalization” steps does not affect $\frac{(\mathbf{u}_i^\top \mathbf{x}_t)^2}{\|\mathbf{x}_t\|_2^2}$ at all. For simplicity, we assume

$$\mathbf{x}_t = \mathbf{T}(\mathbf{I}, \mathbf{x}_{t-1}, \mathbf{x}_{t-1}) = \sum_{i=1}^r \lambda_i (\mathbf{u}_i^\top \mathbf{x}_{t-1})^2 \mathbf{u}_i$$

Observe that

- since $\mathbf{x}_1 = \sum_{i=1}^r \lambda_i (\mathbf{u}_i^\top \mathbf{x}_0)^2 \mathbf{u}_i$, we have

$$(\mathbf{u}_i^\top \mathbf{x}_1)^2 = \lambda_i^2 (\mathbf{u}_i^\top \mathbf{x}_0)^4$$

- since $\mathbf{x}_2 = \sum_{i=1}^r \lambda_i (\mathbf{u}_i^\top \mathbf{x}_1)^2 \mathbf{u}_i$, we have

$$(\mathbf{u}_i^\top \mathbf{x}_2)^2 = \lambda_i^2 (\mathbf{u}_i^\top \mathbf{x}_1)^4 = \lambda_i^6 (\mathbf{u}_i^\top \mathbf{x}_0)^8$$

- since $\mathbf{x}_3 = \sum_{i=1}^r \lambda_i (\mathbf{u}_i^\top \mathbf{x}_2)^2 \mathbf{u}_i$, we have

$$(\mathbf{u}_i^\top \mathbf{x}_3)^2 = \lambda_i^2 (\mathbf{u}_i^\top \mathbf{x}_2)^4 = \lambda_i^{14} (\mathbf{u}_i^\top \mathbf{x}_0)^{16}$$

Proof of Theorem 5.1 (cont.)

By induction, one has

$$(\mathbf{u}_i^\top \mathbf{x}_t)^2 = \lambda_i^{2^{t+1}-2} (\mathbf{u}_i^\top \mathbf{x}_0)^{2^{t+1}}, \quad 1 \leq i \leq r$$

This implies

$$\frac{(\mathbf{u}_i^\top \mathbf{x}_t)^2}{\|\mathbf{x}_t\|_2^2} = \frac{(\mathbf{u}_i^\top \mathbf{x}_t)^2}{\sum_{j=1}^r (\mathbf{u}_j^\top \mathbf{x}_t)^2} = \frac{(\lambda_i \mathbf{u}_i^\top \mathbf{x}_0)^{2^{t+1}}}{\sum_{j=1}^r \left(\frac{\lambda_i}{\lambda_j}\right)^2 (\lambda_j \mathbf{u}_j^\top \mathbf{x}_0)^{2^{t+1}}}$$

and hence

$$\begin{aligned} 1 - \frac{(\mathbf{u}_i^\top \mathbf{x}_t)^2}{\|\mathbf{x}_t\|_2^2} &= \frac{\sum_{j:j \neq i} \left(\frac{\lambda_i}{\lambda_j}\right)^2 (\lambda_j \mathbf{u}_j^\top \mathbf{x}_0)^{2^{t+1}}}{\sum_j \left(\frac{\lambda_i}{\lambda_j}\right)^2 (\lambda_j \mathbf{u}_j^\top \mathbf{x}_0)^{2^{t+1}}} \\ &\leq \frac{\sum_{j:j \neq i} \left(\frac{\lambda_i}{\lambda_j}\right)^2 (\lambda_j \mathbf{u}_j^\top \mathbf{x}_0)^{2^{t+1}}}{(\lambda_i \mathbf{u}_i^\top \mathbf{x}_0)^{2^{t+1}}} \\ &= \lambda_i^2 \sum_{j:j \neq i} \lambda_j^{-2} \left(\frac{\lambda_j \mathbf{u}_j^\top \mathbf{x}_0}{\lambda_i \mathbf{u}_i^\top \mathbf{x}_0} \right)^{2^{t+1}} \end{aligned}$$

General case: reduction to orthogonally decomposable tensors

Suppose $r \leq n$, but $\{\mathbf{u}_i\}$ are not orthonormal

Key idea: use M_2 to find a “whitening matrix” that allows us to orthogonalize $\{\mathbf{u}_i\}$

General case: reduction to orthogonally decomposable tensor

Let \mathbf{W} be a whitening matrix (e.g. $\mathbf{W} = \mathbf{U}\mathbf{\Lambda}^{-1/2}$) obeying

$$\mathbf{W}^\top \mathbf{M}_2 \mathbf{W} = \mathbf{I} \quad (5.1)$$

Then

$$\begin{aligned} M_3(\mathbf{W}, \mathbf{W}, \mathbf{W}) &= \sum_{i=1}^r \lambda_i (\mathbf{W}^\top \mathbf{u}_i) \otimes (\mathbf{W}^\top \mathbf{u}_i) \otimes (\mathbf{W}^\top \mathbf{u}_i) \\ &= \sum_{i=1}^r \lambda_i \tilde{\mathbf{u}}_i \otimes \tilde{\mathbf{u}}_i \otimes \tilde{\mathbf{u}}_i \end{aligned}$$

where $\{\tilde{\mathbf{u}}_i\}$ become **orthonormal vectors**

- use the tensor power method to recover $\{\tilde{\mathbf{u}}_i\}$

Reference

- “*Tensor decompositions for learning latent variable models*,” A. Anandkumar, R. Ge, D. Hsu, S. Kakade, *Journal of machine learning research*, 2014.
- “*Tensor decompositions and applications*,” T. Kolda, B. Bader, *SIAM review*, 2009.
- “*Orthogonal tensor decompositions*,” T. Kolda, *SIAM journal on matrix analysis and applications*, 2001.
- “*Spectral learning on matrices and tensors*,” M. Janzamin, R. Ge, J. Kossaifi and A. Anandkumar, *Foundations and Trends in Machine Learning*, 2019.
- “*On the best rank-1 and rank- (R_1, R_2, \dots, R_n) approximation of higher-order tensors*,” L. De Lathauwer, B. De Moor, J. Vandewalle, *SIAM journal on matrix analysis and applications*, 2000.

Reference

- “*Spectral algorithms for tensor completion*,” A. Montanari, N. Sun, *Communications on pure and applied mathematics*, 2018.
- “*Subspace estimation from unbalanced and incomplete data matrices: $\ell_{2,\infty}$ statistical guarantees*,” C. Cai, G. Li, Y. Chi, H. V. Poor, Y. Chen, *Annals of Statistics*, 2020.