

Optimization for Tensor Models

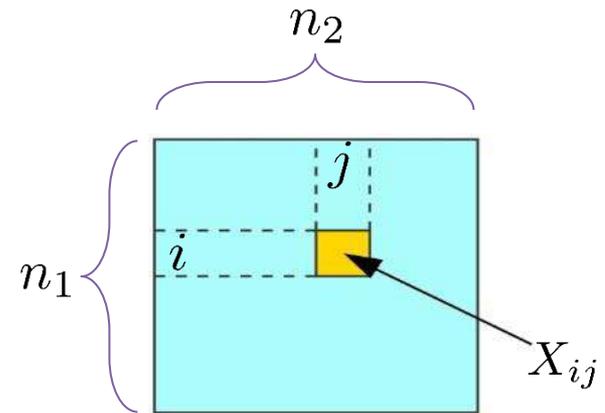
Donald Goldfarb
IEOR Department
Columbia University

UCLA Mathematics Department
Distinguished Lecture Series
May 17 – 19, 2016

Tensors

Matrix

$$\mathbf{X} = (X_{ij}) \in \mathbb{R}^{n_1 \times n_2}$$



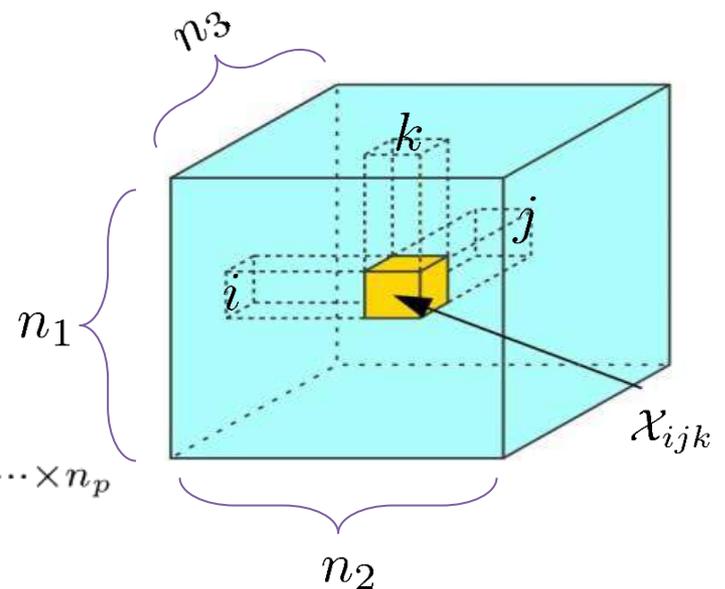
Tensor: higher-order matrix

three-way tensor:

$$\mathcal{X} = (\mathcal{X}_{ijk}) \in \mathbb{R}^{n_1 \times n_2 \times n_3}$$

p-way tensor:

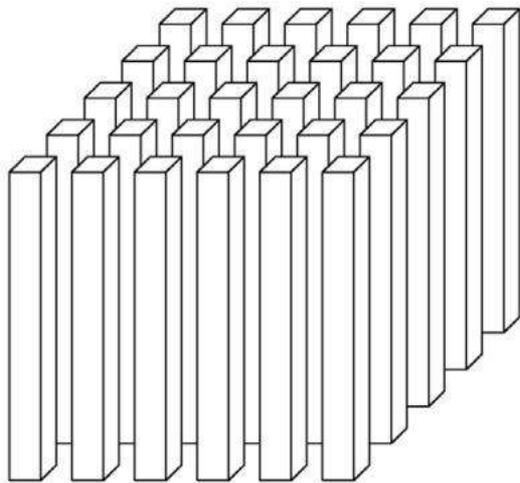
$$\mathcal{X} \in \bigotimes_{i=1}^p \mathbb{R}^{n_i} := (\mathcal{X}_{i_1 i_2 \dots i_k}) \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_p}$$



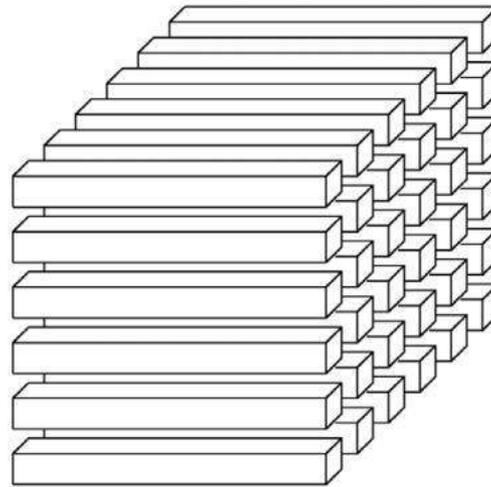
Tensors

Tensor == Vector ? **Fibers?**

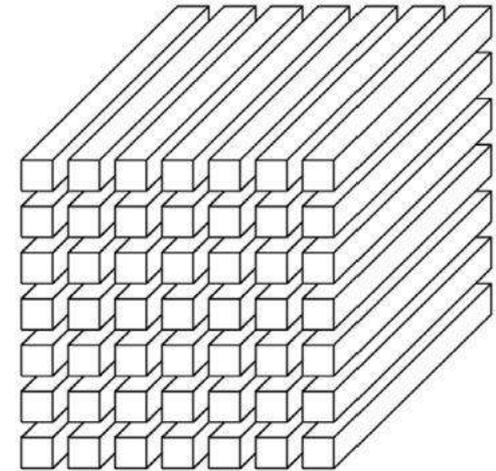
can be viewed as a collection of fibers



(a) Mode-1 (column) fibers: $\mathbf{x}_{:jk}$



(b) Mode-2 (row) fibers: $\mathbf{x}_{i:k}$



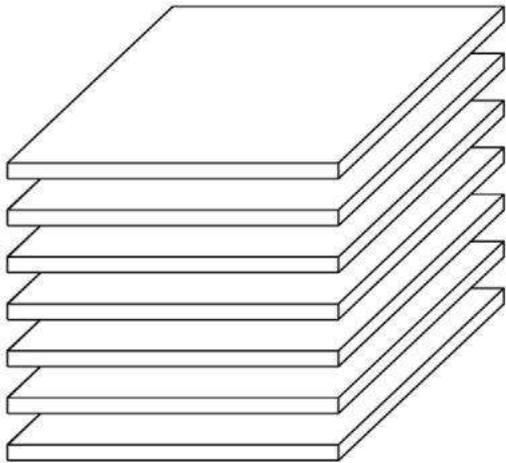
(c) Mode-3 (tube) fibers: $\mathbf{x}_{ij:}$

(Kolda & Bader, 2009)

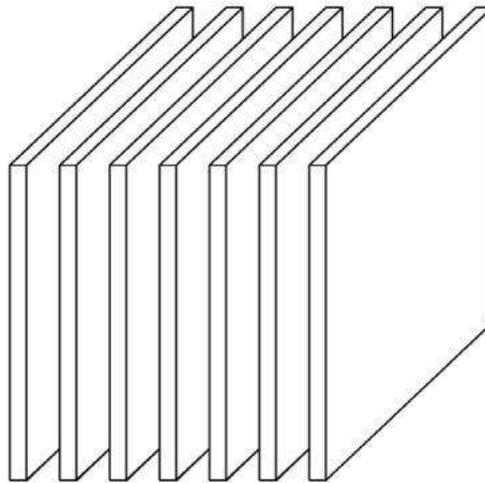
Tensors

Tensor == Matrix ? Slices?

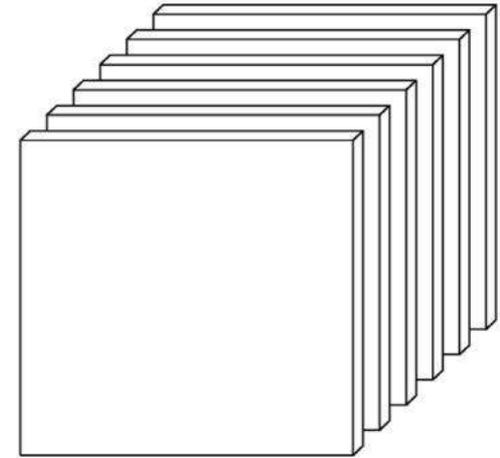
can be viewed as a collection of slices



(a) Horizontal slices: $\mathbf{X}_{i::}$



(b) Lateral slices: $\mathbf{X}_{:,j:}$



(c) Frontal slices: $\mathbf{X}_{::k}$ (or \mathbf{X}_k)

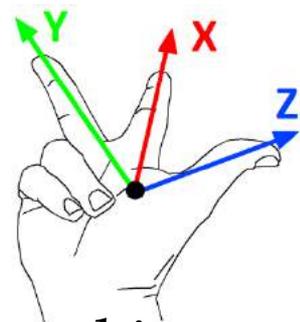
(Kolda & Bader, 2009)

Tensors

Why tensors?

- tensors capture multilinear structure
- more **flexible** and **powerful** models

e.g. parameter estimation in latent variable modelling
(to be discussed shortly)



Tensor: object in its own right

- its own geometrical, statistical and computational issues
- much harder to work with than a matrix

Example: Single Topic Model

Model:

- k topics (dists. over d words)

$$\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^d$$

- sample topic h :

$$\text{prob}[h = i] = w_i \quad (i \in [k])$$

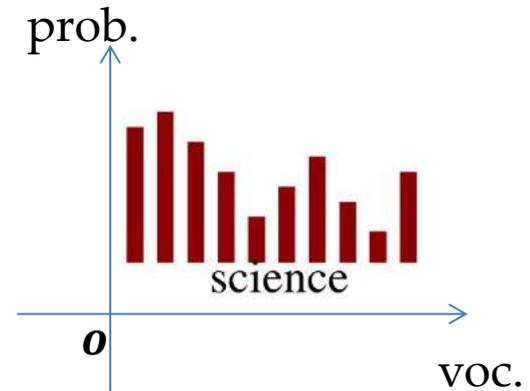
- each document has m words

$$x_1, x_2, \dots, x_m \text{ sampled i.i.d. from } \mu_i$$

Dataset: m -word documents

Goal: learn parameters

$$\theta = \left(\{ \mu_i \}_{i \in [k]}, \{ w_i \}_{i \in [k]} \right)$$



Example: Single Topic Model

Method of moments:



Karl Pearson (1857~1936)

Key idea: find parameters (approx.) consistent with observed moments (Pearson, 1894)

Procedure:

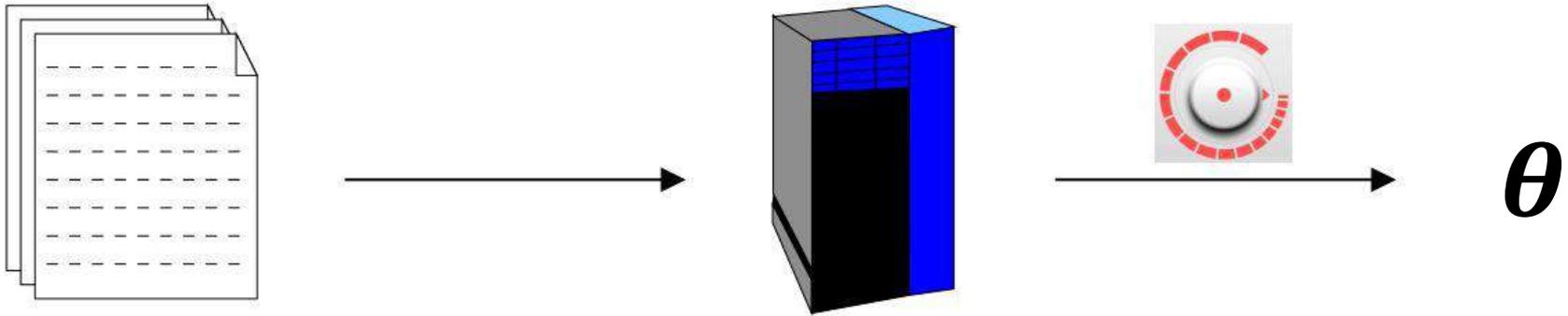
- setup equations

$$\mathbf{Moment}_{\text{population}}(\boldsymbol{\theta}) = \widehat{\mathbf{Moment}}_{\text{sample}}$$

- solve approximately

Example: Single Topic Model

Goal: model the topics of the documents in a given corpus



Samples of documents
generated based on single
topic model with params
 $\theta = (\{\mu_i\}, \{w_i\})$

Unsupervised learning alg.
Method of moments

- setup equations
- solve them approx.

Model parameters

Which moments to use?

1st -order? 2nd -order? 3rd -order? pth -order?

Example: Single Topic Model

Binary encoding:

$\mathbf{x}_t = \mathbf{e}_i \iff$ the t -th word in the document is i -th word in the vocabulary $(t, i) \in [m] \times [d]$

E.g.

topic: animal

3-word document

7-word vocabulary

vocab.	x_1	x_2	x_3
cats	0	0	1
dogs	0	0	0
fear	0	0	0
I	1	0	0
like	0	1	0
raise	0	0	0
want	0	0	0

Example: Single Topic Model

Binary encoding:

$\mathbf{x}_t = \mathbf{e}_i \iff$ the t -th word in the document is i -th word in the vocabulary
 $(t, i) \in [m] \times [d]$

First moment

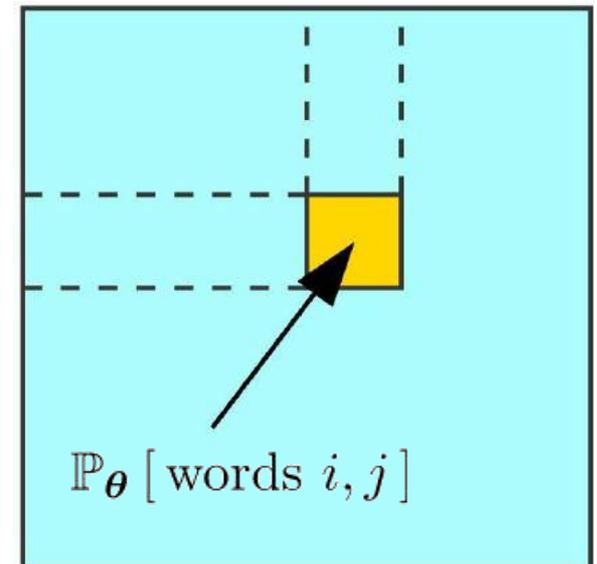
$$\begin{aligned}\mathbb{E}[\mathbf{x}_1] &= \left(\mathbb{P}_{\theta}[\mathbf{x}_1 = \mathbf{e}_i] \right)_{i \in [d]} \\ &= \mathbb{E}_h[\mathbb{E}[\mathbf{x}_1 \mid \text{topic } h]] \\ &= \mathbb{E}_h[\boldsymbol{\mu}_h] \\ &= \sum_{i \in [k]} w_i \boldsymbol{\mu}_i \in \mathbb{R}^d\end{aligned}$$

Not identifiable: only d numbers for $(d + 1)k$ parameters.

Example: Single Topic Model

Second moment

$$\begin{aligned} \mathbf{M} &= \mathbb{E} [\mathbf{x}_1 \otimes \mathbf{x}_2] \\ &= \left(\mathbb{P}_{\boldsymbol{\theta}} [(\mathbf{x}_1, \mathbf{x}_2) = (i, j)] \right)_{(i,j) \in [d] \times [d]} \\ &= \mathbb{E}_h [\mathbb{E} [\mathbf{x}_1 \otimes \mathbf{x}_2 \mid \text{topic } h]] \\ &= \mathbb{E}_h [\boldsymbol{\mu}_h \otimes \boldsymbol{\mu}_h \mid \text{topic } h] \\ &= \sum_{i \in [k]} w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \in \mathbb{R}^{d \times d} \end{aligned}$$



Matrix-mode: still not identifiable even though $\frac{d(d+1)}{2} > (d+1)k$. (Why?)

Example: Single Topic Model

Identifiable? **NO!**

$U \in \mathbb{R}^{n \times k}$ is a solution

$\Leftrightarrow M = UU^T \in \mathbb{R}^{n \times n}$

$\Leftrightarrow M = UQ(UQ)^T$, for any $Q \in \mathcal{O}_k$ ($QQ^T = Q^TQ = I_k$)

$\Leftrightarrow U \in \mathbb{R}^{n \times k}$ is a solution for any $Q \in \mathcal{O}(k)$

\Leftrightarrow AMBIGUITY!

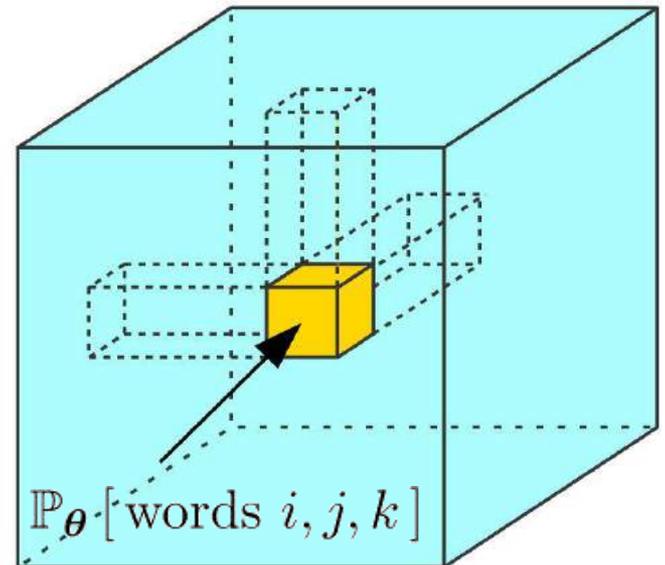
Matrix mode: **INSUFFICIENT** to identify parameters.!

Example: Single Topic Model

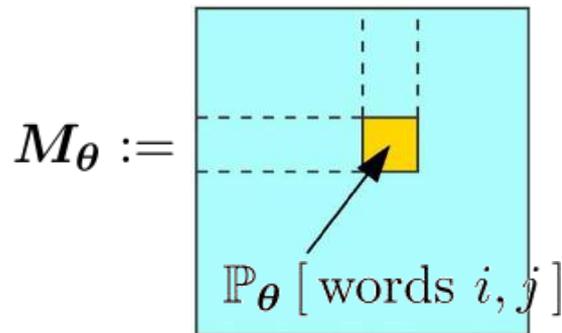
Third-order moment

$$\begin{aligned}\mathcal{T} &= \mathbb{E} [\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3] \\ &= \left(\mathbb{P}_{\boldsymbol{\theta}} [(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) = (i, j, k)] \right)_{(i,j,k) \in [d] \times [d] \times [d]} \\ &= \mathbb{E}_h [\mathbb{E} [\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3 \mid \text{topic } h]] \\ &= \mathbb{E}_h [\boldsymbol{\mu}_h \otimes \boldsymbol{\mu}_h \otimes \boldsymbol{\mu}_h \mid \text{topic } h] \\ &= \sum_{i \in [k]} w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \in \mathbb{R}^{d \times d \times d}\end{aligned}$$

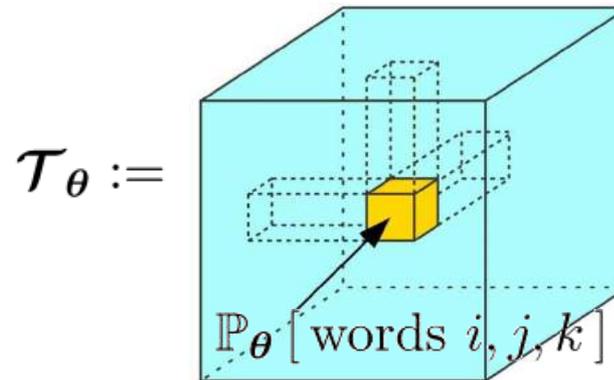
Tensor mode



Example: Single Topic Model



$$M_\theta = \sum_{i \in [k]} w_i \mu_i \otimes \mu_i$$



$$\mathcal{T}_\theta = \sum_{i \in [k]} w_i \mu_i \otimes \mu_i \otimes \mu_i$$

Claim: M_θ & \mathcal{T}_θ uniquely determine the parameters θ

Example: Single Topic Model

Reduction to orthogonal case via whitening

Whiten:

project to k dimensions

transform to orthogonality

$$\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{U}^T$$

$$\mathbf{W} = \mathbf{U}\mathbf{D}^{-1/2}$$

Reduced Eigen-
Decomp.

$$\mathbf{M} = \sum_{i \in [k]} w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i$$

$$\mathcal{T} = \sum_{i \in [k]} w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i$$

apply \mathbf{W} to \mathbf{M} and \mathcal{T}



$$\mathbf{v}_i = \sqrt{w_i} \mathbf{W}^T \boldsymbol{\mu}_i \in \mathbb{R}^k$$

$$\lambda_i = 1/\sqrt{w_i}$$

$$\tilde{\mathbf{M}} = \sum_{i \in [k]} \mathbf{v}_i \otimes \mathbf{v}_i = \mathbf{I}_k$$

$$\tilde{\mathcal{T}} = \sum_{i \in [k]} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i$$

$\{\mathbf{v}_i\}_{i \in [k]}$ forms an orthonormal basis for \mathbb{R}^k

Example: Single Topic Model

Spectral theorem and eigen-decompositions

$\{\mathbf{v}_i\}_{i \in [k]}$ forms an orthonormal basis for \mathbb{R}^k

symmetric matrix

$$\mathbf{X} = \sum_{i \in [k]} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i$$

eigen-decomp. is unique iff
 $\lambda_i \neq \lambda_j$

symmetric tensor

$$\mathcal{X} = \sum_{i \in [k]} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i$$

if such decomp. exists,
then it is always unique
(even if λ_i 's all same)

Uniqueness of orthogonal decomp. \rightarrow \mathbf{M}_θ & \mathcal{T}_θ uniquely determine
the parameters θ

Identifiability issue is resolved via tensor mode!

2nd Example: Mixture of Spherical Gaussians

Model:

k means $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^d$

Sample cluster $h = i$ with probability w_i ($i \in [k]$)

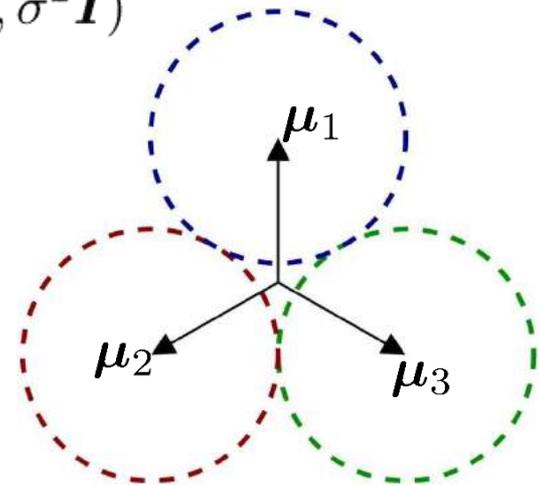
Observe \mathbf{x} , with i.i.d. homogeneous spherical noise

$$\mathbf{x} = \mu_i + \eta, \quad \eta \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

Dataset: multiple points

Goal: learn parameters

$$\theta = \left(\{\mu_i\}_{i \in [k]}, \{w_i\}_{i \in [k]}, \sigma \right)$$



2nd Example: Mixture of Spherical Gaussians

Identifiable using 1st, 2nd and 3rd –order moments together

[Hsu & Kakade, '13]

$$\sigma^2 = \lambda_{\min}(\mathbb{E}[\mathbf{x} \otimes \mathbf{x}])$$

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

$$\begin{aligned} \mathcal{T} &= \mathbb{E}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] \\ &\quad - \sigma^2 \sum_{i \in [d]} (\mathbb{E}[\mathbf{x}] \otimes \mathbf{e}_i \otimes \mathbf{e}_i + \mathbf{e}_i \otimes \mathbb{E}[\mathbf{x}] \otimes \mathbf{e}_i + \mathbf{e}_i \otimes \mathbf{e}_i \otimes \mathbb{E}[\mathbf{x}]) \\ &= \sum_{i \in [k]} w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \end{aligned}$$

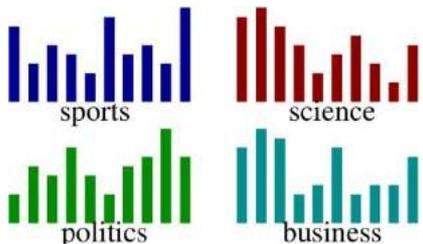
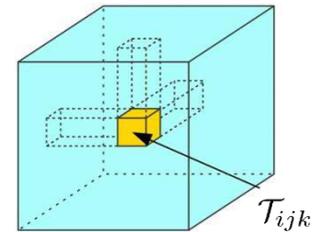
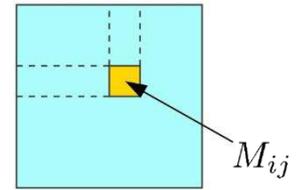
Same approach as simple topic model!

General Principle

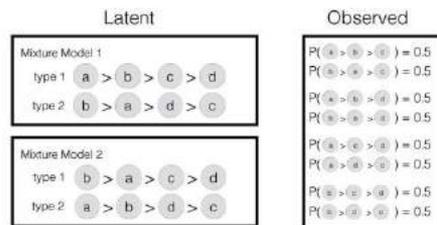
Similar structures prevail in latent variable models

$$\mathbf{M} = f\left(\leq 2\text{nd-order moments}\right) = \sum_{i \in [k]} w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i$$

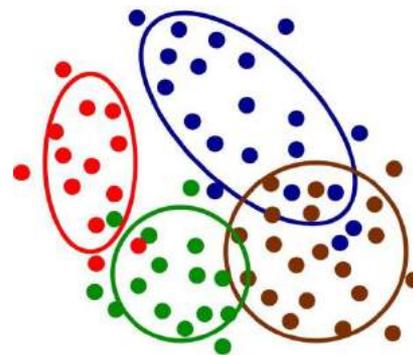
$$\mathcal{T} = g\left(\leq 3\text{rd-order moments}\right) = \sum_{i \in [k]} w_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i$$



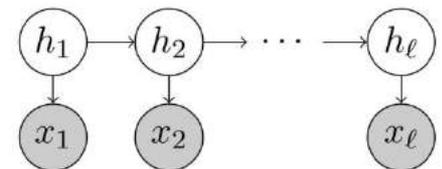
Latent Dirichlet Allocation (LDA)



Mixed Multinomial Logit Model



Mixture of Gaussians (MoG)



Hidden Markov Models (HMMs)

Orthogonal Decomposition

How to find $\{(\lambda_i, \mathbf{v}_i)\}_{i \in [k]}$?

$$(\langle \mathbf{v}_i, \mathbf{v}_j \rangle = \delta_{ij})$$

symmetric matrix

$$\mathbf{X} = \sum_{i \in [k]} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i$$

successive rank-one
approximation (SROA)

symmetric tensor

$$\mathcal{X} = \sum_{i \in [k]} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i$$

generalized SROA?

Successive Rank-One Approximation

Symmetric matrix

$$M = \sum_{i \in [k]} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \quad (\langle \mathbf{v}_i, \mathbf{v}_j \rangle = \delta_{ij})$$

SROA: rank-one approximation + deflation

Repeat k times

$$(\hat{\lambda}, \hat{\mathbf{v}}) \in \arg \min_{\lambda \in \mathbb{R}, \|\mathbf{v}\|=1} \|M - \lambda \mathbf{v} \otimes \mathbf{v}\| \quad \text{(rank-one approx.)}$$

$$M \leftarrow M - \hat{\lambda} \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \quad \text{(deflation)}$$

Recover $\{(\lambda_i, \mathbf{v}_i)\}_{i \in [k]}$ exactly ($\lambda_i \neq \lambda_j$)

Successive Rank-One Approximation

Symmetric orthogonal decomposable (**SOD**) tensor

$$\mathcal{T} = \sum_{i \in [k]} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i \quad (\langle \mathbf{v}_i, \mathbf{v}_j \rangle = \delta_{ij})$$

SROA: **rank-one approximation + deflation**

Repeat k times

$$(\hat{\lambda}, \hat{\mathbf{v}}) \in \arg \min_{\lambda \in \mathbb{R}, \|\mathbf{v}\|=1} \|\mathcal{T} - \lambda \mathbf{v} \otimes \mathbf{v} \otimes \mathbf{v}\|_F \quad \text{(rank-one approx.)}$$

$$\mathcal{T} \leftarrow \mathcal{T} - \hat{\lambda} \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \quad \text{(deflation)}$$

Recover $\{(\lambda_i, \mathbf{v}_i)\}_{i \in [k]}$ exactly up to sign flips [Zhang & Golub, 01]

N.B.: no requirement on λ 's to be distinct

Successive Rank-One Approximation

Sampling error

$$\# \text{ samples} = \infty \quad \rightarrow \quad \hat{\mathcal{T}} = \sum_{i \in [k]} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i$$

$$\# \text{ samples} < \infty \quad \rightarrow \quad \hat{\mathcal{T}} = \sum_{i \in [k]} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i + \boldsymbol{\varepsilon}$$

Other potential sources of perturbation $\boldsymbol{\varepsilon}$

- Model misspecification
- Numerical error
-

Is SROA **robust** to the perturbation?

Recall **matrix perturbation theory** (e.g. Davis-Kahan),
which requires $\| \text{perturbation matrix} \| < \min_{i \neq j} |\lambda_i - \lambda_j|$.

Successive Rank-One Approximation

Perturbed **SOD** tensor

$$\hat{\mathcal{T}} = \sum_{i \in [k]} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i + \boldsymbol{\varepsilon} \quad (\langle \mathbf{v}_i, \mathbf{v}_j \rangle = \delta_{ij})$$

SROA: **rank-one approximation + deflation**

Repeat k times

$$(\hat{\lambda}, \hat{\mathbf{v}}) \in \arg \min_{\lambda \in \mathbb{R}, \|\mathbf{v}\|=1} \left\| \hat{\mathcal{T}} - \lambda \mathbf{v} \otimes \mathbf{v} \otimes \mathbf{v} \right\|_F \quad \text{(rank-one approx.)}$$

$$\hat{\mathcal{T}} \leftarrow \hat{\mathcal{T}} - \hat{\lambda} \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \quad \text{(deflation)}$$

Successive Rank-One Approximation

Input: the perturbed SOD tensor

Repeat k times

$$(\hat{\lambda}, \hat{\mathbf{v}}) \in \arg \min_{\lambda \in \mathbb{R}, \|\mathbf{v}\|=1} \left\| \hat{\mathcal{T}} - \lambda \mathbf{v} \otimes \mathbf{v} \otimes \mathbf{v} \right\|_F \quad \text{(rank-one approx.)}$$

$$\hat{\mathcal{T}} \leftarrow \hat{\mathcal{T}} - \hat{\lambda} \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \quad \text{(deflation)}$$

Output: $\left\{ (\hat{\lambda}_i, \hat{\mathbf{v}}_i) \right\}_{i \in [k]}$

Theorem (MHG, '15)

Output $\left\{ (\hat{\lambda}_i, \hat{\mathbf{v}}_i) \right\}_{i \in [k]}$

$$\|\mathcal{E}\| \leq C \cdot \frac{\lambda_{\min}}{\sqrt{n}}$$



\exists perm. π on $[k]$ s.t. $\forall i \in [k]$

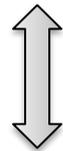
$$\begin{cases} \min_{\pm} \left\{ \left| \lambda_{\pi(i)} \pm \hat{\lambda}_i \right| \right\} \leq 2\varepsilon \\ \min_{\pm} \left\{ \left\| \mathbf{v}_{\pi(i)} \pm \hat{\mathbf{v}}_i \right\| \right\} \leq 20\varepsilon / \lambda_{\pi(i)} \end{cases}$$

N.B.

- generalizes matrix perturbation analysis
- **NO** spectral gap quantity involved

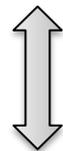
Best Rank-One Tensor Approximation

$$\text{minimize } \|\mathcal{T} - \lambda \underbrace{\mathbf{v} \otimes \mathbf{v} \otimes \cdots \otimes \mathbf{v}}_{p \text{ times}}\|_F^2 \quad \text{subject to } \|\mathbf{v}\| = 1, \quad \lambda \in \mathbb{R}$$



$$\|\mathcal{T} - \lambda \mathbf{v}^{\otimes p}\|_F^2 = \|\mathcal{T}\|_F^2 - 2\lambda \langle \mathcal{T}, \mathbf{v}^{\otimes p} \rangle + \lambda^2$$

$$\min_{\|\mathbf{v}\|=1} \min_{\lambda \in \mathbb{R}} \left(-2\lambda \langle \mathcal{T}, \mathbf{v}^{\otimes p} \rangle + \lambda^2 \right)$$



eliminate λ , then **solve**

$$\min_{\|\mathbf{v}\|=1} \langle \mathcal{T}, \mathbf{v}^{\otimes p} \rangle \quad \& \quad \min_{\|\mathbf{v}\|=1} \langle -\mathcal{T}, \mathbf{v}^{\otimes p} \rangle$$

SDP Relaxation [Jiang, Ma & Zhang, '14]

$$\text{minimize } \langle \mathcal{T}, \mathbf{v} \otimes \mathbf{v} \otimes \mathbf{v} \otimes \mathbf{v} \rangle \quad \text{subject to } \|\mathbf{v}\| = 1 \quad (*)$$

$$\Updownarrow \text{ reshape } \mathcal{T} \in \otimes^4 \mathbb{R}^n \text{ to } \mathbf{T}_{\square} \in \mathbb{R}^{n^2 \times n^2}$$

$$\text{minimize } \left\langle \mathbf{T}_{\square}, \text{vec}(\mathbf{v} \otimes \mathbf{v}) \text{vec}(\mathbf{v} \otimes \mathbf{v})^{\top} \right\rangle \quad \text{subject to } \|\mathbf{v}\| = 1$$

$$\Updownarrow \mathbf{X} = \text{vec}(\mathbf{v} \otimes \mathbf{v}) \text{vec}(\mathbf{v} \otimes \mathbf{v})^{\top} \quad \Downarrow \text{ induced symmetry } \& \|\mathbf{v}\| = 1$$

$$\text{minimize } \langle \mathbf{T}_{\square}, \mathbf{X} \rangle \quad \text{subject to } \mathcal{A}(\mathbf{X}) = \mathbf{b}, \mathbf{X} \succeq \mathbf{0}, \text{rank}(\mathbf{X}) = 1$$

$$\Updownarrow \text{ convex relaxation}$$

$$\text{minimize } \langle \mathbf{T}_{\square}, \mathbf{X} \rangle \quad \text{subject to } \mathcal{A}(\mathbf{X}) = \mathbf{b}, \mathbf{X} \succeq \mathbf{0}$$

(*) The $p = 3$ tensor problem can be transformed into a $p = 4$ tensor problem [JMZ, '14].

SDP Relaxation

Linear constraints $\mathcal{A}(\mathbf{X}) = \mathbf{b}$

E.g. $\mathbf{v} = (v_1, v_2)^T$

$$\mathbf{X} = \text{vec}(\mathbf{v} \otimes \mathbf{v}) \text{vec}(\mathbf{v} \otimes \mathbf{v})^T = \begin{bmatrix} v_1^2 \\ v_1 v_2 \\ v_2 v_1 \\ v_2^2 \end{bmatrix} [v_1^2 \quad v_1 v_2 \quad v_2 v_1 \quad v_2^2]$$

$$= \begin{bmatrix} v_1^4 & v_1^3 v_2 & v_2 v_1^3 & v_1^2 v_2^2 \\ v_1^3 v_2 & v_1^2 v_2^2 & v_1^2 v_2^2 & v_1 v_2^3 \\ v_2 v_1^3 & v_1^2 v_2^2 & v_1^2 v_2^2 & v_1 v_2^3 \\ v_1^2 v_2^2 & v_1 v_2^3 & v_1 v_2^3 & v_2^4 \end{bmatrix}$$

$$\mathcal{A}(\mathbf{X}) = \mathbf{b} \begin{cases} \text{spherical constr.} & 1 = \|\mathbf{v}\|_2^2 = v_1^2 + v_2^2 \\ & = (v_1^2 + v_2^2)^2 = \text{trace}(\mathbf{X}) \\ \text{hyper-symmetry} & \mathbf{X} = \begin{bmatrix} a & d & d & b \\ d & b & b & e \\ d & b & b & e \\ b & e & e & c \end{bmatrix} \end{cases}$$

SDP Relaxation

minimize $\langle \mathbf{T}_{\square}, \mathbf{X} \rangle$ subject to $\mathcal{A}(\mathbf{X}) = \mathbf{b}, \mathbf{X} \succeq \mathbf{0}$

N.B.

- SDP relaxation proposed by [Jiang, Ma & Zhang, '14]
- Square reshaping trick for low-rank tensor recovery [MHWG, '14]
$$\mathbf{T}_{i+(j-1)n, k+(l-1)n} \leftarrow \mathcal{T}_{ijkl}$$
- Equivalent SDP (of reduced size) proposed by [Nie & Wang, '14] using moment-based convex relaxation
- Empirically, $\text{rank}(\mathbf{X}^*) = 1$ observed **almost always!**
i.e., SDP relaxation \rightarrow solves nonconvex problem!

Solving SDP

Semidefinite programming (SDP)

$$\begin{aligned} \min \quad & \langle \mathbf{C}, \mathbf{X} \rangle \\ \text{s.t.} \quad & \mathcal{A}(\mathbf{X}) = \mathbf{b} \\ & \mathbf{X} \succeq \mathbf{0} \end{aligned}$$

Linear constraints:

$$\mathcal{A}(\mathbf{X}) = \mathbf{b}$$

$$\Updownarrow \mathcal{A}(\mathbf{X}) = (\langle \mathbf{A}_1, \mathbf{X} \rangle, \langle \mathbf{A}_2, \mathbf{X} \rangle, \dots, \langle \mathbf{A}_m, \mathbf{X} \rangle)^\top$$

$$\langle \mathbf{A}_i, \mathbf{X} \rangle = b_i, \quad i \in [m]$$

$$\Updownarrow \mathbf{A} = \left(\text{vec}(\mathbf{A}_1), \text{vec}(\mathbf{A}_2), \dots, \text{vec}(\mathbf{A}_m) \right)^\top$$

$$\mathbf{A} \text{vec}(\mathbf{X}) = \mathbf{b}$$

Deriving the dual problem

Lagrangian function

$$\begin{aligned}\mathcal{L}(\mathbf{X}, \mathbf{y}, \mathbf{S}) &= \langle \mathbf{C}, \mathbf{X} \rangle - \mathbf{y}^\top (\mathcal{A}(\mathbf{X}) - \mathbf{b}) - \langle \mathbf{X}, \mathbf{S} \rangle \\ &= \langle \mathbf{C} - \mathcal{A}^*(\mathbf{y}) - \mathbf{S}, \mathbf{X} \rangle + \mathbf{y}^\top \mathbf{b}\end{aligned}$$

$$\mathcal{A}(\mathbf{X}) = \left(\langle \mathbf{A}_i, \mathbf{X} \rangle \right)_{i \in [m]}$$

$$\mathcal{A}^*(\mathbf{y}) = \sum_{i \in [m]} y_i \mathbf{A}_i$$

Dual problem

$$\max_{\mathbf{y}, \mathbf{S} \succeq \mathbf{0}} \min_{\mathbf{X}} \mathcal{L}(\mathbf{X}, \mathbf{y}, \mathbf{S}) = \langle \mathbf{C} - \mathcal{A}^*(\mathbf{y}) - \mathbf{S}, \mathbf{X} \rangle + \mathbf{y}^\top \mathbf{b}$$

$$\max_{\mathbf{y}, \mathbf{S}} \mathbf{b}^\top \mathbf{y}$$

$$\begin{aligned}\text{s.t. } & \mathcal{A}^*(\mathbf{y}) + \mathbf{S} = \mathbf{C} \\ & \mathbf{S} \succeq \mathbf{0}\end{aligned}$$

$$\min_{\mathbf{y}, \mathbf{S}} -\mathbf{b}^\top \mathbf{y}$$

$$\begin{aligned}\text{s.t. } & \mathcal{A}^*(\mathbf{y}) + \mathbf{S} = \mathbf{C} \\ & \mathbf{S} \succeq \mathbf{0}\end{aligned}$$

Augmented Lagrangian Method

$$(D) \quad \min_{\mathbf{y}, \mathbf{S}} -\mathbf{b}^\top \mathbf{y} \quad \text{s.t.} \quad \mathcal{A}^*(\mathbf{y}) + \mathbf{S} = \mathbf{C}, \mathbf{S} \succeq \mathbf{0}$$

Augmented Lagrangian function

$$\mathcal{L}_\mu(\mathbf{X}, \mathbf{y}, \mathbf{S}) = -\mathbf{y}^\top \mathbf{b} + \langle \mathcal{A}^*(\mathbf{y}) + \mathbf{S} - \mathbf{C}, \mathbf{X} \rangle + \frac{1}{2\mu} \|\mathcal{A}^*(\mathbf{y}) + \mathbf{S} - \mathbf{C}\|_F^2$$

Augmented Lagrangian method (ALM)

***k*-th iteration:**

compute $(\mathbf{y}^{k+1}, \mathbf{S}^{k+1}) \in \arg \min_{\mathbf{y}, \mathbf{S} \succeq \mathbf{0}} \mathcal{L}_\mu(\mathbf{X}^k, \mathbf{y}, \mathbf{S})$

update $\mathbf{X}^{k+1} = \mathbf{X}^k + \left(\mathcal{A}^*(\mathbf{y}^{k+1}) + \mathbf{S}^{k+1} - \mathbf{C} \right) / \mu$

➤ minimizing $\mathcal{L}_\mu(\mathbf{X}^k, \mathbf{y}, \mathbf{S})$ jointly over (\mathbf{y}, \mathbf{S}) is hard!

Alternating Direction Method of Multipliers (ADMM)

Remedy: alternating direction

minimize $\mathcal{L}_\mu(\mathbf{X}^k, \mathbf{y}, \mathbf{S})$ along \mathbf{y} -direction and \mathbf{S} -direction alternatively

ADMM

k -th iteration:

$$\mathbf{y}\text{-update: } \mathbf{y}^{k+1} \leftarrow \arg \min_{\mathbf{y}} \mathcal{L}_\mu(\mathbf{X}^k, \mathbf{y}, \mathbf{S}^k)$$

$$\mathbf{S}\text{-update: } \mathbf{S}^{k+1} \leftarrow \arg \min_{\mathbf{S} \succeq \mathbf{0}} \mathcal{L}_\mu(\mathbf{X}^k, \mathbf{y}^{k+1}, \mathbf{S})$$

$$\mathbf{X}\text{-update: } \mathbf{X}^{k+1} = \mathbf{X}^k + \left(\mathcal{A}^*(\mathbf{y}^{k+1}) + \mathbf{S}^{k+1} - \mathbf{C} \right) / \mu$$

Each step is (relatively) easy to compute!

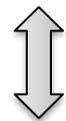
Update \mathbf{y}

$$\mathbf{y}^{k+1} \leftarrow \arg \min_{\mathbf{y}} \mathcal{L}_{\mu}(\mathbf{X}^k, \mathbf{y}, \mathbf{S}^k)$$

$$\mathcal{L}_{\mu}(\mathbf{X}, \mathbf{y}, \mathbf{S}) = -\mathbf{y}^{\top} \mathbf{b} + \langle \mathcal{A}^*(\mathbf{y}) + \mathbf{S} - \mathbf{C}, \mathbf{X} \rangle + \frac{1}{2\mu} \|\mathcal{A}^*(\mathbf{y}) + \mathbf{S} - \mathbf{C}\|_F^2$$

First-order optimality condition:

$$\nabla_{\mathbf{y}} \mathcal{L}_{\mu} = -\mathbf{b} + \mathcal{A}(\mathbf{X}^k) + \frac{1}{\mu} \mathcal{A}(\mathcal{A}^*(\mathbf{y}^{k+1}) + \mathbf{S}^k - \mathbf{C})$$



assume $\mathcal{A} : \mathcal{S}^n \rightarrow \mathbb{R}^m$ is onto (surjective)

$$\mathbf{y}^{k+1} = -(\mathcal{A}\mathcal{A}^*)^{-1} (\mu (\mathcal{A}(\mathbf{x}^k) - \mathbf{b})) + \mathcal{A}(\mathbf{S}^k - \mathbf{C})$$

Update \mathbf{S}

$$\mathcal{L}_\mu(\mathbf{X}, \mathbf{y}, \mathbf{S}) = -\mathbf{y}^\top \mathbf{b} + \langle \mathcal{A}^*(\mathbf{y}) + \mathbf{S} - \mathbf{C}, \mathbf{X} \rangle + \frac{1}{2\mu} \|\mathcal{A}^*(\mathbf{y}) + \mathbf{S} - \mathbf{C}\|_F^2$$

$$\begin{aligned} \mathbf{S}^{k+1} &= \arg \min_{\mathbf{S} \succeq \mathbf{0}} \mathcal{L}_\mu(\mathbf{X}^k, \mathbf{y}^{k+1}, \mathbf{S}) \\ &= \arg \min_{\mathbf{S} \succeq \mathbf{0}} \|\mathbf{S} - \mathbf{V}^{k+1}\|_F^2 \\ &= \mathbf{Q}_+ \boldsymbol{\Sigma}_+ \mathbf{Q}_+^\top \end{aligned}$$

with

$$\begin{aligned} \mathbf{V}^{k+1} &= \mathbf{C} - \mathcal{A}^*(\mathbf{y}^{k+1}) - \mu \mathbf{X}^k \\ &= \begin{bmatrix} \mathbf{Q}_+ & \mathbf{Q}_- \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}_+ & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_- \end{bmatrix} \begin{bmatrix} \mathbf{Q}_+^\top \\ \mathbf{Q}_-^\top \end{bmatrix} \quad (\text{Eigen-Decomp.}) \end{aligned}$$

Update multiplier \mathbf{X}

$$\begin{aligned}\mathbf{X}^{k+1} &= \mathbf{X}^k + \left(\mathcal{A}^*(\mathbf{y}^{k+1}) + \mathbf{S}^{k+1} - \mathbf{C} \right) / \mu \\ &= \frac{1}{\mu} \left(\mathbf{S}^{k+1} - (\mathbf{C} - \mathcal{A}(\mathbf{y}^{k+1}) - \mu \mathbf{X}^k) \right) \\ &= \frac{1}{\mu} (\mathbf{S}^{k+1} - \mathbf{V}^{k+1}) \\ &= -\frac{1}{\mu} \mathbf{Q} \Sigma \mathbf{Q}^\top\end{aligned}$$

Convergence

Semidefinite programming (SDP)

$$\text{minimize } \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{subject to } \mathcal{A}(\mathbf{X}) = \mathbf{b}, \quad \mathbf{X} \succeq \mathbf{0}$$

Assumption:

- $\mathcal{A} : \mathcal{S}^n \rightarrow \mathbb{R}^m$ is onto
- $\exists \mathbf{X}_0$ s.t. $\mathcal{A}(\mathbf{X}_0) = \mathbf{b}$ and $\mathbf{X}_0 \succ \mathbf{0}$ (**Slater's condition**)

Convergence result

Theorem (WGY, '10)

From any starting point,

$$\{(\mathbf{X}^k, \mathbf{y}^k, \mathbf{S}^k)\} \rightarrow \text{a primal and dual solution } \{(\mathbf{X}^*, \mathbf{y}^*, \mathbf{S}^*)\}$$

Another ADMM for Tensor Problem

Tensor robust principal component analysis (T-RPCA)

$$\mathcal{X} = \mathcal{L} + \mathcal{S} \in \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_K}$$

Structural assumptions:

\mathcal{L} : low rank in each mode ($\mathcal{L}_{(k)}$)

i.e. $\text{rank}(\mathcal{L}_{(k)})$ is small for all $k \in [K]$

\mathcal{S} : sparsely supported

i.e. $\text{cardinality}\{\mathcal{S}: \mathcal{S} \neq 0\}$ is small

Problem: Given \mathcal{X} , recover \mathcal{L} and \mathcal{S} .

T-RPCA

Convex surrogates

$$\text{rank}(\mathcal{L}_{(i)}) \implies \|\mathcal{L}_{(i)}\|_* := \sum \sigma_j(\mathcal{L}_{(i)})$$

$$\text{cardinality}(\mathcal{S}) \implies \|\mathcal{S}\|_1 := \sum |\mathcal{S}_{i_1 i_2 \dots i_K}|$$

Convex optimization

$$\min_{\mathcal{L}, \mathcal{S}} \sum \lambda_i \|\mathcal{L}_{(i)}\|_* + \|\mathcal{S}\|_1$$

$$\text{s.t. } \mathcal{L} + \mathcal{S} = \mathcal{X}$$

N.B.

- generalizes matrix robust PCA [CLMW, 11]
- theoretical guarantees [MHWG, '14] [HMGW, '15]

T-RPCA

Variable Splitting

$$\sum_i \lambda_i \|\mathcal{L}^{(i)}\|_*$$



$$\mathcal{L}$$

$\mathcal{L}_1 = \mathcal{L}_2 = \dots = \mathcal{L}_K$

$$\sum_i \lambda_i \|\mathcal{L}_{i,(i)}\|_*$$

T-RPCA

Reformulated problem

$$\begin{aligned} \min_{\{\mathcal{L}_i\}, \mathcal{S}} \quad & \sum \lambda_i \|\mathcal{L}_{i,(i)}\|_* + \|\mathcal{S}\|_1 \\ \text{s.t.} \quad & \mathcal{L}_i + \mathcal{S} = \mathcal{X}, \forall i \in [K] \end{aligned}$$

Augmented Lagrangian function $\mathcal{L}_\mu \left(\{\mathcal{L}_i\}, \mathcal{S}, \{\Lambda_i\} \right)$

$$\sum_{i \in [K]} \lambda_i \|\mathcal{L}_{i,(i)}\|_* + \|\mathcal{S}\|_1 + \sum_{i \in [K]} \left(-\langle \Lambda_i, \mathcal{L}_i + \mathcal{S} - \mathcal{X} \rangle + \frac{1}{2\mu} \|\mathcal{L}_i + \mathcal{S} - \mathcal{X}\|_F^2 \right)$$

ADMM k -th iteration:

$$\mathcal{L}_i \text{-update: } \left\{ \mathcal{L}_i^{k+1} \right\} \leftarrow \arg \min_{\{\mathcal{L}_i\}} \mathcal{L}_\mu \left(\{\mathcal{L}_i\}, \mathcal{S}^k, \{\Lambda_i^k\} \right)$$

$$\mathcal{S}\text{-update: } \mathcal{S}^{k+1} \leftarrow \arg \min_{\mathcal{S}} \mathcal{L}_\mu \left(\left\{ \mathcal{L}_i^{k+1} \right\}, \mathcal{S}, \{\Lambda_i^k\} \right)$$

$$\Lambda_i\text{-update: } \Lambda_i^{k+1} \leftarrow \Lambda_i^k - \frac{1}{\mu} (\mathcal{L}_i^{k+1} + \mathcal{S}^{k+1} - \mathcal{X}) \quad \forall i \in [K]$$

References

- [ZG] Zhang, T. and Golub G. H.. "Rank-one approximation to high order tensors." *SIAM Journal on Matrix Analysis and Applications* 23.2 (2001): 534-550.
- [AGHKT] Anandkumar, A., et al. "Tensor decompositions for learning latent variable models." *The Journal of Machine Learning Research* 15.1 (2014): 2773-2832.
- [HS] Hsu, D. and Sham M. K.. "Learning mixtures of spherical gaussians: moment methods and spectral decompositions." *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*. ACM, 2013.
- [TB] Kolda, T. G. and Bader B. W. Bader.. "Tensor decompositions and applications." *SIAM review* 51.3 (2009): 455-500.
- [JMZ] Jiang, B., Ma, S. and Zhang S.. "Tensor principal component analysis via convex optimization." *Mathematical Programming* 150.2 (2015): 423-457.
- [NW] Nie, J. and Wang, L.. "Semidefinite relaxations for best rank-1 tensor approximations." *SIAM Journal on Matrix Analysis and Applications* 35.3 (2014): 1155-1179.
- [CLMW] Candès, E. J., et al. "Robust principal component analysis?." *Journal of the ACM (JACM)* 58.3 (2011): 11.

References

- [WGY] Wen, Z., Goldfarb, D., & Yin, W. (2010). Alternating direction augmented Lagrangian methods for semidefinite programming. *Mathematical Programming Computation*, 2(3-4), 203-230.
- [MHWG] Mu, C., et al. "Square deal: Lower bounds and improved relaxations for tensor recovery." *Proceedings of the international conference on Machine learning*. ACM, 2014
- [GQ] Goldfarb, D., & Qin, Z.. "Robust low-rank tensor recovery: Models and algorithms." *SIAM Journal on Matrix Analysis and Applications* 35.1 (2014): 225-253.
- [HMGW] Huang, B., et al. "Provable models for robust low-rank tensor completion." *Pacific Journal of Optimization* 11 (2015): 339-364.
- [MHG] Mu, C., Hsu, D., & Goldfarb, D.. "Successive Rank-One Approximations for Nearly Orthogonally Decomposable Symmetric Tensors." *SIAM Journal on Matrix Analysis and Applications* 36.4 (2015): 1638-1659.

