# Estimating Higher-Order Moments Using Symmetric Tensor Decomposition

## Tamara G. Kolda
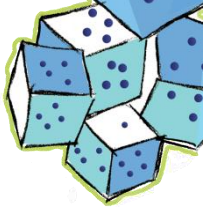
### Sandia National Labs, Livermore, CA
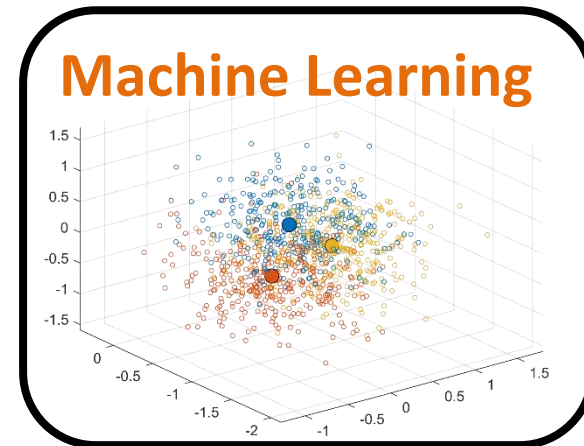### www.kolda.net

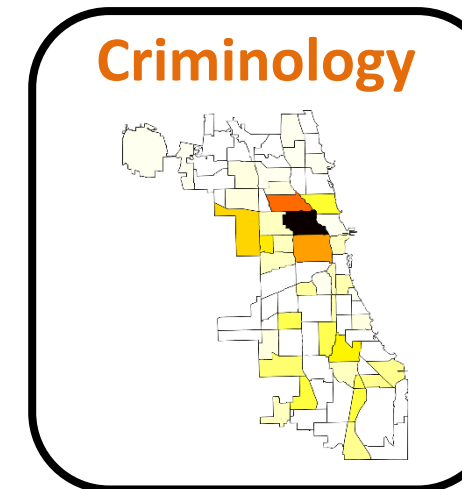Joint work with
## Samantha Sherman
### University of Notre Dame, South Bend, IN

Illustration by Chris Brigman

# Tensors Come From Many Applications

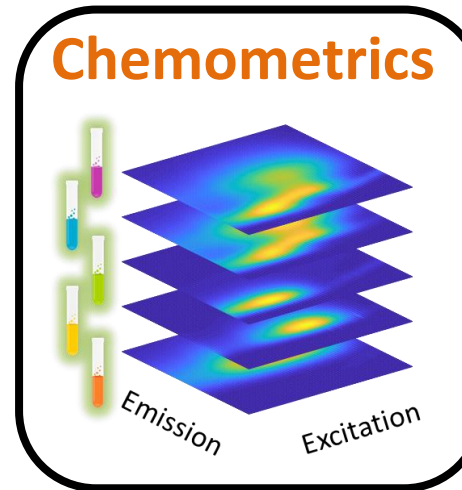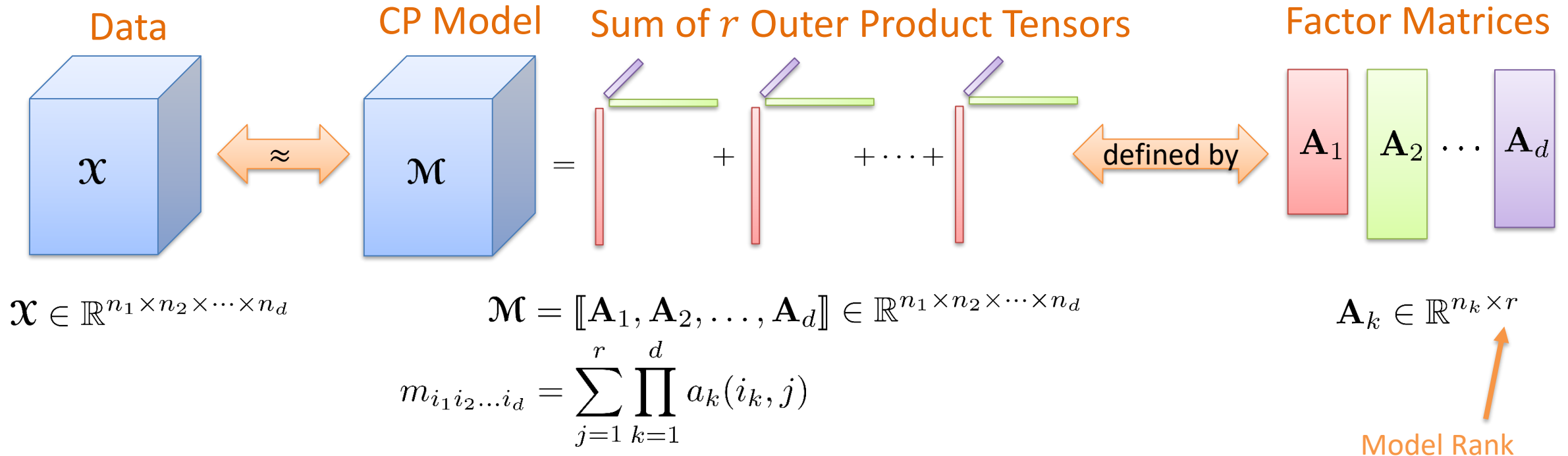- **Chemometrics:** Emission x Excitation x Samples (Fluorescence Spectroscopy)
- **Neuroscience:** Neuron x Time x Trial
- **Criminology:** Day x Hour x Location x Crime (Chicago Crime Reports)
- **Symmetric Higher-order Empirical Moments:** Multivariate Gaussian Distributions in Machine Learning
- **Transportation:** Pickup x Dropoff x Time (Taxis)
- **Sports:** Player x Statistic x Season (Basketball)
- **Cyber-Traffic:** IP x IP x Port x Time
- **Social Network:** Person x Person x Time x Interaction-Type
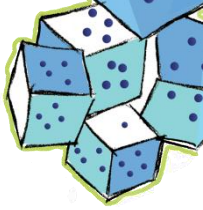- **Symmetric Higher-order Derivatives:** From Optimization

**Chemometrics**



Emission    Excitation

**Neuroscience**



trial

neuron

time

**Criminology**



**Machine Learning**

# CP Tensor Decomposition Identifies Factors



$$\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_d}$$

$$\mathcal{M} = [\![\mathbf{A}_1, \mathbf{A}_2, \ldots, \mathbf{A}_d]\!] \in \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_d}$$

$$\mathbf{A}_k \in \mathbb{R}^{n_k \times r}$$

$$m_{i_1 i_2 \ldots i_d} = \sum_{j=1}^{r} \prod_{k=1}^{d} a_k(i_k, j)$$

Optimization Formulation

$$\min_{\mathbf{A}_1, \ldots, \mathbf{A}_d} \|\mathcal{X} - \mathcal{M}\|^2 = \sum_{i_1=1}^{n_1} \cdots \sum_{i_d=1}^{n_d} (x_{i_1 i_2 \ldots i_d} - m_{i_1 i_2 \ldots i_d})^2$$

# CP First Invented in 1927



Frank Lauren Hitchcock
MIT Professor
(1875–1957)



F. L. Hitchcock, *The Expression of a Tensor or a Polyadic as a Sum of Products*, Journal of Mathematics and Physics, 1927

# CP Independently Reinvented (twice) in 1970
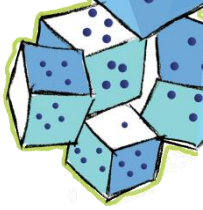
**CANDECOMP: <u>Can</u>onical <u>Decomp</u>osition**

**PARAFAC: <u>Para</u>llel <u>Fac</u>tors**



PSYCHOMETRIKA—VOL. 35, NO. 3
SEPTEMBER, 1970

ANALYSIS OF INDIVIDUAL DIFFERENCES IN MULTIDIMEN-
SIONAL SCALING VIA AN N-WAY GENERALIZATION OF
"ECKART-YOUNG" DECOMPOSITION

J. DOUGLAS CARROLL AND JIH-JIE CHANG

BELL TELEPHONE LABORATORIES
MURRAY HILL, NEW JERSEY

J. Douglas Carroll
Bell Labs
(1939-2011)

Jih-Jie Chang
Bell Labs
(1927-2007)

Richard A. Harshman
Univ. Ontario
(1943-2008)

NOTE: This manuscript was originally published in 1970 and is reproduced here to make it more accessible to interested scholars. The original reference is Harshman, R. A. (1970). Foundations of the PARAFAC procedure: Models and conditions for an "explanatory" multimodal factor analysis. *UCLA Working Papers in Phonetics*, 16, 1-84. (University Microfilms, Ann Arbor, Michigan, No. 10,085).

FOUNDATIONS OF THE PARAFAC PROCEDURE: MODELS AND CONDITIONS
FOR AN "EXPLANATORY" MULTIMODAL FACTOR ANALYSIS

by

Richard A. Harshman

UCLA

*Working Papers in Phonetics*

16

December, 1970

**CP: CANDECOMP/PARAFAC**

In 2000, Henk Kiers proposed this *compromise* name

**CP: Canonical Polyadic**

2010: Pierre Comon, Lieven DeLathauwer, and others reverse-engineered CP, revising some of Hitchcock's terminology

*Many thanks to the following persons for helping me learn about Jih-Jie Chang: Fan Chung, Ron Graham, Shen Lin (husband), May Chang (niece), Lili Bruer (daughter).*

# Tensor Decomposition in Neuroscience

- A. H. Williams et al. **Unsupervised Discovery of Demixed, Low-dimensional Neural Dynamics across Multiple Timescales through Tensor Components Analysis**. *Neuron,* 2018

- D. Hong, T. G. Kolda, J. A. Duersch. **Generalized Canonical Polyadic Tensor Decomposition**. *SIAM Review*, 2020

# Activity of Single Neuron Measured Over Time Produces Vector Data

*Thanks to Schnitzer Group @ Stanford*
Mark Schnitzer, Fori Wang, Tony Kim

Microscope by Inscopix

nVista

mouse in maze

neural activity via calcium imaging

Williams et al., Neuron, 2018

# Activity of Single Neuron Measured Over Time Produces Vector Data

*Thanks to Schnitzer Group @ Stanford*
Mark Schnitzer, Fori Wang, Tony Kim



Microscope by Inscopix

mouse in maze

neural activity via calcium imaging

Williams et al., Neuron, 2018

# Multiple Neurons Measured Over Time Produces Matrix

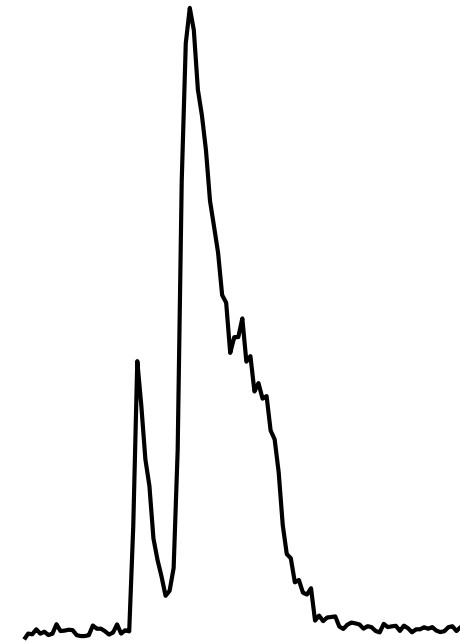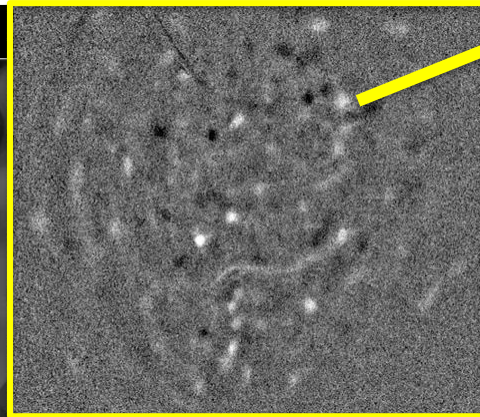*Thanks to Schnitzer Group @ Stanford*
Mark Schnitzer, Fori Wang, Tony Kim

282 neurons × 111 time bins

Microscope by
Inscopix

mouse
in "maze"

neural activity



Williams et al., Neuron, 2018

# Multiple Trials Produces 3-way Tensor

Trial 50    Trial 150    Trial 250

282 neurons × 111 time bins × 300 trials



W    N

wall

- 300 Trials over 5 Days
- Start West
- Conditions Swap Twice
  - ❖ Turn South
  - ❖ Turn North
  - ❖ Turn South

S    E

Williams et al., Neuron, 2018

# Example Neuron Activity

Thin lines show 300 individual trials

Thick line is average



Hong, Kolda, Duersch, SIAM Review, 2020

# CP Tensor for Neuron Activity Data

Data                    CP Model          Sum of $r$ Outer Product Tensors                    Factor Matrices



$$\mathcal{X} \approx \mathcal{M}$$

$$\mathcal{X} \in \mathbb{R}^{282 \times 111 \times 300}$$

$$\mathcal{M} = [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!] \in \mathbb{R}^{282 \times 111 \times 300}$$

defined by

$$\mathbf{A} \quad \mathbf{B} \quad \mathbf{C}$$

# Neuron Factor Vector Visualized as Bar Chart



Hong, Kolda, Duersch, SIAM Review, 2020

Neuron Modes Plotted as a Bar Chart
(Red Lines Correspond to Examples in Prior Slide)



Hong, Kolda, Duersch, SIAM Review, 2020

# Time Factor Vector Visualized as Line

Time (within trial) Plotted as a Line
(Dashed Line is Zero)

trial

neuron

$\mathcal{x}$

time

$\approx$ | $+$ | $+\cdots+$ |

Hong, Kolda, Duersch, SIAM Review, 2020

$c_1$

Rule
Change

Trial Plotted as Scatter Graph
Right turn = Green
Left turn = Orange
Filled = Reward

Rule
Change

Hong, Kolda, Duersch, SIAM Review, 2020

# Visualization of CP Tensor Decomposition Shows the Factors (Vectors)



Hong, Kolda, Duersch, SIAM Review, 2020

# CP Decomposition of Mouse Data

# CP Tensor Decomposition "Sees" Reward

Sam Sherman
Notre Dame

Tammy Kolda
Sandia

# Symmetric CP Tensor Factorization for (Symmetric) Moment Tensors

- S. Sherman, T. G. Kolda. **Estimating Higher-Order Moments Using Symmetric Tensor Decomposition**, revised April 2020, http://arxiv.org/abs/1911.03813

# Symmetric Tensor Entries Invariant Under Permutation of Indices

A tensor is <u>symmetric</u> if its entries are invariant under permutation of the indices

For $d$-way tensor, of dimension $n$, number of unique entries is:

$$\binom{n+d-1}{d} \approx \frac{n^d}{d!}$$

Example 1.2 from Nie (2014)

$3 \times 3 \times 3$ symmetric tensor (10 distinct entries)

$$\mathcal{X} = \begin{pmatrix} 7 & -3 & 9 & -3 & 13 & 20 & 9 & 20 & 19 \\ -3 & 13 & 20 & 13 & -27 & 6 & 20 & 6 & 6 \\ 9 & 20 & 19 & 20 & 6 & 6 & 19 & 6 & 45 \end{pmatrix}$$

$x(1,1,1) = 7 \qquad x(1,3,3) = 19$

$x(1,1,2) = -3 \qquad x(2,2,2) = -27$

$x(1,1,3) = 9 \qquad x(2,2,3) = 6$

$x(1,2,2) = 13 \qquad x(2,3,3) = 6$

$x(1,2,3) = 20 \qquad x(3,3,3) = 45$

# Symmetric CP Tensor Decomposition Has Single Factor Matrix

Symmetric
Data

Symmetric
CP Model

Sum of $r$ Symmetric
Outer Product Tensors

Single
Factor Matrix



$$\mathfrak{X} \in \mathbb{R}^{\underbrace{n \times n \times \cdots \times n}_{\text{order } d}}$$

$$\mathbf{\mathcal{M}} \in \mathbb{R}^{n \times n \times \cdots \times n}$$

defined by

$\boldsymbol{\lambda}$ $\mathbf{A}$

$$\boldsymbol{\lambda} \in \mathbb{R}^r$$
$$\mathbf{A} \in \mathbb{R}^{n \times r}$$

Model Rank

$$m_{i_1 i_2 \ldots i_d} = \sum_{j=1}^{r} \lambda_j \prod_{k=1}^{d} a(i_k, j)$$

# Symmetric Outer Product

Given **a vector**:

$$\mathbf{a} \in \mathbb{R}^n$$

The **outer product** is

$$\mathcal{P} = \mathbf{a}^{\otimes d} \in \mathbb{R}^{n \times n \times \cdots \times n}$$

$$\mathbf{a}^{\otimes 3} \equiv \mathbf{a} \otimes \mathbf{a} \otimes \mathbf{a}$$



$$n \times n \times n$$

$$(\mathbf{a}^{\otimes 3})_{ijk} = a_i a_j a_k$$

# Model Expressed as Sum of Symmetric Outer Products



Symmetric
Data

Symmetric
CP Model

Sum of $r$ Symmetric
Outer Product Tensors

Single
Factor Matrix

$\mathcal{X}$

$\approx$

$\mathcal{M}$

$=$

$+$

$+\cdots+$

defined by

$\boldsymbol{\lambda}$

$\mathbf{A}$

$$\mathcal{M} = \sum_{j=1}^{r} \lambda_j \, \mathbf{a}_j^{\otimes d} \in \mathbb{R}^{n \times n \times \cdots \times n}$$

$\boldsymbol{\lambda} \in \mathbb{R}^r$
$\mathbf{A} \in \mathbb{R}^{n \times r}$

Model Rank

# Symmetric Tensor Rank & Decomposition

Example 1.2 from Nie (2014)

$3 \times 3 \times 3$ symmetric tensor (10 distinct entries)

$$\mathcal{X} = \left( \begin{array}{rrr|rrr|rrr} 7 & -3 & 9 & -3 & 13 & 20 & 9 & 20 & 19 \\ -3 & 13 & 20 & 13 & -27 & 6 & 20 & 6 & 6 \\ 9 & 20 & 19 & 20 & 6 & 6 & 19 & 6 & 45 \end{array} \right)$$
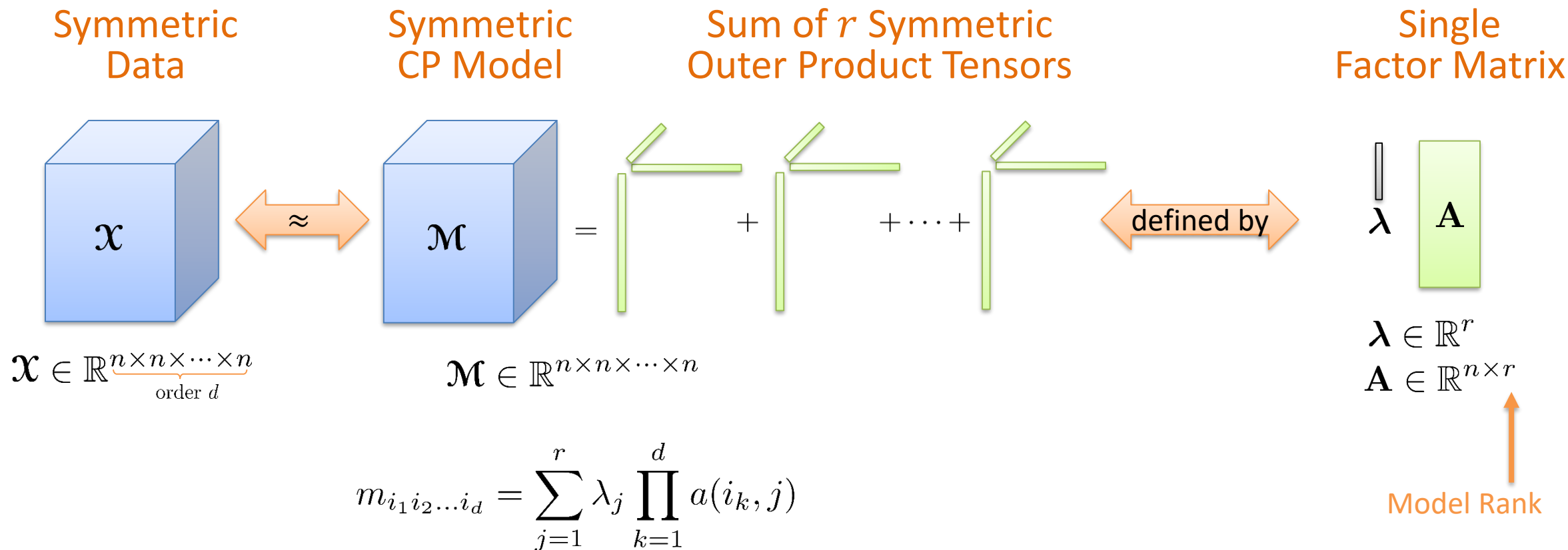
$$\mathrm{rank}(\mathcal{X}) = \min \left\{ r \mid \mathcal{X} = \mathbf{a}_1^{\otimes d} + \cdots + \mathbf{a}_r^{\otimes d} \right\}$$

Rank decomposition

$$\mathcal{X} = 2 \cdot \begin{bmatrix} 1 \\ -2 \\ -1 \end{bmatrix}^{\otimes 3} + 5 \cdot \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}^{\otimes 3} - \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix}^{\otimes 3}$$

- Symmetric tensor rank
  - For any given tensor, NP-hard to compute its rank (Hillar & Lim, 2013)
  - Typical rank known over $\mathbb{C}$ (Comon, Golub, Lim, Mourraine, 2008)
  - In practice, trial and error!
- Symmetric tensor decomposition
  - Waring decomposition (Landsberg, 2012; Oeding & Ottaviani, 2013)
  - Gröbner bases algebraic methods or numerical root-finding method (Nie, 2014)
  - Direct optimization formulation (Kolda, 2015)
  - Subspace power method (Kileel & Pereira, 2019)

# Moment Tensors Arise in Inference of Gaussian Mixture Models (GMMs)



*For ease of illustration, we focus on $n = 2$ dimensions.*
*Generally interested in much higher dimensions, i.e, $n = 500$!*

**Samples from Mixture of 3 Gaussians**

○ Samples

*Given just the samples (point cloud), can we recover the means?*

**PDF for Mixture of 3 Gaussians**

— PDF
● Means

We observe $p$ random vectors of length $n$ coming from a mixture of $r$ Gaussian distributions.
**Can we recover the means of the Gaussians?**

Easy: Means Well Separated

Hard: Means Close Together



For these pictures: $p = 1000, n = 3, r = 3$. Means shown as filled in larger circles. Samples as open circles.
We care about larger values of $n$!

# Moment Structure for Spherical GMMs Corresponds to CP Model

Data Model: $\quad V \sim \mathcal{N}(\boldsymbol{\mu}_\xi, \sigma^2 \mathbf{I}), \quad \xi \sim \mathrm{MULTI}(w_1, \ldots, w_r)$

Multivariate Normal

Probability to select $j$th center is $w_j$

3rd-order Moment:

$\boxed{\mathbb{E}[V^{\otimes 3}]} + O(\sigma^2) = \boxed{\sum_{j=1}^{r} w_j \boldsymbol{\mu}_j^{\otimes 3}}$

*Can also do higher - order moments*

Calculate empirically from data

CP-like Model

$$\mathbf{X} = \frac{1}{p} \sum_{\ell=1}^{p} \mathbf{v}_\ell^{\otimes 3}$$

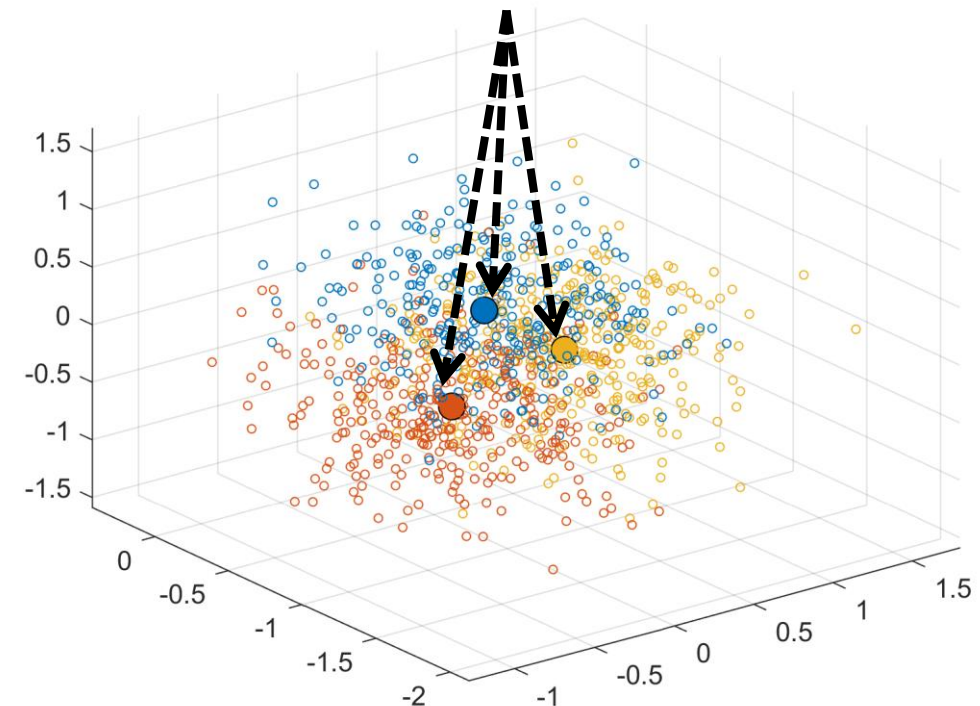$$\mathbf{M} = \sum_{j=1}^{r} \lambda_j \, \mathbf{a}_j^{\otimes 3}$$

Bottlenecks:
$O(pn^d)$ to compute,
$O(n^d)$ to store

Hsu and Kakade, 2013

**Samples from Mixture of 3 Gaussians**



○   Samples

Example: $n$ = 128, $d$ = 4 $\Rightarrow$ storage = 2 GB

Example: $n$ = 512, $d$ = 3 $\Rightarrow$ storage = 1 GB

Simplifying assumptions for this work

$\|\boldsymbol{\mu}_j\|_2 = 1 \;\forall j \in [r]$

$\omega_j = \frac{1}{r} \;\forall j \in [r]$

$$x = \frac{1}{p} \sum_{\ell=1}^{p} \mathbf{v}_\ell^{\otimes d}$$

**Symmetric Data**

**Given Observations**

$$\mathcal{X} = \quad + \quad + \cdots + \quad$$

$$\xleftrightarrow{\text{defined by}}$$

$$\mathbf{V} \in \mathbb{R}^{n \times p}$$

**Symmetric CP Model**
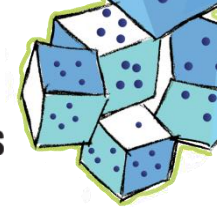
**Want to Find Compact Representation**

$$\mathcal{M} = \quad + \quad + \cdots + \quad$$

$$\xleftrightarrow{\text{defined by}}$$

$$\mathbf{A} \in \mathbb{R}^{n \times r}$$

$$\mathcal{M} = \sum_{j=1}^{r} \lambda_j \, \mathbf{a}_j^{\otimes d}$$

$$r \ll p$$

# Optimization Approach for Symmetric CP of Symmetric Tensor Requires TTSV

Optimization Problem

$$\min_{\boldsymbol{\lambda}, \mathbf{A}} F(\boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{M}}) \equiv \frac{1}{2} \|\boldsymbol{\mathcal{X}} - \boldsymbol{\mathcal{M}}\|^2 \text{ where } \boldsymbol{\mathcal{M}} = \sum_{j=1}^{r} \lambda_j \mathbf{a}_j^{\otimes d}$$

Gradients
$\forall j \in [r]$

$$\frac{\partial F}{\partial \mathbf{a}_j} = -d\lambda_j \boxed{\boldsymbol{\mathcal{X}}\mathbf{a}_j^{d-1}} + d\lambda_j \sum_{k=1}^{r} \lambda_k \langle \mathbf{a}_j, \mathbf{a}_k \rangle^{d-1} \mathbf{a}_k$$

$$\frac{\partial F}{\partial \lambda_j} = -\boldsymbol{\mathcal{X}}\mathbf{a}_j^{d} + \sum_{k=1}^{r} \lambda_k \langle \mathbf{a}_j, \mathbf{a}_k \rangle^{d}$$

*Plug function and gradient into favorite optimization method. My favorite: L-BFGS.*

Bottleneck is TTSV which costs $O(n^d)$



Key Kernel: Tensor Times Single Vector (TTSV)

$$\left(\boldsymbol{\mathcal{X}}\mathbf{a}^{d-1}\right)_{i_1} = \sum_{i_2=1}^{n} \cdots \sum_{i_d=1}^{n} \left( x_{i_1 i_2 \dots i_d} \prod_{k=2}^{d} a_{i_k} \right) \ \forall i_1 \in [n]$$

TTSV Definition:

$$\left(\mathbf{\mathcal{X}}\mathbf{a}^{d-1}\right)_{i_1} = \sum_{i_2=1}^{n} \cdots \sum_{i_d=1}^{n} \left( x_{i_1 i_2 \ldots i_d} \prod_{k=2}^{d} a_{i_k} \right) \; \forall i_1 \in [n]$$

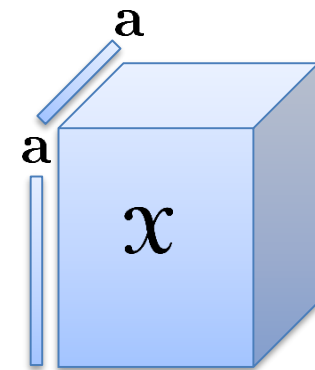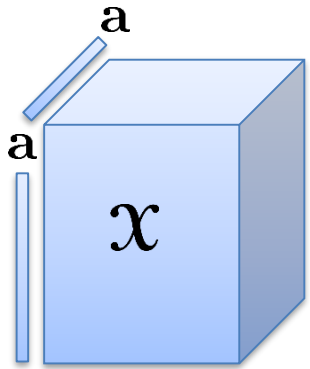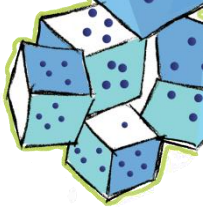**Lemma.** Let $\mathbf{\mathcal{X}} = \dfrac{1}{p} \sum_{\ell=1}^{p} \mathbf{v}_\ell^{\otimes d}$ and $\mathbf{V} = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_p \end{bmatrix}$, then

*Entry-wise Power*

$$\mathbf{\mathcal{X}}\mathbf{a}^{d-1} = \frac{1}{p} \mathbf{V} \left[\mathbf{V}^\mathsf{T}\mathbf{a}\right]^{d-1}$$

$O(n^d)$

$O(pn)$

**Left column:**

1: **function** FG_EXPLICIT$(\mathcal{X}, \boldsymbol{\lambda}, \mathbf{A}, \alpha)$

2:     **for** $j = 1, \ldots, r$, **do** $\mathbf{y}_j = \mathcal{X}\mathbf{a}_j^{d-1}$, **end**

3:     **for** $j = 1, \ldots, r$, **do** $w_j = \mathbf{a}_j^T \mathbf{y}_j$, **end**

4:     $\mathbf{B} = \mathbf{A}^T \mathbf{A}$

5:     $\mathbf{C} = [\mathbf{B}]^{d-1}$

6:     $\mathbf{u} = (\mathbf{B} * \mathbf{C})\boldsymbol{\lambda}$

7:     $f = \alpha + \boldsymbol{\lambda}^T \mathbf{u} - 2\mathbf{w}^T \boldsymbol{\lambda}$

8:     $\mathbf{g}_{\boldsymbol{\lambda}} = -2(\mathbf{w} - \mathbf{u})$

9:     $\mathbf{G}_{\mathbf{A}} = -2d(\mathbf{Y} - \mathbf{A}\mathbf{D}_{\boldsymbol{\lambda}}\mathbf{C})\mathbf{D}_{\boldsymbol{\lambda}}$

10:     **return** $f, \mathbf{g}_{\boldsymbol{\lambda}}, \mathbf{G}_{\mathbf{A}}$

11: **end function**

**Right column:**

1: **function** FG_IMPLICIT$(\mathbf{V}, \boldsymbol{\lambda}, \mathbf{A}, \alpha)$

2:     $\mathbf{Y} = \frac{1}{p}\mathbf{V}[\mathbf{V}^\mathsf{T}\mathbf{A}]^{d-1}$

3:     **for** $j = 1, \ldots, r$, **do** $w_j = \mathbf{a}_j^T \mathbf{y}_j$, **end**

4:     $\mathbf{B} = \mathbf{A}^T \mathbf{A}$

5:     $\mathbf{C} = [\mathbf{B}]^{d-1}$

6:     $\mathbf{u} = (\mathbf{B} * \mathbf{C})\boldsymbol{\lambda}$

7:     $f = \alpha + \boldsymbol{\lambda}^T \mathbf{u} - 2\mathbf{w}^T \boldsymbol{\lambda}$

8:     $\mathbf{g}_{\boldsymbol{\lambda}} = -2(\mathbf{w} - \mathbf{u})$

9:     $\mathbf{G}_{\mathbf{A}} = -2d(\mathbf{Y} - \mathbf{A}\mathbf{D}_{\boldsymbol{\lambda}}\mathbf{C})\mathbf{D}_{\boldsymbol{\lambda}}$

10:     **return** $f, \mathbf{g}_{\boldsymbol{\lambda}}, \mathbf{G}_{\mathbf{A}}$

11: **end function**

# Experimental Difference in Per-Iteration Cost of Implicit versus Explicit

Rank-$r$ Symmetric CP Tensor Factorization
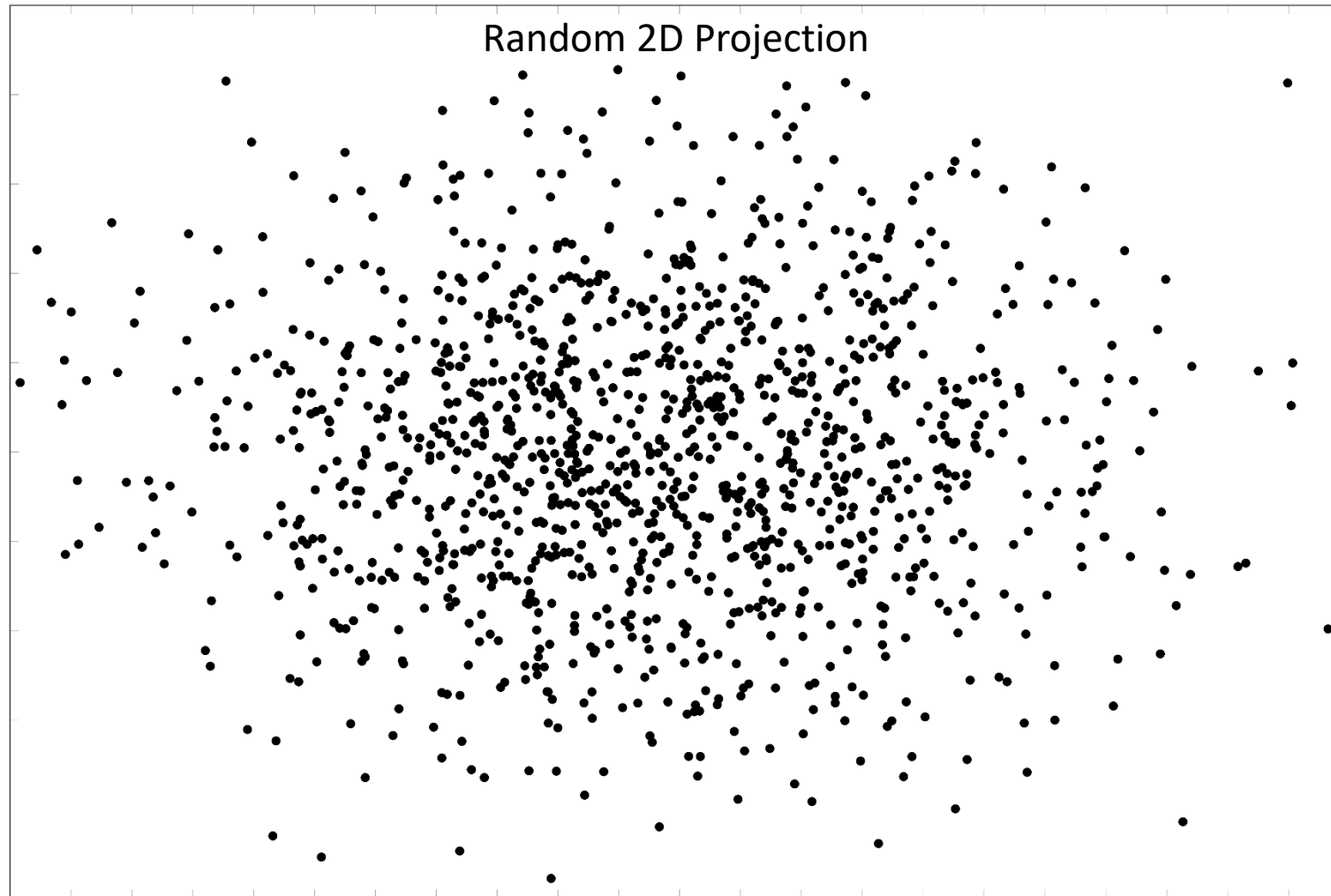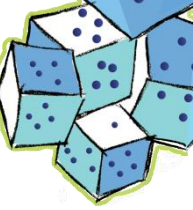for $d$-way tensor of size $n$

$$r < n < p$$

| Method | Storage | Per-Iteration |
|---|---|---|
| Explicit | $O(n^d)$ | $O(rn^d)$ |
| Implicit | $O(pn)$ | $O(pnr)$ |

*Implicit cheaper if* $p < O(n^{d-1})$

Average cost per iteration for $r = 5$ over 10 runs

| $d$ | $n$ | $p$ | $n^{d-1}$ | Explicit | Implicit |
|---|---|---|---|---|---|
| 3 | 75 | 3750 | 5625 | 5e-4 | 8e-4 |
| 3 | 375 | 3750 | 140625 | 2e-2 | 5e-3 |
| 4 | 75 | 3750 | 421875 | 1e-2 | 9e-4 |

For $d = 3$, explicit method requires 1 GB storage

For $d = 4$, explicit method requires 500 GB storage



Random 2D Projection

Random 2D Projection, Color-Coded by Component

# GMM Example with $r$=5 (mixtures), $n$=500 (dimension) and $p$=750 (observations)

$$\boldsymbol{\mu}_j \in \mathbb{R}^{500}$$

$$\|\boldsymbol{\mu}_j\|_2 = 1$$

$$\forall j \in [r]$$

$$\boldsymbol{\mu}_j^T \boldsymbol{\mu}_k = 0.5$$

$$\forall j \neq k$$

Shown here:
$$\sigma = 0.1$$



Random 2D Projection, Color-Coded by Component, With Means Denoted

# Choosing Starting Guess Within Range of Observations is Key!

Randomized Range Finder (RRF): $\quad \mathbf{A}_0 = \mathbf{V\Omega}, \; \mathbf{\Omega} \sim \mathcal{N}(0,1)^{p \times \hat{r}}$

Random: $\quad \mathbf{A}_0 \sim \mathcal{N}(0,1)^{n \times \hat{r}}$

[with columns normalized in both cases]

Results of computing
$\hat{r} = 3$ approximation for
moment tensor of order $d = 3$,
with
$r = 3$ components,
$n = 500$ dimensions, and
$p = 750$ observations
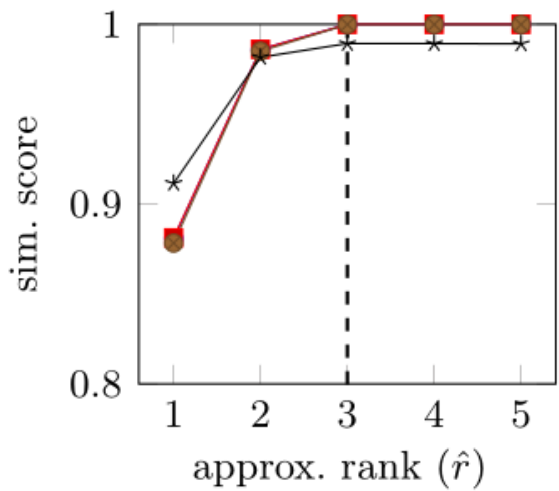
Similarity Score for Best Rel. Error of 10 Runs for $n = 500$ (Dimension), Varying Other Parameters

# Identified Factors for $\hat{r}$=5
## with $r = 5, p = 750, n = 500, \sigma = 0.1$

$\boldsymbol{\mu}_j \in \mathbb{R}^{500}$

$\|\boldsymbol{\mu}_j\|_2 = 1$

$\forall j \in [r]$

$\boldsymbol{\mu}_j^T \boldsymbol{\mu}_k = 0.5$

$\forall j \neq k$

Shown here:
$\sigma = 0.1$

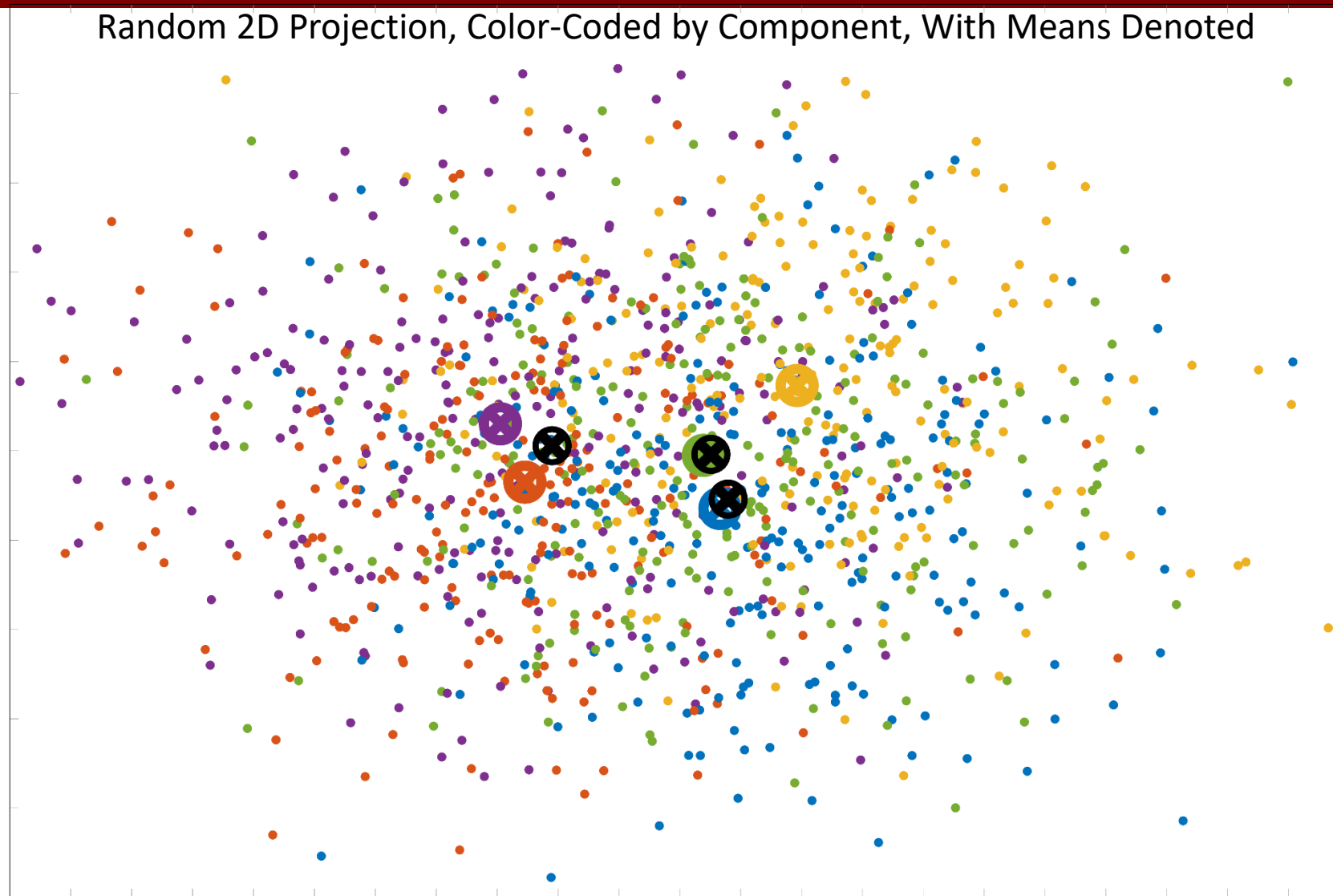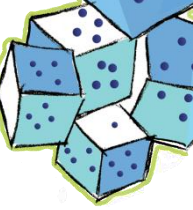Random 2D Projection, Color-Coded by Component, With Means Denoted

Random 2D Projection, Color-Coded by Component, With Means Denoted
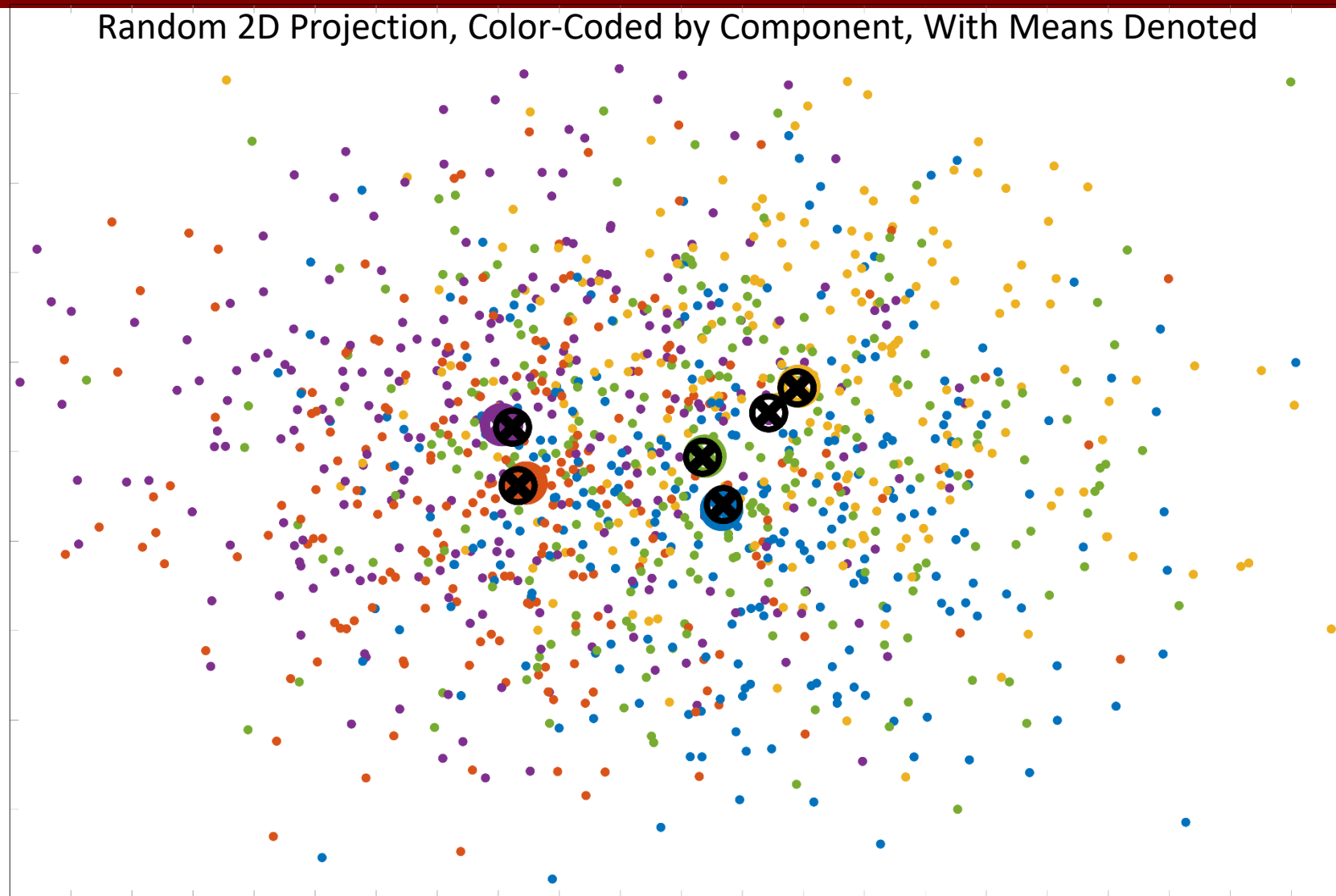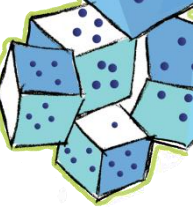
# Identified Factors for $\hat{r}$=4
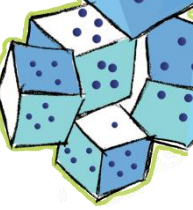## with $r = 5, p = 750, n = 500, \sigma = 0.1$



Random 2D Projection, Color-Coded by Component, With Means Denoted

Random 2D Projection, Color-Coded by Component, With Means Denoted

Random 2D Projection, Color-Coded by Component, With Means Denoted

# For Massive Numbers of Observations, Use Stochastic Variants

$$V \in \mathbb{R}^{n \times p}$$

Sample columns with replacement

$$\tilde{V} \in \mathbb{R}^{n \times s}$$

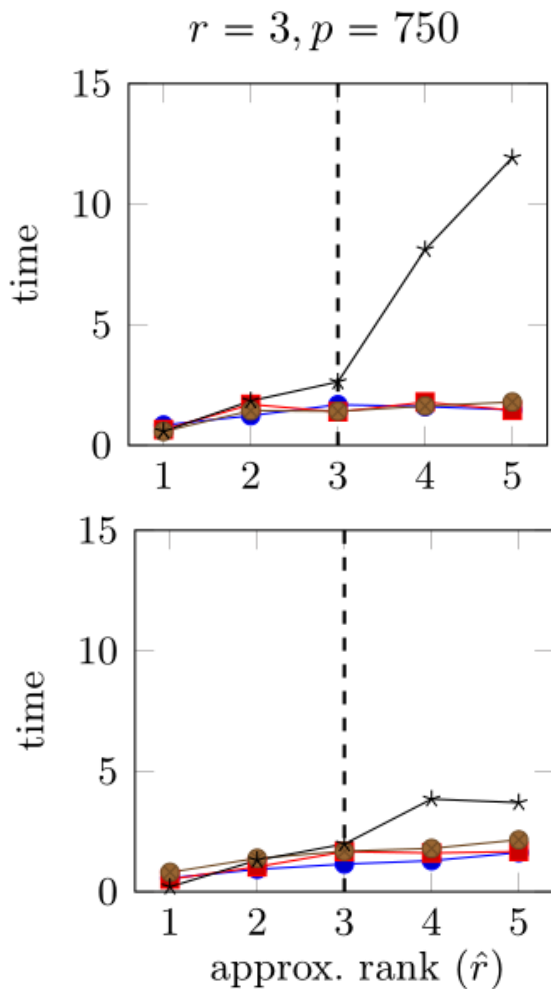$$\mathfrak{X} = \frac{1}{p} \sum_{\ell=1}^{p} v_\ell^{\otimes d}$$

$$\tilde{\mathfrak{X}} = \frac{1}{s} \sum_{\ell=1}^{s} \tilde{v}_\ell^{\otimes d}$$

$$\Rightarrow \qquad \mathbb{E}[\tilde{\mathfrak{X}}a^{d-1}] = \mathfrak{X}a^{d-1}$$
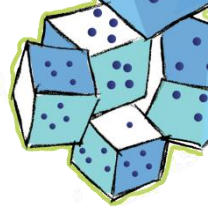
**Example Results**

$\hat{r} = r = 10, n = 500,$
$\sigma = 0.1, d = 3$
$p = 100{,}000$

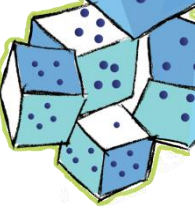| Method | Best $f$ (shifted) | Sim. Score | Total Time (s) |
|---|---|---|---|
| standard | $-0.2471$ | $0.9998$ | $2166.70$ |
| Adam, s=10 | $-0.2209$ | $0.9225$ | $8.03$ |
| Adam, s=100 | $-0.2427$ | $0.9929$ | $10.48$ |
| Adam, s=1000 | $-0.2464$ | $0.9990$ | $41.00$ |

**Best Runs (of 10)**

$$\hat{r} = r = 10, n = 500, \sigma = 0.1, d = 3, p = 100{,}000$$

# Conclusions and Future Work

- In data analysis, $d$th-order moment is expensive to compute – instead work with implicit moment
  - Reduces storage from $O(n^d)$ to $O(np)$
  - Reduces computation per iteration from $O(rn^d)$ to $O(rnp)$
- Shows promise for fitting spherical GMMs
  - Example with $n = 500$ (dimension), $r \in \{3,5,10\}$ (components), $p = 250r$, $\hat{r} \in \{r - 2, \ldots, r + 2\}$, and $d = 3,4$ (orders)
  - Future work will incorporate lower-order terms, different $\sigma$ for each component, multiple values for $d$ simultaneously, etc.
- Many extensions possible, e.g., for subspace power method
- Reference: S. Sherman, T. G. Kolda. **Estimating Higher-Order Moments Using Symmetric Tensor Decomposition**, submitted for publication, 2019, arXiv:1911.03813