



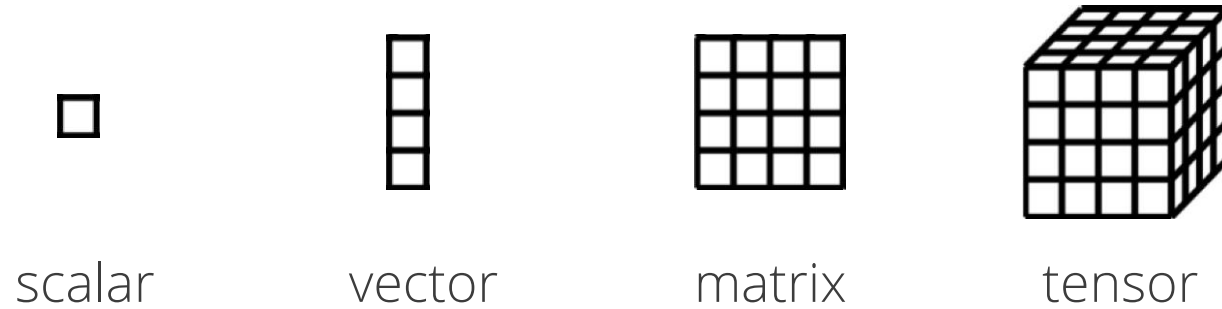
Exceptional service in the national interest

# THE POISSON CANONICAL POLYADIC TENSOR MODEL AS A LATENT-VARIABLE MODEL

Carlos Llosa, Daniel M. Dunlavy, Richard B. Lehoucq,  
Oscar Lopez, Arvind Prasad

August 2024, JSM

# OVERVIEW



- **Goal:** Understand key relationships in tensor data
- **Current approach**
  - Low-rank tensor models (canonical polyadic, Tucker, tensor train, ...)
  - Parameter inference via maximum likelihood estimation
- **Our approach**
  - Latent-variable model formulation
  - Parameter inference via complete-data loglikelihood
    - EM algorithms for maximum likelihood estimation
    - Fisher information matrix

# FROM DENSE-CONTINUOUS TO SPARSE-DISCRETE TENSOR DATA ANALYSIS

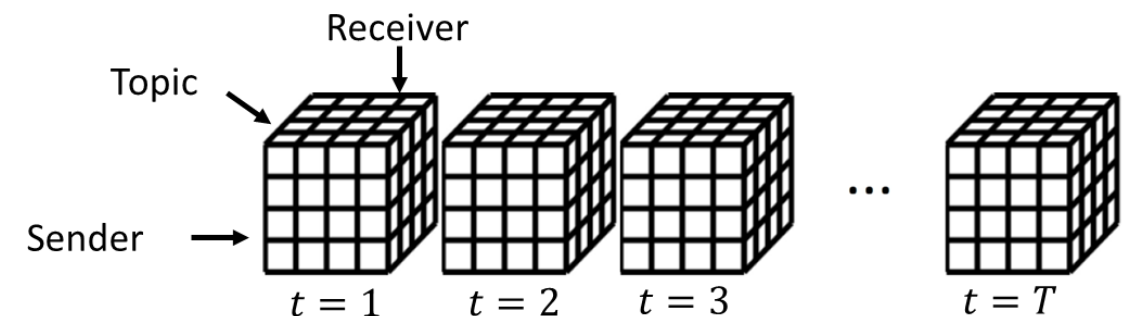
## Netflix Prize [1]

- Sparse data: only ~1% of entries are observed
- Ranking data (1-5)
- Winner algorithm used matrix factorization techniques
- This led to increased interest in non-Gaussian matrix factorization

		Movies									
		1	2	3	4	5	6	7	8	...	17,700
Users	1		4		2	4					
	2	3		3				3			
	3		3			1					
	4			2			4		1		5
	...										
480,000			2					3		1	

## ICEWS Database [2]

- Countries as receivers and senders
- Events such as threats or aid
- Count data: number of times an event happens from a receiver to a sender
- Sparse data, low-count data



[1] Bennett and Lanning, *The Netflix Prize*, Proc. of KDD Cup and Workshop, 2007.

[2] O'Brien, *Crisis Early Warning and Decision Support: Contemporary Approaches and Thoughts on Future Research*, ISR, 2010.

# POISSON CANONICAL POLYADIC (PCP) TENSOR MODEL [3][4]

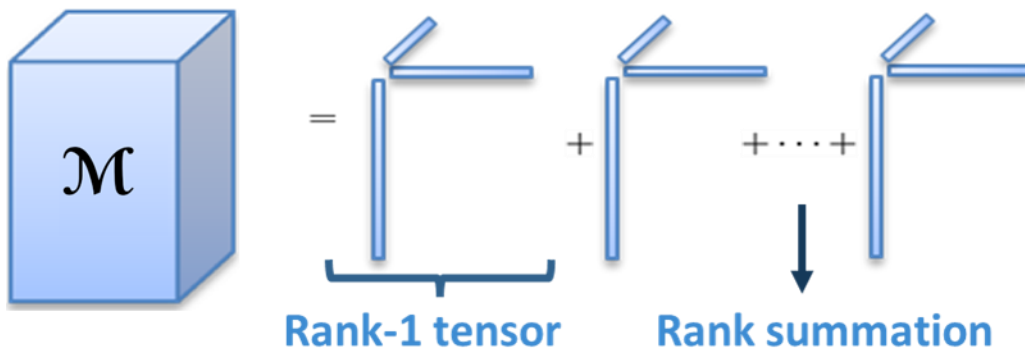
The count tensor follows a Poisson distribution element-wise

$$\mathcal{X} \sim \text{Poisson}(\mathcal{M}) \iff \mathcal{X}_{i,j,k} \stackrel{\text{indep.}}{\sim} \text{Poisson}(\mathcal{M}_{i,j,k})$$

We present results for 3-way tensors, but our work generalizes to arbitrary D-way tensors

$$\ell(\theta) = \sum_{i,j,k} [\mathcal{X}_{i,j,k} * \log(\mathcal{M}_{i,j,k}) - \mathcal{M}_{i,j,k}] + \text{constant}$$

The parameter tensor is imposed a Canonical Polyadic tensor structure



$$\theta = [\text{vec}(A)' \text{vec}(B)' \text{vec}(C)']'$$

$$\mathcal{M}_{i,j,k} = \sum_{r=1}^R A_{i,r} B_{j,r} C_{k,r}$$

[3] Lee and Seung, *Algorithms for Non-negative Matrix Factorization*, NeurIPS 2000.

[4] Chi and Kolda, *On Tensors, Sparsity, and Nonnegative Factorizations*, SIMAX 2012.

# PCP TENSOR MODEL: CHALLENGES AND APPROACHES

$$\ell(\boldsymbol{\theta}) = \sum_{i,j,k} [\mathbf{x}_{i,j,k} * \log(\mathcal{M}_{i,j,k}) - \mathcal{M}_{i,j,k}] + \text{constant}$$

$$\mathcal{M}_{i,j,k} = \sum_{r=1}^R A_{i,r} B_{j,r} C_{k,r}$$

## Optimization Approach [3,4]

How to efficiently optimize the loglikelihood?

- MM optimization [3,4]

Their MM algorithms are actually EM algorithms!

- Higher order methods

Our Fisher Info can be used for Fisher scoring optimization!

## Our Approach

PCP is a latent-variable model!

- EM algorithms
- Fisher information

Evaluate model and fit

- Well-posed statistical problems
- Evaluate convergence of algorithm

## Probabilistic Approach [5]

How many entries do I need to recover  $\mathcal{M}$ ?

- Matrix/tensor completion
- Typically an upper bound on MSE

Our Fisher Info can be used for Cramer Rao bounds on MSE!

[3] Lee and Seung, *Algorithms for Non-negative Matrix Factorization*, NeurIPS 2000.

[4] Chi and Kolda, *On Tensors, Sparsity, and Nonnegative Factorizations*, SIMAX 2012.

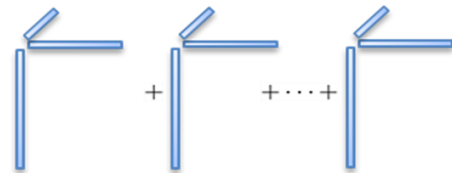
[5] Cao and Xie, *Poisson Matrix Recovery and Completion*, IEEE TSP 2016

# FIRST MAIN RESULT: PCP IS A LATENT-VARIABLE MODEL

Random Variable

$$\mathcal{X} = \text{3D grid}$$

$$\mathcal{X} \sim \text{Poisson} \left( \sum_{r=1}^R A[:, r] \circ B[:, r] \circ C[:, r] \right)$$



Latent mechanism  $\mathcal{X} \stackrel{d}{=} \sum_{r=1}^R \mathcal{Z}_r$

Latent Random Variable

$$\mathcal{Z} = \left( \mathcal{Z}_1 \quad \mathcal{Z}_2 \quad \dots \quad \mathcal{Z}_R \right)$$

$$\mathcal{Z}_r \sim \text{Poisson} (A[:, r] \circ B[:, r] \circ C[:, r])$$



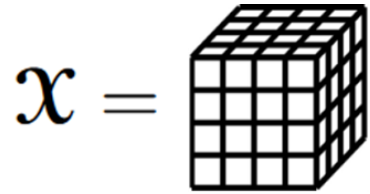
$\mathcal{X}$  follows a rank- $R$  PCP model



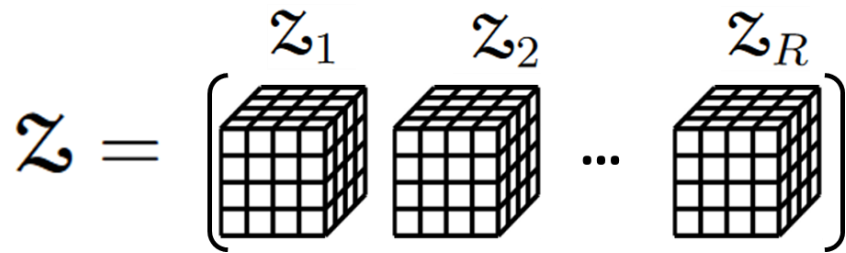
$\mathcal{X}$  is the sum of  $R$  independent  $\mathcal{Z}_r$  each following a rank-1 PCP model

# PCP AS A LATENT-VARIABLE MODEL

Observed Data  $\mathbf{x} = \text{vec}(\mathcal{X})$



Complete Data  $\mathbf{z} = \text{vec}(\mathcal{Z})$



Latent mechanism  $\mathbf{x} \stackrel{d}{=} \sum_{r=1}^R \mathcal{Z}_r$

$$\boldsymbol{\theta} = [\text{vec}(A)' \text{vec}(B)' \text{vec}(C)']'$$

Loglikelihood  $x_{i,j,k} \sim \text{Poisson} \left( \sum_r A_{i,r} B_{j,r} C_{k,r} \right)$

$$\begin{aligned} \ell(\mathbf{x}|\boldsymbol{\theta}) &:= \log p(\mathbf{x}|\boldsymbol{\theta}) \\ &= \underbrace{\log p(\mathbf{z}|\boldsymbol{\theta})}_{\ell_c(\mathbf{z}|\boldsymbol{\theta})} - \underbrace{\log p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})}_{\ell_m(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} \end{aligned}$$

Complete loglikelihood  
 $z_{r,i,j,k} \sim \text{Poisson}(A_{i,r} B_{j,r} C_{k,r})$

Missing loglikelihood  
 $(z_{1,i,j,k} \cdots z_{R,i,j,k}) | x_{i,j,k} \sim \text{Multinomial}(n = x_{i,j,k}, p_1, \dots, p_R)$   

$$p_r = \frac{A_{i,r} B_{j,r} C_{k,r}}{\sum_r A_{i,r} B_{j,r} C_{k,r}}$$

## SECOND MAIN RESULT: EXISTING MM ALGORITHMS CAN BE DERIVED AS EM ALGORITHMS

E-step  $Q(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}) := \mathbb{E}_{\mathbf{z}|\mathbf{x}, \bar{\boldsymbol{\theta}}}(\ell_c(\boldsymbol{\theta}))$

$$= \sum_{r,i,j,k} \left[ \bar{\mathbf{z}}_{r,i,j,k} * \log(A_{i,r} B_{j,r} C_{k,r}) - (A_{i,r} B_{j,r} C_{k,r}) \right]$$

$$\bar{\mathbf{z}}_{r,i,j,k} := \mathbb{E}_{\mathbf{z}|\mathbf{x}, \bar{\boldsymbol{\theta}}}(\mathbf{z}_{r,i,j,k}) = \mathbf{x}_{i,j,k} \frac{\bar{A}_{i,r} \bar{B}_{j,r} \bar{C}_{k,r}}{\sum_r \bar{A}_{i,r} \bar{B}_{j,r} \bar{C}_{k,r}}$$

← This is the mean of a multinomial

CM-step Here  $\boldsymbol{\theta}$  is split into 3 blocks, corresponding to  $A$ ,  $B$  and  $C$ .

Solve  $\underbrace{\operatorname{argmax}_A Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)})}_{\text{This is EM algorithm!}} = \left\{ \sum_{j,k} \bar{\mathbf{z}}_{r,i,j,k}^{(t)} \right\}_{i,r} = A * \underbrace{([\mathbf{x}_{(1)} \oslash (A(C \odot B)')] (C \odot B))}_{\text{This is the update used in [3] and [4]!!}}$

[3] Lee and Seung, *Algorithms for Non-negative Matrix Factorization*, NeurIPS 2000.

[4] Chi and Kolda, *On Tensors, Sparsity, and Nonnegative Factorizations*, SIMAX 2012.



# A PATH TOWARDS THE FISHER INFORMATION

Latent mechanism  $\mathbf{x} \stackrel{d}{=} \sum_{r=1}^R \mathbf{z}_r$

## The Missing Information Principle [6]

“observed information” equals the “complete information” minus the “missing information”

$$\mathcal{I}(\boldsymbol{\theta}; \mathbf{x}) := - \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \ell(\mathbf{x} | \boldsymbol{\theta})$$

$$= \underbrace{\mathbb{E} \left[ - \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \ell_c(\mathbf{z}, \mathbf{x} | \boldsymbol{\theta}) \middle| \mathbf{x} \right]}_{\mathcal{I}_c(\boldsymbol{\theta}; \mathbf{x})} - \underbrace{\mathbb{E} \left[ - \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \ell_m(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}) \middle| \mathbf{x} \right]}_{\mathcal{I}_m(\boldsymbol{\theta}; \mathbf{x})}$$

observed information

complete information

missing information

Poisson  $\left( \sum_r A_{i,r} B_{j,r} C_{k,r} \right)$

Poisson  $(A_{i,r} B_{j,r} C_{k,r})$

Multinomial  $(n = x_{i,j,k}, p_1, \dots, p_R)$

# RANK 1 PCP CASE: FISHER INFORMATION MATRIX

When  $R=1$ , the complete data is observed  $\mathbf{x} = \mathbf{z}$ . No information is lost.

Parameter  
Vector

$$\boldsymbol{\theta} = [\mathbf{a}'\mathbf{b}'\mathbf{c}']'$$

Model

$$\mathbf{x} \sim \text{Poisson}(\mathbf{a} \circ \mathbf{b} \circ \mathbf{c})$$

Loglikelihood

$$\ell(\boldsymbol{\theta}) = \sum_{i,j,k} [\mathbf{x}_{i,j,k} * \log(\mathbf{a}_i \mathbf{b}_j \mathbf{c}_k) - \mathbf{a}_i \mathbf{b}_j \mathbf{c}_k] + \text{constant}$$

Fisher  
Information

$$\mathcal{I}(\boldsymbol{\theta}) = \begin{bmatrix} \text{diag}(\mathbf{a}^{-1}) & \mathbf{1}\mathbf{1}' & \mathbf{1}\mathbf{1}' \\ \mathbf{1}\mathbf{1}' & \lambda \text{diag}(\mathbf{b}^{-1}) & \lambda \mathbf{1}\mathbf{1}' \\ \mathbf{1}\mathbf{1}' & \lambda \mathbf{1}\mathbf{1}' & \lambda \text{diag}(\mathbf{c}^{-1}) \end{bmatrix}$$

- Above we are parameterizing  $\boldsymbol{\theta}$  so that  $\lambda = \mathbf{a}'\mathbf{1}$  and  $\mathbf{1} = \mathbf{b}'\mathbf{1} = \mathbf{c}'\mathbf{1}$ . Any parameterization follows from this
- For  $\mathbf{x}$  of size  $N_1 \times N_2 \times N_3$ ,  $\mathcal{I}(\boldsymbol{\theta})$  is square with  $N_1 + N_2 + N_3$  rows and columns
- $\mathcal{I}(\boldsymbol{\theta})$  is singular of rank  $N_1 + N_2 + N_3 - 2$
- The FIM is nonsingular if you remove any one entry from  $\mathbf{b}$  and any one entry from  $\mathbf{c}$

# GENERAL RANK PCP CASE, AND OAKES' THEOREM

- For the general rank case, we have missing information.
- Unlike Gaussian CP, the Fisher information for  $R=1$  is not a special case of general rank
- Direct differentiation of the loglikelihood is challenging
- We can leverage the missing information principle

$$\mathcal{I}(\boldsymbol{\theta}; \mathbf{x}) = \mathcal{I}_c(\boldsymbol{\theta}; \mathbf{x}) - \mathcal{I}_m(\boldsymbol{\theta}; \mathbf{x})$$



Similar to  
 $R=1$  case



Challenging

Latent mechanism

$$\mathbf{x} \stackrel{d}{=} \sum_{r=1}^R \mathbf{z}_r$$

- Many techniques exist for obtaining/estimating  $\mathcal{I}_m(\boldsymbol{\theta}; \mathbf{x})$  from the complete loglikelihood
- Most popular is Louis' method [7], but can only be evaluated at the MLE.
- I use Oakes' method [8], which is more general:  $\mathcal{I}_m(\boldsymbol{\theta}; \mathbf{x}) = \left[ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \bar{\boldsymbol{\theta}}'} Q(\boldsymbol{\theta}, \bar{\boldsymbol{\theta}}) \right]_{\bar{\boldsymbol{\theta}}=\boldsymbol{\theta}}$   $\mathbb{E}_{\mathbf{z}|\mathbf{x}, \bar{\boldsymbol{\theta}}}(\ell_c(\boldsymbol{\theta}))$

[7] Louis, *Finding the Observed Information Matrix when Using the EM Algorithm*, JRSSB 1982

[8] Oakes, *Direct Calculation of the Information Matrix via the EM Algorithm*, JRSSB 1999

# GENERAL RANK PCP: FISHER INFORMATION MATRIX

## Parameter Vector

Model

$$\mathbf{x}_{a,b,c} \stackrel{\text{indep.}}{\sim} \text{Poisson}(\mathcal{M}_{a,b,c})$$

$$\boldsymbol{\theta} = [\text{vec}(A)' \text{vec}(B)' \text{vec}(C)']'$$

Loglikelihood

$$\ell(\boldsymbol{\theta}) = \sum_{a,b,c} [\mathbf{x}_{a,b,c} * \log(\mathcal{M}_{a,b,c}) - \mathcal{M}_{a,b,c}]$$

$$\mathcal{M}_{a,b,c} = \sum_{r=1}^R A[a,r]B[b,r]C[c,r]$$

Fisher Information

$$\mathcal{I}(\boldsymbol{\theta}) = \left\{ \left\{ G_{k,l}^{r,s} \right\}_{r,s=1,\dots,R} \right\}_{k,l=1,2,3} \quad G_{k,l}^{r,s} = \begin{cases} \text{diagonal matrix} & k = l \\ \text{dense matrix} & k \neq l \end{cases}$$

- $\mathcal{I}(\boldsymbol{\theta})$  is a 3 x 3 block matrix, where each block is itself a  $R \times R$  block matrix
- Above is for arbitrary parameterization of  $\boldsymbol{\theta}$
- For  $\mathbf{x}$  of size  $N_1 \times N_2 \times N_3$ ,  $\mathcal{I}(\boldsymbol{\theta})$  is square with  $R(N_1 + N_2 + N_3)$  rows and columns
- $\mathcal{I}(\boldsymbol{\theta})$  is singular of rank  $R(N_1 + N_2 + N_3 - 2)$
- The FIM is nonsingular if you remove one entry from each column of B, and one entry from each column of C

# BIAS-VARIANCE TRADE-OFF AND THE CRAMER-RAO LOWER BOUND

Consider a true parameter vector  $\theta_o = [\text{vec}(A)^\top \text{vec}(B)^\top \text{vec}(C)^\top]^\top$  and its estimated counterpart  $\hat{\theta}$

Call

1. Bias<sup>2</sup>       $\|\mathbb{E}(\hat{\theta}) - \theta_o\|^2$
2. Variance:     $\mathbb{E}\|\hat{\theta} - \mathbb{E}(\hat{\theta})\|^2$
3. MSE:         $\mathbb{E}\|\hat{\theta} - \theta_o\|^2$
4. CRLB:        $\text{tr}(\mathcal{I}^\dagger(\theta_o))$

Then:

- $\text{MSE} = \text{Bias}^2 + \text{Variance}$
- $\text{Bias}^2 = 0 \implies \text{CRLB} \leq \text{Variance}$

## Monte-Carlo Study

- Draw  $\mathcal{X}_k^*$  from  $\text{Poisson}(\mathcal{M}(\theta_o))$
- Estimate  $\hat{\theta}_k^*$  from  $\mathcal{X}_k^*$ .
- Repeat for some large  $K$   
 $k = 1, 2, \dots, K$

Monte-Carlo approximations:

1. Mean:         $\bar{\theta}_{MC} = K^{-1} \sum_k \hat{\theta}_k^*$
2. Bias2:        $\|\bar{\theta}_{MC} - \theta_o\|^2$
3. Variance:     $K^{-1} \sum_k \|\hat{\theta}_k^* - \bar{\theta}_{MC}\|^2$
4. MSE:         $K^{-1} \sum_k \|\hat{\theta}_k^* - \theta_o\|^2$

# BIAS-VARIANCE TRADE-OFF AND THE CRAMER-RAO LOWER BOUND

## Simulation Setup

Visualize the bias-variance trade-off for tensors  $\mathcal{M}(\theta_0)$  with varying:

- Entry-wise mean  $s = 1, 2, 3, 4$   
 $s = \text{mean}(\mathcal{M})$
- Sizes  $N = 10, 15, 20$   
 $\mathcal{M} \in \mathbb{R}^{N \times N \times N}$
- Rank  $R = 1, 2, \dots, 16$

## Cramer-Rao Lower Bound

$$\left( \frac{\partial}{\partial \theta} \mathbb{E}(\hat{\theta}) \right) \mathcal{I}^\dagger(\theta) \left( \frac{\partial}{\partial \theta} \mathbb{E}(\hat{\theta}) \right)^\top \leq \text{Var}(\hat{\theta})$$

↓

Jacobian identity

$$\frac{\partial}{\partial \theta} \mathbb{E}(\hat{\theta}) = \text{Cov} \left( \hat{\theta}, \frac{\partial}{\partial \theta} \log f(\mathbf{x}; \theta) \right)$$

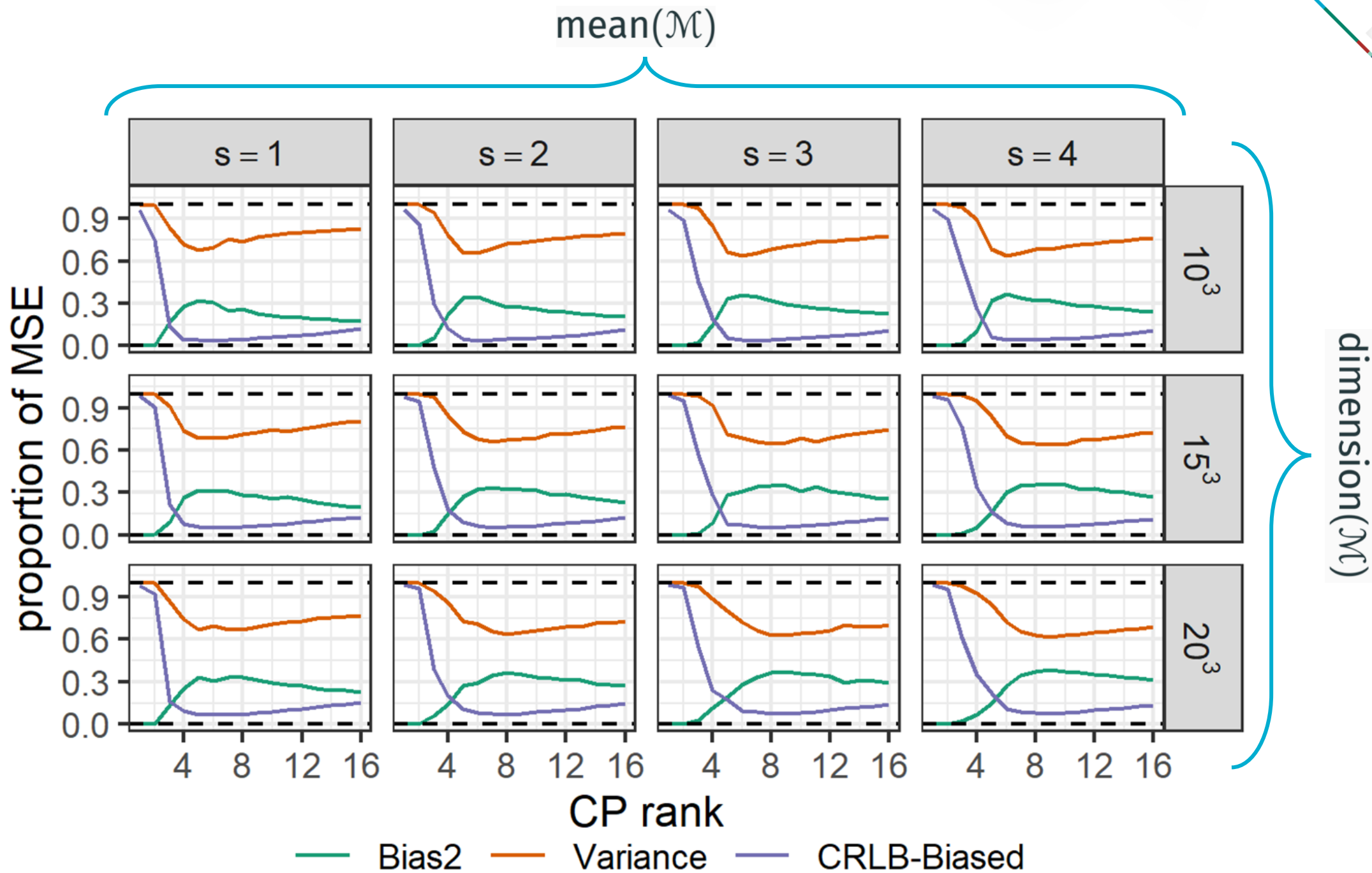
↓

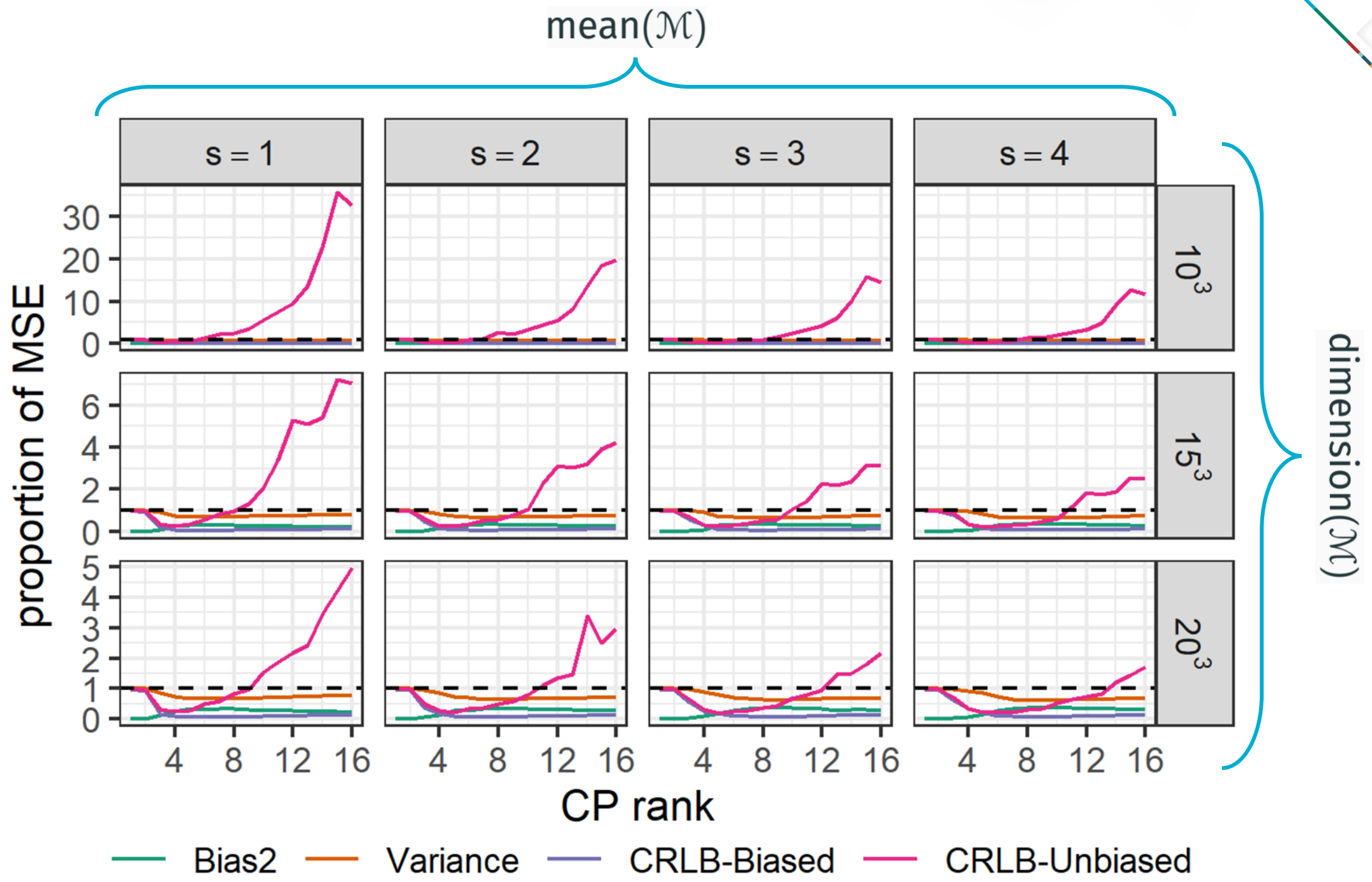
Score simplification

$$\frac{\partial}{\partial \theta} \log f(\mathbf{x}; \theta) = \begin{bmatrix} \text{vec}[\mathcal{S}_{(1)}(\mathbf{C} \odot \mathbf{B})] \\ \text{vec}[\mathcal{S}_{(2)}(\mathbf{C} \odot \mathbf{A})] \\ \text{vec}[\mathcal{S}_{(3)}(\mathbf{B} \odot \mathbf{A})] \end{bmatrix}$$

↓

$$\mathcal{S}_i = \mathcal{X}_i / \mathcal{M}_i - 1$$







## CONCLUSIONS AND PATH FORWARD

- PCP is a very popular method, we demonstrate it's a latent variable model
  - Applications in topic modeling, document clustering and classification, poll analysis, etc.
  - We allow for parameter inference through complete loglikelihood
- We rediscover popular estimating algorithms as instances of EM algorithms
  - Shed light on the properties of existing algorithms
  - Help bridge two fields of machine learning research
- Derived novel Fisher information matrix, using the missing information principle
  - Can be used to propose new Fisher scoring algorithms, Cramer Rao inequalities
  - Allows us to gauge the conditions for a well-posed parameter inference problem
- Variance-trade off simulation study
  - Bias-variance trade-off of PCP
  - Comparison against CRLB



# THANK YOU!

## The Poisson Canonical Polyadic Tensor Model as a Latent-Variable Model

Carlos Llosa [cjllosa@sandia.gov](mailto:cjllosa@sandia.gov)