



THE UNIVERSITY OF UTAH

Bayesian Streaming Sparse Tucker Decomposition

Shikai Fang, Robert M. Kirby, Shandian Zhe

Presenter: Shikai Fang

School of computing, The University of Utah

For UAI 2021



THE UNIVERSITY OF UTAH

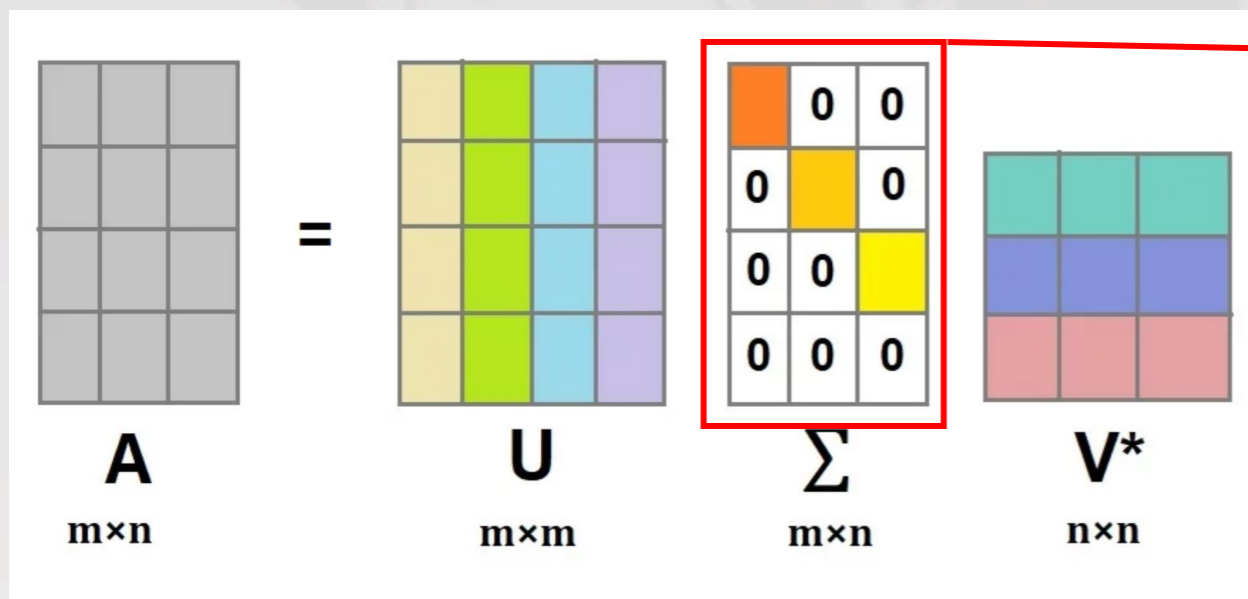
Outline

1. Background
2. Motivation
3. Bayesian Sparse Tucker Model
4. One-shot & Streaming Inference
5. Experiments on Real-world Data



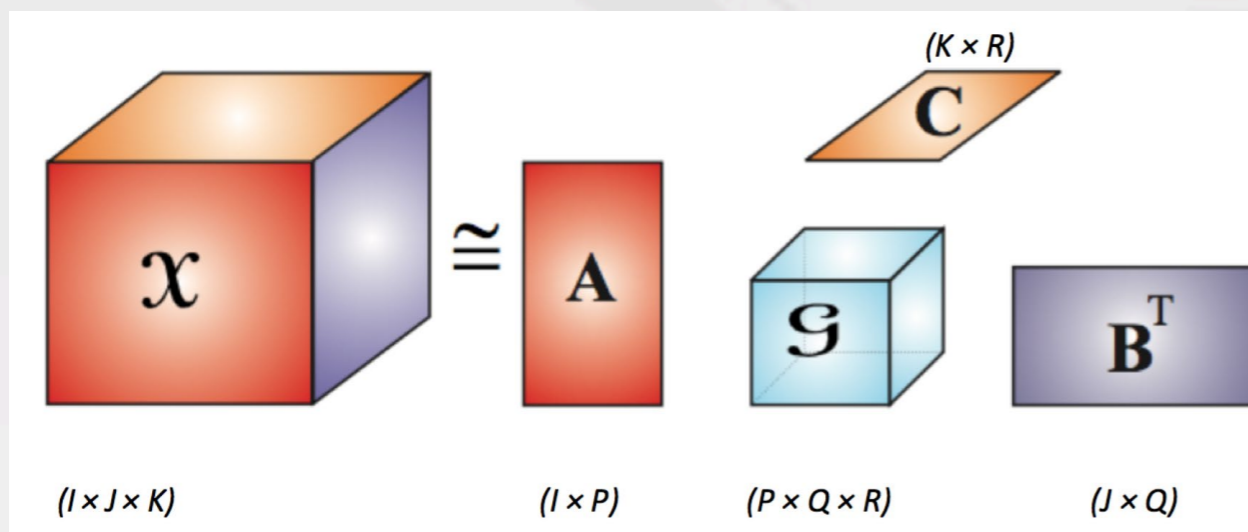
1. Background

- Tucker tensor decomposition: Generalization of the matrix SVD
(also called *Higher Order Singular Value Decomposition* ([HOSVD](#)))



SVD core:

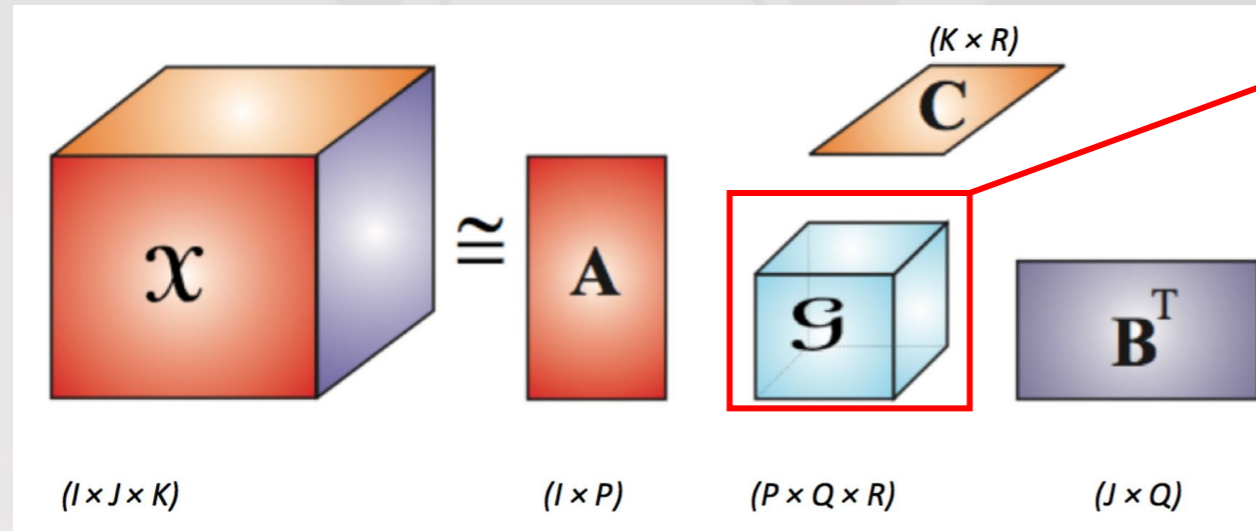
- 2-D diagonal matrix
- only model interactions of embeddings on same dim





1. Background

- Tucker tensor decomposition: Generalization of the matrix SVD
(also called *Higher Order Singular Value Decomposition* ([HOSVD](#)))



- Tucker core (3-mode example):
- 3-d dense tensor
 - model all possible interactions of embeddings at every dim

- Element-wise form for a K-mode tensor \mathcal{Y} :

$$y_i \approx \sum_{r_1=1}^{R_1} \cdots \sum_{r_K=1}^{R_K} \left[\underbrace{w(r_1, \dots, r_K)}_{\text{Core tensor element: interaction weight}} \cdot \prod_{k=1}^K \underbrace{u_{i_k, r_k}^k}_{\text{Embeddings}} \right]$$

Traverse each dim per mode

One interaction



1. Background

- Probabilistic/Bayesian version of tucker decomposition:
everything is random variable (distribution)
- For uncertainty measure and robustness
- Element-wise form for a K-mode tensor \mathcal{Y} :

$$y_i \approx \sum_{r_1=1}^{R_1} \cdots \sum_{r_K=1}^{R_K} \left[w(r_1, \dots, r_K) \cdot \prod_{k=1}^K u_{i_k, r_k}^k \right].$$

All random variables:
place priors and do inference!



2. Motivation:

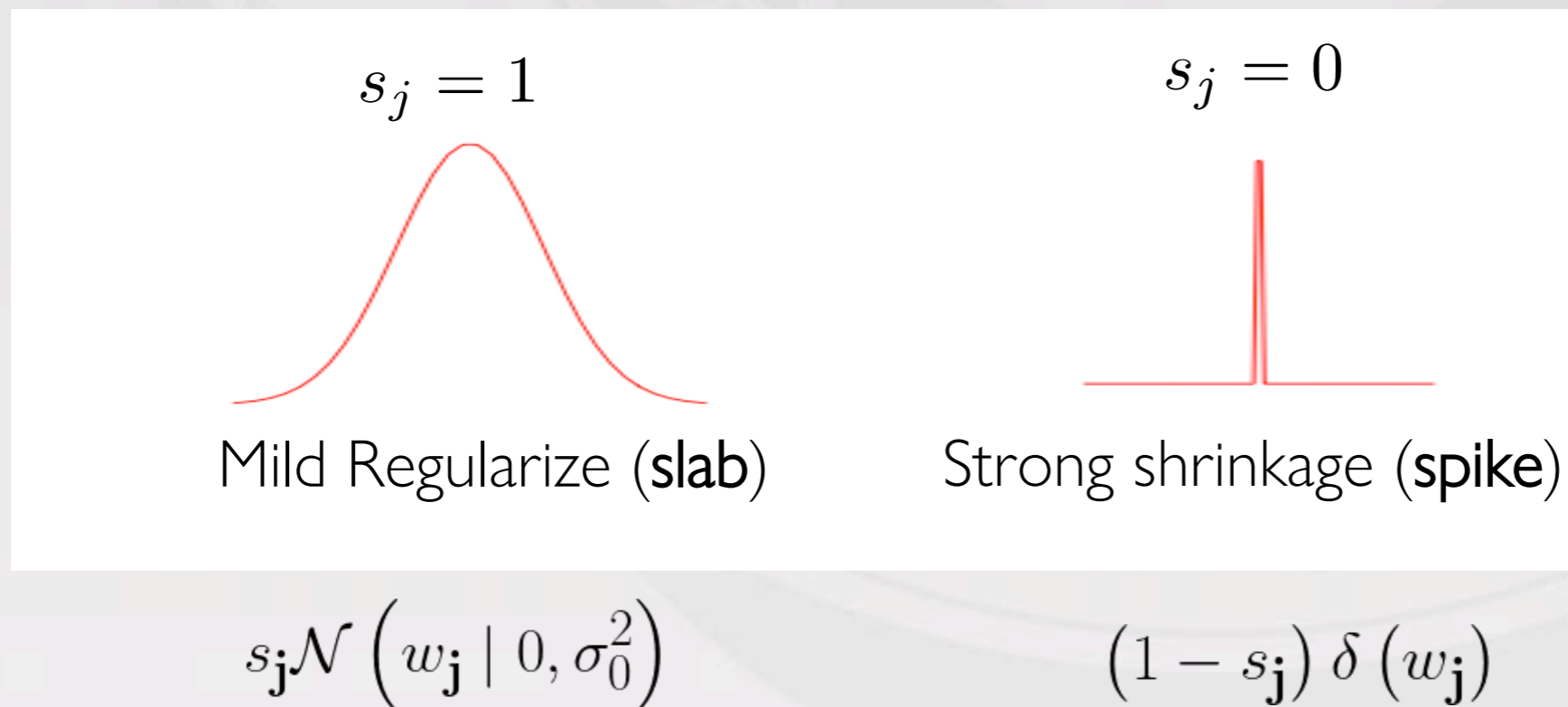
- Classical Tucker tensor decomposition is featured as
 - I. Flexible: model all possible interactions
 - II. Interpretable: core tensor indicates interaction strengthsbut suffers from:
 - I. Estimating core-tensor is memory & computationally intensive
 - II. Overparameterizing and Overfitting risk, esp. for sparse data
 - III. The two problems are more severe for streaming data!

- Goal
 - I. Alleviate over-parameterization: automatic selecting meaningful interactions
 - II. Efficient streaming posterior inference



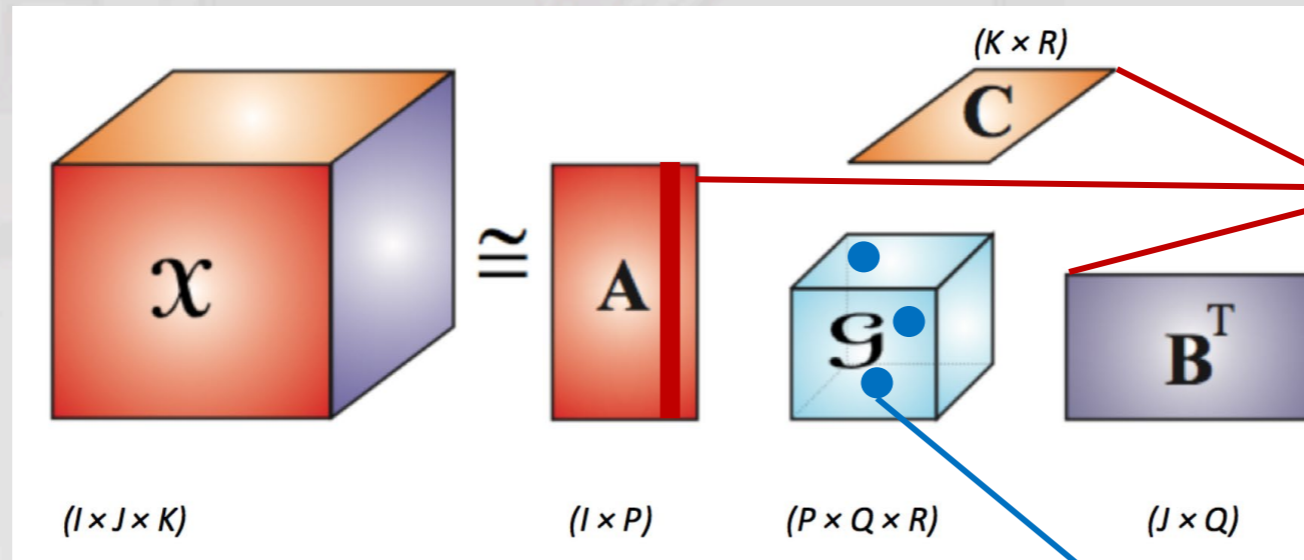
3. Bayesian Sparse Tucker Model

- Spike and Slab priors :
Introduce binary selection indicators
on each core tensor element ! s_1, s_2, \dots, s_d





3. Bayesian Sparse Tucker Model



S&S Prior over each core tensor element

$$p(\mathcal{S}, \mathcal{W}, \mathcal{U}, \mathcal{Y}, \tau) = \prod_k \prod_j \mathcal{N}(\mathbf{u}_j^k | \mathbf{0}, \mathbf{I}) \text{Gam}(\tau | a_0, b_0)$$

$$\cdot \prod_j \text{Bern}(s_j | \rho_0) \left(s_j \mathcal{N}(w_j | 0, \sigma_0^2) + (1 - s_j) \delta(w_j) \right)$$

Binary indicators

Selective shrinkage priors on core based on indicator

$$\cdot \prod_{\mathbf{i} \in \mathcal{F}} \mathcal{N}\left(y_{\mathbf{i}} | \mathcal{W} \times_1 \left(\mathbf{u}_{i_1}^1\right)^\top \times_2 \dots \times_K \left(\mathbf{u}_{i_K}^K\right)^\top, \tau^{-1}\right)$$



3. Bayesian Sparse Tucker Model

- Exact posterior distribution: Intractable!
- Approximation with distributions in exponential family:

$$q_{\text{cur}}(\mathcal{W}, \mathcal{U}, \tau) \propto p(\mathcal{S}) \xi(\mathcal{W}, \mathcal{S}) \cdot \prod_{k=1}^n \prod_{j=1}^{a_k} \mathcal{N}(\mathbf{u}_j^k \mid \boldsymbol{\mu}_j^k, \mathbf{V}_j^k) \cdot \mathcal{N}(\text{vec}(\mathcal{W}) \mid \mathbf{m}, \boldsymbol{\Sigma}) \text{Gam}(\tau \mid a, b)$$

Approximation of SS priors

Approximation of data likelihood

where:

$$\begin{aligned} \xi(\mathcal{W}, \mathcal{S}) &= \prod_j \xi_j(w_j, s_j) \\ &= \prod_j \text{Bern}(s_j \mid c(\rho_j)) \mathcal{N}(w_j \mid m_j, \eta_j) \propto p(\mathcal{W} \mid \mathcal{S}), \end{aligned}$$



3. Streaming & One-shot inference

- Streaming: data come, model update, data drops
- Incremental Bayesian rule:

Exact / Approx posterior on old data

$$p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}} \cup \mathcal{D}_{\text{new}}) \propto p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}}) p(\mathcal{D}_{\text{new}} \mid \boldsymbol{\theta})$$



3. Streaming & One-shot inference

- Streaming: data come, model update, data drops
- Incremental Bayesian rule:

Exact / Approx posterior on old data

$$p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}} \cup \mathcal{D}_{\text{new}}) \propto p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}}) p(\mathcal{D}_{\text{new}} \mid \boldsymbol{\theta})$$

Data likelihood on current model



3. Streaming & One-shot inference

- Streaming: data come, model update, data drops
- Incremental Bayesian rule:

Exact / Approx posterior on old data

$$p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}} \cup \mathcal{D}_{\text{new}}) \propto p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}}) p(\mathcal{D}_{\text{new}} \mid \boldsymbol{\theta})$$

Exact / Approx posterior on all data

Data likelihood on current model



3. Streaming & One-shot inference

- Streaming: data come, model update, data drops
- Incremental Bayesian rule:

Exact / Approx posterior on old data

$$p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}} \cup \mathcal{D}_{\text{new}}) \propto p(\boldsymbol{\theta} \mid \mathcal{D}_{\text{old}}) p(\mathcal{D}_{\text{new}} \mid \boldsymbol{\theta})$$

Exact / Approx posterior on all data

Data likelihood on current model

- Classical ADF: integrating data points one by one via **moment matching** --- inefficient (esp for core tensor update); many approximations
- Our goal: assimilating a **batch of streaming** data points at a time:
-- more efficient and improve the quality



3. Streaming & One-shot inference

- For tractable moment in ADF, we made 3 tech-contributions

I. Conditional Expectation Propagation(CEP)

$$\mathbb{E}_{\tilde{p}}[\phi(\mathcal{W})] = \mathbb{E}_{\tilde{p}(\Theta_{\setminus w})} \left[\mathbb{E}_{\tilde{p}(\mathcal{W}|\Theta_{\setminus w})} \phi(\mathcal{W}) \mid \Theta_{\setminus w} \right] \text{ Tractable Conditional Moment!}$$

II. Delta method: Expectation on first-order Taylor approximation

$$\mathbb{E}_{q_{\text{cur}}}(\Theta_{\setminus \mathcal{W}}) \left[\mathbf{h}(\Theta_{\setminus \mathcal{W}}) \right] \approx \mathbf{h} \left(\mathbb{E}_{q_{\text{cur}}} \left[\Theta_{\setminus \mathcal{W}} \right] \right) \text{ h : first-order approx. at the mean}$$

III. Repeated update of S&S prior approx. to ensure sparsity inducing effect

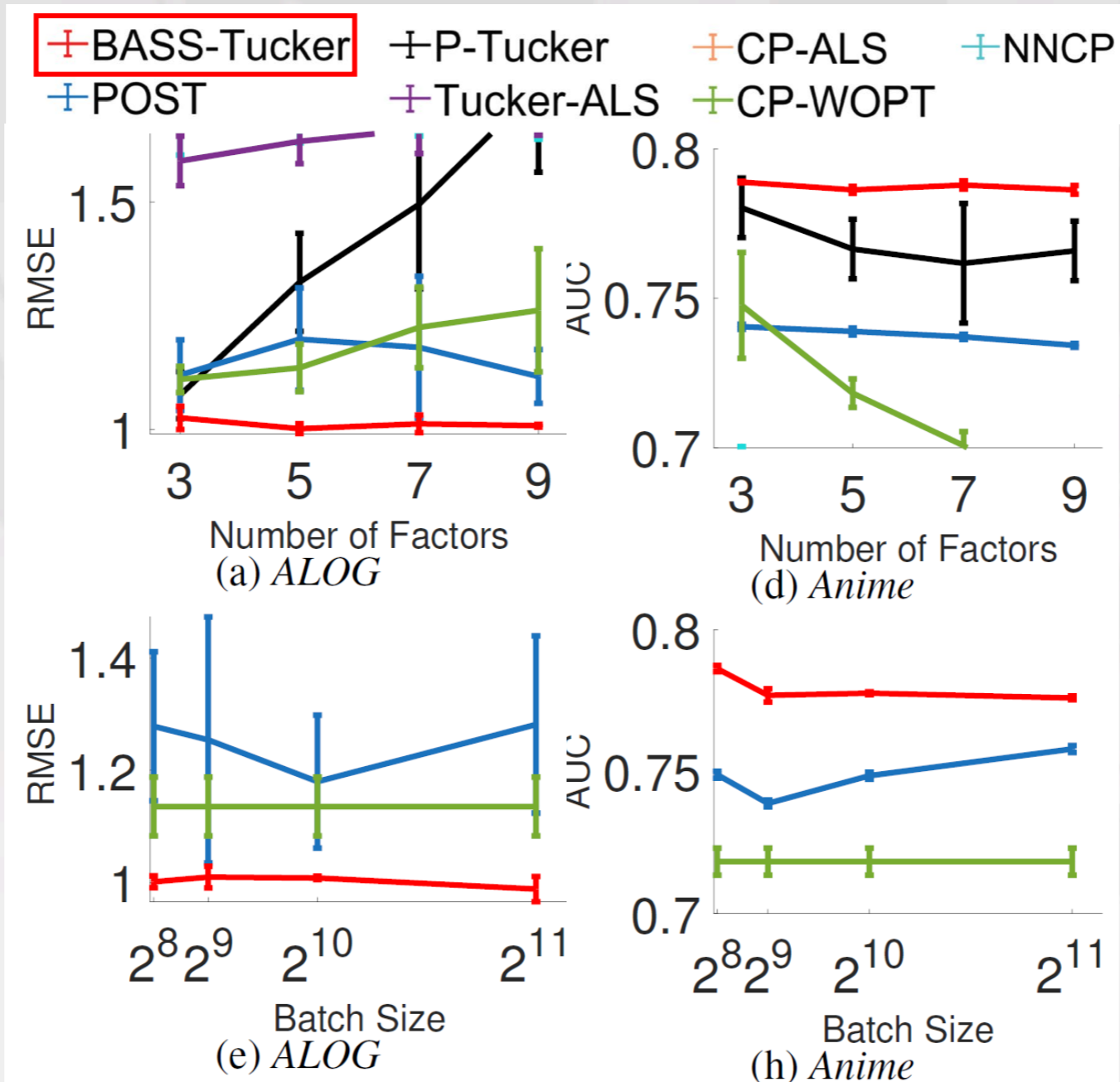
$$q^*(w_j, s_j) = \text{Bern} \left(s_j \mid c \left(\rho_j^* \right) \right) \mathcal{N} \left(w_j \mid \mu_j^*, v_j^* \right)$$



5. Experiments on real-world data

- Predictive performance on large real-world datasets
- With different factors / streaming batch size

Our method!

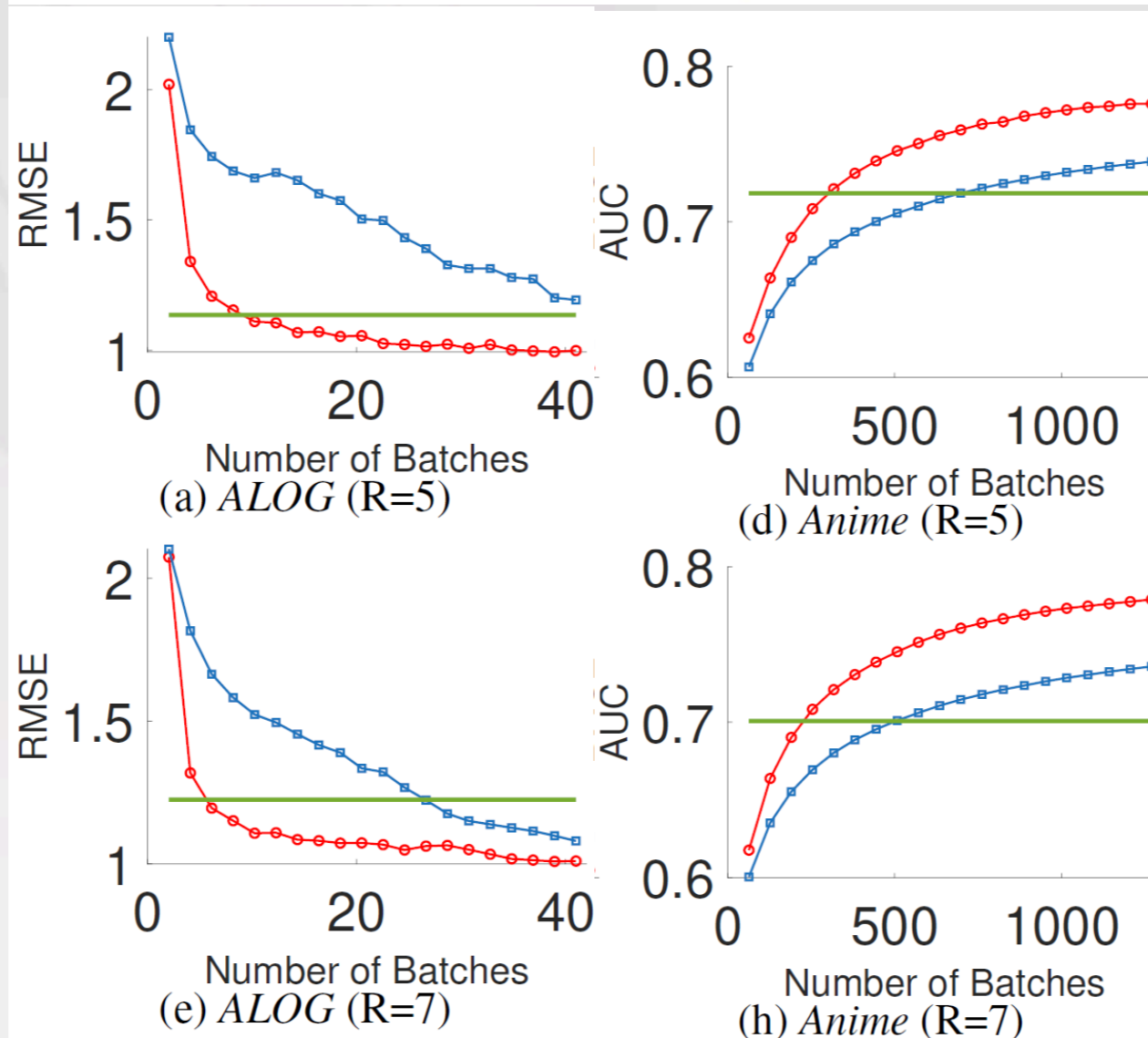




5. Experiments on real-world data

- Running prediction large real-world datasets
- With different number of factors

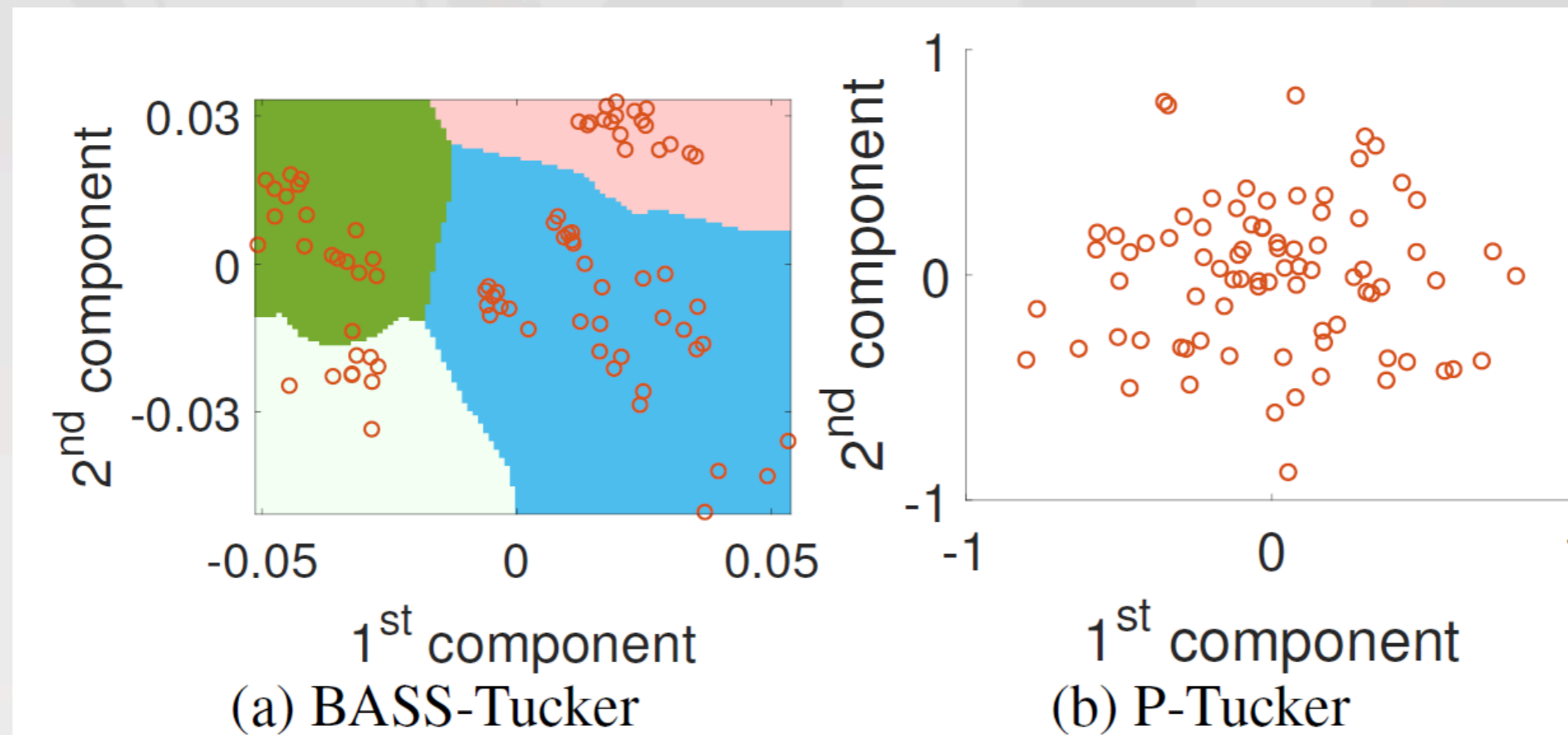
Our method!





5. Experiments on Real-world Data

- More significant Core tensor structure





THE UNIVERSITY OF UTAH

Thanks for attention Q&A Time

Authors' email: shikai.fang@utah.edu, {kirby, zhe}@cs.utah.edu

Focus: Probabilistic model, Bayesian machine learning and its application