



THE UNIVERSITY OF UTAH

Bayesian Tensor Decomposition: Dynamics and Sparsity

Presenter: 方楷楷 Shikai Fang

School of computing, The University of Utah

For ZJU talk



Outline

- 1. Background**
2. Tensor learning via Bayesian Inference
3. Dynamics in Tensor (ICML 2022 oral paper)
4. Sparsity in Tensor (UAI & ICML 2021 paper)



THE UNIVERSITY OF UTAH

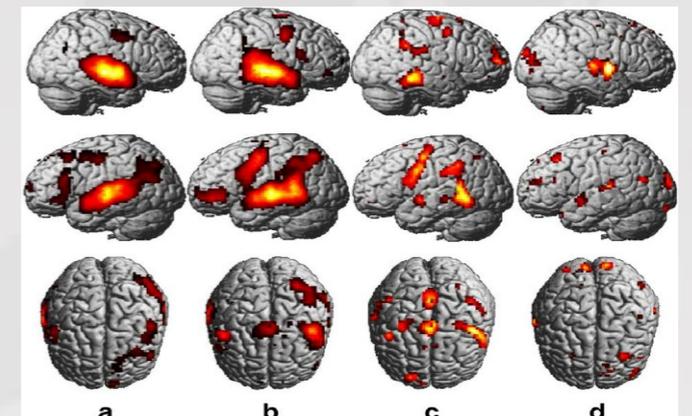
Tensor Data: Widely Used High-Order Data Structures to Represent **Interactions** of **Multiple Objects/Entities**



(user, movie, episode)



(user, advertisement, page-section)



(subject, voxel, electrode)



(location, region, time, climate)



(user, location, message-type)

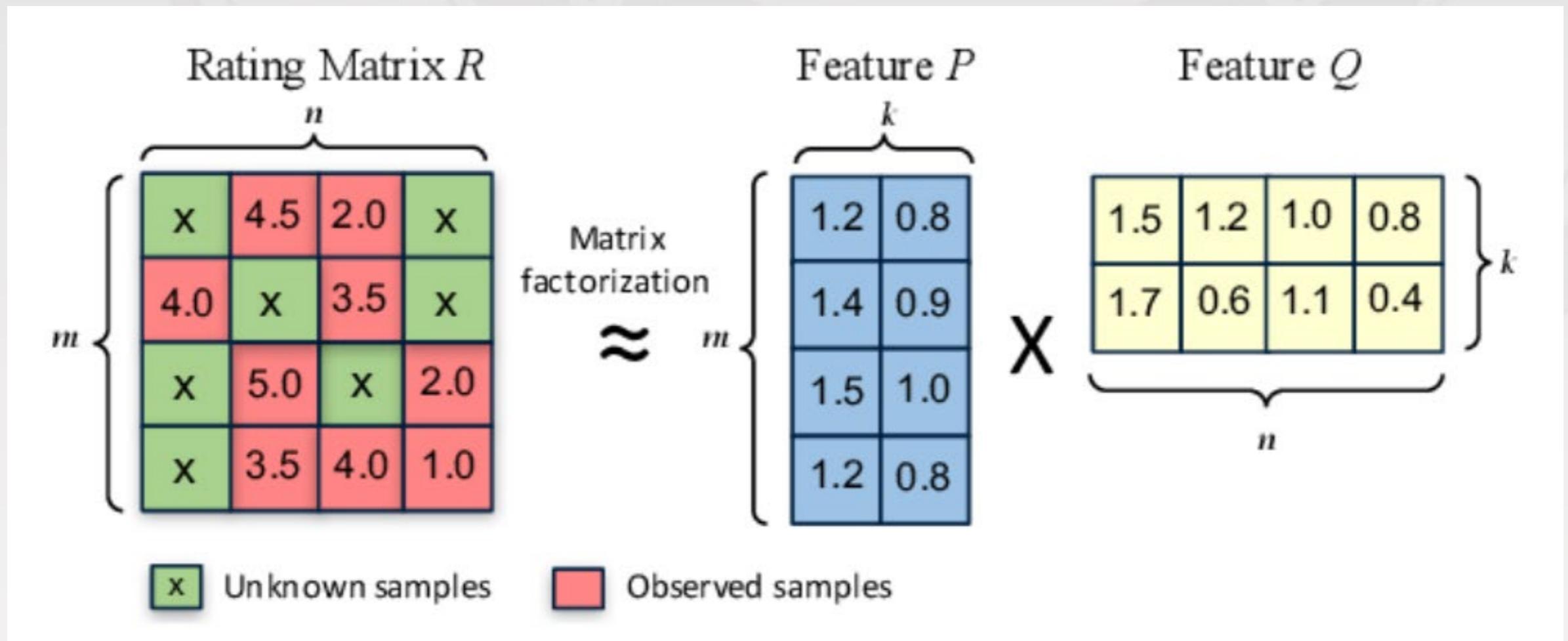


(patient, gene, condition)



Tensor Decomposition: estimate latent factors to reconstruct tensor with observed entries

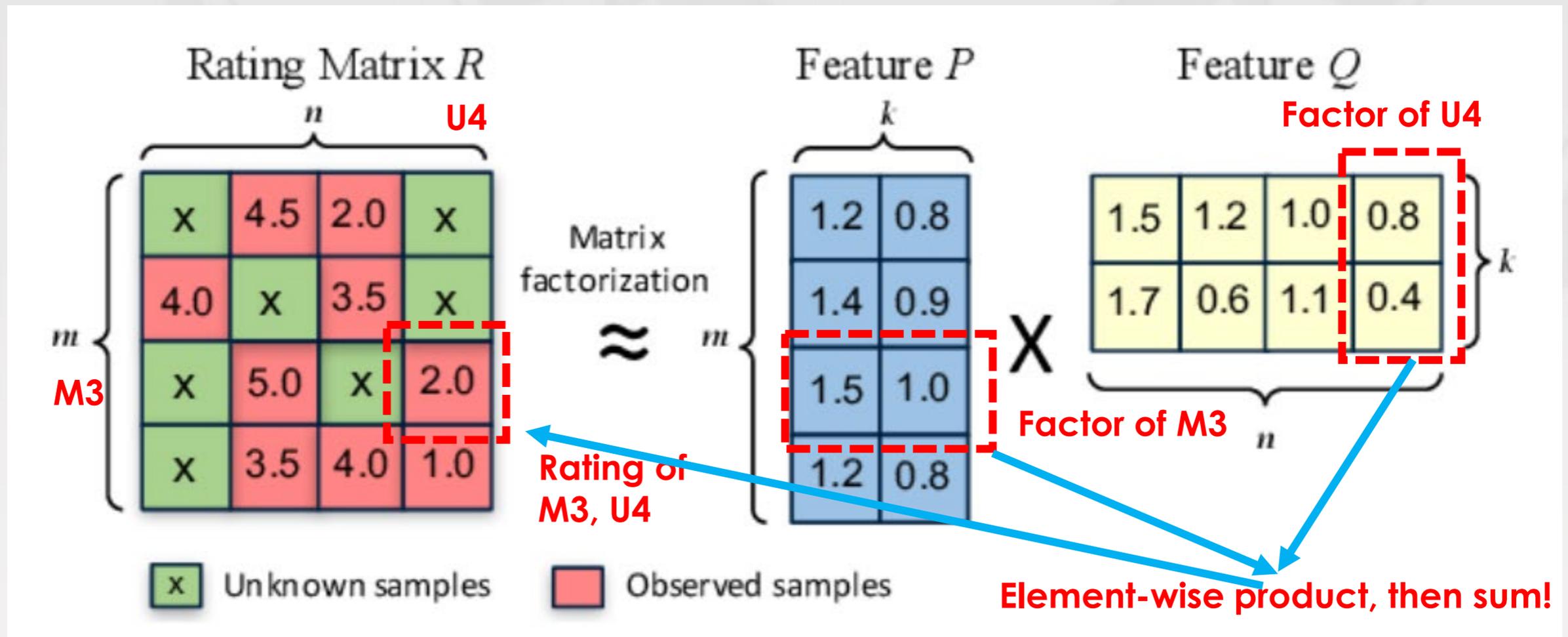
- Simple case:
Collaborative Filtering (Matrix Factorization)





Tensor Decomposition: estimate latent factors to reconstruct tensor with observed entries

- Simple case:
Collaborative Filtering (Matrix Factorization)





From matrix to tensor: CP decomposition

- Example of three-mode tensor:

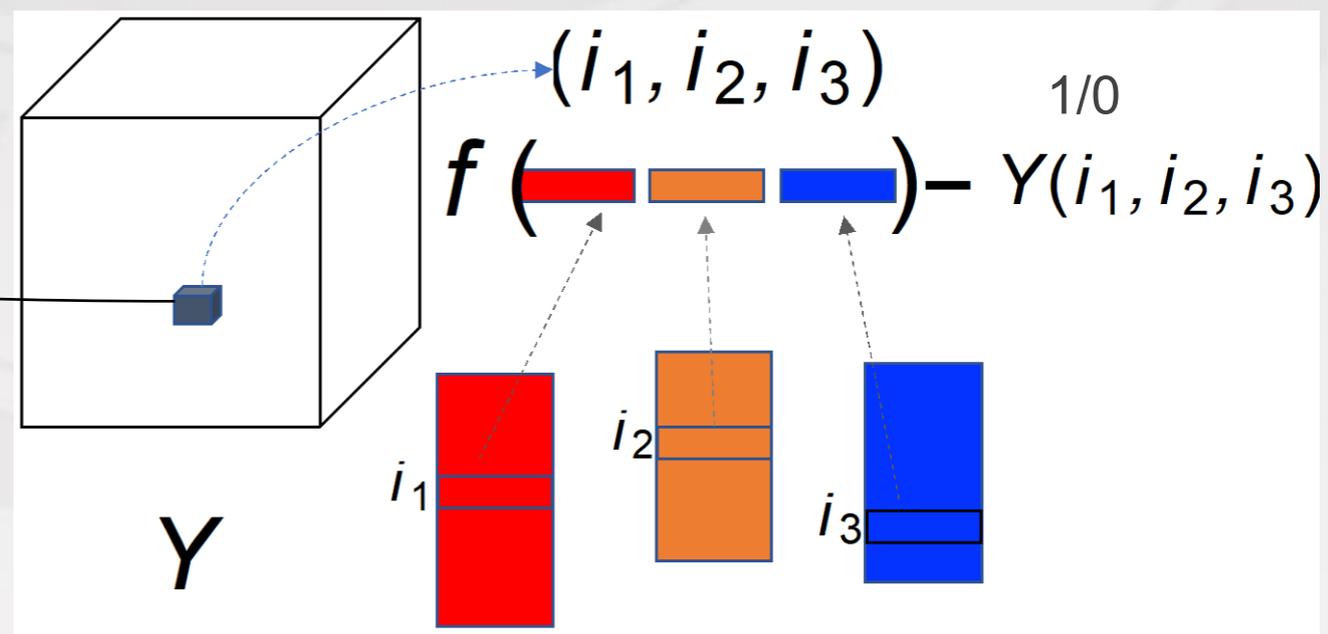
$$\mathcal{X} \in \mathbb{R}^{I \times J \times K}, A \in \mathbb{R}^{I \times R}, B \in \mathbb{R}^{J \times R}, C \in \mathbb{R}^{K \times R}$$

$$x_{ijk} \approx \sum_{r=1}^R A_{ir} B_{jr} C_{kr} \text{ for } i = 1, \dots, I, j = 1, \dots, J, k = 1, \dots, K$$

Element-wise product, then sum!

Interaction Records

user	item	page	purchase
100	25	35	1
23	21	56	0
100	25	35	1
..
32	33	46	0





THE UNIVERSITY OF UTAH

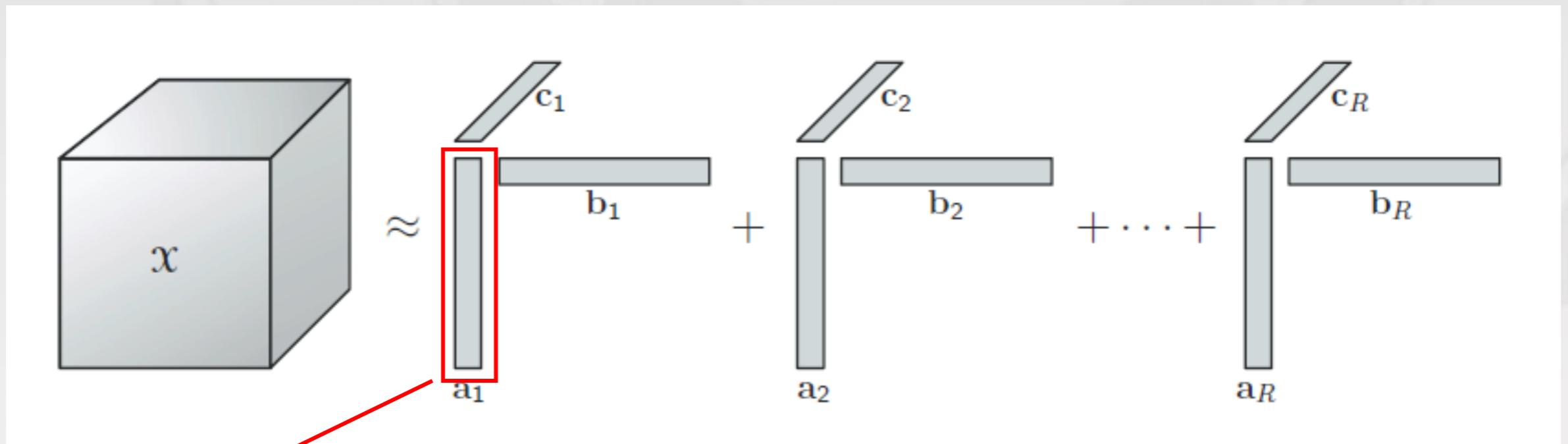
Outline

1. Background
- 2. Tensor learning via Bayesian Inference**
3. Dynamics in Tensor (ICML 2022 oral paper)
4. Sparsity in Tensor (UAI & ICML 2021 paper)



- Make latent factor from variable to a **distribution**
- Turns Deterministic model -> **Bayesian model**

Uncertainty really counts!



Deterministic version:

$$a_k \in \mathbb{R}^{d_k}$$

Probabilistic version:

$$a_k \sim \mathcal{N}(\mathbf{u}_s^k \mid \mathbf{u}_s^k, v\mathbf{I}), \mathbf{u}_s^k \in \mathbb{R}^{d_k}$$



- Bayesian version for CP decomposition,
- the **joint distribution** is:

Observed data, factors, noise Prior of the noise Gauss Prior of latent factors

$$p(\{y_i\}_{i \in S}, \mathcal{U}, \tau) = \text{Gam}(\tau \mid a_0, b_0) \prod_{k=1}^K \prod_{s=1}^{d_k} \mathcal{N}(\mathbf{u}_s^k \mid \mathbf{m}_s^k, v\mathbf{I})$$
$$\prod_{i \in S} \mathcal{N}(y_i \mid \mathbf{1}^\top (\mathbf{u}_{i_1}^1 \circ \dots \circ \mathbf{u}_{i_K}^K), \tau^{-1})$$

Gauss likelihood of the prediction

Final goal : Get the **exact posterior distribution** :

$$p(\mathcal{U}, \tau \mid \{y_i\}_{i \in S}) = \frac{p(\{y_i\}_{i \in S}, \mathcal{U}, \tau)}{p(\{y_i\}_{i \in S})}$$

Distribution of data, Constant, we never know!!



- Exact posterior distribution p is always intractable
- Approximation it by a **tractable distribution: q**

$$p(\mathcal{U}, \tau | \{y_i\}_{i \in S}) \approx q(\mathcal{U}, \tau) = q(\tau) \prod_{k=1}^K \prod_{s=1}^{d_k} q(\mathbf{u}_s^k)$$
$$= \text{Gamma}(\tau | a^*, b^*) \prod_{k=1}^K \prod_{s=1}^{d_k} \mathcal{N}(\mathbf{u}_s^k | \boldsymbol{\mu}_s^{k*}, \boldsymbol{\Sigma}_s^{k*})$$

**Parameters needed to be optimized
to make the approximation accurate !!**

Become a problem of **distribution match !**



- Key point in Bayesian machine learning

$$p(\{y_i\}_{i \in S}, \mathcal{U}, \tau) = \text{Gam}(\tau \mid a_0, b_0) \prod_{k=1}^K \prod_{s=1}^{d_k} \mathcal{N}(\mathbf{u}_s^k \mid \mathbf{m}_s^k, v\mathbf{I})$$

$$\prod_{i \in S} \mathcal{N}(y_i \mid f(\mathbf{u}_{i_1}^1, \dots, \mathbf{u}_{i_K}^K), \tau^{-1})$$

$$p(\mathcal{U}, \tau \mid \{y_i\}_{i \in S}) \approx q(\mathcal{U}, \tau) = q(\tau) \prod_{k=1}^K \prod_{s=1}^{d_k} q(\mathbf{u}_s^k)$$

$$= \text{Gamma}(\tau \mid a^*, b^*) \prod_{k=1}^K \prod_{s=1}^{d_k} \mathcal{N}(\mathbf{u}_s^k \mid \boldsymbol{\mu}_s^{k*}, \boldsymbol{\Sigma}_s^{k*})$$

- Complete toolbox, but no silver bullet.
- Depend on the model **f** picked for the likelihood



Most-frequently used tool:

Variational inference(VI): minimize the **KL divergence** of two distributions

$$\mathcal{L} = \int q^*(\mathcal{U}, \tau) \log \frac{p(\{y_i\}_{i \in S_t} | \mathcal{U}, \tau) q(\mathcal{U}, \tau)}{q^*(\mathcal{U}, \tau)} d\mathcal{U} d\tau$$

- **Expectation propagation(EP)**:
moment match when L (reverse order) has **closed form solution**
- **Assumed density filtering(ADF)**:
moment match when having **tractable normalization term**
- **Reparameterization trick (SVI)**:
Probabilistic version SGD when L is **totally intractable**
...(sampling based methods)



Series of work from our group under such architecture.

Short Names	Likelihood Model & Prior	Inference Method	Publication
POST[1]	CP decomposition/multi-linear	EP	ICDM 2018
POND[2]	Neural-kernel Gaussian process	SVI	ICDM 2020
SBTD[3]	Bayesian neural network + Sparse	ADF (Streaming)	ICML 2021
BASS[4]	Tucker decomposition + Sparse	ADF (Streaming)	UAI 2021
BCTT[5]	Tucker decomposition + Dynamics	CEP + KF/RTS	ICML 2022

[1] Du, Yishuai, et al. "Probabilistic streaming tensor decomposition." *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2018.

[2] Conor Tillinghast, **Shikai Fang**, Kai Zheng, and Shandian Zhe, "Probabilistic Neural-Kernel Tensor Decomposition", *IEEE International Conference on Data Mining (ICDM)*, 2020.

[3] **Fang, Shikai**, et al. "Streaming Probabilistic Deep Tensor Factorization." *The Thirty-eighth International Conference on Machine Learning (ICML)*, 2021

[4] **Fang, Shikai**, et al. "Bayesian Streaming Sparse Tucker decomposition." *Conference on Uncertainty in Artificial Intelligence*. UAI, 2021.

[5] **Fang, Shikai**, et al. "Bayesian Continuous-Time Tucker Decomposition." *The Thirty-ninth International Conference on Machine Learning (ICML)*, 2022



THE UNIVERSITY OF UTAH

Outline

1. Background
2. Tensor learning via Bayesian Inference
- 3. Dynamics in Tensor (ICML 2022 oral paper)**
4. Sparsity in Tensor (UAI & ICML 2021 paper)

Shikai Fang, Akil Narayan, Robert M. Kirby, and Shandian Zhe, “Bayesian Continuous-Time Tucker Decomposition“ (Oral), The 39th International Conference on Machine Learning (ICML), 2022

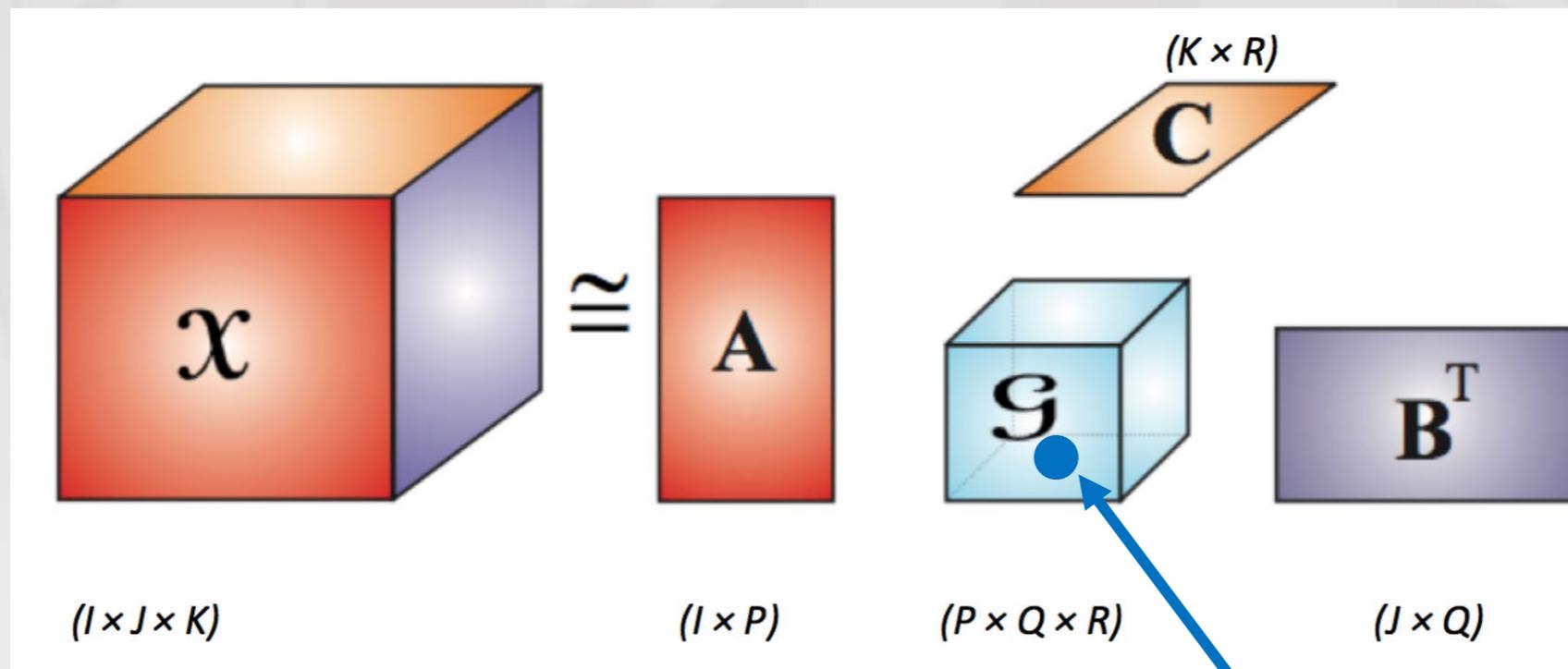


Tucker Decomposition

- 2-D matrix \Rightarrow N-D tensor
- Element-wise interaction \Rightarrow all possible interactions



Tucker Decomposition



One interaction weight

One Interaction of latent factors

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

$$y_i \approx \text{vec}(\mathcal{W})^T \left(\mathbf{u}_{i_1}^1 \otimes \dots \otimes \mathbf{u}_{i_K}^K \right)$$

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

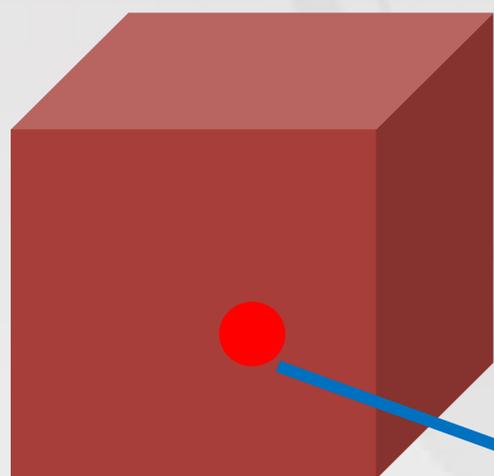


Challenge: Temporal info in Tensor

What about each entry is time-dependent?

Straightforward Solution:

- Drop time or
- Augment tensor with **time-step mode**

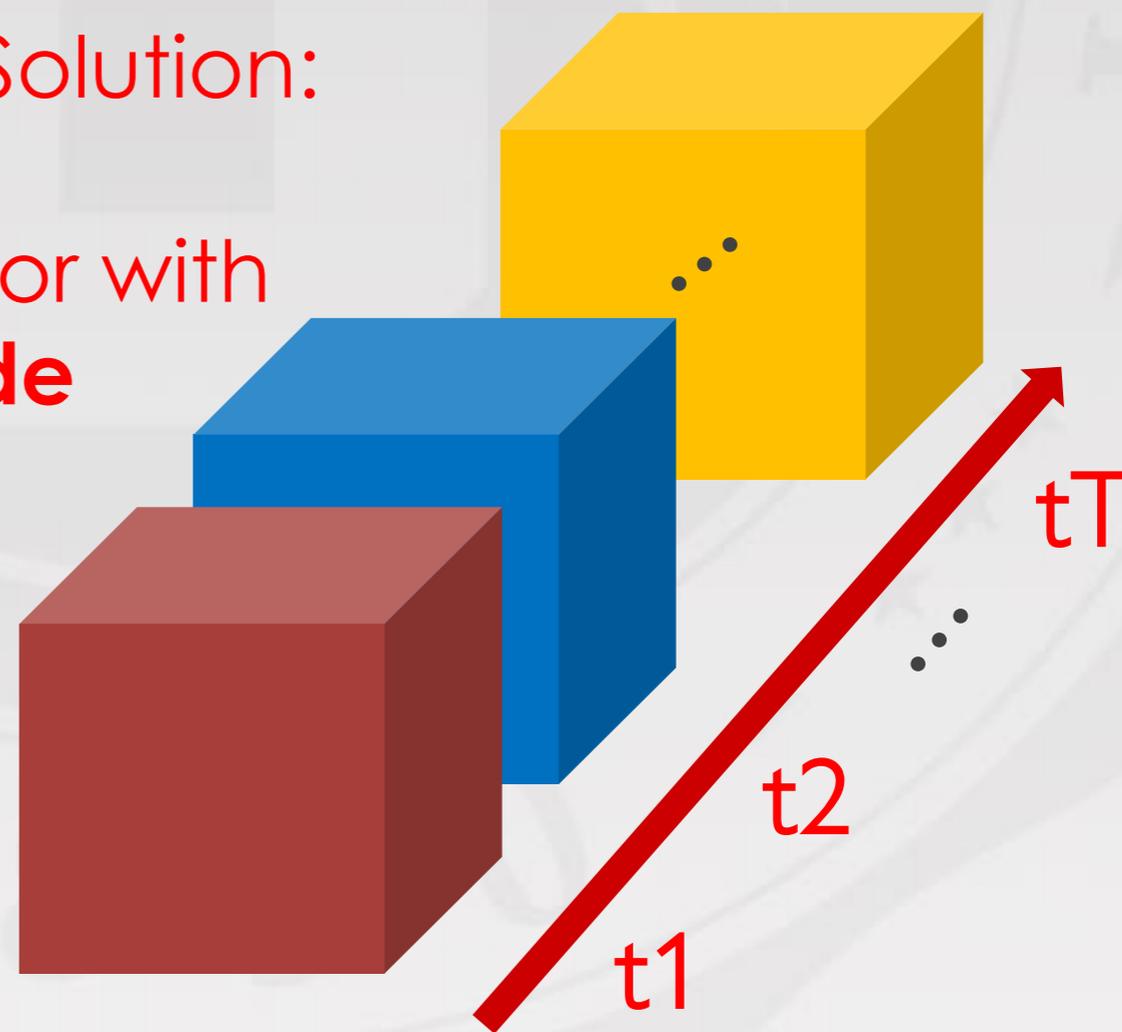


$$X_{ijk}(t)$$

$$(I \times J \times K)$$

Problem:

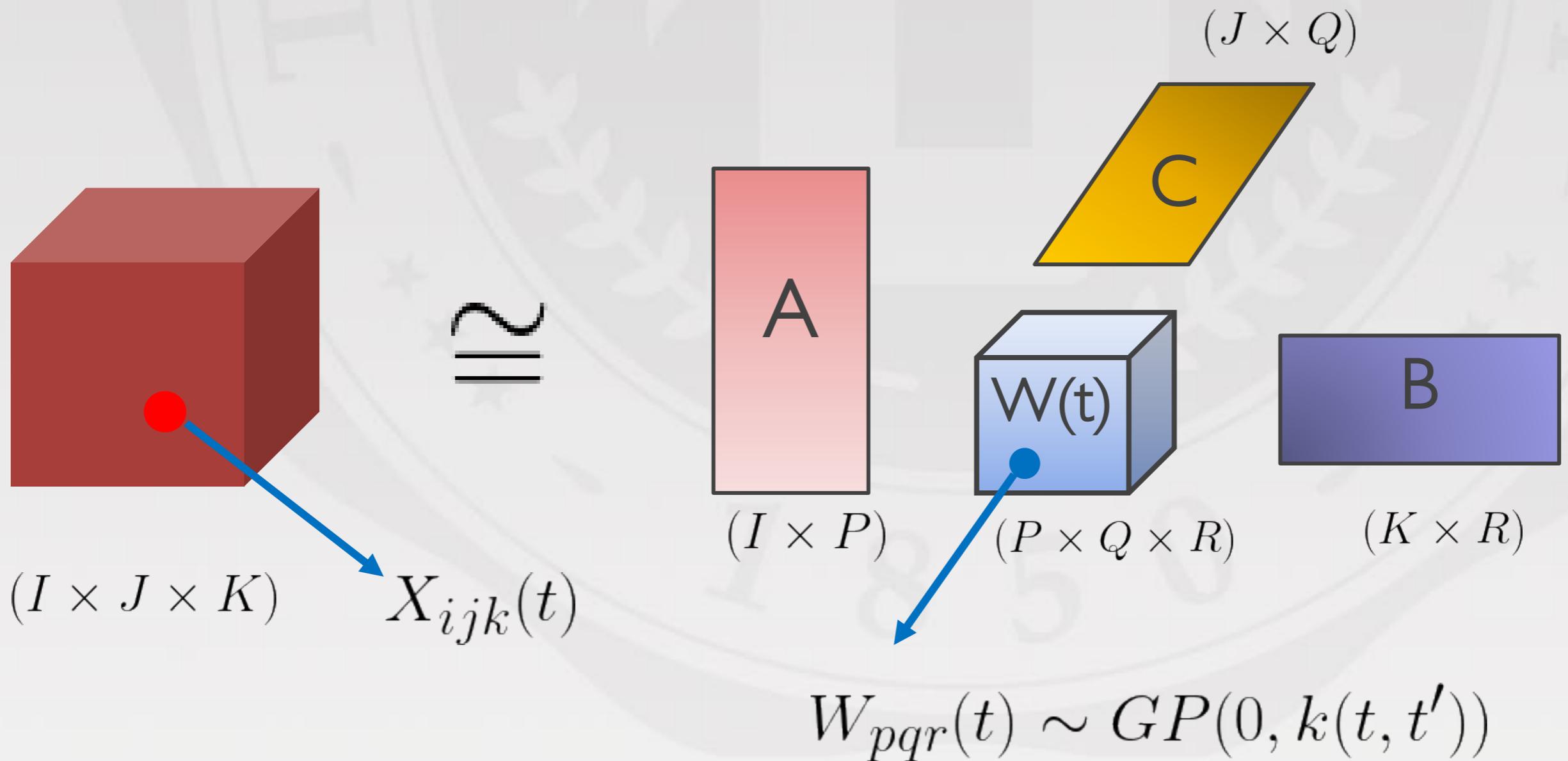
1. Too Sparse
2. Ignore the temporary continuity



$$(I \times J \times K \times T)$$



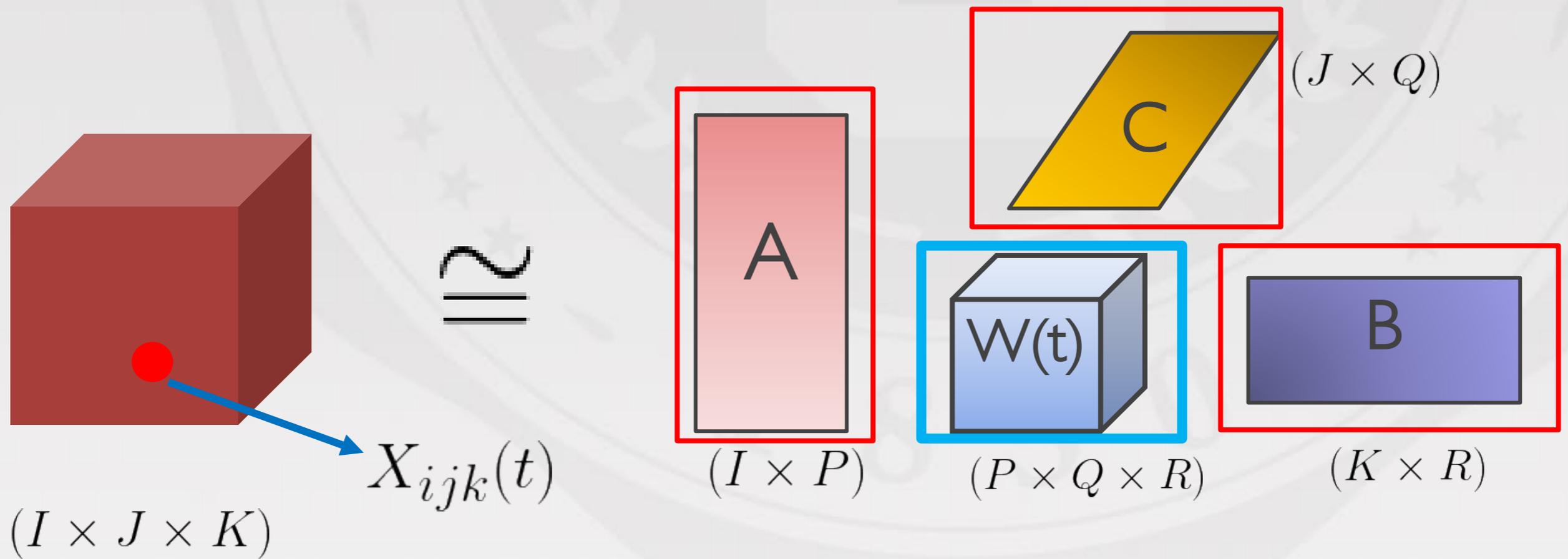
Our Solution: Modeling Dynamic Tucker Core by Temporal Gaussian Processes





High-level Motivation:

Decouple the **representation learning of factors** and the **capture of dynamic pattern**





Joint Probability:

$$p(\mathcal{U}, \{\mathbf{w}_r\}_r, \tau, \mathbf{y}) =$$

$$\text{Gam}(\tau \mid b_0, c_0) \prod_{k=1}^K \prod_{j=1}^{d_k} \mathcal{N}(\mathbf{u}_j^k \mid \mathbf{0}, \mathbf{I})$$

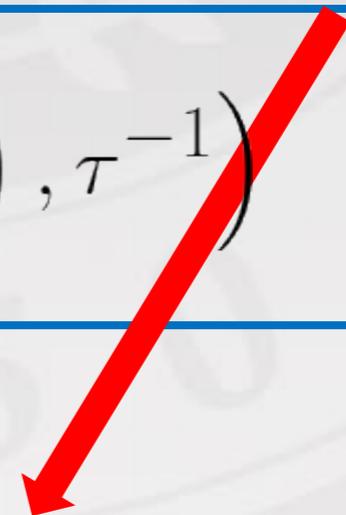
Priors of factors and noise

$$\times \prod_{\mathbf{r}=(1,\dots,1)}^{(R_1,\dots,R_K)} \mathcal{N}(\mathbf{w}_r \mid \mathbf{0}, \mathbf{K}_r)$$

Temporal GPs on Tucker Core

$$\prod_{n=1}^N \mathcal{N}(y_n \mid \text{vec}(\mathcal{W}(t_n))^\top (\mathbf{u}_{i_{n_1}}^1 \otimes \dots \otimes \mathbf{u}_{i_{n_K}}^K), \tau^{-1})$$

Gaussian Likelihood



Computational challenge: $\mathcal{O}(N^3)$ cost of full GPs



THE UNIVERSITY OF UTAH

To **avoid low-rank/sparse approx.** (low quality),
but enjoy **linear-cost inference of full GPs,**

We apply a crucial fact:

Temporal GPs



with stationary kernel

Linear Time-Invariant(LTI) SDE



discrete form on $(t_1, t_2 \dots t_N)$

State Space Model (Gauss Markov Chain)

Can be solved by
**Kalman filtering &
RTS Smoothing in $O(N)^*$**

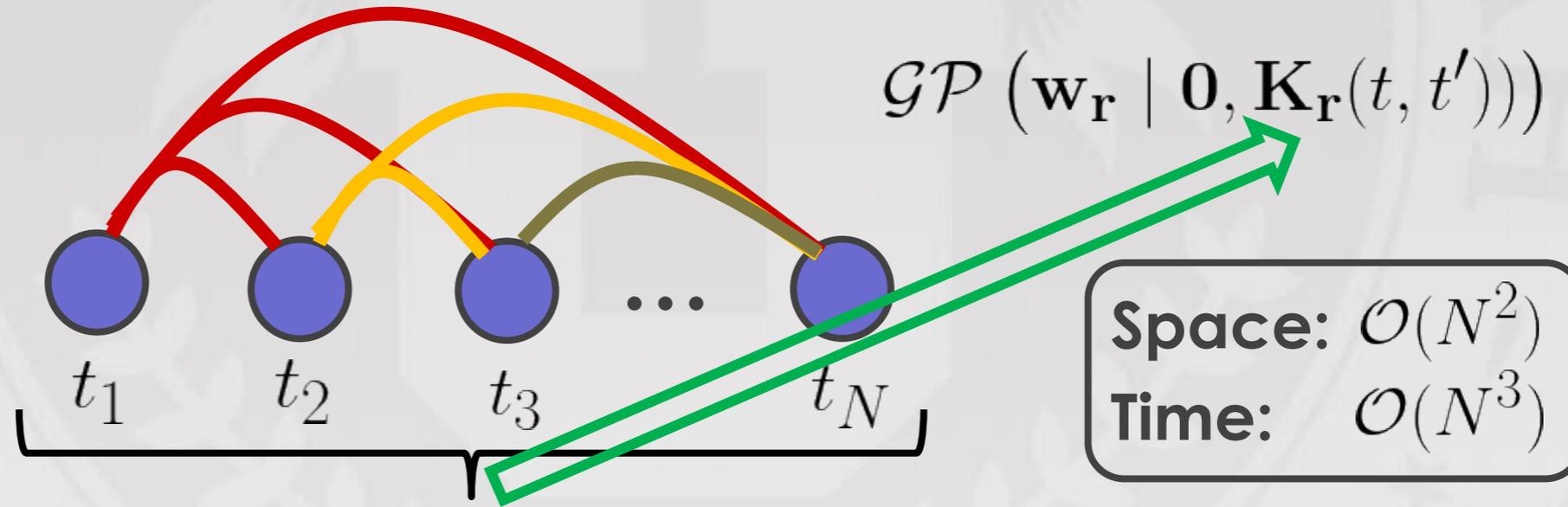


*: it holds with linear emission/observed likelihood, if with non-linear, we could apply non-linear filter and smoothing



Illustration of computation cost:

Temporal GPs



Temporal States:



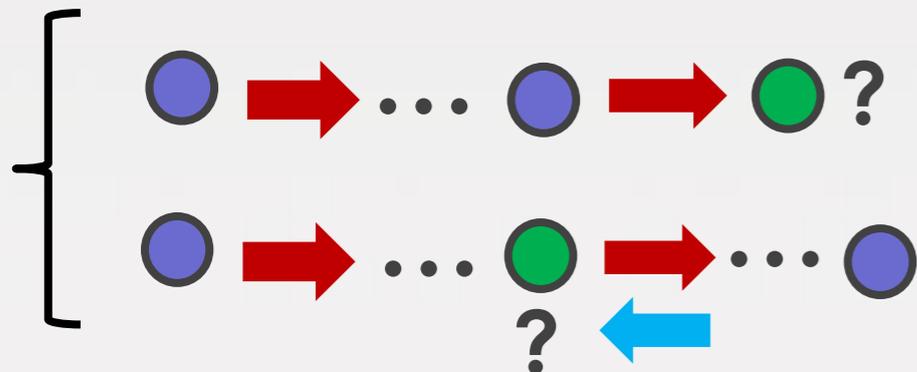
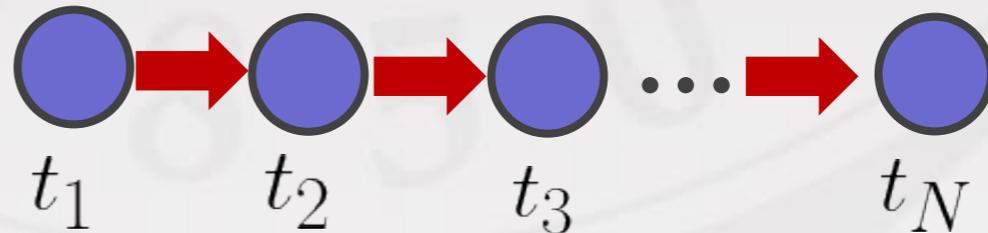
LTI-SDE

$$\frac{d\gamma_{\mathbf{r}}(t)}{dt} = \mathbf{F}\gamma_{\mathbf{r}} + \mathbf{L}\xi(t)$$

discrete form

State Space Model

$$p(\gamma_{\mathbf{r}}(t_{n+1}) \mid \gamma_{\mathbf{r}}(t_n)) = \mathcal{N}(\gamma_{\mathbf{r}}(t_{n+1}) \mid \mathbf{A}_n \gamma_{\mathbf{r}}(t_n), \mathbf{Q}_n)$$



Kalman Filter

RTS Smoothing

Space: $\mathcal{O}(N)$
Time: $\mathcal{O}(N)$



Specifically,

Temporal GPs

$$\prod_{\mathbf{r}=(1,\dots,1)}^{(R_1,\dots,R_K)} \mathcal{N}(\mathbf{w}_{\mathbf{r}} \mid \mathbf{0}, \mathbf{K}_{\mathbf{r}})$$



with stationary kernel
(e.g., Matern25)

$$\begin{cases} \gamma_{\mathbf{r}}(t) = \left(w_{\mathbf{r}}, \frac{dw_{\mathbf{r}}}{dt} \right)^{\top} \\ \frac{d\gamma_{\mathbf{r}}(t)}{dt} = \mathbf{F}\gamma_{\mathbf{r}} + \mathbf{L}\xi(t) \end{cases}$$

Linear Time-Invariant(LTI) SDE



discrete form on $(t_1, t_2 \dots t_N)$

$$\begin{cases} p(\gamma_{\mathbf{r}}(t_1)) = \mathcal{N}(\gamma_{\mathbf{r}}(t_1) \mid \mathbf{0}, \mathbf{P}_{\infty}) \\ p(\gamma_{\mathbf{r}}(t_{n+1}) \mid \gamma_{\mathbf{r}}(t_n)) = \mathcal{N}(\gamma_{\mathbf{r}}(t_{n+1}) \mid \mathbf{A}_n\gamma_{\mathbf{r}}(t_n), \mathbf{Q}_n) \end{cases}$$

State Space Model (Gauss Markov Chain)



Recall: linear-cost solver - KF, RTS



Reformulate **Tucker core** with **State Space Priors**

$$p(\bar{\gamma}_1) \prod_{n=1}^{N-1} p(\bar{\gamma}_{n+1} \mid \bar{\gamma}_n)$$

We post **Gaussian-Gamma Approx.** to fit each data-llk

$$\mathcal{N}\left(y_n \mid (\mathbf{H}\bar{\gamma}_n)^\top \left(\mathbf{u}_{i_{n1}}^1 \otimes \dots \otimes \mathbf{u}_{i_{nK}}^K\right), \tau^{-1}\right) \approx$$

$$Z_n \prod_{k=1}^K \mathcal{N}\left(\mathbf{u}_{i_{nk}}^k \mid \mathbf{m}_{i_{nk}}^{k,n}, \mathbf{V}_{i_{nk}}^{k,n}\right) \cdot \text{Gam}(\tau \mid b_n, c_n)$$

Approx. Msg of Factors & noise

$$\times \mathcal{N}(\mathbf{H}\bar{\gamma}_n \mid \boldsymbol{\beta}_n, \mathbf{S}_n)$$

Approx. Msg of SDE states /Tucker core

Substitute these into joint prob.



The proposed approx. posterior is:

$$q(\mathcal{U}, \{\bar{\gamma}_n\}, \tau) \propto \prod_{k=1}^K \prod_{j=1}^{d_k} \mathcal{N}(\mathbf{u}_j^k \mid \mathbf{0}, \mathbf{I}) \text{Gam}(\tau \mid b_0, c_0)$$

Standard moment match? Infeasible!

$$\prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(\mathbf{u}_{i_{n_k}}^k \mid \mathbf{m}_{i_{n_k}}^{k,n}, \mathbf{V}_{i_{n_k}}^{k,n}) \text{Gam}(\tau \mid b_n, c_n)$$

$$p(\bar{\gamma}_1) \mathcal{N}(\mathbf{H}\bar{\gamma}_1 \mid \boldsymbol{\beta}_1, \mathbf{S}_1) \prod_{n=1}^{N-1} p(\bar{\gamma}_{n+1} \mid \bar{\gamma}_n) \mathcal{N}(\mathbf{H}\bar{\gamma}_n \mid \boldsymbol{\beta}_n, \mathbf{S}_n)$$

SDE states: Solve by KF and RTS

Apply conditional moment matching and delta method!



- **Conditional Moment Match**

$$\mathbb{E}_{\tilde{p}}[\phi(\boldsymbol{\eta}_n)] = \mathbb{E}_{\tilde{p}(\Theta \setminus \eta_n)} \left[\mathbb{E}_{\tilde{p}(\boldsymbol{\eta}_n | \Theta \setminus \eta_n)} [\phi(\boldsymbol{\eta}) | \Theta \setminus \eta_n] \right]$$

- **Delta method:**

$$\mathbb{E}_q(\Theta \setminus \eta_n) [\boldsymbol{\rho}_n] \approx \rho_n \left(\mathbb{E}_q [\Theta \setminus \eta_n] \right)$$

Enable **tractable moment matching**

to update **approx. probability terms** under
Expectation Propagation (EP) framework



Algorithm 1 BCTT

Input: $\mathcal{D} = \{(\mathbf{i}_1, t_1, y_1), \dots, (\mathbf{i}_N, t_N, y_N)\}$, kernel hyper-parameters l, σ^2

Initialize approximation terms in (10) for each likelihood.

repeat

Run KF and RTS smoothing to compute each $q(\bar{\gamma}_n)$

for $n = 1$ **to** N **in parallel do**

Simultaneously update $\mathcal{N}(\mathbf{H}\bar{\gamma}_n | \beta_n, \mathbf{S}_n)$,
Gam($\tau | b_n, c_n$) and $\left\{ \mathcal{N} \left(\mathbf{u}_{i_{n_k}}^k | \mathbf{m}_{i_{n_k}}^{k,n}, \mathbf{V}_{i_{n_k}}^{k,n} \right) \right\}_k$
in (10) with conditional moment matching and multi-variate delta method.

end for

until Convergence

Return: $\{q(\mathcal{W}(t_n))\}_{n=1}^N, \{q(\mathbf{u}_j^k)\}_{1 \leq k \leq K, 1 \leq j \leq d_k}, q(\tau)$

Time cost: $\mathcal{O}(N\bar{R})$ **Space cost:** $\mathcal{O} \left(N \left(\bar{R}^2 + \sum_{k=1}^K R_k^2 \right) \right)$



➡ Kalman Filtering (forward) ➡ RTS Smoothing (backward) ➡ Conditional Moment Matching (parallel)

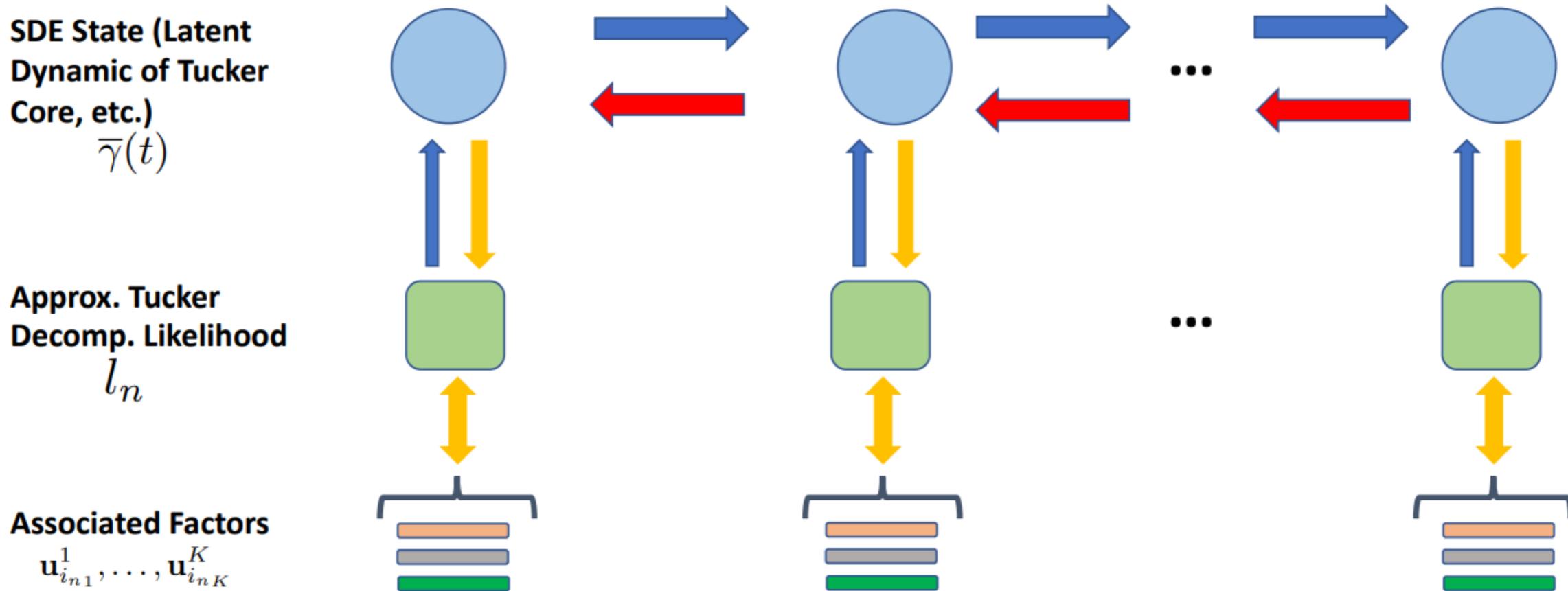
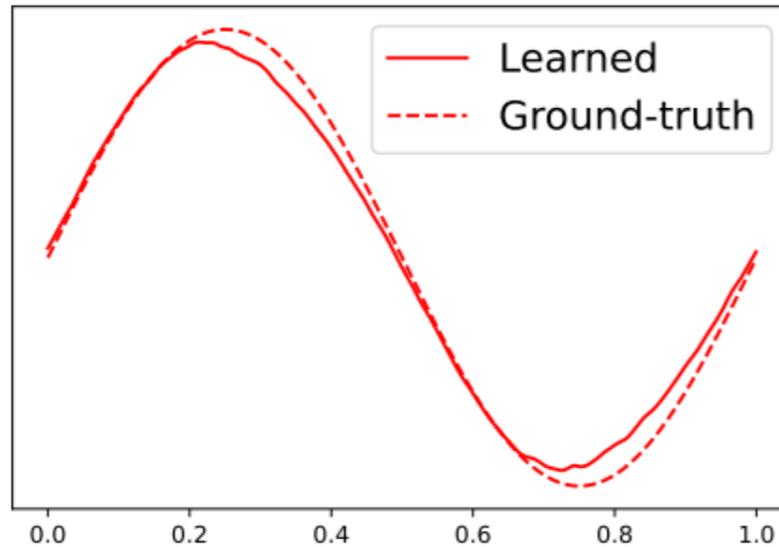


Figure 1. Graphical illustration of the message-passing inference algorithm.

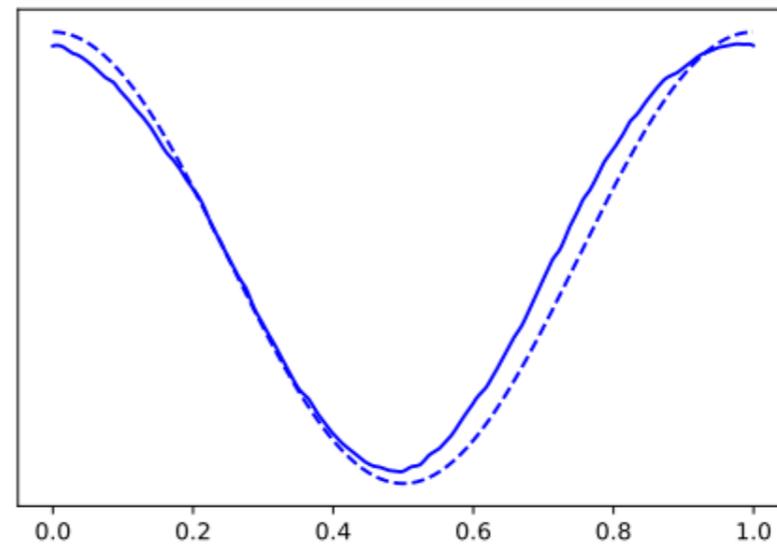


Can BCTT capture the temporal patterns in tensor?

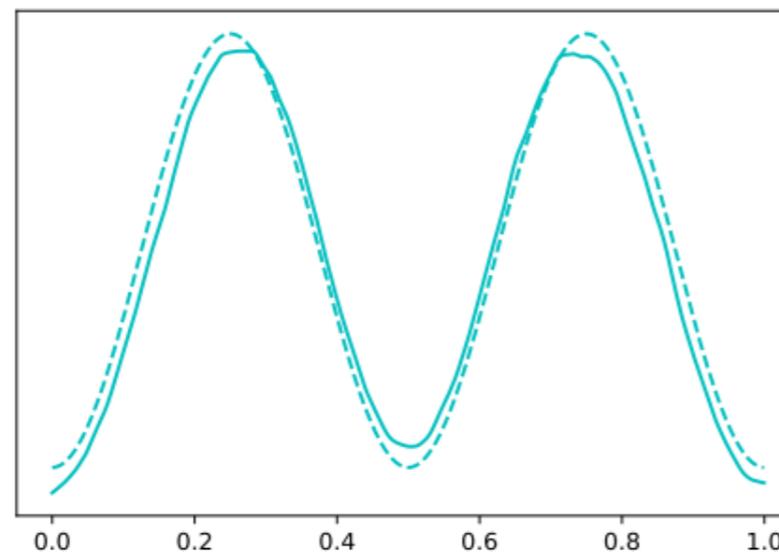
- Exp on **simulation data**
- Plot the dynamics of Tucker core



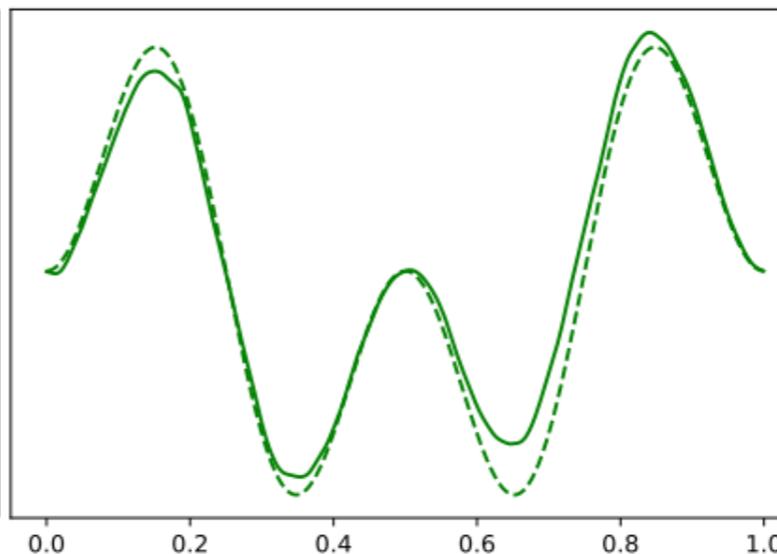
(a) $w_{(1,1)}(t)$



(b) $w_{(1,2)}(t)$



(c) $w_{(2,1)}(t)$



(d) $w_{(2,2)}(t)$



Can BCTT capture the temporal patterns in tensor?

- Exp on **real-world data (DBLP dataset)**
- Scatter low-rank structures of Tucker core

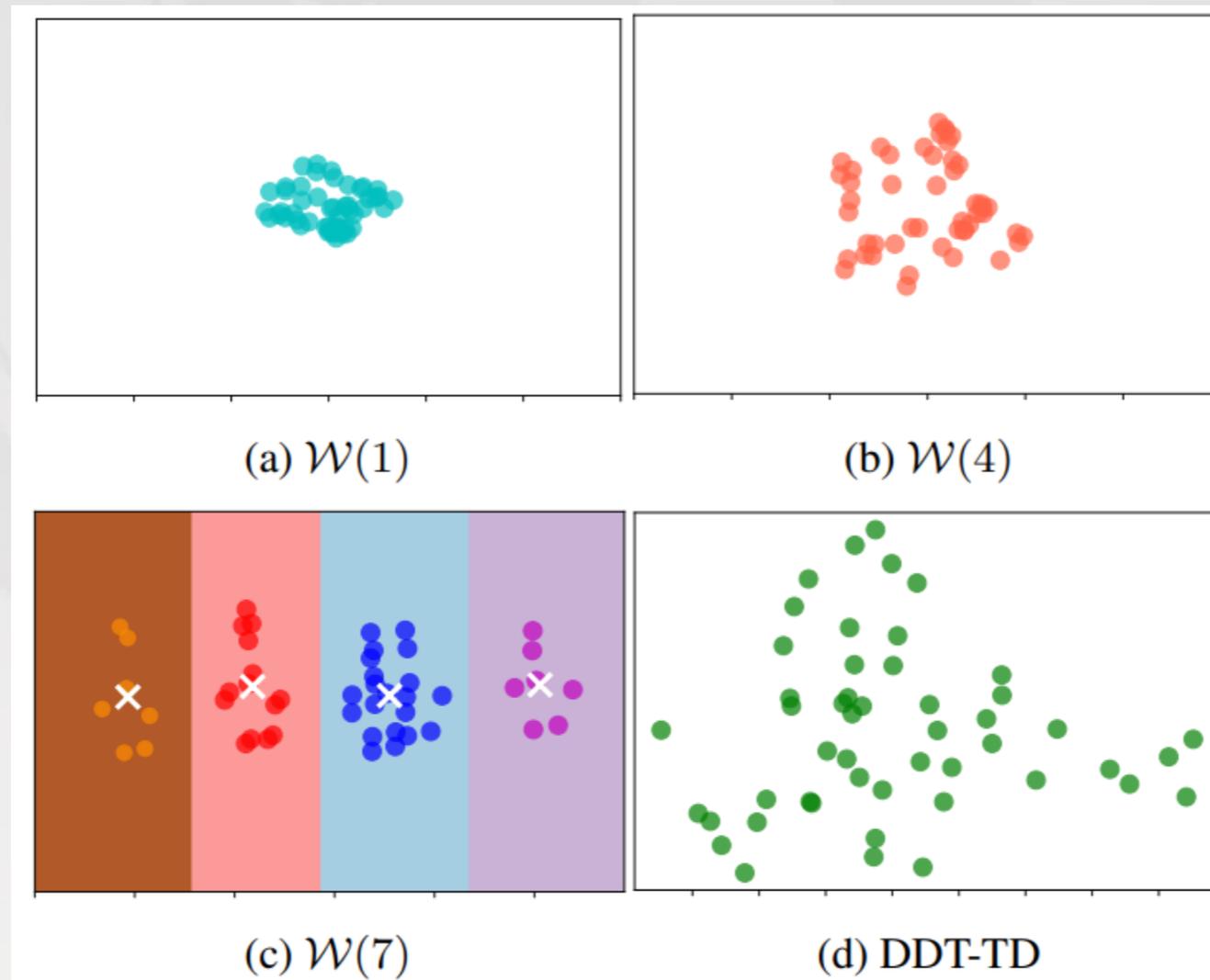


Figure 4. The structures of learned tensor-core at different time points by BCTT (a-c) and the static tensor-score learned by dynamic discrete-time Tucker decomposition (DDT-TD).



Prediction with BCTT

- Prediction performance of BCTT on 3 real-world data

RMSE	<i>MovieLens</i>	<i>AdsClicks</i>	<i>DBLP</i>
CT-CP	1.113 ± 0.004	1.337 ± 0.013	0.240 ± 0.007
CT-GP	0.949 ± 0.008	1.422 ± 0.008	0.227 ± 0.009
DT-GP	0.963 ± 0.008	1.436 ± 0.015	0.227 ± 0.007
DDT-GP	0.957 ± 0.008	1.437 ± 0.010	0.225 ± 0.006
DDT-CP	1.022 ± 0.003	1.420 ± 0.020	0.245 ± 0.004
DDT-TD	1.059 ± 0.006	1.401 ± 0.022	0.232 ± 0.09
BCTT	0.922 ± 0.002	1.322 ± 0.012	0.214 ± 0.009

MAE

CT-CP	0.788 ± 0.004	0.787 ± 0.006	0.105 ± 0.001
CT-GP	0.714 ± 0.004	0.891 ± 0.011	0.092 ± 0.004
DT-GP	0.722 ± 0.008	0.893 ± 0.008	0.084 ± 0.003
DDT-GP	0.720 ± 0.003	0.894 ± 0.009	0.083 ± 0.001
DDT-CP	0.755 ± 0.002	0.901 ± 0.011	0.114 ± 0.002
DDT-TD	0.742 ± 0.006	0.866 ± 0.012	0.101 ± 0.001
BCTT	0.698 ± 0.002	0.777 ± 0.016	0.084 ± 0.001

(a) $R = 3$

RMSE	<i>MovieLens</i>	<i>AdsClicks</i>	<i>DBLP</i>
CT-CP	1.165 ± 0.008	1.324 ± 0.013	0.263 ± 0.006
CT-GP	0.965 ± 0.019	1.410 ± 0.015	0.227 ± 0.007
DT-GP	0.949 ± 0.007	1.425 ± 0.015	0.225 ± 0.008
DDT-GP	0.948 ± 0.005	1.421 ± 0.012	0.220 ± 0.006
DDT-CP	1.141 ± 0.007	1.623 ± 0.013	0.282 ± 0.011
DDT-TD	0.944 ± 0.003	1.453 ± 0.035	0.312 ± 0.072
BCTT	0.895 ± 0.007	1.304 ± 0.018	0.202 ± 0.009

MAE

CT-CP	0.835 ± 0.006	0.792 ± 0.007	0.128 ± 0.001
CT-GP	0.717 ± 0.012	0.883 ± 0.016	0.092 ± 0.002
DT-GP	0.714 ± 0.005	0.886 ± 0.012	0.084 ± 0.001
DDT-GP	0.707 ± 0.004	0.882 ± 0.015	0.082 ± 0.003
DDT-CP	0.843 ± 0.003	1.082 ± 0.013	0.141 ± 0.004
DDT-TD	0.712 ± 0.002	0.903 ± 0.024	0.221 ± 0.047
BCTT	0.679 ± 0.001	0.785 ± 0.010	0.080 ± 0.001

(b) $R = 7$



THE UNIVERSITY OF UTAH

Outline

1. Background
2. Tensor learning via Bayesian Inference
3. Dynamics in Tensor (ICML 2022 oral paper)
- 4. Sparsity in Tensor(UAI & ICML 2021 paper)**

Shikai Fang, Zheng Wang, Zhimeng Pan, Ji Liu, and Shandian Zhe, “Streaming Bayesian Deep Tensor Factorization”, The Thirty-eighth International Conference on Machine Learning (ICML), 2021

Shikai Fang, Robert. M. Kirby, and Shandian Zhe, “Bayesian Streaming Sparse Tucker Decomposition”, The 37th Conference on Uncertainty in Artificial Intelligence (UAI), 2021



- Sparsity in tensor data requires Sparsity in model

Tensor-Datasets	Size	#Observed entries	Observed Ratio
Gowalla	18737*1000*32510	821,931	0.0001%
SG	2321*5596*1600	105,764	0.0005%
ACC	3000*150*30000	1,220,000	0.1%
Movielens1M	6000*3700	1,000,000	4%

- otherwise, **overfitting** risk, especially for complex model like NN
- How non-Bayesian people get sparsity?
 - L1 regular terms
 - overall sparse, not accurate enough



How Bayesian people get sparsity?

- **Spike and Slab Priors**

Introduce binary selection indicators s_1, s_2, \dots, s_d
on each parameter ! -- Element-wise sparse control!

$$s_j = 1$$



Mild Regularize
(slab)

$$s_j \mathcal{N}(w_j | 0, \sigma_0^2)$$

$$s_j = 0$$



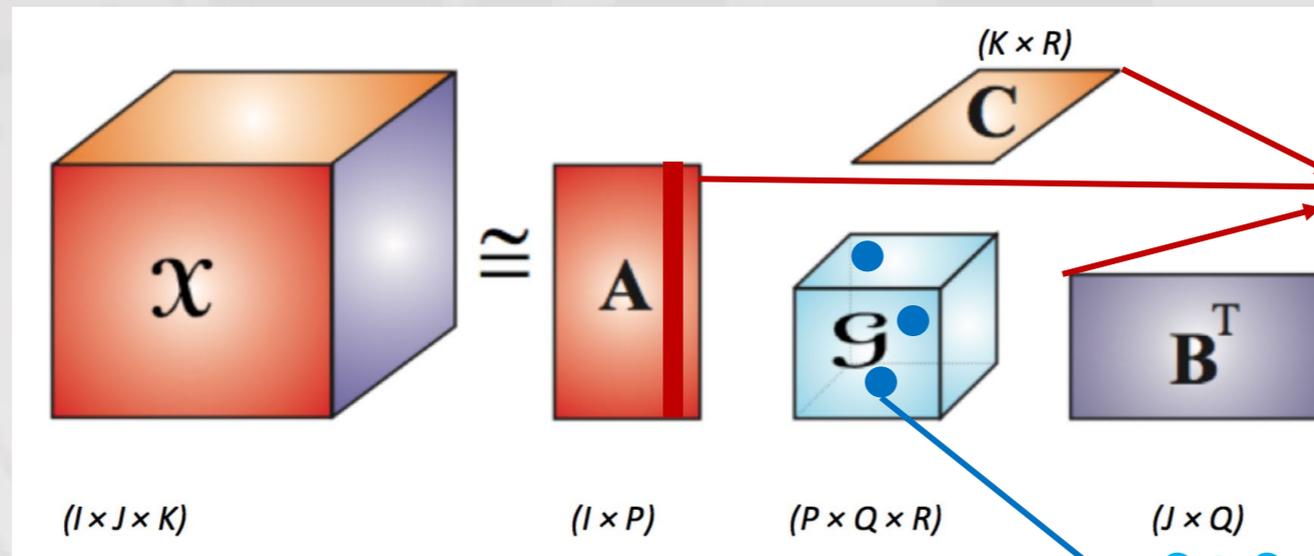
Strong shrinkage
(spike)

$$(1 - s_j) \delta(w_j)$$



Where to put such sparsity?

- **On Tucker Core– build a sparse core (BASS[4])**



Standard Gaussian prior

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

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

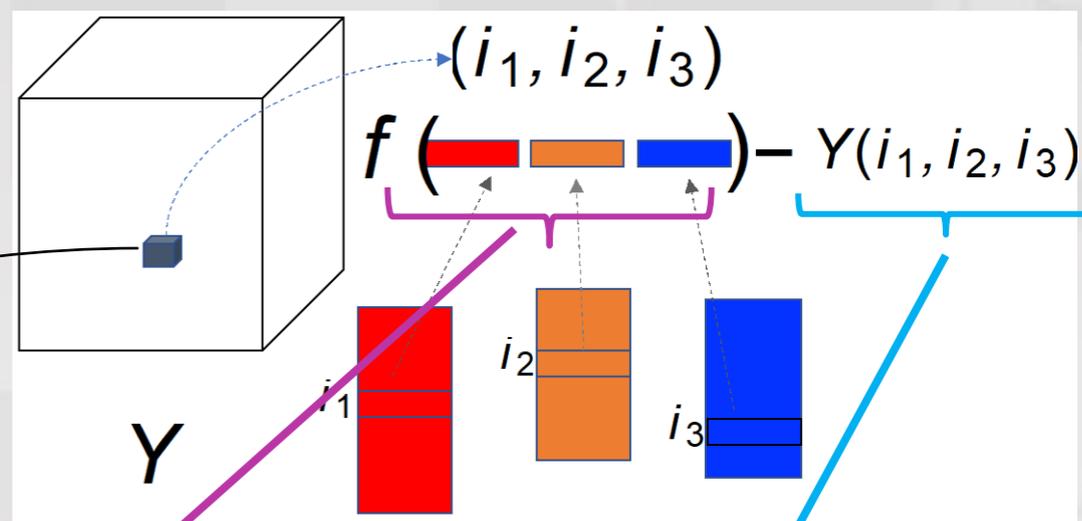


Where to put such sparsity?

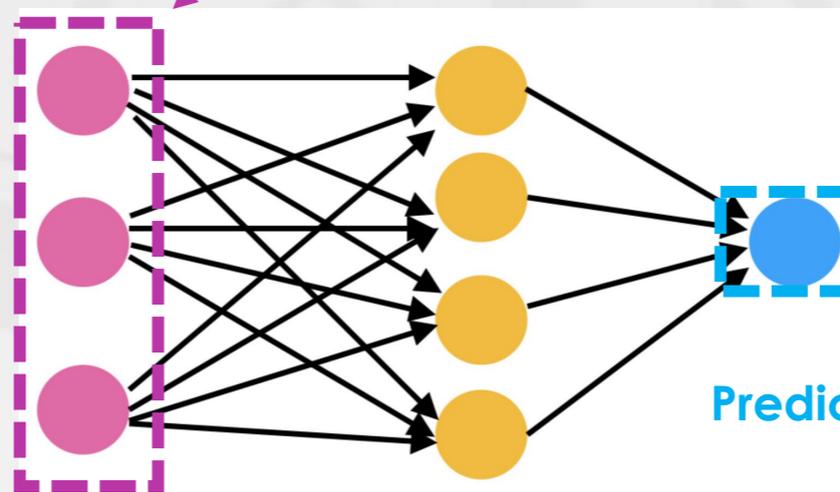
- **On NN weights – build a sparse BNN (SBDT[4])**

Interaction Records

use r	item	page	purchase
100	25	35	1
23	21	56	0
100	25	35	1
..
32	33	46	0



Latent factors at each mode as NN input



Predicted Entry as NN output

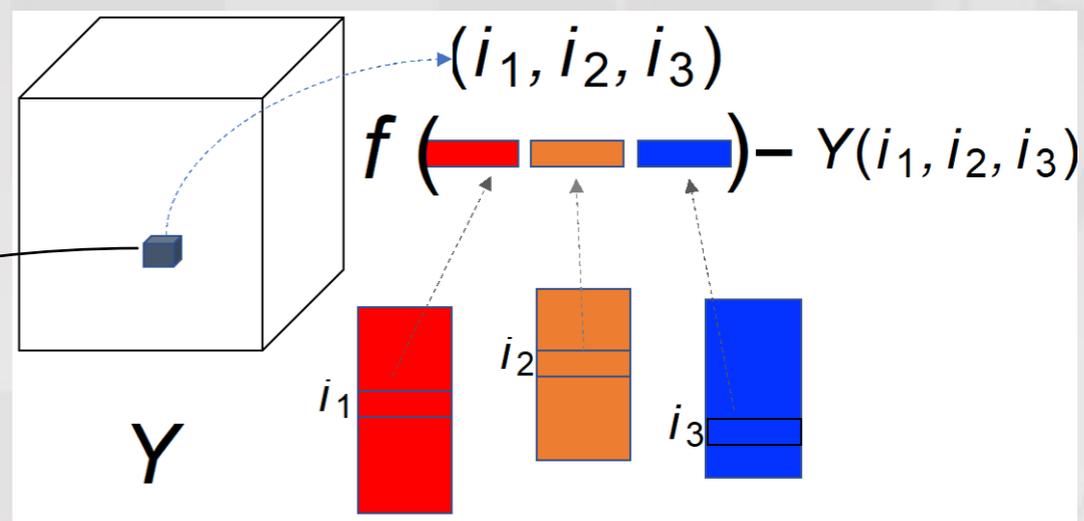


Where to put such sparsity?

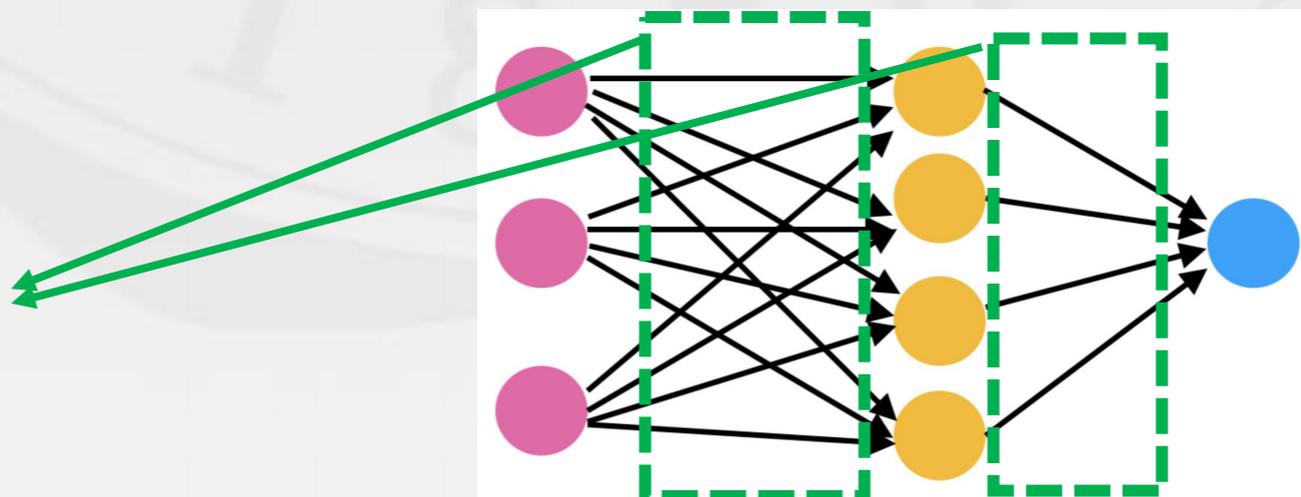
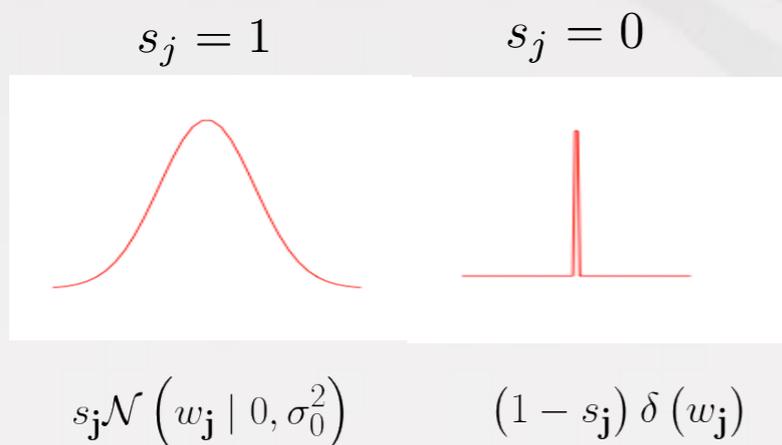
- **On NN weights – build a sparse BNN (SBDT[4])**

Interaction Records

use r	item	page	purchase
100	25	35	1
23	21	56	0
100	25	35	1
..
32	33	46	0



Assign s&s priors over each NN weights !





How the final sparsity exactly look like?

- Approx. of S&S Priors in exponential family:
Gaussian + Bernoulli

$$\begin{aligned} p(w_{mjt} | s_{mjt}) &\propto A(w_{mjt}, s_{mjt}) \\ &= \text{Bern}(s_{mjt} | c(\rho_{mjt})) \mathcal{N}(w_{mjt} | \mu_{mjt}^0, v_{mjt}^0) \end{aligned}$$

**Select posterior prob of each weight
<0.5: unselected <=> sparse**

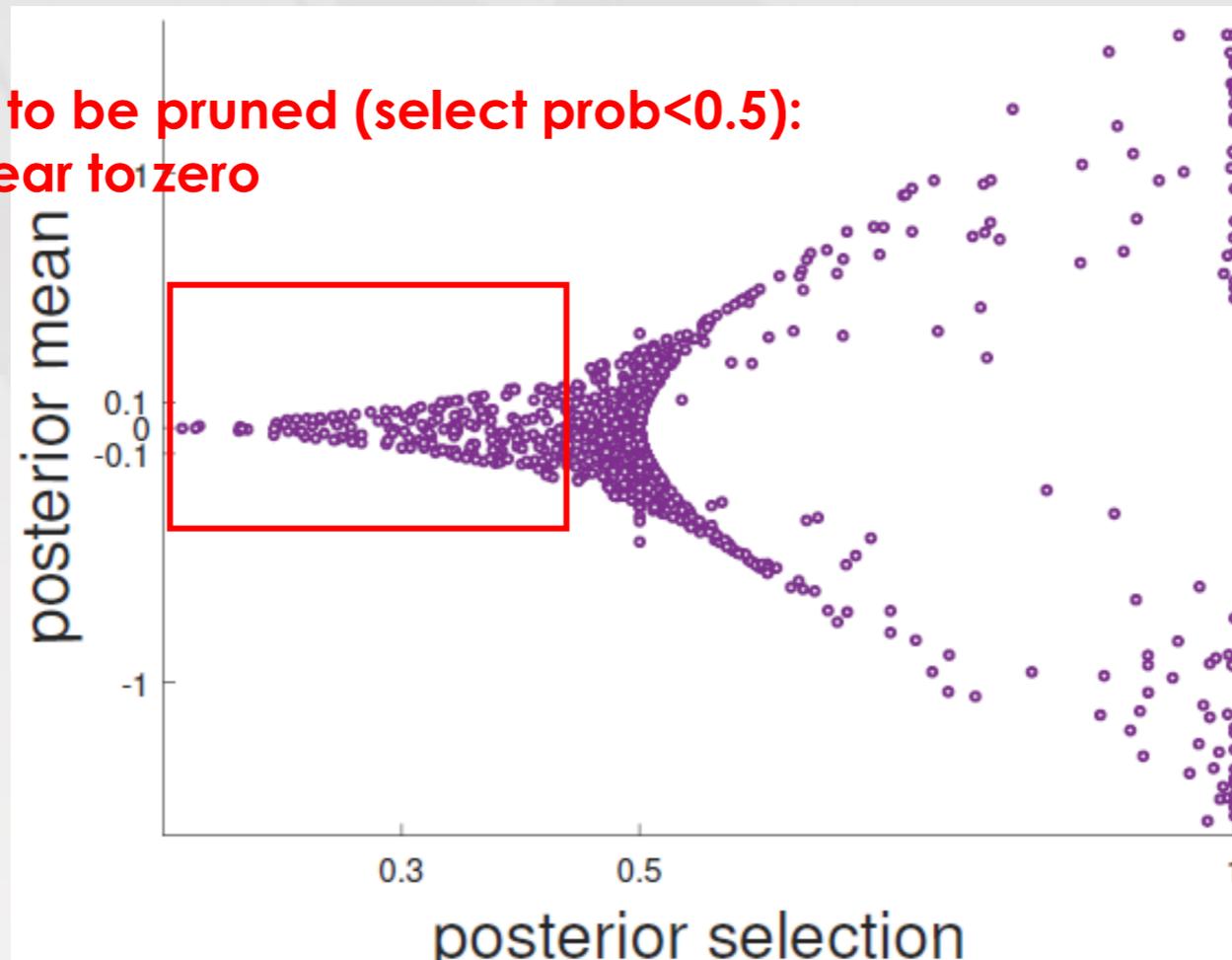


- **How sparse can model get? – Light model**

For SBDT work (Sparse BNN as factorization model)

- Plot of sparse-BNN weights after training
- Each weight has its posterior mean, var and selection prob.

**Weights to be pruned (select prob < 0.5):
Mean near to zero**

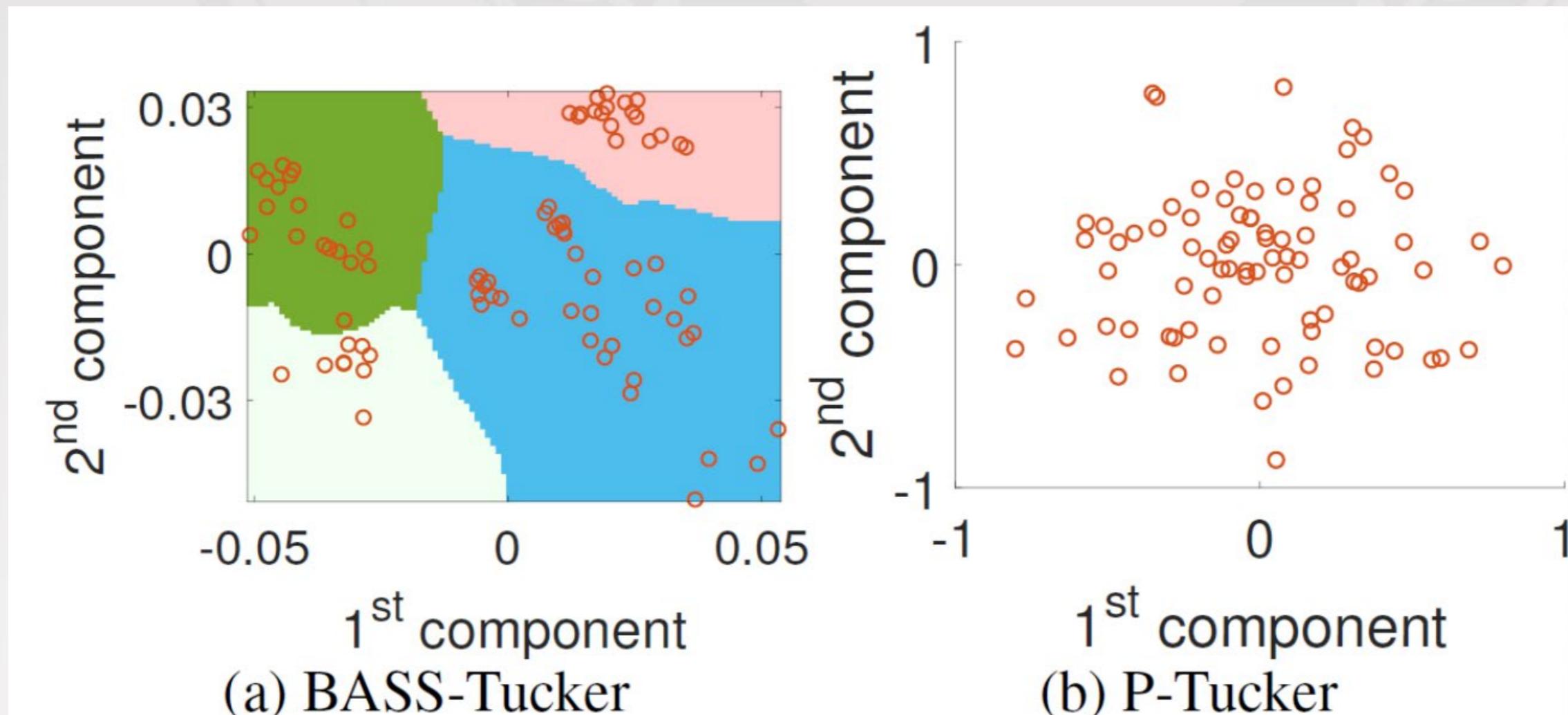




- **How sparse can model get? – Interpretability**

For BASS work (Tucker with Sparse core as factorization model)

- Plot of projected Tucker core elements with sparsity
- Significant structure of interactions





How to make use of quantized uncertainty?

An example in ad-recommend system

- plots of model predictions on **CTR(click-through-rate) tensor dataset**
- **exploration and exploitation**, optimization policy for down steaming tasks

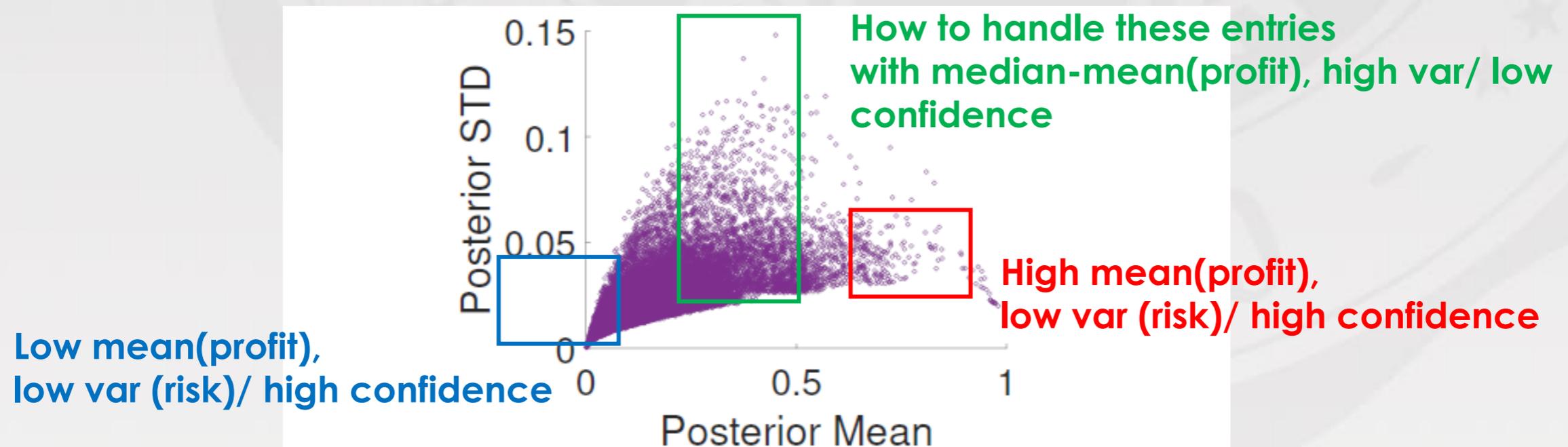


Figure 4: the posterior mean v.s. the posterior standard deviation (STD) of the click probability prediction.



Open questions for Cooperation

- Domain knowledge embedded in prior
- Make good use of the uncertainty measure
- Challenges and inspiration from real world
- ...

Domain model(PDE/SDE) + AI4Science + new algos...



THE UNIVERSITY OF UTAH

Thanks for attention Q&A Time

Presenter' email: shikai.fang@utah.edu

Webpage: <https://www.cs.utah.edu/~shikai/>

Focus: Bayesian machine learning, tensor learning

知乎: 方轩固