



Near-Universal Multiplicative Updates for Nonnegative Einsum Factorization

John Hood
University of Chicago

Aaron Schein
University of Chicago



ICML
International Conference
On Machine Learning

Motivation: Large Family of Tensor Methods

Many powerful algorithms/models/tensor methods:



but extending/modifying/changing is hard/model-specific.

Tensor Decomposition as Einsum Factorization

Simple matrix factorization assumes:

$$y_{i,j} \approx \sum_{q=1}^Q \psi_{iq} \beta_{jq}$$

Indices appearing on the RHS but not LHS are **contracted**

observed indices vs latent indices

Written more compactly in terms of only the indices:

$$ij \leftarrow iq, jq$$

Tensor Contractions are Einsum Factorizations

$$y_{i \rightarrow j}^{(t)} \approx \sum_{c=1}^C \sum_{d=1}^C \sum_{k=1}^K \sum_{r=1}^R \psi_{ic}^{(\rightarrow)} \psi_{jd}^{(\leftarrow)} \phi_{ak} \theta_{t,r} \lambda_{c \rightarrow d}^{(r)}$$

$$ij \atop a \rightarrow j \leftarrow ic, jd, ak, tr, cdkr$$

4-mode Tucker: Contraction of 4 matrices, 1 tensor

$$ij \atop a \rightarrow j \leftarrow ic, jd, ak, tr, cp, dp, kp, rp, p$$

4-mode Tucker with low-rank core tensor

$$ij \atop a \rightarrow j \leftarrow itc, jtd, ak, tr, cdkr$$

4-mode Tucker with time-varying factor matrices

...

This Work: Nonnegative Einsum Factorization (NNEinFact)

```

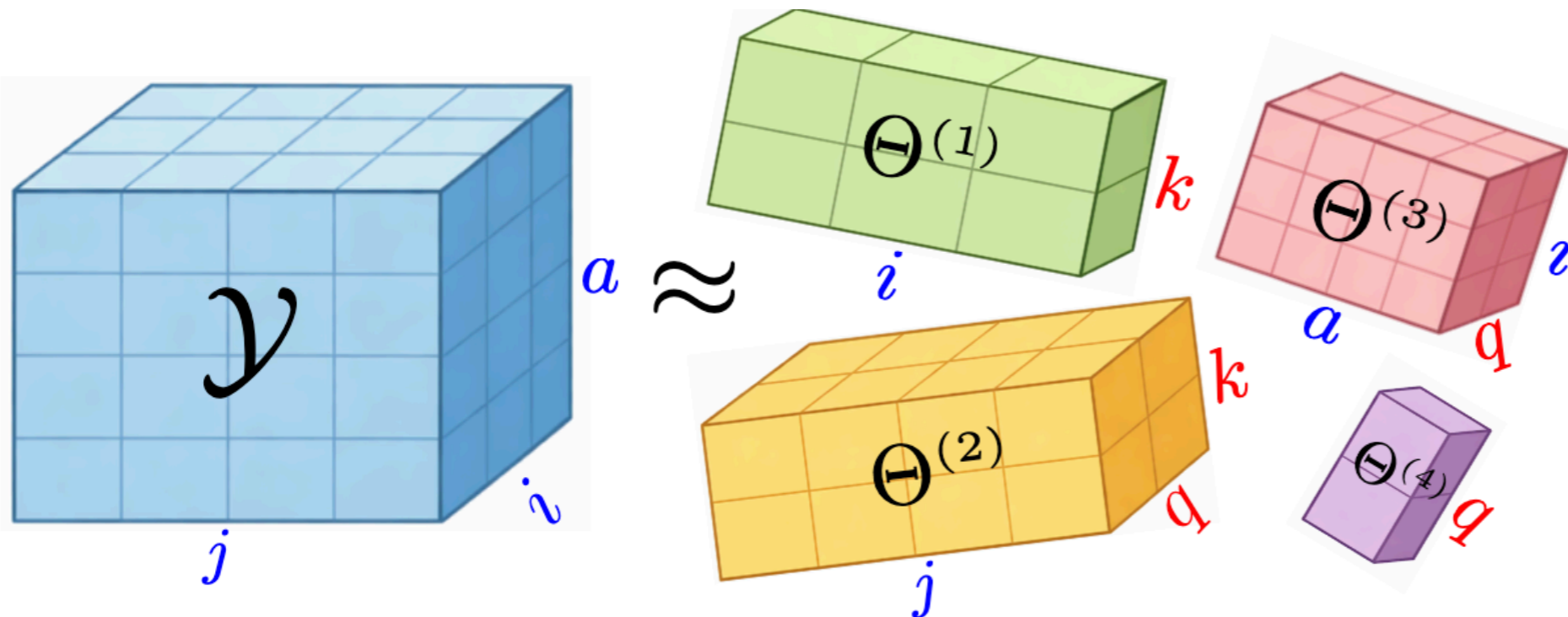
from einfact import NNEinFact

model = NNEinFact(
    model_str='ik,jqk,aiq,q->ija',
    loss='alpha-beta',  $\alpha=0.1$ ,  $\beta=1.1$ )

model.fit(data=Y_ija)

 $\theta_{ik}$ ,  $\theta_{jqk}$ ,  $\theta_{aiq}$ ,  $\theta_q$  = model.get_params()
    
```

NMF	$ik, jk \rightarrow ij$
CP	$ik, jk, ak \rightarrow ija$
Tucker	$ic, jd, ak, cdk \rightarrow ija$
Tensor train	$ic, jcd, ad \rightarrow ija$
CANDELINC	$ic, jd, ak, cq, dq, kq \rightarrow ija$
Custom	$ic, jd, ack, kd, k \rightarrow ija$



Loss examples

- Euclidean distance
- Forward KL divergence
- Reverse KL divergence
- Itakura-Saito divergence
- Hellinger distance
- Neyman χ^2 divergence
- Pearson χ^2 divergence
- α -divergence
- β -divergence
- (α, β) -divergence

One Multiplicative Update Algorithm for All

$$\hat{Y} \leftarrow \text{einsum}(\text{model_str}, \{\Theta^{(\ell)}\}_{\ell=1}^L) \quad \text{loss} \quad \min_{\Theta^{(1)}, \dots, \Theta^{(L)}} \mathcal{L}(Y, \hat{Y}) = \sum_{\mathbf{i}} \mathcal{L}(y_{\mathbf{i}}, \hat{y}_{\mathbf{i}}), \quad \Theta^{(\ell)} \geq \varepsilon$$

“near-universal” multiplicative update rule

$$\Theta^{(\ell)} \leftarrow \max \left(\varepsilon, \Theta^{(\ell)} \odot g^{-1} \left(\frac{\sum_{\mathbf{i}} [\nabla_{\Theta^{(\ell)}} \hat{y}_{\mathbf{i}}] a(y_{\mathbf{i}}, \hat{y}_{\mathbf{i}})}{\sum_{\mathbf{i}} [\nabla_{\Theta^{(\ell)}} \hat{y}_{\mathbf{i}}] b(y_{\mathbf{i}}, \hat{y}_{\mathbf{i}})} \right) \right)$$

implementation just involves calls to einsum

4: while not converged do

5: for $\ell = 1, \dots, L$ do

6: $\hat{Y} \leftarrow \text{einsum}(\text{model_str}, \{\Theta^{(\ell)}\}_{\ell=1}^L)$

7: $A \leftarrow \text{einsum}(\text{instr}_{\ell}, \Theta^{(1)}, \dots, \Theta^{(\ell-1)}, a(Y, \hat{Y}), \Theta^{(\ell+1)}, \dots, \Theta^{(L)})$

8: $B \leftarrow \text{einsum}(\text{instr}_{\ell}, \Theta^{(1)}, \dots, \Theta^{(\ell-1)}, b(Y, \hat{Y}), \Theta^{(\ell+1)}, \dots, \Theta^{(L)})$

9: $\Theta^{(\ell)} \leftarrow \max(\varepsilon, \Theta^{(\ell)} \odot g^{-1}(\frac{A}{B}))$

Theory and Computation

Theory

- ▶ Convergence of the objective values via MM
- ▶ Convergence of the iterates $\Theta^{(1)}, \dots, \Theta^{(L)}$ to a stationary point via BCD theory

Computation

- ▶ einfact.py file is less than 200 lines
- ▶ Factorization structure
 - ▶ To change the model, change the string
 - ▶ Rapid iteration and refinement
- ▶ Loss function
 - ▶ α and β : tunable parameters that shape the loss
 - ▶ $0 < \alpha < 1$: shrinks large entries y_i — robustness to outliers
 - ▶ $\alpha > 1$: magnifies large entries — robustness to missing data input as zeros
 - ▶ $\beta = 1$: least-squares, $\beta = 0$: KL divergence

```
from einfact import NNEinFact

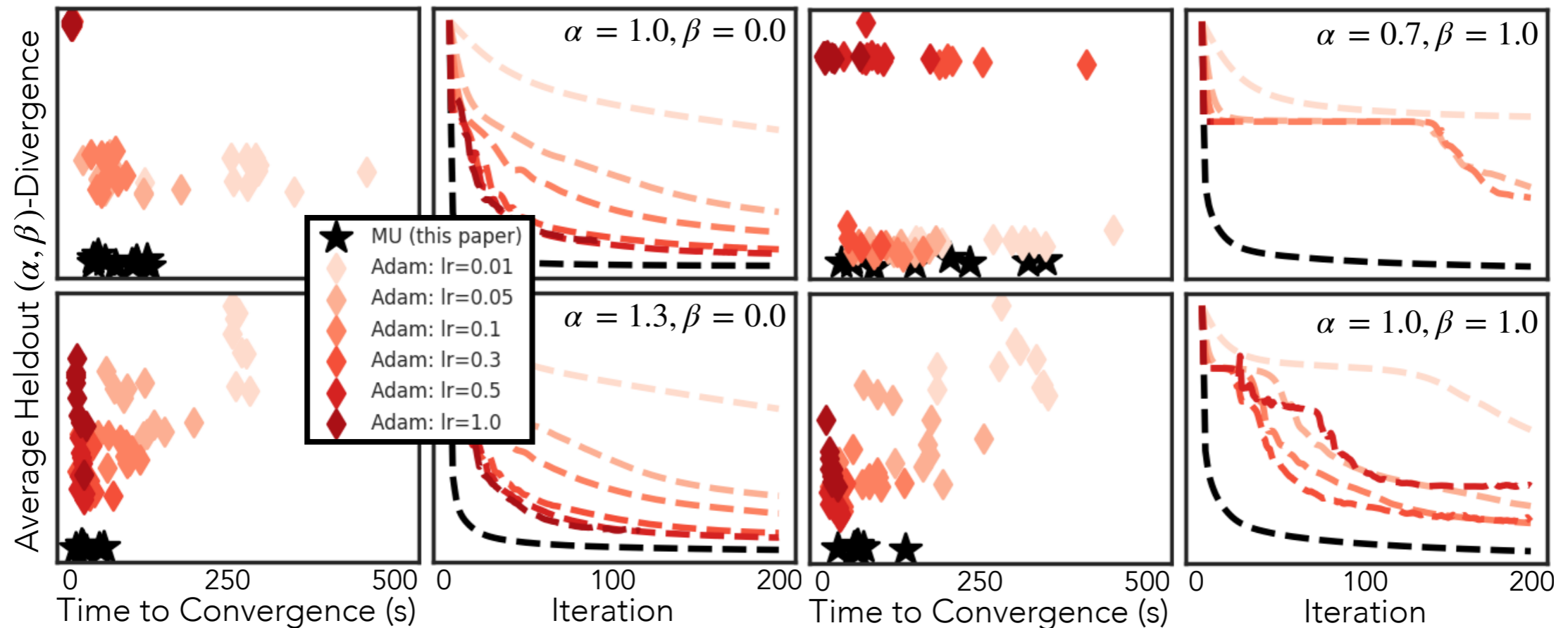
model = NNEinFact(
    model_str='ik,jqk,aiq,q->ija',
    loss='alpha-beta', alpha=0.1, beta=1.1)

model.fit(data=Y_ija)

theta_ik, theta_jqk, theta_aiq, theta_q = model.get_params()
```

Fast Convergence of Multiplicative Updates

To fit custom models, the only alternative options are automatic-differentiation gradient-based methods, such as Adam.

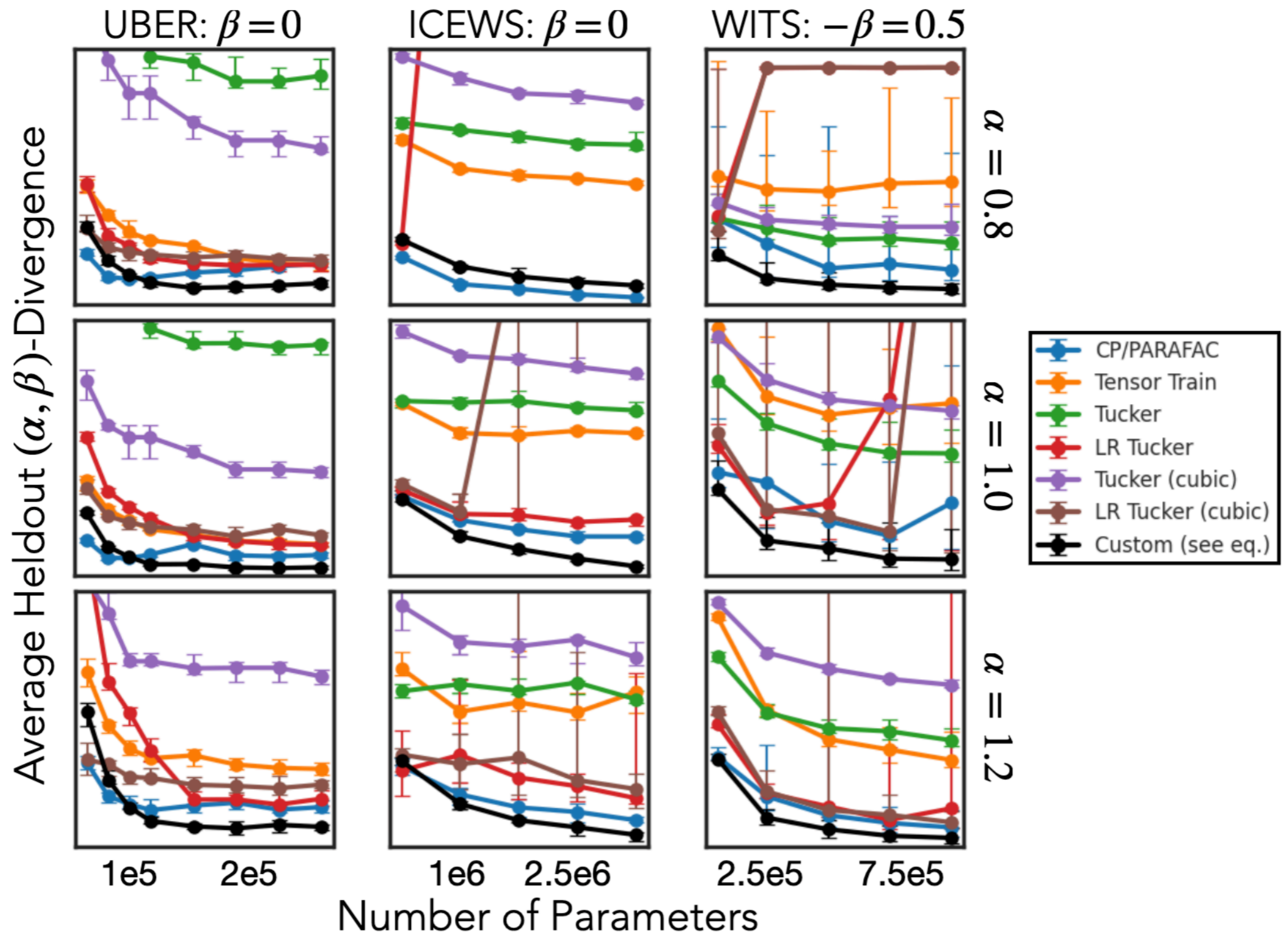


Compared to Adam, NNEinFact achieves:

- ▶ Faster convergence
- ▶ Lower heldout loss (better fit to the data)

Custom Models Often Fit the Data Better

Custom models \implies more parameter-efficient representations of the data



Flexible Parts-Based Scientific Discovery

Case study: New York City Uber pickup data, April to September of 2014

$$\mathcal{Y} \in \mathbb{N}_0^{27 \times 7 \times 24 \times 400 \times 400}$$

5-mode tensor:

- ▶ week
- ▶ day of week
- ▶ hour of day
- ▶ latitude
- ▶ longitude

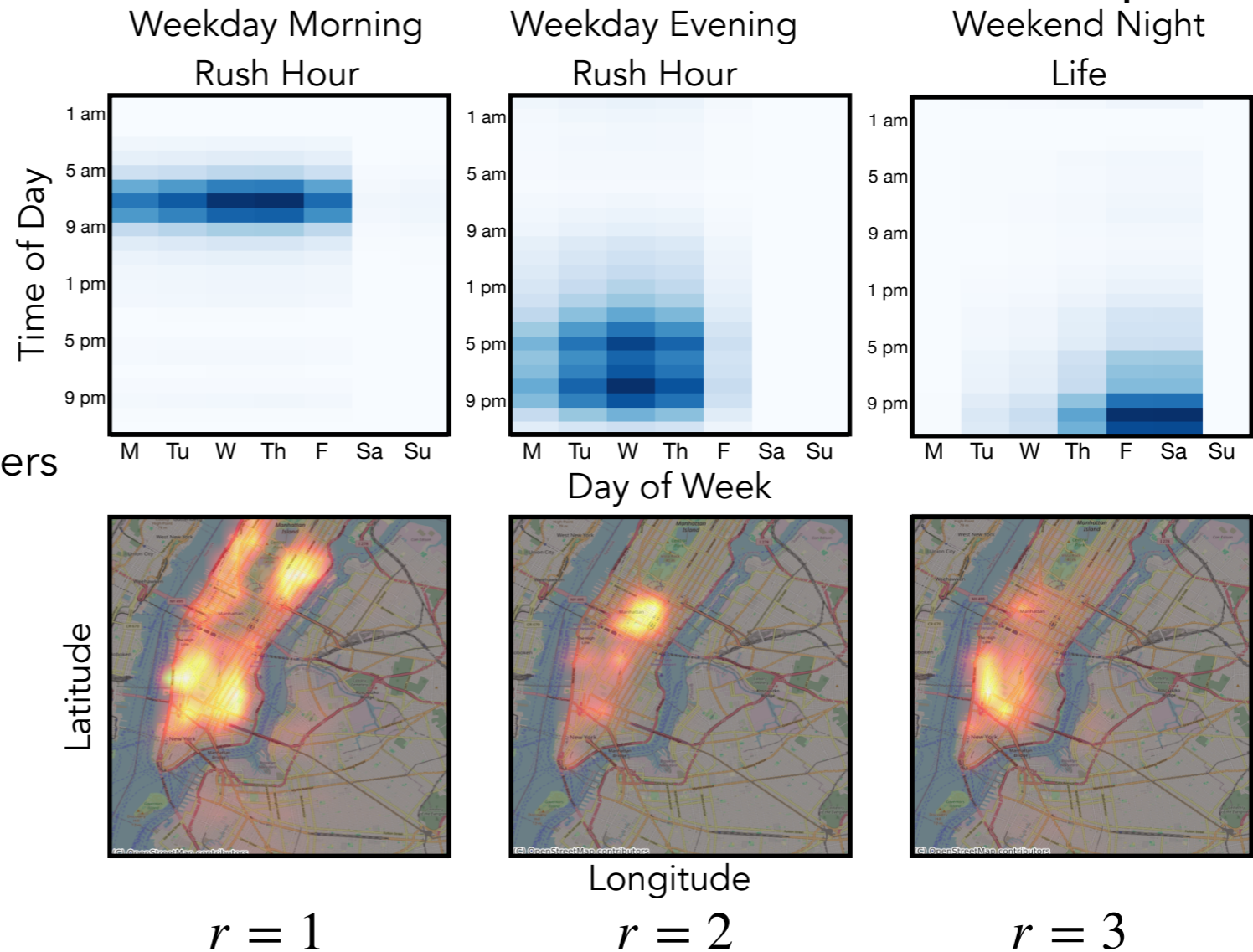
Model: CP-like, except extra parameters for spatial structure

$$\hat{y}_{wdhij} = \sum_{r=1}^R \theta_{wr}^{(1)} \theta_{dr}^{(2)} \theta_{hr}^{(3)} \underbrace{\sum_{k=1}^K \theta_{irk}^{(4)} \theta_{jrk}^{(5)}}_{\phi_{ijr}}$$



$wr, dr, hr, ikr, jkr \rightarrow wdhij$

Inferred Parameters in a Custom Model for NYC Uber Pickup Data



NNEinFact: Summary

Einsum representation: Any tensor network, decomposition, or factorization that can be represented as a tensor contraction can be represented as an einsum factorization.

Multiplicative update algorithm: NNEinFact is a general-purpose multiplicative update algorithm for *any* nonnegative einsum factorization under a wide family of loss functions (e.g., Gaussian likelihood, Poisson likelihood, robust losses, ...).

- ▶ Guaranteed to converge
- ▶ Flexible tool with simple, user-friendly code

Empirical advantages:

- ▶ Converges faster than gradient-based approaches to lower loss values
- ▶ Allows for expressive, parameter-efficient factorizations tailored to your data