# **Tensorial Change Analysis using Probabilistic** AAAI-19 #2614 **Tensor Regression**



Tsuyoshi Idé ("Ide-san", 井手 剛, tide@us.ibm.com) IBM Research, T. J. Watson Research Center



## **Problem setting**



### **Prior work**

## Alternating least squares (ALS; [Zhou+ 13]) $\rightarrow$ Non-probabilistic

Linear model with CP (canonical polyadic) expansion  $y \sim (\boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{A}}) \equiv \sum_{i_1, \dots, i_M} \mathcal{X}_{i_1, \dots, i_M} \mathcal{A}_{i_1, \dots, i_M}$  $\boldsymbol{\mathcal{A}} = \sum_{r=1}^R \boldsymbol{a}^{1,r} \circ \boldsymbol{a}^{2,r} \circ \dots \circ \boldsymbol{a}^{M,r}, \text{ or } \mathcal{A}_{i_1, i_2, \dots, i_M} = \sum_{r=1}^R a_{i_1}^{1,r} a_{i_2}^{2,r} \cdots a_{i_M}^{M,r}$ Factor out one component using tensor algebra  $(\boldsymbol{\mathcal{X}},\boldsymbol{\mathcal{A}}) = \operatorname{Tr}\left[ (\boldsymbol{\mathsf{A}}^{l})^{\top} \boldsymbol{\mathsf{X}}_{(l)} (\boldsymbol{\mathsf{A}}^{M} \odot \cdots \boldsymbol{\mathsf{A}}^{l+1} \odot \boldsymbol{\mathsf{A}}^{l-1} \odot \cdots \boldsymbol{\mathsf{A}}^{1}) \right]$  $\mathbf{A}^{l} \equiv [\mathbf{a}^{l,1}, \dots, \mathbf{a}^{l,R}], \odot = \text{Khatri-Rao product}$  $\mathbf{X}_{(l)} \equiv \text{mode-}l \text{ matricization of } \mathbf{X},$ •Find it, given the others. Repeat.

## Gaussian process regression [Zhao+ 14]. $\rightarrow$ Much computational cost on testing

•Use standard GPR formula (K: kernel matrix) (predictive mean) =  $\boldsymbol{k}(\boldsymbol{\mathcal{X}})^{\top} [\boldsymbol{\mathsf{K}} + \sigma^2 \boldsymbol{\mathsf{I}}_N]^{-1} \boldsymbol{y}_N$ (predictive variance) =  $k_0 - \mathbf{k}(\mathbf{\mathcal{X}})^{\top} [\mathbf{K} + \sigma^2 \mathbf{I}_N]^{-1} \mathbf{k}(\mathbf{\mathcal{X}})$  $[\mathbf{K}]_{n,n'} = k(\boldsymbol{\mathcal{X}}^{(n)}, \boldsymbol{\mathcal{X}}^{(n')})$ 

Define proper kernel function for tensors • Note: You can **NOT** use the Frobenius norm as it reduces to the simple vectorization approach

 $\|\boldsymbol{\mathcal{X}}^{(n)} - \boldsymbol{\mathcal{X}}^{(n')}\|_{\mathrm{F}}^2 = \|\mathrm{vec}(\boldsymbol{\mathcal{X}}^{(n)}) - \mathrm{vec}(\boldsymbol{\mathcal{X}}^{(n')})\|_2^2$ 

## **Bayesian method [Guhaniyogi+ 17].** $\rightarrow$ Needs complex Monte Carlo sampling

- Observation model with CP expanded coefficients  $p(y \mid \mathcal{X}, \mathcal{A}, \lambda) = \mathcal{N}(y \mid (\mathcal{A}, \mathcal{X}), \lambda^{-1})$
- Prior distributions allowing sparse model by shrinkage  $p(\{\boldsymbol{a}^{l,r}\}) = \prod \prod \mathcal{N}(\boldsymbol{a}^{l,r} \mid \boldsymbol{0}, (\phi^r \tau) \operatorname{diag}(\boldsymbol{w}^{l,r})), \ p(\lambda) = \mathcal{G}(\lambda \mid a_{\lambda}, b_{\lambda})$  $p(\boldsymbol{w}^{l,r}) = \prod_{i=1}^{a_l} \operatorname{Exp}(w_i^{l,r} \mid \nu^{l,r}), \ p(\sqrt{\nu^{l,r}}) = \mathcal{G}(\nu^{l,r} \mid a_{\nu}, b_{\nu})$  $p(\boldsymbol{\phi}) = \operatorname{Dir}(\boldsymbol{\phi} \mid \boldsymbol{\alpha}), \ p(\tau) = \mathcal{G}(\tau \mid \boldsymbol{a}_{\tau}, \boldsymbol{b}_{\tau})$

Inference is done by Markov chain Monte Carlo

## **Probabilistic tensor regression**

## **Proposed Bayesian model**

Observation model: Gaussian + CP expansion  $p(y \mid \mathcal{X}, \mathcal{A}, \lambda) = \mathcal{N}(y \mid (\mathcal{A}, \mathcal{X}), \lambda^{-1})$  Gaussian distribution  $\boldsymbol{\mathcal{A}} = \sum^{N} \boldsymbol{a}^{1,r} \circ \boldsymbol{a}^{2,r} \circ \cdots \circ \boldsymbol{a}^{M,r} \quad (\circ = \text{direct product})$ 

Transform to linear model by factoring out each

## Variational EM framework for inference



component given the others

$$(\boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{A}}) = \sum_{r=1}^{R} (\boldsymbol{\mathcal{X}}, \boldsymbol{a}^{1,r} \circ \dots \circ \boldsymbol{a}^{M,r}) = \sum_{r=1}^{R} (\boldsymbol{a}^{l,r})^{\top} \boldsymbol{\phi}^{l,r} (\boldsymbol{\mathcal{X}}),$$
$$\boldsymbol{\phi}^{l,r} = \boldsymbol{\mathsf{X}}_{(l)} (\boldsymbol{a}^{M,r} \otimes \dots \otimes \boldsymbol{a}^{l+1,r} \otimes \boldsymbol{a}^{l-1,r} \otimes \dots \otimes \boldsymbol{a}^{1,r}),$$
$$\text{``Mode-! matricization''} \quad \otimes = (\text{Kronecker product})$$

Prior distribution: Gauss-Gamma

$$p(\boldsymbol{a}^{l,r} \mid \boldsymbol{b}^{l,r}) = \mathcal{N}(\boldsymbol{a}^{l,r} \mid \boldsymbol{0}, (\boldsymbol{b}^{l,r})^{-1} \mathbf{I}_{d_l}),$$
$$p(\boldsymbol{b}^{l,r} \mid \alpha_0, \beta_0) = \frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} (\boldsymbol{b}^{l,r})^{\alpha_0 - 1} \mathrm{e}^{-\beta_0 \boldsymbol{b}^{l,r}} \mathrm{e}^{-\beta_0 \boldsymbol{b}^{l,r}}$$
gamma distribution

### **Change scores**

**Overall change: Use predictive distribution** 

Outliner score for a new sample

$$c(y, \boldsymbol{\mathcal{X}}) \equiv -\ln p(y \mid \boldsymbol{\mathcal{X}}, \mathbf{D})$$
$$= \frac{\{y - \mu(\boldsymbol{\mathcal{X}})\}^2}{2\sigma^2(\boldsymbol{\mathcal{X}})} + \frac{1}{2}\ln\{2\pi\sigma^2(\boldsymbol{\mathcal{X}})\}$$

Change-point score: smoothed version of it

 $Q(\{a^{l,r}\}) = \prod q^{1,r}(a^{1,r})q^{2,r}(a^{2,r})\cdots q^{M,r}(a^{M,r})$ 

-Iterate between VB for  $\{ {m a}^{l,r} \}$  and  $\{ b^{l,r} \}$  and point-estimate of  $\lambda$ 

Results:  $\circ \text{Posterior} \quad q_1^{l,r}(\boldsymbol{a}^{l,r}) = \mathcal{N}(\boldsymbol{a}^{l,r} \mid \boldsymbol{\mu}^{l,r}, \boldsymbol{\Sigma}^{l,r}), \quad q_2^{l,r}(\boldsymbol{b}^{l,r}) = \mathcal{G}(\boldsymbol{b}^{l,r} \mid \boldsymbol{\alpha}^{l,r}, \boldsymbol{\beta}^{l,r})$  $\circ \operatorname{Precision} \quad \lambda^{-1} = \frac{1}{N} \sum_{n=1}^{N} \left\{ y^{(n)} - \sum_{r=1}^{R} (\mathcal{X}^{(n)}, \boldsymbol{\mu}^{1, r} \circ \cdots \circ \boldsymbol{\mu}^{M, r}) \right\}^{2} + (\operatorname{higher order term})$ Predictive distribution  $p(y \mid \boldsymbol{x}, D) = \mathcal{N}(y \mid \mu(\boldsymbol{\mathcal{X}}), \sigma^2(\boldsymbol{\mathcal{X}}))$  $\mu(\boldsymbol{\mathcal{X}}) = \eta + \sum_{r=1}^{R} (\boldsymbol{\mathcal{X}}, \boldsymbol{\mu}^{1,r} \circ \cdots \circ \boldsymbol{\mu}^{M,r}), \quad \sigma^{2}(\boldsymbol{\mathcal{X}}) = \lambda^{-1} + \sum_{r=1}^{R} \sum_{l=1}^{M} (\boldsymbol{\varphi}^{l,r})^{\top} \boldsymbol{\Sigma}^{l,r} \boldsymbol{\varphi}^{l,r}$  $\boldsymbol{\varphi}^{l,r} = \boldsymbol{\mathsf{X}}_{(l)}(\boldsymbol{\mu}^{M,r} \otimes \cdots \otimes \boldsymbol{\mu}^{l+1,r} \otimes \boldsymbol{\mu}^{l-1,r} \otimes \cdots \otimes \boldsymbol{\mu}^{1,r}).$ 

#### Mode- and dimension-wise change: Use posterior learned

Change analysis score of the *i*-th dimension of the *l*-th mode

$$c_{i}^{l}(\tilde{\mathbf{D}},\mathbf{D}) \equiv \frac{1}{R} \sum_{r=1}^{R} \int d\boldsymbol{a}^{l,r} q^{l,r}(\boldsymbol{a}^{l,r}) \ln \frac{q^{l,r}(a_{i}^{l,r} \mid \boldsymbol{a}_{-i}^{l,r})}{\tilde{q}^{l,r}(a_{i}^{l,r} \mid \boldsymbol{a}_{-i}^{l,r})}$$

VB plays the key role in change analysis: It is the factorized assumption that makes possible mode- and dimension-wise scoring



Application

#### **Change diagnosis for semiconductor etching tool**



"golden period" (or reference period)

