# (Some) Challenges in Tensor Mining

**Evrim Acar**

Sandia National Labs., Livermore, CA
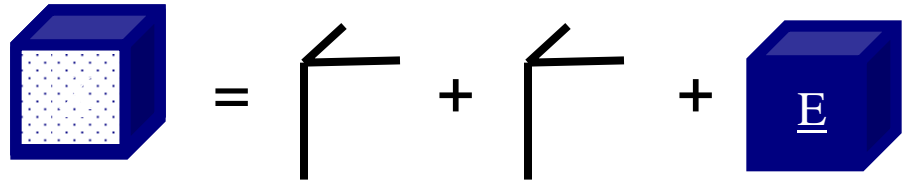
Sandia National Laboratories

# Tensor Mining

**unsupervised**

**Parafac**

dense or sparse

$=$ + + $\underline{E}$

**Tucker**

$=$ + $\underline{E}$

**supervised**

$\underline{X}_{train}$

$\underline{X}_{test}$

$y$

$\underline{X}_{train}$ $\approx$

$y$

# App I: Social Networks Analysis
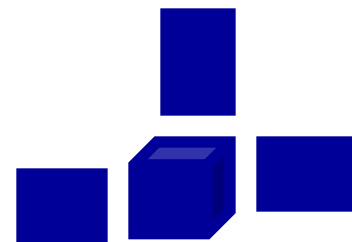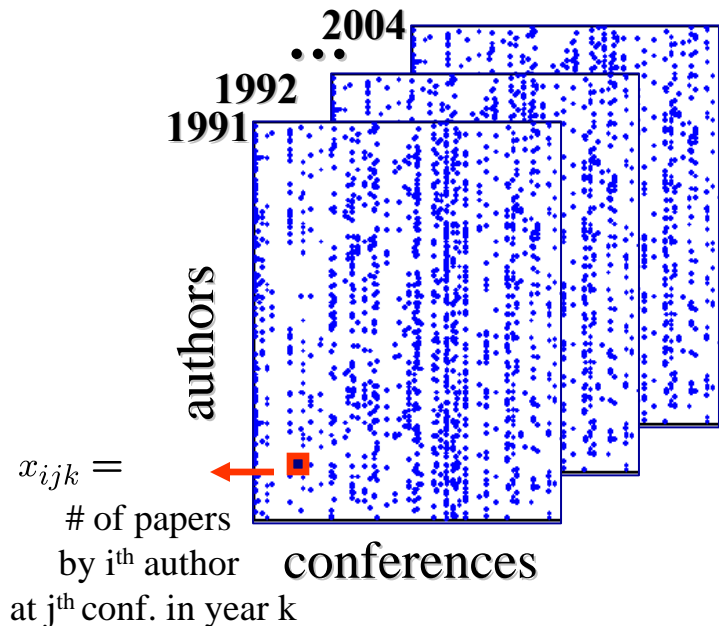
*Joint work with*
*T.G. Kolda and D. M. Dunlavy*

- In social networks, we are interested in modeling relationships (links) evolving over time.

- Example:
  - DBLP dataset: Authors x Conferences x Years (10K x 2K x 14: ~0.1% dense)



$x_{ijk} =$

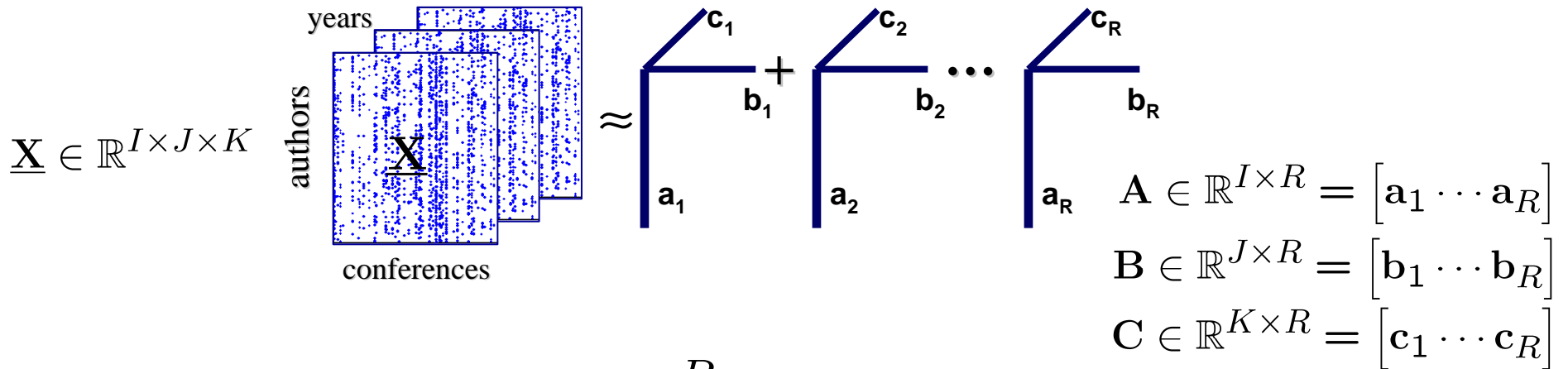\# of papers by i[th] author at j[th] conf. in year k

**Q1: Can we use tensor decompositions to model the data and extract meaningful underlying factors?**

Q2: Can we predict who is going to publish at which conferences in future?

*(Link Prediction in time)*

*SIAM CS&E*
*March 2-6, 2009*

# Modeling DBLP using PARAFAC

$$\underline{\mathbf{X}} \in \mathbb{R}^{I \times J \times K}$$



$$\mathbf{A} \in \mathbb{R}^{I \times R} = \begin{bmatrix} \mathbf{a}_1 \cdots \mathbf{a}_R \end{bmatrix}$$

$$\mathbf{B} \in \mathbb{R}^{J \times R} = \begin{bmatrix} \mathbf{b}_1 \cdots \mathbf{b}_R \end{bmatrix}$$

$$\mathbf{C} \in \mathbb{R}^{K \times R} = \begin{bmatrix} \mathbf{c}_1 \cdots \mathbf{c}_R \end{bmatrix}$$

$$\underline{\mathbf{X}} \approx \sum_{r=1}^{R} \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

$$\min_{\mathbf{A},\mathbf{B},\mathbf{C}} \left\| \underline{\mathbf{X}} - \left( \sum_{r=1}^{R} \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r \right) \right\|^2$$

$$\| \underline{\mathbf{X}} \| = \sum_{k=1}^{K} \sum_{j=1}^{J} \sum_{i=1}^{I} x_{ijk}^2$$

- Solve using a gradient-based optimization approach
- Initialization:
  - first two modes using svd, $R \leq I, J$
  - last mode: random, $R > K$

# Components make sense!

# What if data is a **Sparse** tensor with **Missing** entries?

- Sparse Data:

$$\min_{A,B,C} \left\| \underline{X} - \left( \sum_{r=1}^{R} a_r \circ b_r \circ c_r \right) \right\|^2$$

$$A \in \mathbb{R}^{I \times R} = \begin{bmatrix} a_1 \cdots a_R \end{bmatrix}$$

$$B \in \mathbb{R}^{J \times R} = \begin{bmatrix} b_1 \cdots b_R \end{bmatrix}$$

**Success with 70% randomly missing data**



[Tomasi&Bro, 2005]

- Missing Data [Kiers, 1997; Tomasi & Bro, 2005] :

$$\min_{A,B,C} \left\| \underline{W} * \left( \underline{X} - \left( \sum_{r=1}^{R} a_r \circ b_r \circ c_r \right) \right) \right\|$$

$$w_{ijk} = \begin{cases} 1, & \text{if } x_{ijk} \text{ not missing,} \\ 0, & \text{if } x_{ijk} \text{ missing.} \end{cases}$$

- Sparse & Missing:

$$\min_{A,B,C} \left\| \underline{W} * \left( \underline{X} - \left( \sum_{r=1}^{R} a_r \circ b_r \circ c_r \right) \right) \right\|^2 + \text{???}$$

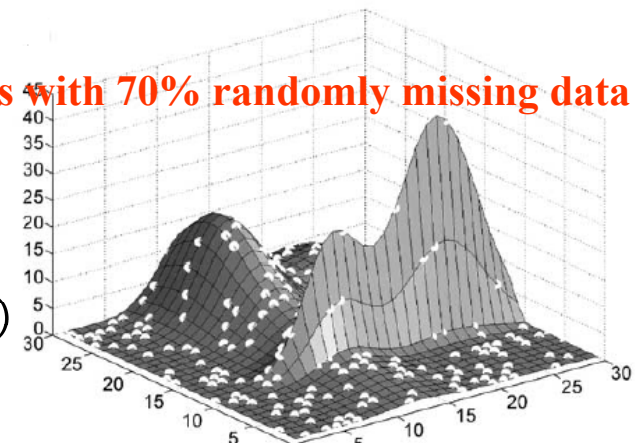$$w_{ijk} = \begin{cases} 1, & \text{if } x_{ijk} \text{ not missing,} \\ 0, & \text{if } x_{ijk} \text{ missing.} \end{cases}$$
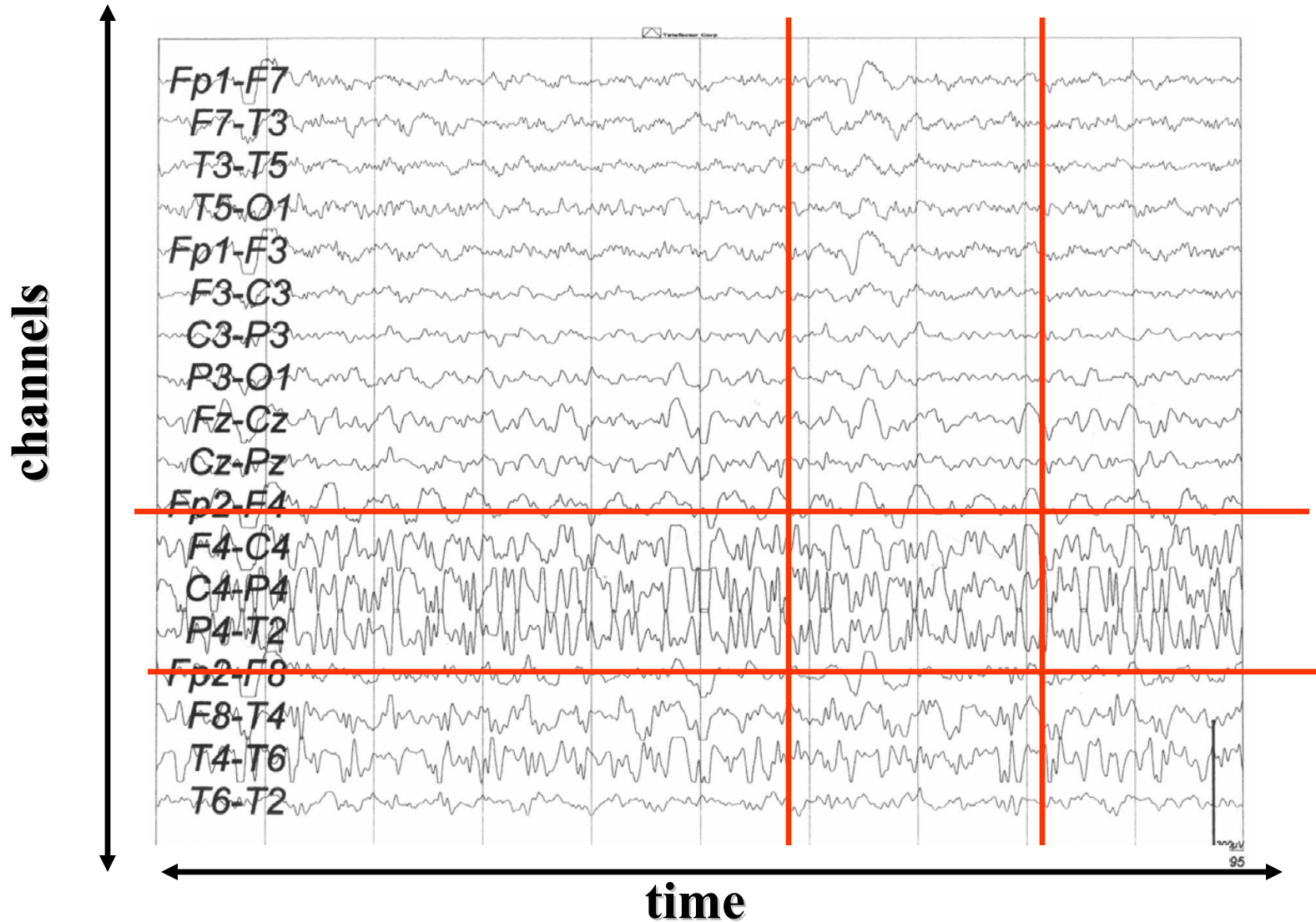
# App II: Understanding Epileptic Seizures

*Joint work with*
*R. Bro, B. Yener, C. A. Bingol. H. Bingol*

Sandia National Laboratories



channels

time

# Epilepsy Tensors

$x_{ij}$: Electrical potential at $i^{th}$ sample $j^{th}$ channel

**CWT**

$x_{ijk}$: Power of a wavelet coeff. at $i^{th}$ sample $j^{th}$ scale $k^{th}$ channel

**Time samples**

**Channels**

**Time samples**

**Channels**

**Scales (freq.)**

• Data rearranged as a three-way array using continuous wavelet transform (CWT):

  • Let $c_{ijk}$ be the wavelet coefficient at time sample $i$ at scale $j$ for the $k^{th}$ channel.

  • An Epilepsy Tensor is a three-way array, $\underline{X}$, where each entry $x_{ijk}$ is computed as:

$$x_{ijk} = |c_{ijk}|^2$$

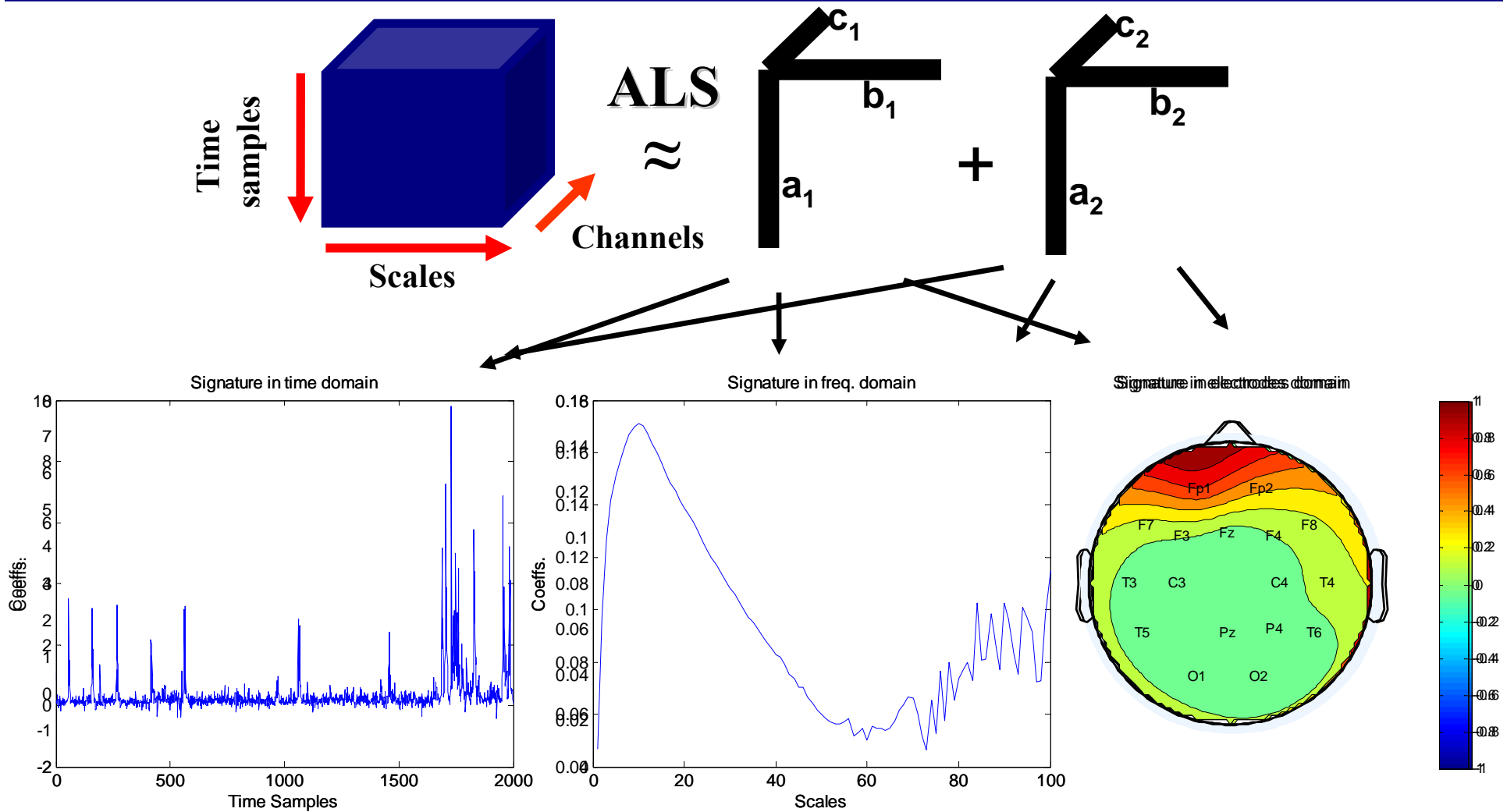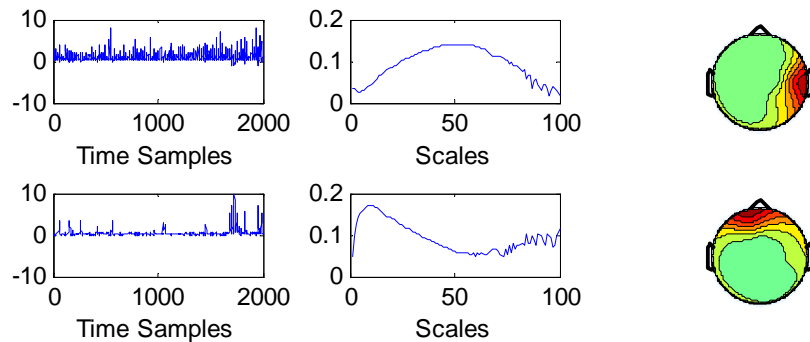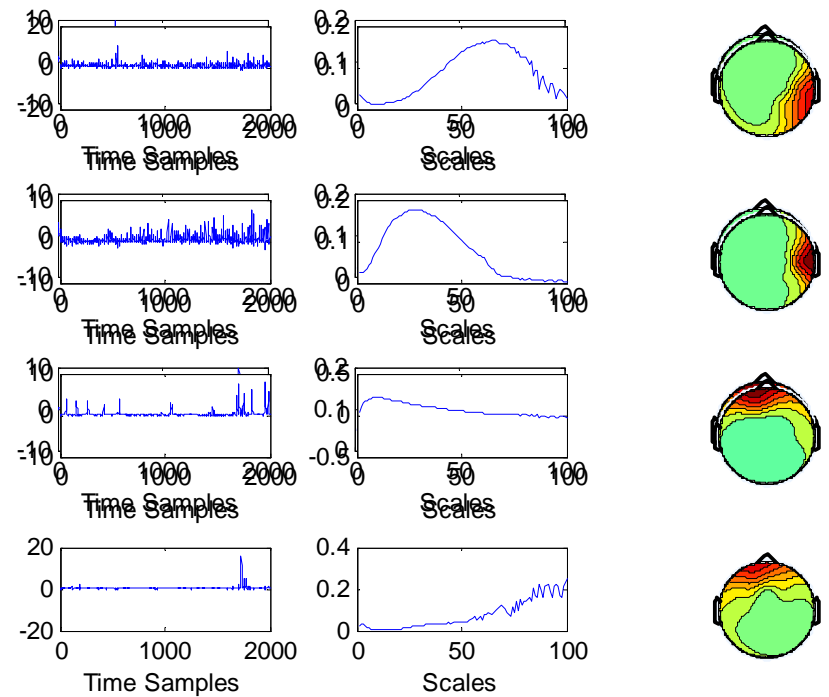# Epilepsy Focus Localization



Acar et al'07, De Vos et al'07

# How many components?

$$\underline{\mathbf{X}} \approx \sum_{r=1}^{R} \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

$R = 2$

$R = 4$

# How to initialize?

$$\underline{X} \approx \sum_{r=1}^{2} \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

**HOSVD**

**RANDOM**

# Understanding Epileptic Seizures

# Epilepsy Feature Tensor

- Construction of an Epilepsy Feature Tensor from multi-channel EEG



$$\begin{bmatrix} f_1(s) \\ f_2(s) \\ \\ f_n(s) \end{bmatrix}$$

$x_{ij}$: Electrical potential at $i^{th}$ channel $j^{th}$ time sample

**Time samples**

**Channels**

**Channels**

**Time epochs**

**Epilepsy Feature Tensor**

**Features**

$x_{ijk}$: Value of $j^{th}$ feature at $i^{th}$ epoch recorded at $k^{th}$ channel

# Seizure Recognition

**Training Set**
- Build a model using the training set **X** and the labels **y**.

**Test Set**
- Predict the labels of new recordings.

**X**          **y**$_{train}$

Time epochs

Pre$_1$
Seizure$_1$ → seizure
Post$_1$

Pre$_2$
Seizure$_2$
Post$_2$

Pre$_3$ → non-seizure
Seizure$_3$
Post$_3$

**y**$_{test}$

**Test**   **?**

# Multiway Classification(?)

- Potential Approaches
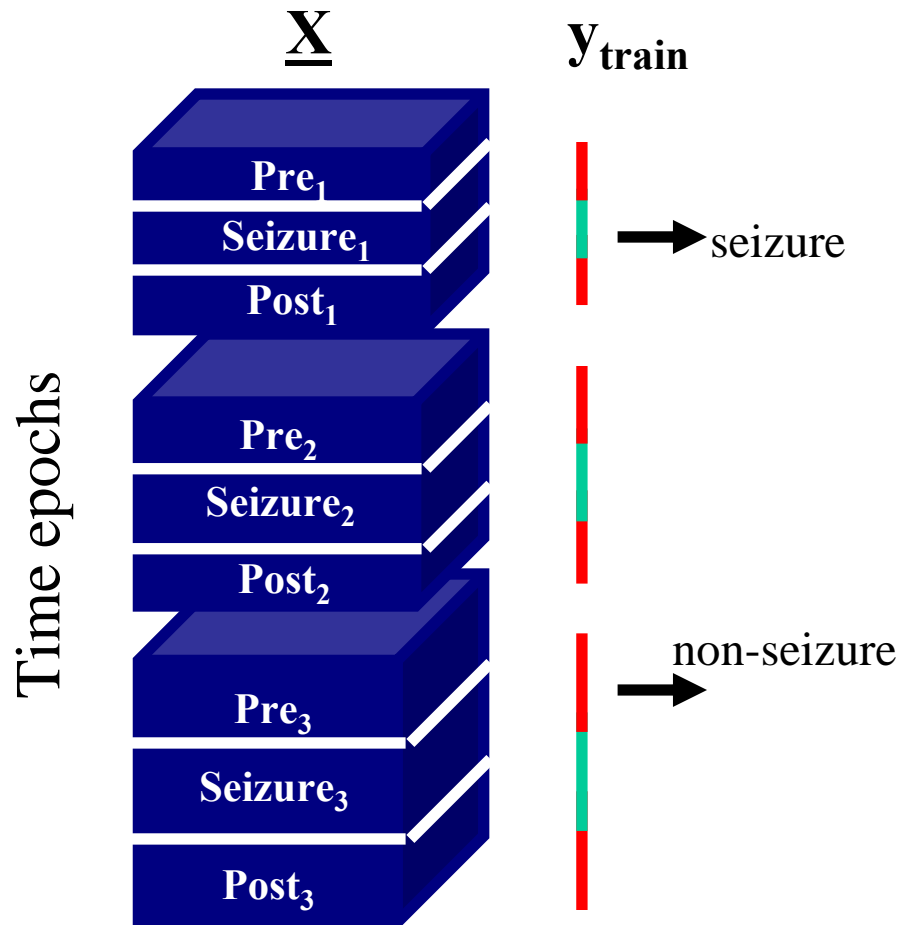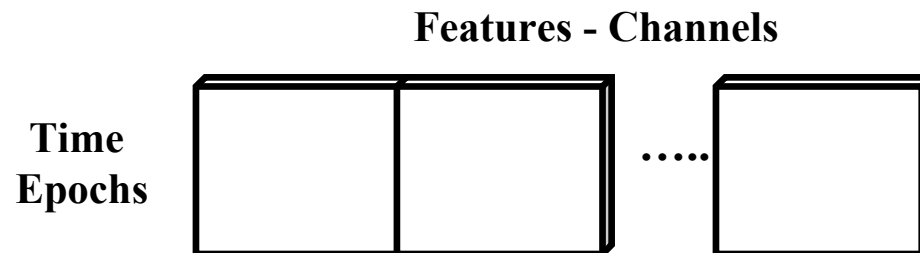  - Modify multiway regression models, e.g., multilinear PLS [Bro, 1996; Bro et al., 2001], as classifiers.

$$\max_{\mathbf{w}^J, \mathbf{w}^K} [cov(\mathbf{t}, \mathbf{y}) | t_i = (\mathbf{w}^J)^T \mathbf{X}_i \mathbf{w}^K]$$



  - Unfold the data and apply two-way classification, e.g., SVM.

**Features - Channels**

# Some challenges are …

- **Handling Sparse Data with Missing Entries:**
  - We need models to capture the underlying sparse factors in sparse tensors with missing entries.

- **Determining the Rank:**
  - Important also in practice.

- **Initialization:**
  - Algorithms suffer from the local minima problem. In practice, we may end up interpreting our results differently.

- **Supervised learning on tensors:**
  - We need classification models for tensors as good as the state-of-the-art two-way classification approaches such as SVMs.

# Thank you!

- **References:**
  - **Social Networks Analysis**: [Tensor toolbox & Poblano toolbox (by Sandia)]
    - Acar, Kolda and Dunlavy, An Optimization Approach for Fitting Canonical Tensor Decompositions, SAND2009-0857, Feb. 2009.
  - **Understanding Epileptic Seizures:** [PLS toolbox (by Eigenvector Research)]
    - Acar, Bingol, Bingol, Bro and Yener, Multiway Analysis of Epilepsy Tensors, *Bioinformatics*, 23(13): i10-i18, 2007.
    - Acar, Bingol, Bingol, Bro and Yener, Seizure Recognition on Epilepsy Feature Tensor, *Proc. 29th Int. Conf. IEEE Engineering in Medicine and Biology Society*, 2007.
  - **Survey:**
    - Acar and Yener, Unsupervised Multiway Data Analysis: A Literature Survey, *IEEE Transactions on Knowledge and Data Engineering*, 21(1): 6-20, 2009.

- **Contact:**

  Evrim Acar, Sandia National Laboratories, **eacarat@sandia.gov**