

IBM Research


# Sparse Gaussian Markov Random Field Mixtures for Anomaly Detection

Tsuyoshi Idé (“**Idé-san**”), Ankush Khandelwal\*, Jayant Kalagnanam  
IBM Research, T. J. Watson Research Center  
(\*Currently with University of Minnesota)

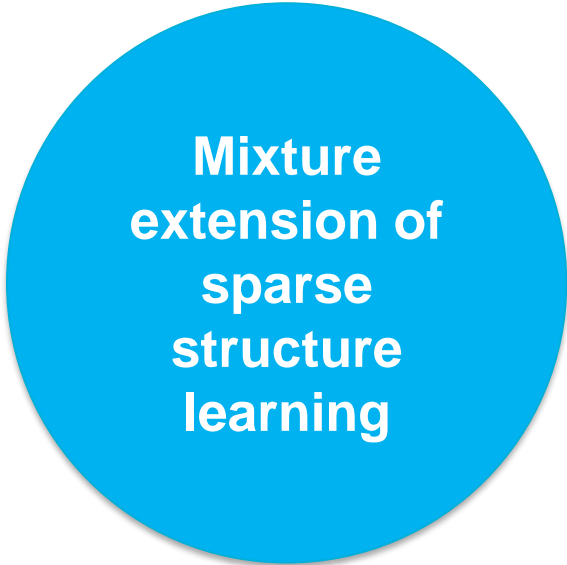
Proceedings of the 2016 IEEE International Conference on Data Mining ([ICDM 16](#)), Dec. 13-15, 2016, pp.955-960.

## Summary: Gaussian mixture + anomaly detection

Newly added features:



**Principled  
variable-wise  
scoring**



**Mixture  
extension of  
sparse  
structure  
learning**

## Define variable-wise anomaly score as conditional log-loss

- Anomaly score for the  $i$ -th variable of a new sample  $\mathbf{x} \in \mathbb{R}^M$

$$a_i(\mathbf{x}) = -\ln p(x_i \mid \mathbf{x}_{-i}, \mathcal{D})$$

conditional predictive  
distribution

- $x_i$ : the  $i$ -th variable
- $\mathbf{x}_{-i}$ : the rest
- $\mathcal{D}$ : training data

- c.f. overall score [Yamanishi+ 00]

$$a(\mathbf{x}) = -\ln p(\mathbf{x} \mid \mathcal{D})$$

predictive  
distribution for  $\mathbf{x}$

“ $a$  will be large if  $\mathbf{x}$  falls in the area where  $p(\mathbf{x} / \mathcal{D})$  is small”

## (For ref.) Why negative log p? It reproduces Mahalanobis distance in the single Gaussian case

Gaussian with mean  $\mu$  and covariance  $\Sigma$

$$\begin{aligned} a(\mathbf{x}) &= -\ln \mathcal{N}(\mathbf{x} \mid \mu, \Sigma) \\ &= \text{const.} + \frac{1}{2} \underbrace{(\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu)}_{\text{Mahalanobis distance}} \end{aligned}$$

# Use mixture of Gaussian Markov random fields (GMRF) for the conditional distribution

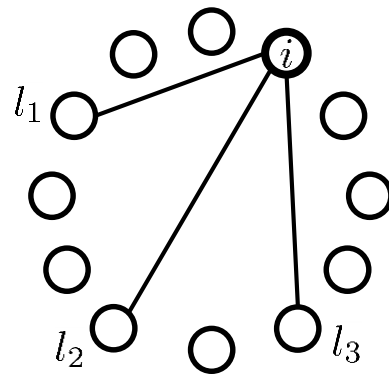
Variable-wise mixture weight

$$p(x_i | \mathbf{x}_{-i}, \mathcal{D}) = \sum_{k=1}^K \underbrace{g_k^i(\mathbf{x})}_{\text{GMRF}} \mathcal{N}(x_i | u_i^k, w_i^k),$$

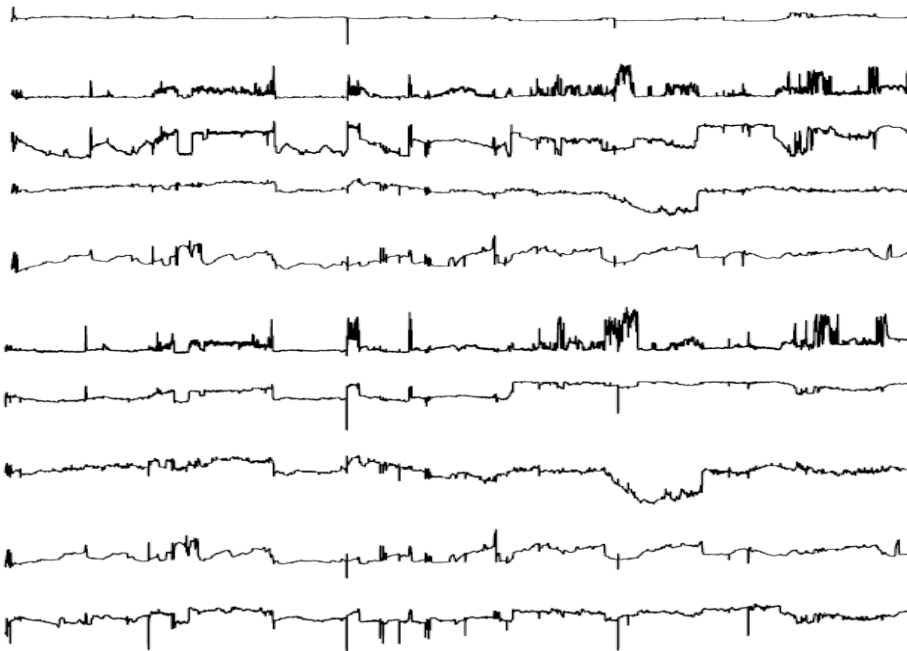
$$u_i^k = m_i^k - \frac{1}{A_{i,i}^k} \sum_{l \neq i}^M A_{i,l}^k (x_l - m_l^k)$$

$$w_i^k = \frac{1}{A_{i,i}^k}.$$

- GMRF is characterized as conditional distribution of Gaussian  $\mathcal{N}(\mathbf{x} | \mathbf{m}^k, (\mathbf{A}^k)^{-1})$
- GMRF describes dependency among variables



## Mixture model with i.i.d. assumption on data is a practical compromise: Compressor data example



This is **normal** operation data

It looks piece-wise stationary (with heavy noise)

Time-series modeling looks too hard

# Tackling noisy data of complex systems: Strategy of designing inference algorithm



## How to choose $K$ ?

- Give a large enough  $K$  first. Let the algorithm decide on an optimal value.



## How to handle heavy noise? How to stably learn the model?

- Use sparsity-enforcing prior to leverage sparse structure learning
- 

## Two-step approach to GMRF mixture learning: Model

$$p(x_i | \mathbf{x}_{-i}, \mathcal{D}) = \sum_{k=1}^K g_k^i(\mathbf{x}) \mathcal{N}(x_i | u_i^k, w_i^k),$$

### Step 1: Find GMRF parameters

- Observation model

$$p(\mathbf{x} | \mathbf{z}, \boldsymbol{\mu}, \Lambda) \equiv \prod_{k=1}^K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}^k, (\Lambda^k)^{-1})^{z_k},$$

- Priors:

Gauss-Laplace for  $(\boldsymbol{\mu}^k, \Lambda^k)$

Categorical for  $\{\mathbf{z}_k\}$  (each sample)

### Step 2: Find variable-wise weights given GMRF parameters

- Observation model

$$p(x_i | \mathbf{x}_{-i}, h^i) = \prod_{k=1}^K \mathcal{N}(x_i | u_i^k, w_i^k)^{h_k^i}$$

- Priors:

Categorical-Dirichlet for  $\{h_k^i\}$

(each sample)



## Two-step approach to GMRF mixture learning: Inference

- Use variational Bayes (VB) for the 1<sup>st</sup> and 2<sup>nd</sup> steps
- The 1<sup>st</sup> step achieves sparsity over variable dependency and mixture components
  - Variable dependency: (iteratively) solve graphical lasso [Friedman+ 08]

$$\bar{\Lambda}^k \leftarrow \arg \max_{\Lambda^k} \left\{ \ln |\Lambda^k| - \text{Tr}(\Lambda^k Q^k) - \frac{\rho}{N^k} \|\Lambda^k\|_1 \right\}.$$

- Mixture components: (iteratively) point-estimated for ARD (automated relevance determination) [Corduneanu+ 01]
- Details → paper

## Overview of the approach (for multivariate noisy sensor data)

### ■ Initialize

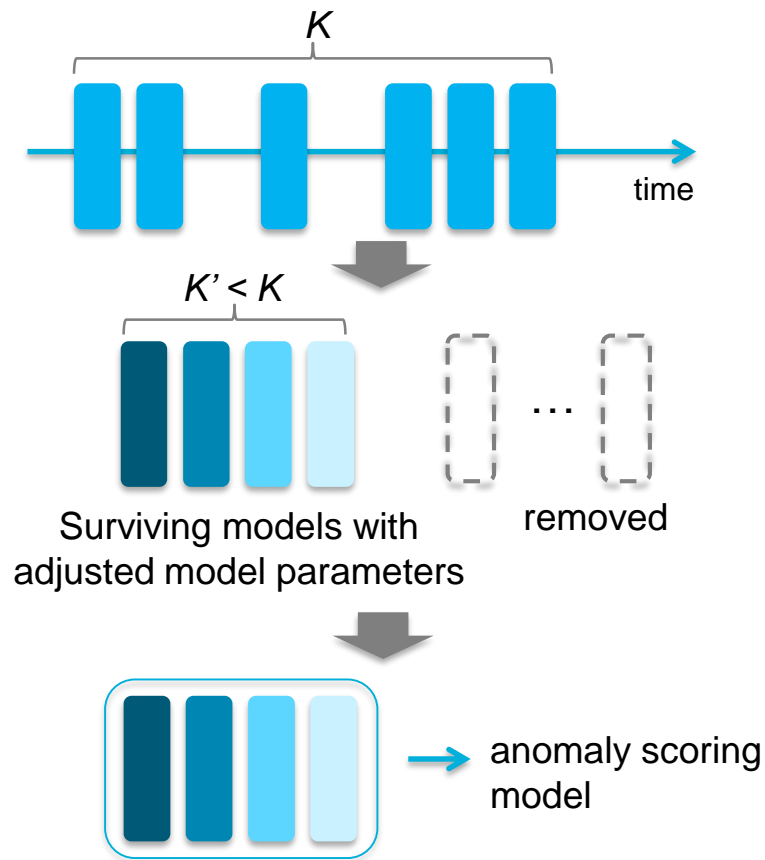
- Randomly pick time-series blocks with a large enough  $K$
- Run graphical lasso separately to initialize  $\{(\mu^k, \Lambda^k)\}$

### ■ Step 1

- Iteratively update  $\{(\mu^k, \Lambda^k)\}$  and
- remove clusters with zero weight

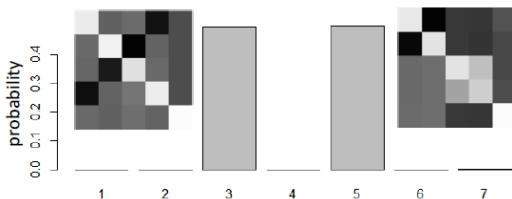
### ■ Step 2

- Compute variable-wise mixture weights
- Produce anomaly scoring model

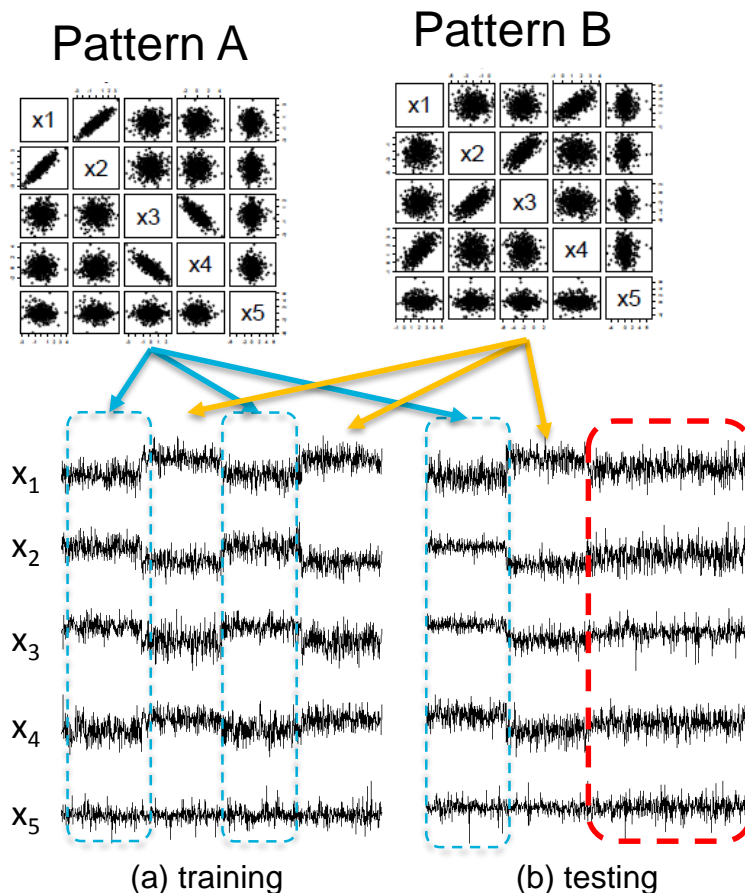


## Results: Synthetic data (see paper for real application)

- Data generated
  - Training: A-B-A-B
  - Testing: A-B-(anomaly)
- Results
  - Successfully recovered 2 major patterns starting from  $K=7$



- Achieved better performance in anomaly detection (in terms of AUC)



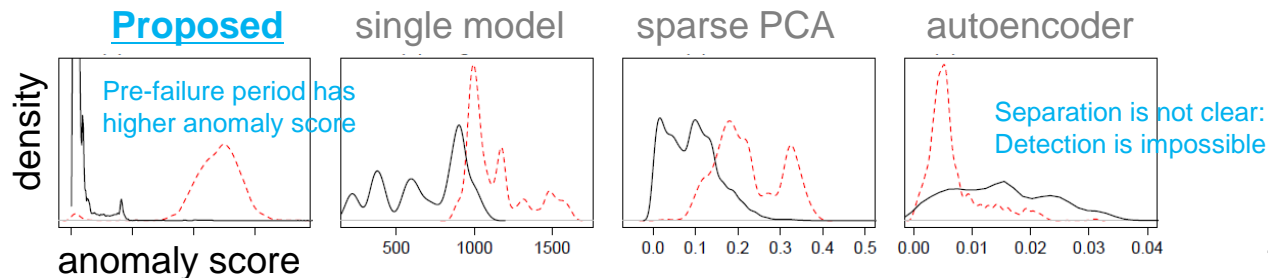
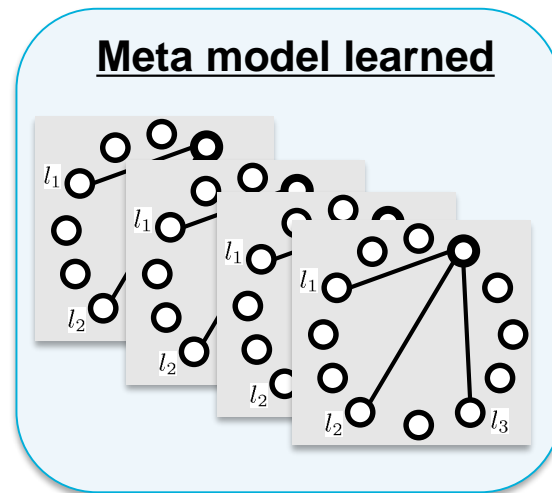
## Conclusion

- Proposed a new outlier detection method, the sparse GMRF mixture
- Our method is capable of handling multiple operational states in the normal condition and variable-wise anomaly scores.
- Derived variational Bayes iterative equations based on the Gauss-delta posterior model

**Thank you!**

## Results: Detecting pump surge failures

- Computed anomaly score for  $x_{14}$ , which is a flow-rate variable, based on the normal state model learned
- Compared anomaly score between the pre-failure region (24h) and several normal periods
  - Black: normal
  - red: pre-failure period
- Clearly outperformed alternative methods including neural network (autoencoder)



## Leveraging variational Bayes method for inference

- Assumption: posterior distribution is factorized

$$p(\mathbf{z}^{(1)}, \dots, \boldsymbol{\mu}^1, \dots, \Lambda^1, \dots) = \prod_{n=1}^N q(\mathbf{z}^{(n)}) \prod_{k=1}^K q(\boldsymbol{\mu}^k, \Lambda^k)$$

- Posterior is determined so that the KL divergence between the factorized form and the full posterior
  - Full posterior is proportional to the complete likelihood

## Two-step approach to GMRF mixture learning: Inference

- Use variational Bayes (VB) for the 1<sup>st</sup> and 2<sup>nd</sup> steps
- The 1<sup>st</sup> step achieves sparsity over variable dependency and mixture components
  - Variable dependency: (iteratively) solve graphical lasso
  - Mixture components: point-estimated for ARD (automated relevance determination)
    - ✓ [Corduneanu-Bishop 01]

### VB iteration for the 1<sup>st</sup> step

$$N^k \leftarrow \sum_{n=1}^N r_k^{(n)}, \quad \pi_k \leftarrow \frac{N^k}{N},$$

Point-estimated cluster weight

$$\bar{\mathbf{x}}^k \leftarrow \frac{1}{N^k} \sum_{n=1}^N r_k^{(n)} \mathbf{x}^{(n)},$$

$$\Sigma^k \leftarrow \frac{1}{N^k} \sum_{n=1}^N r_k^{(n)} (\mathbf{x}^{(n)} - \bar{\mathbf{x}}^k)(\mathbf{x}^{(n)} - \bar{\mathbf{x}}^k)^\top,$$

$$\lambda^k \leftarrow \lambda_0 + N^k, \quad \mathbf{m}^k \leftarrow \frac{1}{\lambda^k} (\lambda_0 \mathbf{m}_0 + N^k \bar{\mathbf{x}}^k),$$

$$\mathbf{Q}^k \leftarrow \Sigma^k + \frac{\lambda_0}{\lambda^k} (\bar{\mathbf{x}}^k - \mathbf{m}_0)(\bar{\mathbf{x}}^k - \mathbf{m}_0)^\top,$$

$$\bar{\Lambda}^k \leftarrow \arg \max_{\Lambda^k} \left\{ \ln |\Lambda^k| - \text{Tr}(\Lambda^k \mathbf{Q}^k) - \frac{\rho}{N^k} \|\Lambda^k\|_1 \right\}$$

graphical lasso [Friedman+ 08]