IBM **Research**

# Multi-task Multi-modal Models for Collective Anomaly Detection

**Tsuyoshi Ide** ("**Ide-san**"), Dzung T. Phan, J. Kalagnanam
PhD, Senior Technical Staff Member
IBM Thomas J. Watson Research Center

This slides are available at ide-research.net.

# Outline

- Problem setting

- Modeling strategy

- Model inference approach

- Experimental results

# Wish to build a <u>collective</u> monitoring solution

System 1
(in New
Orleans)

$\vdots$

System $s$

$\vdots$

System $S$
(in New York)

- You have many similar but not identical industrial assets

- You want to build an anomaly detection model for each of the assets

- Straightforward solutions have serious limitations
  - 1. Treat the systems separately. Create each model individually
    - ✓ Suffers from lack of fault examples
  - 2. Build one universal model by disregarding individuality
    - ✓ Model fit is not good

# Practical requirements: Need to capture both commonality and individuality
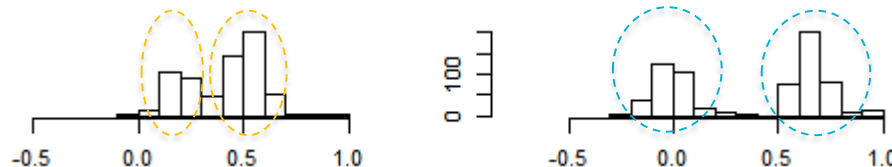
System 1
(in New
Orleans)

System $s$
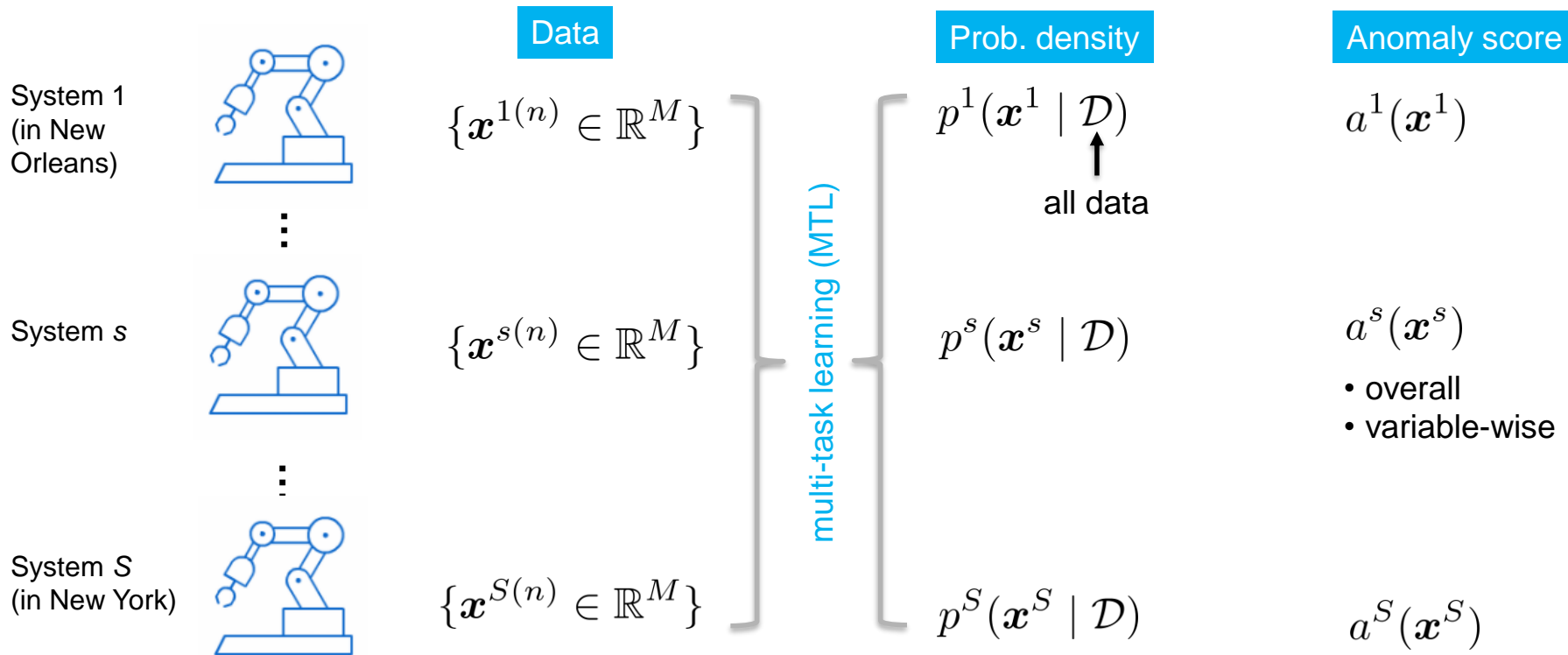
System $S$
(in New York)

- Capture both individuality and commonality

- Automatically capture multiple operational states
  o Real-world is not single-peaked (single-modal)

- Be robust to noise

- Be highly interpretable for diagnosis purposes

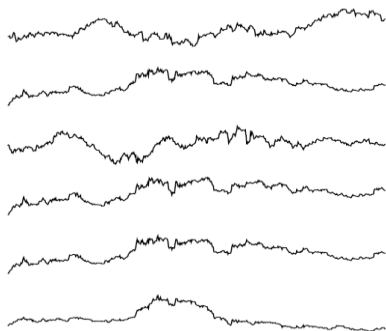# Formalizing the problem as multi-task density estimation for anomaly detection

| | Data | multi-task learning (MTL) | Prob. density | Anomaly score |
|---|---|---|---|---|

System 1 (in New Orleans)

$$\{\boldsymbol{x}^{1(n)} \in \mathbb{R}^M\}$$

$$p^1(\boldsymbol{x}^1 \mid \mathcal{D})$$

↑ all data

$$a^1(\boldsymbol{x}^1)$$

⋮

System $s$

$$\{\boldsymbol{x}^{s(n)} \in \mathbb{R}^M\}$$

$$p^s(\boldsymbol{x}^s \mid \mathcal{D})$$

$$a^s(\boldsymbol{x}^s)$$

• overall
• variable-wise

⋮

System $S$ (in New York)

$$\{\boldsymbol{x}^{S(n)} \in \mathbb{R}^M\}$$

$$p^S(\boldsymbol{x}^S \mid \mathcal{D})$$

$$a^S(\boldsymbol{x}^S)$$

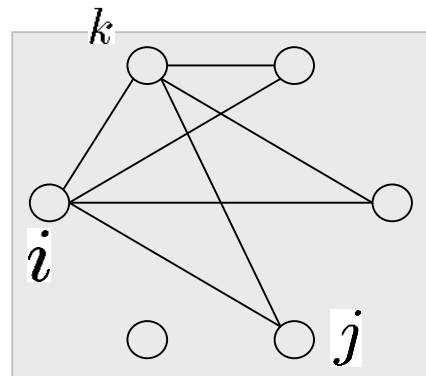multi-task learning (MTL)

# Outline

- Problem setting

- Modeling strategy

- Model inference approach

- Experimental results

# Use Gaussian graphical model (GGM)-based anomaly detection approach as the basic building block

**Multi-variate data**

**Sparse graphical model**

**Anomaly score**



$$\max_{\Lambda} \left\{ \ln \det \Lambda - \mathrm{tr}(\Sigma \Lambda) - \rho ||\Lambda||_1 \right\}$$

sample covariance

$$a(\boldsymbol{x}) = \begin{cases} -\ln p(\boldsymbol{x} \mid \mathcal{D}) \\ \quad \text{Overall score} \\ -\ln p(x_i \mid \boldsymbol{x}_{-i}, \mathcal{D}) \end{cases}$$

Variable-wise score

training data

[Ide+ SDM09] [Ide+ ICDM16]

# Basic modeling strategy: Combine common pattern dictionary with individual weights



**Individual sparse weights**

**Common dictionary of sparse graphs**
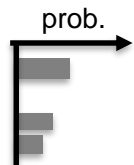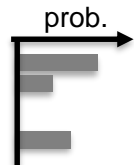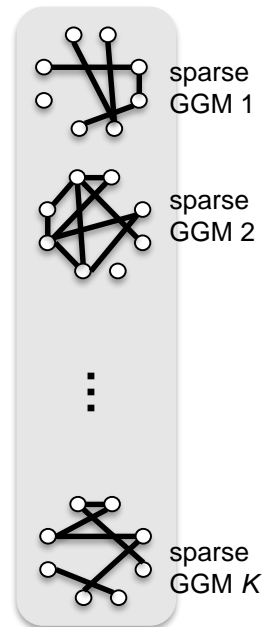
System 1 (in New Orleans)

System $s$

System $S$ (in New York)

prob.

prob.

prob.

sparse GGM 1

sparse GGM 2

sparse GGM $K$

Monitoring model for System 1

Monitoring model for System 2

Monitoring model for System $S$

GGM=Gaussian Graphical Model

8

# Basic modeling strategy: Resulting model will be a sparse mixture of sparse GGM

System *s*



prob.

$\times$

sparse GGM 1

sparse GGM 2

⋮

sparse GGM *K*

GGM=Gaussian Graphical Model

Monitoring model for System s

Gaussian mixture

$$= \sum_{k=1}^{K} \pi_k^s \, \mathcal{N}(\boldsymbol{x}^s \mid \boldsymbol{\mu}^k, (\Lambda^k)^{-1})$$

**Sparse mixture weights**

(= automatic determination of the number of patterns)

**Sparse Gaussian graphical model**

# Outline

- Problem setting

- Modeling strategy

- Model inference approach

- Experimental results

# Employing a Bayesian model for multi-modal MTL

- **Observation model (for the *s*-th task)**
  - Gaussian mixture with task-dependent weight

$$\prod_{k=1}^{K} \mathcal{N}(\boldsymbol{x}^s \mid \boldsymbol{\mu}^k, (\Lambda^k)^{-1})^{z_k^s}$$

- **Sparsity enforcing priors (non-conjugate)**
  - Laplace prior for the precision matrix
  - Bernoulli prior for the mixture weights

$$p(\Lambda^k) = \left(\frac{\rho}{4}\right)^{M^2} \exp\left(-\frac{\rho}{2}\|\Lambda^k\|_1\right)$$

$$p(\boldsymbol{\pi}^s) = p_0^{\|\boldsymbol{\pi}^s\|_0}(1-p_0)^{G-\|\boldsymbol{\pi}^s\|_0}$$

- **Conjugate prior on $\{\boldsymbol{\mu}^k\}$ and $\{\boldsymbol{z}^s\}$**

# Maximizing log likelihood using variational Bayes combined with point-estimation

- Log likelihood

$$L = \sum_{s=1}^{S} \sum_{n=1}^{N_s} \sum_{k=1}^{K} \ln \mathcal{N}(\boldsymbol{x}^{s(n)} \mid \boldsymbol{\mu}^k)^{z^{s(n)}} + \sum_{k=1}^{K} \mathrm{Lap}(\Lambda^k \mid \rho) p(\boldsymbol{\mu}^k \mid \Lambda^k) + \sum_{s=1}^{S} z^{s(n)} \ln \pi_k^s + \sum_{s=1}^{S} \ln p(\boldsymbol{\pi}^s)$$

Likelihood by the obs. model          Prior distributions

- Use VB for $\{\boldsymbol{\mu}^k\}, \{\boldsymbol{z}^{s(n)}\}$

- Use point-estimate for $\{\Lambda^k\}, \{\boldsymbol{\pi}^s\}$
  - Results in two convex optimization problems

# Maximizing log likelihood using variational Bayes combined with point-estimation

- Update sample weights

- Update cluster weights

- Update precision matrices

- Update other parameters

Use new semi-closed form solution

$$\max_{\boldsymbol{\pi}^s} \left\{ \sum_{k=1}^{K} c_k^s \ln \pi_k^s - \tau \|\boldsymbol{\pi}^s\|_0 \right\}$$

$$\text{s.t.} \quad \|\boldsymbol{\pi}^s\|_1 = 1.$$

The ratio of samples assigned to the *k*-th cluster

Solved by graphical lasso [Friedman 08]

$$\max_{\Lambda^k} \left\{ \ln |\Lambda^k| - \text{Tr}(\Lambda^k Q^k) - \frac{\rho}{N_k} \|\Lambda^k\|_1 \right\}$$

total # of samples assigned to the *k*-th cluster

# Solving the L0-regularized optimization problem for mixture weights

- What is the problem of the conventional VB approach?
  - Simply differentiate w.r.t. $\pi_k^s$
  - Claims to get a sparse solution [Corduneanu+ 01]
  - But mathematically $\pi_k^s$ cannot be zero due to logarithm

$$\max_{\boldsymbol{\pi}^s} \left\{ \sum_{k=1}^{K} c_k^s \ln \pi_k^s \right\}$$

$$\text{s.t.} \ \ \|\boldsymbol{\pi}^s\|_1 = 1.$$

- We re-formalized the problem as a convex mixed-integer programming
  - A semi-closed form solution can be derived (→ see paper)

$$\max_{\boldsymbol{\pi}^s, \boldsymbol{y}^s} \sum_{k=1}^{K} \{c_k^s \ln \pi_k^s - \tau y_k^s\} \ \ \text{s.t.} \ \ \sum_{k=1}^{K} \pi_k^s = 1,$$

$$y_k^s \geq \pi_k^s - \epsilon, \ \ y_k^s \in \{0,1\} \ \ \text{for} \ \ k = 1, \dots, K,$$

# Comparison with possible alternatives

| | | Interpretability | Noise reduction | Fleet-readiness | Multi-modality |
|---|---|---|---|---|---|
| **Our work [Ide et al. ICDM 17]** | | **Yes** | **Yes** | **Yes** | **Yes** |
| (single) sparse GGM | [Ide et al. SDM 2009, Ide et al. ICDM 2016] | **Yes** | **Yes** | No | No |
| Gaussian mixtures | [Yamanishi et al., 2000; Zhang and Fung, 2013; Gao et al., 2016] | Limited | Limited | No | **Yes** |
| Multi-task sparse GGM | [Varoquaux et al., 2010; Honorio and Samaras, 2010; Chiquet et al., 2011; Danaher et al., 2014; Gao et al., 2016; Peterson et al., 2015]. | **Yes** | **Yes** | **Yes** | No |
| Multi-task learning anomaly detection | [Bahadori et al., 2011; He et al., 2014; Xiao et al., 2015] | No | (depends) | **Yes** | No |

# Outline

- Problem setting

- Modeling strategy

- Model inference approach

- Experimental results

# Experiment (1): Learning sparse mixture weights

- Conventional ARD approach sometimes get stuck with local minima
  - ARD = automatic relevance determination
  - Often less sparser than the proposed convex L0 approach

- Typical result of log likelihood vs VB iteration count →

Proposed convex L0 approach gives better likelihood

Conventional ARD approach gets stuck with a local minimum



17

# Experiment (2): Learning GGMs and detecting anomalies

- "Anuran Calls" (frog voice) data in UCI Archive
  - o Multi-modal (multi-peaked)
  - o Voice signal + attributes (species, etc.)
- Created 10-variate, 3-task dataset
  - o Use the species of "Rhinellagranulosa" as the anomaly
- Results
  - o Two non-empty GGMs are automatically detected starting from K=9
  - o Clearly outperformed single-modal MTL alternative in anomaly detection
    - ✓ Group graphical lasso, fused graphical lasso

Example of variable-wise distribution



Automatically learned GGMs

# Conclusion

- Developed multi-task density estimation framework that can handle multi-modality
    - Featuring double sparsity: mixture weights, variable dependency

- Demonstrated the utility in the context of condition-based asset management

**Thank you!**

# Integrated monitoring tool allows sharing rare anomaly data across different assets

Anomalous: 0.2%

Normal: 99.8%

- In condition-based monitoring, big data may not be really big
  - Anomalous samples account for less than 0.2% in a metal smelting process

- Coverage of anomalies and thus accuracy can be limited due to lack of data

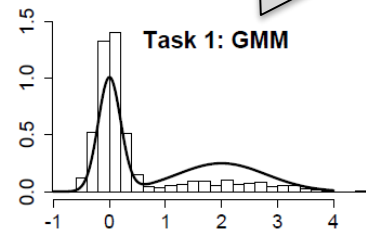# Existing methods cannot handle multi-modality



Comparing the proposed multi-task multi-modal (MTL-MM) model with standard Gaussian mixture (GMM) and multi-task learning (MTL) models