

# Statistical Machine Learning Approaches to Change Detection

Song Liu, JSPS DC2,  
Sugiyama Lab, Tokyo Institute of Technology  
2014/3/10

# Song Liu (Tokyo Tech)

- Graduating PhD student (3<sup>rd</sup> year) from Tokyo Tech
  - Thesis: Statistical Machine Learning Approaches to Change Detection
  - Fellow of Japan Society for Promotion of Science (JSPS)
- Visiting at National Institute of Informatics, JP
- MSc in Advanced Computing (University of Bristol)
- BEng in Computer Science (Soochow University, CN)

# Learning from Data

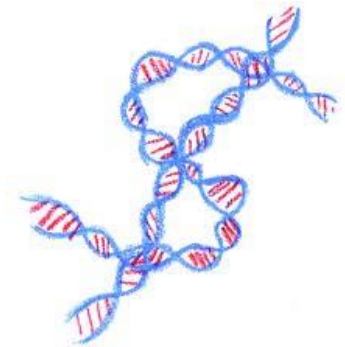
- The era of big data is coming.



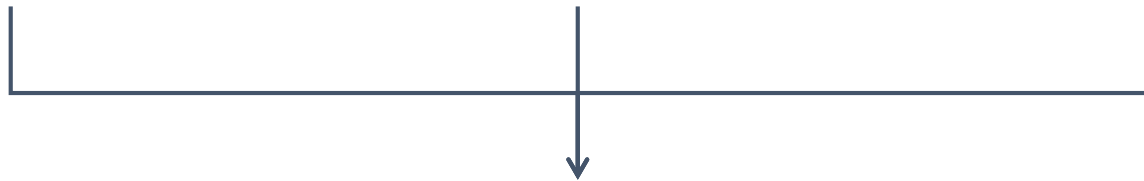
1 hour video per second



1M transaction per hour



2.5 petabytes per 1000 person

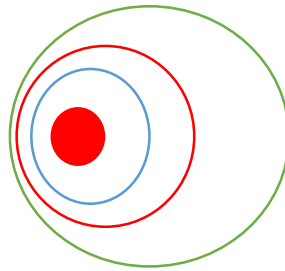


- Without interpretation, raw data make no sense to human.
- Machine learning helps make sense out of data.

# The Human Learning in Astronomy

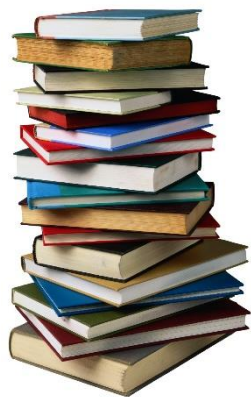
## A brief history

- Learning Constellations
  - Patterns of Stars
  - Predicting warfare, harvest
- Patterns of Planetary Movements
  - Tycho Brahe's data
  - Kepler's law
- Newton Laws
  - Predicting the movements of new planets



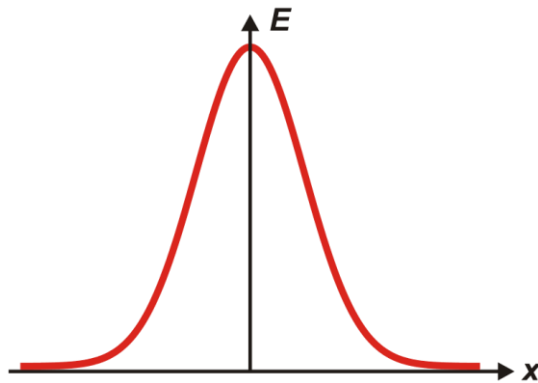
$$F = ma$$

# Statistical Machine Learning



A photograph of a piece of paper with the handwritten equation 
$$m = \frac{m_0 c^2}{\sqrt{1 - \frac{v^2}{c^2}}}$$
 written in black ink.

Data is so big, so we look for compressive rules (**patterns**).



Data involving uncertainty, we prefer **statistical** patterns .

# However, Data is Always Changing!

- Patterns learned today may be useless tomorrow.



Everyday, 20% Google queries have never been seen before<sup>1</sup>.

- **Learning pattern changes** are also important.

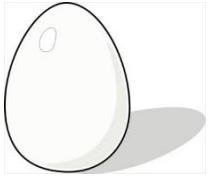


Red spots are detected changes on NASA satellite images, reflecting the heaviest damages after super typhoon Haiyan struck Philippines on Nov. 8 2013<sup>2</sup>.

- We need a **new paradigm** of machine learning.

1. <http://certifiedknowledge.org/blog/are-search-queries-becoming-even-more-unique-statistics-from-google>  
2. <http://www.nasa.gov/content/goddard/haiyan-northwestern-pacific-ocean/#.Uq3u6vQW2So>

# Two Paradigms of Machine Learning



**Static Learning:** Data do not change

(p.5)

**Supervised**

Predict Future:  $p(y|x)$

**Classification, Regression**

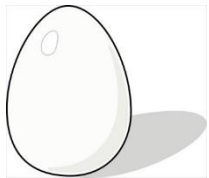
(Duda et al., 2001, Bishop, 2006)

**Unsupervised**

Learn Patterns:  $p(x)$

**Clustering, Anomaly Detection**

(Murphy, 2012, Shaw-Taylor & Cristianini, 2004)



**Dynamic Learning:** Data are shifting

**Supervised**

Predict future under data shifts

$$p(y'|x'), p \neq p'$$

e.g. **Transfer learning** (Sugiyama et al., 2012)

**Unsupervised**

Learn changes of patterns

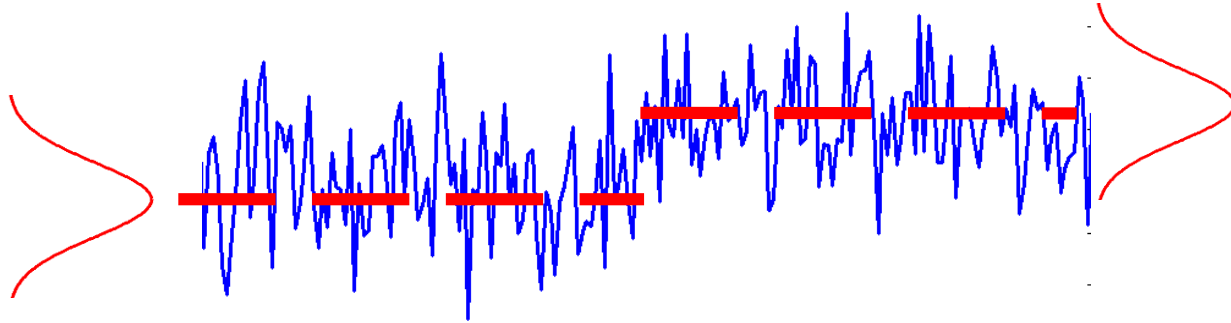
$$p/q \text{ or } p - q, p \neq q$$

**Change detection**

# Two Ways of Learning Changes

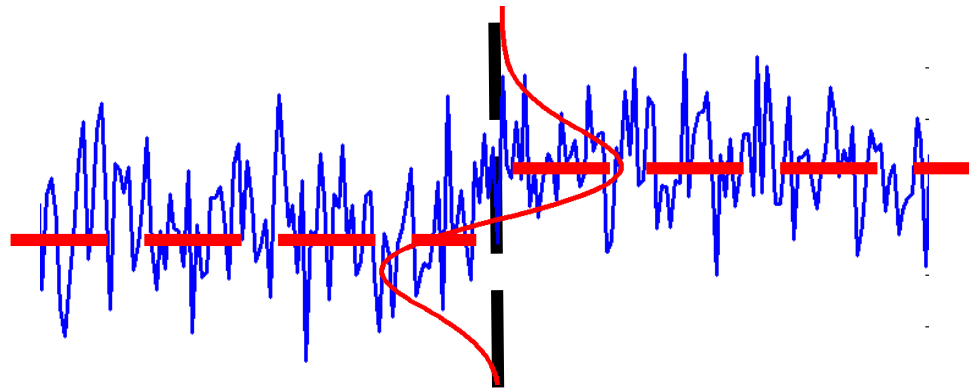
## Separated Learning Approach

①



## Direct Learning Approach

②



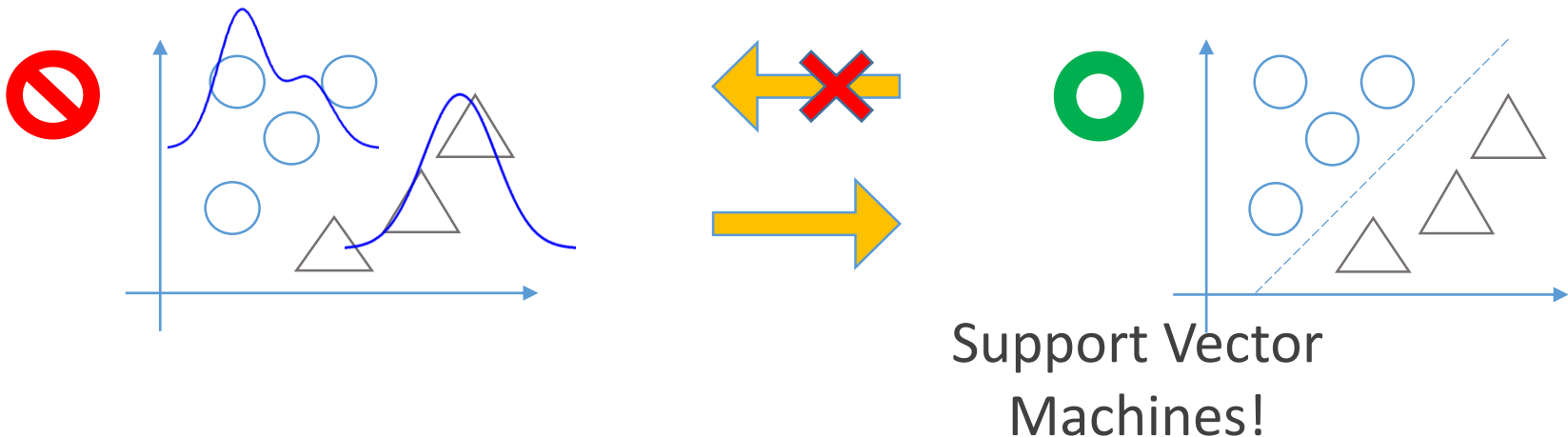
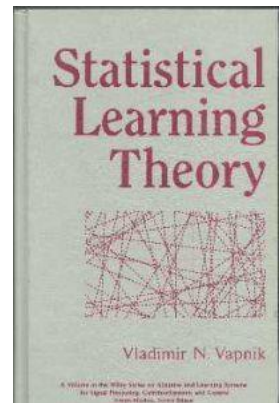
Which one is the best ?



# Vapnik's Principle

(Vapnik, Statistical Learning Theory, 1998)

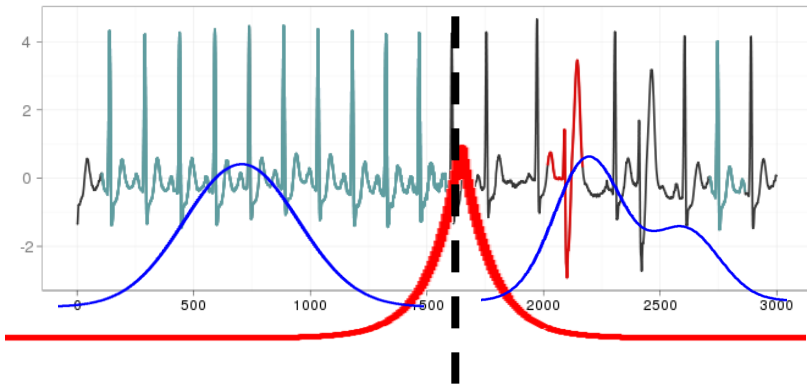
- Generally, we prefer ② to ① because of **Vapnik's principle**.
- “When solving a problem of interest, one should not solve a more general problem as an intermediate step .”



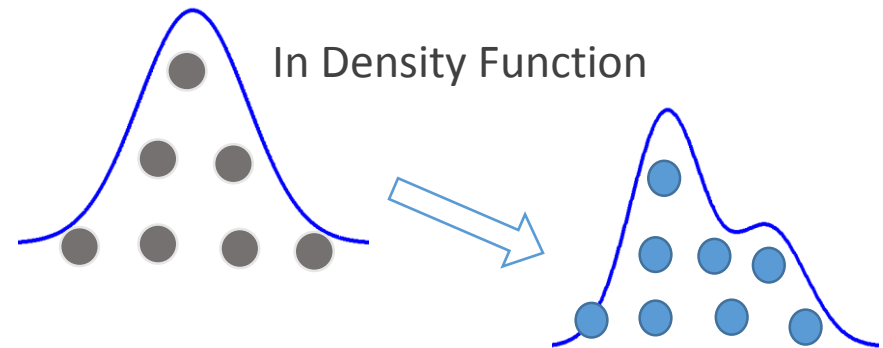
# Unsupervised Change Detection

(p.10)

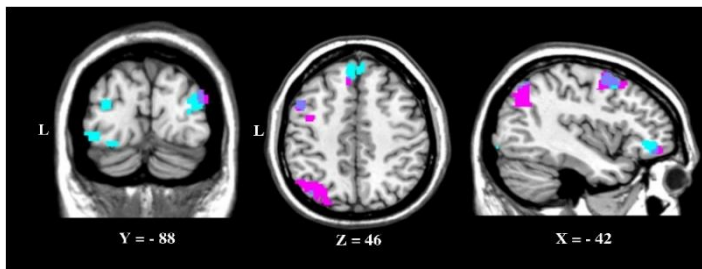
- Has it been changed? **Distributional Change Detection**



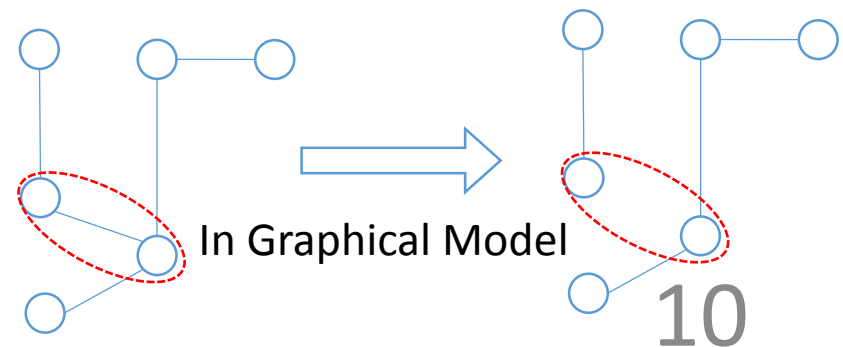
Heartbeat disorder detection, from JMOTIF project.



- What has been changed? **Structural Change Detection**



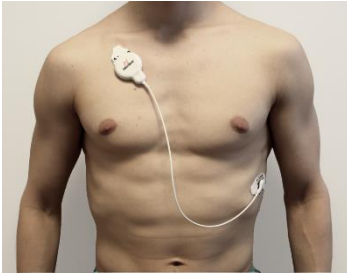
Changes in brain activation (Cléry et al., 2006)



# Chapter 2, Distributional Change Detection

1. Background & Existing Methods
2. Problem Formulation
3. Divergence based Change-point Score
4. Experiments

# Changes in Time-series



## Health monitoring

e.g. Kawahara et al., 2012



## Engineering fault detection

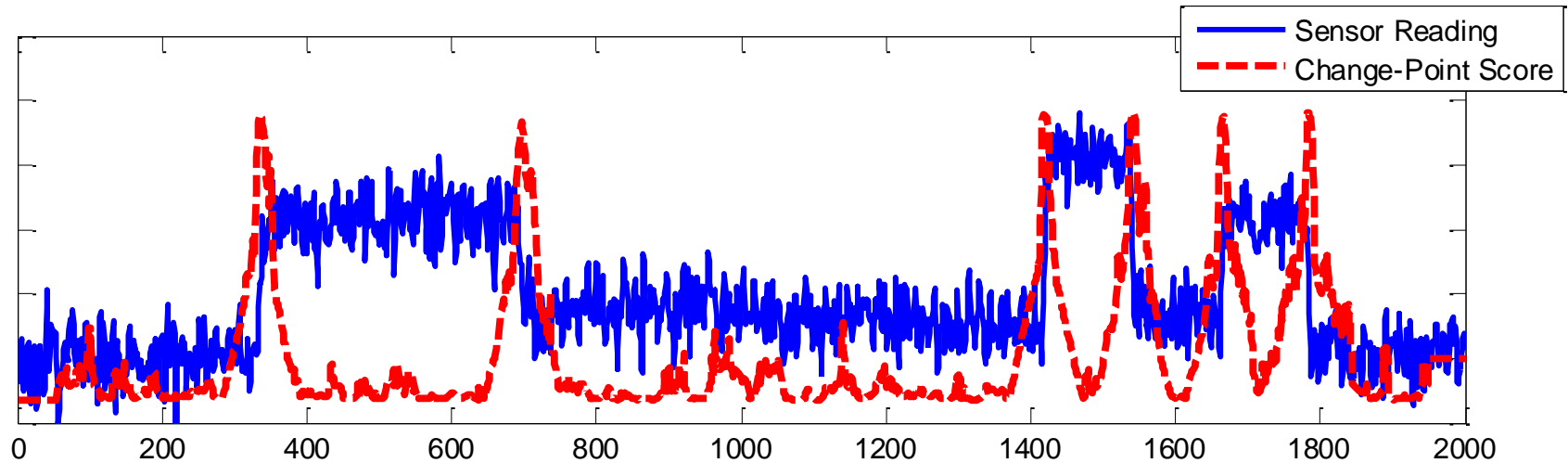
e.g. Keogh et al., 2005



## Network intrusion detection

e.g. Takeuchi et al., 2006

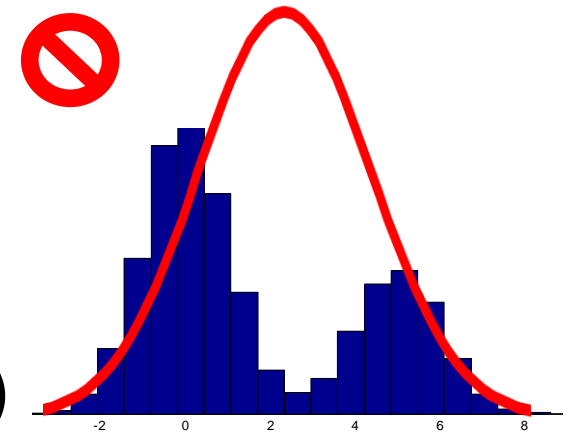
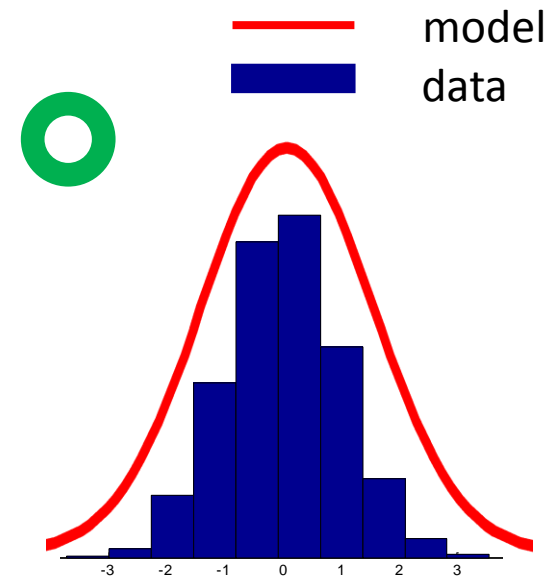
# Distributional Change Detection in Time-series



- **Objective:** Detecting **abrupt changes** lying among time-series data
- **Change-point score:** Plausibility of changes that have happened
- **Metrics:** True positive rate, false positive rate, and detection delay

# Existing Methods

- Model based Methods:
  - **Auto Regressive (AR)**
    - (Takeuchi and Yamanishi, 2006)
  - Autoregressive model with Gaussian noise.
  - **Singular Spectrum Transformation (SST)**
    - (Moskvina and Zhigljavsky, 2003)
  - **Subspace Identification (SI)**
    - (Kawahara et al., 2007)
  - State Space model with Gaussian noise.
- Model-free Methods:
  - **One-class Support Vector Machine (OSVM)**
    - (Desobry et al., 2005)
  - **Density Ratio based Methods (KLIEP)**
    - KLIEP (Kawahara & Sugiyama, 2009)
  - No model assumption.



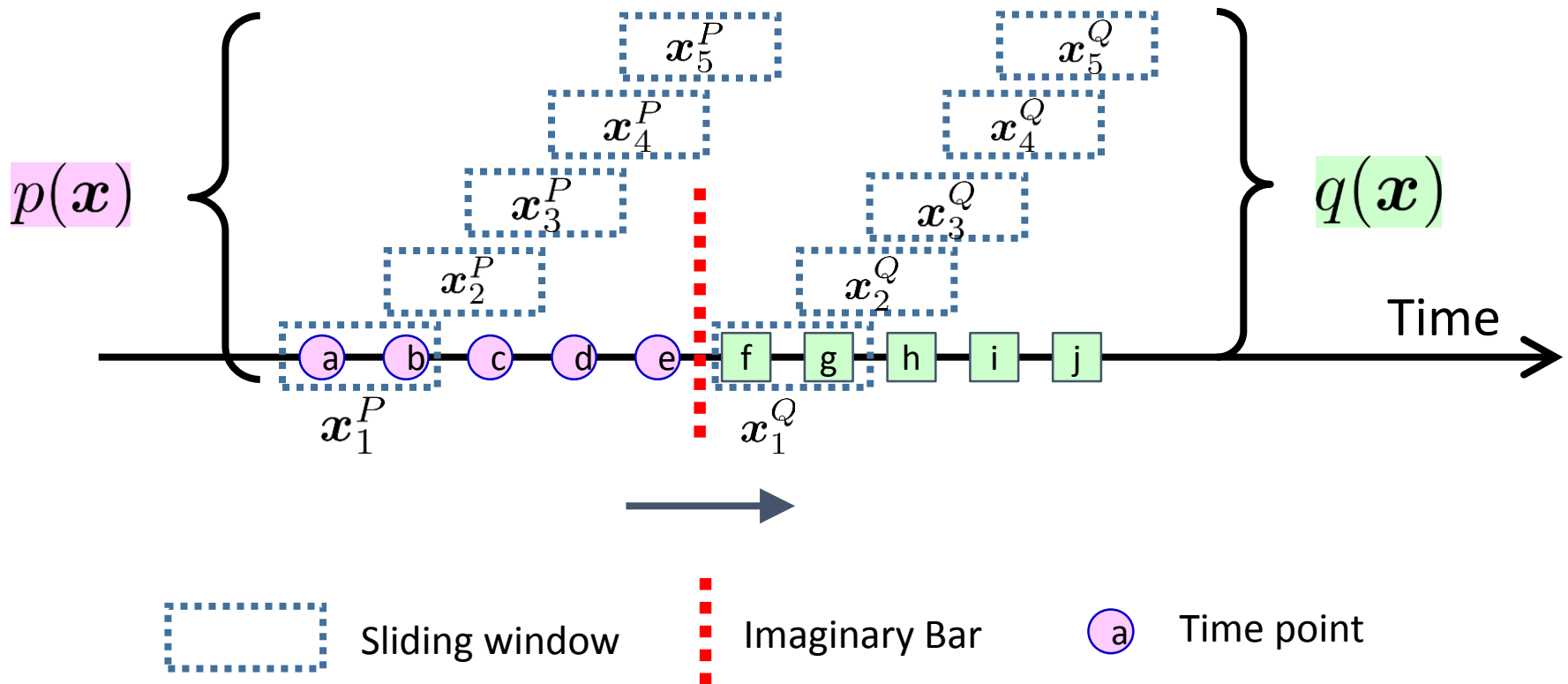
Model based methods fail to characterize data when model mismatches.

# Chapter 2, Distributional Change Detection

1. Background & Existing Methods
2. Problem Formulation
3. Proposed Method
4. Experiments

# Formulate Problem from Time-series

- Construct samples by using sliding window.
- An imaginary bar in the middle divides samples into two groups.
- Assume two groups of samples are from  $p$  and  $q$ .



How to generate a change point score?



# Chapter 2, Distributional Change Detection

1. Background & Existing Methods
2. Problem Formulation
3. Proposed Method
4. Experiments

# Divergence based Change Point Detection

- Need a **measure of change** between  $p(\mathbf{x})$  and  $q(\mathbf{x})$ .
- $f$ -divergence is a **distance** from  $p$  to  $q$ .

(Ali and Silvey, 1966)

$$D[p||q] = \int f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) q(\mathbf{x}) \, d\mathbf{x}$$

$f$  is convex  
 $f(1) = 0$

- $f$  is a function of a ratio between two **probability density functions** ( $p$  and  $q$ ).
- $D[p||q]$  is asymmetric, thus we symmetrized it:

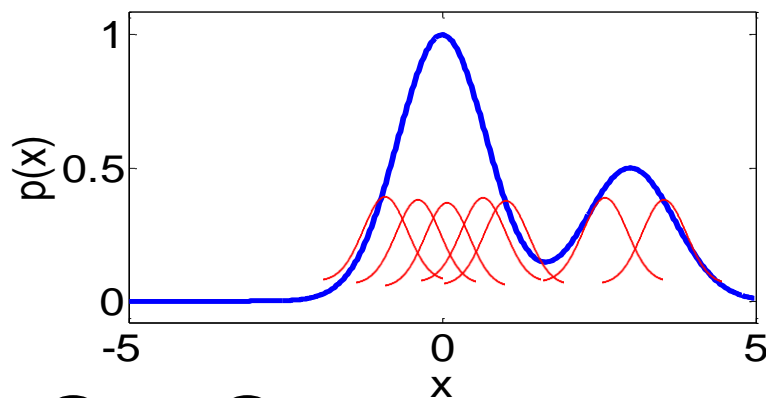
$$D[p||q] + D[q||p]$$

# Density Estimation


- A naive way is a two step approach:

$$\left. \begin{array}{l} \textcircled{1} p(\mathbf{x}) \approx \hat{p}(\mathbf{x}) \\ \textcircled{2} q(\mathbf{x}) \approx \hat{q}(\mathbf{x}) \end{array} \right\} \frac{p(\mathbf{x})}{q(\mathbf{x})} \approx \frac{\hat{p}(\mathbf{x})}{\hat{q}(\mathbf{x})}$$

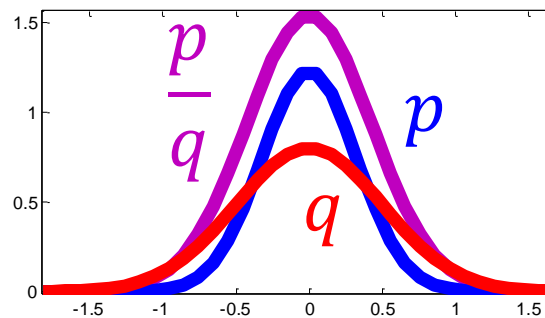
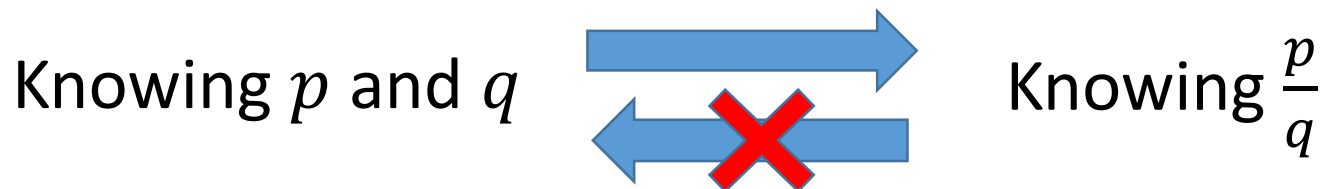
- e.g. Kernel Density Estimation (Csörgo & Horváth 1988)



$$\hat{p}(\mathbf{x}; \sigma) = \frac{1}{n} \sum_{i=1}^n K_{\sigma}(\mathbf{x}, \mathbf{x}_i)$$
$$K_{\sigma}(\mathbf{x}, \mathbf{x}') = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{-2\sigma^2}\right)$$

-   $\textcircled{1}$  or  $\textcircled{2}$  is carried out without taking care of the ratio.
  - Error may be magnified after the division.

# Density Ratio Estimation



Vapnik's Principle!

- **Direct** Density Ratio Estimation (Sugiyama et al., 2012).

$$\frac{p(\mathbf{x})}{q(\mathbf{x})} \approx \hat{r}(\mathbf{x}; \boldsymbol{\theta}) = \sum_i^n \theta_i K_\sigma(\mathbf{x}, \mathbf{x}_i)$$

- By plugging in estimated density ratio, we may obtain various  $f$ -divergences.

# Kullback-Leibler Divergence as Change-point Score (p.22)


We may have different choices for change-point score:

- Kullback-Leibler Divergence (Kullback and Leibler, 1951)

$$\text{KL}[p||q] = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

- Has been applied, and promising results were obtained.  
Kawahara & Sugiyama, SDM 2009
- Density ratio can be obtained via **Kullback-leibler Importance Estimation Procedure (KLIEP)** (Sugiyama, et al., NIPS 2008).

 Optimization for density ratio is convex!

 Slow, less robust

# Pearson Divergence as Change-point Score (pp.24)

- The log term in Kullback-leibler divergence can go infinity if  $\frac{p(\mathbf{x})}{q(\mathbf{x})} = 0$ .
- Pearson Divergence (Pearson, 1900)

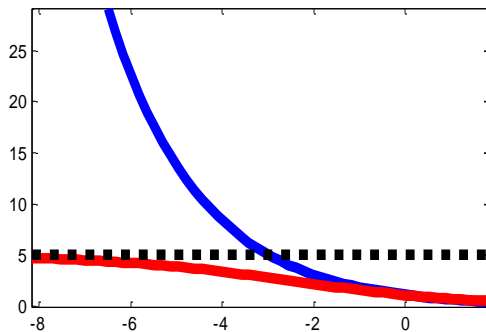
$$\text{PE}[p||q] = \frac{1}{2} \int q(\mathbf{x}) \left( \frac{p(\mathbf{x})}{q(\mathbf{x})} - 1 \right)^2 d\mathbf{x}$$

- Density ratio can be obtained via **unconstrained Least Square Importance Fitting (uLSIF)** (Kanamori, et al., JMLR 2009).
- Solution has an analytic form.

 Fast, relatively robust!

# Relative Pearson Divergence as Change-point Score (pp.26)

- The density ratio may go unbounded.
- Bounded Density Ratio  $\rightarrow$  Relative Density Ratio



$$\text{---} \frac{p}{q} \quad \text{---} \frac{p}{\alpha p + (1 - \alpha)q} \leq \frac{1}{\alpha}$$

- Relative Pearson Divergence (Yamada, et al. , NIPS 2011)

$$\begin{aligned} \text{PE}_\alpha[p||q] &= \text{PE}[p||\alpha p + (1 - \alpha)q] \\ &= \frac{1}{2} \int (\alpha p(\mathbf{x}) + (1 - \alpha)q(\mathbf{x})) \left( \frac{p(\mathbf{x})}{\alpha p(\mathbf{x}) + (1 - \alpha)q(\mathbf{x})} - 1 \right)^2 d\mathbf{x} \end{aligned}$$

More robust!

# Divergences as Change-point Scores (Summary) (pp.22-28)

- **Kullback-Leibler Divergence (via KLIEP) :**

$$\widehat{\text{KL}} := \frac{1}{n} \sum_{i=1}^n \log \hat{r}(\mathbf{x}_i)$$

Existing method

Previous Work: Kawahara & Sugiyama, 2009

- **Pearson Divergence (via uLSIF):**

$$\widehat{\text{PE}} := -\frac{1}{2n} \sum_{i=1}^n \hat{r}(\mathbf{x}'_i)^2 + \frac{1}{n} \sum_{i=1}^n \hat{r}(\mathbf{x}_i) - \frac{1}{2}$$

- **Relative Pearson Divergence (via RuLSIF):**

$$\widehat{\text{PE}}_\alpha = -\frac{\alpha}{2n} \sum_{i=1}^n \hat{r}(\mathbf{x}_i)^2 - \frac{1-\alpha}{2n} \sum_{i=1}^n \hat{r}(\mathbf{x}'_i)^2 + \frac{1}{n} \sum_{i=1}^n \hat{r}(\mathbf{x}_i) - \frac{1}{2}$$

Robust  
Change  
Detection



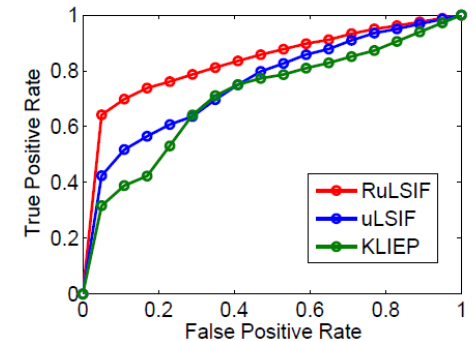
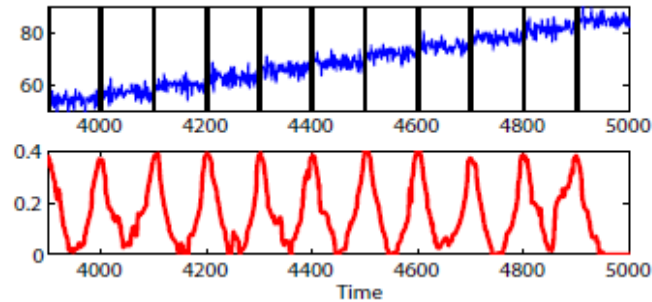
# Chapter 2, Distributional Change Detection

1. Background & Existing Methods
2. Problem Formulation
3. Proposed Method
4. Experiments

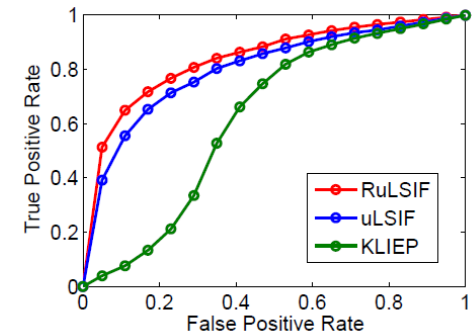
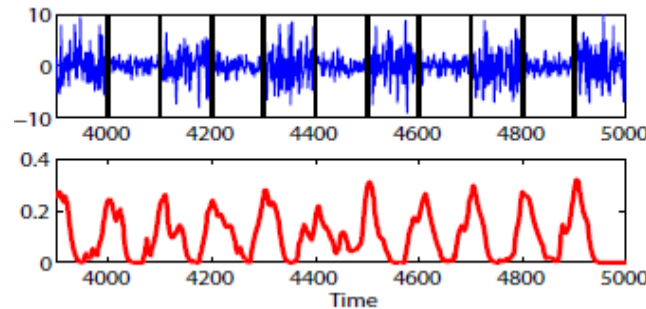
# Experiments (Synthetic)

— original signal  
— score  
— ground truth

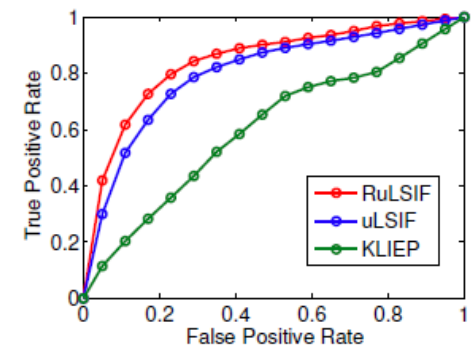
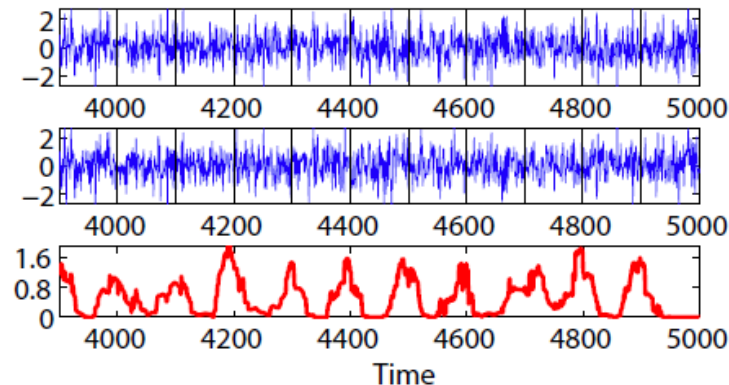
## Mean Shift



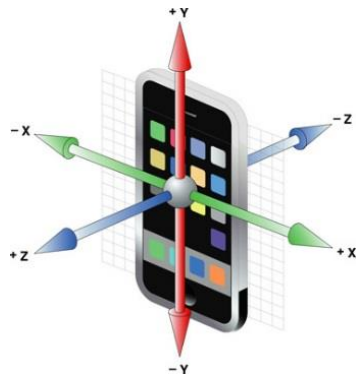
## Variance Scaling



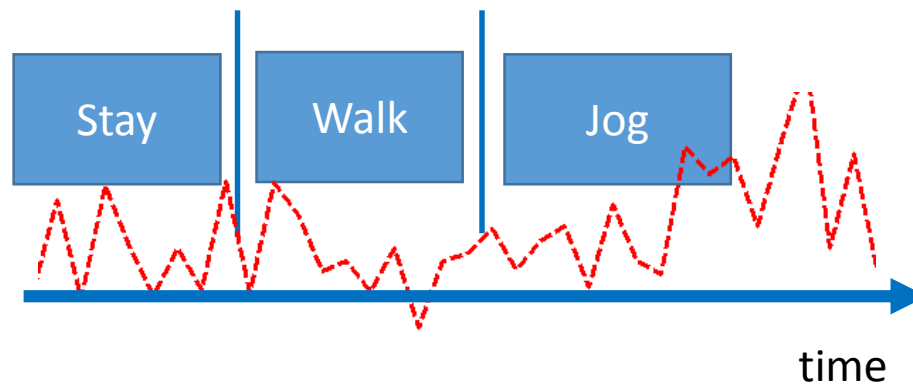
## Covariance Switching



# Experiments (Sensor)

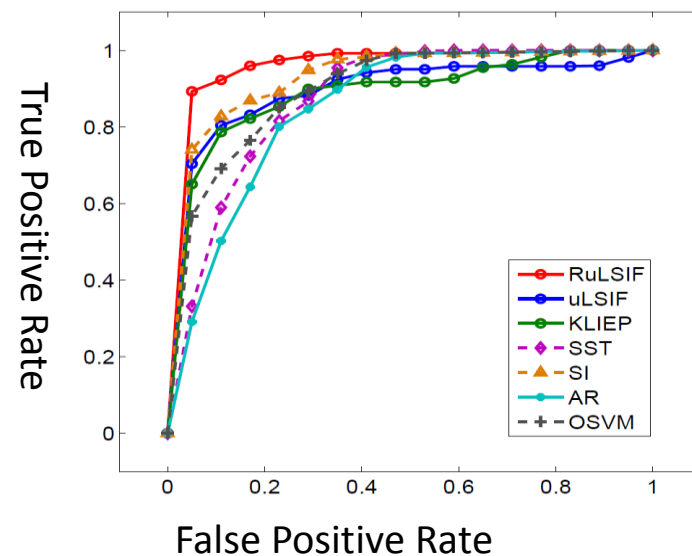
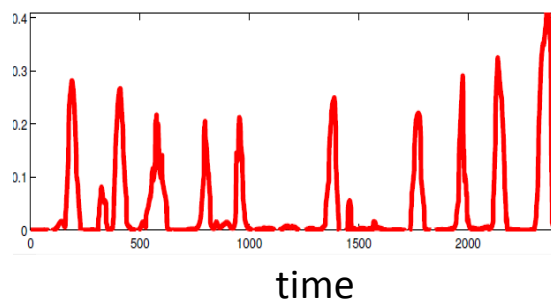
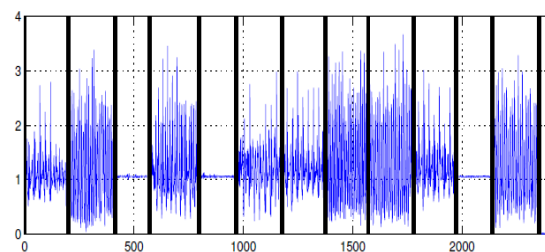


HASC 2011 Challenge

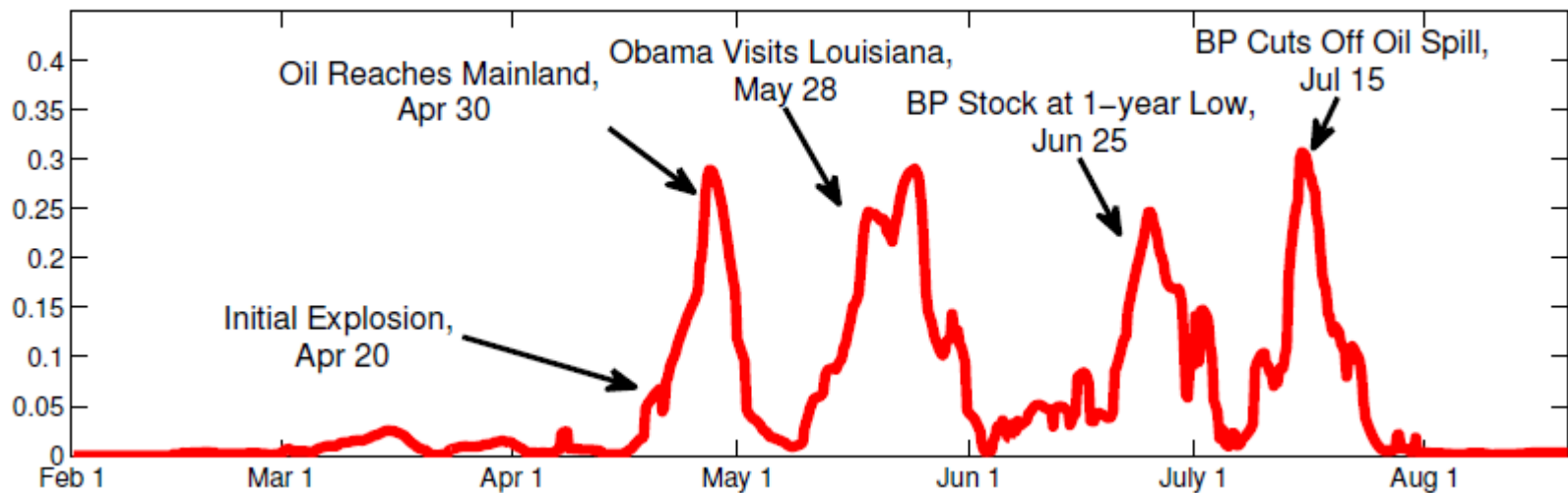
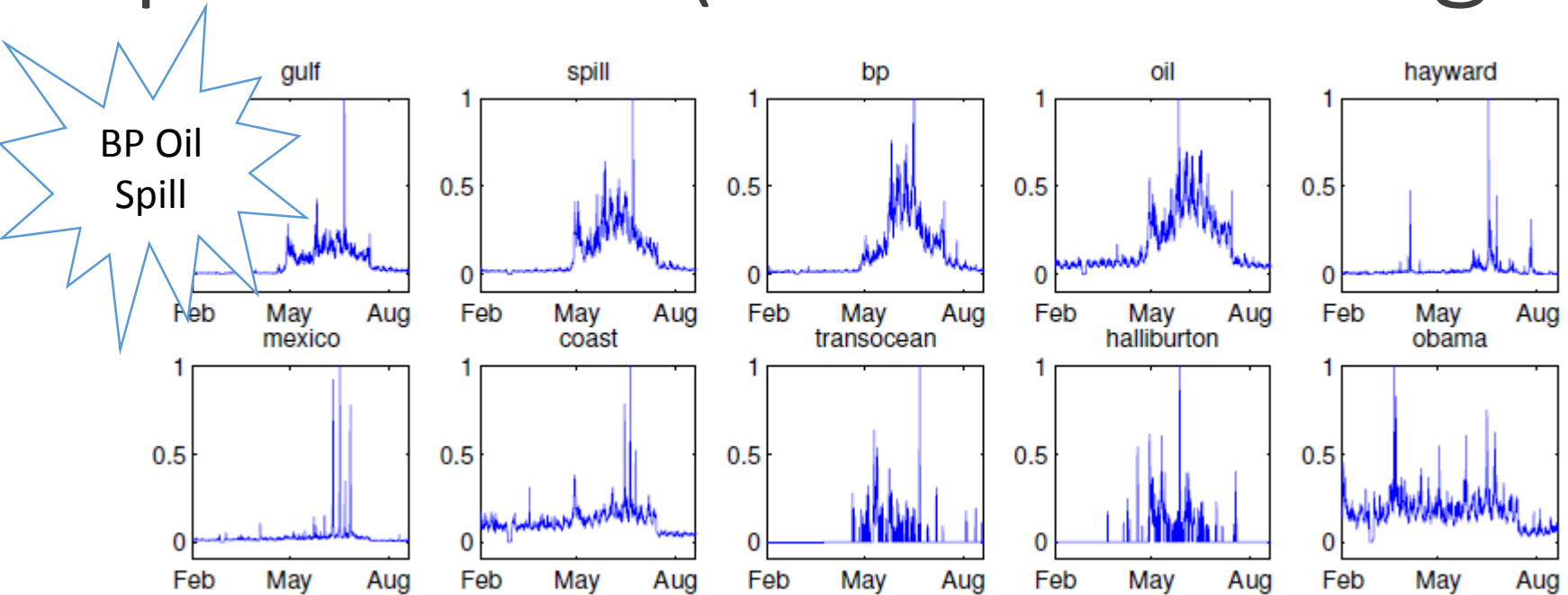


**Mobile *Sensor* data**, measured in 20 seconds, recording human activity.

**Task:** Segmenting the activities.  
e.g. walk, jogging, stairs up/down...



# Experiments (Twitter Messages)



# Conclusion

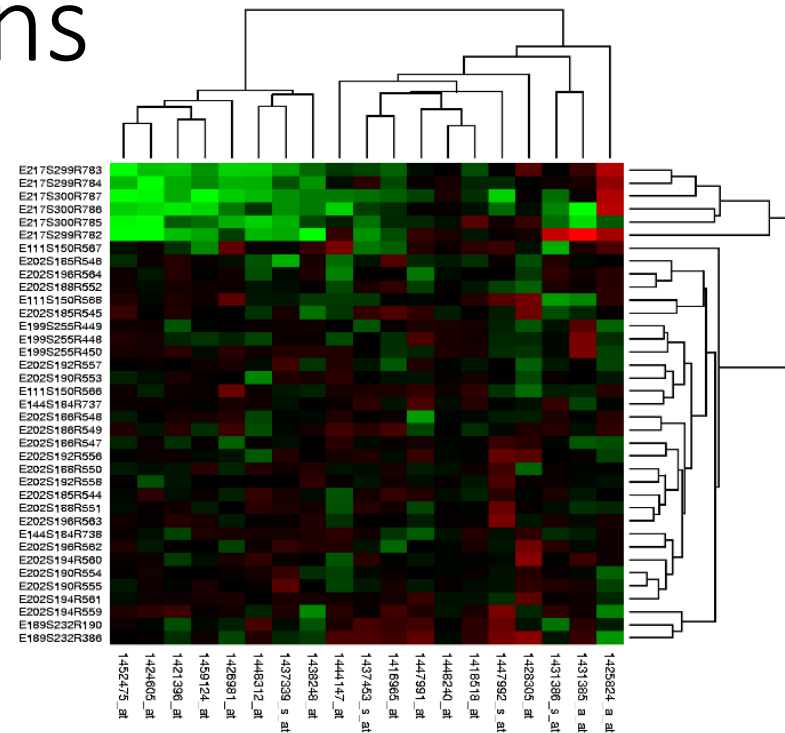
- A **robust** and **model-free** method for distributional change detection in time-series data was proposed.
- The proposed method obtained the **leading** performance in various artificial and real-world experiments.

# Chapter 3, Structural Change Detection

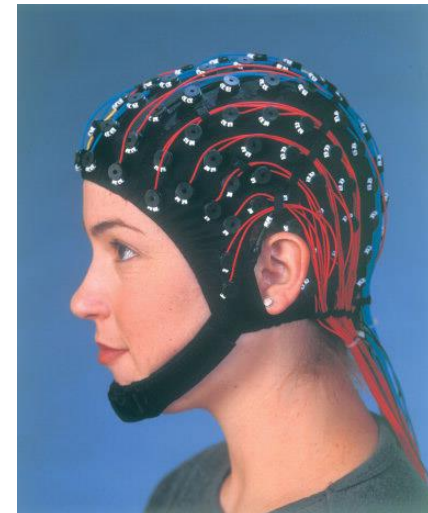
1. Background
2. Separate Estimation Methods
3. Proposed Methods
4. Experiments

# Patterns of Interactions

- Genes regulate each other via gene network.
- Brain EEG signals may be synchronized in a certain pattern.
- However, such patterns of interactions may be **changing!**



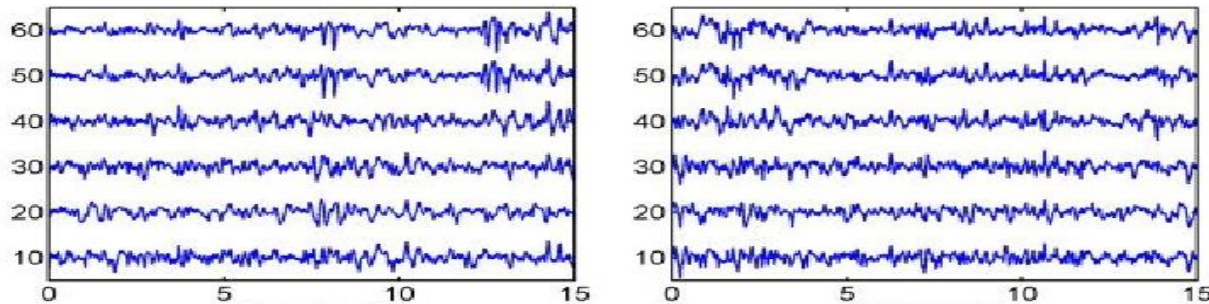
(Wikipedia)



# Structural Change Detection

- Interactions between features may change.  $\mathbf{x} = (x^{(1)}, \dots, x^{(d)})^T$

$$\{\mathbf{x}_i^P\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}) \qquad \{\mathbf{x}_i^Q\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} q(\mathbf{x})$$



(Williamson et al., 2012)

The change of brain signal correlation at two different experiment intervals.

- e.g. Interactions between some genes may be activated,
  - but only under some conditions.
- “apple” may co-occur with “banana” quite often in cookbook,
  - but **not in IT news**.



# Chapter 3, Structural Change Detection

1. Background
2. Existing Methods
3. Proposed Methods
4. Experiments

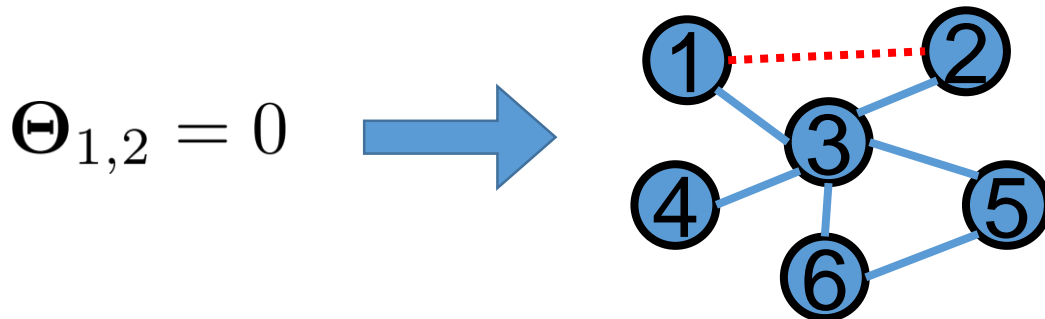
# Gaussian Markov Networks (GMNs)

- The interactions between random variables can be modelled by **Markov Networks**.
- **Markov Networks (MNs)** are **undirected** Graphical Models.
- The simplest example of MN is a Gaussian MN:

$$p(\mathbf{x}; \Theta) = \frac{\det(\Theta)^{1/2}}{(2\pi)^{d/2}} \exp \left( -\frac{1}{2} \mathbf{x}^\top \Theta \mathbf{x} \right)$$

$\mathbf{x} = (x^{(1)}, x^{(2)}, \dots, x^{(d)})$        $\Theta$  is the inverse covariance matrix

- We can visualize the above MN using an undirected graph.



# Estimating Separate GMN

Tibshirani, JRSS 1996; Friedman et al., Biostatistics 2008

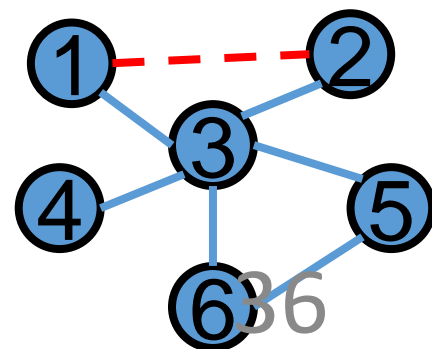
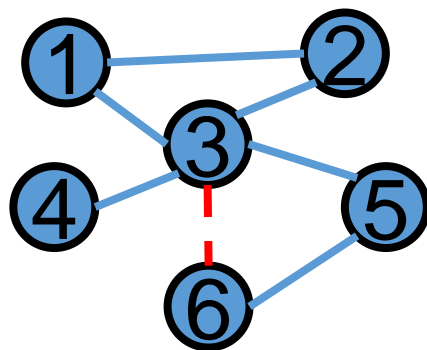
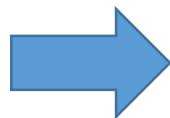
- Recall, we would like to detect changes between MNs.
- Changes can be found once the structure of two **separate** GMNs are known to us.
- Estimating Sparse GMNs can be done via **Graphical Lasso (Glasso)**.

$$\max_{\Theta^P} \sum_{i=1}^n \log p(\mathbf{x}_i^P; \Theta^P) - \lambda^P \|\Theta^P\|_1$$

• Idea: Twice Glasso, take the parameter difference.



$$\Theta_{u,v}^P - \Theta_{u,v}^Q$$



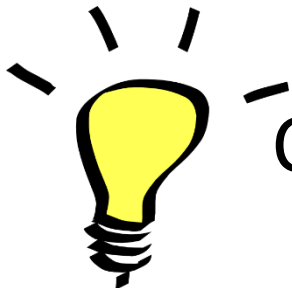
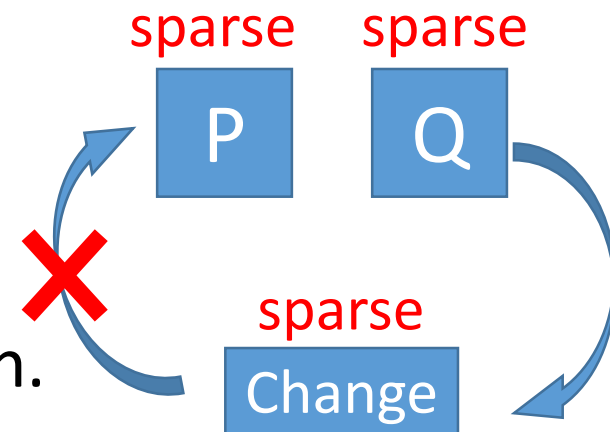
# Glasso: Pros and Cons

## 😊 Pros:

- Statistical properties are well studied.
- Off-the-shelf software can be used.
- **Sparse change is produced.**

## 😞 Cons:

- Cannot detect high-order correlation.
- Not clear how to choose hyper-parameters.
- **Does not work if  $p$  or  $q$  is dense.**



Can we combine two optimizations into one?

# Semi-direct Approach (Fused-lasso)

- We can impose sparsity directly on  $\Theta_{u,v}^P - \Theta_{u,v}^Q$ , using **Fused-lasso (Flasso)** (Tibshirani et al., 2005).

- Consider the following objective:

$$\max_{\theta} \ell(\mathbf{x}^P, \Theta^P) + \ell(\mathbf{x}^Q, \Theta^Q) - \lambda \|\Theta^P - \Theta^Q\|_1$$

Gaussian

Gaussian

(Zhang & Wang, **UAI2010**)

- ☺ We don't have to assume  $p$  or  $q$  is sparse.
- ☺ Sparsity control is much easier than Glasso.
- ☹ Gaussianity is still assumed.

# Nonparanormal (NPN) Extension

Liu et al., JMLR 2009

- Gaussian assumption is too restrictive.
- However, considering general non-Gaussian model can be computationally challenging.
- We may consider a method half-way between Gaussian and non-Gaussian.
- Non-paranormal (NPN).

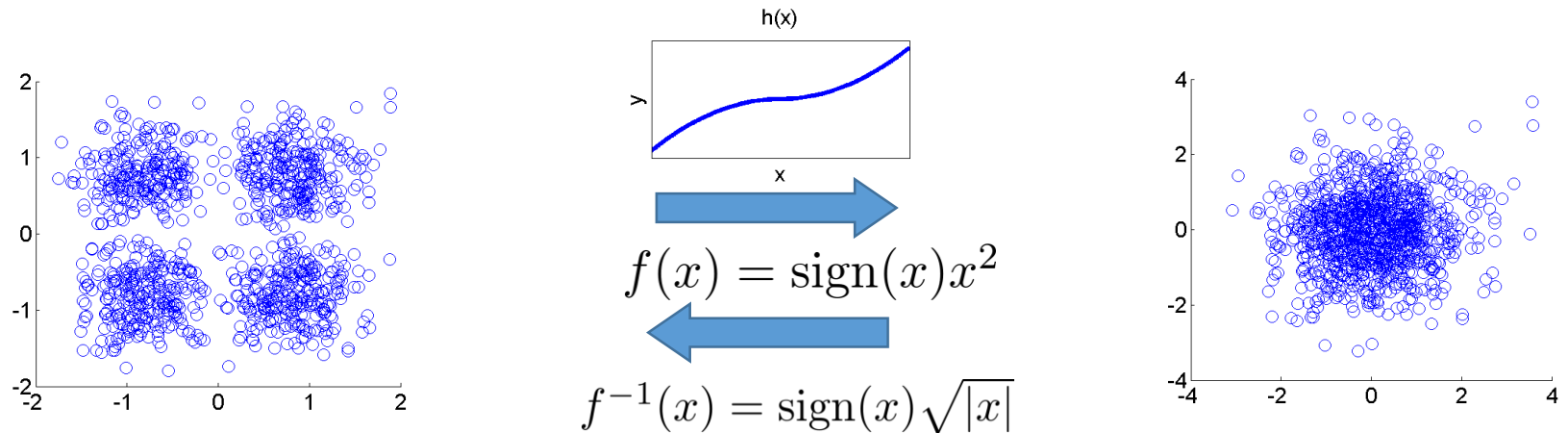
$$p(\mathbf{x}; \Theta) = \frac{\det(\Theta)^{1/2}}{(2\pi)^{d/2}} \exp \left( -\frac{1}{2} \mathbf{f}(\mathbf{x})^\top \Theta \mathbf{f}(\mathbf{x}) \right) \prod_{i=1}^d |f'_i(x^{(i)})|$$

$$\mathbf{x} = (x^{(1)}, x^{(2)}, \dots, x^{(d)})$$

$$\mathbf{f}(\mathbf{x}) = (f_1(x^{(1)}), f_2(x^{(2)}), \dots, f_d(x^{(d)}))$$

$f_k$  : Monotone, differentiable function

# Nonparanormal (NPN) Extension



😊 More flexible than Gaussian methods, still tractable.

😞 However, NPN extension is still restrictive.

# Non-Gaussian Log-linear Model

- Pairwise Markov Network

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left( \sum_{u \geq v} \boldsymbol{\theta}_{u,v}^\top \mathbf{f}(x^{(u)}, x^{(v)}) \right)$$

- $\mathbf{f}$  are feature vectors.

$$\mathbf{f}: \mathcal{R}^2 \rightarrow \mathcal{R}^b$$

Gaussian:  $f_{gau}(x, y) = xy$

Nonparanormal:  $f_{npn}(x, y) = f(x)f(y)$

Polynomial:  $\mathbf{f}_{poly}(x, y) = [x^k, y^k, x^{k-1}y, \dots, x, y, 1]$

- The normalization term  $Z(\boldsymbol{\theta})$  is generally intractable.
  - Gaussian or Nonparanormal models are **exceptions**.



# The Normalization Issue

$$Z(\boldsymbol{\theta}) = \int \exp \left( \sum_{u \geq v} \boldsymbol{\theta}_{u,v}^\top \mathbf{f}(x^{(u)}, x^{(v)}) \right) d\mathbf{x} \quad !$$

- However, for a generalized Markov Network, there is no closed-form for  $Z(\boldsymbol{\theta})$ .
- Importance Sampling (IS):

$$Z(\boldsymbol{\theta}) = \int p_{\text{inst}}(\mathbf{x}) \frac{\exp \left( \sum_{u \geq v} \boldsymbol{\theta}_{u,v}^\top \mathbf{f}(x^{(u)}, x^{(v)}) \right)}{p_{\text{inst}}(\mathbf{x})} d\mathbf{x}$$

- IS can be used to extend both Glasso and Flasso.
  - IS-Glasso, IS-Flasso.

# Chapter 3, Structural Change Detection


1. Background
2. Existing Methods
3. Proposed Methods
4. Experiments

# Modeling Changes Directly

- Recall, our interest is:  $\beta_{u,v} := \theta_{u,v}^P - \theta_{u,v}^Q$
- The ratio of two MNs naturally incorporates the  $\theta_{u,v}^P - \theta_{u,v}^Q$ !

$$\frac{p(x; \theta^P)}{q(x; \theta^Q)} \propto \exp \left( \sum_{u \geq v} (\theta_{u,v}^P - \theta_{u,v}^Q)^\top \mathbf{f}(x^{(u)}, x^{(v)}) \right)$$




$$\exp \left( \sum_{u \geq v} \beta_{u,v}^\top \mathbf{f}(x^{(u)}, x^{(v)}) \right)$$

So, model the **ratio** directly!

# Modeling Changes Directly

- Density ratio model:

$$r(\mathbf{x}; \boldsymbol{\beta}) = \frac{1}{N(\boldsymbol{\beta})} \exp \left( \sum_{u \geq v} \boldsymbol{\beta}_{u,v}^\top \mathbf{f}(x^{(u)}, x^{(v)}) \right)$$

- The normalization term is:

$$N(\boldsymbol{\beta}) = \int q(\mathbf{x}) \exp \left( \sum_{u \geq v} \boldsymbol{\beta}_{u,v}^\top \mathbf{f}(x^{(u)}, x^{(v)}) \right) d\mathbf{x}$$

To ensure:

$$\hat{p} = q(\mathbf{x}) r(\mathbf{x}, \hat{\boldsymbol{\beta}})$$



$$\int r(\mathbf{x}; \boldsymbol{\beta}) q(\mathbf{x}) d\mathbf{x} = 1$$



**Sample average approximation**

Also works when **integral has no closed form!**


$$\frac{1}{n} \sum_i \exp \left( \sum_{u \geq v} \boldsymbol{\beta}_{u,v}^\top \mathbf{f}(x^{(u)Q}, x^{(v)Q}) \right)$$

# Estimating Density Ratio

Sugiyama et al., NIPS 2007

- Kullback-Leibler Importance Estimation Procedure (KLIEP):

$$\begin{aligned}\operatorname{argmin}_{\theta} \text{KL}[p||\hat{p}] &= \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{\hat{p}(\mathbf{x})} d\mathbf{x} \\ &= \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})\hat{r}(\mathbf{x})} d\mathbf{x} \quad \hat{p}(\mathbf{x}) = q(\mathbf{x})\hat{r}(\mathbf{x}) \\ &= \text{Const.} - \int p(\mathbf{x}) \log \hat{r}(\mathbf{x}; \beta) d\mathbf{x}\end{aligned}$$


$$\operatorname{argmax}_{\theta} \ell_{\text{KLIEP}}(\beta) = \frac{1}{n} \sum_{i=1}^n \log \hat{r}(\mathbf{x}_i; \beta)$$

**Unconstrained convex optimization!**

Tsuboi et al, JIP 2009

# Sparse Regularization

L2 regularizers    Group lasso regularizer

$$\max_{\boldsymbol{\beta}} \left[ \ell_{\text{KLIEP}}(\boldsymbol{\beta}) - \lambda_1 \underbrace{\|\boldsymbol{\beta}\|_2^2}_{\text{Elastic Net}} - \lambda_2 \underbrace{\sum_{u \geq v} \|\boldsymbol{\beta}_{u,v}\|_2}_{\text{Group lasso regularizer}} \right]$$

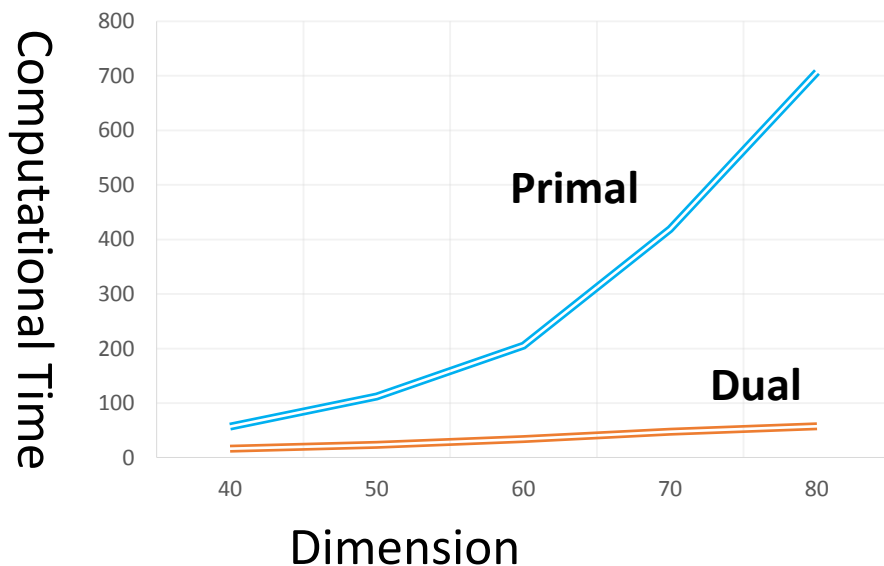
- Impose sparsity constraints on each factor  $\boldsymbol{\beta}_{u,v}$ .
  - equals to impose sparsity on changes.
- So finally, we can obtain a  $\boldsymbol{\beta}$  with group sparsity!

# The Dual Formulation p. 57

- When dimensionality is high, the dual formulation is preferred.

$$\min_{\alpha = (\alpha_1, \dots, \alpha_{n_Q})^\top} \sum_{i=1}^{n_Q} \alpha_i \log \alpha_i + \frac{1}{\lambda_1} \sum_{u \geq v} \max(0, \|\xi_{u,v}\| - \lambda_2)^2$$

subject to  $\alpha_1, \dots, \alpha_{n_Q} \geq 0$  and  $\sum_{i=1}^{n_Q} \alpha_i = 1,$



$$\theta_{u,v} = \begin{cases} \frac{1}{\lambda_1} \left(1 - \frac{\lambda_2}{\|\xi_{u,v}\|}\right) \xi_{u,v} & \text{if } \|\xi_{u,v}\| > \lambda_2, \\ 0 & \text{if } \|\xi_{u,v}\| \leq \lambda_2. \end{cases}$$

$$\xi_{u,v} = g_{u,v} - H_{u,v} \alpha,$$

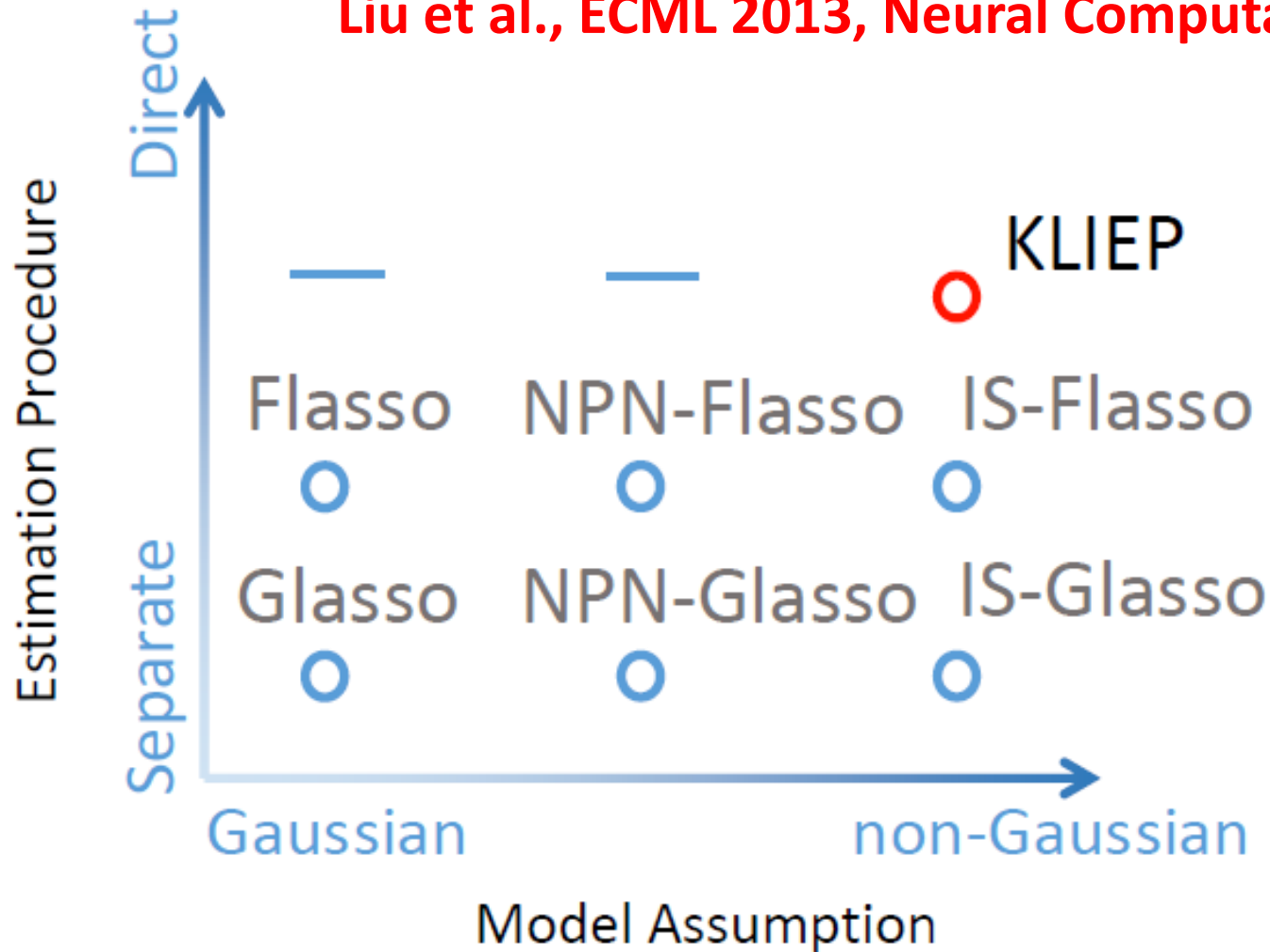
$$H_{u,v} = [\mathbf{f}(x_1^{(u)Q}, x_1^{(v)Q}), \dots, \mathbf{f}(x_{n_Q}^{(u)Q}, x_{n_Q}^{(v)Q})],$$

$$g_{u,v} = \frac{1}{n_P} \sum_{i=1}^{n_P} \mathbf{f}(x_i^{(u)P}, x_i^{(v)P}).$$

We can recover primal solution! 48

# Summary: Comparison between Methods

Liu et al., ECML 2013, Neural Computation 2014



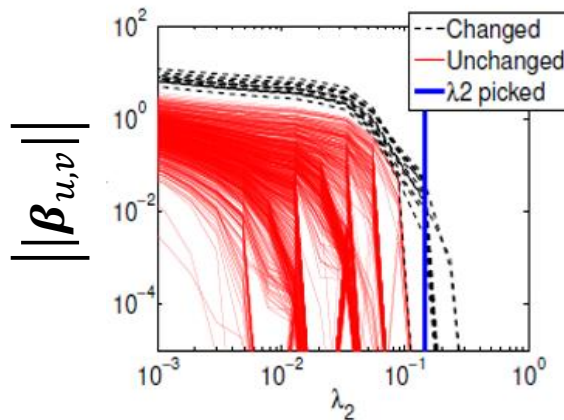
- The Illustration of methods on **model assumption** and **estimation procedure**



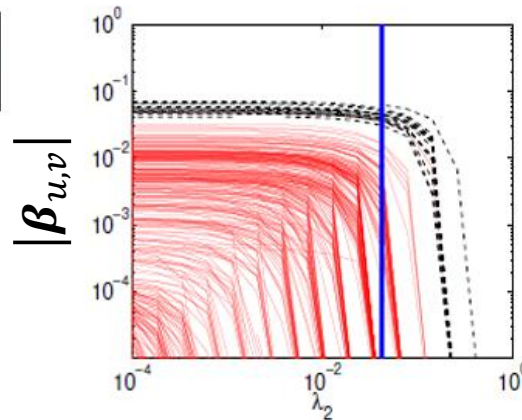
# Chapter 3, Structural Change Detection

1. Background
2. Existing Methods
3. Proposed Methods
4. Experiments

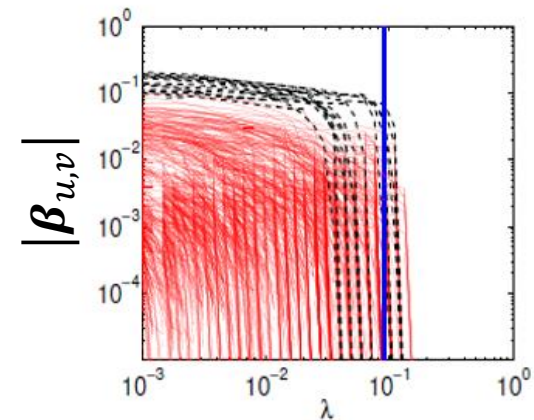
# Gaussian Distribution ( $n = 100, d = 40$ )



(a) KLIEP,  $n = 100$



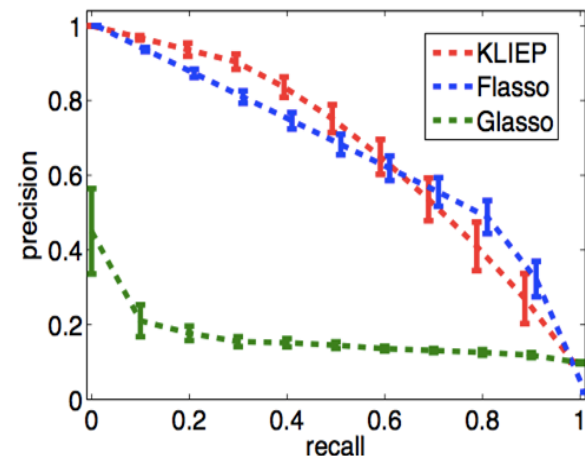
(b) Flasso,  $n = 100$



(c) Glasso,  $n = 100$

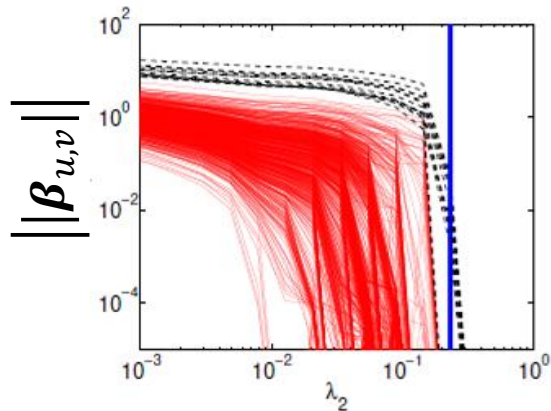
**Regularization path**

- Start from 40 dimensional GMN with random correlations.
- Randomly drops 15 edges.
- Precision and Recall curves are averaged over 20 runs.

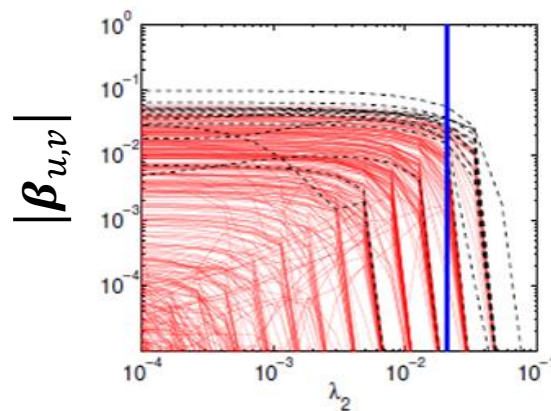


**P-R curve**

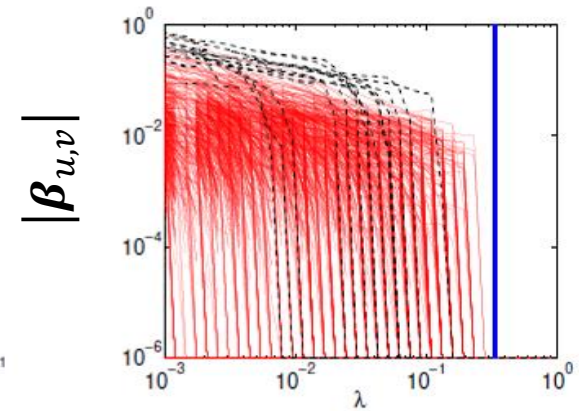
# Gaussian Distribution ( $n = 50, d = 40$ )



(d) KLIEP,  $n = 50$



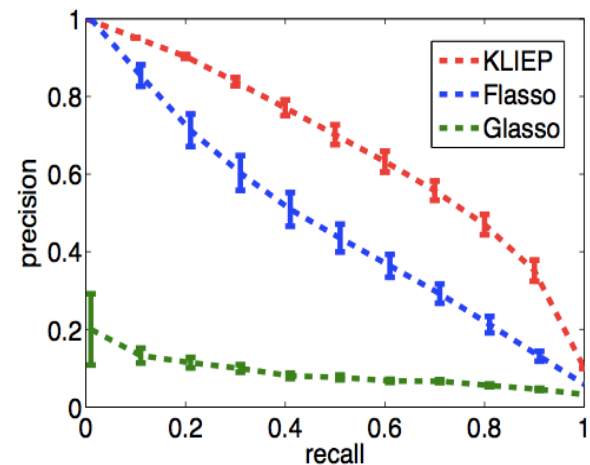
(e) Flasso,  $n = 50$



(f) Glasso,  $n = 50$

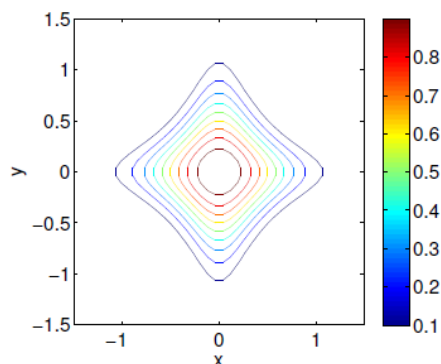
Regularization path

- Start from 40 dimensional GMN with random correlations.
- Randomly drops 15 edges.
- Precision and Recall curves are averaged over 20 runs.

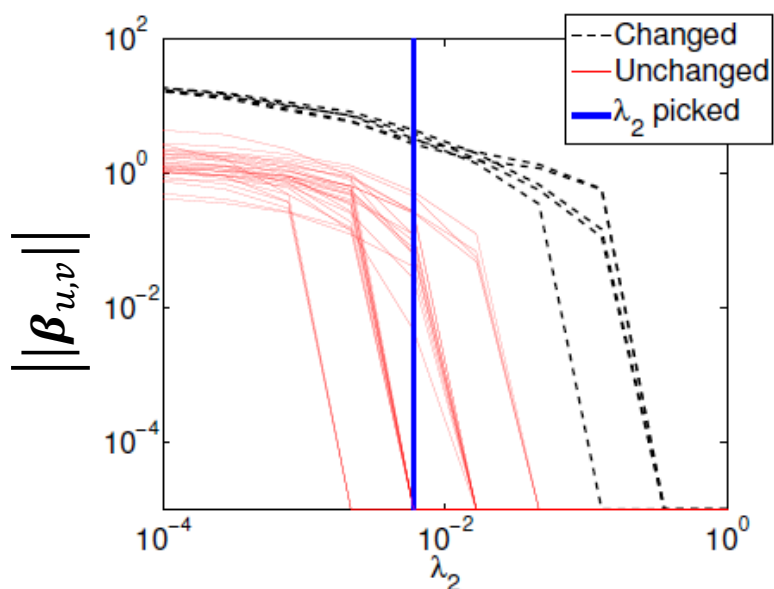


P-R curve

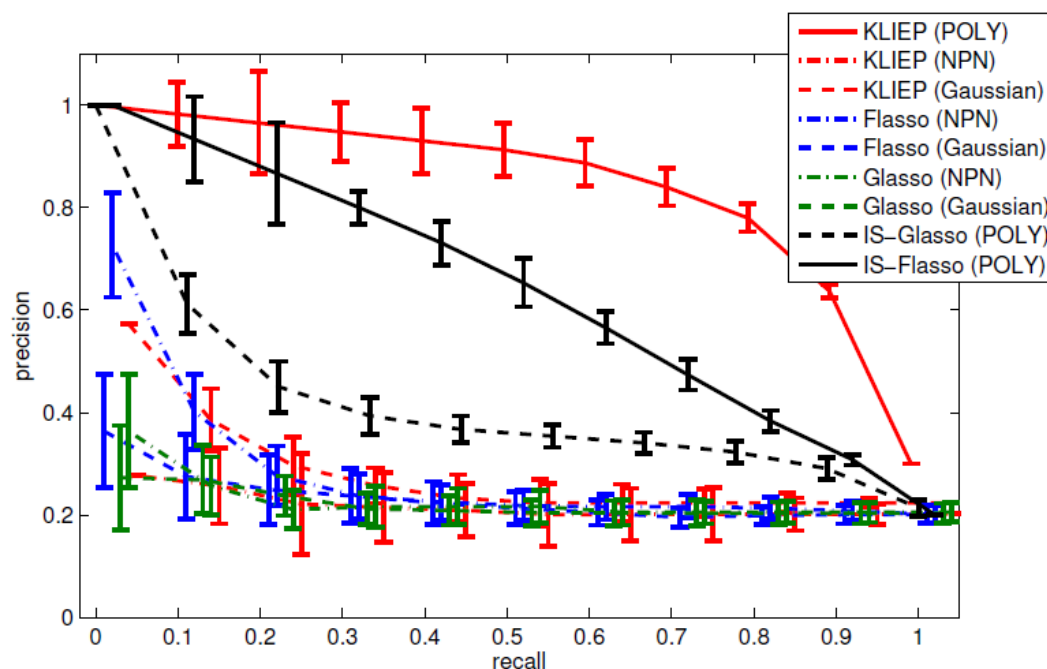
# Diamond Distribution ( $n = 5000, d = 9$ )



$$p(\mathbf{x}) \propto \exp \left( - \sum_i 2x_i^2 - \sum_{i,j \in E} 20x_i^2 x_j^2 \right)$$



Regularization path



P-R curve

- The proposed method has the leading performance.

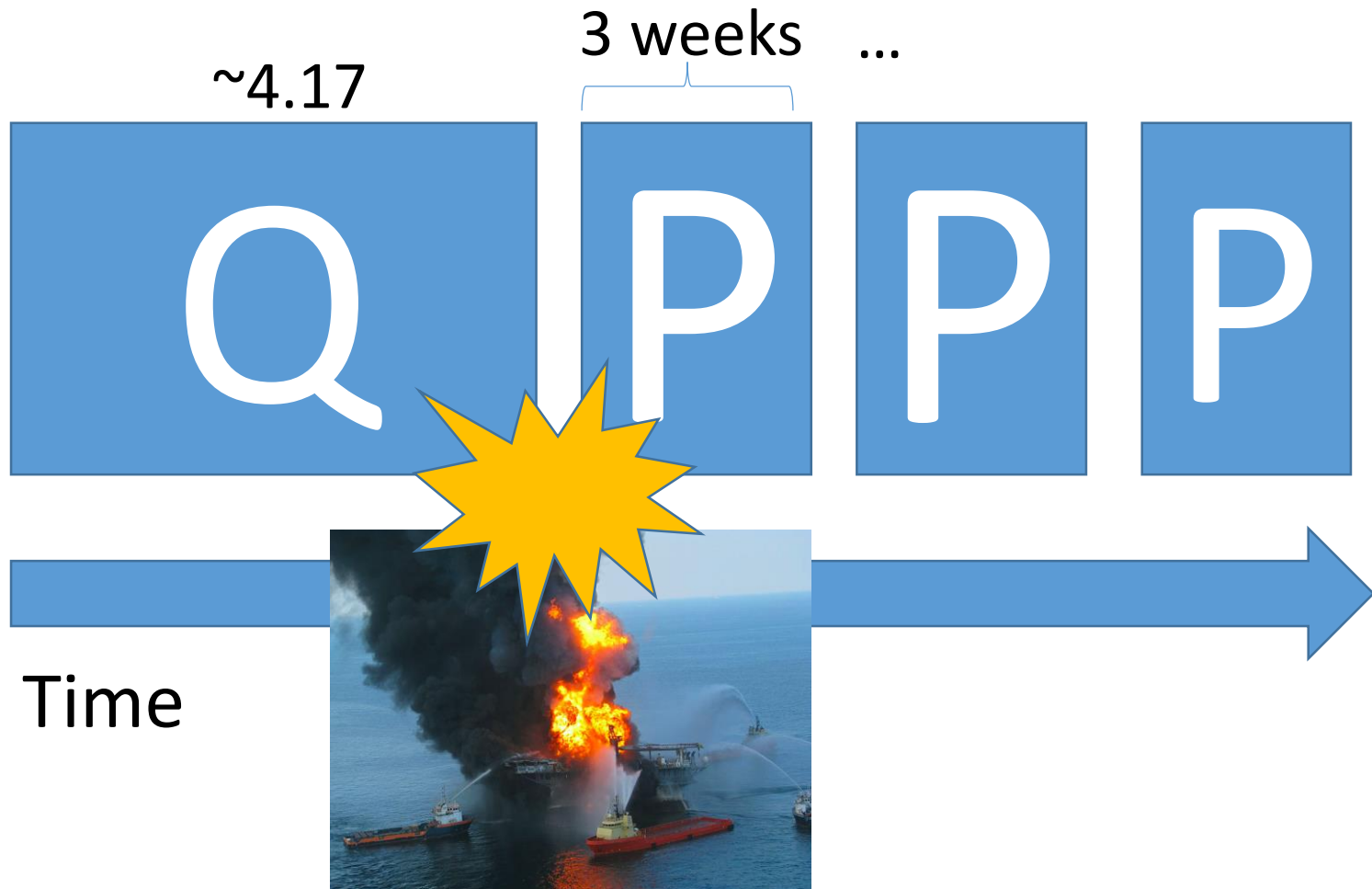
# Twitter Dataset



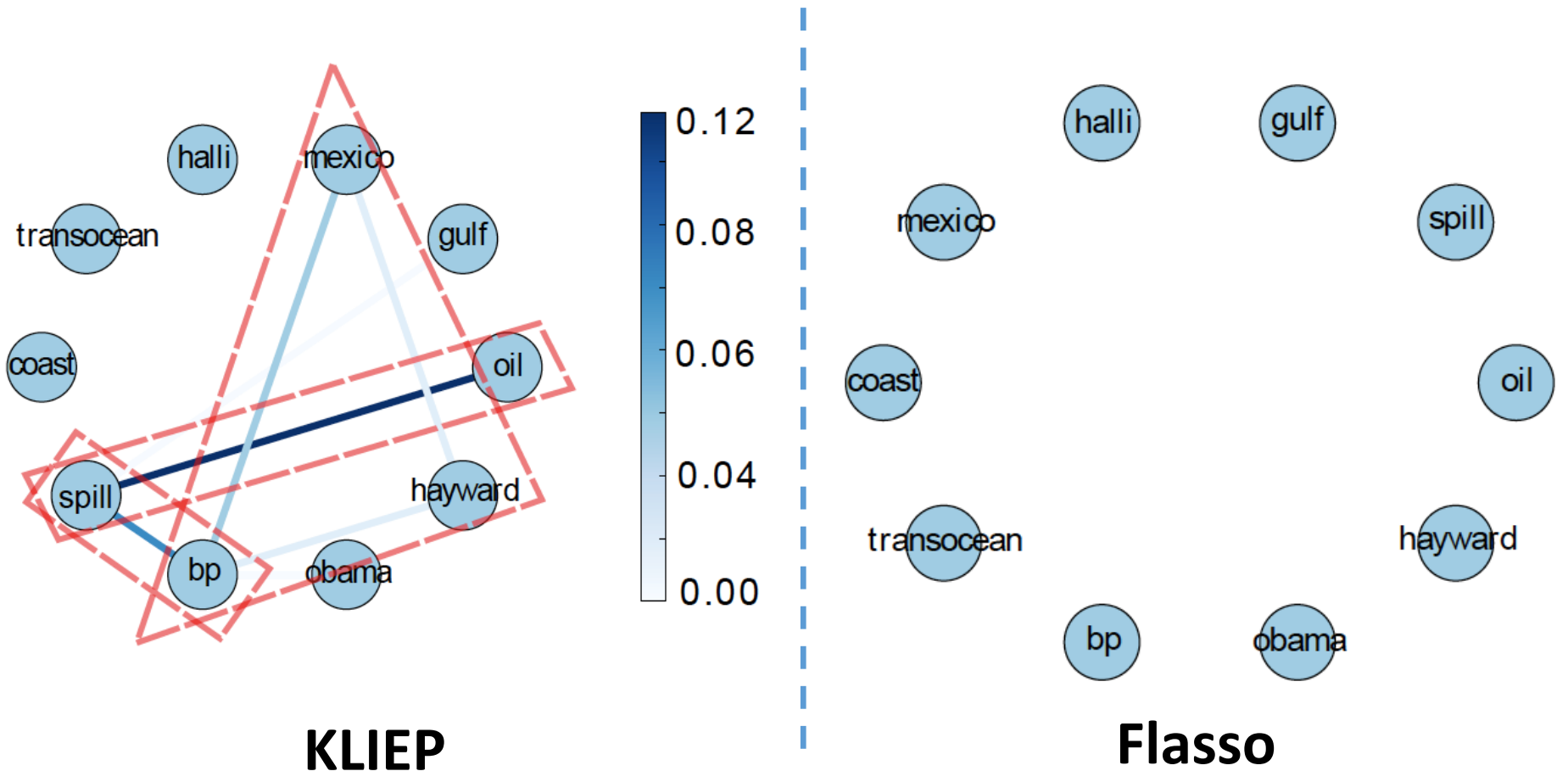
source: Wikipedia

- We choose the Deepwater Horizon oil spill as the target event.
  - Samples are the frequencies of 10 related keywords over time.
  - Detecting the change of **co-occurrences on keywords** before and after a certain event.

# Twitter Dataset (pp. 68-72)



# The Change of Correlation From 7.26-9.14



Small # of samples, severe over-fitting of Flasso!

# Conclusion

- We proposed a **direct** structural change detection method for Markov Networks.
- The proposed method can handle **non-Gaussian data**, and can be efficiently solved using **primal** or **dual** objective.
- The proposed method showed the **leading** performance on both artificial and real-world datasets