

MLSS, Beijing

Jun. 16, 2014

Unsupervised Change Detection

Masashi Sugiyama

Tokyo Institute of Technology, Japan

sugi@cs.titech.ac.jp

<http://sugiyama-www.cs.titech.ac.jp/~sugi/>



Abstract

2

We consider the problem of detecting change in two sets of samples, and explore two approaches: distributional and structural change detection.

Distributional change detection is aimed at estimating a divergence between the probability densities behind the two sets of samples. We first explain that the two-step approach of first estimating the probability densities and then computing the divergence from the estimated densities results in systematic under-estimation of the divergence. Then we introduce methods to directly estimate the ratio of densities and the difference of densities, which are shown to be more reliable than the density estimation approach.

Abstract (cont.)

3

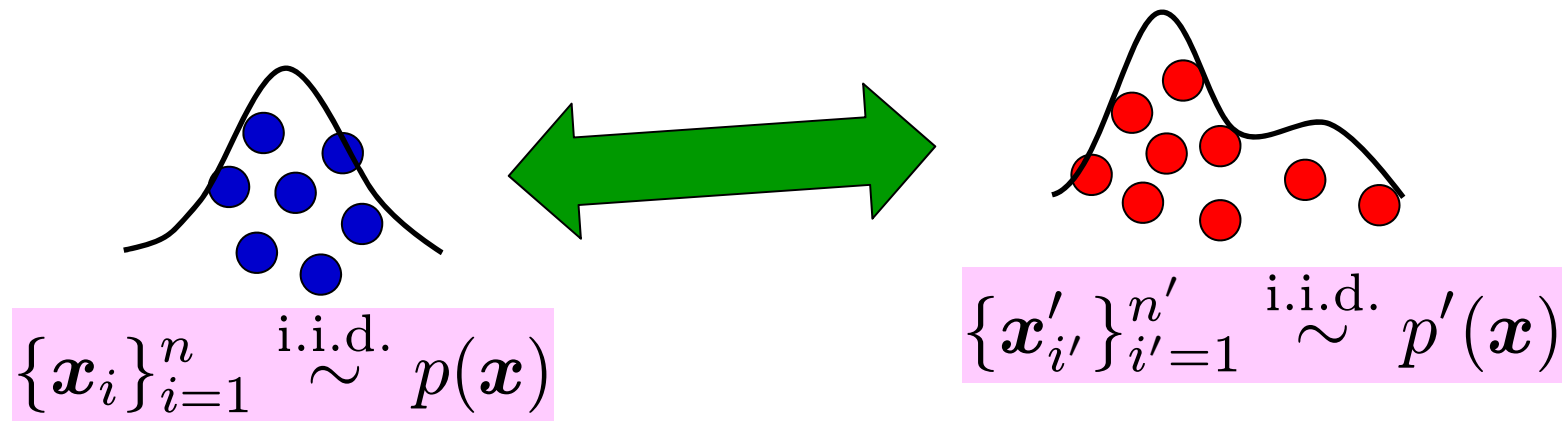
Structural change detection tries to identify change in element-wise dependency structure in multi-dimensional samples. We first consider the Gaussian sparse covariance selection setup and introduce approaches based on LASSO and fused-LASSO. Then we extend our discussion to non-Gaussian Markov networks, which generally suffer computational intractability of the normalization term, and introduce the importance sampling technique and the score matching method. Finally, we cover a method to directly compare two Markov networks for change detection.

No solid background on change detection is necessary, but basic knowledge of elementary statistics, linear algebra, and optimization is assumed.

Change Detection

4

- **Goal:** Given two sets of samples, we want to compare the probability distributions behind



- Two approaches:

- **Distributional change detection:** Flexible and robust
- **Structural change detection:** Interpretable



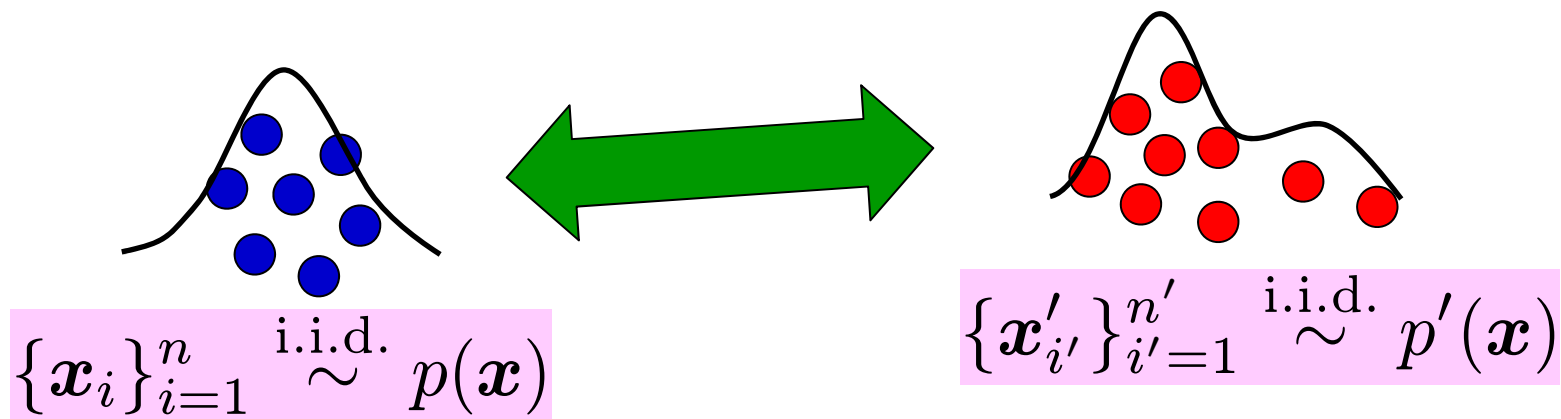
Contents

5

1. Distributional change detection
 - A) Problem setup and motivating examples
 - B) Distances
 - C) Distance approximation
2. Structural change detection

Distributional Change Detection 6

- **Goal:** Detect change in probability distributions behind two sets of samples **through divergence**



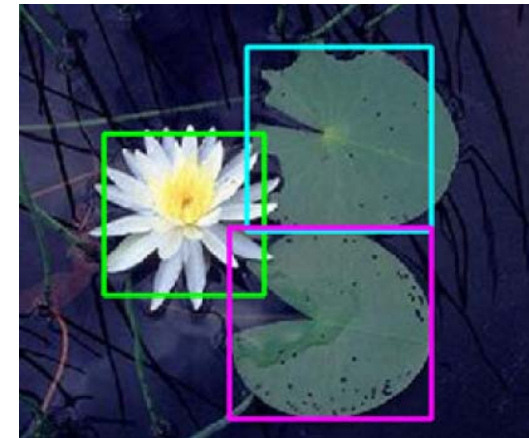
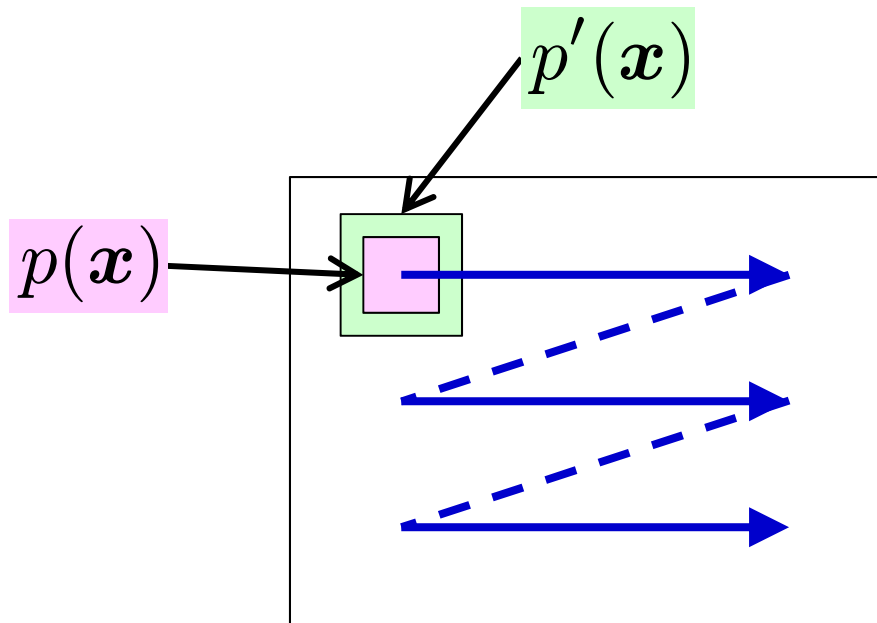
$$\text{Divergence}(p, p') < \varepsilon ?$$

Motivating Example 1

7

■ Region-of-interest detection in images:

- $p(\mathbf{x})$ and $p'(\mathbf{x})$ are significantly different when a visually salient object is included inside.

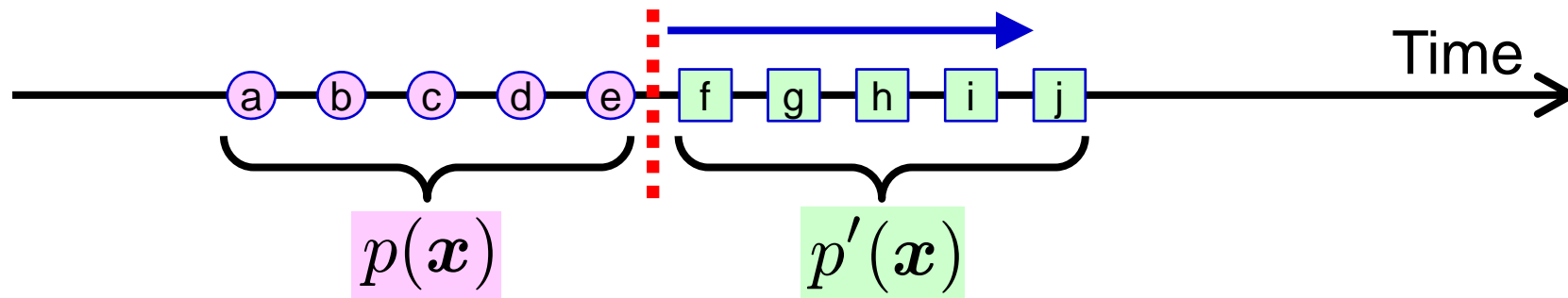


Motivating Example 2

8

■ Event detection in movies:

- $p(x)$ and $p'(x)$ are significantly different when an irregular event occurs.



52



57



62



67



72



77



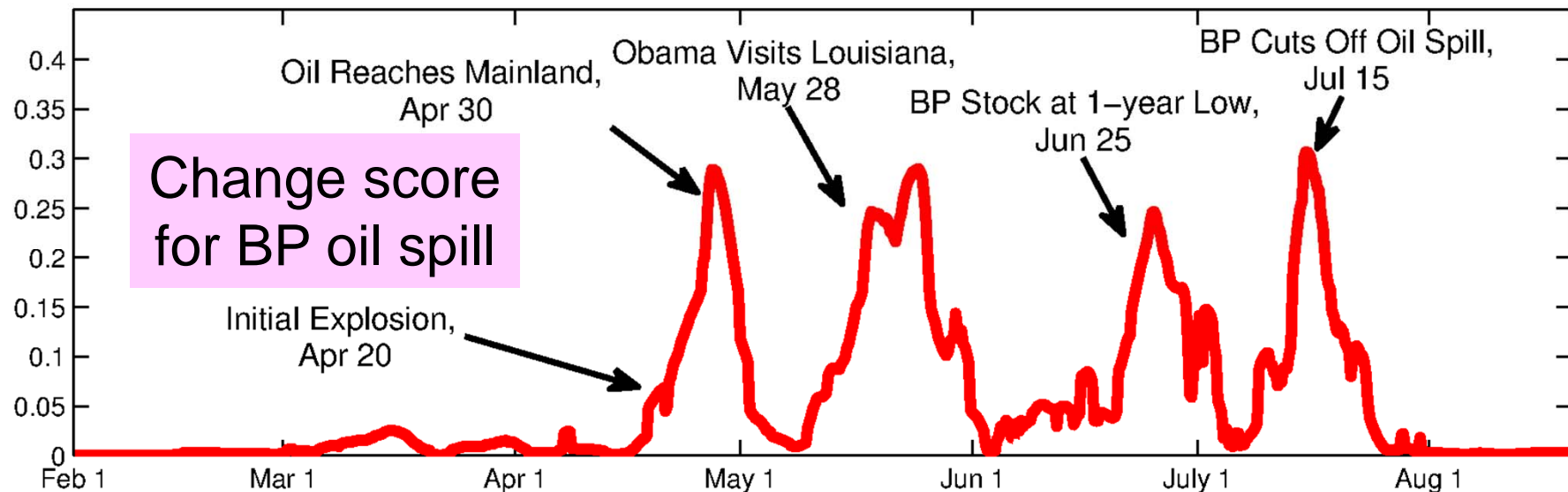
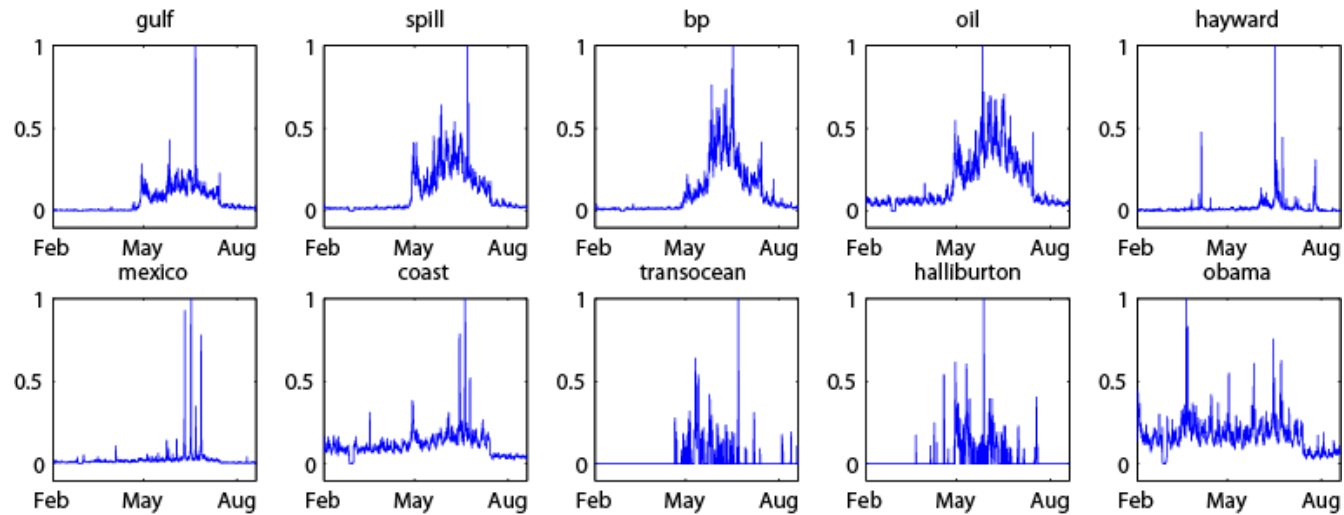
82



87

Motivating Example 3

Event detection from Twitter:





Contents

10

1. **Distributional change detection**
 - A) Problem setup and motivating examples
 - B) **Distances**
 - I. Density-ratio divergences
 - II. Density-difference distances
 - C) Distance approximation
2. Structural change detection

Distances and Divergences

11

Distance:

- Non-negativity:

$$\forall x, y, \quad d(x, y) \geq 0$$

- Non-degeneracy:

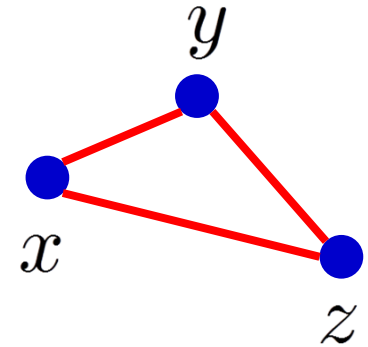
$$d(x, y) = 0 \iff x = y$$

- Symmetry:

$$\forall x, y, \quad d(x, y) = d(y, x)$$

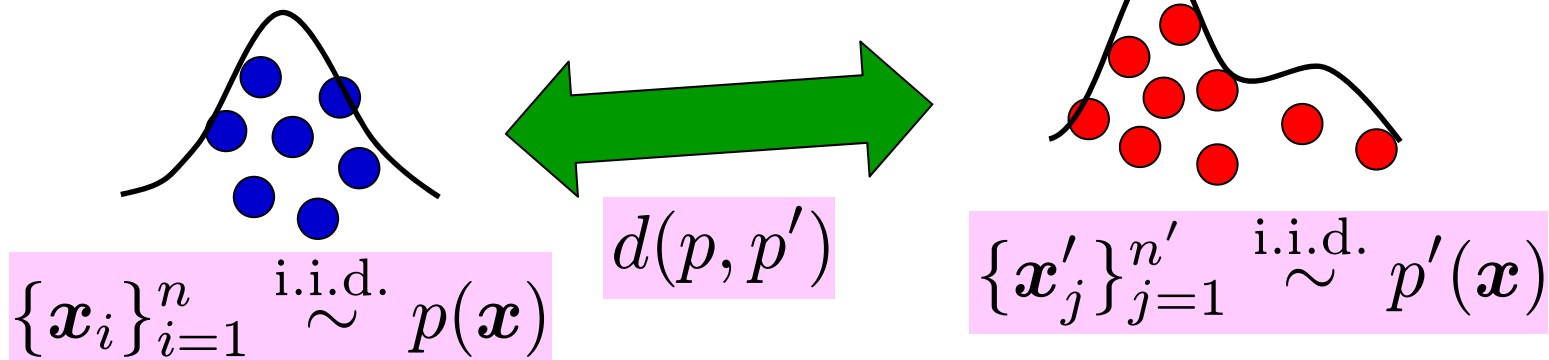
- Triangularity:

$$\forall x, y, z \quad d(x, z) \leq d(x, y) + d(y, z)$$



■ A **divergence** is a pseudo-distance.

■ We consider distances/divergences between **probability densities**.





Contents

12

1. **Distributional change detection**
 - A) Problem setup and motivating examples
 - B) **Distances**
 - I. **Density-ratio divergences**
 - II. Density-difference distances
 - C) Distance approximation
2. Structural change detection

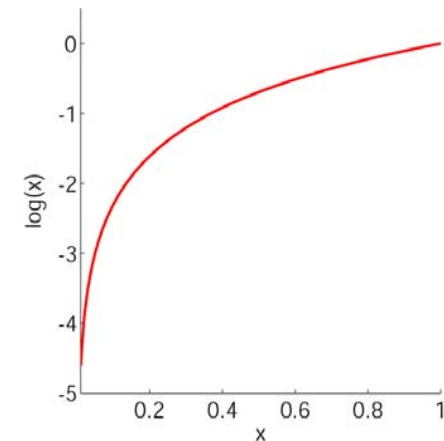
Kullback-Leibler Divergence

13

Kullback & Leibler (1951)

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

- ☺ Compatible with **maximum likelihood**.
- ☺ **Invariant** under input transformation. $\frac{p(\mathbf{x})}{p'(\mathbf{x})}$
(Jacobians cancel in the density ratio)
- ☹ Doesn't satisfy symmetry and triangularity.
- ☹ **Sensitive to outliers**
(due to log and ratio).



f-Divergences

14

Ali & Silvey (1966), Csiszár (1967)

$$F(p||p') = \int p'(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{p'(\mathbf{x})}\right) d\mathbf{x}$$

f : Convex function such that $f(1) = 0$

- $f(t) = t \log t$ yields the KL-divergence:

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

- To avoid the log function, let us use

$$f(t) = (t - 1)^2$$

Pearson (PE) Divergence

15

Pearson (1900)

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

- 😊 Compatible with **least-squares**.
- 😊 **Invariant** under input transformation.
- 😞 Doesn't satisfy symmetry and triangularity.
- 😞 **Sensitive to outliers** (no log, but still ratio).

$$\frac{p(\mathbf{x})}{p'(\mathbf{x})}$$

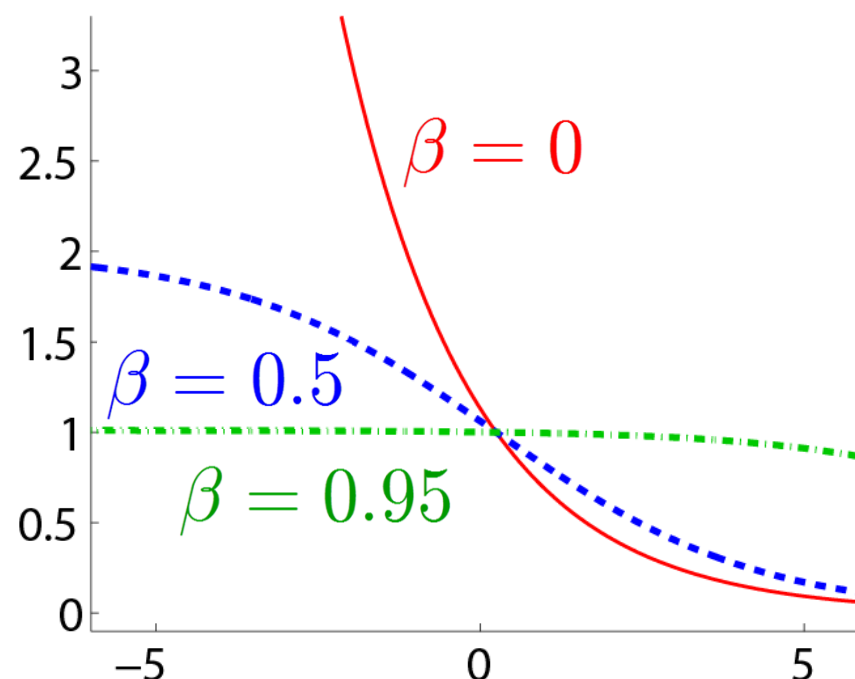
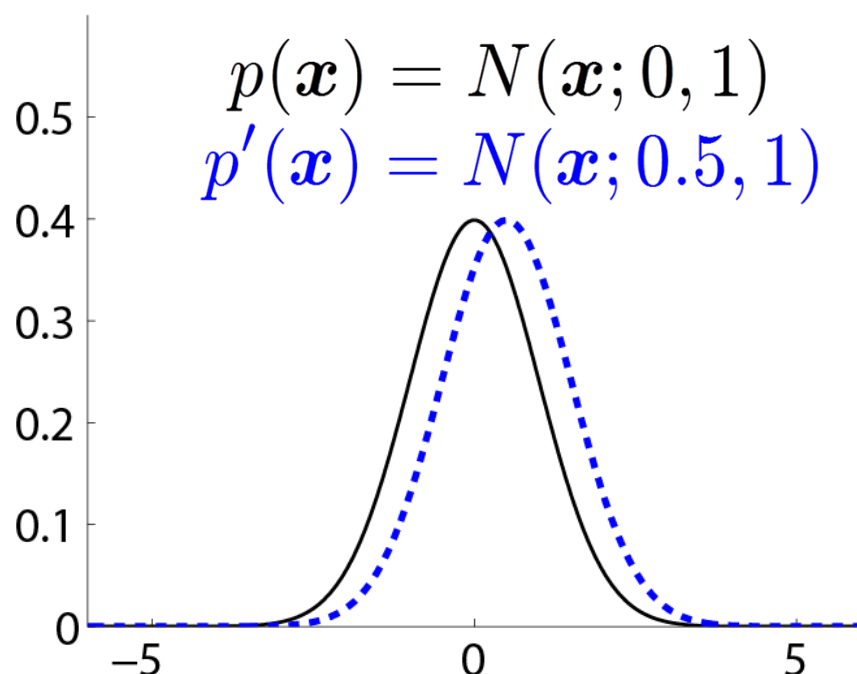
Relative Density Ratio

16

Yamada *et al.* (NIPS2011, NeCo2013)

- Density ratio $\frac{p(\mathbf{x})}{p'(\mathbf{x})}$ can diverge to **infinity**.
- **Relative density ratio** is always bounded:

$$\frac{p(\mathbf{x})}{\beta p(\mathbf{x}) + (1 - \beta)p'(\mathbf{x})} < \frac{1}{\beta} \quad 0 \leq \beta < 1$$



Relative Pearson (rPE) Divergence¹⁷

$$\text{rPE}(p||p') = \text{PE}(p||p_\beta) = \int p_\beta(\mathbf{x}) \left(\frac{p(\mathbf{x})}{p_\beta(\mathbf{x})} - 1 \right)^2 d\mathbf{x}$$

$$0 \leq \beta < 1$$

$$p_\beta(\mathbf{x}) = \beta p(\mathbf{x}) + (1 - \beta)p'(\mathbf{x})$$

- 😊 Compatible with **least-squares**.
- 😊 **Invariant** under input transformation.
- 😊 **Robust against outliers**.
- 😞 Doesn't satisfy symmetry and triangularity.
- 😞 **Not clear how to choose β** .



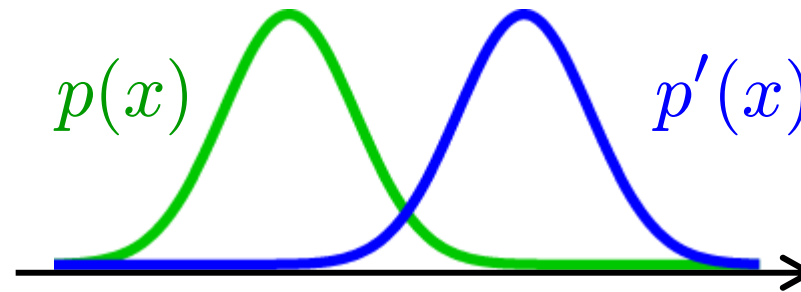
Contents

18

1. **Distributional change detection**
 - A) Problem setup and motivating examples
 - B) **Distances**
 - I. Density-ratio divergences
 - II. **Density-difference distances**
 - C) Distance approximation
2. Structural change detection

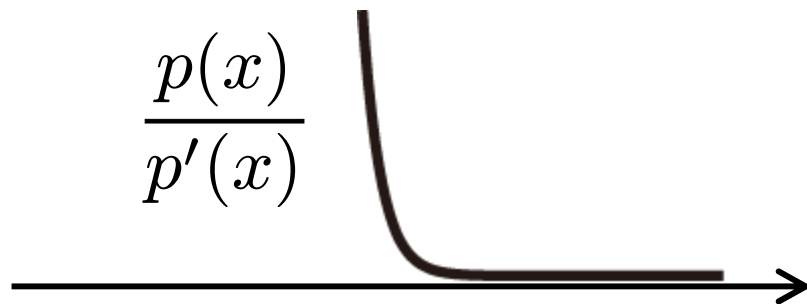
Density Difference

19



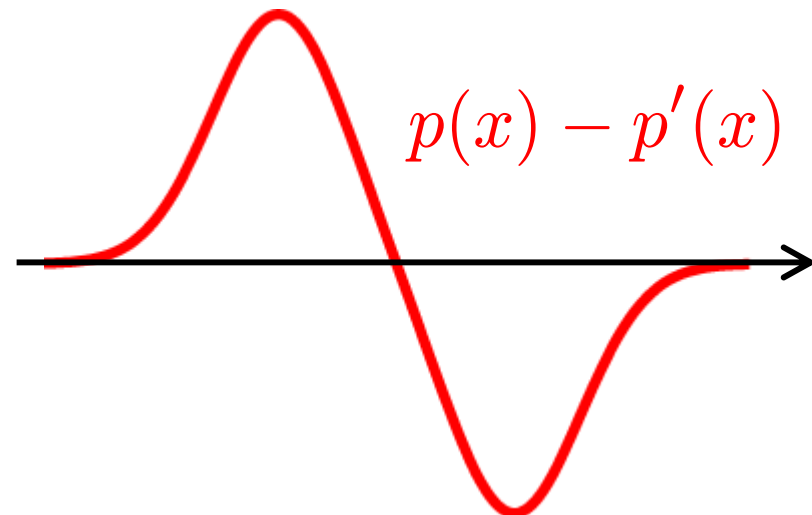
■ Density ratio based distance:

- Is the ratio 1?



■ Density difference based distance:

- Is the difference 0?



L^t -Distance

$$L^t(p, p') = \int |p(\mathbf{x}) - p'(\mathbf{x})|^t d\mathbf{x} \quad t \geq 0$$

☺ Proper distance.

☺ **Robust against outliers** (no ratio).

■ When $t = 2$:

☺ Compatible with **least-squares**.

☹ **Not invariant** under input transformation.

■ When $t = 1$:

☺ **Invariant** under input transformation

(because f-div).

$$f(t) = |t - 1|$$

$$L^1(p, p') = \int p'(\mathbf{x}) \left| \frac{p(\mathbf{x})}{p'(\mathbf{x})} - 1 \right| d\mathbf{x}$$

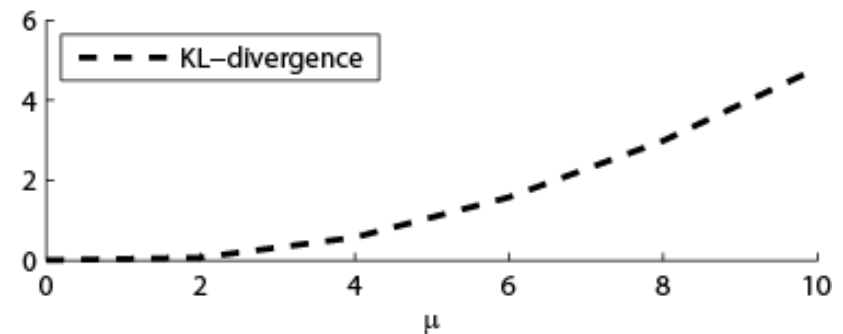
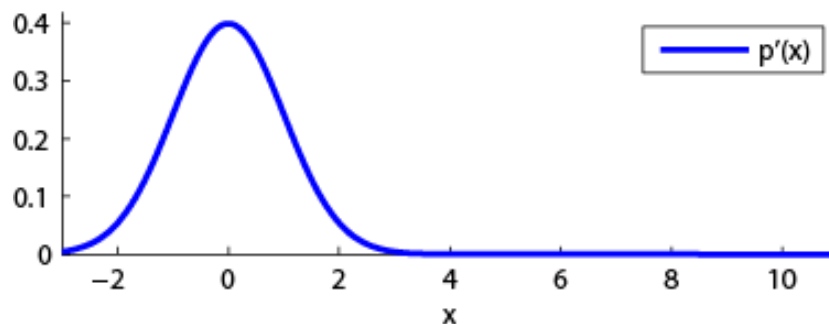
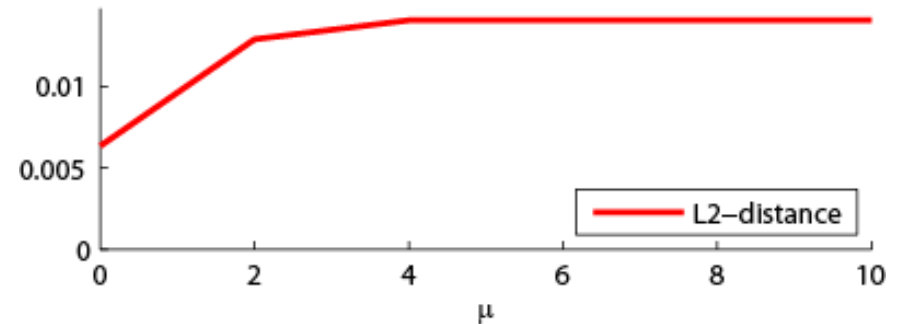
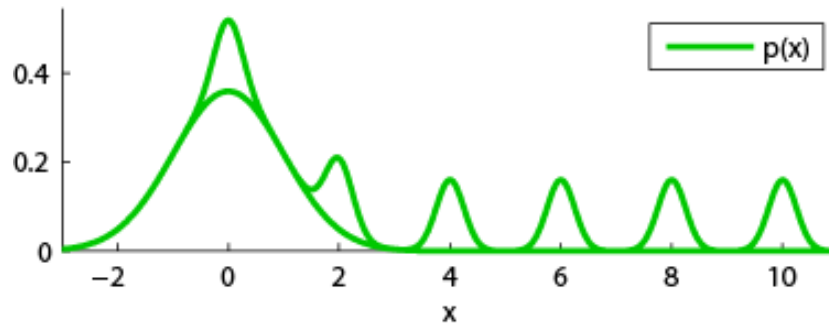
KL vs. L^2 for Outliers

21

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

$$L^2(p, p') = \int (p(\mathbf{x}) - p'(\mathbf{x}))^2 d\mathbf{x}$$

$$p(x) = 0.9p'(x) + 0.1q(x - \mu)$$



- L^2 -distance is bounded.
- KL-divergence is unbounded.



Contents

22

1. Distributional change detection
 - A) Problem setup and motivating examples
 - B) Distances
 - C) **Distance approximation**
 - I. Density-ratio divergences
 - II. Density-difference distances
2. Structural change detection

Distance Approximation via Density Estimation

1. **Estimate densities** $p(\mathbf{x}), p'(\mathbf{x})$ from samples:

$$\{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}) \quad \{\mathbf{x}'_{i'}\}_{i'=1}^{n'} \stackrel{\text{i.i.d.}}{\sim} p'(\mathbf{x})$$

- Maximum likelihood estimation
- Bayes estimation
- Kernel density estimation
- Nearest-neighbor density estimation.

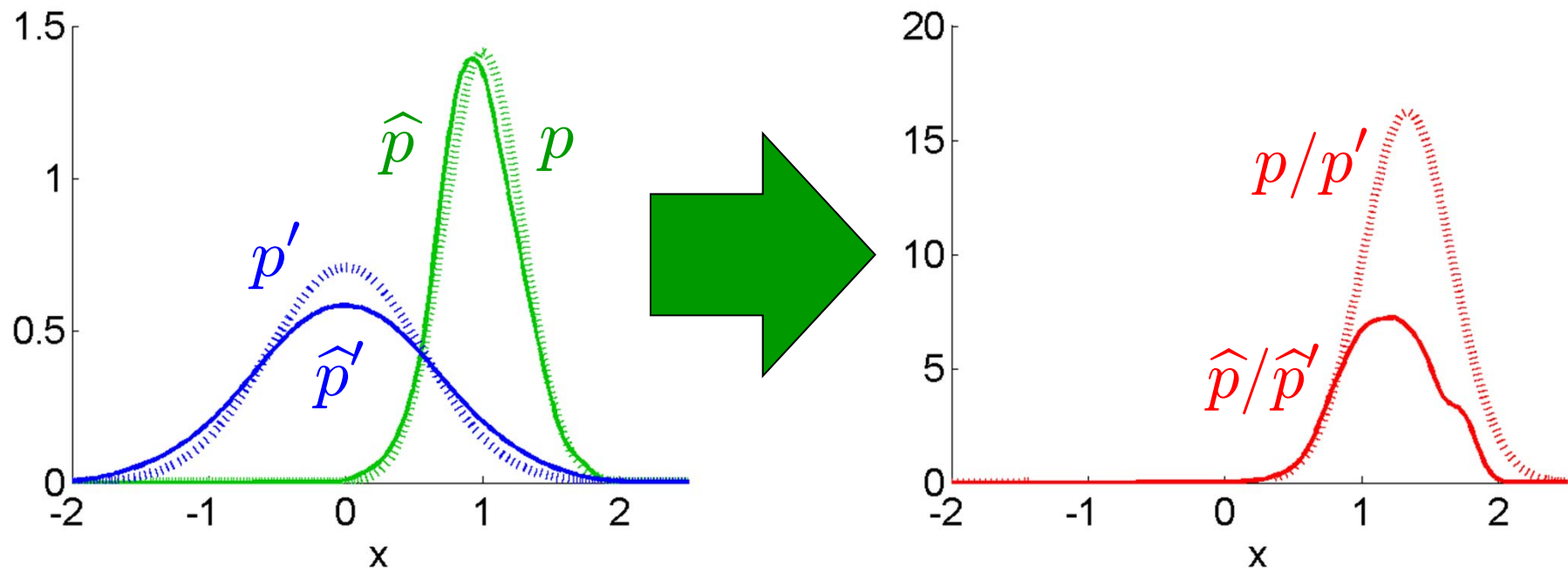
2. **Plug-in** the estimated densities $\hat{p}(\mathbf{x}), \hat{p}'(\mathbf{x})$:

$$\widehat{\text{KL}}(p||p') = \int \hat{p}(\mathbf{x}) \log \frac{\hat{p}(\mathbf{x})}{\hat{p}'(\mathbf{x})} d\mathbf{x}$$

$$\widehat{L}^2(p, p') = \int \left(\hat{p}(\mathbf{x}) - \hat{p}'(\mathbf{x}) \right)^2 d\mathbf{x}$$

Drawback of Plug-In Density Estimation Approach

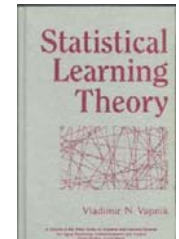
- Densities are estimated **without regard to taking their ratio later**.
- Division by \hat{p}' **magnifies** estimation error in \hat{p} .



Guiding Principle

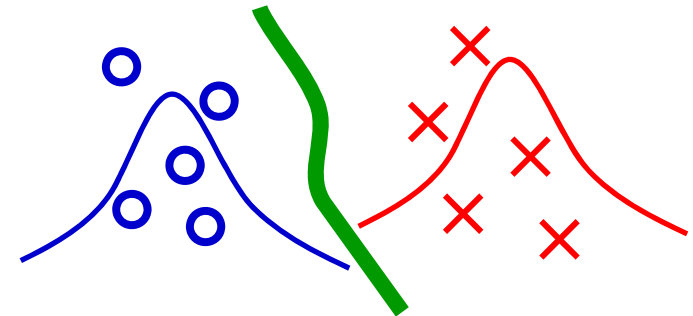
- **Vapnik's principle:** Vapnik (Wiley 1998)

When solving a problem of interest, one should not solve a more general problem as an intermediate step



- **Support vector machine** avoids general density estimation and directly learns the boundary.

Cortes & Vapnik (MLJ1995)



- Let's avoid separately estimating $p(x)$ and $p'(x)$, and **directly compare the densities!**

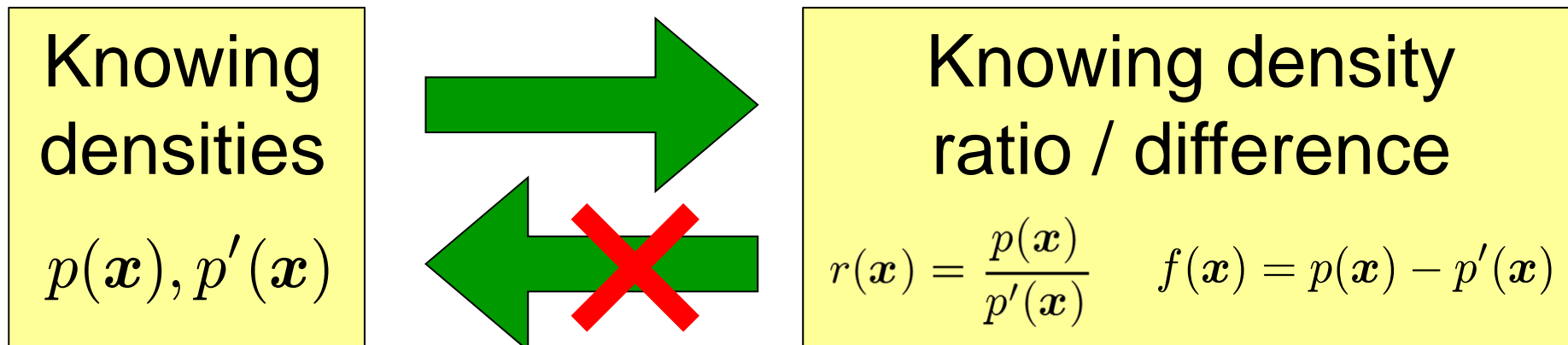
Vapnik's Principle in Distance Approximation

$$\text{KL}(p\|p') = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{p'(\mathbf{x})} d\mathbf{x} \quad L^2(p, p') = \int \left(p(\mathbf{x}) - p'(\mathbf{x}) \right)^2 d\mathbf{x}$$

- Directly estimate the **density ratio / difference**:

$$r(\mathbf{x}) = \frac{p(\mathbf{x})}{p'(\mathbf{x})} \quad f(\mathbf{x}) = p(\mathbf{x}) - p'(\mathbf{x})$$

without estimating each density $p(\mathbf{x}), p'(\mathbf{x})$.





Contents

27

1. Distributional change detection
 - A) Problem setup and motivating examples
 - B) Distances
 - C) Distance approximation
 - I. Density-ratio divergences
 - II. Density-difference distances
2. Structural change detection

KL-Divergence Approximation

28

Nguyen *et al.* (NIPS2007, IEEE-IT2010)

Sugiyama *et al.* (NIPS2007, AISM2008)

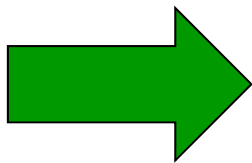
$$\text{KL}(p||p') = \int p(\mathbf{x}) \log r(\mathbf{x}) d\mathbf{x}$$

$$r(\mathbf{x}) = \frac{p(\mathbf{x})}{p'(\mathbf{x})}$$

- Directly approximate the **density ratio** with **log-loss**:

$$\hat{r} = \underset{\tilde{r}}{\text{argmin}} \text{KL}(p||\tilde{r} \cdot p')$$

$$\text{subject to } \int \tilde{r}(\mathbf{x}) p'(\mathbf{x}) d\mathbf{x} = 1 \quad \text{and} \quad \tilde{r} \geq 0$$



$$\text{KL}(p||p') \approx \int p(\mathbf{x}) \log \hat{r}(\mathbf{x}) d\mathbf{x}$$

- Expectation is approximated by empirical average.

Solution for Linear Model

29

Linear-in-parameter model:

$$r_{\alpha}(\mathbf{x}) = \sum_{j=1}^b \alpha_j \phi_j(\mathbf{x}) = \alpha^{\top} \phi(\mathbf{x})$$

$\alpha = (\alpha_1, \dots, \alpha_b)^{\top}$: Fixed basis functions
 $\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_b(\mathbf{x}))^{\top}$: Parameters

Empirical optimization problem:

$$\hat{\alpha} = \operatorname{argmax}_{\alpha} \frac{1}{n} \sum_{i=1}^n \log r_{\alpha}(\mathbf{x}_i)$$

$$\text{subject to } \frac{1}{n'} \sum_{i'=1}^{n'} r_{\alpha}(\mathbf{x}'_{i'}) = 1 \text{ and } \alpha \geq \mathbf{0}$$

- The solution tends to be **sparse** due to $\alpha \geq \mathbf{0}$.

Solution for Linear Model

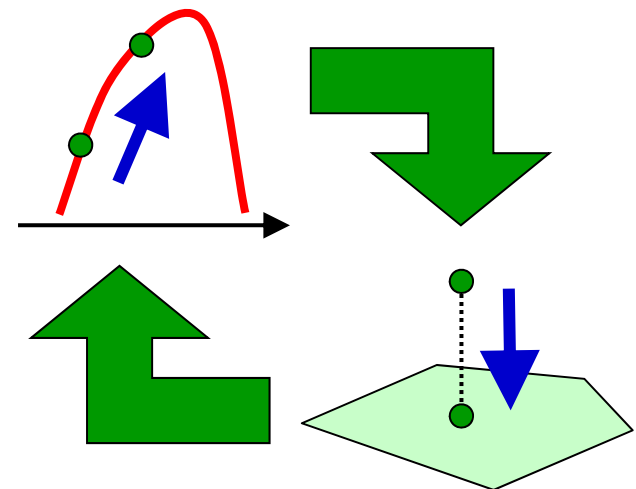
30

$$\hat{\alpha} = \operatorname{argmax}_{\alpha} \frac{1}{n} \sum_{i=1}^n \log \alpha^{\top} \phi(\mathbf{x}_i)$$

$$\text{subject to } \frac{1}{n'} \sum_{i'=1}^{n'} \alpha^{\top} \phi(\mathbf{x}'_{i'}) \text{ and } \alpha \geq \mathbf{0}$$

- Thanks to **convexity**, global optimal solution can be obtained by simple **gradient-projection**.
- Resulting KL-divergence approximator:

$$\text{KL}(p||p') \approx \frac{1}{n} \sum_{i=1}^n \log \hat{\alpha}^{\top} \phi(\mathbf{x}_i)$$



Other Models

Kernel model:

☺ Nonparametric

$$r_{\alpha}(\mathbf{x}) = \sum_{j=1}^n \alpha_j K(\mathbf{x}, \mathbf{x}_j)$$

Log-linear model:

☺ Always positive

$$r_{\alpha}(\mathbf{x}) = \exp \left(\sum_{j=1}^b \alpha_j \phi_j(\mathbf{x}) \right)$$

☺ Compatible with Markov networks

Gaussian mixture model:

☺ More flexible

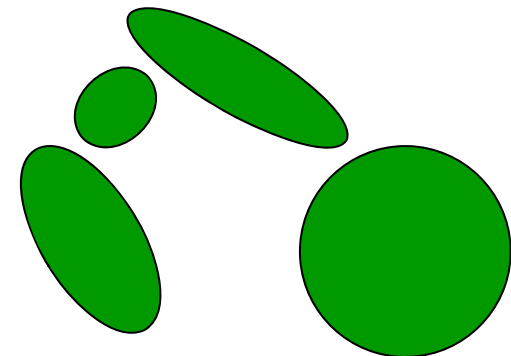
☹ Non-convex optimization

$$r_{\alpha}(\mathbf{x}) = \sum_{j=1}^b \alpha_j N(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

Probabilistic PCA mixture:

☺ Local dimension reduction

☹ Non-convex optimization

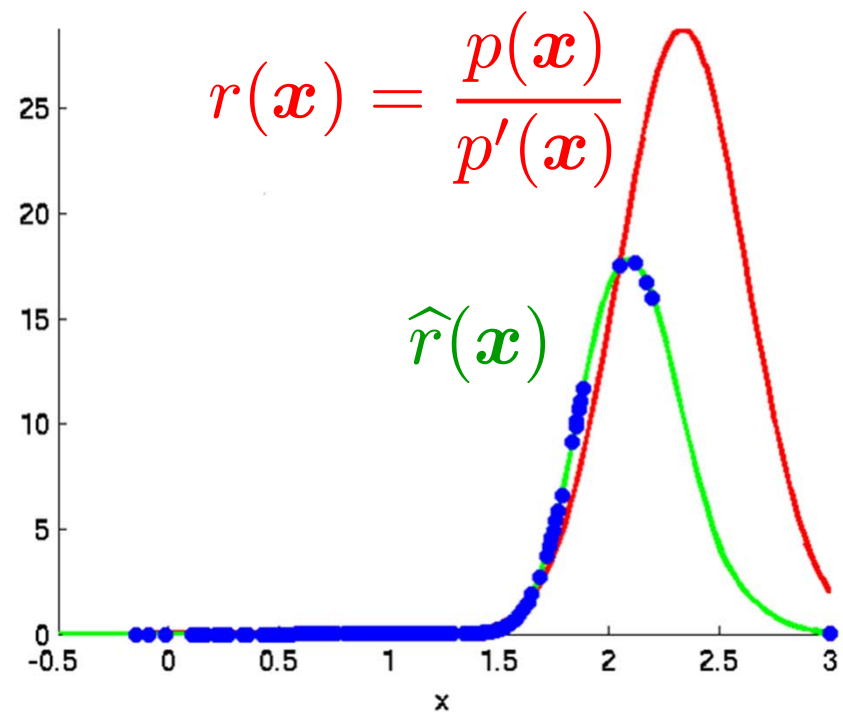
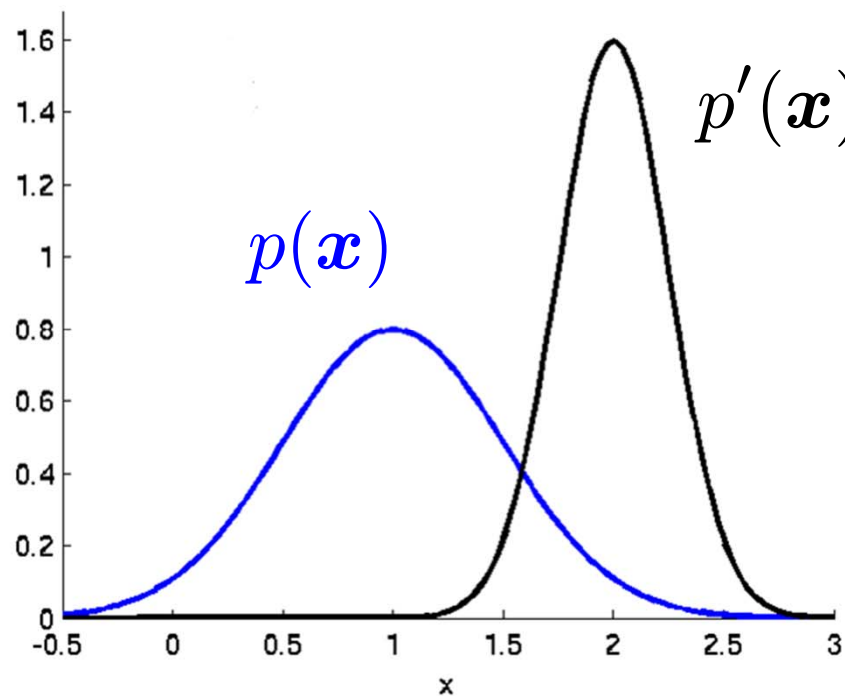


Numerical Example

32

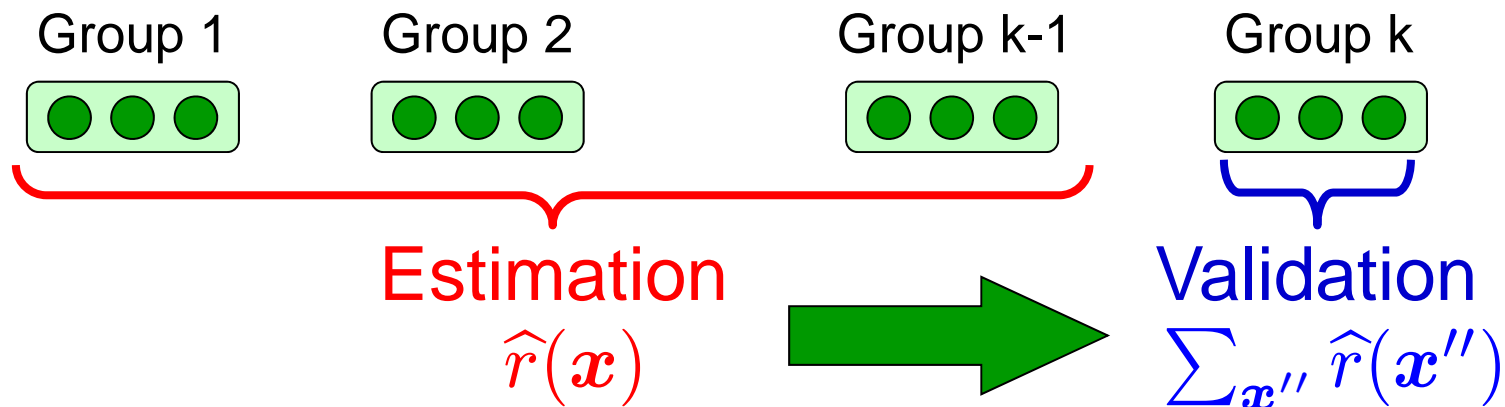
■ Gaussian kernel model:

$$r_{\alpha}(\mathbf{x}) = \sum_{j=1}^n \alpha_j \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$



Model Selection

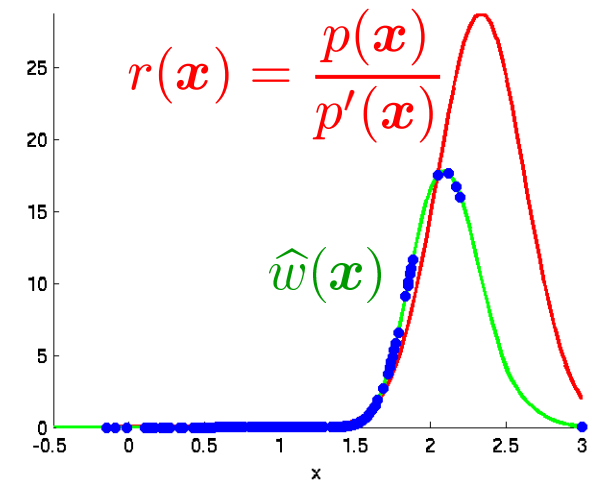
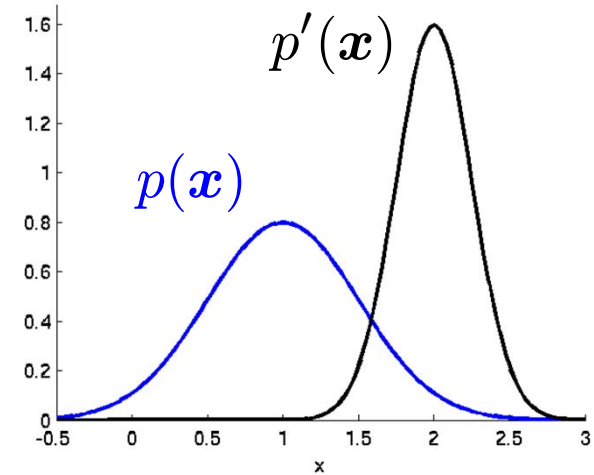
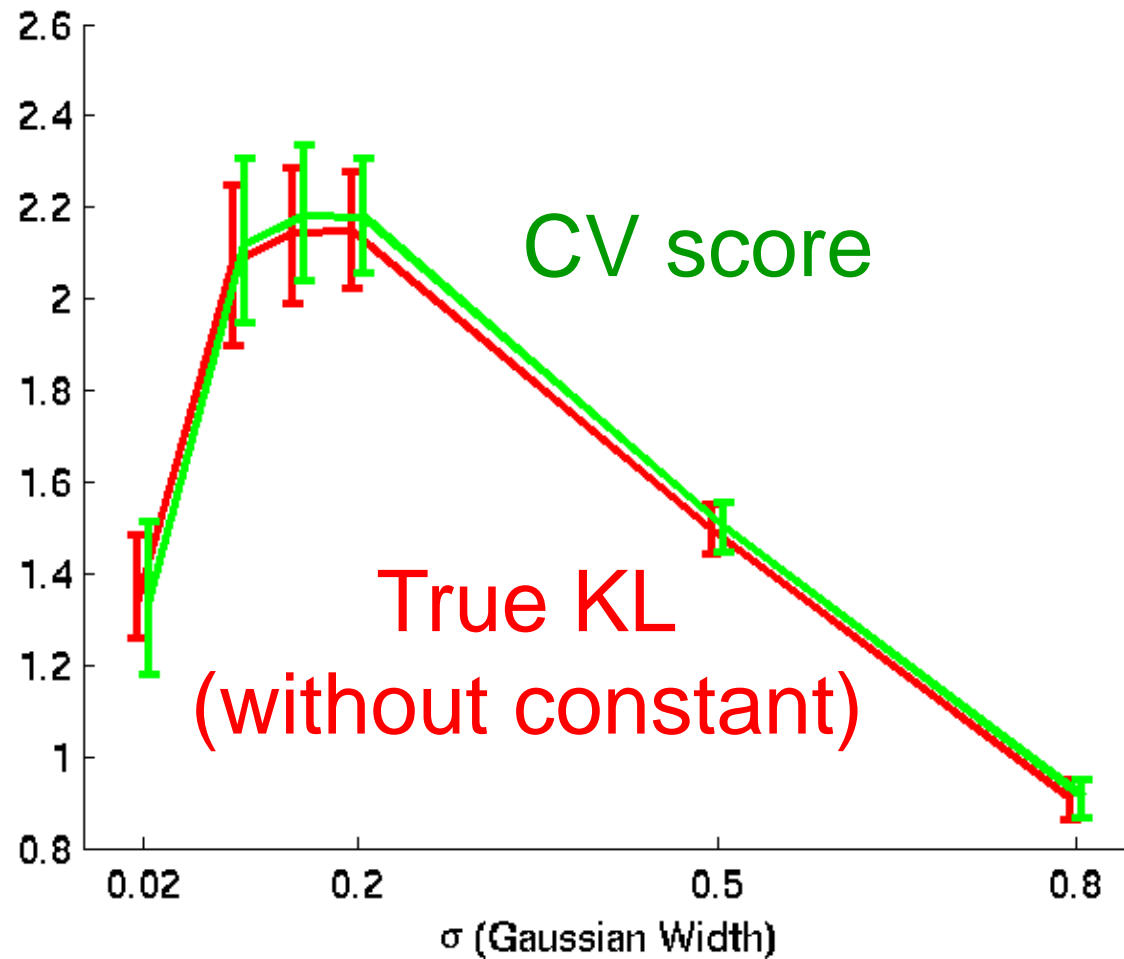
- Choice of the Gaussian bandwidth affects the performance.
- **Cross-validation (CV):**
 - Split $\{\mathbf{x}_i\}_{i=1}^n$ into estimation and validation subsets.



- Repeat this estimation-validation process for all combinations
- **CV gives an almost unbiased estimator of KL.**

Numerical Example

34



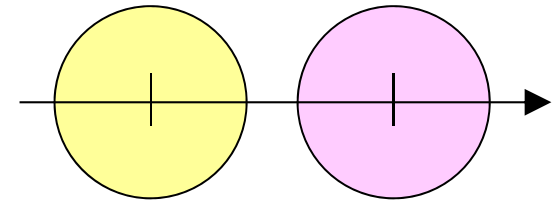
■ Model selection by CV works.

Comparison with KDE

35

■ d-dimensional Gaussians with covariance identity and

- **Denominator:** mean $(0,0,0,\dots,0)$
- **Numerator:** mean $(1,0,0,\dots,0)$



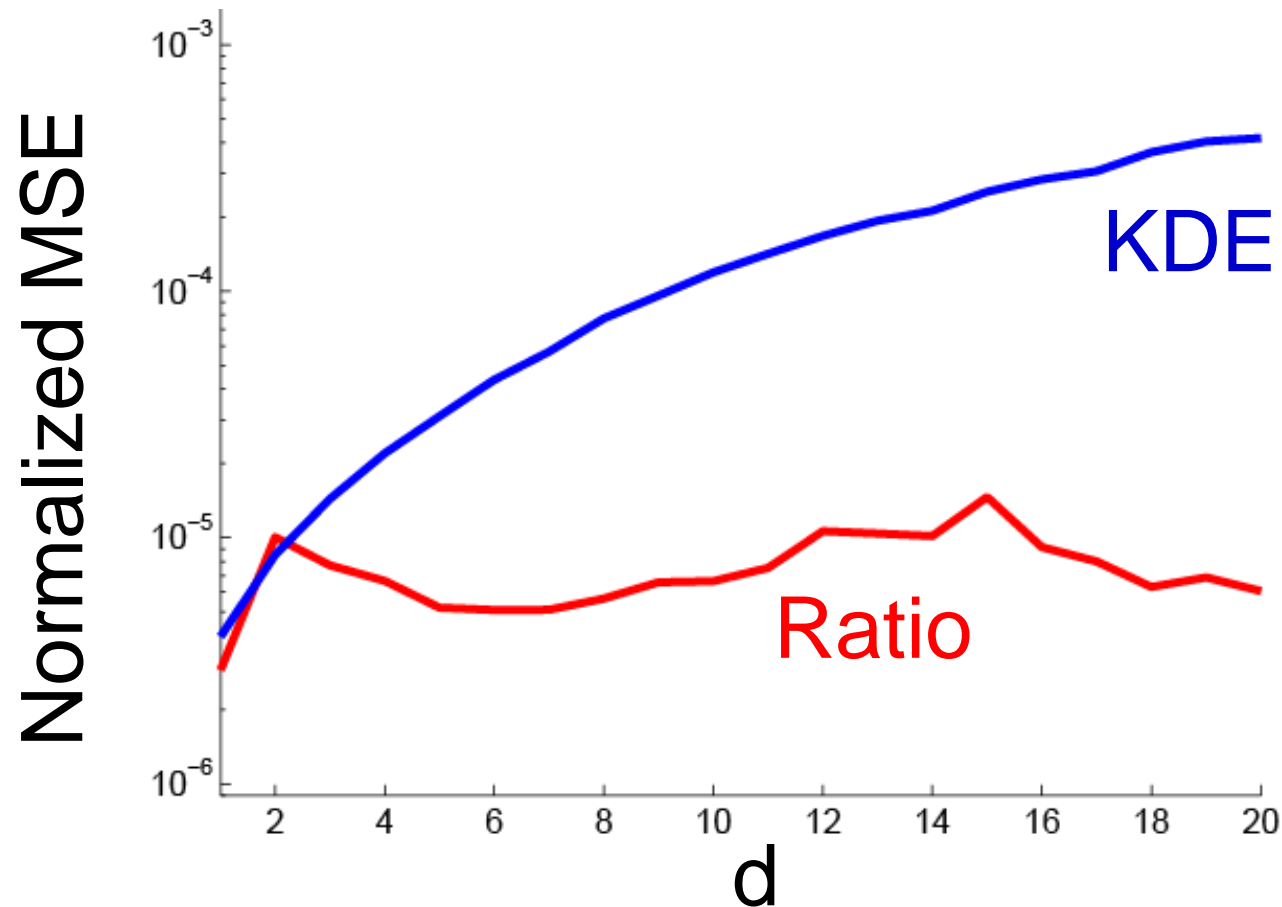
■ Kernel density estimation (KDE):

- Estimate two densities separately and take ratio.
- Gaussian widths are chosen by CV.

■ Ratio:

- Estimate the density ratio directly.
- Gaussian width is chosen by CV.

Accuracy as a Function of Input Dimensionality



- Density ratio approach works better.

PE-Divergence Approximation 37

Kanamori *et al.* (NIPS2008, JMLR2009)

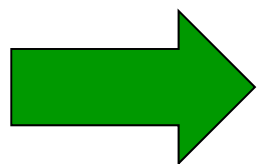
$$\text{PE}(p||p') = \int p'(\mathbf{x}) \left(r(\mathbf{x}) - 1 \right)^2 d\mathbf{x} = \int p(\mathbf{x}) r(\mathbf{x}) d\mathbf{x} - 1$$

- Directly approximate the density ratio by **least-squares**:

$$\hat{r} = \operatorname{argmin}_{\tilde{r}} \int p'(\mathbf{x}) \left(\tilde{r}(\mathbf{x}) - r(\mathbf{x}) \right)^2 d\mathbf{x}$$

$$r(\mathbf{x}) = \frac{p(\mathbf{x})}{p'(\mathbf{x})}$$

$$= \operatorname{argmin}_{\tilde{r}} \int p'(\mathbf{x}) \left(\tilde{r}(\mathbf{x}) \right)^2 d\mathbf{x} - 2 \int p(\mathbf{x}) r(\mathbf{x}) d\mathbf{x}$$



$$\text{PE}(p||p') \approx \int p(\mathbf{x}) \hat{r}(\mathbf{x}) d\mathbf{x} - 1$$

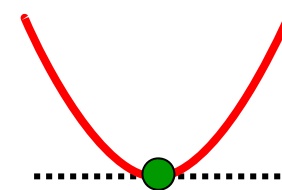
- Expectation is approximated by empirical average.

PE-Divergence Approximation for Linear Model

38

$$\hat{\alpha} = \operatorname{argmin}_{\alpha} \left[\frac{1}{n'} \sum_{i'=1}^{n'} r_{\alpha}(\mathbf{x}'_{i'})^2 - \frac{2}{n} \sum_{i=1}^n r_{\alpha}(\mathbf{x}_i) \right]$$

$$r_{\alpha}(\mathbf{x}) = \alpha^{\top} \phi(\mathbf{x})$$



- Solution is given **analytically**:

$$\hat{\alpha} = \operatorname{argmin}_{\alpha} \left[\alpha^{\top} \hat{G} \alpha - 2 \hat{h}^{\top} \alpha + \lambda \alpha^{\top} \alpha \right]$$

$$= (\hat{G} + \lambda \mathbf{I})^{-1} \hat{h}$$

$$\hat{G} = \frac{1}{n'} \sum_{i'=1}^{n'} \phi(\mathbf{x}'_{i'}) \phi(\mathbf{x}'_{i'})^{\top}$$

$$\hat{h} = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i)$$

- Resulting PE-divergence approximator:

$$\text{PE}(p \| p') \approx \frac{1}{n} \sum_{i=1}^n \hat{\alpha}^{\top} \phi(\mathbf{x}_i) - 1 = \hat{h}^{\top} (\hat{G} + \lambda \mathbf{I})^{-1} \hat{h} - 1$$

MATLAB Implementation for Gauss Kernel Model

$$r_{\alpha}(\mathbf{x}) = \sum_{j=1}^n \alpha_j \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_j\|^2}{2\sigma^2}\right) \quad \text{PE}(p\|p') \approx \hat{\mathbf{h}}^{\top} (\hat{\mathbf{G}} + \lambda \mathbf{I})^{-1} \hat{\mathbf{h}} - 1$$

$$\hat{h}_j = \frac{1}{n} \sum_{i=1}^n \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right) \quad \hat{G}_{j,j'} = \frac{1}{n'} \sum_{i'=1}^{n'} \exp\left(-\frac{\|\mathbf{x}'_{i'} - \mathbf{x}_j\|^2}{2\sigma^2}\right) \exp\left(-\frac{\|\mathbf{x}'_{i'} - \mathbf{x}_{j'}\|^2}{2\sigma^2}\right)$$

```
n=1000; x=randn(n,1); y=randn(n,1)+1/2;
x2=x.^2; xx=repmat(x2,1,n)+repmat(x2',n,1)-2*x*x'; s=exp(-xx);
y2=y.^2; yx=repmat(y2,1,n)+repmat(x2',n,1)-2*y*x'; t=exp(-yx);
PE=mean(s*((t'*t/n+eye(n))\((mean(s,2)))))-1;
```

- **Relative density ratio** can also be estimated in the almost same way.

f-Divergences and Duality

40

$$F(p||p') = \int p'(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{p'(\mathbf{x})}\right) d\mathbf{x}$$

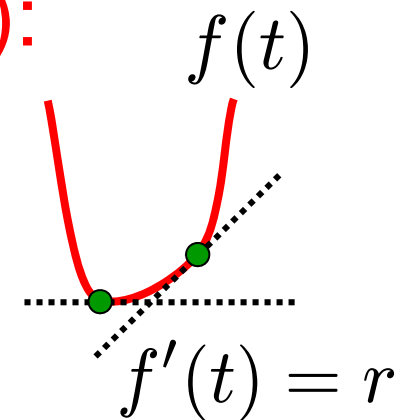
f : Convex function such that $f(1) = 0$

■ Fenchel transform (convex conjugate):

$$f^*(r) = -\inf_t [f(t) - rt]$$

■ Conjugate of conjugate:

$$f(t) = -\inf_r [f^*(r) - rt]$$



● KL-divergence: $f(t) = t \log t$ $f^*(r) = \exp(r - 1)$

● PE-divergence: $f(t) = (t - 1)^2$ $f^*(r) = r^2 / 2 + r$

Lower Bound of f-Divergences 41

Nguyen *et al.* (NIPS2007, IEEE-IT2010)

$$F(p||p') = \int p'(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{p'(\mathbf{x})}\right) d\mathbf{x}$$

$$f(t) = -\inf_r [f^*(r) - rt]$$

■ Lower bound of f-divergences:

$$F(p||p') = -\inf_r \left[\int p'(\mathbf{x}) f^*(r(\mathbf{x})) d\mathbf{x} - \int p(\mathbf{x}) r(\mathbf{x}) d\mathbf{x} \right]$$

■ Sample approximation gives

$$\hat{F}(p||p') = -\min_{\alpha} \left[\frac{1}{n'} \sum_{i'=1}^{n'} f^*(r_{\alpha}(\mathbf{x}'_{i'})) - \frac{1}{n} \sum_{i=1}^n r_{\alpha}(\mathbf{x}_i) \right]$$



Contents

42

1. Distributional change detection
 - A) Problem setup and motivating examples
 - B) Distances
 - C) Distance approximation
 - I. Density-ratio divergences
 - II. Density-difference distances
2. Structural change detection

L²-Distance Approximation

43

Kim & Scott (IEEE-TPAMI2010)

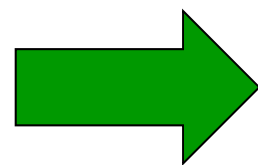
Sugiyama *et al.* (NIPS2012, NeCo2013)

$$L^2(p, p') = \int f(\mathbf{x})^2 d\mathbf{x} \quad f(\mathbf{x}) = p(\mathbf{x}) - p'(\mathbf{x})$$

- Directly approximate the density difference by LS:

$$\hat{f} = \operatorname{argmin}_{\tilde{f}} \int \left(\tilde{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2 d\mathbf{x}$$

$$= \operatorname{argmin}_{\tilde{f}} \int \left(\tilde{f}(\mathbf{x}) \right)^2 d\mathbf{x} - 2 \int f(\mathbf{x}) \tilde{f}(\mathbf{x}) d\mathbf{x}$$



$$L^2(p, p') \approx \int \hat{f}(\mathbf{x})^2 d\mathbf{x}$$

- Expectation is approximated by empirical average.

Solution for Linear Model

44

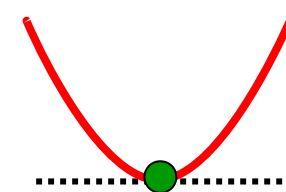
$$\hat{\alpha} = \operatorname{argmin}_{\alpha} \int \left(f_{\alpha}(\mathbf{x}) \right)^2 d\mathbf{x} + \frac{1}{n'} \sum_{i'=1}^{n'} f_{\alpha}(\mathbf{x}'_{i'})^2 - \frac{1}{n} \sum_{i=1}^n f_{\alpha}(\mathbf{x}_i)$$

$$f_{\alpha}(\mathbf{x}) = \sum_{j=1}^b \alpha_j \phi_j(\mathbf{x}) = \alpha^{\top} \phi(\mathbf{x})$$

- (Regularized) solution is given analytically:

$$\hat{\alpha} = \operatorname{argmin}_{\alpha} \left[\alpha^{\top} \mathbf{G} \alpha - 2 \hat{\mathbf{h}}^{\top} \alpha + \lambda \alpha^{\top} \alpha \right]$$

$$= (\mathbf{G} + \lambda \mathbf{I})^{-1} \hat{\mathbf{h}}$$



$$\mathbf{G} = \int \phi(\mathbf{x}) \phi(\mathbf{x})^{\top} d\mathbf{x}$$

$$\hat{\mathbf{h}} = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) - \frac{1}{n'} \sum_{i'=1}^{n'} \phi(\mathbf{x}'_{i'})$$

Resulting L²-Distance Approximator ⁴⁵

- Two ways to approximate the L²-distance by density-difference estimation:

- $L^2(p, p') = \int f(\mathbf{x})^2 d\mathbf{x} \approx \hat{\alpha}^\top \mathbf{G} \hat{\alpha}$

$$\hat{\alpha} = (\mathbf{G} + \lambda \mathbf{I})^{-1} \hat{\mathbf{h}} \quad f(\mathbf{x}) = p(\mathbf{x}) - p'(\mathbf{x}) \approx \hat{\alpha}^\top \phi(\mathbf{x})$$

- $L^2(p, p') = \int (p(\mathbf{x}) - p'(\mathbf{x})) f(\mathbf{x}) d\mathbf{x} \approx \hat{\mathbf{h}}^\top \boldsymbol{\alpha}$

$$\mathbf{G} = \int \phi(\mathbf{x}) \phi(\mathbf{x})^\top d\mathbf{x}$$

$$\hat{\mathbf{h}} = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) - \frac{1}{n'} \sum_{i'=1}^{n'} \phi(\mathbf{x}'_{i'})$$

Bias Reduction

- Consider their **linear combination**:

$$\kappa \hat{\mathbf{h}}^\top \hat{\boldsymbol{\alpha}} + (1 - \kappa) \hat{\boldsymbol{\alpha}}^\top \mathbf{G} \hat{\boldsymbol{\alpha}} \quad \kappa \in \mathbb{R}$$

- For small λ ,

$$\kappa \hat{\mathbf{h}}^\top \hat{\boldsymbol{\alpha}} + (1 - \kappa) \hat{\boldsymbol{\alpha}}^\top \mathbf{G} \hat{\boldsymbol{\alpha}}$$

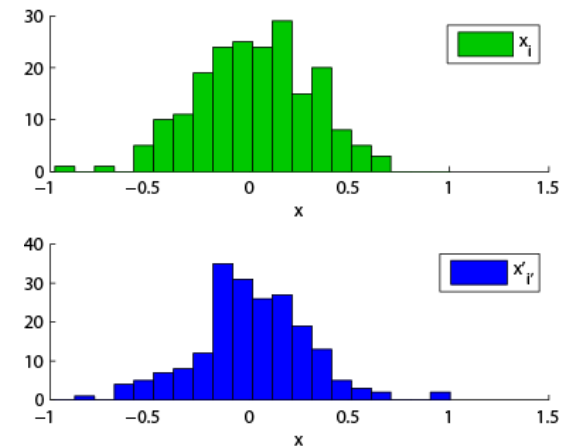
$$= \hat{\mathbf{h}}^\top \mathbf{G}^{-1} \hat{\mathbf{h}} - \lambda(2 - \kappa) \hat{\mathbf{h}}^\top \mathbf{G}^{-2} \hat{\mathbf{h}} + o_p(\lambda)$$

- $\kappa = 2$ removes the **regularization-induced bias**:

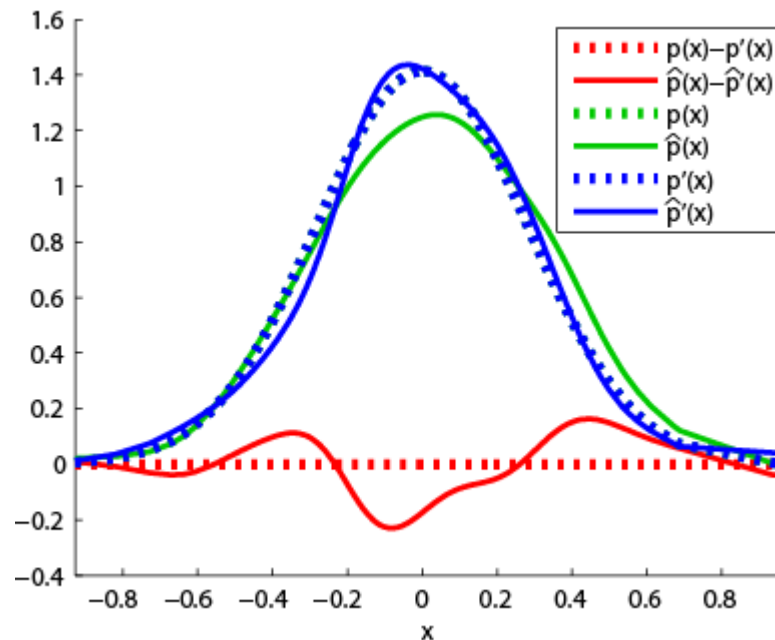
$$\hat{L}^2(\mathcal{X}, \mathcal{X}') = 2 \hat{\mathbf{h}}^\top \hat{\boldsymbol{\alpha}} - \hat{\boldsymbol{\alpha}}^\top \mathbf{G} \hat{\boldsymbol{\alpha}}$$

Density-Difference Estimation (1)⁴⁷

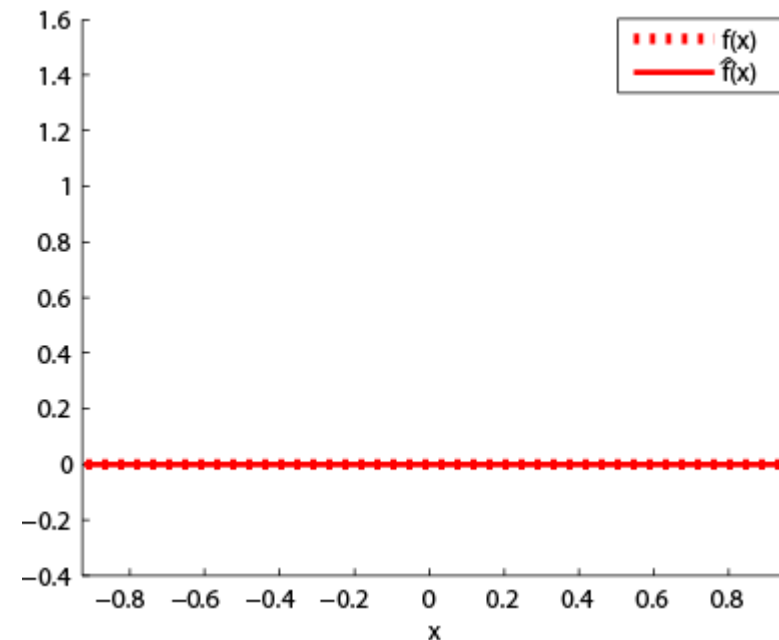
■ $p(x) = p'(x) = N(x; 0, (4\pi)^{-1})$
 $n = n' = 200$



Difference of kernel
Density estimators (KDE)



Least-squares density
-difference estimation (LSDD)

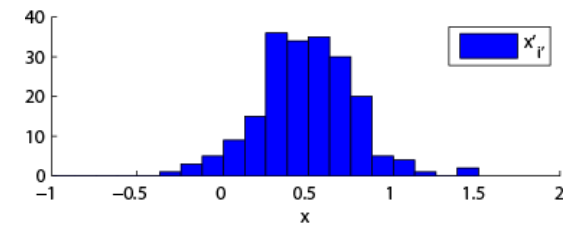
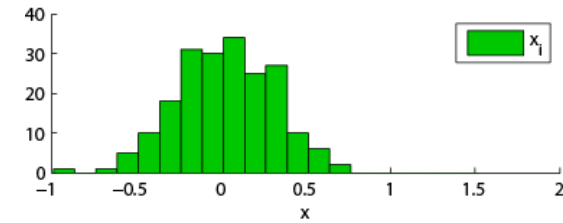


Density-Difference Estimation (2)⁴⁸

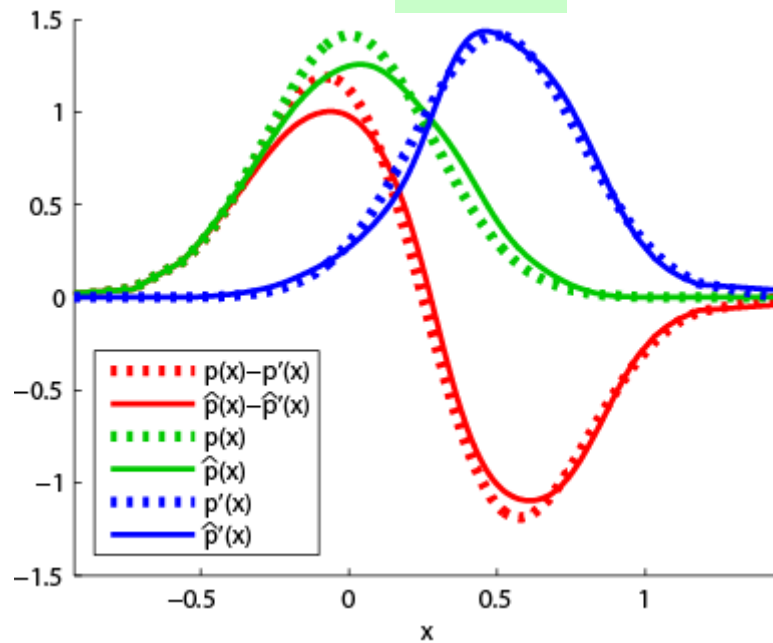
■ $p(x) = N(x; 0, (4\pi)^{-1})$

■ $p'(x) = N(x; 0.5, (4\pi)^{-1})$

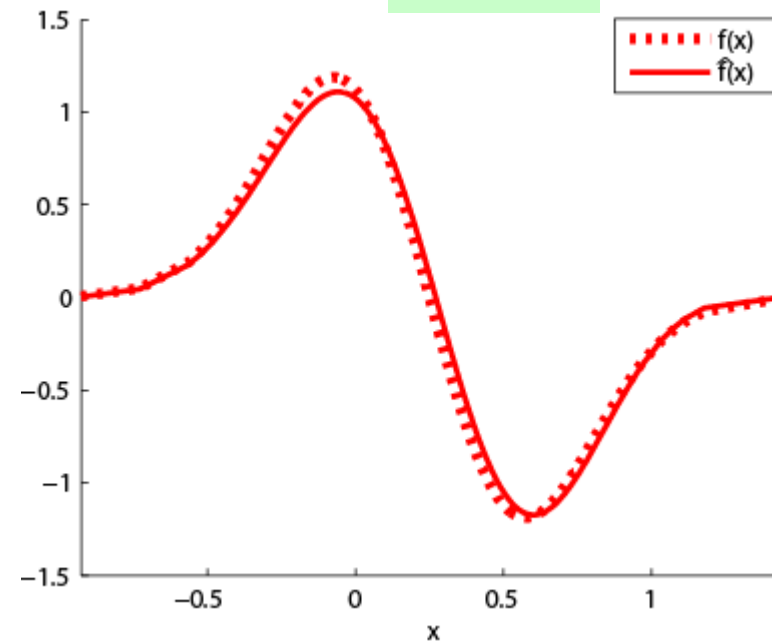
$n = n' = 200$



KDE



LSDD

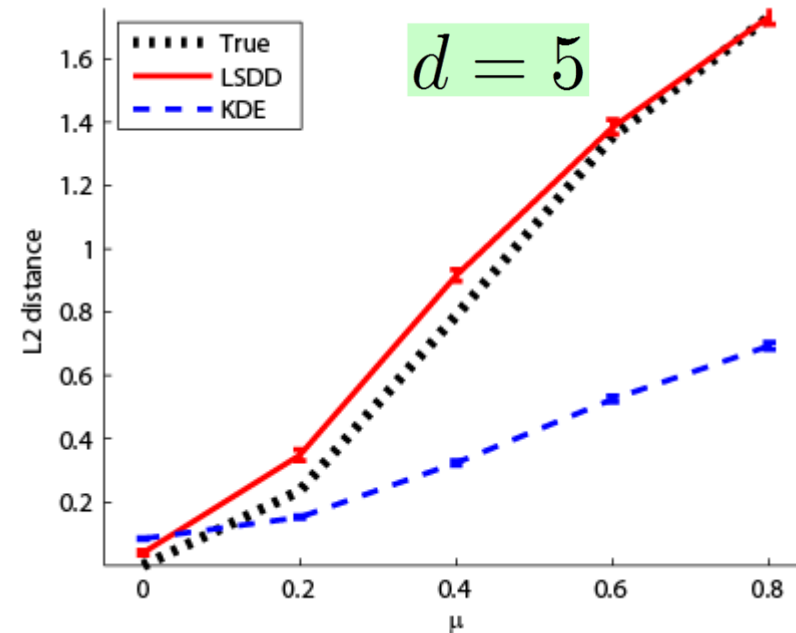
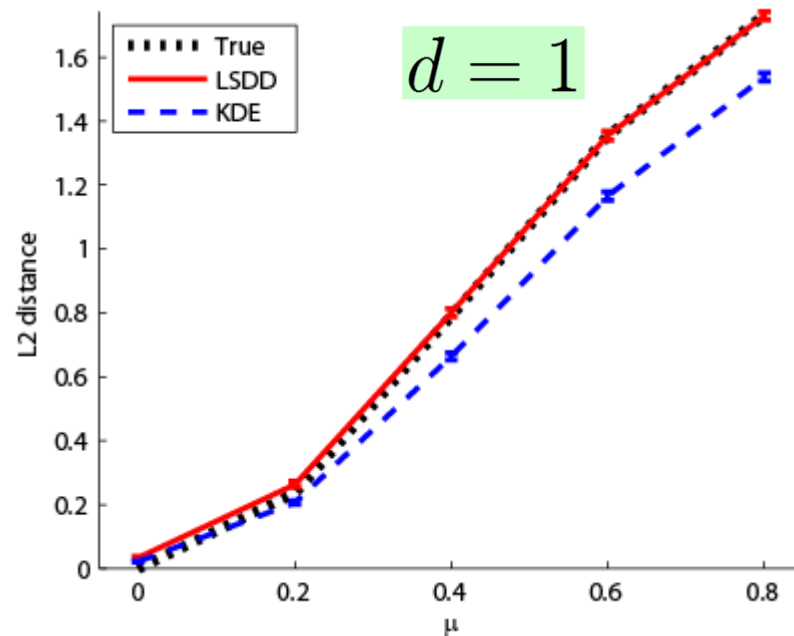


L²-Distance Approximation

49

■ $p(\mathbf{x}) = N(\mathbf{x}; (\mu, 0, \dots, 0)^\top, (4\pi)^{-1} \mathbf{I}_d)$ $n = n' = 100$

■ $p'(\mathbf{x}) = N(\mathbf{x}; (0, 0, \dots, 0)^\top, (4\pi)^{-1} \mathbf{I}_d)$

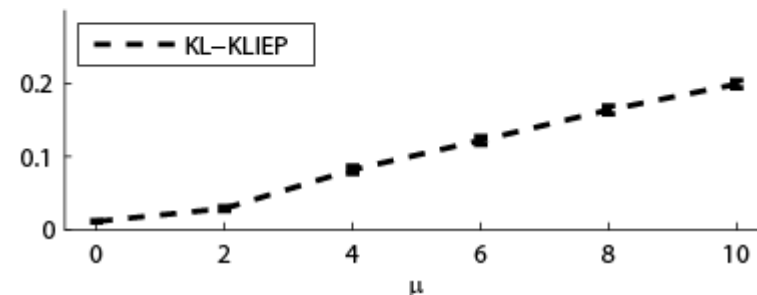
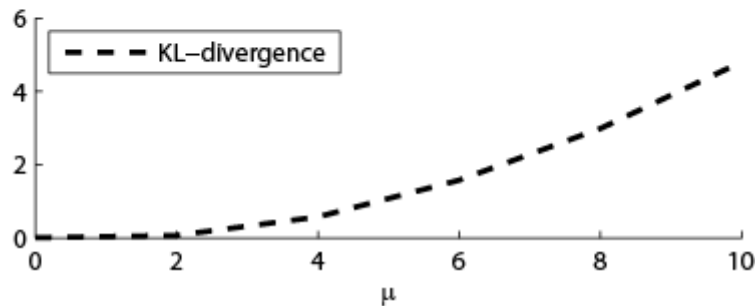
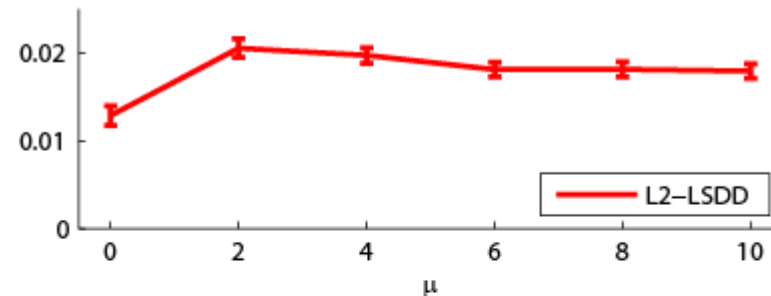
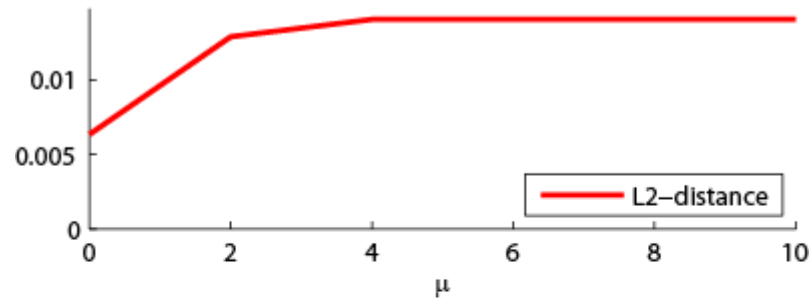
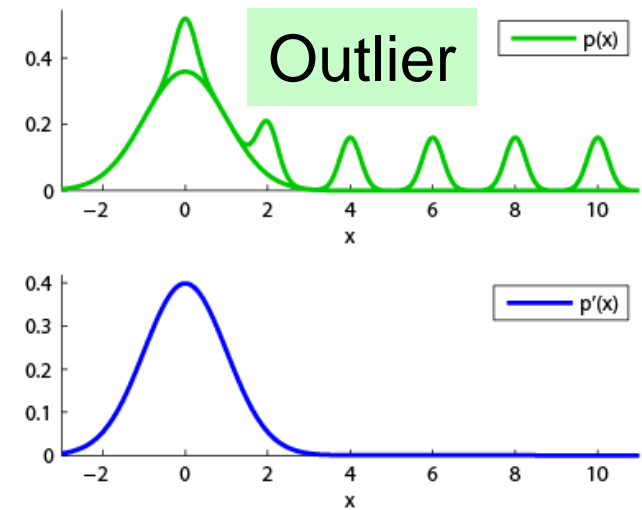


- KDE significantly under-estimates.
- LSDD slightly over-estimates.

L²-Distance vs. KL-Divergence 50

$$L^2(p, p') = \int (p(\mathbf{x}) - p'(\mathbf{x}))^2 d\mathbf{x}$$

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



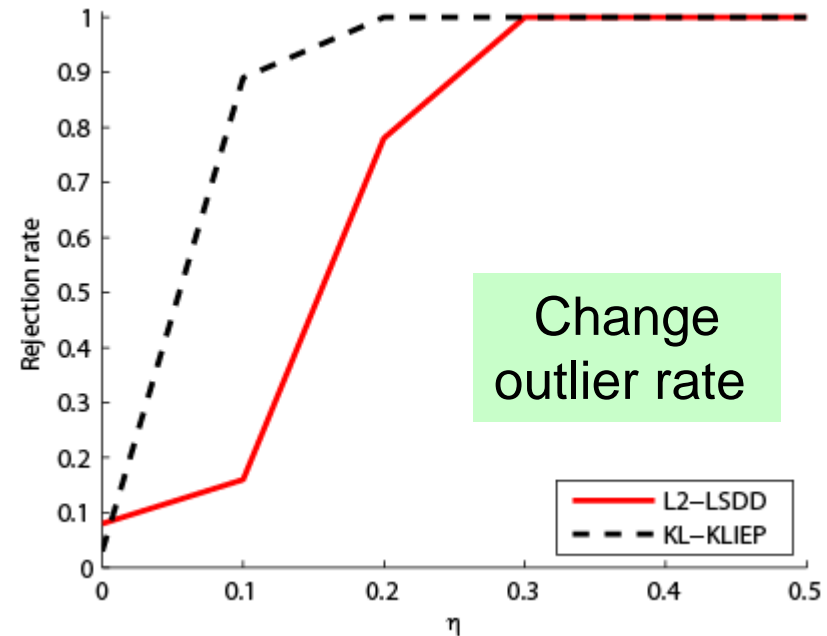
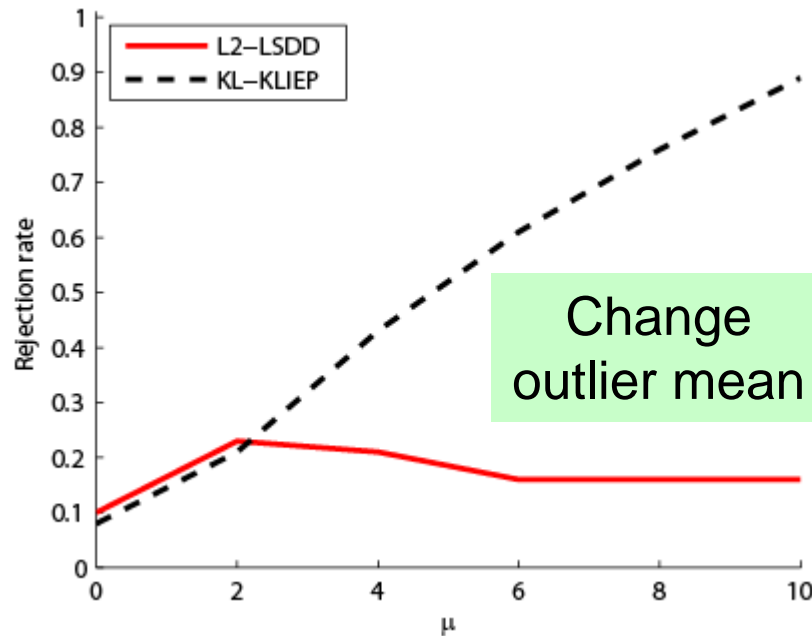
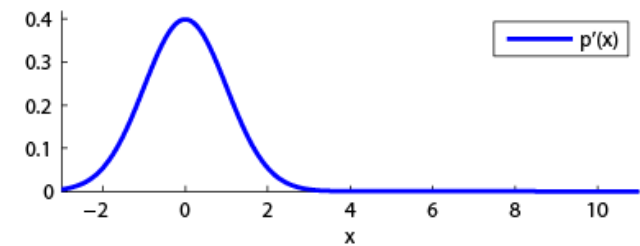
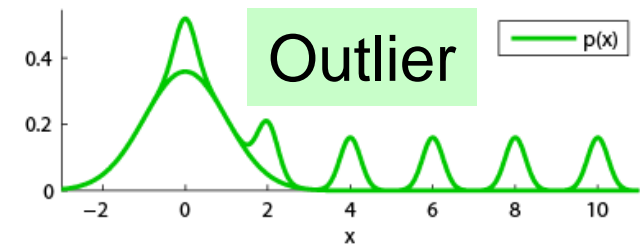
L²-distance is less sensitive to outliers.

Robust Two-Sample Test

51

■ **Two-sample test:** Are two distributions the same?

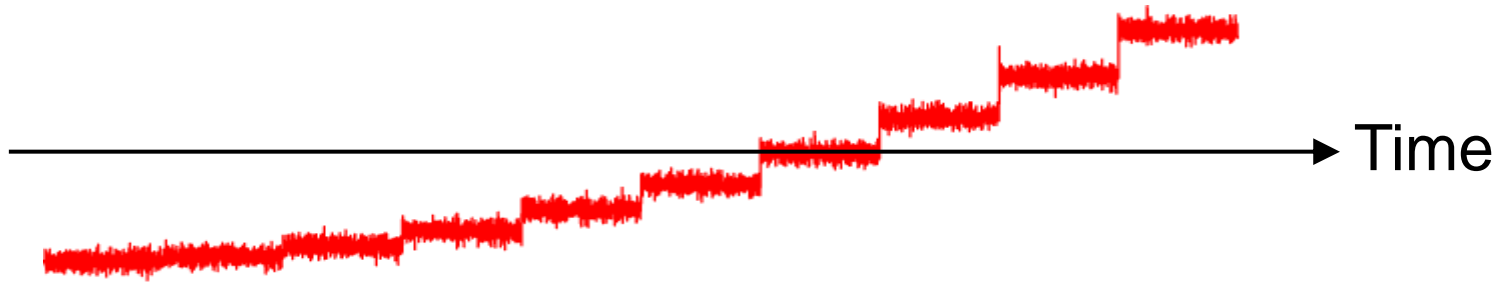
- **Null:** Two are the same
- **Alternative:** Two are different



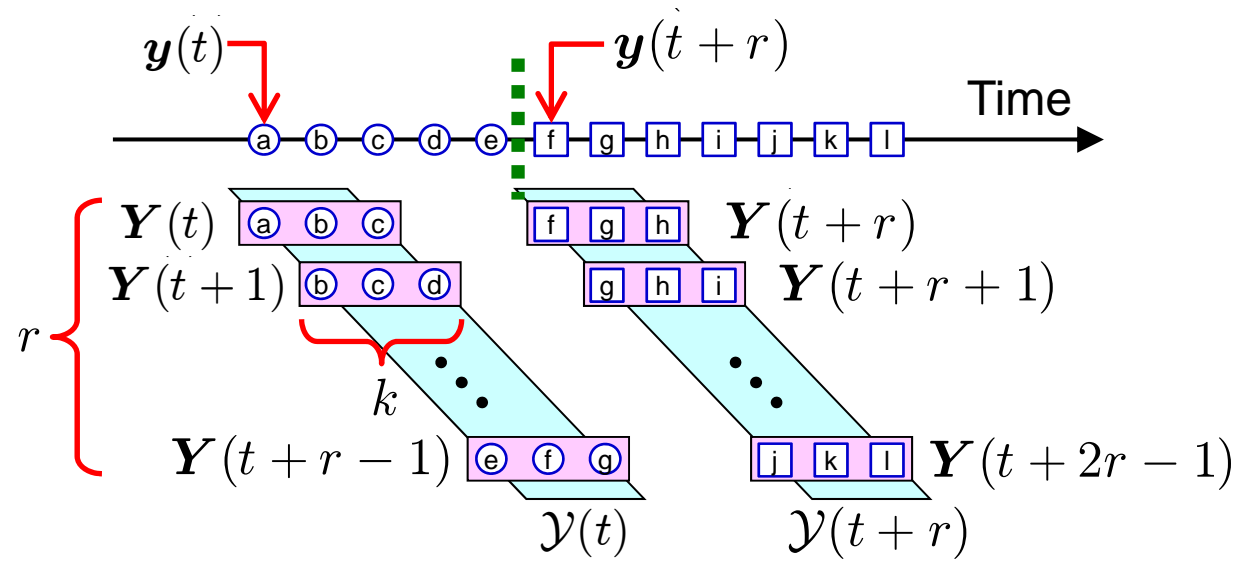
L²-based test is more robust against outliers.

Unsupervised Change Detection ⁵²

- Identify change points in time-series:



- Use the distance between the distributions of sliding-windowed past and current data.

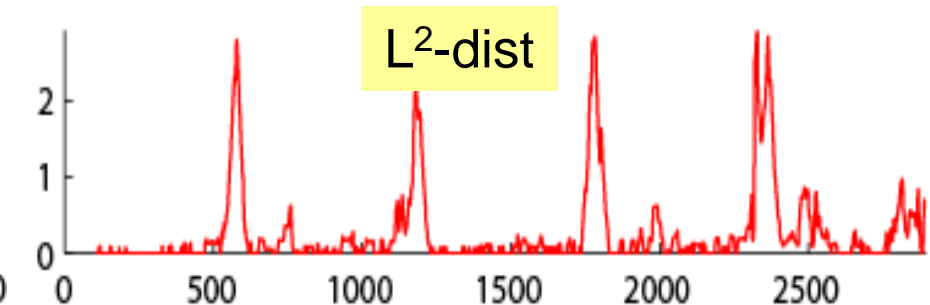
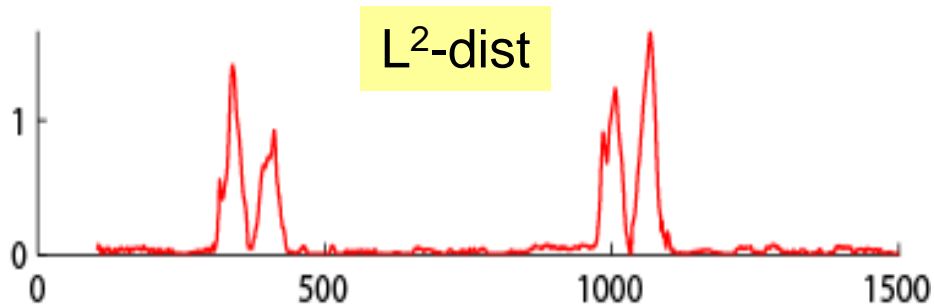
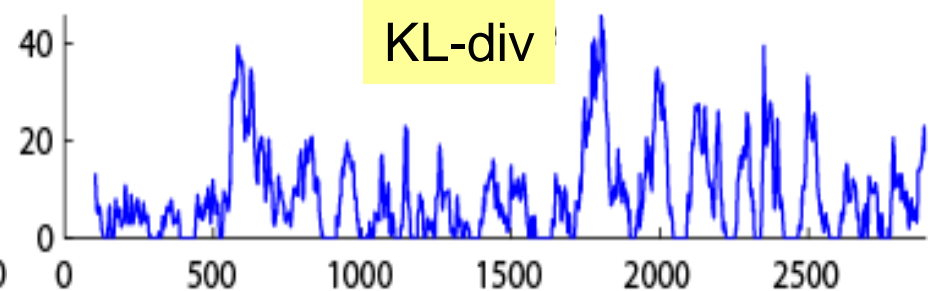
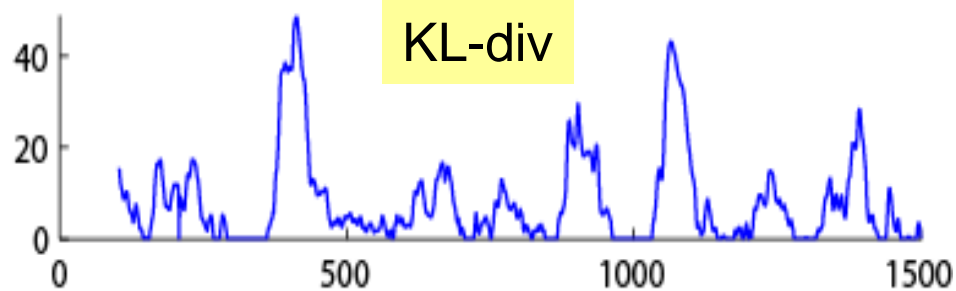
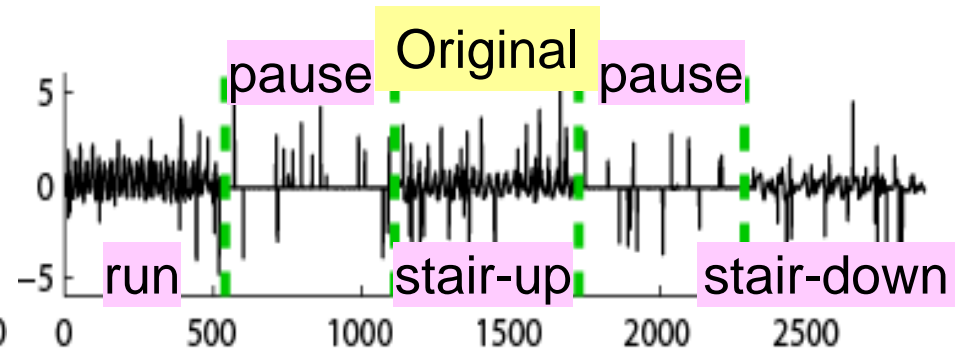
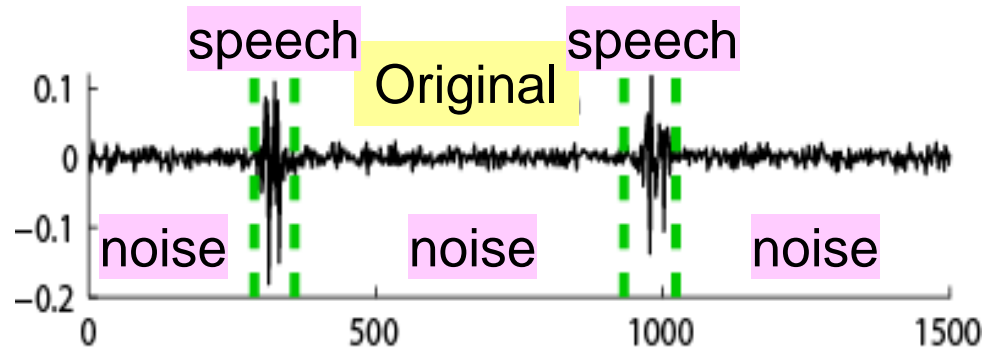


Results

53

CENSREC Speech Data

HASC Accelerometer Data



L² is more robust against noise.

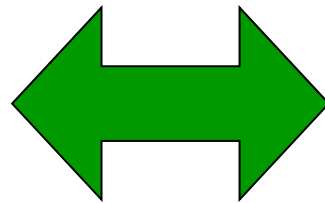
Mutual Information

$$\text{MI} = \iint p(\mathbf{x}, \mathbf{y}) \log \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} d\mathbf{x}d\mathbf{y}$$

- **Mutual information** is the KL-divergence from the joint density $p(\mathbf{x}, \mathbf{y})$ to the product of marginal densities $p(\mathbf{x})p(\mathbf{y})$.
- **Independence** can be measured:

- $\text{MI} \geq 0$

- $\text{MI} = 0$



x and y are independent

Mutual Information Approximation⁵⁵

$$\text{MI} = \iint p(\mathbf{x}, \mathbf{y}) \log \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} d\mathbf{x}d\mathbf{y}$$

- Estimation of density ratio $r(\mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})}$

from $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, \mathbf{y})$

$\{(\mathbf{x}_i, \mathbf{y}_{i'})\}_{i,i'=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})p(\mathbf{y})$

gives an MI approximator:

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

Variations of MI

- Squared-loss MI (Pearson divergence):

$$\text{SMI} = \iint p(\mathbf{x})p(\mathbf{y}) \left(\frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} - 1 \right)^2 d\mathbf{x}d\mathbf{y}$$

- Relative SMI:

$$\text{rSMI} = \iint p_{\beta}(\mathbf{x}, \mathbf{y}) \left(\frac{p(\mathbf{x}, \mathbf{y})}{p_{\beta}(\mathbf{x}, \mathbf{y})} - 1 \right)^2 d\mathbf{x}d\mathbf{y}$$

$$0 \leq \beta < 1 \quad p_{\beta}(\mathbf{x}, \mathbf{y}) = \beta p(\mathbf{x}, \mathbf{y}) + (1 - \beta)p(\mathbf{x})p(\mathbf{y})$$

- Quadratic MI:

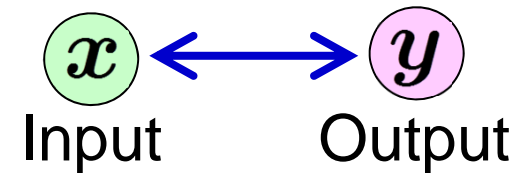
$$\text{QMI} = \iint \left(p(\mathbf{x}, \mathbf{y}) - p(\mathbf{x})p(\mathbf{y}) \right)^2 d\mathbf{x}d\mathbf{y}$$

Usages of MI Approximator

57

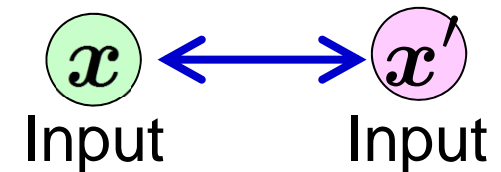
■ MI between input and output:

- Feature selection/extraction
- Clustering



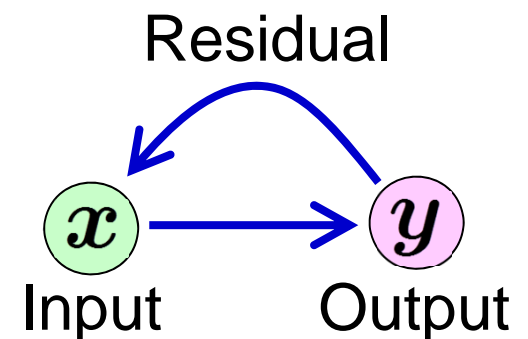
■ MI between inputs:

- Independent component analysis
- Higher-order canonical correlation analysis
- Object matching



■ MI between input and residual:

- Causal direction inference





Summary of Distributional Change Detection

- Compute a **divergence** between distributions:
 - Separate density estimation does not work well, because Vapnik's principle is violated.
 - **Direct estimation of density ratio/difference** seems more sensible.
- Don't simply use KL as a divergence measure just because it is popular.
 - **Relative PE** and L^2 could be more robust against outliers and computationally more efficient.
- **MI** can also be approximated in the same way.



A Little Break: Artist Agent

59

Ning *et al.* (ICML2012)

Brush movement learning by **reinforcement learning**.





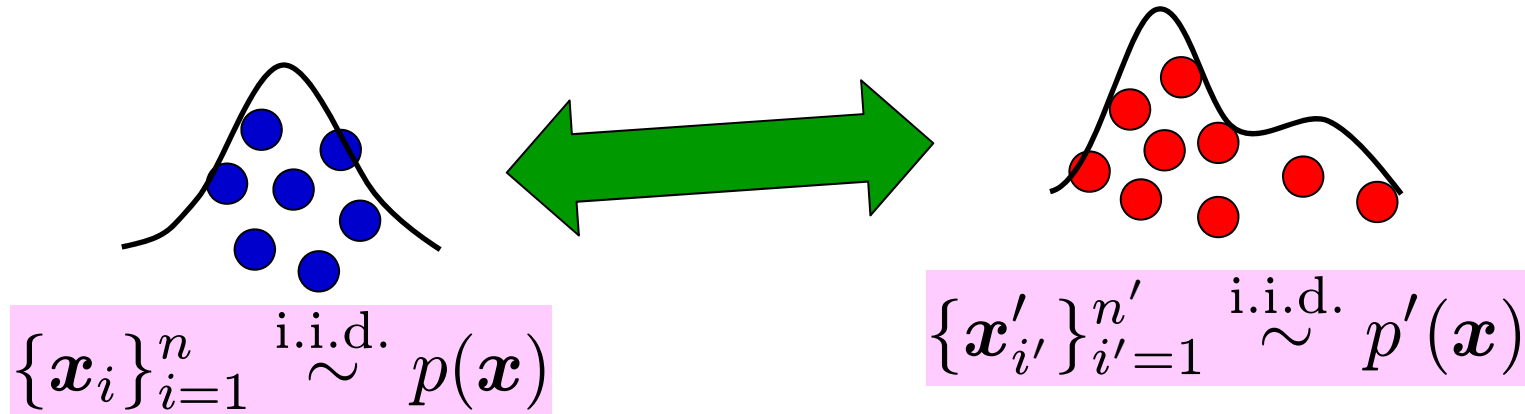


Contents

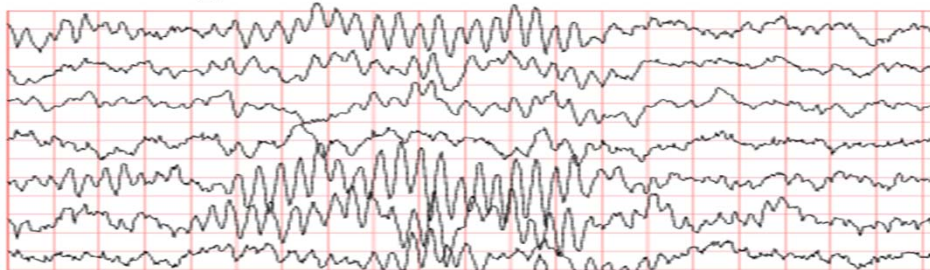
61

1. Distributional change detection
2. **Structural change detection**
 - A) Density estimation approach
 - B) Density-ratio estimation approach

From Distributional Change to Structural Change



- Through distance estimation, **distributional change** can be detected.
- We investigate how distributions are changed through **interaction between variables**.

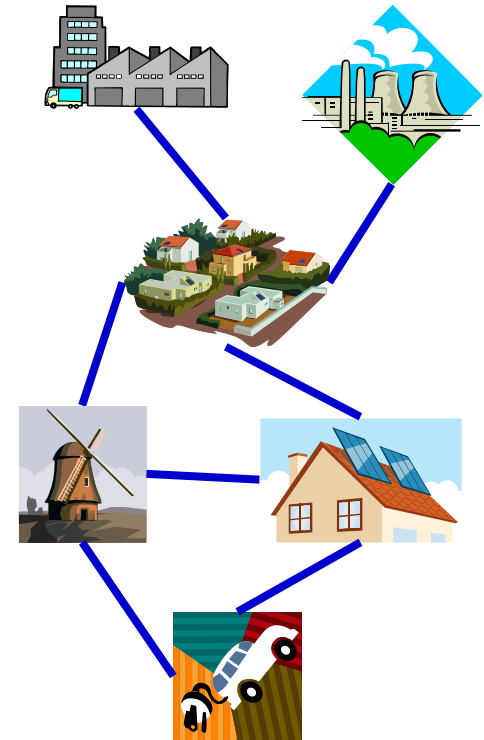
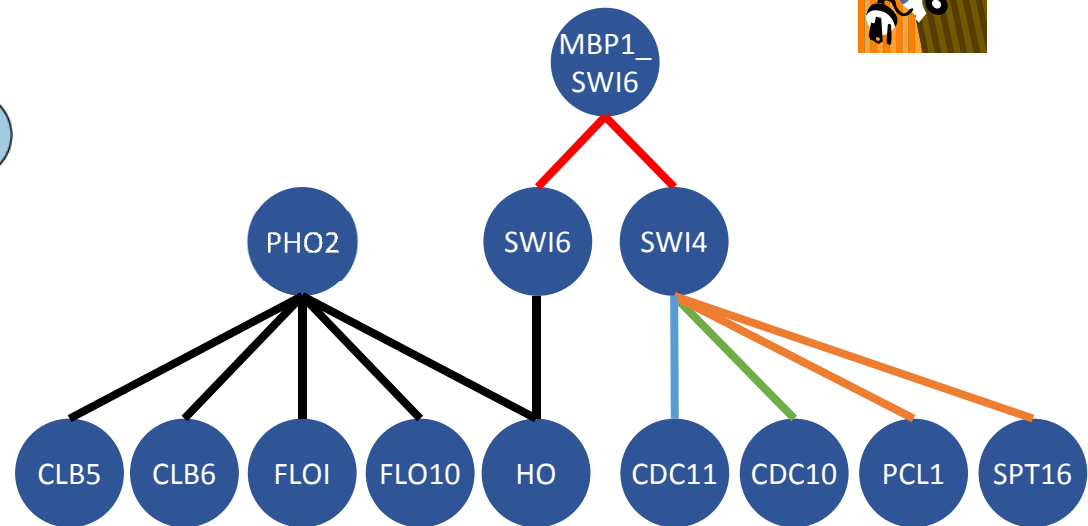
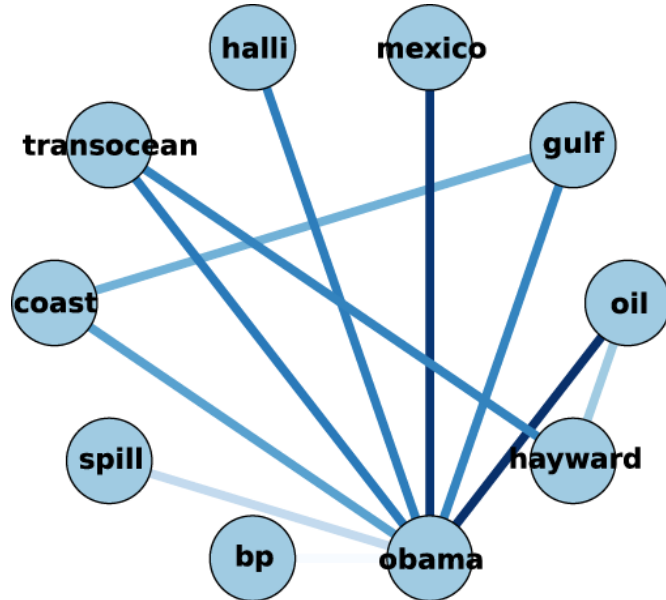


$$\mathbf{x} = (x^{(1)}, \dots, x^{(d)})^\top$$

Motivating Examples

63

- Word co-occurrence in Twitter
- Gene regulatory networks
- Fraud detection in smart grid





Contents

64

1. Distributional change detection
2. Structural change detection
 - A) Density estimation approach
 - I. Gauss models
 - II. Non-Gauss models
 - B) Density-ratio estimation approach

Gauss Model

65

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

Θ : (sparse) inverse covariance matrix

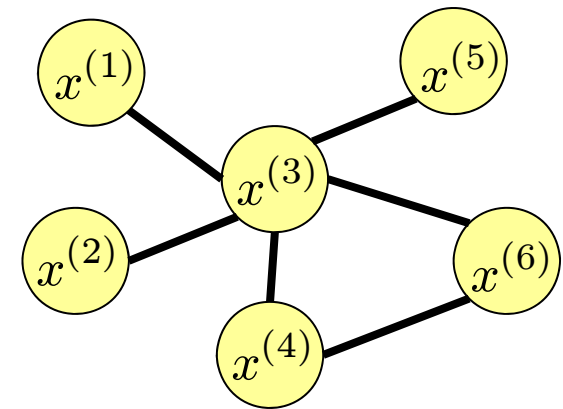
■ **Conditional independence:**

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

$$\Theta_{k,k'} = 0 \iff x^{(k)} \perp\!\!\!\perp x^{(k')} \mid \{x^{(\ell)}\}_{\ell \neq k,k'}$$

■ **Graphical representation:**

- **Node:** Each variable
- **Edge:** Exists if $\Theta_{i,j} \neq 0$
- **Only connected variables affect!**



$$x^{(1)} \perp\!\!\!\perp x^{(2)} \mid x^{(3)}$$

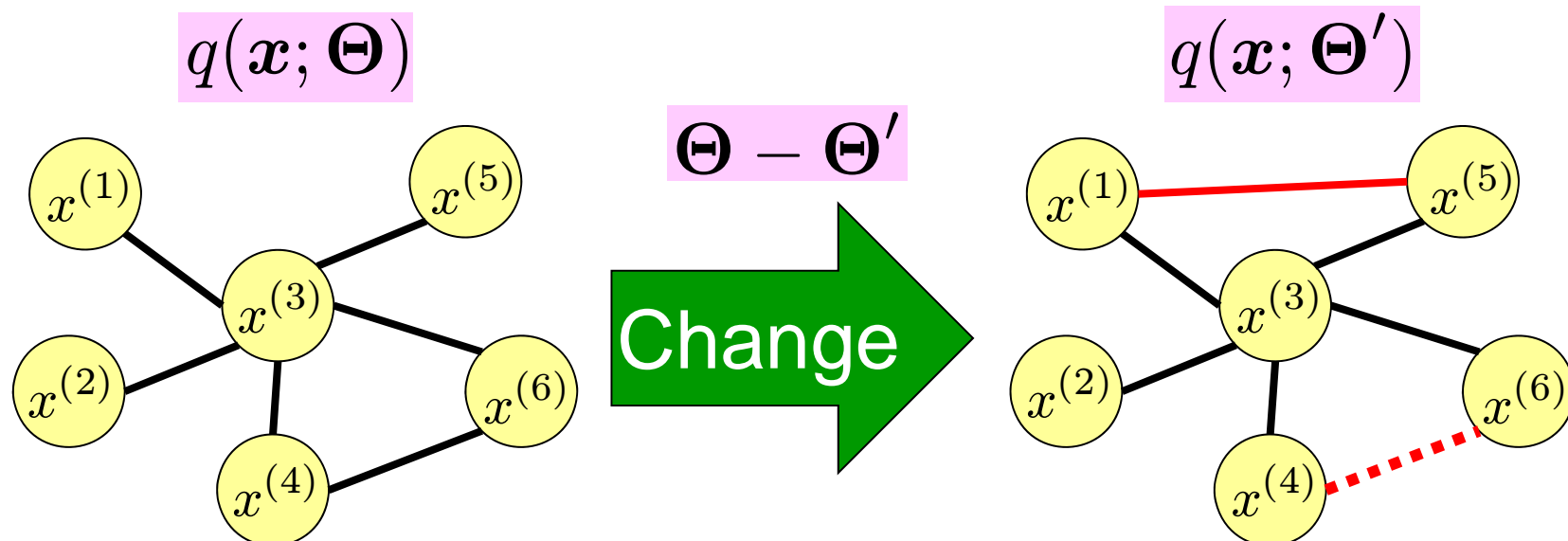
Structural Change Detection with Gauss Models

66

- Use Gauss models for $p(x)$ and $p'(x)$:

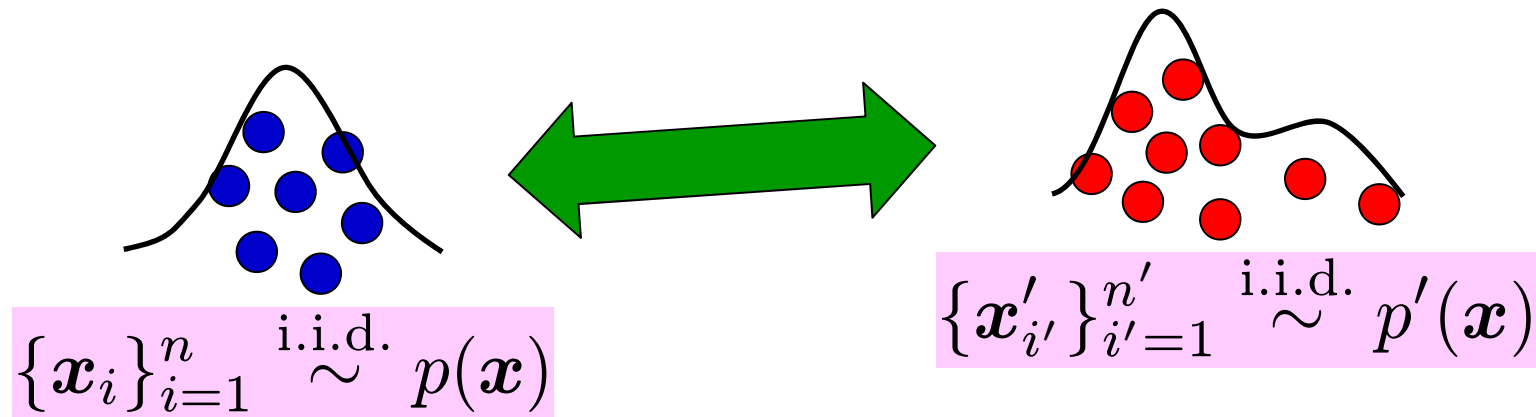
$$q(x; \Theta) = \frac{\det(\Theta)^{1/2}}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2}x^\top \Theta x\right) \quad q(x; \Theta')$$

- Detect **sparse change** in covariance structure:



Structural Change Detection by Graphical Lasso (Glasso)

Tibshirani (JRSS1996), Friedman *et al.* (Biostat2008)



■ Sparse maximum likelihood estimation:

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

$$\max_{\Theta'} \sum_{i'=1}^{n'} \log q(\mathbf{x}'_{i'}; \Theta') - \lambda' \|\Theta'\|_1$$

$$q(\mathbf{x}; \Theta) = \frac{\det(\Theta)^{1/2}}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} \mathbf{x}^\top \Theta \mathbf{x}\right) \quad \lambda, \lambda' \geq 0$$

Structural Change Detection by Glasso

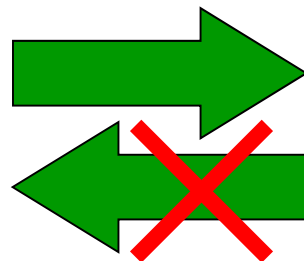
68

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

$$\max_{\Theta'} \sum_{i'=1}^{n'} \log q(\mathbf{x}'_{i'}; \Theta') - \lambda' \|\Theta'\|_1$$

- 😊 Scalable to high-dimensional datasets.
- 😊 Statistical properties have been well studied.
- 😞 Does not work if true Θ and Θ' are dense.
- 😞 Choice of λ and λ' is not straightforward.

Both Θ and Θ'
are sparse



Change $\Theta - \Theta'$
is sparse

Structural Change Detection by Fused Lasso (Flasso)

69

Tibshirani *et al.* (JRSS2005)
Zhang & Wang (UAI2010)

- Directly penalize **the difference of parameters** to be sparse:

$$\max_{\Theta, \Theta'} \sum_{i=1}^n \log q(\mathbf{x}_i; \Theta) + \sum_{i'=1}^{n'} \log q(\mathbf{x}'_{i'}; \Theta') - \gamma \|\Theta - \Theta'\|_1$$

$$\gamma \geq 0$$

- ☺ Scalable to high-dimensional datasets.
- ☺ Work well even if true Θ and Θ' are dense.



Contents

70

1. Distributional change detection
2. Structural change detection
 - A) Density estimation approach
 - I. Gauss models
 - II. Non-Gauss models
 - B) Density-ratio estimation approach

Correlation and Dependence

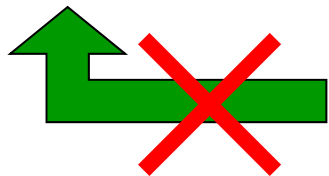
71

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

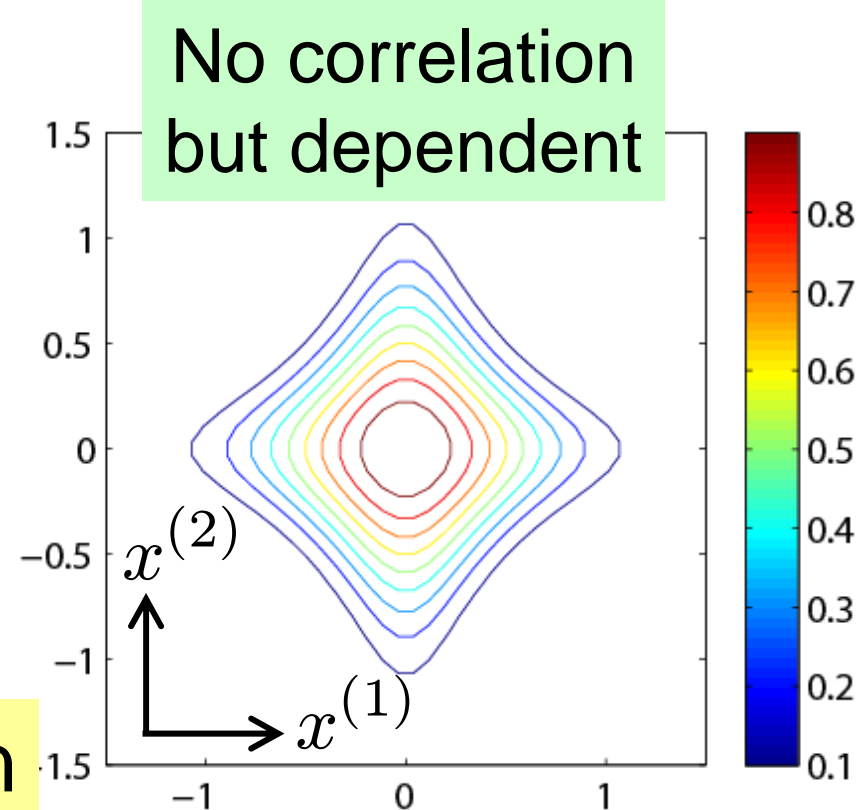
Θ : (sparse) inverse covariance matrix

- Gauss models cannot capture **higher-order correlations**.
- No correlation does not imply independence.

Independence



No correlation



Nonparanormal Models

72

Liu *et al.* (JMLR2009)

- Gaussian after **element-wise transformation**:

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

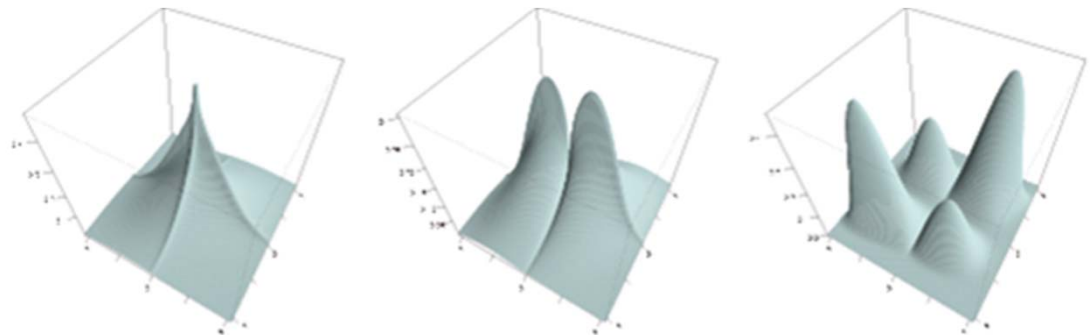
$$\mathbf{f}(\mathbf{x}) = (f_1(x^{(1)}), \dots, f_d(x^{(d)}))^\top$$

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

f_k : Monotone and differentiable function

😊 More flexible than ordinary Gauss models.

☹ Still restrictive
in representation



Pairwise Markov Networks

73

$$q(\mathbf{x}; \boldsymbol{\theta}) = \frac{\bar{q}(\mathbf{x}; \boldsymbol{\theta})}{Z(\boldsymbol{\theta})}$$

$$\bar{q}(\mathbf{x}; \boldsymbol{\theta}) = \exp \left(\sum_{k \geq k'} \boldsymbol{\theta}_{k,k'}^\top \mathbf{f}(x^{(k)}, x^{(k')}) \right)$$

$\mathbf{f}(x, x')$: feature vector

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

$$\boldsymbol{\theta} = (\boldsymbol{\theta}_{1,1}^\top, \dots, \boldsymbol{\theta}_{d,d}^\top)^\top$$

■ Gaussian: $\mathbf{f}(x, x') = xx'$

■ Nonparanormal: $\mathbf{f}(x, x') = f(x)f(x')$

■ Polynomial: $\mathbf{f}(x, x') = [x^t, x^{t-1}x', \dots, x, x', 1]^\top$

☺ High representation capability.

☹ Normalization $Z(\boldsymbol{\theta}) = \int \bar{q}(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}$ is intractable.

Importance Sampling

74

1. Draw **pseudo-samples** from a **proposal density**:

$$\{\mathbf{x}_{i''}''\}_{i''=1}^{n''} \stackrel{\text{i.i.d.}}{\sim} p''(\mathbf{x}) \quad (\text{e.g., Gaussian})$$

2. Approximate the integration by **importance-weighted** sample average:

$$Z(\boldsymbol{\theta}) = \int \bar{q}(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} = \int \frac{\bar{q}(\mathbf{x}; \boldsymbol{\theta})}{p''(\mathbf{x})} p''(\mathbf{x}) d\mathbf{x}$$

$$\approx \frac{1}{n''} \sum_{i''=1}^{n''} \frac{\bar{q}(\mathbf{x}_{i''}''; \boldsymbol{\theta})}{p(\mathbf{x}_{i''}''')} \xrightarrow{n'' \rightarrow \infty} \int \bar{q}(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}$$

- ☺ Law of large numbers guarantees **consistency**.
- ☹ Unstable due to **large variance**.

Score Matching

75

Hyvärinen (JMLR2005)

- Learn unnormalized density model $\bar{q}(\mathbf{x}; \boldsymbol{\theta})$ by least-squares matching of **score functions**:

$$\min_{\boldsymbol{\theta}} \int p(\mathbf{x}) \|\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}) - \nabla_{\mathbf{x}} \log p(\mathbf{x})\|^2 d\mathbf{x}$$

$$\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}) = \nabla_{\mathbf{x}} \log \bar{q}(\mathbf{x}; \boldsymbol{\theta}) \quad \nabla_{\mathbf{x}} = (\partial_{x^{(1)}}, \dots, \partial_{x^{(d)}})^{\top}$$

- Empirical version (use **integration-by-parts**):

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n S(\mathbf{x}_i; \boldsymbol{\theta}) \quad S(\mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^d \left(\psi_k(\mathbf{x}; \boldsymbol{\theta})^2 + 2\partial_{x^{(k)}} \psi_k(\mathbf{x}; \boldsymbol{\theta}) \right)$$

$$\int \psi_k(\mathbf{x}; \boldsymbol{\theta}) \partial_{x^{(k)}} p(\mathbf{x}) d\mathbf{x} = - \int \partial_{x^{(k)}} \psi_k(\mathbf{x}; \boldsymbol{\theta}) p(\mathbf{x}) d\mathbf{x}$$

☺ No normalization is needed.



Contents

76

1. Distributional change detection
2. Structural change detection
 - A) Density estimation approach
 - B) Density-ratio estimation approach

Avoiding Density Estimation

77

■ Fused lasso + Score matching:

$$\gamma \geq 0$$

$$\max_{\Theta, \Theta'} \sum_{i=1}^n S(\mathbf{x}_i; \boldsymbol{\theta}) + \sum_{i'=1}^{n'} S(\mathbf{x}'_{i'}; \boldsymbol{\theta}') - \gamma \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_1$$

- 😊 Work well even if true Θ and Θ' are dense.
- 😊 Higher-order correlations can be captured.
- ☹ Still need explicit modeling of $p(\mathbf{x})$ and $p'(\mathbf{x})$.

■ Vapnik's principle:

*Don't solve
a more general problem*



Direct Change Modeling in Markov Networks

Liu *et al.* (ECML2013, NeCo2014)

- Without separately modeling $p(\mathbf{x})$ and $p'(\mathbf{x})$, we directly model the **density ratio** $p(\mathbf{x})/p'(\mathbf{x})$:

$$r(\mathbf{x}) = \frac{p(\mathbf{x})}{p'(\mathbf{x})} \approx \frac{q(\mathbf{x}; \boldsymbol{\theta})}{q(\mathbf{x}; \boldsymbol{\theta}')} \propto \exp \left(\sum_{k \geq k'} (\boldsymbol{\theta}_{k,k'} - \boldsymbol{\theta}'_{k,k'})^\top \mathbf{f}(x^{(k)}, x^{(k')}) \right)$$

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

- Individual parameters** $\boldsymbol{\theta}, \boldsymbol{\theta}'$ are not necessary, but their **difference** $\boldsymbol{\alpha} = \boldsymbol{\theta} - \boldsymbol{\theta}'$ is enough.

Ratio of Markov Network Models⁷⁹

$$r_{\boldsymbol{\alpha}}(\mathbf{x}) = \frac{1}{N(\boldsymbol{\alpha})} \exp \left(\sum_{k \geq k'} \boldsymbol{\alpha}_{k,k'}^{\top} \mathbf{f}(x^{(k)}, x^{(k')}) \right)$$

■ Normalization:

$$\boldsymbol{\alpha} = (\boldsymbol{\alpha}_{1,1}^{\top}, \dots, \boldsymbol{\alpha}_{d,d}^{\top})^{\top}$$

$$r(\mathbf{x}) = \frac{p(\mathbf{x})}{p'(\mathbf{x})} \implies \int p'(\mathbf{x}) r(\mathbf{x}) d\mathbf{x} = \int p(\mathbf{x}) d\mathbf{x} = 1$$

☺ Simple sample averaging is consistent:

$$N(\boldsymbol{\alpha}) = \int p'(\mathbf{x}) \exp \left(\sum_{k \geq k'} \boldsymbol{\alpha}_{k,k'}^{\top} \mathbf{f}(x^{(k)}, x^{(k')}) \right) d\mathbf{x}$$

$$\approx \frac{1}{n'} \sum_{i'=1}^{n'} \exp \left(\sum_{k \geq k'} \boldsymbol{\alpha}_{k,k'}^{\top} \mathbf{f}(x'_{i'}^{(k)}, x'_{i'}^{(k')}) \right)$$

Sparse Density-Ratio Estimation⁸⁰

Sugiyama *et al.* (NIPS2007, AISM2008)

- Density-ratio matching under KL-divergence:

$$\min_{\alpha} \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{p'(\mathbf{x})r_{\alpha}(\mathbf{x})} d\mathbf{x}$$

$$r_{\alpha}(\mathbf{x}) \approx \frac{p(\mathbf{x})}{p'(\mathbf{x})}$$

- Sample approximation gives

$$\min_{\alpha} \log \frac{1}{n'} \sum_{i'=1}^{n'} \exp \left(\sum_{k \geq k'} \alpha_{k,k'}^{\top} \mathbf{f}(x_{i'}^{(k)}, x_{i'}^{(k')}) \right) - \frac{1}{n} \sum_{i=1}^n \sum_{k \geq k'} \alpha_{k,k'}^{\top} \mathbf{f}(x_i^{(k)}, x_i^{(k')})$$

- Tractable for any feature $\mathbf{f}(x^{(k)}, x^{(k')})$.
- Add a smoothing regularizer: $+ \eta \|\alpha\|^2$
- Add a **group-sparsity** regularizer: $+ \gamma \sum_{k \geq k'} \|\alpha_{k,k'}\|$

Primal Optimization

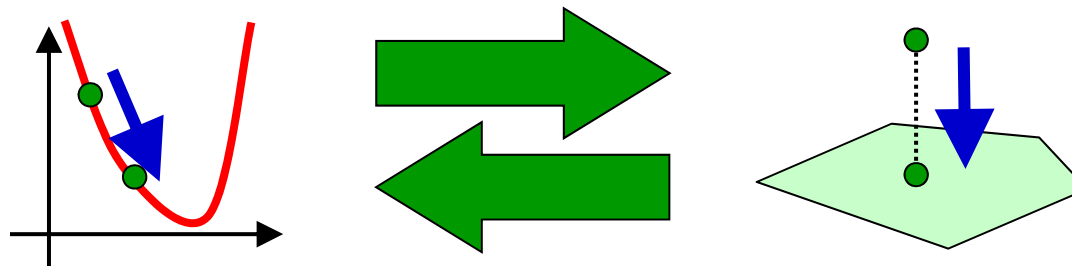
81

$$\min_{\alpha} \log \frac{1}{n'} \sum_{i'=1}^{n'} \exp \left(\sum_{k \geq k'} \alpha_{k,k'}^{\top} \mathbf{f}(x_{i'}^{(k)}, x_{i'}^{(k')}) \right)$$

$$- \frac{1}{n} \sum_{i=1}^n \sum_{k \geq k'} \alpha_{k,k'}^{\top} \mathbf{f}(x_i^{(k)}, x_i^{(k')}) + \eta \|\alpha\|^2$$

$$\text{subject to } \sum_{k \geq k'} \|\alpha_{k,k'}\| \leq C_{\gamma}$$

- Simple gradient-projection gives the global solution.
- Efficient when more samples than parameters.

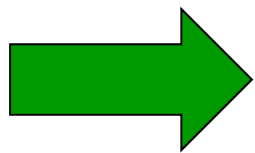


Dual Optimization

$$\min_{\beta} \sum_{i'=1}^{n'} \beta_{i'}^{\top} \log \beta_{i'} + \frac{1}{2\eta} \sum_{k \geq k'} \max(0, \|\mathbf{m}_{k,k'}\| - \gamma)^2$$

$$\text{subject to } \beta_1, \dots, \beta_{n'} \geq 0, \sum_{i'=1}^{n'} \beta_{i'} = 1$$

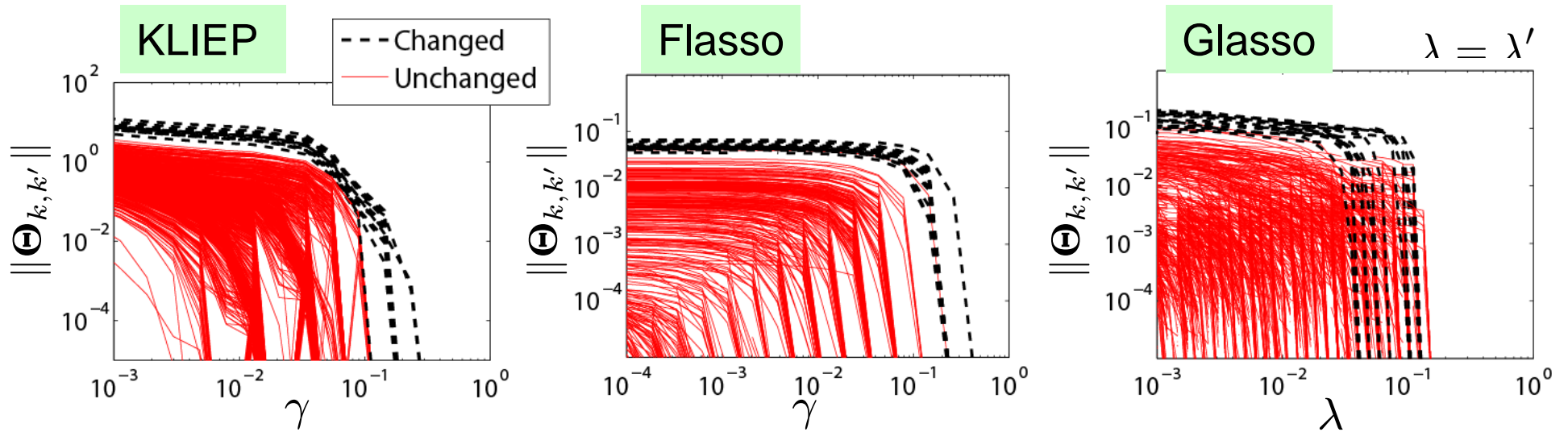
$$\mathbf{m}_{k,k'} = \frac{1}{n} \sum_{i=1}^n \mathbf{f}(\mathbf{x}_i^{(k)}, \mathbf{x}_i^{(k')}) - \frac{1}{n'} \sum_{i'=1}^{n'} \beta_{i'} \mathbf{f}(\mathbf{x}'_{i'}^{(k)}, \mathbf{x}'_{i'}^{(k')})$$



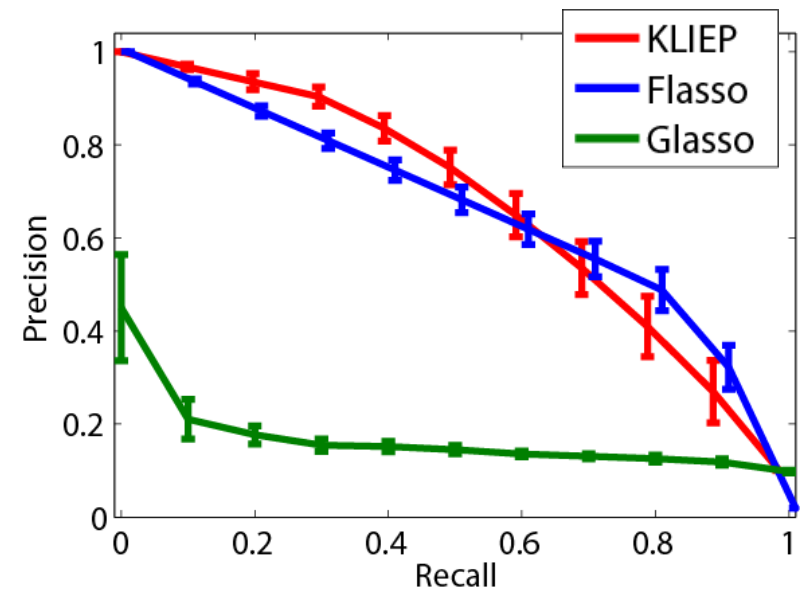
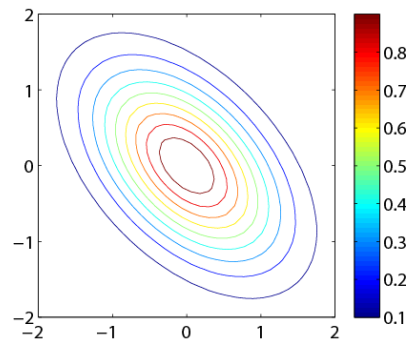
$$\alpha_{k,k'} = \max(0, \|\mathbf{m}_{k,k'}\| - \gamma) \frac{\mathbf{m}_{k,k'}}{\eta \|\mathbf{m}_{k,k'}\|}$$

- Simple gradient-projection gives the global solution.
- Efficient when more parameters than samples.

($d=40$, $n=n'=100$, Change in 15 Edges)



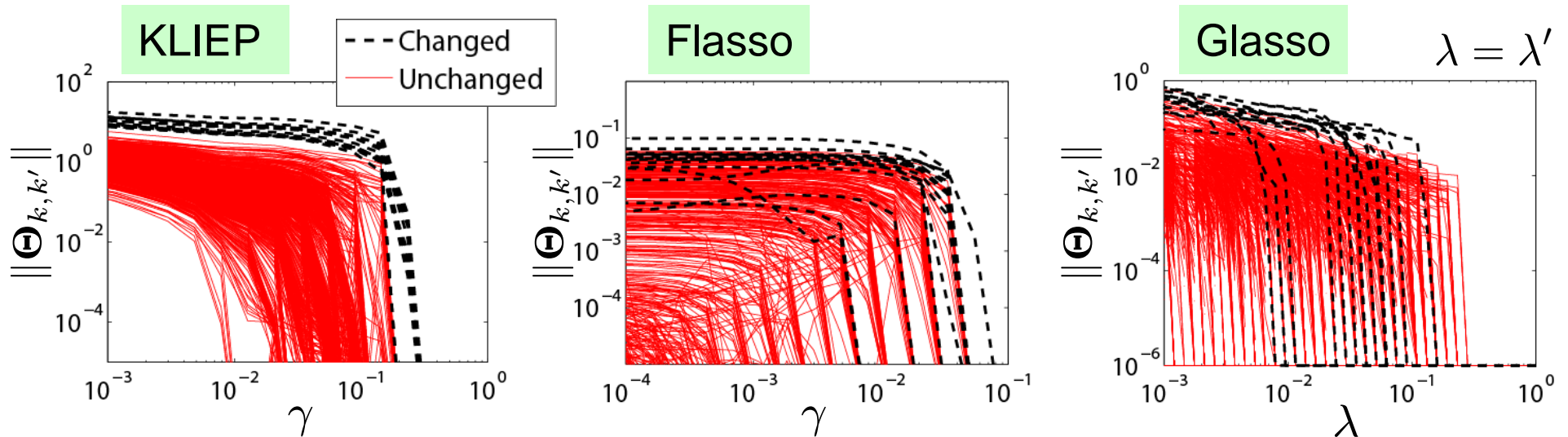
- All use the Gaussian model.
- KLIEP and Flasso work well.



Gaussian Data

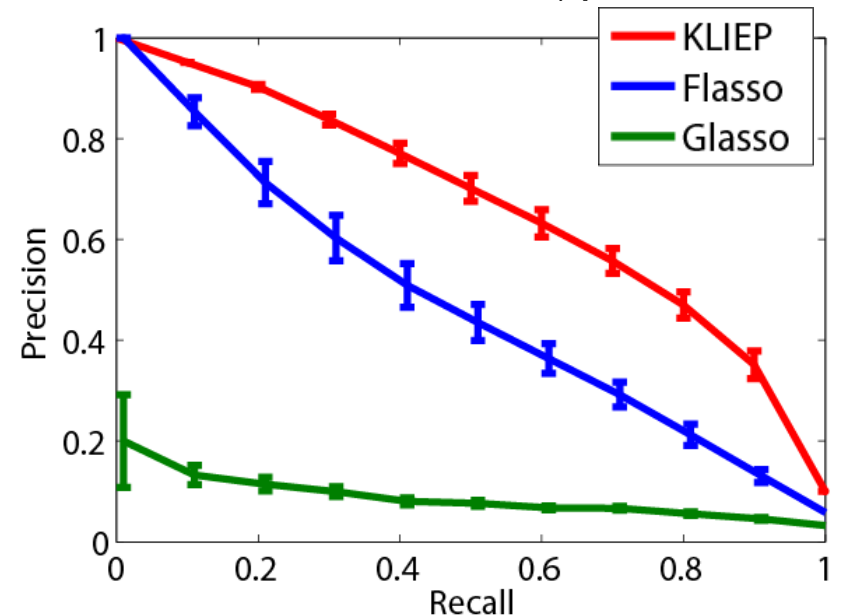
84

($d=40$, $n=n'=50$, Change in 15 Edges)



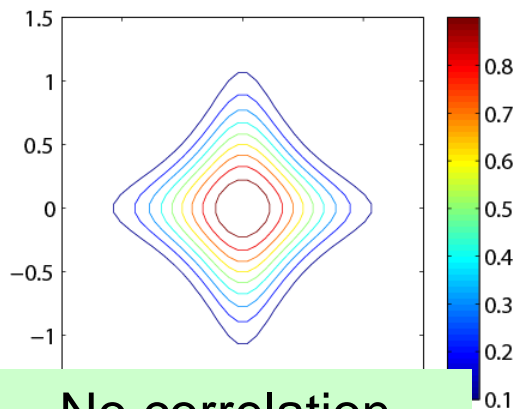
■ KLIEP works well even with small samples.

$$\alpha = \theta - \theta'$$



Non-Gaussian Data

($d=9$, $n=n'=5000$, Change in 7 Edges)

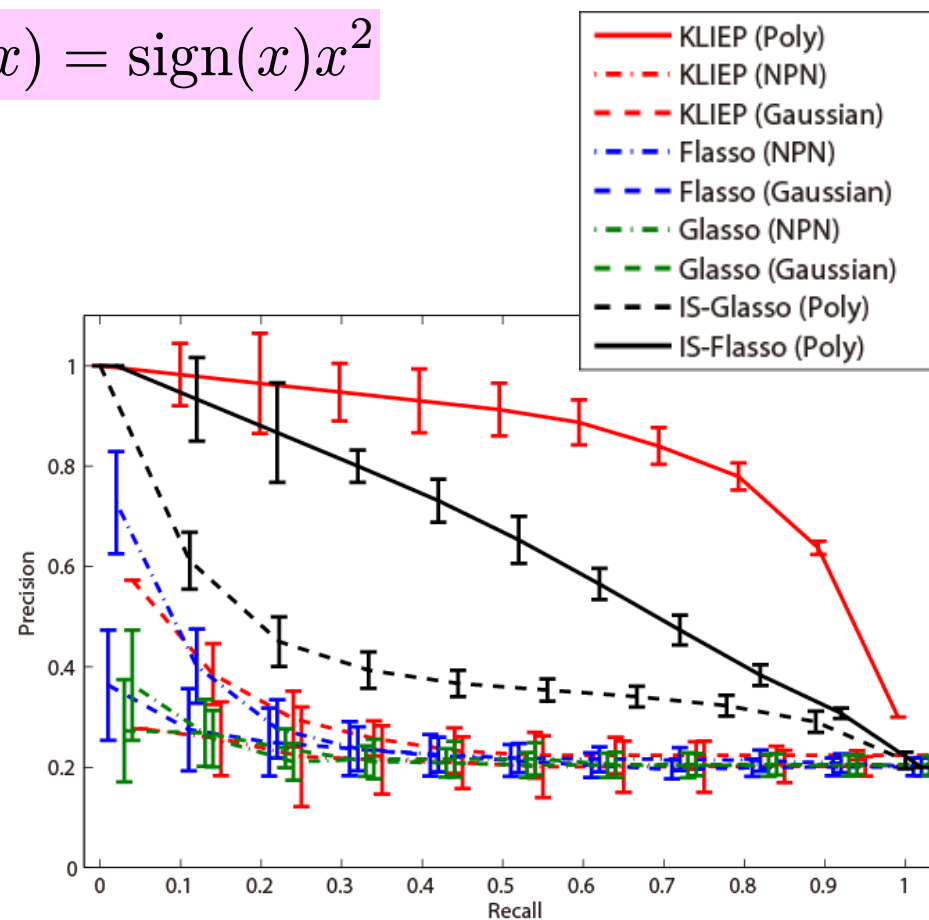
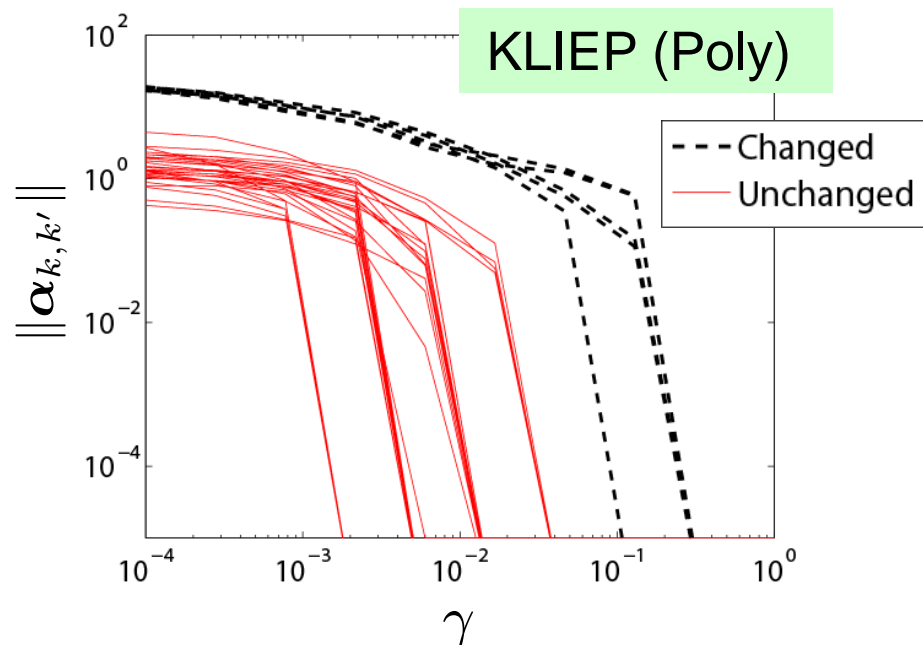


No correlation,
no nonparanormal

■ KLIEP (Poly) works well.

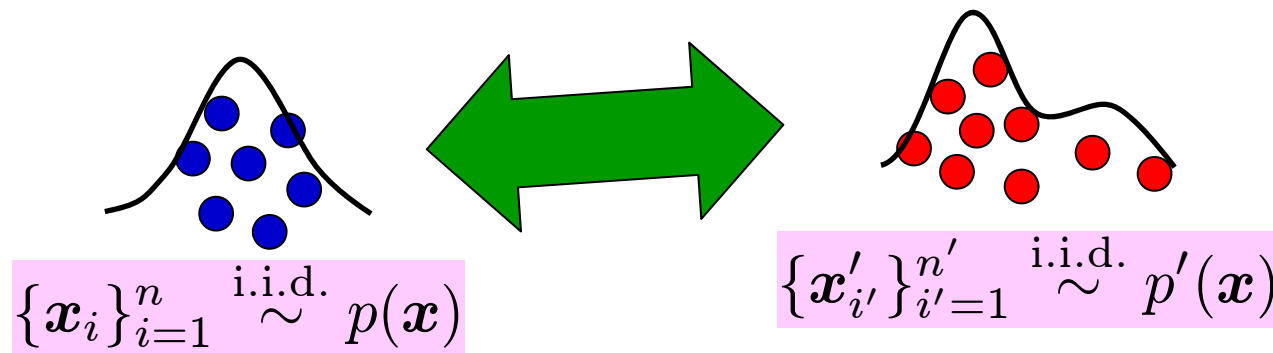
● Poly: $f(x, x') = [x^t, x^{t-1}x', \dots, x, x', 1]^T$

● NPN: $f(x) = \text{sign}(x)x^2$



Take-Home Messages

86



Don't solve
a more general
problem



■ Directly learn the change:

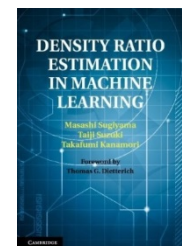
- Flexible and robust distributional change detection by direct density-ratio/density-difference estimation
- Interpretable and tractable structural change detection by group-sparse density-ratio estimation

■ Software: <http://sugiyama-www.cs.titech.ac.jp/~sugi/software/>

Schölkopf *et al.* (eds.),
*Empirical Inference, Festschrift
in Honor of Vladimir N. Vapnik*,
Springer, 2013



Sugiyama *et al.*,
*Density Ratio Estimation
in Machine Learning*,
Cambridge University Press, 2012



PE-Divergence Approximation

89

Kanamori *et al.* (NIPS2008, JMLR2009)

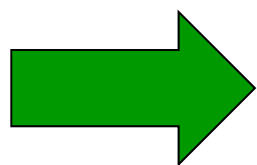
$$\text{PE}(p||p') = \int p'(\mathbf{x}) \left(r(\mathbf{x}) - 1 \right)^2 d\mathbf{x} = \int p(\mathbf{x}) r(\mathbf{x}) d\mathbf{x} - 1$$

- Directly approximate the density ratio by **least-squares**:

$$\hat{r} = \operatorname{argmin}_{\tilde{r}} \int p'(\mathbf{x}) \left(\tilde{r}(\mathbf{x}) - r(\mathbf{x}) \right)^2 d\mathbf{x}$$

$$r(\mathbf{x}) = \frac{p(\mathbf{x})}{p'(\mathbf{x})}$$

$$= \operatorname{argmin}_{\tilde{r}} \int p'(\mathbf{x}) \left(\tilde{r}(\mathbf{x}) \right)^2 d\mathbf{x} - 2 \int p(\mathbf{x}) r(\mathbf{x}) d\mathbf{x}$$



$$\text{PE}(p||p') \approx \int p(\mathbf{x}) \hat{r}(\mathbf{x}) d\mathbf{x} - 1$$

- Expectation is approximated by empirical average.



Contents

90

1. **Distributional change detection**
 - A) Problem setup and motivating examples
 - B) **Distance approximation**
 - I. Kullback-Leibler divergence
 - II. Pearson divergence
 - III. **Relative Pearson divergence**
 - IV. L^2 -distance
2. Structural change detection

rPE-Divergence Approximation 91

Yamada *et al.* (NIPS2011, NeCo2013)

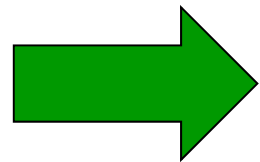
$$\text{rPE}(p||p') = \int p_\beta(\mathbf{x}) \left(\frac{p(\mathbf{x})}{p_\beta(\mathbf{x})} - 1 \right)^2 d\mathbf{x} = \int p(\mathbf{x}) \frac{p(\mathbf{x})}{p_\beta(\mathbf{x})} d\mathbf{x} - 1$$

$$p_\beta(\mathbf{x}) = \beta p(\mathbf{x}) + (1 - \beta)p'(\mathbf{x}) \quad 0 \leq \beta < 1$$

- Directly approximate the relative density ratio by LS:

$$\hat{r}_\beta = \underset{\tilde{r}}{\operatorname{argmin}} \int p_\beta(\mathbf{x}) \left(\tilde{r}(\mathbf{x}) - \frac{p(\mathbf{x})}{p_\beta(\mathbf{x})} \right)^2 d\mathbf{x}$$

$$= \underset{\tilde{r}}{\operatorname{argmin}} \int p_\beta(\mathbf{x}) (\tilde{r}(\mathbf{x}))^2 d\mathbf{x} - 2 \int p(\mathbf{x}) \tilde{r}(\mathbf{x}) d\mathbf{x}$$



$$\text{rPE}(p||p') \approx \int p(\mathbf{x}) \hat{r}_\beta(\mathbf{x}) d\mathbf{x} - 1$$

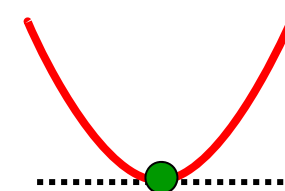
- Expectation is approximated by empirical average.

Solution for Linear Model

92

$$\hat{\alpha}_\beta = \operatorname{argmin}_\alpha \frac{\beta}{n'} \sum_{i'=1}^{n'} r_\alpha(\mathbf{x}'_{i'})^2 + \frac{1-\beta}{n} \sum_{i=1}^n r_\alpha(\mathbf{x}_i)^2 - \frac{2}{n} \sum_{i=1}^n r_\alpha(\mathbf{x}_i)$$

$$r_\alpha(\mathbf{x}) = \alpha^\top \phi(\mathbf{x})$$



- (Regularized) solution is given analytically:

$$\hat{\alpha}_\beta = \operatorname{argmin}_\alpha \left[\alpha^\top \hat{\mathbf{G}}_\beta \alpha - 2\hat{\mathbf{h}}^\top \alpha + \lambda \alpha^\top \alpha \right]$$

$$\hat{\mathbf{h}} = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i)$$

$$= (\hat{\mathbf{G}}_\beta + \lambda \mathbf{I})^{-1} \hat{\mathbf{h}}$$

$$\hat{\mathbf{G}}_\beta = \frac{\beta}{n'} \sum_{i'=1}^{n'} \phi(\mathbf{x}'_{i'}) \phi(\mathbf{x}'_{i'})^\top + \frac{1-\beta}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top$$

- Resulting rPE-divergence approximator:

$$\operatorname{rPE}(p \| p') \approx \frac{1}{n} \sum_{i=1}^n \hat{\alpha}_\beta^\top \phi(\mathbf{x}_i) - 1 = \hat{\mathbf{h}}^\top (\hat{\mathbf{G}}_\beta + \lambda \mathbf{I})^{-1} \hat{\mathbf{h}} - 1$$