

CHAPTER CONTENTS

Distributional Change Detection	469
KL Divergence	470
Pearson Divergence	470
L_2 -Distance	471
L_1 -Distance	474
Maximum Mean Discrepancy (MMD)	476
Energy Distance	477
Application to Change Detection in Time Series	477
Structural Change Detection	478
Sparse MLE	478
Sparse Density Ratio Estimation	482

The objective of *change detection* is to investigate whether change exists between two data sets $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$. In this chapter, two statistical approaches to change detection, *distributional change detection* and *structural change detection*, are explored. Below, $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ are assumed to be drawn independently from the probability distributions with densities $p(\mathbf{x})$ and $p'(\mathbf{x})$, respectively.

39.1 DISTRIBUTIONAL CHANGE DETECTION

Distributional change detection is aimed at identifying change in probability distributions behind $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$. This can be achieved by estimating a *distance* or a *divergence* between $p(\mathbf{x})$ and $p'(\mathbf{x})$. As explained in Section 14.2, a distance satisfies four conditions: *non-negativity*, *symmetry*, *identity*, and the *triangle inequality*. On the other hand, a divergence is a pseudodistance that still acts like a distance, but it may violate some of the above conditions. In this section, divergence and distance measures between probability distributions that are useful for change detection and their approximators from samples $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ are introduced.

39.1.1 KL DIVERGENCE

The most popular divergence measure in statistics and machine learning would be the KL divergence (see Section 14.2), because of its compatibility with MLE (see Chapter 12):

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

Due to the log function which is sharp near zero, the KL divergence may allow sensitive detection of small change in probability distributions. Another advantage of the KL divergence is that it is *invariant* under input metric change, i.e., the value of the KL divergence does not change even if \mathbf{x} is transformed to $\tilde{\mathbf{x}}$ by *any* mapping [5].

The KL divergence can be approximated from samples $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ by the *direct density ratio estimator* introduced in Section 38.3. More specifically, an estimator $\widehat{w}(\mathbf{x})$ of the density ratio function,

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

can be obtained by *KL density ratio estimation*. Then, the KL divergence can be immediately approximated as

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

However, due to the log function, estimation of the KL divergence is prone to be sensitive to *outliers*.

39.1.2 PEARSON DIVERGENCE

Ali-Silvey-Csiszár divergences [2, 35], which is also known as *f-divergences*, are generalization of the KL divergence defined as

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

where $f(t)$ is a *convex* function (see Fig. 8.3) such that $f(1) = 0$. It can be easily confirmed that, with $f(t) = t \log t$, the Ali-Silvey-Csiszár divergence is reduced to the KL divergence. Note that all Ali-Silvey-Csiszár divergences are invariant under input metric change.

The *Pearson divergence* [79] is a squared-loss variant of the KL divergence defined as the Ali-Silvey-Csiszár divergence with the squared function $f(t) = (t-1)^2$:

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

Since the Pearson divergence does not include the log function, it would be more robust against outliers. However, it still includes the density ratio function $p(\mathbf{x})/p'(\mathbf{x})$, which tends to be a sharp function and is possibly unbounded. Therefore, its accurate approximation is not straightforward in practice. As discussed in Section 33.2.2, this problem can be mitigated by considering the *relative density ratio* for $\beta \in [0, 1]$:

$$w_\beta(\mathbf{x}) = \frac{p(\mathbf{x})}{\beta p(\mathbf{x}) + (1 - \beta)p'(\mathbf{x})}.$$

The Pearson divergence extended using the relative density ratio is called the *relative Pearson divergence* [121]:

$$\begin{aligned} \text{rPE}(p\|p') &= \text{PE}(p\|\beta p + (1 - \beta)p') \\ &= \int (\beta p(\mathbf{x}) + (1 - \beta)p'(\mathbf{x})) (w_\beta(\mathbf{x}) - 1)^2 d\mathbf{x}. \end{aligned}$$

The relative Pearson divergence can be approximated from samples $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ by the *direct relative density ratio estimator* introduced in Section 33.2.4. More specifically, an estimator $\widehat{w}_\beta(\mathbf{x})$ of the relative density ratio function $w_\beta(\mathbf{x})$ can be obtained by *LS relative density ratio estimation*. Then, the relative Pearson divergence can be approximated as

$$\widehat{\text{rPE}} = \frac{1}{n} \sum_{i=1}^n \widehat{w}_\beta(\mathbf{x}_i) - 1,$$

which comes from the following expression of the relative Pearson divergence:

$$\text{rPE}(p\|p') = \int p(\mathbf{x}) w_\beta(\mathbf{x}) d\mathbf{x} - 1.$$

The tuning parameter $\beta \in [0, 1]$ controls the trade-off between sensitivity and robustness of the divergence measure, which should be appropriately chosen in practice.

39.1.3 L_2 -DISTANCE

The L_2 -distance is another standard distance measure between probability distributions:

$$L_2(p, p') = \int f(\mathbf{x})^2 d\mathbf{x},$$

where

$$f(\mathbf{x}) = p(\mathbf{x}) - p'(\mathbf{x}).$$

The L_2 -distance is a proper distance measure, and thus it is symmetric and satisfies the triangle inequality unlike the KL divergence and the (relative) Pearson divergence. Furthermore, the density difference $f(\mathbf{x})$ is always bounded as long as each

density is bounded. Therefore, the L_2 -distance is stable, without the need of tuning any control parameter such as β in the relative Pearson divergence.

The L_2 -distance can be approximated from samples $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ by *LS density difference estimation* [103] that directly estimates the density difference function $f(\mathbf{x})$ without estimating $p(\mathbf{x})$ and $p'(\mathbf{x})$ (see Section 37.4 for direct estimation of $p(\mathbf{x}, y) - p(\mathbf{x})p(y)$). More specifically, let us consider the following Gaussian density difference model:

$$f_{\alpha}(\mathbf{x}) = \sum_{j=1}^{n+n'} \alpha_j \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_j\|^2}{2h^2}\right),$$

where $h > 0$ denotes the Gaussian bandwidth and

$$(\mathbf{c}_1, \dots, \mathbf{c}_n, \mathbf{c}_{n+1}, \dots, \mathbf{c}_{n+n'}) = (\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}'_1, \dots, \mathbf{x}'_{n'})$$

are the Gaussian centers. The parameters $\alpha = (\alpha_1, \dots, \alpha_{n+n'})^\top$ are estimated by LS fitting to the true density difference function $f(\mathbf{x})$:

$$\min_{\alpha} \int (f_{\alpha}(\mathbf{x}) - f(\mathbf{x}))^2 d\mathbf{x}.$$

Its empirical criterion where an irrelevant constant is ignored and the expectation is approximated by the sample average is given by

$$\min_{\alpha} [\alpha^\top U \alpha - 2\alpha^\top \hat{\mathbf{v}} + \lambda \|\alpha\|^2],$$

where the ℓ_2 -regularizer $\lambda \|\alpha\|^2$ is included. U is the $(n + n') \times (n + n')$ matrix with the (j, j') th element defined by

$$\begin{aligned} U_{j,j'} &= \int \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_j\|^2}{2h^2}\right) \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_{j'}\|^2}{2h^2}\right) d\mathbf{x} \\ &= (\pi h^2)^{d/2} \exp\left(-\frac{\|\mathbf{c}_j - \mathbf{c}_{j'}\|^2}{4h^2}\right), \end{aligned}$$

where d denotes the dimensionality of \mathbf{x} . $\hat{\mathbf{v}}$ is the $(n + n')$ -dimensional vector with the j th element defined by

$$\hat{v}_j = \frac{1}{n} \sum_{i=1}^n \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{c}_j\|^2}{2h^2}\right) - \frac{1}{n'} \sum_{i'=1}^{n'} \exp\left(-\frac{\|\mathbf{x}'_{i'} - \mathbf{c}_j\|^2}{2h^2}\right).$$

This is a convex optimization problem, and the global optimal solution $\hat{\alpha}$ can be obtained *analytically* as

$$\hat{\alpha} = (U + \lambda I)^{-1} \hat{\mathbf{v}}.$$

```

n=200; x=randn(n,1); y=randn(n,1)+1;
hhs=2*[0.5 1 3].^2; ls=10.^[-2 -1 0]; m=5;
x2=x.^2; xx= repmat(x2,1,n)+ repmat(x2',n,1)-2*x*x';
y2=y.^2; yx= repmat(y2,1,n)+ repmat(x2',n,1)-2*y*x';
u=mod(randperm(n),m)+1; v=mod(randperm(n),m)+1;

for hk=1:length(hhs)
    hh=hhs(hk); k=exp(-xx/hh); r=exp(-yx/hh);
    U=(pi*hh/2)^(1/2)*exp(-xx/(2*hh));
    for i=1:m
        vh=mean(k(u~=i,:))'-mean(r(v~=i,:))';
        z=mean(k(u==i,:))-mean(r(v==i,:));
        for lk=1:length(ls)
            l=ls(lk); a=(U+l*eye(n))\vh; g(hk,lk,i)=a'*U*a-2*z*a;
        end, end, end
    [gl,ggl]=min(mean(g,3),[],2); [ghl,gghl]=min(gl);
    L=ls(ggl(gghl)); HH=hhs(gghl);
    k=exp(-xx/HH); r=exp(-yx/HH); vh=mean(k)'-mean(r)';
    U=(pi*HH/2)^(1/2)*exp(-xx/(2*HH));
    a=(U+L*eye(n))\vh; s=[k;r]*a; L2=a'*vh;
    figure(1); clf; hold on; plot([x;y],s,'rx');

```

FIGURE 39.1

MATLAB code for LS density difference estimation.

The Gaussian width h and the regularization parameter λ may be optimized by cross validation with respect to the squared error criterion.

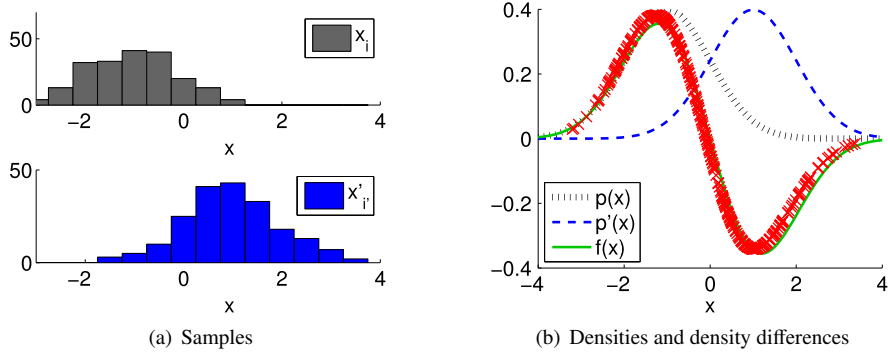
Finally, the L_2 -distance can be approximated by

$$\widehat{L}^2 = \frac{1}{n} \sum_{i=1}^n f_{\widehat{\alpha}}(\mathbf{x}_i) - \frac{1}{n'} \sum_{i'=1}^{n'} f_{\widehat{\alpha}}(\mathbf{x}'_{i'}) = \widehat{\mathbf{v}}^\top \widehat{\alpha}, \quad (39.1)$$

which comes from the following expression of the L_2 -distance:

$$L_2(p, p') = \int (p(\mathbf{x}) - p'(\mathbf{x})) f(\mathbf{x}) d\mathbf{x}.$$

A MATLAB code for LS density difference estimation is provided in [Fig. 39.1](#), and its behavior is illustrated in [Fig. 39.2](#). This shows that the density difference function can be accurately estimated.

**FIGURE 39.2**

Example of LS density difference estimation. \times 's in the right plot show estimated density difference values at $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$.

39.1.4 L_1 -DISTANCE

The L_2 -distance can be generalized to the class of L_s -distances for $s > 0$:

$$L_t(p, p') = \int |p(\mathbf{x}) - p'(\mathbf{x})|^s d\mathbf{x}.$$

Among the class of L_s -distances, the L_1 -distance is also a member of the class of the Ali-Silvey-Csiszár divergences with the absolute function $f(t) = |t - 1|$ (see Section 39.1.2):

$$L_1(p, p') = \int |p(\mathbf{x}) - p'(\mathbf{x})| d\mathbf{x} = \int p'(\mathbf{x}) \left| \frac{p(\mathbf{x})}{p'(\mathbf{x})} - 1 \right| d\mathbf{x}.$$

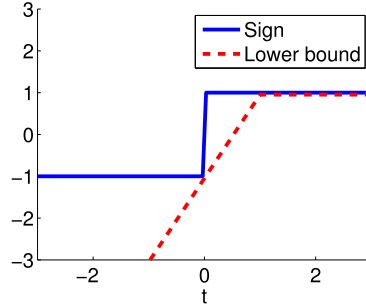
This implies that the L_1 -distance is also invariant under input metric change.

Approximation of the L_1 -distance from samples $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ can actually be performed by the (weighted) *support vector machine* introduced in Chapter 27 [75]. More specifically, since $|t| = \text{sign}(t)t$, the L_1 -distance can be expressed as

$$L_1(p, p') = \int \text{sign}(p(\mathbf{x}) - p'(\mathbf{x})) (p(\mathbf{x}) - p'(\mathbf{x})) d\mathbf{x},$$

where $\text{sign}(t)$ denotes the sign function:

$$\text{sign}(t) = \begin{cases} 1 & (t > 0), \\ 0 & (t = 0), \\ -1 & (t < 0). \end{cases}$$

**FIGURE 39.3**

Lower bound of $\text{sign}(t)$ by $-2 \max(0, 1 - t) + 1$.

For a density difference model $f_\alpha(\mathbf{x})$ with parameter α , the L_1 -distance $L_1(p, p')$ can be lower-bounded as

$$\begin{aligned} L_1(p, p') &\geq \max_{\alpha} \int \text{sign}(f_\alpha(\mathbf{x})) (p(\mathbf{x}) - p'(\mathbf{x})) d\mathbf{x} \\ &= \max_{\alpha} \left[\int \text{sign}(f_\alpha(\mathbf{x})) p(\mathbf{x}) d\mathbf{x} + \int \text{sign}(-f_\alpha(\mathbf{x})) p'(\mathbf{x}) d\mathbf{x} \right], \end{aligned}$$

where the last term is due to $\text{sign}(t) = -\text{sign}(-t)$. As plotted in Fig. 39.3, the sign function can be lower-bounded by

$$-2 \max(0, 1 - t) + 1 = \begin{cases} 1 & (t > 1), \\ 2t - 1 & (t \leq 1). \end{cases}$$

Based on this, the L_1 -distance $L_1(p, p')$ can be further lower-bounded as

$$\begin{aligned} L_1(p, p') &\geq 2 - 2 \min_{\alpha} \left[\int \max(0, 1 - f_\alpha(\mathbf{x})) p(\mathbf{x}) d\mathbf{x} \right. \\ &\quad \left. + \int \max(0, 1 + f_\alpha(\mathbf{x})) p'(\mathbf{x}) d\mathbf{x} \right]. \end{aligned}$$

Let us employ a *linear-in-parameter* density difference model:

$$f_\alpha(\mathbf{x}) = \sum_{j=1}^b \alpha_j \psi_j(\mathbf{x}) = \alpha^\top \boldsymbol{\psi}(\mathbf{x}).$$

Then the empirical version of the above maximization problem (without irrelevant multiplicative and additive constants) is given by

$$\min_{\alpha} \left[\frac{1}{n} \sum_{i=1}^n \max(0, 1 - \alpha^\top \boldsymbol{\psi}(\mathbf{x}_i)) + \frac{1}{n'} \sum_{i'=1}^{n'} \max(0, 1 + \alpha^\top \boldsymbol{\psi}(\mathbf{x}'_{i'})) \right].$$

Let us assign *class labels* $y_i = +1$ to \mathbf{x}_i for $i = 1, \dots, n$ and $y_{i'} = -1$ to $\mathbf{x}'_{i'}$ for $i' = 1, \dots, n'$. If $n = n'$, the above optimization problem agrees with *hinge loss minimization* (see Section 27.6):

$$\min_{\alpha} \left[\sum_{i=1}^n \max(0, 1 - y_i \alpha^\top \psi(\mathbf{x}_i)) + \sum_{i'=1}^{n'} \max(0, 1 - y_{i'} \alpha^\top \psi(\mathbf{x}'_{i'})) \right].$$

If $n \neq n'$, L_1 -distance approximation corresponds to weighted hinge loss minimization with weight $1/n$ for $\{\mathbf{x}_i\}_{i=1}^n$ and $1/n'$ for $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$.

The above formulation shows that the support vector machine is actually approximating the sign of the density difference. More specifically, let $p_+(\mathbf{x})$ and $p_-(\mathbf{x})$ be the probability density functions of samples in the positive class and negative class, respectively. Then, the support vector machine approximates

$$\text{sign}(p_+(\mathbf{x}) - p_-(\mathbf{x})),$$

which is the optimal decision function. Thus, support vector classification can be interpreted as directly approximating the optimal decision function without estimating the densities $p_+(\mathbf{x})$ and $p_-(\mathbf{x})$.

39.1.5 MAXIMUM MEAN DISCREPANCY (MMD)

MMD [17] measures the distance between embeddings of probability distributions in a *reproducing kernel Hilbert space* [9].

More specifically, the MMD between p and p' is defined as

$$\begin{aligned} \text{MMD}(p, p') &= \mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim p} [K(\mathbf{x}, \tilde{\mathbf{x}})] + \mathbb{E}_{\mathbf{x}', \tilde{\mathbf{x}}' \sim p'} [K(\mathbf{x}', \tilde{\mathbf{x}}')] \\ &\quad - 2\mathbb{E}_{\mathbf{x} \sim p, \mathbf{x}' \sim p'} [K(\mathbf{x}, \mathbf{x}')], \end{aligned}$$

where $K(\mathbf{x}, \mathbf{x}')$ is a reproducing kernel, $\mathbb{E}_{\mathbf{x} \sim p}$ denotes the expectation with respect to \mathbf{x} following density p . $\text{MMD}(p, p')$ is always non-negative, and $\text{MMD}(p, p') = 0$ if and only if $p = p'$ when $K(\mathbf{x}, \mathbf{x}')$ is a *characteristic kernel* [45] such as the Gaussian kernel.

An advantage of MMD is that it can be directly approximated using samples as

$$\frac{1}{n^2} \sum_{i, \tilde{i}=1}^n K(\mathbf{x}_i, \mathbf{x}_{\tilde{i}}) + \frac{1}{n'^2} \sum_{i', \tilde{i}'=1}^{n'} K(\mathbf{x}'_{i'}, \mathbf{x}'_{\tilde{i}'}) - \frac{2}{nn'} \sum_{i=1}^n \sum_{i'=1}^{n'} K(\mathbf{x}_i, \mathbf{x}'_{i'}).$$

Thus, no estimation is involved when approximating MMD from samples. However, it is not clear how to choose kernel functions in practice. Using the Gaussian kernel with bandwidth set at the median distance between samples is a popular heuristic [49], but this does not always work well in practice [50].

39.1.6 ENERGY DISTANCE

Another useful distance measure is the *energy distance* [106] introduced in Section 33.3.2:

$$D_E(p, p') = \int_{\mathbb{R}^d} \|\varphi_p(\mathbf{t}) - \varphi_{p'}(\mathbf{t})\|^2 \left(\frac{\pi^{\frac{d+1}{2}}}{\Gamma(\frac{d+1}{2})} \|\mathbf{t}\|^{d+1} \right)^{-1} d\mathbf{t},$$

where $\|\cdot\|$ denotes the Euclidean norm, φ_p denotes the *characteristic function* (see Section 2.4.3) of p , $\Gamma(\cdot)$ is the *gamma function* (see Section 4.3), and d denotes the dimensionality of \mathbf{x} .

An important property of the energy distance is that it can be expressed as

$$D_E(p, p') = 2\mathbb{E}_{\mathbf{x} \sim p, \mathbf{x}' \sim p'} \|\mathbf{x} - \mathbf{x}'\| - \mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim p} \|\mathbf{x} - \tilde{\mathbf{x}}\| - \mathbb{E}_{\mathbf{x}', \tilde{\mathbf{x}}' \sim p'} \|\mathbf{x}' - \tilde{\mathbf{x}}'\|,$$

where $\mathbb{E}_{\mathbf{x} \sim p}$ denotes the expectation with respect to \mathbf{x} following density p . This can be directly approximated using samples as

$$\frac{2}{nn'} \sum_{i=1}^n \sum_{i'=1}^{n'} \|\mathbf{x}_i - \mathbf{x}'_{i'}\| - \frac{1}{n^2} \sum_{i, \tilde{i}=1}^n \|\mathbf{x}_i - \mathbf{x}_{\tilde{i}}\| - \frac{1}{n'^2} \sum_{i', \tilde{i}'=1}^{n'} \|\mathbf{x}'_{i'} - \mathbf{x}'_{\tilde{i}'}\|.$$

Thus, no estimation and no tuning parameters are involved when approximating the energy distance from samples, which is a useful properties in practice.

Actually, it was shown [91] that the energy distance is a special case of MMD. Indeed, MMD with kernel function defined as

$$K(\mathbf{x}, \mathbf{x}') = -\|\mathbf{x} - \mathbf{x}'\| + \|\mathbf{x}\| + \|\mathbf{x}'\|$$

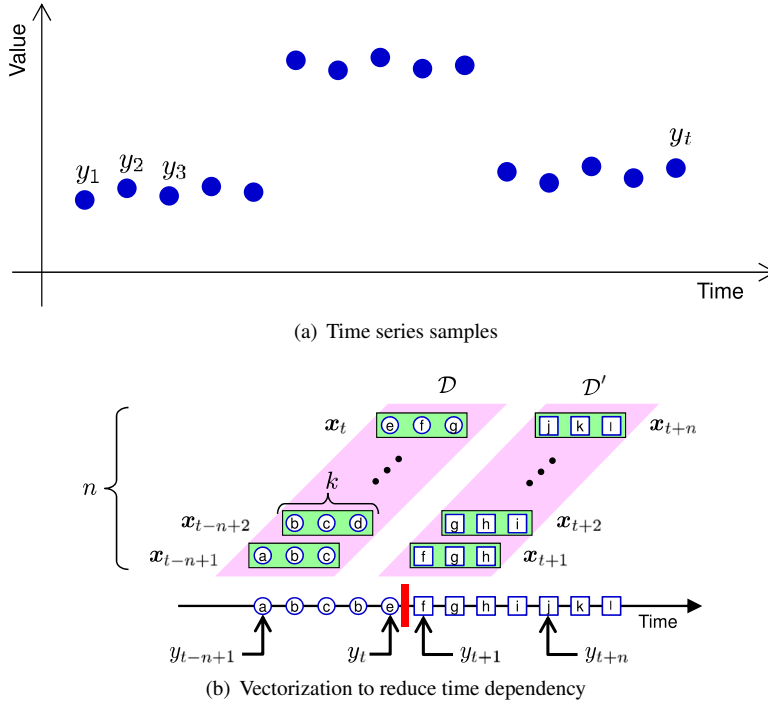
agrees with the energy distance.

39.1.7 APPLICATION TO CHANGE DETECTION IN TIME SERIES

Let us consider the problem of *change detection in time series* (Fig. 39.4(a)). More specifically, given time series samples $\{y_i\}_{i=1}^N$, the objective is to identify whether change in probability distributions exists between y_t and y_{t+1} for some t . This problem can be tackled by estimating a distance (or a divergence) between the probability distributions of $\{y_i\}_{i=t-n+1}^t$ and $\{y_i\}_{i=t+1}^{t+n}$.

A challenge in change detection in time series is that samples $\{y_i\}_{i=1}^N$ are often dependent over time, which violates the presumption in this chapter. A practical approach to mitigate this problem is to *vectorize* data [60], as illustrated in Fig. 39.4(b). That is, instead of handling time series sample y_i as it is, its vectorization with k consecutive samples $\mathbf{x}_i = (y_i, \dots, y_{i+k-1})^\top$ is considered, and a distance (or a divergence) is estimated from $\mathcal{D} = \{\mathbf{x}_i\}_{i=t-n+1}^t$ and $\mathcal{D}' = \{\mathbf{x}_i\}_{i=t+1}^{t+n}$.

A MATLAB code for change detection in time series based on the energy distance is provided in Fig. 39.5, and its behavior is illustrated in Fig. 39.6. This shows that the energy distance well captures the distributional change in time series.

**FIGURE 39.4**

Change detection in time series.

39.2 STRUCTURAL CHANGE DETECTION

Distributional change detection introduced in the previous section focused on investigating whether change exists in probability distributions. The aim of *structural change detection* introduced in this section is to analyze change in the *dependency structure* between elements of d -dimensional variable $\mathbf{x} = (x^{(1)}, \dots, x^{(d)})^\top$.

39.2.1 SPARSE MLE

Let us consider a *Gaussian Markov network*, which is a d -dimensional Gaussian model with expectation zero (Section 6.2):

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

where not the variance-covariance matrix, but its inverse called the *precision matrix* is parameterized by Θ . If Θ is regarded as an *adjacency matrix*, the Gaussian Markov network can be visualized as a *graph* (see Fig. 39.7). An advantage of this precision-based parameterization is that the connectivity governs conditional independence.

```

N=300; k=5; n=10; m=N-k+1; E=nan(1,N);
y=zeros(1,N); y(101:200)=3; y=y+randn(1,N);
%y=sin([1:N]/2); y(101:200)=sin([101:200]);
x=toeplitz(y); x=x(1:k,1:m); x2=sum(x.^2);
D=sqrt(repmat(x2',1,m)+repmat(x2,m,1)-2*x'*x);
for t=n:N-n-k+1
    a=[t-n+1:t]; b=[t+1:t+n];
    E(t)=2*mean(mean(D(a,b)))-mean(mean(D(a,a))) ...
        -mean(mean(D(b,b)));
end

figure(1); clf; hold on; plot(y,'b-'); plot(E,'r--');
legend('Time series','Energy distance')

```

FIGURE 39.5

MATLAB code for change detection in time series based on the energy distance.

For example, in the Gaussian Markov network illustrated in the left-hand side of Fig. 39.7, $x^{(1)}$ and $x^{(2)}$ are connected via $x^{(3)}$. This means that $x^{(1)}$ and $x^{(2)}$ are conditionally independent given $x^{(3)}$.

Suppose that $\{x_i\}_{i=1}^n$ and $\{x'_{i'}\}_{i'=1}^{n'}$ are drawn independently from the Gaussian Markov networks with precision matrices Θ and Θ' , respectively. Then analyzing $\Theta - \Theta'$ allows us to identify the change in Markov network structure (see Fig. 39.7 again).

A sparse estimate of Θ may be obtained by MLE with the ℓ_1 -constraint (see Chapter 24):

$$\max_{\Theta} \sum_{i=1}^n \log q(x_i; \Theta) \quad \text{subject to } \|\Theta\|_1 \leq R^2,$$

where $R \geq 0$ is the radius of the ℓ_1 -ball. This method is also referred to as the *graphical lasso* [44].

The derivative of $\log q(x; \Theta)$ with respect to Θ is given by

$$\frac{\partial \log q(x; \Theta)}{\partial \Theta} = \frac{1}{2} \Theta^{-1} - \frac{1}{2} x x^\top,$$

where the following formulas are used for its derivation:

$$\frac{\partial \log \det(\Theta)}{\partial \Theta} = \Theta^{-1} \quad \text{and} \quad \frac{\partial x^\top \Theta x}{\partial \Theta} = x x^\top.$$

A MATLAB code of a gradient-projection algorithm of ℓ_1 -constraint MLE for Gaussian Markov networks is given in Fig. 39.8, where projection onto the ℓ_1 -ball is computed by the method developed in [39].

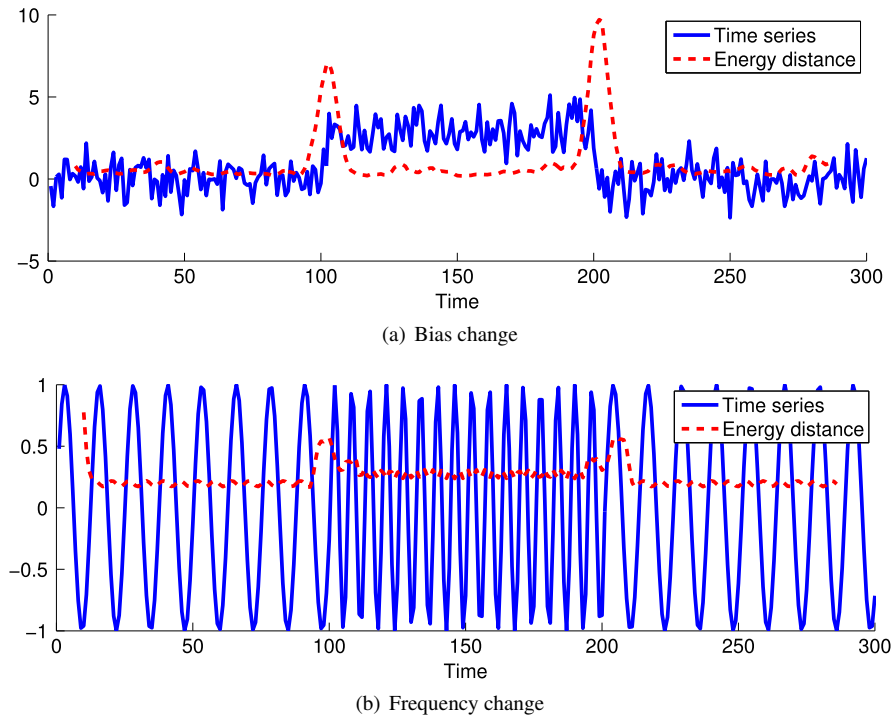


FIGURE 39.6

Examples of change detection in time series based on the energy distance.

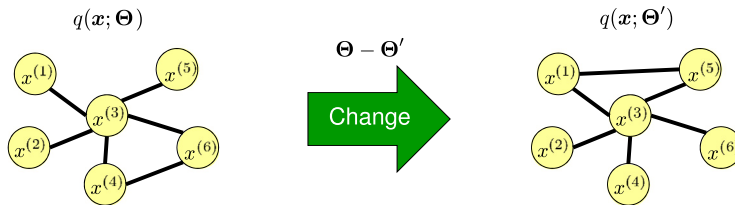


FIGURE 39.7

Structural change in Gaussian Markov networks.

For the true precision matrices

$$\Theta = \begin{pmatrix} 2 & 0 & 1 \\ 0 & 2 & 0 \\ 1 & 0 & 2 \end{pmatrix} \quad \text{and} \quad \Theta' = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{pmatrix},$$

```

TT=[2 0 1; 0 2 0; 1 0 2];
%TT=[2 0 0; 0 2 0; 0 0 2];
%TT=[2 1 0; 1 2 1; 0 1 2];
%TT=[2 0 1; 0 2 1; 1 1 2];
d=3; n=50; x=TT^(-1/2)*randn(d,n); S=x*x'/n;
T0=eye(d); C=5; e=0.1;
for o=1:100000
    T=T0+e*(inv(T0)-S);
    T(:)=L1BallProjection(T(:),C);
    if norm(T-T0)<0.00000001, break, end
    T0=T;
end
T, TT

```

```

function w=L1BallProjection(x,C)

u=sort(abs(x),'descend'); s=cumsum(u);
r=find(u>(s-C)./(1:length(u))',1,'last');
w=sign(x).*max(0,abs(x)-max(0,(s(r)-C)/r));

```

FIGURE 39.8

MATLAB code of a gradient-projection algorithm of ℓ_1 -constraint MLE for Gaussian Markov networks. The bottom function should be saved as “L1BallProjection.m.”

sparse MLE gives

$$\widehat{\Theta} = \begin{pmatrix} 1.382 & 0 & 0.201 \\ 0 & 1.788 & 0 \\ 0.201 & 0 & 1.428 \end{pmatrix} \quad \text{and} \quad \widehat{\Theta}' = \begin{pmatrix} 1.617 & 0 & 0 \\ 0 & 1.711 & 0 \\ 0 & 0 & 1.672 \end{pmatrix}.$$

Thus, the true sparsity patterns of Θ and Θ' (in off-diagonal elements) can be successfully recovered. Since

$$\Theta - \Theta' = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \widehat{\Theta} - \widehat{\Theta}' = \begin{pmatrix} -0.235 & 0 & 0.201 \\ 0 & 0.077 & 0 \\ 0.201 & 0 & -0.244 \end{pmatrix},$$

change in sparsity patterns (in off-diagonal elements) can be correctly identified.

On the other hand, when the true precision matrices are

$$\Theta = \begin{pmatrix} 2 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 2 \end{pmatrix} \quad \text{and} \quad \Theta' = \begin{pmatrix} 2 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix},$$

sparse MLE gives

$$\hat{\Theta} = \begin{pmatrix} 1.303 & 0.348 & 0 \\ 0.348 & 1.157 & 0.240 \\ 0 & 0.240 & 1.365 \end{pmatrix} \quad \text{and} \quad \hat{\Theta}' = \begin{pmatrix} 1.343 & 0 & 0.297 \\ 0 & 1.435 & 0.236 \\ 0.297 & 0.236 & 1.156 \end{pmatrix}.$$

Thus, the true sparsity patterns of Θ and Θ' can still be successfully recovered. However, since

$$\Theta - \Theta' = \begin{pmatrix} 0 & 1 & -1 \\ 1 & 0 & 0 \\ -1 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \hat{\Theta} - \hat{\Theta}' = \begin{pmatrix} -0.040 & 0.348 & -0.297 \\ 0.348 & -0.278 & 0.004 \\ -0.297 & 0.004 & 0.209 \end{pmatrix},$$

change in sparsity patterns was not correctly identified (although 0.004 is reasonably close to zero). This shows that, when a nonzero unchanged edge exists, say $\Theta_{k,k'} = \Theta'_{k,k'} > 0$ for some k and k' , it is difficult to identify this unchanged edge because $\hat{\Theta}_{k,k'} \approx \hat{\Theta}'_{k,k'}$ does not necessarily hold by separate sparse MLE from $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$.

39.2.2 SPARSE DENSITY RATIO ESTIMATION

As illustrated above, sparse MLE can perform poorly in structural change detection. Another limitation of sparse MLE is the Gaussian assumption. A Gaussian Markov network can be extended to a non-Gaussian model as

$$q(\mathbf{x}; \theta) = \frac{\bar{q}(\mathbf{x}; \theta)}{\int \bar{q}(\mathbf{x}; \theta) d\mathbf{x}},$$

where, for a *feature vector* $\mathbf{f}(x, x')$,

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

This model is reduced to the Gaussian Markov network if

$$\mathbf{f}(x, x') = -\frac{1}{2}xx',$$

while higher-order correlations can be captured by considering higher-order terms in the feature vector. However, applying sparse MLE to non-Gaussian Markov networks is not straightforward in practice because the normalization term $\int \bar{q}(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}$ is often computationally intractable.

To cope with these limitations, let us handle the change in parameters, $\boldsymbol{\theta}_{k,k'} - \boldsymbol{\theta}'_{k,k'}$, directly via the following density ratio function:

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

Based on this expression, let us consider the following density ratio model:

$$r(\mathbf{x}; \boldsymbol{\alpha}) = \frac{\exp \left(\sum_{k \geq k'} \boldsymbol{\alpha}_{k,k'}^\top \mathbf{f}(x^{(k)}, x^{(k')}) \right)}{\int p'(\mathbf{x}) \exp \left(\sum_{k \geq k'} \boldsymbol{\alpha}_{k,k'}^\top \mathbf{f}(x^{(k)}, x^{(k')}) \right) d\mathbf{x}}, \quad (39.2)$$

where $\boldsymbol{\alpha}_{k,k'}$ is the difference of parameters:

$$\boldsymbol{\alpha}_{k,k'} = \boldsymbol{\theta}_{k,k'} - \boldsymbol{\theta}'_{k,k'}.$$

$p'(\mathbf{x})$ in the denominator of Eq. (39.2) comes from the fact that $r(\mathbf{x}; \boldsymbol{\alpha})$ approximates $p(\mathbf{x})/p'(\mathbf{x})$ and thus the normalization constraint,

$$\int r(\mathbf{x}; \boldsymbol{\alpha}) p'(\mathbf{x}) d\mathbf{x} = 1,$$

is imposed.

Let us learn the parameters $\{\boldsymbol{\alpha}_{k,k'}\}_{k \geq k'}$ by a *group-sparse* variant (see Section 24.4.4) of *KL density ratio estimation* explained in Section 38.3 [69]:

$$\begin{aligned} \min_{\{\boldsymbol{\alpha}_{k,k'}\}_{k \geq k'}} \quad & \log \frac{1}{n'} \sum_{i'=1}^{n'} \exp \left(\sum_{k \geq k'} \boldsymbol{\alpha}_{k,k'}^\top \mathbf{f}(x_{i'}^{(k)}, x_{i'}^{(k')}) \right) \\ & - \frac{1}{n} \sum_{i=1}^n \sum_{k \geq k'} \boldsymbol{\alpha}_{k,k'}^\top \mathbf{f}(x_i^{(k)}, x_i^{(k')}) \\ \text{subject to} \quad & \sum_{k \geq k'} \|\boldsymbol{\alpha}_{k,k'}\| \leq R^2, \end{aligned}$$

where $R \geq 0$ controls the sparseness of the solution.

A MATLAB code of a gradient-projection algorithm of sparse KL density ratio estimation for Gaussian Markov networks is given in Fig. 39.9. For the true precision matrices

$$\boldsymbol{\Theta} - \boldsymbol{\Theta}' = \begin{pmatrix} 2 & 0 & 1 \\ 0 & 2 & 0 \\ 1 & 0 & 2 \end{pmatrix} - \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix},$$

```

Tp=[2 0 1; 0 2 0; 1 0 2]; Tq=[2 0 0; 0 2 0; 0 0 2];
Tp=[2 1 0; 1 2 1; 0 1 2]; Tq=[2 0 1; 0 2 1; 1 1 2];
d=3; n=50; xp=Tp^(-1/2)*randn(d,n); Sp=xp*xp'/n;
xq=Tq^(-1/2)*randn(d,n); A0=eye(d); C=1; e=0.1;
for o=1:1000000
    U=exp(sum((A0*xq).*xq));
    A=A0-e*((repmat(U,[d 1]).*xq)*xq'/sum(U)-Sp);
    A(:)=L1BallProjection(A(:),C);
    if norm(A-A0)<0.00000001, break, end
    A0=A;
end
-2*A, Tp-Tq

```

FIGURE 39.9

MATLAB code of a gradient-projection algorithm of ℓ_1 -constraint KL density ratio estimation for Gaussian Markov networks. “L1BallProjection.m” is given in [Fig. 39.8](#).

sparse KL density ratio estimation gives

$$\begin{pmatrix} 0 & 0 & 1.000 \\ 0 & 0 & 0 \\ 1.000 & 0 & 0 \end{pmatrix}.$$

This implies that the change in sparsity patterns can be correctly identified.

Even when nonzero unchanged edges exist as

$$\Theta - \Theta' = \begin{pmatrix} 2 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 2 \end{pmatrix} - \begin{pmatrix} 2 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix} = \begin{pmatrix} 0 & 1 & -1 \\ 1 & 0 & 0 \\ -1 & 0 & 0 \end{pmatrix},$$

sparse KL density ratio estimation gives

$$\begin{pmatrix} 0 & 0.707 & -0.293 \\ 0.707 & 0 & 0 \\ -0.293 & 0 & 0 \end{pmatrix}.$$

Thus, the change in Markov network structure can still be correctly identified.