

Change-Point Detection in Time-Series Data by Relative Density-Ratio Estimation

Song Liu¹, Makoto Yamada², Nigel Collier³, and Masashi Sugiyama¹

¹ Tokyo Institute of Technology
2-12-1 O-okayama, Meguro-ku, Tokyo 152-8552, Japan
{song@sg.,sugi@}cs.titech.ac.jp

² NTT Communication Science Laboratories
2-4, Hikaridai, Seika-cho, Kyoto, Japan 619-0237
yamada.makoto@lab.ntt.co.jp

³ National Institute of Informatics
2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo 101-8430, Japan
collier@nii.ac.jp

Abstract. The objective of change-point detection is to discover abrupt property changes lying behind time-series data. In this paper, we present a novel statistical change-point detection algorithm that is based on non-parametric divergence estimation between two retrospective segments. Our method uses the relative Pearson divergence as a divergence measure, and it is accurately and efficiently estimated by a method of direct density-ratio estimation. Through experiments on real-world human-activity sensing, speech, and Twitter datasets, we demonstrate the usefulness of the proposed method.

Keywords: change-point detection, distribution comparison, relative density-ratio estimation, kernel methods, time-series data.

1 Introduction

Detecting abrupt changes in time-series data, called *change-point detection*, has attracted researchers in the statistics and data mining communities for decades [1–6].

Some pioneer works demonstrated good change-point detection performance by comparing the probability distributions of time-series samples over past and present intervals [1]. As both the intervals move forward, a typical strategy is to issue an alarm for a change point when the two distributions are becoming significantly different. Various change-point detection methods follow this strategy, for example, the *cumulative sum* [1], the *generalized likelihood-ratio method* [2], and the *change finder* [3].

Another group of methods that have attracted high popularity in recent years is the *subspace methods* [4, 5]. By using a pre-designed time-series model, a subspace is discovered by principle component analysis from trajectories in past

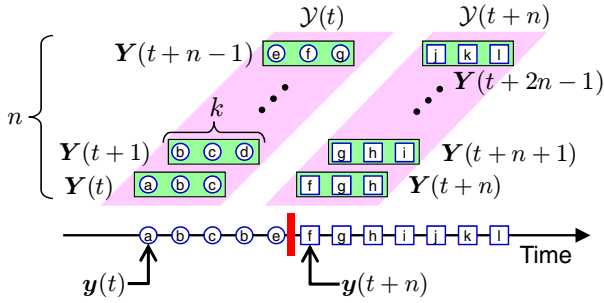


Fig. 1. Notation

and present intervals, and their dissimilarity is measured by the distance between the subspaces. One of the major approaches is called *subspace identification* [5], which compares the subspaces spanned by the columns of an *extended observability matrix* generated by a state-space model with system noise.

However, the methods explained above rely on pre-designed parametric models such as underlying probability distributions [1, 2], auto-regressive models [3], and state-space models [4, 5], for tracking some specific statistics such as the mean, the variance, and the spectrum. Thus, they are not robust against different types of changes, which significantly limits the range of applications in practice. To cope with this problem, non-parametric estimation methods such as *kernel density estimation* may be used. However, non-parametric methods tend to be less accurate in high-dimensional problems because of the so-called *curse of dimensionality*.

To overcome this difficulty, a new strategy was introduced recently which estimates the *ratio* of probability densities directly without going through density estimation [7]. In the context of change-point detection, a direct density-ratio estimation method called the *Kullback-Leibler importance estimation procedure* (KLIEP) [8] was reported to outperform competitive approaches [6] such as the *one-class support vector machine* [9] and *singular-spectrum analysis* [4].

The goal of this paper is to further advance this line of research. More specifically, our contributions in this paper are two folds.

- We apply a recently-proposed density-ratio estimation method called the *unconstrained least-squares importance fitting* (uLSIF) [10] to change-point detection. Notable advantages of uLSIF are that an analytical solution can be obtained, it achieves the optimal non-parametric convergence rate, it has optimal numerical stability, and it has higher robustness [7].
- We further improve the uLSIF-based change-point detection method by employing a state-of-the-art extension of uLSIF called *relative uLSIF* (RuLSIF) [11], which was proved to have an even better non-parametric convergence property than plain uLSIF [11], with other advantages of uLSIF maintained.

2 Problem Formulation

In this section, we formulate our change-point detection problem (see Figure 1).

Let $\mathbf{y}(t) \in \mathbb{R}^d$ be a d -dimensional time-series sample at time t . Let

$$\mathbf{Y}(t) := [\mathbf{y}(t)^\top, \mathbf{y}(t+1)^\top, \dots, \mathbf{y}(t+k-1)^\top]^\top \in \mathbb{R}^{dk}$$

be a subsequence of time series at time t with length k , where $^\top$ represents the transpose. Following the previous work [6], we treat the subsequence $\mathbf{Y}(t)$ as a sample, instead of a single point $\mathbf{y}(t)$, by which time-dependent information can be incorporated naturally. Let $\mathcal{Y}(t)$ be a set of n retrospective subsequence samples starting at time t :

$$\mathcal{Y}(t) := \{\mathbf{Y}(t), \mathbf{Y}(t+1), \dots, \mathbf{Y}(t+n-1)\}.$$

For change-point detection, let us consider two consecutive segments $\mathcal{Y}(t)$ and $\mathcal{Y}(t+n)$. Our strategy is to compute a certain dissimilarity measure between $\mathcal{Y}(t)$ and $\mathcal{Y}(t+n)$, and use it as the plausibility of change points. More specifically, the larger the dissimilarity is, the more likely the point is a change point.

Now the problems that need to be addressed are what kind of dissimilarity measure we should use and how we estimate it from data. We will discuss these issues in the next section.

3 Change-Point Detection via Density-Ratio Estimation

In this section, we first define our dissimilarity measure, and then show methods for estimating the dissimilarity measure.

3.1 Divergence-Based Dissimilarity Measure

We use a dissimilarity measure of the following form:

$$D(P_t \| P_{t+n}) + D(P_{t+n} \| P_t), \tag{1}$$

where P_t and P_{t+n} are probability distributions of samples in $\mathcal{Y}(t)$ and $\mathcal{Y}(t+n)$, respectively. $D(P \| P')$ denotes the f -divergence [12, 13]:

$$D(P \| P') := \int p'(\mathbf{Y}) f\left(\frac{p(\mathbf{Y})}{p'(\mathbf{Y})}\right) d\mathbf{Y},$$

where f is a convex function such that $f(1) = 0$, and $p(\mathbf{Y})$ and $p'(\mathbf{Y})$ are probability density functions of P and P' , respectively. Because the f -divergence is not symmetric, we use a symmetrized divergence in Eq.(1).

The f -divergence includes various popular divergences such as the *Kullback-Leibler (KL) divergence* by $f(t) = t \log t$ and the *Pearson (PE) divergence* by $f(t) = \frac{1}{2}(t - 1)^2$:

$$\text{KL}(P \| P') := \int p(\mathbf{Y}) \log \frac{p(\mathbf{Y})}{p'(\mathbf{Y})} d\mathbf{Y} \quad \text{and} \quad \text{PE}(P \| P') := \frac{1}{2} \int p'(\mathbf{Y}) \left(\frac{p(\mathbf{Y})}{p'(\mathbf{Y})} - 1\right)^2 d\mathbf{Y}.$$

In the rest of this section, we explain three methods of directly estimating the density ratio $\frac{p(\mathbf{Y})}{p'(\mathbf{Y})}$ from samples $\{\mathbf{Y}_i\}_{i=1}^n$ and $\{\mathbf{Y}'_j\}_{j=1}^n$ drawn from $p(\mathbf{Y})$ and $p'(\mathbf{Y})$: the *KL importance estimation procedure* (KLIEP) [8] in Section 3.2, *unconstrained least-squares importance fitting* (uLSIF) [10] in Section 3.3, and *relative uLSIF (RuLSIF)* [11] in Section 3.4.

3.2 Kullback-Leibler Importance Estimation Procedure (KLIEP)

KLIEP [8] is a direct density-ratio estimation algorithm that is suitable for estimating the KL divergence.

Density-Ratio Model: Let us model the density ratio $\frac{p(\mathbf{Y})}{p'(\mathbf{Y})}$ by the following kernel model:

$$g(\mathbf{Y}; \boldsymbol{\theta}) := \sum_{\ell=1}^n \theta_\ell K(\mathbf{Y}, \mathbf{Y}_\ell), \tag{2}$$

where $\boldsymbol{\theta} := (\theta_1, \dots, \theta_n)^\top$ are parameters to be learned from data samples, and $K(\mathbf{Y}, \mathbf{Y}')$ is a kernel basis function. In practice, we use the Gaussian kernel and the kernel width is chosen by cross-validation (see [8] for details).

Learning Algorithm: The parameters $\boldsymbol{\theta}$ in the model $g(\mathbf{Y}; \boldsymbol{\theta})$ are determined so that the empirical KL divergence from $p(\mathbf{Y})$ to $g(\mathbf{Y}; \boldsymbol{\theta})p'(\mathbf{Y})$ is minimized:

$$\max_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \log \left(\sum_{\ell=1}^n \theta_\ell K(\mathbf{Y}_i, \mathbf{Y}_\ell) \right) \text{ s.t. } \frac{1}{n} \sum_{j=1}^n \sum_{\ell=1}^n \theta_\ell K(\mathbf{Y}'_j, \mathbf{Y}_\ell) = 1, \quad \theta_1, \dots, \theta_n \geq 0.$$

The equality constraint is for normalization purposes because $g(\mathbf{Y}; \boldsymbol{\theta})p'(\mathbf{Y})$ should be a probability density function. The inequality constraint comes from the non-negativity of the density-ratio function. Since this is a convex optimization problem, the unique global optimal solution $\hat{\boldsymbol{\theta}}$ can be simply obtained, for example, by a gradient-projection iteration. Finally, a density-ratio estimator is given as

$$\hat{g}(\mathbf{Y}) = \sum_{\ell=1}^n \hat{\theta}_\ell K(\mathbf{Y}, \mathbf{Y}_\ell). \tag{3}$$

KLIEP was shown to achieve the optimal non-parametric convergence rate [8].

Change-Point Detection by KLIEP: Given a density-ratio estimator $\hat{g}(\mathbf{Y})$, an approximator of the KL divergence is given as

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

In the previous work [6], this KLIEP-based KL-divergence estimator was applied to change-point detection and demonstrated to be promising in experiments.

3.3 Unconstrained Least-Squares Importance Fitting (uLSIF)

Recently, another direct density-ratio estimator called uLSIF was proposed [10], which is suitable for estimating the PE divergence.

Learning Algorithm: In uLSIF, the same density-ratio model $g(\mathbf{Y}; \boldsymbol{\theta})$ as KLIEP (see Eq.(2)) is used. However, its training criterion is different; the density-ratio model is fitted to the true density ratio under the squared loss. More specifically, the parameter $\boldsymbol{\theta}$ in the model $g(\mathbf{Y}; \boldsymbol{\theta})$ is determined so that the following squared loss $J(\mathbf{Y})$ is minimized:

$$\begin{aligned}
 J(\mathbf{Y}) &:= \frac{1}{2} \int \left(\frac{p(\mathbf{Y})}{p'(\mathbf{Y})} - g(\mathbf{Y}; \boldsymbol{\theta}) \right)^2 p'(\mathbf{Y}) \, d\mathbf{Y} \\
 &= \frac{1}{2} \int \frac{p(\mathbf{Y})^2}{p'(\mathbf{Y})} p'(\mathbf{Y}) \, d\mathbf{Y} - \int p(\mathbf{Y})g(\mathbf{Y}; \boldsymbol{\theta}) \, d\mathbf{Y} + \frac{1}{2} \int g(\mathbf{Y}; \boldsymbol{\theta})^2 p'(\mathbf{Y}) \, d\mathbf{Y}.
 \end{aligned}$$

Since the first term is a constant, we focus on the last two terms. By approximating the expectations by the empirical averages, the uLSIF optimization problem is given as follows:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^n} \left[\frac{1}{2} \boldsymbol{\theta}^\top \widehat{\mathbf{H}} \boldsymbol{\theta} - \widehat{\mathbf{h}}^\top \boldsymbol{\theta} + \frac{\lambda}{2} \boldsymbol{\theta}^\top \boldsymbol{\theta} \right], \tag{4}$$

where the penalty term $\frac{\lambda}{2} \boldsymbol{\theta}^\top \boldsymbol{\theta}$ is included for regularization purposes, and λ (≥ 0) denotes the regularization parameter, which is chosen by cross validation. (see [10] for details). $\widehat{\mathbf{H}}$ is the $n \times n$ matrix and $\widehat{\mathbf{h}}$ is the n -dimensional vector defined as

$$\widehat{H}_{\ell, \ell'} := \frac{1}{n} \sum_{j=1}^n K(\mathbf{Y}'_j, \mathbf{Y}_\ell) K(\mathbf{Y}'_j, \mathbf{Y}_{\ell'}) \quad \text{and} \quad \widehat{h}_\ell := \frac{1}{n} \sum_{i=1}^n K(\mathbf{Y}_i, \mathbf{Y}_\ell).$$

It is easy to confirm that the solution $\widehat{\boldsymbol{\theta}}$ of (4) can be analytically obtained as

$$\widehat{\boldsymbol{\theta}} = (\widehat{\mathbf{H}} + \lambda \mathbf{I}_n)^{-1} \widehat{\mathbf{h}}, \tag{5}$$

where \mathbf{I}_n denotes the n -dimensional identity matrix. Finally, a density-ratio estimator is given by Eq.(3) with Eq.(5).

Change-Point Detection by uLSIF: Given a density-ratio estimator $\widehat{g}(\mathbf{Y})$, an approximator of the PE divergence can be constructed as

$$\widehat{\text{PE}} := -\frac{1}{2n} \sum_{j=1}^n \widehat{g}(\mathbf{Y}'_j)^2 + \frac{1}{n} \sum_{i=1}^n \widehat{g}(\mathbf{Y}_i) - \frac{1}{2}.$$

This approximator is derived from the following expression of the PE divergence:

$$\text{PE}(P \| P') = -\frac{1}{2} \int \left(\frac{p(\mathbf{Y})}{p'(\mathbf{Y})} \right)^2 p'(\mathbf{Y}) \, d\mathbf{Y} + \int \left(\frac{p(\mathbf{Y})}{p'(\mathbf{Y})} \right) p(\mathbf{Y}) \, d\mathbf{Y} - \frac{1}{2}. \tag{6}$$

Notable advantages of uLSIF are that its solution can be computed analytically, it possesses the optimal non-parametric convergence rate, it has the optimal numerical stability, and it has higher robustness [7]. As experimentally demonstrated in our supplementary technical report [14], uLSIF-based change-point detection compares favorably with the KLIEP-based method.

3.4 Relative uLSIF (RuLSIF)

Depending on the condition of the denominator density $p'(\mathbf{Y})$, the density-ratio value $\frac{p(\mathbf{Y})}{p'(\mathbf{Y})}$ can be unbounded (i.e., they can be infinity). This is actually problematic because the non-parametric convergence rate of uLSIF is governed by the “sup”-norm of the true density-ratio function: $\max_{\mathbf{Y}} \frac{p(\mathbf{Y})}{p'(\mathbf{Y})}$. To overcome this problem, *relative density-ratio estimation* was introduced [11].

Relative PE Divergence: Let us consider the α -relative PE-divergence for $0 \leq \alpha < 1$:

$$\text{PE}_\alpha(P\|P') := \text{PE}(P\|\alpha P + (1 - \alpha)P') = \int p'_\alpha(\mathbf{Y}) (r_\alpha(\mathbf{Y}) - 1)^2 d\mathbf{Y},$$

where $p'_\alpha(\mathbf{Y}) = \alpha p(\mathbf{Y}) + (1 - \alpha)p'(\mathbf{Y})$ and $r_\alpha(\mathbf{Y}) = \frac{p(\mathbf{Y})}{p'_\alpha(\mathbf{Y})}$. We refer to $r_\alpha(\mathbf{Y})$ as the α -relative density ratio. The α -relative density ratio is reduced to the plain density ratio if $\alpha = 0$, and it tends to be “smoother” as α gets larger. Indeed, the α -relative density ratio is bounded above by $1/\alpha$ for $\alpha > 0$, even when the plain density ratio $\frac{p(\mathbf{Y})}{p'(\mathbf{Y})}$ is unbounded. This was proved to contribute to improving the estimation accuracy [11].

Learning Algorithm: In the same way as the uLSIF method, the parameter θ of the model $g(\mathbf{Y}; \theta)$ is learned by minimizing the squared difference between true and estimated ratios:

$$\begin{aligned} J(\mathbf{Y}) &= \frac{1}{2} \int p'_\alpha(\mathbf{Y})(r_\alpha(\mathbf{Y}) - g(\mathbf{Y}; \theta))^2 d\mathbf{Y} \\ &= \frac{1}{2} \int p'_\alpha(\mathbf{Y})r_\alpha^2(\mathbf{Y})d\mathbf{Y} - \int p(\mathbf{Y})r_\alpha(\mathbf{Y})g(\mathbf{Y}; \theta) d\mathbf{Y} \\ &\quad + \frac{\alpha}{2} \int p(\mathbf{Y})g(\mathbf{Y}; \theta)^2 d\mathbf{Y} - \frac{1 - \alpha}{2} \int p'(\mathbf{Y})g(\mathbf{Y}; \theta)^2 d\mathbf{Y}, \end{aligned}$$

where the first term is a constant term. Note that we still use the same kernel model (2) as $g(\mathbf{Y}; \theta)$ for approximating the α -relative density ratio.

Again, by ignoring the constant and approximating the expectations by empirical averages, the α -relative density ratio can be learned in the same way as the plain density ratio. Indeed, the optimization problem of a relative variant of uLSIF, called RuLSIF, is given as the same form as uLSIF; the only difference is the definition of the matrix $\widehat{\mathbf{H}}$, which is now given by

$$\widehat{H}_{\ell, \ell'} := \frac{\alpha}{n} \sum_{i=1}^n K(\mathbf{Y}_i, \mathbf{Y}_\ell) K(\mathbf{Y}_i, \mathbf{Y}_{\ell'}) + \frac{(1-\alpha)}{n} \sum_{j=1}^n K(\mathbf{Y}'_j, \mathbf{Y}_\ell) K(\mathbf{Y}'_j, \mathbf{Y}_{\ell'}).$$

RuLSIF inherits the advantages of uLSIF, i.e., its solution can be computed analytically, it has the superior numerical stability, and it has higher robustness; furthermore, RuLSIF possesses an even better non-parametric convergence property than uLSIF [11].

Change-Point Detection by RuLSIF: By using an estimator $\widehat{g}(\mathbf{Y})$ of the α -relative density ratio, the α -relative PE divergence can be approximated as

$$\widehat{PE}_\alpha := -\frac{\alpha}{2n} \sum_{i=1}^n \widehat{g}(\mathbf{Y}_i)^2 - \frac{1-\alpha}{2n} \sum_{j=1}^n \widehat{g}(\mathbf{Y}'_j)^2 + \frac{1}{n} \sum_{i=1}^n \widehat{g}(\mathbf{Y}_i) - \frac{1}{2}.$$

As experimentally demonstrated in our supplementary technical report [14], the RuLSIF-based change-point detection performs even better than the plain uLSIF-based method. Thus, we focus on RuLSIF in the experiments in Section 4.

4 Experiments

In this section, we experimentally investigate the performance of the proposed and existing change-point detection methods.

First, we use a human activity dataset and a speech dataset. The human activity dataset is a subset of the *Human Activity Sensing Consortium (HASC) challenge 2011*, which provides human activity information collected by portable three-axis accelerometers. The speech dataset is the *IPSJ SIG-SLP Corpora and Environments for Noisy Speech Recognition (CENSREC)* dataset provided by National Institute of Informatics (NII), which records human voice in a noisy

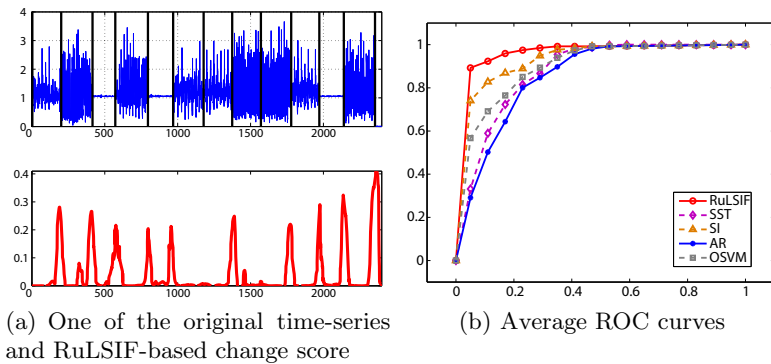


Fig. 2. HASC human-activity dataset (<http://hasc.jp/hc2011/>)

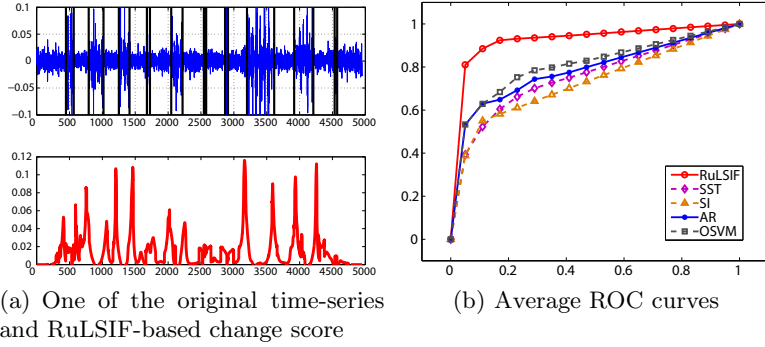


Fig. 3. NII speech dataset (<http://research.nii.ac.jp/src/eng/list/index.html>)

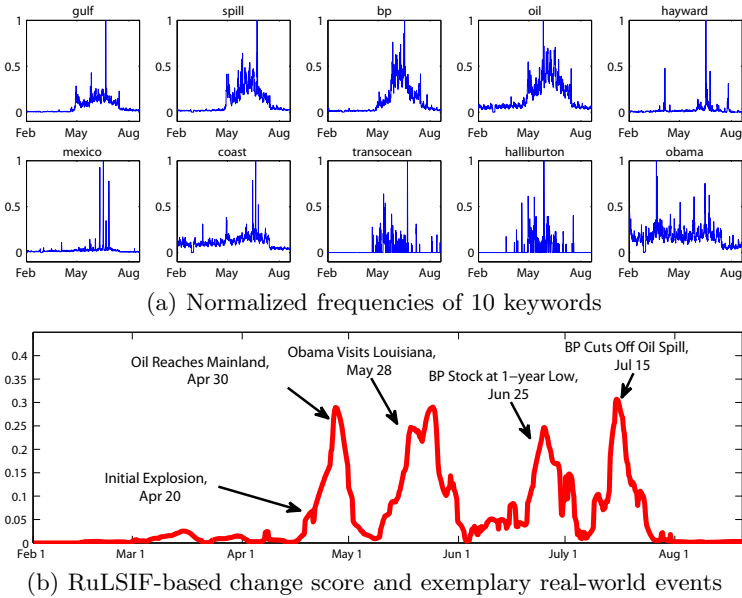


Fig. 4. Twitter dataset (<http://www.ark.cs.cmu.edu/tweets/>)

environment. We compare our RuLSIF-based method with several state-of-the-art methods: *Singular spectrum transformation* (SST) [4], *subspace identification* (SI) [5], *auto regressive* (AR) [3], and *one-class support vector machine* (OSVM) [9]. Examples of RuLSIF-based change score and ROC curves over 10 datasets are plotted in Figures 2 and 3, showing that the proposed RuLSIF-based method outperforms other methods.

Finally, we apply the proposed change-point detection method to the *CMU Twitter dataset*, which is an archive of Twitter messages that have been collected from April 2010 to October 2010 via the Twitter API. Here we track the degree

of popularity of a given topic by monitoring the frequency of selected keywords. More specifically, we focus on events related to “*Deepwater Horizon oil spill in the Gulf of Mexico*” which occurred on April 20, 2010, and was widely broadcast among the Twitter community. We use the frequency of 10 keywords: “*gulf*”, “*spill*”, “*bp*”, “*oil*”, “*hayward*”, “*mexico*”, “*coast*”, “*transocean*”, “*halliburton*”, and “*obama*” (see Figure 4(a)). For quantitative evaluation, we referred to the Wikipedia entry “Timeline of the Deepwater Horizon oil spill” as a real-world event source. The change-point score obtained by the proposed RuLSIF-based method is plotted in Figure 4(b), where four occurrences of important real-world events show the development of this news story.

As we can see from Figure 4(b), the change-point score increases immediately after the initial explosion of the deepwater horizon oil platform and soon reaches the first peak when oil was found on the sea shore of Louisiana on April 30. Shortly after BP announced its preliminary estimation on the amount of leaking oil, the change-point score rises quickly again and reaches its second peak at the end of May, at which time President Obama visited Louisiana to assure local residents of the federal government’s support. On June 25, the BP stock was at its one year’s lowest price, while the change-point score spikes at the third time. Finally, BP cuts off the spill on July 15, as the score reaches its last peak.

5 Conclusion

We extended the existing KLIEP-based change detection method and proposed to use uLSIF or RuLSIF as a building block. Through experiments, we demonstrated that the RuLSIF-based change detection method is promising.

SL was supported by NII internship fund and the JST PRESTO program. MY and MS were supported by the JST PRESTO program. NC was supported by NII Grand Challenge project fund.

References

1. Basseville, M., Nikiforov, I.V.: *Detection of Abrupt Changes: Theory and Application*. Prentice-Hall, Inc., Upper Saddle River (1993)
2. Gustafsson, F.: The marginalized likelihood ratio test for detecting abrupt changes. *IEEE Transactions on Automatic Control* 41(1), 66–78 (1996)
3. Takeuchi, Y., Yamanishi, K.: A unifying framework for detecting outliers and change points from non-stationary time series data. *IEEE Transactions on Knowledge and Data Engineering* 18(4), 482–489 (2006)
4. Moskvina, V., Zhigljavsky, A.: Change-point detection algorithm based on the singular-spectrum analysis. *Communications in Statistics: Simulation and Computation* 32, 319–352 (2003)
5. Kawahara, Y., Yairi, T., Machida, K.: Change-point detection in time-series data based on subspace identification. In: *Proceedings of the 7th IEEE International Conference on Data Mining*, pp. 559–564 (2007)
6. Kawahara, Y., Sugiyama, M.: Sequential change-point detection based on direct density-ratio estimation. *Statistical Analysis and Data Mining* 5(2), 114–127 (2012)

7. Sugiyama, M., Suzuki, T., Kanamori, T.: Density Ratio Estimation in Machine Learning. Cambridge University Press, Cambridge (2012)
8. Sugiyama, M., Suzuki, T., Nakajima, S., Kashima, H., von Buenau, P., Kawanabe, M.: Direct importance estimation for covariate shift adaptation. *Annals of the Institute of Statistical Mathematics* 60(4), 699–746 (2008)
9. Desobry, F., Davy, M., Doncarli, C.: An online kernel change detection algorithm. *IEEE Transactions on Signal Processing* 53(8), 2961–2974 (2005)
10. Kanamori, T., Hido, S., Sugiyama, M.: A least-squares approach to direct importance estimation. *Journal of Machine Learning Research* 10, 1391–1445 (2009)
11. Yamada, M., Suzuki, T., Kanamori, T., Hachiya, H., Sugiyama, M.: Relative density-ratio estimation for robust distribution comparison. *Advances in Neural Information Processing Systems* 24, 594–602 (2011)
12. Ali, S.M., Silvey, S.D.: A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society, Series B* 28(1), 131–142 (1966)
13. Csiszár, I.: Information-type measures of difference of probability distributions and indirect observation. *Studia Scientiarum Mathematicarum Hungarica* 2, 229–318 (1967)
14. Liu, S., Yamada, M., Collier, N., Sugiyama, M.: Change-point detection in time-series data by relative density-ratio estimation. *arXiv 1203.0453* (2012)