

Multiscale parametric approach for change point detection

A. Suvorikova, V. Spokoiny, N. Buzun
{suvorikova,spokoiny,buzun}@wias-berlin.de

¹ Weierstrass Institute for Applied Analysis and Stochastics

² Institute for Information Transmission Problems

Abstract. This work presents a novel algorithm for change point detection, that can be applied for analysis of data of unknown nature. It is based on likelihood-ratio test statistics, as its behaviour can be described in terms of χ^2 -distribution even in case of model misspecification. To discover change point in the quickest way, statistics is calculated in a set of running windows of different scales. Algorithm is self-tuned: critical values are justified by data and calculated with multiplier bootstrap procedure. To make the method more robust for outliers, the concept of change-point patterns is presented.

Keywords: change point detection, multiscale inference

1 Introduction

The problem of change point detection has a wide range of applications, that varies from life-critical to pure scientific ones. It appears each time one needs to explore a set of random data and make a decision about homogeneity of its structure. In other words, the problem can be stated as two following questions: were there any structural changes in the nature of observed data? At which moments, if so? These and similar questions arise in many areas of theoretical and engineering research. For example, algorithms of change point detection are used for identification and elimination of faults of aeroplane's navigation system, so as to perform better geolocation Nikiforov [2003]. There are many other examples of real-world applications, such as analysis of stock markets Lavielle and Teyssiere [2006] or anomaly detection in computer traffic Tartakovsky et al. [2006], Casas et al. [2010]. The present work mainly focuses on the *sequential* or *online* change point detection. In this case the data is aggregated from running random process. Let Y_τ be the observation at the current moment τ , $\tau > 1$. The moment τ is a *change point*, if stochastic properties of observed signal have undergone some changes:

$$\begin{cases} Y_t \sim \mathbb{P}_1 & t < \tau, \\ Y_t \sim \mathbb{P}_2 & t \geq \tau. \end{cases}$$

The goal is to find such structural breaks as soon as possible. The problem arises across many scientific areas: quality control Lai [1995], cybersecurity Blazek and Kim [2001], Wang et al. [2004], econometrics Spokoiny [2009], Mikosch and Starica [2004], geodesy e.t.c. Shiryaev [1963] describes classical results in change point detection theory. Overview of the state-of-art methods is presented in Polunchenko and Tartakovsky [2011] or Shiryaev [2010]. The problem in hand can be easily reduced to the problem of hypothesis testing in a rolling window. Let t' be a candidate for a change point and let $(Y_{t'-h}, \dots, Y_{t'+h-1})$ be observed data in the rolling window of size $2h$, then

$$H_0 : Y_t \sim \mathbb{P}_1, \quad t' - h \leq t \leq t' + h - 1$$

$$H_1 : \begin{cases} Y_t \sim \mathbb{P}_1, & t' - h \leq t \leq t' - 1, \\ Y_t \sim \mathbb{P}_2, & t' \leq t \leq t' + h - 1. \end{cases}$$

One popular solution is likelihood-ratio test (LRT). Its application for change point detection goes back at least as far as Shewhart [1931]. This work presents the concept of *control chart* for quality control. The work Quandt [1960] proposes application of LRT for detection of breaks in linear regression model. It was further developed by many authors, e.g. Kim and Siegmund [1989], Haccou et al. [1987], Srivastava and Worsley [1986]. Liu et al. [2008], Zou et al. [2007] investigate LRT for change point detection for nonparametric case. In general, nonparametric approaches need more information for change point detection than their parametric alternatives. Introduction of *parametric assumption*: $\mathbb{P}_1, \mathbb{P}_2 \in (\mathbb{P}(\theta), \theta \in \Theta \subseteq \mathbb{R}^p)$ allows to reduce average number of observations. The state-of-the-art review of parametric models based on LRT and its

application to economics and bio-informatics are presented by Chen and Gupta [2012]. The paper Gombay [2000] explores how LRT can be used for sequential change point detection in case $\mathbb{P}(\theta)$ is exponential family. Lai [1995] proposes 'window-limited LRT schemes': test statistics is calculated in rolling window. This concept naturally expands to on-line change point detection. Many works are dedicated to asymptotic behaviour of LRT, e.g. Jandhyala and Fotopoulos [1999] obtains lower and upper bounds for distribution of asymptotic maximum likelihood estimator. The work Kim [1994] provides a very detailed study of its asymptotic behaviour in linear regression models. Similar results for change in mean of a Gaussian process are in Fotopoulos et al. [2010]. As far as the authors know, the most comprehensive study of the LRT behaviour is done by Fan et al. [2001]. It shows that LRT is asymptotically χ^2 distributed. The idea of the proof is based on the *Wilks's phenomenon* Wilks [1938], Boucheron and Massart [2011]. The aim of the present paper is to describe the LRT behaviour in finite-sample case using non-asymptotic Wilks and Fischer theorems Spokoiny [2012]. It is shown that the distribution of LRT is similar to ordinary χ^2 under H_0 . In case non-homogeneous sample inside a rolling window, the systematic drift of LRT appears. Thus, under H_1 the test statistic behaves like non-central χ^2 random value. This drift is referred to as a *change-point pattern*. The result holds for both correct and misspecified parametric models. The cornerstone of the new change point detection procedure is the concept of change-point pattern. The geometry of a pattern depends on a type of switch between distributions the data obeys before and after a change respectively. Three examples are presented at the Fig. 1. The triangle pattern appears in case of an abrupt switch from $\mathbb{P}_{\theta_1^*}$ to $\mathbb{P}_{\theta_2^*}$. A smooth transition between two regimes entails the trapezium change-point pattern. And an inverted triangle pattern appears due to a change in coefficients of linear regression. The control of a change-point pattern instead of a single LRT-value allows to reduce false-alarm rate to zero.

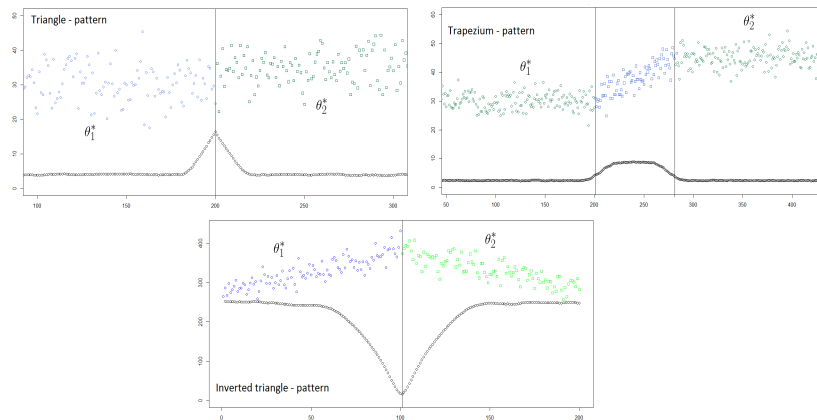


Fig. 1. Type of change point and the geometry of change-point pattern

The existence of change-point patterns is the corollary of the Theorem 2. However, the description of their theoretical properties is beyond the scope of the present work.

Any procedure of change point detection exploits the concept of "memory". It is the number of observations to be analyzed to detect a change point. Introduction of *multiscale* approach provides automatic optimal choice of memory parameter. The results are described in Section 3. This technique is popular, e.g. Frick et al. [2014], Spokoiny [2009] and performs analysis of the data on different scales simultaneously. The procedure proposed below computes the test statistics in rolling windows of different sizes and controls change-point pattern at each level. The greater number of scales at which a change-point pattern is detected, the more sure algorithm is, that there is a change point.

Under some assumptions on the frequency of change points provided in Section 3, the method is applied to the *multiple* change point problem. The survey on existing models can be found in Chib [1998]. The last, but not at all the least feature of the proposed algorithm is that critical values for test statistic are computed in a data-driven way. The idea is to use the multiplier bootstrap procedure Chernozhukov et al. [2013]. The work of Spokoiny and Zhilova [2014] shows that it can be used for the construction of confidence intervals even in case of a misspecified parametric model. Despite the fact, that theoretical properties of data-driven critical values are beyond the scope of this paper, the procedure of computation is presented in Algorithm 3.

The paper is organized as follows. Section 2 presents the description of the algorithm. Theoretical properties of the procedure are discussed in Section 3. Section 4 compares the new algorithm with existing methods using simulated data sets. It also illustrates the performance of the method on a real data set.

2 Algorithm

This section provides the description of the proposed algorithm. Let $(\mathbb{P}(\theta), \theta \in \Theta \subseteq \mathbb{R}^p)$ be a local parametric assumption about the nature of data inside a window $(Y_{t-h}, \dots, Y_{t+h-1})$. Here and further we assume, that the observations are independent and identically distributed. The generalised likelihood ratio test is

$$T_h(t) = \sup_{\theta \in \Theta} L(\theta; Y_{t-h}, \dots, Y_{t-1}) + \sup_{\theta \in \Theta} L(\theta; Y_t, \dots, Y_{t+h-1}) \\ - \sup_{\theta \in \Theta} \{L(\theta; Y_{t-h}, \dots, Y_{t-1}) + L(\theta; Y_t, \dots, Y_{t+h-1})\},$$

where $L(\theta; \cdot)$ is a log-likelihood function. To control a change point pattern, the procedure monitors $2h$ values of the LRT simultaneously:

$$\mathbb{T}_h(t) = (\sqrt{2T_h(t-h)}, \dots, \sqrt{2T_h(t+h-1)}).$$

The test statistics in hand is a convolution of $\mathbb{T}_h(t)$ with a predefined change-point pattern $P_h \in \mathbb{R}^{2h}$.

$$\widehat{\mathbb{T}}_h(t) = \langle \mathbb{T}_h(t), P_h \rangle.$$

Under *online* framework, the algorithm marks a time moment τ at a scale h as a change point, if the test statistic $\widehat{\mathbb{T}}_h(t)$ exceeds critical value $z(h)$ at the moment $t = \tau + h$:

$$\{\tau : \widehat{\mathbb{T}}_h(\tau + h) > z(h)\}.$$

Under *offline* setting, τ is marked as a change point if

$$\{\tau = \operatorname{argmax}_{t \in \{1, \dots, M\}} \sum_{h \in H} w_h \widehat{\mathbb{T}}_h(t), \quad \exists h \in H : \widehat{\mathbb{T}}_h(\tau) > z(h)\},$$

where M is the number of observations, and $\{w_h\}_{h \in H}$ are weights for window size preferences.

In both cases the procedure repeats itself simultaneously on different scales $H = \{h_i\}$. The greater element position in ordered sequence H with which τ is marked as change point, the more sure algorithm is, that τ the *true* change point is. The multiplicity correction introduced in Algorithm 3 allows to avoid an increase in false-alarm rate with the growth of H size.

Algorithm 1, 2 summarises above ideas for sequential case and case with preloaded fixed data. Here the current moment is supposed to be $\tau + 2h_N - 2$ and a candidate for the change point is τ . Designation $(t_1 : t_2)$ means range of natural values $t_1, t_1 + 1, \dots, t_2$.

Algorithm 3 presents the procedure of calculation of a critical value z_h for a fixed window size $2h$. Let $\mathbb{Y} = (Y_1, \dots, Y_M)$ be a training set (the described procedure is applicable only for independent data \mathbb{Y}). Let weighted likelihood

function be a convolution of independent likelihood components and a weight vector (u_1, \dots, u_M) :

$$L^b(\theta; Y_1, \dots, Y_M) = \sum_{m=1}^M u_m l(\theta, Y_m), \quad (\text{Lb})$$

where $\{u_m\}_{m=1}^M$ are i.i.d. and $\mathbb{E}u_m = \text{Var } u_m = 1$. Then bootstrapped generalised likelihood ratio test is

$$\begin{aligned} T_h^b(t) &= \sup_{\theta \in \Theta} L^b(\theta; Y_{t-h}, \dots, Y_{t-1}) + \sup_{\theta \in \Theta} L^b(\theta; Y_t, \dots, Y_{t+h-1}) \quad (\text{Tb}) \\ &\quad - \sup_{\theta \in \Theta} \{L^b(\theta; Y_{t-h}, \dots, Y_{t-1}) + L^b(\theta + \hat{\theta}_{12}; Y_t, \dots, Y_{t+h-1})\}, \end{aligned}$$

$$\hat{\theta}_{12} = \underset{\theta}{\text{argmax}} L(\theta; Y_t, \dots, Y_{t+h-1}) - \underset{\theta}{\text{argmax}} L(\theta; Y_{t-h}, \dots, Y_{t-1}).$$

Parameter $\hat{\theta}_{12}$ allows bootstrap calibration with corrected bias, which described in remark 4 in Section 3.

Algorithm 3 use multiplicity correction for multiple hypothesis testing: $H_h : \max_{\tau} \hat{\mathbb{T}}_h^b(\tau) < z(h)$ for each h . Let $z(h, \alpha)$ be α quantile of variable $\max_{\tau} \hat{\mathbb{T}}_h^b(\tau)$. The probability that at least one hypothesis is false equals to

$$\mathbb{P}(\{\exists h : \max_{\tau} \hat{\mathbb{T}}_h^b(\tau) - z(h, \alpha) > 0\}) = \mathbb{P}(\{\exists h : \text{p-value}(\max_{\tau} \hat{\mathbb{T}}_h^b(\tau)) < \alpha\}) \geq \alpha.$$

One may decrease above probability by confidence reduction:

$$\mathbb{P}(\{\exists h : \text{p-value}(\max_{\tau} \hat{\mathbb{T}}_h^b(\tau)) < \alpha - \alpha'\}) = \alpha.$$

$Q_h(t) = 0$ – change point signals;
 H – window sizes set;
 get $z(h)$ by Algorithm 3;
foreach *window position* t **do**
 foreach h **do**
 add $T_h(t)$ to \mathbb{T}_h ;
 $\widehat{\mathbb{T}}_h = \langle \mathbb{T}_h(t-h), P_h \rangle$;
 if $\widehat{\mathbb{T}}_h > z(h)$ **and**
 $Q_{(1:h)}(t-2h:t) = 0$ **then**
 $Q_h(t) = 1$;
 end
 end
 if $\max_h Q_h(t) = 1$ **then**
 t is change point;
 end
end

Algorithm 1: LRTOonline.

S – change points set; H – window sizes set;
 w_j – window size weights;
function FindCP(Y_1, \dots, Y_M):
 get $z(h)$ by Algorithm 3;
foreach h **do**
 foreach *window position* t **do**
 compute $T_h(t)$;
 end
 foreach τ **do**
 $\widehat{\mathbb{T}}_h(\tau) = \langle \mathbb{T}_h(\tau), P_h \rangle$;
 end
end
 $\tau = \operatorname{argmax}_\tau \sum_{h \in H} w_h \widehat{\mathbb{T}}_h(\tau)$;
if $\exists h : \mathbb{T}_h(\tau) > z(h)$ **then**
 add τ to S ;
 FindCP(Y_1, \dots, Y_τ);
 FindCP(Y_τ, \dots, Y_M);
end

Algorithm 2: LRTOoffline.

Data: (Y_1, \dots, Y_M), h , P_h , S – weights generation counts
Result: f_h^b – bootstrap distribution of maximal convolution inside the window
for $s = 1$ **to** S **do**
 generate $u = (u_1, \dots, u_M)$;
 foreach *window position* t **do**
 compute $T_h^b(t)$;
 end
 foreach τ **do**
 $\widehat{\mathbb{T}}_h^b(\tau) = \langle T_h^b(\tau), P_h \rangle$;
 end
 add $\max_\tau \widehat{\mathbb{T}}_h^b(\tau)$ to f_h^b ;
end

Data: $H = (h_1, \dots, h_N)$, f_h^b , α – confidence
Result: critical values $z(h)$
 Multiplicity correction:
for $s = 1$ **to** S **do**
 generate $u = (u_1, \dots, u_M)$;
 add
 $\min_h \text{p-value}(\max_\tau \widehat{\mathbb{T}}_h^b(\tau), f_h^b)$ to empirical distribution \mathbb{P}_f
end
 find α' from condition
 $\mathbb{P}_f(\min_h \text{p-value}(\cdot) < \alpha - \alpha') = \alpha$;
foreach h *in* H **do**
 $z(h) = \text{quantile}(f_h^b, \alpha - \alpha')$;
end

Algorithm 3: Critical values calibration

3 Theoretical results

3.1. LRT statistic

This section presents main results that describe theoretical properties of the likelihood-ratio statistics (LRT). They are essential for the proposed algorithm of change point detection. Further assume that log-likelihood function $L(\theta) = L(Y, \theta)$ has rather precise approximation by its quadratic part in local region $\Theta_0(r)$ of θ^* , $\Theta_0(r) \subseteq \mathbb{R}^p$, where

$$\theta^* = \operatorname{argmax}_{\theta} \mathbb{E}L(\theta), \quad \hat{\theta} = \operatorname{argmax}_{\theta} L(\theta)$$

and $\Theta_0(r) = \{\|D(\theta - \theta^*)\| < r\}$. Spokoiny [2012] provides required conditions for justified quadratic approximation and parameter concentration in the local region. Approximation error involves the next variables for its estimation:

$$\begin{aligned} \alpha(\theta, \theta_0) &= L(\theta) - L(\theta_0) - (\theta - \theta_0)^T \nabla L(\theta_0) + \frac{1}{2} \|D(\theta - \theta_0)\|^2, \\ \chi(\theta, \theta_0) &= D^{-1} \nabla \alpha(\theta, \theta_0) = D^{-1} (\nabla L(\theta) - \nabla L(\theta_0)) + D(\theta - \theta_0). \end{aligned}$$

Let in region $\Theta_0(r)$ with probability $1 - e^{-x}$:

$$\frac{|\alpha(\theta, \theta^*)|}{\|D(\theta - \theta^*)\|} \leq \diamond(r, x), \quad \|\chi(\theta, \theta^*)\| \leq \diamond(r, x), \quad (\text{A})$$

where $\diamond(r, x) = (\delta(r) + 6\nu_0 z_H(x)\omega)r$,

$$D^2(\theta) = -\nabla^2 \mathbb{E}L(\theta), \quad D = D(\theta^*), \quad (\text{D})$$

$$\|D^{-1} D^2(\theta) D^{-1} - I_p\| \leq \delta(r), \quad (\text{dD})$$

$$\forall \lambda \leq g, \gamma_1 \gamma_2 \in \mathbb{R}^p : \quad \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \frac{\gamma_1^T \nabla^2 \overset{\circ}{L}(\theta) \gamma_2}{\|D\gamma_2\| \|D\gamma_2\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad (\text{ED2})$$

$$z_H(x) = \sqrt{H} + \sqrt{2x} + \frac{g^{-2}x + 1}{g} H, \quad H = 6p.$$

Condition (dD) ensures quadratic approximation of $\mathbb{E}L(\theta)$ and (ED2) ensures linear approximation of centered likelihood $\overset{\circ}{L}(\theta) = L(\theta) - \mathbb{E}L(\theta)$.

Firstly, to provide a simple non-strict explanation of what kind of distribution the main statistic T_h is supposed to have, review T_h as

$$T_h = L(\hat{\theta}) - L(\hat{\theta}_{H_0}), \quad L(\theta_1, \theta_2) = L_1(\theta_1) + L_2(\theta_2),$$

$$L_1 = L(Y_1, \dots, Y_h), \quad L_2 = L(Y_h, \dots, Y_{2h}),$$

where $\hat{\theta}_{H_0}$ is argmax of L under condition $H_0 : \theta_1^* = \theta_2^*$. Then due to quadratic approximation T_h corresponds to Taylor equation with point $\hat{\theta}$:

$$T_h \approx \frac{1}{2} \|D(\hat{\theta} - \hat{\theta}_{H_0})\|^2.$$

If $\hat{\theta}$ and $\hat{\theta}_{H_0}$ tend to be Normal and H_0 is true then their difference are close to a centered Normal variable. If H_0 is false – the Normal variable will have mean that is equal to $D(\theta^* - \theta_{H_0}^*)$.

The next equations describes strict equation for LRT statistic distribution in quadratic model case. Sum of two quadratic functions $L_1(\theta) + L_2(\theta)$ is a quadratic function with central point $\hat{\theta}$.

$$\begin{aligned} L(\theta) &= L_1(\theta) + L_2(\theta) \\ &= L_1(\hat{\theta}_1) + L_2(\hat{\theta}_2) - \frac{1}{2}(\theta - \hat{\theta}_1)^T D_1^2(\theta - \hat{\theta}_1) - \frac{1}{2}(\theta - \hat{\theta}_2)^T D_2^2(\theta - \hat{\theta}_2) \\ &= L(\hat{\theta}) - \frac{1}{2}(\theta - \hat{\theta})^T D^2(\theta - \hat{\theta}), \\ \hat{\theta} &= D^{-2}(D_1^2 \hat{\theta}_1 + D_2^2 \hat{\theta}_2), \quad D^2 = D_1^2 + D_2^2. \end{aligned}$$

$$\begin{aligned} T_h &= L_1(\hat{\theta}_1) + L_2(\hat{\theta}_2) - L(\hat{\theta}) \\ &= \frac{1}{2}(\hat{\theta} - \hat{\theta}_1)^T D_1^2(\hat{\theta} - \hat{\theta}_1) + \frac{1}{2}(\hat{\theta} - \hat{\theta}_2)^T D_2^2(\hat{\theta} - \hat{\theta}_2). \\ \hat{\theta} - \hat{\theta}_1 &= D^{-2}(D_1^2 \hat{\theta}_1 + D_2^2 \hat{\theta}_2) - \hat{\theta}_1 = D^{-2} D_2^2(\hat{\theta}_2 - \hat{\theta}_1), \\ \hat{\theta} - \hat{\theta}_2 &= D^{-2}(D_1^2 \hat{\theta}_1 + D_2^2 \hat{\theta}_2) - \hat{\theta}_2 = D^{-2} D_1^2(\hat{\theta}_1 - \hat{\theta}_2). \\ 2T_h &= (\hat{\theta}_2 - \hat{\theta}_1)^T \Sigma^2(\hat{\theta}_2 - \hat{\theta}_1), \end{aligned}$$

where

$$\Sigma^2 = D_2^2 D^{-2} D_1^2 D^{-2} D_2^2 + D_1^2 D^{-2} D_2^2 D^{-2} D_1^2 = D_1^2 D^{-2} D_2^2 \approx \frac{1}{4} D^2, \quad (\text{S})$$

$$D_1^2 = -\nabla^2 \mathbb{E}L(\theta_1^*), \quad D_2^2 = -\nabla^2 \mathbb{E}L(\theta_2^*), \quad D_1 \approx D_2.$$

In quadratic model following equations provides replacement of $\hat{\theta}_2$, $\hat{\theta}_1$ in the equation for T_h with regard to condition $\chi(\theta, \theta^*) = 0$:

$$D_1(\hat{\theta}_1 - \theta_1^*) = \xi_1 = D_1^{-1} \nabla L(\theta_1^*), \quad D_2(\hat{\theta}_2 - \theta_2^*) = \xi_2 = D_2^{-1} \nabla L(\theta_2^*).$$

The next theorem concludes these considerations to generalized result for non-quadratic model.

Theorem 1. Assume condition (L^*) and quadratic Laplace approximation (A) of L_1 and L_2 are fulfilled with probability $1 - 2e^{-x}$, additionally with probability $1 - 2e^{-x}$

$$\|\xi_i\| \leq z(x), \quad z^2(x) = \max_i p_{B_i} + 6\lambda_{B_i} x,$$

$$B_i = D_i^{-1} \text{Var}(\nabla L_i(\theta_i^*)) D_i^{-1}, \quad p_B = \text{tr}(B), \quad \lambda_B = \lambda_{\max}(B). \quad (\text{B})$$

Then in the local region with probability $1 - 8e^{-x}$

$$2T_h = \|\xi_{12} + \theta_{12}^*\|^2 + O(\{r + z(x)\} \diamond(r, x)),$$

where

$$\xi_{12} = \Sigma(D_2^{-1} \xi_2 - D_1^{-1} \xi_1), \quad \theta_{12}^* = \Sigma(\theta_2^* - \theta_1^*).$$

Remark 1. In increasing sample size $n \rightarrow \infty$ the stochastic component tends to Normal distribution:

$$\xi_{12} \rightarrow \mathcal{N}(0, B_1 + B_2).$$

Remark 2. Both $L_1(\hat{\theta})$ and $L_2(\hat{\theta})$ should have an opportunity to be presented in quadratic form in local regions $\Theta_1(r) = \{\theta : \|D_1(\theta - \theta_1^*)\|\}$ and $\Theta_2(r) = \{\theta : \|D_2(\theta - \theta_2^*)\|\}$. For the condition $\hat{\theta} \in \Theta_1(r) \cap \Theta_2(r)$ the restriction of the parameter θ^* variability is required

$$\|D(\theta_1^* - \theta_2^*)\| \leq r. \quad (\text{L}^*)$$

Proof of a similar statement (theorem 1) for statistic $\sqrt{2T_h}$ one could get from condition (A). With probability $1 - 2e^{-x}$

$$\begin{aligned} \left| T_h(\hat{\theta}_1, \hat{\theta}_2) - \frac{1}{2} \|\Sigma(\hat{\theta}_2 - \hat{\theta}_1)\|^2 \right| &\leq 2\|D_1(\hat{\theta}_1 - \hat{\theta})\| \diamond(r, x) + 2\|D_2(\hat{\theta}_2 - \hat{\theta})\| \diamond(r, x) \\ &\leq 4\|\Sigma(\hat{\theta}_2 - \hat{\theta}_1)\| \diamond(r, x). \end{aligned}$$

Inequality $|a - b| \leq |a^2 - b^2|/b$, $b > 0$ converts the previous term to

$$\left| \sqrt{2T_h(\hat{\theta}_1, \hat{\theta}_2)} - \|\Sigma(\hat{\theta}_2 - \hat{\theta}_1)\| \right| \leq 8\diamond(r, x).$$

Replacement $(\hat{\theta}_1, \hat{\theta}_2)$ with $(D_1^{-1}\xi_1 + \theta_1^*, D_2^{-1}\xi_2 + \theta_2^*)$ results in

$$\begin{aligned} &\left| \|\Sigma(\hat{\theta}_2 - \hat{\theta}_1)\| - \|\xi_{12} + \theta_{12}^*\| \right| \\ &\leq \|\Sigma(\hat{\theta}_1 - \theta_1^*) - \Sigma D_1^{-1}\xi_1\| + \|\Sigma(\hat{\theta}_2 - \theta_2^*) - \Sigma D_2^{-1}\xi_2\| \leq 2\diamond(r, x). \end{aligned}$$

The next theorem summarizes the statements above.

Theorem 2. Assume condition (L*) and quadratic Laplace approximation (A) with probability $1 - 2e^{-x}$ are fulfilled. Then with probability $1 - 4e^{-x}$ in the local region $\Theta_1(r) \cap \Theta_2(r)$ took place

$$\left| \sqrt{2T_h} - \|\xi_{12} + \theta_{12}^*\| \right| \leq 10\diamond(r, x).$$

where ξ_{12} and θ_{12}^* are defined in theorem 1.

Remark 3. The constant near $\diamond(r, x)$ could be decreased, expanding series of $L_1(\theta)$, $L_2(\theta)$ and $L(\theta)$ in the local regions around θ_1^* , θ_2^* and θ^* instead of MLE values:

$$\begin{aligned} 2T_h &= -\|\xi\|^2 + \|\xi_1\|^2 + \|\xi_2\|^2 - 2\xi_1^T D_1 D^{-2} D_2^2 (\theta_2^* - \theta_1^*) + 2\xi_2^T D_2 D^{-2} D_1^2 (\theta_2^* - \theta_1^*) \\ &\quad + \|D_1 D^{-2} D_2^2 (\theta_2^* - \theta_1^*)\|^2 + \|D_2 D^{-2} D_1^2 (\theta_2^* - \theta_1^*)\|^2 \pm (2\diamond(r, x)r + 2\delta(r)r^2) \\ &= -\|\xi\|^2 + \|\xi_1\|^2 + \|\xi_2\|^2 + 2(D_2^{-1}\xi_2 - D_1^{-1}\xi_1)^T \Sigma^2 (\theta_2^* - \theta_1^*) + \|\Sigma(\theta_2^* - \theta_1^*)\|^2 \\ &\quad \pm (2\diamond(r, x)r + 2\delta(r)r^2). \end{aligned}$$

Referring to condition A, $\|D^{-1}(D_1\xi_1 + D_2\xi_2)\|^2 \pm 2\langle(r, x)z(x)$ replaces $\|\xi\|^2$.

$$-\|\xi\|^2 + \|\xi_1\|^2 + \|\xi_2\|^2 = \|\Sigma(D_2^{-1}\xi_2 - D_1^{-1}\xi_1)\|^2 \pm 2\langle(r, x)z(x).$$

That leads to result

$$\left|2T_h - \|\xi_{12} + \theta_{12}^*\|^2\right| \leq (4\langle(r, x)r + 2\delta(r)r^2).$$

Remark 4. Weighted LRT statistic (Tb) has similar approximation:

$$2T_h^b \approx \|D(\widehat{\theta}^b - \widehat{\theta}_{H_0}^b)\|^2 = \|\xi_{12}^b\|^2.$$

where $\widehat{\theta}_{H_0}^b$ is argmax of L^b under condition $H_0 : \widehat{\theta}_2 - \widehat{\theta}_1 = \widehat{\theta}_{12}$, which is true. That's why the mean of difference $(\widehat{\theta}^b - \widehat{\theta}_{H_0}^b)$ is zero.

3.2. Optimal window size

The change point detection algorithm described above has rather meaningful parameter window size (h) that determines sample sizes on which MLE $(\widehat{\theta}_1, \widehat{\theta}_2)$ will be compared. One may find out the required sample size from condition

$$h\mathcal{KL}(\theta_1^*, \theta_2^*) > h\mathcal{KL}(\widehat{\theta}_1, \theta_1^*) + h\mathcal{KL}(\widehat{\theta}_2, \theta_2^*), \quad (1)$$

which ensures that distance between distributions $\mathcal{P}(\theta_1^*)$ and $\mathcal{P}(\theta_2^*)$ is greater than their fluctuations caused by error in θ_i^* estimation. Wilks theorem (reg. Spokoiny [2012]) gives upper approximation with probability $1 - 10e^{-x}$

$$h\mathcal{KL}(\widehat{\theta}_1, \theta_1^*) + h\mathcal{KL}(\widehat{\theta}_2, \theta_2^*) \leq 2r\langle(r, x) + \frac{\|\xi_1\|^2}{2} + \frac{\|\xi_2\|^2}{2},$$

where with probability $1 - 4e^{-x}$

$$\frac{\|\xi_1\|^2}{2} + \frac{\|\xi_2\|^2}{2} \leq z^2(x) = p_B + 6\lambda_B x, \quad p_B = \frac{p_{B_1} + p_{B_2}}{2}, \quad \lambda_B = \frac{\lambda_{B_1} + \lambda_{B_2}}{2}.$$

In case with

$$r\langle(r, x) = \sqrt{\frac{C(p_B + x)^3}{h}}, \quad h > C(p_B + x),$$

lower bound for parameter change from initial condition 1 results in estimation

$$h\mathcal{KL}(\theta_1^*, \theta_2^*) > 3p_B + (6\lambda_B + 2)x.$$

Optimal h is finite. Increasing a sample size one decreases an impact of stochastic part of $\|\xi_{12} + \theta_{12}^*\|$, since $\|\theta_{12}^*\|$ grows. But at the same time $\|\theta_{12}^*\|$ will not be changed by window replacement when $h \rightarrow \infty$.

The previous estimation of h doesn't take into account convolutions with patterns. Next reasoning estimates h from comparison of pattern convolution with statistic $T_h(i)$ with and without growing static term $\theta_{12}^*(i)$, where i indicates

sliding window position. Note that angle of $\|\theta_{12}^*(i)\|$ growth decreases with h . The optimal window size is the smallest one that is sufficient to overcome random fluctuations in convolution of $\|\xi_{12}(i) + \theta_{12}^*(i)\|$ with linear function $f(i) = i$ (corresponds to half of triangle pattern). Define new variables

$$b = \|\theta_{12}^*\| = \sqrt{h}b_0, \quad b_i = \frac{i}{h}b, \quad i > 0, \quad \xi_i = \xi_{12}(i).$$

Optimal window size for online change point detection is to be derived from the following inequality.

$$\sum_{i=1}^h i\|\xi_i + b_i\| \geq \sum_{i=1}^h i \left(\|\xi_i\| + 10\diamond(r, x) \right).$$

Theorem 4.1 from paper Spokoiny and Zhilova [2013] ensures following inequality with probability $1 - 2e^{-x}$

$$\begin{aligned} \|\xi_i + b_i\| &\geq \sqrt{\|\xi_i\|^2 + \|b_i\|^2 - 2\|b_i\|} - 2\delta_1(x) \geq \\ &\geq \|b_i\| - 2 - \sqrt{4 + 2\delta_1(x)}. \end{aligned}$$

With probability $1 - 4e^{-x}$ under condition that statement from theorem 2 is true one comes to a final estimation of the minimal sufficient window size:

$$h \geq \frac{9(2 + \sqrt{4 + 2\delta_1(x)} + z(x) + 10\diamond(r, x))^2}{4b_0^2} \sim \frac{c_1 + c_2p}{b_0^2}.$$

Consequently smaller h , in case data size is sufficient, is more preferable and has greater weight on offline mode. Due to b parameter is unknown in advance, the proposed multiscale method executes change point detection with different scale parameters $h \in H$.

4 Experiments

4.1. Experiments with synthetic data

This section presents results of the comparison of the proposed algorithm of change point detection (referred as *LRTOnline* or *LRTOffline*) with two other methods: *Bayesian online changepoint detection (BOCPD)* Adams and MacKay [2007] and *cpt.meanvar(PELT,...)* (RMeanVar) from RPa [2014]. The first method is constructed for online inference, but so far as it returns CP location with each CP signal, it is also applicable for offline testing scenario. The idea of this method is predictive filtering: its forecasts a new data point using only the information have been observed already, where the distribution family is fixed (Normal for the tests in this paper). Bayesian inference calculates the length of the observed data (from the last CP). The second algorithm also uses preliminary specified model. Its design focuses into finding multiple changes in mean and variance in Normally (another distributions also supported) distributed data. The returned set of change points is the result of sequential testing H_0 (existing number of change points) against H_1 (one extra change point) applying the likelihood ratio statistic of the whole data coupled with the penalty for CP count. Originally the method has offline change point detection interface, but one could adapt it for online case by buffering incoming data elements and clearing the buffer when at least one CP have been observed in the buffered data. RMeanVar performs better than well known method CUSUM due to synchronous changes in both data parameters mean and variance. In total, each of these two algorithms has modification in the way that allows one to use it in both online and offline testing mode.

LRTOffline configuration:

window sizes $(h_1, \dots, h_W) = (10, 20, 40, 70)$; confidence for the upper bound of convolution with pattern = 0.1; window weights $(u_1, \dots, u_W) = (1.0, 2.0, 0.5, 0.2)$.

LRTOnline configuration:

window sizes $(h_1, \dots, h_W) = (30, 50, 70)$; confidence = 0.1.

Quality of measurements uses three following metrics: Normalised Mutual Information (NMI), Delay (average time interval in which CP have been detected after it had taken place), Precision and Recall. The next equation defines NMI measure of two partitions (X, Y) of time range by change points

$$\text{NMI}(X, Y) = 2 \frac{H(X) + H(Y) - H(X, Y)}{H(X) + H(Y)}.$$

$H(X)$ and $H(X, Y)$ and entropy functions. Higher NMI values (they are in $[0, 1]$) correspond to better quality. Quality comparison in offline case apply NMI measure, while for online mode involves Delay, Precision and Recall.

Synthetic test data have been generated for different values of difference norm of the data distribution parameter. Such values are denoted as *delta*. Each delta corresponds to 10 sampled data sequences over which one compute measure average. In online mode each data sequence could have one or none change

points, in offline mode – two, one or none change points. The data has two distributions: Normal ($\mathcal{N}(\theta(1), \theta(2))$) and Poisson ($Po(\theta)$). Parametric assumption for all methods is $\mathcal{N}(\theta(1), \theta(2))$, so Poisson data corresponds to misspecification scenario.

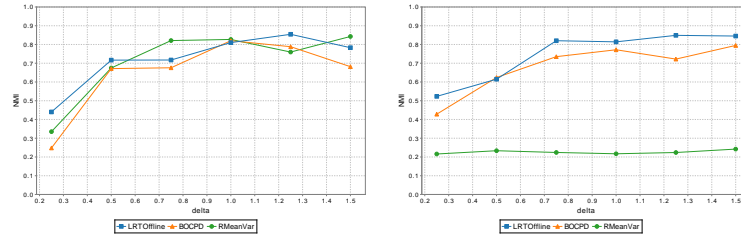


Fig. 2. Offline mode. First data: $\mathcal{N}(\theta(1), \theta(2))$, second data: $Po(\theta)$, $\delta = \|\theta_{12}^*\|$, data size = 340, parametric assumption for all methods is $\mathcal{N}(\theta(1), \theta(2))$, NMI – Normalized Mutual Information between predicted and reference partitions of time interval with change points, tests per $\delta = 10$, change point per test = $\{0, 1, 2\}$.

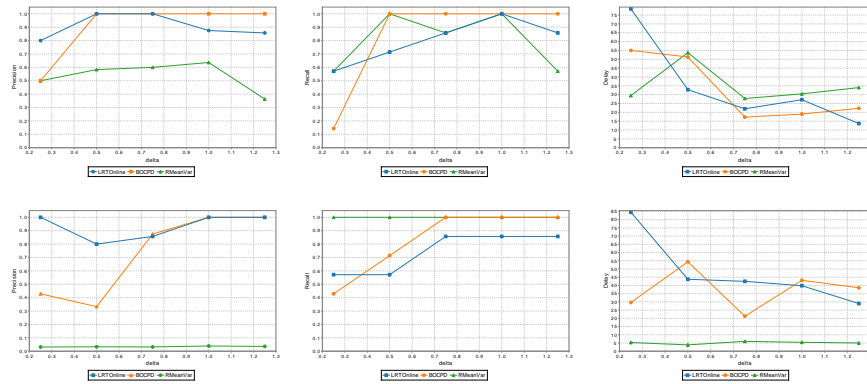


Fig. 3. Online mode. First row data: $\mathcal{N}(\theta(1), \theta(2))$, second row data: $Po(\theta)$, $\delta = \|\theta_{12}^*\|$, data size = 340, PA for all methods is $\mathcal{N}(\theta(1), \theta(2))$, tests per $\delta = 10$, change point per test = $\{0, 1\}$.

In the offline tests with Normal data all the methods achieves similar NMI scores, nonetheless LRTOOffline is more stable for decreasing strength of CP (δ). In the tests with Poisson data (misspecification) RMeanVar has relatively low quality. The online tests characterizes the proposed method (LR-

TOnline) as more stable along different delta values what is accomplished by multiscale heuristic.

The experiments reveal following meaningful properties of the proposed method configuration:

1. Quality is sensitive to selection of time interval $\tau \in [t_1, t_2]$ for upper bounds $(z(h), h \in H)$ calibration of convolution $\max_{\tau \in [t_1, t_2]} \widehat{\mathbb{T}}_h^b(\tau)$ in offline mode. For example in data $\mathcal{N}(0, 1).\text{sample}(100) \cup \mathcal{N}(1, 2).\text{sample}(100)$ is preferable to use only slice of 0 to 100 elements for calibration, because of lower $\text{Var } \xi_{12}$. Generally according to remark 1 from Section 3 one should run calibration in the range with the lowest possible values of $\text{tr}(B_1 + B_2)$. This improvement additionally requires approximation for the convolution maximum in expanding data ranges.
2. It is influenced to find out the minimal h sufficient for bootstrap usage. Bootstrap measure used for $z(h)$ calibration becomes closer to real measure with increasing h . Small h improves Delay (which also predicted theoretically in 3.2. Subsection) but makes unable to keep high level of Precision and Recall in online mode.

4.2. Experiments with real data

Here data from 1972-75 Dow Jones Returns Adams and MacKay [2007] describes several major events with potential macroeconomic effects (the most significant among them are the Watergate affair and the OPEC oil embargo). Convolutions plot $(\widehat{\mathbb{T}}_h^b(\tau)$ as function from $\tau)$ with its upper bounds $z(h)$ on this dataset appeared to be a nice illustration of multiscale detection importance: CP near $t = 325$ is better perceptible when window size is equal to 30 and CP near $t = 600$ has more perceptible convolution when window size is equal to 70. Two plots presented below includes convolutions with Static and Fitted Patterns, where one could remark better separability of convolution peaks in fitted case.

4.3. Sources

Demo of the LRTOonline method is available by link localcpdetector.shinyapps.io/localCP

Scala project with LRTOoffline and LRTOonline methods could be cloned from github.com/nazarblch/cpd which also includes testing system for abrupt change points detection applications and generated data used in the experiments.

5 Conclusion

This paper presented new change point detection method that works in offline and online modes. The method accounts properties of LRT statistic, which

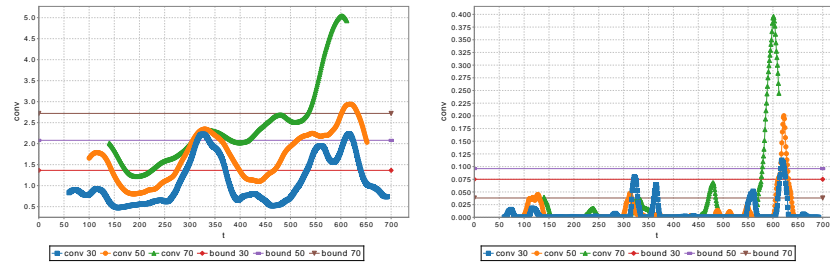


Fig. 4. Data: daily returns of the Dow Jones Industrial Average from July 3, 1972 to June 30, 1975. Left plot – convolution with static triangle pattern; right plot – convolution with fitted triangle pattern. The time axis is in business days, conv 30 (50, 70) corresponds to pattern with window size 30 (50, 70), bound 30 (50, 70) corresponds to convolution upper bound. Three reference CP are presented: the conviction of G. Gordon Liddy and James W. McCord, Jr. on January 30, 1973 ($t = 142$); the beginning of the OPEC embargo against the United States on October 19, 1973 ($t = 325$); the resignation of President Nixon on August 9, 1974 ($t = 548$).

has shifted χ^2 -distribution. Bootstrap technique appeared to be rather effective for LRT critical values calibration. Experiments and quality measurements confirm stability of the proposed algorithm in possibility to detect change points with different significance. The introduced concept of patterns allows to reveal different types of change point and filter regions with outliers.

Bibliography

- I. Nikiforov. A lower bound for the detection/isolation delay in a class of sequential tests. *IEEE Trans. Inf. Theory*, 49:3037–3047, 2003.
- M. Lavielle and G. Teyssiere. Detection of multiple change-points in multivariate time series. *Lithuanian Math. J.*, 46:287–306, 2006.
- A. Tartakovsky, B. Rozovskii, R. Blazek, and H. Kim. Detection of intrusions in information systems by sequential change-point methods. *Stat. Methodology*, 3:252–340, 2006.
- Pedro Casas, Sandrine Vaton, Lionel Fillatre, and Igor Nikiforov. Optimal volume anomaly detection and isolation in large-scale ip networks using coarse-grained measurements. *Computer Networks*, 54(11):1750–1766, 2010.
- Tze Leung Lai. Sequential changepoint detection in quality control and dynamical systems. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 613–658, 1995.
- R. Blazek and H. Kim. A novel approach to detection of denial-of-service attacks via adaptive sequential and batch-sequential change-point detection methods. In David Kurlander, Marc Brown, and Ramana Rao, editors, *Proc. of IEEE Workshop on Systems, Man, and Cybernetics Information Assurance*, pages 41–50. ACM Press, June 2001.
- Haining Wang, Danlu Zhang, and Kang G Shin. Change-point monitoring for the detection of dos attacks. *Dependable and Secure Computing, IEEE Transactions on*, 1(4):193–208, 2004.
- V. Spokoiny. Multiscale local change point detection with applications to value-at-risk. *Ann. of Stat.*, 2009.
- T. Mikosch and C. Starica. Changes of structure in financial time series and the garch model. *Econometrics* 0412003, EconWPA, 2004.
- Albert N Shiryaev. On optimum methods in quickest detection problems. *Theory of Probability & Its Applications*, 8(1):22–46, 1963.
- A. Polunchenko and A. Tartakovsky. State-of-the-art in sequential change-point detection. *Methodol. Comput. Appl. Probab.*, 14:649–684, 2011.
- A.N. Shiryaev. Quickest detection problems: Fifty years later. *Sequential Anal.: Design Methods and Applicat.*, 29:345–385, 2010.
- Walter Andrew Shewhart. *Economic control of quality of manufactured product*. ASQ Quality Press, 1931.
- Richard E Quandt. Tests of the hypothesis that a linear regression system obeys two separate regimes. *Journal of the American statistical Association*, 55(290): 324–330, 1960.
- Hyune-Ju Kim and David Siegmund. The likelihood ratio test for a change-point in simple linear regression. *Biometrika*, 76(3):409–423, 1989.
- Patsy Haccou, Evert Meelis, and Sara Van De Geer. The likelihood ratio test for the change point problem for exponentially distributed random variables. *Stochastic processes and their applications*, 27:121–139, 1987.

- MS Srivastava and Keith J Worsley. Likelihood ratio tests for a change in the multivariate normal mean. *Journal of the American Statistical Association*, 81(393):199–204, 1986.
- Yukun Liu, Changliang Zou, and Runchu Zhang. Empirical likelihood ratio test for a change-point in linear regression model. *Communications in Statistics—Theory and Methods*, 37(16):2551–2563, 2008.
- Changliang Zou, Yukun Liu, Peng Qin, and Zhaojun Wang. Empirical likelihood ratio test for the change-point problem. *Statistics & probability letters*, 77(4):374–382, 2007.
- J. Chen and A.K. Gupta. *Parametric Statistical Change Point Analysis: With Applications to Genetics, Medicine, and Finance*. Springer, 2012. ISBN 0817648003. URL <http://www.amazon.com/Parametric-Statistical-Change-Point-Analysis/dp/0817648003>.
- Edit Gombay. Sequential change-point detection with likelihood ratios. *Statistics & probability letters*, 49(2):195–204, 2000.
- BK Jandhyala and Stergios B Fotopoulos. Capturing the distributional behaviour of the maximum likelihood estimator of a changepoint. *Biometrika*, 86(1):129–140, 1999.
- Hyune-Ju Kim. Tests for a change-point in linear regression. *Lecture Notes-Monograph Series*, pages 170–176, 1994.
- Stergios B Fotopoulos, Venkata K Jandhyala, and Elena Khapalova. Exact asymptotic distribution of change-point mle for change in the mean of gaussian sequences. *The Annals of Applied Statistics*, pages 1081–1104, 2010.
- Jianqing Fan, Chunming Zhang, and Jian Zhang. Generalized likelihood ratio statistics and wilks phenomenon. *Annals of Statistics*, pages 153–193, 2001.
- Samuel S Wilks. The large-sample distribution of the likelihood ratio for testing composite hypotheses. *The Annals of Mathematical Statistics*, 9(1):60–62, 1938.
- Stéphane Boucheron and Pascal Massart. A high-dimensional wilks phenomenon. *Probability theory and related fields*, 150(3-4):405–433, 2011.
- V. Spokoiny. Penalized maximum likelihood estimation and effective dimension. *eprint arXiv:1205.0498*, 2012.
- Klaus Frick, Axel Munk, and Hannes Sieling. Multiscale change point inference. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(3):495–580, 2014.
- Siddhartha Chib. Estimation and comparison of multiple change-point models. *Journal of econometrics*, 86(2):221–241, 1998.
- V. Chernozhukov, D. Chetverikov, and K. Kato. Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *Ann. of Stat.*, 41:2786–2819, 2013.
- V. Spokoiny and M. Zhilova. Bootstrap confidence sets under a model misspecification. *Preprint no. 1992, WIAS*, 2014.
- V. Spokoiny and M. Zhilova. Sharp deviation bounds for quadratic forms. *Mathematical Methods of Statistics*, 22(2):100–113, 2013.
- R.P Adams and D.J.C. MacKay. Bayesian online changepoint detection. 2007. changepoint: An r package for changepoint analysis, 2014.