

DETECTING A CHANGE IN REGRESSION: FIRST-ORDER OPTIMALITY

BY BENJAMIN YAKIR,¹ ABBA M. KRIEGER AND MOSHE POLLAK²

Hebrew University and University of Pennsylvania

Observations are generated according to a regression with normal error as a function of time, when the process is in control. The process potentially changes at some unknown point of time and then the ensuing observations are normal with the same mean function plus an arbitrary function under suitable regularity conditions. The problem is to obtain a stopping rule that is optimal in the sense that the rule minimizes the expected delay in detecting a change subject to a constraint on the average run length to a false alarm. A bound on the expected delay is first obtained. It is then shown that the cusum and Shiriyayev–Roberts procedures achieve this bound to first order.

1. Introduction. A new product has entered the market recently (e.g., cellular phones). Industry sales are measured monthly. If the product catches on, one may anticipate a sharp and sustained increase in the growth of monthly sales. For a variety of reasons (e.g., production, inventory control and advertising strategy), it is important to determine as soon as possible that a marked increase has occurred in the growth of monthly sales. For similar reasons, it is costly to claim that such a change has occurred when it has not.

The problem can be modeled by a sequence of independent, homoscedastic, normally distributed observations Y_i (where i indexes time). Initially, the mean grows according to a function $\beta(i)$. If a change occurs at time ν , then the mean of Y_i , $i \geq \nu$, grows according to $\beta(i) + \gamma(i - \nu + 1)$.

In the first part of this paper, it is assumed that the variance σ^2 of Y_i is known and that $\beta(i)$ and $\gamma(i)$ are also known. (Without loss of generality, $\beta(i)$ can be assumed to be 0 and σ^2 can be assumed to be 1. Some regularity conditions on γ will be imposed.) In the latter part of the paper, the case of unknown parameters is studied.

In formal terms, a detection scheme is characterized by a stopping time N at which an alarm is raised. The problem addressed in this paper is that of minimizing the expected delay until detection, $E^{(\nu)}(N - \nu + 1 | N \geq \nu)$, subject to a lower bound on the average run length (ARL) to false alarm $E^{(\infty)}N$. The expectation $E^{(\infty)}$ corresponds to $\nu = \infty$, or, in other words, the case of no change.

Received December 1996; revised May 1999.

¹Supported by a grant from the Israel Science Foundation.

²Supported by a grant from the Israel Science Foundation.

AMS 1991 subject classifications. 62L10, 62N10.

Key words and phrases. Change point detection, regression, stopping rules, information bound.

Much has been written on this problem in the context of a change from one fixed distribution to another fixed distribution. In that context, Lorden (1971) found an asymptotic lower bound on the expected delay to detection as a function of the ARL to false alarm (as the latter tends to ∞). Other and more precise optimality results were obtained later [Pollak (1985), Pollak and Siegmund (1985), Moustakides (1986), Ritov (1990), Lai (1995), Yakir (1997)]. Yao (1993) extended Lorden’s results to detecting a change in regression in the case of bounded information per observation (i.e., essentially $\sup_{1 \leq i < \infty} E^{(1)} \log[dP^{(1)}(Y_i)/dP^{(\infty)}(Y_i)] < \infty$). Here we study the case of possibly unbounded information, which includes the problem of detecting a change of a slope with respect to time.

There are two basic approaches to proving optimality statements. One approach is Bayesian or decision theoretic [cf. Pollak (1985), Ritov (1990), Yakir (1997)]. The other approach is classical and uses the theory of optimal stopping [cf. Lorden (1971), Moustakides (1986), Yao (1993), Lai (1995)]. It is this latter approach that is employed here, though our method of proof is different. [See also Lai (1993) and Yakir (1996).]

This paper has three main results. In Section 2 a lower bound on the expected delay is developed. In Section 3, it is shown that, asymptotically, the cusum and the Shirayayev–Roberts procedures achieve this bound and are therefore asymptotically optimal (to first order). In Section 4, a procedure which asymptotically attains these bounds is developed for the more practical case where baseline and post-change parameters are unknown. The paper is concluded with remarks and a discussion of extensions.

2. A lower bound. Our main concern in this section is the investigation of the optimal rate of detection for various post-change structures, and it will be assumed that it is known that $\sigma^2 = 1$ and $\beta(\cdot) \equiv 0$. The distribution of the sequence of observations is denoted by $P^{(\nu)}$. If $\nu = k$, then

$$Y_i \sim N(\gamma(i - k + 1) \mathbb{1}(i \geq k), 1)$$

for some given function $\gamma(\cdot)$.

Denote by $l^{(k)}(i)$ the log-likelihood ratio of the observation Y_i for the $P^{(k)}$ relative to the $P^{(\infty)}$ measures when $\nu = k$. It follows that if $i < k$, then $l^{(k)}(i) = 0$, and if $i \geq k$, then

$$l^{(k)}(i) = \gamma(i - k + 1)Y_i - \gamma^2(i - k + 1)/2.$$

Let $S^{(k)}(n) = \sum_{i=k}^n l^{(k)}(i)$ be the log-likelihood ratio of the first n observations.

Let N be any change point detection policy. In the following theorem the rate of detection of N , $E^{(k)}(N - k + 1 \mid N \geq k)$, is related to the information on $\gamma(\cdot)$ carried by the sequence of observations. The Kullback–Leibler information function is $E^{(k)}S^{(k)}(n) = \sum_{i=k}^n \gamma^2(i - k + 1)/2$. In particular, when $k = 1$, the information becomes

$$I(n) = \sum_{i=1}^n \frac{\gamma^2(i)}{2}.$$

For the sake of convenience, regard $I(x)$ as a piecewise linear continuous function of real x . Moreover, for the sake of clarity of exposition, we will assume that $I(x)$ is strictly increasing in x . We show that if $I(n)$ behaves as a power function in n and if the ARL to false alarm of N is not less than A , then the rate of detection is bounded from below by the inverse of the information function, calculated at $\log A$.

THEOREM 1. *Assume that $I(n)$ is of the form $L(n)n^r$, for some $r > 0$ and some slowly changing function L . Then*

$$(2.1) \quad \inf_{\{N: \mathbf{E}^{(\infty)}N \geq A\}} \sup_{1 \leq k < \infty} \mathbf{E}^{(k)}(N - k + 1 \mid N \geq k) \geq I^{-1}(\log A)(1 + o(1)),$$

where $o(1) \rightarrow 0$ as $A \rightarrow \infty$.

Note that $I(n)$ must be monotone in n . Recall that $L(n)$ is slowly changing if $\lim_{n \rightarrow \infty} L(bn)/L(n) = 1$, for all $b > 0$.

A key ingredient in the proof of the theorem is a lemma that relates the distribution of the detection policy N to the distribution of the stopping time of a one-sided sequential probability ratio test (SPRT). Given a boundary $b > 0$, the one-sided SPRT of $H_0: \nu = \infty$ versus $H_1: \nu = k$ is defined by

$$M_b^{(k)} = \inf\{n: S^{(k)}(n) \geq b\}.$$

LEMMA 1. *Let A, c and t be positive numbers such that $a = \log A > 2c$. If N is a stopping time for which $\mathbf{E}^{(\infty)}N \geq A$, then there exists an integer k such that*

$$\mathbf{P}^{(k)}(N - k + 1 < M_{a-2c}^{(k)} - k + 1 \mid N \geq k) \leq \mathbf{P}^{(1)}(M_{a-2c}^{(1)} > t) + (8t + 9)e^{-c}.$$

PROOF. Define the (truncated) Shiriyayev–Roberts statistics by

$$R(n; t, m) = \sum_{k=(n-t) \vee m}^n \exp\{S^{(k)}(n)\}.$$

For any $j \geq 0$,

$$\begin{aligned} \mathbf{P}^{(\infty)}(N \geq jA/2) &= \sum_{n=jA/2}^{\infty} \mathbf{P}^{(\infty)}(N = n) \\ &= \sum_{n=jA/2}^{\infty} \sum_{k=(n-t) \vee jA/2}^n \mathbf{E}^{(k)} \left[\frac{\mathbb{1}(N = n)}{R(n; t, jA/2)} \right] \\ &= \sum_{k=jA/2}^{\infty} \sum_{n=k}^{k+t} \mathbf{E}^{(k)} \left[\frac{\mathbb{1}(N = n)}{R(n; t, jA/2)} \right] \\ &= \sum_{k=jA/2}^{\infty} \mathbf{E}^{(k)} \left[\frac{\mathbb{1}(k \leq N \leq k+t)}{R(N; t, jA/2)} \right] \end{aligned}$$

$$\begin{aligned} &\geq \sum_{k=jA/2}^{(j+1)A/2} \mathbf{E}^{(k)} \left[\frac{\mathbb{1}(N \leq k+t)}{R(N; t, jA/2)} \mid N \geq k \right] \mathbf{P}^{(\infty)}(N \geq k) \\ &\geq \mathbf{P}^{(\infty)}(N \geq (j+1)A/2) \\ &\quad \times \sum_{k=jA/2}^{(j+1)A/2} \mathbf{E}^{(k)} \left[\frac{\mathbb{1}(N \leq k+t)}{R(N; t, jA/2)} \mid N \geq k \right]. \end{aligned}$$

Define, for A , the given N and all $0 \leq j < \infty$, the conditional probability

$$\gamma_j = \mathbf{P}^{(\infty)}(N \geq (j+1)A/2 \mid N \geq jA/2) = \frac{\mathbf{P}^{(\infty)}(N \geq (j+1)A/2)}{\mathbf{P}^{(\infty)}(N \geq jA/2)}.$$

It follows that

$$(2.2) \quad \frac{1}{\gamma_j} \geq \frac{A}{2} \min_{jA/2 \leq k \leq (j+1)A/2} \mathbf{E}^{(k)} \left[\frac{\mathbb{1}(N \leq k+t)}{R(N; t, jA/2)} \mid N \geq k \right].$$

Since $A \leq \mathbf{E}^{(\infty)} N \leq [A/2] \sum_{n=0}^{\infty} \prod_{j=0}^{n-1} \gamma_j$ (with $\prod_{j=0}^{-1} \gamma_j = 1$) it can be concluded that there must exist j for which $\gamma_j \geq 1/2$. Hence, there must exist $k, jA/2 \leq k \leq (j+1)A/2$, such that

$$(2.3) \quad \mathbf{E}^{(k)} \left[\frac{\mathbb{1}(N \leq k+t)}{R(N; t, jA/2)} \mid N \geq k \right] \leq 4/A.$$

In the rest of the proof we fix these j and k .

Let $c > 0$ be given and consider the log-likelihood ratio $S^{(k)}(N)$. It is straightforward to show that

$$(2.4) \quad \begin{aligned} \mathbf{P}^{(k)}(S^{(k)}(N) < \log A - 2c; N \leq k+t \mid N \geq k) \\ \leq \mathbf{P}^{(k)}(R(N; t, jA/2) \leq e^{-c} A/4; N \leq k+t \mid N \geq k) \end{aligned}$$

$$(2.5) \quad \begin{aligned} + \mathbf{P}^{(k)}(R(N; t, jA/2) \exp(-S^{(k)}(N)) \\ \geq e^c/4; N \leq k+t \mid N \geq k). \end{aligned}$$

Inequality (2.3) can be used to show for (2.4) that

$$\mathbf{P}^{(k)}(R(N; t, jA/2) \leq e^{-c} A/4; N \leq k+t \mid N \geq k) \leq e^{-c}.$$

In order to bound the term in (2.5) notice that over the event $\{k \leq N \leq k+t\}$,

$$R(N; t, jA/2) \exp(-S^{(k)}(N)) \leq \max_{k \leq n \leq k+t} \sum_{i=(k-t) \vee jA/2}^n \exp(S^{(i)}(n) - S^{(k)}(n)).$$

Hence,

$$\begin{aligned} \mathbf{P}^{(k)}(R(N; t, jA/2) \exp(-S^{(k)}(N)) \geq e^c/4; N \leq k+t \mid N \geq k) \\ \leq \mathbf{P}^{(k)} \left(\max_{k \leq n \leq k+t} \sum_{i=(k-t) \vee jA/2}^n \exp(S^{(i)}(n) - S^{(k)}(n)) \geq e^c/4 \mid N \geq k \right) \end{aligned}$$

$$\begin{aligned} &\leq \mathbb{P}^{(k)}\left(\max_{k \leq n \leq k+t} \sum_{i=(k-t) \vee jA/2}^n \exp(S^{(i)}(n) - S^{(k)}(n)) \geq e^c/4 \mid N \geq jA/2\right) \\ &\quad \times \mathbb{P}^{(\infty)}(N \geq jA/2) / \mathbb{P}^{(\infty)}(N \geq (j+1)A/2) \\ &= \frac{1}{\gamma_j} \mathbb{P}^{(k)}\left(\max_{k \leq n \leq k+t} \sum_{i=(k-t) \vee jA/2}^n \exp(S^{(i)}(n) - S^{(k)}(n)) \geq e^c/4\right) \\ &\leq 8(t+1)e^{-c}, \end{aligned}$$

since $\gamma_j \geq 1/2$ and since Doob's inequality can be applied to the $\mathbb{P}^{(k)}$ -submartingale

$$\sum_{i=(k-t) \vee jA}^n \exp(S^{(i)}(n) - S^{(k)}(n)), \quad k \leq n \leq k+t.$$

To conclude the above discussion,

$$(2.6) \quad \mathbb{P}^{(k)}(S^{(k)}(N) < a - 2c; N \leq k+t \mid N \geq k) \leq (8t+9)e^{-c}.$$

It follows from (2.6) that

$$\mathbb{P}^{(k)}(N < M_{a-2c}^{(k)}; N \leq k+t \mid N \geq k) \leq (8t+9)e^{-c}.$$

Therefore,

$$\mathbb{P}^{(k)}(N < M_{a-2c}^{(k)}; M_{a-2c}^{(k)} \leq k+t \mid N \geq k) \leq (8t+9)e^{-c}.$$

The event $\{N \geq k\}$ is independent of the stopping time $M_{a-2c}^{(k)}$ and the $\mathbb{P}^{(k)}$ -distribution of $M_{a-2c}^{(k)} - k + 1$ is independent of k . Hence

$$\begin{aligned} \mathbb{P}^{(k)}(N < M_{a-2c}^{(k)} \mid N \geq k) &\leq \mathbb{P}^{(k)}(M_{a-2c}^{(k)} - k + 1 > t) + (8t+9)e^{-c} \\ &= \mathbb{P}^{(1)}(M_{a-2c}^{(1)} > t) + (8t+9)e^{-c}. \quad \square \end{aligned}$$

PROOF OF THEOREM 1. Let I^{-1} be the inverse function of I . Observe first that for b fixed, $b > 0$, and for a , $a \rightarrow \infty$,

$$I(bI^{-1}(a)) = ab^r \frac{L(bI^{-1}(a))}{L(I^{-1}(a))} \sim ab^r,$$

since L is slowly changing. It follows that $I^{-1}(ba) \sim b^{1/r} I^{-1}(a)$.

Let $\varepsilon > 0$ be given. It will be shown that for any A large,

$$(2.7) \quad \inf_{\{N: \mathbb{E}^{(\infty)}N \geq A\}} \sup_{1 \leq \nu < \infty} \mathbb{E}^{(\nu)}(N - \nu + 1 \mid N \geq \nu) \geq e^{-\varepsilon} I^{-1}(\log A).$$

Define $t = \lceil I^{-1}(\log A) \rceil = \lceil I^{-1}(a) \rceil$, and for a small, but positive ε_1 , define $c = \varepsilon_1 a$. It follows from Lemma 1 and from Chebyshev's inequality that for

any stopping rule N , for which $E^{(\infty)}N \geq A$, there exists an integer k such that

$$\begin{aligned} &P^{(k)}(N - k + 1 < M_{(1-2\varepsilon_1)a}^{(k)} - k + 1 \mid N \geq k) \\ &\leq P^{(1)}(M_{(1-2\varepsilon_1)a}^{(1)} > t) + (8t + 9) \exp\{-\varepsilon_1 a\} \\ &\leq P^{(1)}(S^{(1)}(t) < (1 - 2\varepsilon_1)a) + (8I^{-1}(a) + 9) \exp\{-\varepsilon_1 a\} \\ &\leq \frac{1}{2\varepsilon_1^2 a} + (8I^{-1}(a) + 9) \exp\{-\varepsilon_1 a\}, \end{aligned}$$

since $E^{(1)}(S^{(1)}(t)) = I(t) \sim a$ and $\text{Var}^{(1)}(S^{(1)}(t)) = 2I(t)$. This expression converges to zero as $a \rightarrow \infty$.

Moreover,

$$\begin{aligned} &P^{(k)}(M_{(1-2\varepsilon_1)a}^{(k)} - k + 1 < \exp(-\varepsilon_1)I^{-1}((1 - 2\varepsilon_1)a)) \\ &\leq P^{(1)}\left(\sup_{n \leq \exp(-\varepsilon_1)I^{-1}((1-2\varepsilon_1)a)} S^{(1)}(n) > (1 - 2\varepsilon_1)a\right). \end{aligned}$$

However, from Doob's inequality, the monotonicity of $I(\cdot)$ and for a_1 , such that $a_1 = (1 - 2\varepsilon_1)a$,

$$\begin{aligned} &P^{(1)}\left(\sup_{n \leq \exp(-\varepsilon_1)I^{-1}(a_1)} S^{(1)}(n) > a_1\right) \\ &\leq P^{(1)}\left(\sup_{n \leq \exp(-\varepsilon_1)I^{-1}(a_1)} (S^{(1)}(n) - I(n))^2 > (a_1 - I(\exp(-\varepsilon_1)I^{-1}(a_1)))^2\right) \\ &\leq \frac{2I(\exp(-\varepsilon_1)I^{-1}(a_1))}{(a_1 - I(\exp(-\varepsilon_1)I^{-1}(a_1)))^2}, \end{aligned}$$

which, again, converges to zero as $a \rightarrow \infty$.

Therefore,

$$P^{(k)}(N - k + 1 \geq \exp(-\varepsilon_1)I^{-1}((1 - 2\varepsilon_1)a) \mid N \geq k) \geq \exp(-\varepsilon/2),$$

provided that a is large. A choice of a small enough ε_1 would lead to inequality (2.7) since $\exp(-\varepsilon_1)I^{-1}((1 - 2\varepsilon_1)a) \sim \exp(-\varepsilon_1)(1 - 2\varepsilon_1)^{1/r}I^{-1}(a)$. The proof of the theorem thus follows. \square

3. Asymptotically optimal detection policies. In this section we relax the assumption that γ is known and consider the construction of asymptotically optimal detection stopping times. Natural candidates, such as the cusum or the Shiriyayev–Roberts procedure, are based on log-likelihood ratios. These log-likelihood ratios, however, involve unknowns—the function γ —which can be estimated. Martingale consideration would suggest estimation which is based only on observations which are prior to the current one. A simple estimate of $\gamma(i - k + 1)$, which is the $P^{(k)}$ -mean of the observation Y_i , is Y_{i-1} (or

0 if $i = k$). Substituting estimates leads to the (approximated) log-likelihood ratios

$$\tilde{S}^{(k)}(n) = \sum_{i=k+1}^n (Y_{i-1}Y_i - Y_{i-1}^2/2).$$

Note that $\tilde{S}^{(k)}(n)$ is itself a log-likelihood ratio.

Consider the Shiriyayev–Roberts stopping rule

$$N_{\text{SR}} = \inf \left\{ n: \sum_{k=1}^n \exp\{\tilde{S}^{(k)}(n)\} \geq A \right\},$$

and the cusum stopping rule

$$N_{\text{CS}} = \inf \left\{ n: \max_{k \leq n} \tilde{S}^{(k)}(n) \geq \log A \right\}.$$

THEOREM 2. *Let N be either N_{SR} or N_{CS} .*

(i) $\mathbf{E}^{(\infty)}N \geq A$.

(ii) *If $\gamma(\cdot)$ is a positive and increasing function and if $I(n) = L(n)n^r$ for some $r > 1$ and some slowly changing function L , then*

$$\sup_{1 \leq k < \infty} \mathbf{E}^{(k)}(N - k + 1 \mid N \geq k) = I^{-1}(\log A)(1 + o(1)),$$

where $o(1) \rightarrow 0$ as $A \rightarrow \infty$.

REMARK. Note that $\gamma(n) \rightarrow \infty$, $\gamma(n)/\gamma(n-1) \rightarrow 1$ and $Y_n/\gamma(n-k+1) \rightarrow_p 1$ under $\mathbf{P}^{(k)}$ as $n \rightarrow \infty$, providing for intuitive plausibility of the success of N_{SR} and N_{CS} .

PROOF OF THEOREM 2(i). The stopping time N_{CS} dominates the stopping time N_{SR} . Hence it is sufficient to prove the claim for $N = N_{\text{SR}}$.

If $\lim_{m \rightarrow \infty} \mathbf{P}^{(\infty)}(N < m) < 1$, then $\mathbf{E}^{(\infty)}N = \infty$ and the claim is trivial. Assume that the probability converges to 1. Obviously, $\mathbf{E}^{(\infty)}N \geq \mathbf{E}^{(\infty)}(N \wedge m)$, for $1 \leq m < \infty$. The process $\sum_{k=1}^n \exp\{\tilde{S}^{(k)}(n)\} - n$ is a $\mathbf{P}^{(\infty)}$ -martingale. Hence, by the optional sampling theorem,

$$\mathbf{E}^{(\infty)}(N \wedge m) = \mathbf{E}^{(\infty)} \sum_{k=1}^{N \wedge m} \exp\{\tilde{S}^{(k)}(N \wedge m)\}.$$

The statistic $\sum_{k=1}^{N \wedge m} \exp\{\tilde{S}^{(k)}(N \wedge m)\}$ is positive and is larger than A on the event $\{N < m\}$. Therefore,

$$\mathbf{E}^{(\infty)}N \geq A\mathbf{P}^{(\infty)}(N < m).$$

The proof follows since $\mathbf{P}^{(\infty)}(N < m) \rightarrow_{m \rightarrow \infty} 1$ by assumption. \square

PROOF OF THEOREM 2(ii). Let N be either N_{SR} or N_{CS} . For any $k, 1 \leq k < \infty$, let $M_{\log A}^{(k)}$ be the SPRT stopping time as in the previous section with $S^{(k)}(n)$ replaced by $\tilde{S}^{(k)}(n)$. It is easy to see that

$$\begin{aligned} \mathbf{E}^{(k)}(N - k + 1 \mid N \geq k) &\leq \mathbf{E}^{(k)}(M_{\log A}^{(k)} - k + 1 \mid N \geq k) \\ &= \mathbf{E}^{(1)}M_{\log A}^{(1)}. \end{aligned}$$

In the sequel we bound this last expectation.

From the definition of $M_{\log A}^{(1)} = M$ it follows that

$$\begin{aligned} (3.1) \quad \log A &\geq \sum_{i=2}^{M-1} (Y_{i-1}Y_i - Y_{i-1}^2/2) \\ &= I(M - 1) + \sum_{i=2}^{M-1} (Y_{i-1}Y_i - Y_{i-1}^2/2 - \gamma^2(i)/2) - \frac{\gamma^2(1)}{2}. \end{aligned}$$

The expectation of the second term in the last line of (3.1) can be bounded from below by

$$\begin{aligned} &\mathbf{E}^{(1)} \sum_{i=2}^{M-1} (Y_{i-1}Y_i - Y_{i-1}^2/2 - \gamma^2(i)/2) \\ &= \mathbf{E}^{(1)} \sum_{i=2}^{\infty} (Y_{i-1}Y_i - Y_{i-1}^2/2 - \gamma^2(i)/2) \mathbb{1}(M \geq i + 1) \\ &\geq -\mathbf{E}^{(1)} \sum_{i=2}^{\infty} |Y_{i-1}Y_i - Y_{i-1}^2/2 - \gamma^2(i)/2| \mathbb{1}(M \geq i + 1) \\ &\geq -\mathbf{E}^{(1)} \sum_{i=2}^{\infty} |Y_{i-1}Y_i - Y_{i-1}^2/2 - \gamma^2(i)/2| \mathbb{1}(M \geq i - 1) \\ &= -\sum_{i=2}^{\infty} \mathbf{E}^{(1)} |Y_{i-1}Y_i - Y_{i-1}^2/2 - \gamma^2(i)/2| \mathbf{P}^{(1)}(M \geq i - 1) \\ &= -\mathbf{E}^{(1)} \sum_{i=2}^{M+1} \mathbf{E}^{(1)} |Y_{i-1}Y_i - Y_{i-1}^2/2 - \gamma^2(i)/2|. \end{aligned}$$

Straightforward calculations and the assumptions of the theorem show that $\gamma^2(n)/2 = o(I(n))$, implying $I(n - 1) = (1 + o(1))I(n)$, and that

$$\begin{aligned} &|Y_{i-1}Y_i - Y_{i-1}^2/2 - \gamma^2(i)/2| \\ &\leq |Y_{i-1} - \gamma(i - 1)| |Y_i - \gamma(i)| + \gamma(i) |Y_i - \gamma(i)| \\ &\quad + (Y_i - \gamma(i - 1))^2/2 + (\gamma^2(i) - \gamma^2(i - 1))/2, \end{aligned}$$

so that

$$\sum_{i=2}^{n+1} \mathbf{E}^{(1)} |Y_{i-1}Y_i - Y_{i-1}^2/2 - \gamma^2(i)/2| = o(I(n)).$$

Therefore, for any $\varepsilon > 0$ a constant C can be found such that

$$I(n - 1) - \sum_{i=2}^{n+1} \mathbf{E}^{(1)} \left| Y_{i-1} Y_i - Y_{i-1}^2/2 - \gamma^2(i)/2 \right| - \frac{\gamma^2(1)}{2} \geq (1 - \varepsilon)I(n) - C,$$

for all $n \geq 1$. Moreover, the function $I(n)$ is convex since the function $\gamma(\cdot)$ is increasing. It can be concluded that

$$\begin{aligned} \log A &\geq (1 - \varepsilon)\mathbf{E}^{(1)}I(M) - C \\ &\geq (1 - \varepsilon)I(\mathbf{E}^{(1)}M) - C. \end{aligned}$$

The proof of the theorem thus follows. \square

4. A case of unknown baseline parameters. Often, β and σ^2 are unknown and need to be estimated from the data. In this section, we construct a first-order asymptotically optimal procedure for this case. In order for the change point problem to be well defined, β should be structured, since γ is not. In this section we assume that $\beta(i) = \beta_0 + i\beta_1$ for some unknown scalars β_0 and β_1 . Given the first n observations, the standard estimators of $(\beta_0, \beta_1, \sigma^2)$ are (when all observations are prechange)

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_{i=1}^n iY_i - ((n + 1)/2) \sum_{i=1}^n Y_i}{n(n^2 - 1)/12} \\ \hat{\beta}_0 &= \frac{\sum_{i=1}^n Y_i}{n} - \hat{\beta}_1 \frac{n + 1}{2} \\ \hat{\sigma}^2 &= \frac{1}{n - 2} \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 i)^2. \end{aligned}$$

Note that $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$ and $\hat{\sigma}^2$ implicitly depend on n , the number of observations at the current inspection period. Define $\hat{l}^{(k)}(i, n) = \hat{Y}_{i-1, n} \hat{Y}_{i, n} - \hat{Y}_{i-1, n}^2/2$, for all $3 < k < i \leq n$, with $\hat{Y}_{i, n} = (Y_i - \hat{\beta}_0 - \hat{\beta}_1 i)/\hat{\sigma}$. Consider the pseudo-log-likelihood ratio

$$(4.1) \quad \hat{S}^{(k)}(n) = \sum_{i=k+1}^n \hat{l}^{(k)}(i, n),$$

and define the (window truncated) Shiriyayev–Roberts stopping rule,

$$N_{\text{SR}} = \inf \left\{ n: \sum_{k=n-t}^n \exp\{\hat{S}^{(k)}(n)\} \geq dA \right\},$$

and the (window truncated) cusum stopping rule,

$$N_{\text{CS}} = \inf \left\{ n: \max_{n-t \leq k \leq n} \hat{S}^{(k)}(n) \geq \log(dA) \right\},$$

for some $t = t(A)$, $t = o(A)$, $t/\log A \rightarrow \infty$, and $d > 2$. This d , together with $c > 1$, solves the equation

$$A = \frac{1/2 + [1/4 - (2t)/(dA)]^{1/2}}{1 - \{(cA - 1)/(cA)\} \{1/2 + [1/4 - (2t)/(dA)]^{1/2}\}^{1/t}}.$$

Note that $c = c(A)$ and $d = d(A)$ can be taken to be bounded. Let G be a geometric random variable, independent of everything else, such that $E(G) = cA$ and denote $N = N_{SR} \wedge G$ or $N = N_{CS} \wedge G$.

It will be shown below that, under mild regularity conditions, N attains the optimal rate of detection, at least when the change does not occur too soon after initiation of the monitoring. (This restriction is unavoidable; for example, if the change is a change of slope, the state {no change} is indistinguishable from the state {change from the very beginning}.) Before considering the rate of detection of a policy N , however, one needs to demonstrate that $E^{(\infty)}N \geq A$. It should be noted that $\hat{S}^{(k)}(n)$ is no longer a log-likelihood ratio, hence standard martingale results cannot be applied. Nonetheless, as shown in the proof of the following lemma, $\hat{S}^{(k)}(n)$ is smaller than another (true) log-likelihood ratio. It follows that the expectation constraint on the rate of false alarms is satisfied for the proposed procedure.

LEMMA 2. *Let N be defined as above. Then $E^{(\infty)}N \geq A$.*

The optimality claim is stated next as a theorem. The proof of both this theorem and of Lemma 2 is relegated to the Appendix.

THEOREM 3. *Let N be defined as above. If $I(n) = L(n)n^r$ for some $r > 1$ and some slowly changing function L and if $\sum_{i=1}^n [\gamma(i) - \gamma(i - 1)]^2 = o(I(n))$, then there exists a constant η such that*

$$\sup_{\eta \log A \leq k < \infty} E^{(k)}(N - k + 1 | N \geq k) = I^{-1}(\log A)(1 + o(1)).$$

REMARK. Note that one need not know the value of r to apply this procedure and that this procedure is first-order asymptotically optimal whatever the value of $r > 1$ is.

5. Comments and extensions. In the first part of this paper, the first-order optimality of the cusum and the Shirayayev–Roberts procedures was established. This result is in a regression context where the observations are independent and normally distributed with known variance and the mean level grows according to a known function of time before a change has taken place and according to another unknown function once a change has occurred. It is interesting to note that an optimal procedure is still based on the likelihood ratio’s crossing of a constant boundary. We conjecture that this may be the case in a wide variety of contexts, as hinted by the following considerations.

Regard a sequential hypothesis testing problem of a simple H_0 versus a simple H_1 based on a sequence of (not necessarily independent) observations

Y_1, Y_2, \dots . Without loss of generality, regard procedures based on L_1, L_2, \dots , where L_n is the likelihood ratio of the first n observations. Note that L_n equals the likelihood ratio of (L_1, \dots, L_n) . Letting f_n be the density of L_n , one can rephrase H_0 and H_1 as $H_0: f_n = f_n^0; n = 1, 2, \dots$ and $H_1: f_n = f_n^1; n = 1, 2, \dots$. Note that $f_n^1(y)/f_n^0(y) = y$. For a power one test, the operating characteristics of interest are $\alpha = P_0(N < \infty)$ and $E_1 N$, where N is a stopping time at which (if $N < \infty$) one stops and rejects H_0 . The optimal stopping problem is to minimize $E_1 N$ subject to a given α .

Consider stopping times of the form $N = \min\{n: L_n \geq C_n\}$ ($N = \infty$ if no such n exists), where $\{C_n\}$ is a sequence of constants. Regard the surrogate problem of minimizing the H_1 -expected number of times L_n is below C_n , subject to a fixed value δ of the H_0 -expected number of visits of L_n above C_n . Intuitively, the H_1 -expected number of times L_n is below C_n differs from $E_1 N$ by an additive constant, and δ differs from α by a multiplicative constant, so that an optimal procedure for the surrogate problem is almost (asymptotically, $\alpha \rightarrow 0$) optimal for the original one. Now use a Lagrangian multiplier argument to solve the surrogate problem,

$$\frac{\partial}{\partial C_j} \left[\sum_{n=1}^{\infty} P_1(L_n \leq C_n) + \lambda \left(\sum_{n=1}^{\infty} P_0(L_n > C_n) - \delta \right) \right] = f_j^1(C_j) - \lambda f_j^0(C_j)$$

which equal zero iff $f_n^1(C_n)/f_n^0(C_n) = \lambda$ for each n . Since $f_n^1(C_n)/f_n^0(C_n) = C_n$, it follows that the rule which calls for stopping the first time that the likelihood ratio exceeds a constant is optimal for the surrogate problem. The close relationship between power one tests and change point problems leads one to conjecture that the analogous rule (cusum, or Shiriyayev–Roberts) is asymptotically almost optimal. [Actually, these considerations lead one to suspect that the optimality is much stronger than that claimed in Sections 2 and 3; $I^{-1}(\log A)(1 + o(1))$ may perhaps turn out to be $I^{-1}(\log A) + O(1)$. However, we have so far been unable to provide a full proof along these lines.]

The results of this paper apply to cases where the information contained in the first n observations increases according to $L(n)n^r$. Note that a change in a linear trend can be related to $r = 3$. The example of the success of a new product (mentioned in the introduction) may correspond to $r > 3$. The classical change point problem (change from one constant mean to another constant mean level) corresponds to $r = 1$ and the case $r < 1$ corresponds to a situation where the mean tends to revert back to its original prechange level (at not too fast a rate), and the results of Section 2 apply, although those of Section 3 and 4 do not. The results of this paper, however, do not apply to regression problems where the information accumulates exponentially.

Parallel results can be obtained when the errors are not normally distributed. It is important to note that the main lemma in Section 2 is not based on normality. In this context it is useful to compare our results with the results in Robbins and Zhang (1993). They considered a change point detection problem in an exponential family context with a probability constraint imposed on the rate of false alarm. Mixture-type stopping rules were inves-

tigated in terms of their ARL to detection. It can be shown, however, that mixture-type stopping rule are suboptimal when $r > 1$.

APPENDIX

This Appendix contains the proofs of Lemma 2 and Theorem 3.

The basic idea in proving Lemma 2 is to show that $\hat{S}^{(k)}(n)$ is dominated by a true log-likelihood. This is done in five steps:

1. A true log-likelihood ratio, which involves invariant statistics, is defined. The likelihood ratio is formed by dividing the joint density of these statistics under $P^{(k)}$ by their joint density under $P^{(\infty)}$.
2. It is shown that the true log-likelihood ratio is greater than a computable expression.
3. The above computable expression is presented as a sum of $\hat{S}^{(k)}(n)$ and an additional term.
4. The additional term in (3) is proven to be positive.
5. The stopping time involving the true log-likelihood ratio of (1) is shown to satisfy the constraint on false alarms.

Step 1. Initially, let $Y_i, i \geq 1$, be independent homoscedastic normally distributed observations. It will be convenient to use vector and matrix notation. Let m be the number of observations that is known to be before the change has taken place. (Note that in the paper m is assumed to be 3.) Let

$$\mathbf{Y}' = (\mathbf{Y}_{-1}, \mathbf{Y}_0, \mathbf{Y}_1)' = ((Y_1, \dots, Y_m), (Y_{m+1}, \dots, Y_{k-1}), (Y_k, \dots, Y_n))$$

and let $\mathbf{X} = (\mathbf{X}_{-1}, \mathbf{X}_0, \mathbf{X}_1)$, where

$$\mathbf{X}' = \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 2 & 3 & \dots & n \end{pmatrix},$$

is the associated design matrix. Note that in the text $m = 3$ observations are assumed to be prechange. Let $\hat{\mathbf{b}}$ and $\hat{\nu}$ be estimates of $(\beta_0, \beta_1)'$ and σ^2 , based on \mathbf{Y}_{-1} ,

$$\begin{aligned} \hat{\mathbf{b}} &= (\mathbf{X}'_{-1}\mathbf{X}_{-1})^{-1}\mathbf{X}'_{-1}\mathbf{Y}_{-1}, \\ \hat{\nu} &= \mathbf{Y}'_{-1}[\mathbf{I} - \mathbf{X}_{-1}(\mathbf{X}'_{-1}\mathbf{X}_{-1})^{-1}\mathbf{X}'_{-1}]\mathbf{Y}_{-1}/(m - 2). \end{aligned}$$

Define the statistics

$$(A.1) \quad \mathbf{T} = (\mathbf{T}_0, \mathbf{T}_1) = \hat{\nu}^{-1/2}((\mathbf{Y}_0, \mathbf{Y}_1) - (\mathbf{X}_0, \mathbf{X}_1)\hat{\mathbf{b}}).$$

Regard \mathbf{Y}_{-1} as a learning sample, taken from the prechange distribution. A change of the regression parameters, if it occurs at all, is assumed to occur at some (unknown) time $\nu = k > m$, in which case $\mathcal{L}(Y_i | Y_1, \dots, Y_{i-1}) = N(Y_{i-1}, \sigma^2)$ for $i \geq k + 1$. (If there is no change at all, denote $\nu = \infty$.) Monitoring for a change will be based on the the invariant sequence of $|\mathbf{T}|$'s (where $|\mathbf{T}|$ is the vector of absolute values of the coordinates of \mathbf{T}). Let $S^{(k)}(\mathbf{T})$ be the marginal log-likelihood ratio of the sequence of invariant statistics.

The distributions of both $\hat{S}^{(k)}(n)$ and $S^{(k)}(\mathbf{T})$ are invariant with respect to an affine transformation of time and multiplication by a positive scalar. Hence, one may assume, without loss of generality, that $\beta_0 = \beta_1 = 0$ and $\sigma^2 = 1$.

Step 2. Denote: $\Theta = (\hat{\mathbf{b}}, \hat{v})$. Let $S^{(k)}(\mathbf{T} | \Theta)$ be the conditional log-likelihood ratio of \mathbf{T} , for $\{\nu = k\}$ versus $\{\nu = \infty\}$. Let f_k denote the density under $\nu = k$. Now,

$$\begin{aligned} S^{(k)}(t) &\stackrel{\text{def}}{=} \log \frac{f_k(t)}{f_\infty(t)} = \log \frac{\int f_k(t, \theta) d\theta}{\int f_\infty(t, \theta) d\theta} \\ &= \log \int \frac{f_k(t, \theta)}{f_\infty(t, \theta)} \left[\frac{f_\infty(t, \theta)}{\int f_\infty(t, u) du} \right] d\theta \\ &= \log \int \frac{f_k(t|\theta)}{f_\infty(t|\theta)} f_\infty(\theta|t) d\theta \\ &\geq \int \log \left[\frac{f_k(t|\theta)}{f_\infty(t|\theta)} \right] f_\infty(\theta|t) d\theta. \quad \square \end{aligned}$$

It follows that $\tilde{S}^{(k)}(\mathbf{T}) = E_\infty[S^{(k)}(\mathbf{T} | \Theta) | \mathbf{T}]$ is smaller than the (true) log-likelihood ratio $S^{(k)}(\mathbf{T})$.

Step 3. Here it is shown that

$$(A.2) \quad \tilde{S}^{(k)}(\mathbf{T}) = \hat{\sigma}^{-2}(\mathbf{Y}_1 - \mathbf{X}_1\hat{\beta})' \mathbf{A}(\mathbf{Y}_1 - \mathbf{X}_1\hat{\beta}) + \text{trace}(\mathbf{X}'_1 \mathbf{A} \mathbf{X}_1 (\mathbf{X}' \mathbf{X})^{-1}),$$

with

$$\begin{aligned} \hat{\beta} &= (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}, \\ \hat{\sigma}^2 &= \frac{\mathbf{Y}' \mathbf{Y} - \mathbf{Y}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}}{n - 2}. \end{aligned}$$

Note that the first term on the right-hand side of (A.2) is what is denoted by $\hat{S}^{(k)}(n)$. Throughout this step the superscript k is omitted in order to simplify the notation.

Observe that

$$\begin{aligned} \hat{\mathbf{b}} &\sim N(0, (\mathbf{X}'_{-1} \mathbf{X}_{-1})^{-1}), \\ \hat{v} &\sim \Gamma((m - 2)/2, (m - 2)/2) \end{aligned}$$

are independent. Furthermore,

$$\mathcal{L}(\mathbf{T} | \theta) = \mathcal{L}(\mathbf{T} | \hat{\mathbf{b}}, \hat{v}) = N(-\hat{v}^{-1/2}(\mathbf{X}_0, \mathbf{X}_1)\hat{\mathbf{b}}, \hat{v}^{-1}\mathbf{I}).$$

Standard Bayesian arguments can be used to show that

$$(A.3) \quad \mathcal{L}(\hat{\mathbf{b}} | \hat{v}, \mathbf{T}) = N(-\hat{v}^{1/2}(\mathbf{X}' \mathbf{X})^{-1}(\mathbf{X}_0, \mathbf{X}_1)' \mathbf{T}, (\mathbf{X}' \mathbf{X})^{-1}),$$

$$(A.4) \quad \mathcal{L}(\hat{v} | \mathbf{T}) = \Gamma\left(\frac{n - 2}{2}, \frac{\mathbf{T}' \mathbf{T} + m - 2 - \mathbf{T}'(\mathbf{X}_0, \mathbf{X}_1)(\mathbf{X}' \mathbf{X})^{-1}(\mathbf{X}_0, \mathbf{X}_1)' \mathbf{T}}{2}\right).$$

Let $\mathbf{A} = \mathbf{A}_{k, n}$ be, for general k and n , an $(n - k + 1)$ by $(n - k + 1)$ matrix of the following form (demonstrated for $n - k + 1 = 5$):

$$\mathbf{A} = \begin{pmatrix} -0.5 & 1 & 0 & 0 & 0 \\ 0 & -0.5 & 1 & 0 & 0 \\ 0 & 0 & -0.5 & 1 & 0 \\ 0 & 0 & 0 & -0.5 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

Recall that if a change occurs at time $\nu = k$, then $\mathcal{L}(Y_i | Y_1, \dots, Y_{i-1}) = N(Y_{i-1}, 1)$ for $i \geq k + 1$. Thus the log-likelihood ratio of \mathbf{Y} for $\{\nu = k\}$ versus $\{\nu = \infty\}$ is $\sum_{i=k+1}^n (Y_i Y_{i-1} - Y_{i-1}^2/2) = \mathbf{Y}'_1 \mathbf{A} \mathbf{Y}_1$.

The term $S(\mathbf{T} | \hat{\mathbf{b}}, \hat{\nu})$ is the conditional log-likelihood ratio of \mathbf{T} for $\nu = k$ versus $\nu = \infty$. From (A.1) it follows that $S(\mathbf{T} | \hat{\mathbf{b}}, \hat{\nu})$ equals the conditional log-likelihood ratio of $\mathbf{Y}_1 = \hat{\nu}^{1/2} \mathbf{T}_1 + \mathbf{X}_1 \hat{\mathbf{b}}$, so that

$$S(\mathbf{T} | \hat{\mathbf{b}}, \hat{\nu}) = [\hat{\nu}^{1/2} \mathbf{T}_1 + \mathbf{X}_1 \hat{\mathbf{b}}]' \mathbf{A} [\hat{\nu}^{1/2} \mathbf{T}_1 + \mathbf{X}_1 \hat{\mathbf{b}}].$$

Regarding this as a function of $\hat{\mathbf{b}}$, obtain from (A.3) that

$$E_\infty[S(\mathbf{T} | \hat{\mathbf{b}}, \hat{\nu}) | \hat{\nu}, \mathbf{T}] = \hat{\nu}(\mathbf{T}_1 - \mathbf{X}_1 \tilde{\beta})' \mathbf{A} (\mathbf{T}_1 - \mathbf{X}_1 \tilde{\beta}) + \text{trace}(\mathbf{X}'_1 \mathbf{A} \mathbf{X}_1 (\mathbf{X} \mathbf{X})^{-1}),$$

where $\tilde{\beta} = (\mathbf{X} \mathbf{X})^{-1}(\mathbf{X}_0, \mathbf{X}_1)' \mathbf{T}$. Therefore,

$$(A.5) \quad \tilde{S}(\mathbf{T}) = E_\infty[S(\mathbf{T} | \hat{\mathbf{b}}, \hat{\nu}) | \mathbf{T}] = E_\infty[E_\infty(S(\mathbf{T} | \hat{\mathbf{b}}, \hat{\nu}) | \hat{\nu}, \mathbf{T}) | \mathbf{T}]$$

$$(A.6) \quad = \tilde{\sigma}^{-2}(\mathbf{T}_1 - \mathbf{X}_1 \tilde{\beta})' \mathbf{A} (\mathbf{T}_1 - \mathbf{X}_1 \tilde{\beta}) + \text{trace}(\mathbf{X}'_1 \mathbf{A} \mathbf{X}_1 (\mathbf{X} \mathbf{X})^{-1}),$$

where

$$\tilde{\sigma}^2 = \frac{\mathbf{T}' \mathbf{T} + m - 2 - \mathbf{T}'(\mathbf{X}_0, \mathbf{X}_1)(\mathbf{X} \mathbf{X})^{-1}(\mathbf{X}_0, \mathbf{X}_1)' \mathbf{T}}{n - 2}.$$

However, from

$$(\mathbf{T}_1 - \mathbf{X}_1 \tilde{\beta}) = \frac{\mathbf{Y}_1 - \mathbf{X}_1 \hat{\mathbf{b}}}{\hat{\nu}^{1/2}} - \frac{\mathbf{X}_1 (\mathbf{X} \mathbf{X})^{-1} (\mathbf{X}_0, \mathbf{X}_1)' [(\mathbf{Y}_0, \mathbf{Y}_1) - (\mathbf{X}_0, \mathbf{X}_1) \hat{\mathbf{b}}]}{\hat{\nu}^{1/2}}$$

obtain that [since $(\mathbf{X}'_{-1} \mathbf{X}_{-1}) \hat{\mathbf{b}} = \mathbf{X}'_{-1} \mathbf{Y}_{-1}$]

$$\begin{aligned} & \hat{\nu}^{1/2}(\mathbf{T}_1 - \mathbf{X}_1 \tilde{\beta}) \\ &= \mathbf{Y}_1 - \mathbf{X}_1 \{(\mathbf{X} \mathbf{X})^{-1}(\mathbf{X}_{-1}, \mathbf{X}_0, \mathbf{X}_1)'(\mathbf{Y}_{-1}, \mathbf{Y}_0, \mathbf{Y}_1)\} \\ & \quad + \mathbf{X}_1 \{(\mathbf{X} \mathbf{X})^{-1}(\mathbf{X}_0, \mathbf{X}_1)'(\mathbf{X}_0, \mathbf{X}_1) \hat{\mathbf{b}} + (\mathbf{X} \mathbf{X})^{-1} \mathbf{X}'_{-1} \mathbf{Y}_{-1} - \hat{\mathbf{b}}\} \\ &= \mathbf{Y}_1 - \mathbf{X}_1 \{(\mathbf{X} \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}\} \\ & \quad + \mathbf{X}_1 \{(\mathbf{X} \mathbf{X})^{-1}(\mathbf{X}_0, \mathbf{X}_1)'(\mathbf{X}_0, \mathbf{X}_1) \hat{\mathbf{b}} + (\mathbf{X} \mathbf{X})^{-1} \mathbf{X}'_{-1} \mathbf{Y}_{-1} - \hat{\mathbf{b}}\} \\ &= \mathbf{Y}_1 - \mathbf{X}_1 \hat{\beta} + \mathbf{X}_1 \{(\mathbf{X} \mathbf{X})^{-1}(\mathbf{X} \mathbf{X}) \hat{\mathbf{b}} - \hat{\mathbf{b}}\} \\ &= \mathbf{Y}_1 - \mathbf{X}_1 \hat{\beta} \end{aligned}$$

and likewise

$$\hat{v}^{1/2}(\mathbf{T}_0 - \mathbf{X}_0\tilde{\beta}) = \mathbf{Y}_0 - \mathbf{X}_0\hat{\beta}.$$

Therefore, this step will be completed if it can be shown that $\hat{v}\hat{\sigma}^2 = \hat{\sigma}^2$. Indeed,

$$\begin{aligned} (n-2)\hat{\sigma}^2 &= \mathbf{T}'_0(\mathbf{T}_0 - \mathbf{X}_0\tilde{\beta}) + \mathbf{T}'_1(\mathbf{T}_1 - \mathbf{X}_1\tilde{\beta}) + \hat{v}^{-1}\hat{v}(m-2) \\ &= \hat{v}^{-1}(\mathbf{Y}_0 - \mathbf{X}_0\hat{\mathbf{b}})'(\mathbf{Y}_0 - \mathbf{X}_0\hat{\beta}) + \hat{v}^{-1}(\mathbf{Y}_1 - \mathbf{X}_1\hat{\mathbf{b}})'(\mathbf{Y}_1 - \mathbf{X}_1\hat{\beta}) \\ &\quad + \hat{v}^{-1}(\mathbf{Y}'_{-1}\mathbf{Y}_{-1} - \mathbf{Y}'_{-1}\mathbf{X}_{-1}\hat{\mathbf{b}}) \end{aligned}$$

so that [since $\mathbf{Y}'_{-1}\mathbf{X}_{-1} = \hat{\mathbf{b}}'(\mathbf{X}'_{-1}\mathbf{X}_{-1})$]

$$\begin{aligned} \hat{v}(n-2)\hat{\sigma}^2 &= \mathbf{Y}'\mathbf{Y} - \mathbf{Y}'\mathbf{X}\hat{\beta} - (\hat{\mathbf{b}}'\mathbf{X}'\mathbf{Y} - \hat{\mathbf{b}}'\mathbf{X}'_0\mathbf{X}_0\hat{\beta} - \hat{\mathbf{b}}'\mathbf{X}'_1\mathbf{X}_1\hat{\beta} - \mathbf{Y}'_{-1}\mathbf{X}_{-1}\hat{\beta}) \\ &= (n-2)\hat{\sigma}^2 - (\hat{\mathbf{b}}'\mathbf{X}'\mathbf{Y} - \hat{\mathbf{b}}'\mathbf{X}'\mathbf{X}\hat{\beta}) \\ &= (n-2)\hat{\sigma}^2. \end{aligned} \quad \square$$

Step 4. $\text{trace}(\mathbf{X}'_1\mathbf{A}\mathbf{X}_1(\mathbf{X}'\mathbf{X})^{-1}) > 0$. Straightforward calculations yield

$$\mathbf{X}'_1\mathbf{A}\mathbf{X}_1 = \begin{pmatrix} \frac{n-k}{2} & \frac{(n-k)(n+k+3)}{4} \\ \frac{(n-k)(n+k-1)}{4} & \frac{n-k}{12}(2n^2 + 2kn + 2k^2 + 3n + 3k - 5) \end{pmatrix}$$

and

$$(\mathbf{X}'\mathbf{X})^{-1} = \begin{pmatrix} \frac{n(n+1)(2n+1)}{6} & -\frac{n(n+1)}{2} \\ -\frac{n(n+1)}{2} & n \end{pmatrix} \frac{1}{Q},$$

where $Q = n^2(n+1)(2n+1)/6 - n^2(n+1)^2/4$. Algebra yields (for $n \geq k \geq 2$)

$$\text{trace}(\mathbf{X}'_1\mathbf{A}\mathbf{X}_1(\mathbf{X}'\mathbf{X})^{-1}) = \frac{(n-k)n}{12Q}[n^2 - nk + 2k^2 - 7] > 0. \quad \square$$

Step 5. Now the proof of Lemma 2 is complete. From Steps 3 and 4 obtain that $S^{(k)}(\mathbf{T}) \geq \hat{S}^{(k)}(n)$. Let N_{SR}^S be defined as N_{SR} with $S^{(k)}(\mathbf{T})$ instead of $\hat{S}^{(k)}(n)$. It follows from the above that N_{SR}^S is stochastically smaller than N_{SR} .

We turn to the investigation of the $P^{(\infty)}$ -distribution of N_{SR}^S . For any $j \geq 2$,

$$\begin{aligned} &P^{(\infty)}((j-1)t < N_{\text{SR}}^S \leq jt) \\ \text{(A.7)} \quad &= P^{(\infty)}\left((j-1)t < N_{\text{SR}}^S, \max_{(j-1)t < n \leq jt} \sum_{k=n-t}^n \exp(S^{(k)}(n)) \geq dA\right) \end{aligned}$$

$$(A.8) \quad \leq P^{(\infty)}\left((j-2)t < N_{SR}^S, \max_{(j-1)t < n \leq jt} \sum_{k=(j-2)t+1}^n \exp(S^{(k)}(n)) \geq dA\right)$$

$$(A.9) \quad \leq \frac{2t}{dA} P^{(\infty)}((j-2)t < N_{SR}^S),$$

where the last inequality follows from Doob’s inequality, after noticing that

$$\exp\{S^{(k)}(\mathbf{T})\} = \frac{f_k(T_4, \dots, T_n)/f_k(T_4, \dots, T_{(j-2)t})}{f_\infty(T_4, \dots, T_n)/f_\infty(T_4, \dots, T_{(j-2)t})}$$

is a likelihood ratio (conditional on $T_4, \dots, T_{(j-2)t}$), making the sum in (A.8) a conditional submartingale.

For all $1 \leq n \leq t$, note that $P^{(\infty)}(N_{SR}^S > n) \geq 1 - n/dA$. Divide (A.7) and (A.9) by $P^{(\infty)}(N_{SR}^S > (j-1)t)$. Apply induction to get the bound

$$P^{(\infty)}(N_{SR}^S > jt \mid N_{SR}^S > (j-1)t) \geq 1/2 + [1/4 - (2t)/(dA)]^{1/2} \approx 1 - (2t)/(dA),$$

which leads to the relation

$$(A.10) \quad P^{(\infty)}(N_{SR}^S > n) \geq \{1/2 + [1/4 - (2t)/(dA)]^{1/2}\}^{n/t+1}.$$

It follows from the definition of c and d that $E^{(\infty)}N \geq E^{(\infty)}(G \wedge N_{SR}^S) \geq A$. This completes the proof of Lemma 2. \square

PROOF OF THEOREM 3. The distribution of $\hat{S}^{(k)}(n)$ is invariant with respect to an affine transformation of time and multiplication by a positive scalar. Hence, one may assume, without loss of generality, that $\beta_0 = \beta_1 = 0$ and $\sigma^2 = 1$.

Fix k . Consider the statistic $\mathbf{Z}_k = (\sum_{i=1}^{k-1} Y_i, \sum_{i=1}^{k-1} iY_i, \sum_{i=1}^{k-1} Y_i^2)$ and define the event B_k as

$$\left\{ \left| \sum_{i=1}^{k-1} Y_i \right| \leq \varepsilon(k-1), \left| \sum_{i=1}^{k-1} iY_i \right| \leq \varepsilon(k-1)^2, \left| \sum_{i=1}^{k-1} Y_i^2 - (k-1) \right| \leq \varepsilon(k-1) \right\}.$$

Note that $P^{(k)}(\bar{B}_k) \leq \exp(-\tau(k-1))$, for some $\tau > 0$, where \bar{B}_k is the complement of B_k . Let $k < n$ and consider $\hat{S}^{(k)}(n)$. Note that given \mathbf{Z}_k , the event $\{N \geq k\}$ and the statistic $\hat{S}^{(k)}(n)$ are independent.

Let $n_1 = k - 1 + \lfloor (1 + \varepsilon)I^{-1}(\log A) \rfloor$ and $n_2 = k + \lceil c_1 I^{-1}(\log A) \rceil$ for some small ε and some large c_1 (to be determined later). It follows that

$$\begin{aligned} & E^{(k)}(N - k + 1 \mid N \geq k) \\ & \leq (1 + \varepsilon)I^{-1}(\log A) + c_1 I^{-1}(\log A) \max_{\mathbf{Z}_k \in B_k} P^{(k)}(\hat{S}^{(k)}(n_1) < \log A \mid \mathbf{Z}_k) \\ & \quad + E(G) \frac{P^{(\infty)}(\bar{B}_k)}{P^{(\infty)}(N \geq k)} + E(G) \max_{\mathbf{Z}_k \in \bar{B}_k} P^{(k)}(\hat{S}^{(k)}(n_2) < \log A \mid \mathbf{Z}_k). \end{aligned}$$

The proof will be complete when, on the right-hand side of this inequality, all terms but the first will be shown to be negligible.

Begin with the term $E(G)P^{(\infty)}(\bar{B}_k)/P^{(\infty)}(N \geq k)$. Note that the ratio of the probabilities is less than $\exp\{t/A - (k - 1)[\tau - (\zeta + 1/c)/A]\}$, for some $\zeta > 0$. Thus, if $k \gg (1/\tau) \log A$ then this term is $o(1/A)$, as $A \rightarrow \infty$.

For the other two terms we need to investigate $\hat{S}^{(k)}(n)$, both when $n = n_1$ and when $n = n_2$. However, it can be shown that

$$\hat{S}^{(k)}(n) = \frac{1}{2\hat{\sigma}^2} \left[\sum_{i=k+1}^n (Y_i - (\hat{\beta}_0 + i\hat{\beta}_1))^2 - \sum_{i=k+1}^n (Y_i - Y_{i-1} - \hat{\beta}_1)^2 \right].$$

The distribution of $\hat{S}^{(k)}(n)$ depends on \mathbf{Z}_k only via $(\hat{\beta}_0, \hat{\beta}_1, \hat{\sigma}^2)$. We first show that the conditional probability of the event

$$C_{k,n} = \{\hat{\sigma}^2 > 1 + 2\varepsilon\} \cup \{|\hat{\beta}_0 + |n\hat{\beta}_1| > 2\varepsilon\}$$

is $o(1)$ when $n = n_1$ and is $o(1/A)$ when $n = n_2$. It is then straightforward to see that the probability of the event $\{\hat{S}^{(k)}(n) < \log A\} \cap \bar{C}_{k,n}$, is $o(1)$ when $n = n_1$ and is $o(1/A)$ when $n = n_2$, which is all that is needed in order to prove that the terms in question are small. Indeed,

$$\begin{aligned} P^{(k)}(\hat{\sigma}^2 > 1 + 2\varepsilon | \mathbf{Z}_k) &\leq P^{(k)}\left(\frac{\sum_{i=1}^{k-1} Y_i^2 + \sum_{i=k}^n Y_i^2}{n - 2} > 1 + \varepsilon | \mathbf{Z}_k\right) \\ &\leq P^{(k)}(\sum_{i=k}^n Y_i^2 \geq \varepsilon(n - 2)), \end{aligned}$$

since $\mathbf{Z}_k \in B_k$. The $P^{(k)}$ -distribution of $\sum_{i=k}^n Y_i^2$ is a noncentral χ^2 with $n - k + 1$ degrees of freedom and $2I(n - k + 1)$ as the parameter of noncentrality. An exponential Markov inequality can be used in order to establish the necessary bounds. Similar arguments can be used for the conditional probability of the event $\{|\hat{\beta}_0 + |n\hat{\beta}_1| > \varepsilon\}$. \square

REFERENCES

LAI, T. L. (1993). Information bounds and quick detection of parameter changes in stochastic systems. Technical report, Stanford Univ.
 LAI, T. L. (1995). Sequential changepoint detection in quality control and dynamical systems. *J. Roy. Statist. Soc. Ser. B* **57** 613–658.
 LORDEN, G. (1971). Procedures for reacting to a change in distribution. *Ann. Math. Statist.* **42** 1897–1908.
 MOUSTAKIDES, G. V. (1986). Optimal stopping times for detecting changes in distributions. *Ann. Statist.* **14** 1379–1387.
 POLLAK, M. (1985). Optimal detection of a change in distribution. *Ann. Statist.* **13** 206–227.
 POLLAK, M. and SIEGMUND, D. (1985). A diffusion process and its application to detecting a change in the drift of a Brownian motion. *Biometrika* **72** 267–280.
 RITOV, Y. (1990). Decision theoretic optimality of the CUSUM procedure. *Ann. Statist.* **18** 1464–1469.
 ROBBINS, H. and ZHANG, C.-H. (1993). A change point problem with some applications. Technical report, Rutgers Univ.
 YAKIR, B. (1996). A lower bound on the ARL to detection with a probability constraint on false alarm. *Ann. Statist.* **24** 431–435.

- YAKIR, B. (1997). A note on optimal detection of a change in distribution. *Ann. Statist.* **25** 2117-2126.
- YAO, Q. (1993). Asymptotically optimal detection of a change in a linear model. *Sequential Anal.* **12** 201-210.

DEPARTMENT OF STATISTICS
HEBREW UNIVERSITY
JERUSALEM
ISRAEL 91905
E-MAIL: msby@mssc.huji.ac.il
msmp@mssc.huji.ac.il

DEPARTMENT OF STATISTICS
UNIVERSITY OF PENNSYLVANIA
PHILADELPHIA, PENNSYLVANIA 19104-6302
E-MAIL: abba@compstat.wharton.upenn.edu