University of Alberta

# Sign Test for Change-Point Problem

by

Xiongsheng Jin     ©

A thesis

submitted to the Faculty of Graduate Studies and Research in

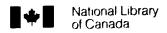partial fulfillment of the requirements for the degree of

Master of Science

in

Statistics

Department of Mathematical Sciences

Edmonton, Alberta

Fall 1996

Canada

# University of Alberta
# Library Release Form

**Name of Auther:**  Xiongsheng Jin

**Title of Thesis:**  Sign Test for Change-Point Problem

**Degree:**  Master of Science

**Year this Degree Granted:**  Fall 1996

Permission is hereby granted to the University of Alberta Library to re produce single copies of this thesis and to lend or sell such copies for private, scholarly, or scientific research purposes only.

The auther reserves all other publication and other rights in association with the copyright in the thesis, and except as hereinbefore provided, neither the thesis nor any substantial portion thereof may be printed or otherwise reproduced in any material from whatever without the auther's prior written permission.

_____

Department of Mathematical Sciences

University of Alberta

Edmonton, Alberta

Canada, T6G 2G1

Date: September 11, 1996

University of Alberta

Faculty of Graduate Studies and Research

The undersigned certify that they have read, and recommend to the Faculty of Graduate Studies and Research for acceptance, a thesis entitled **Sign Test for Change-Point Problem** submitted by **Xiongsheng Jin** in partial fulfillment of the requirements for the degree of **Master of Science** in Statistics.

Dr. **Edit Gombay** (Supervisor)

Dr. **A.N.Al-Hussaini**

(Committee Chair and Examiner)

Dr. **Connie K. Varnhagen**

Date: August 8, 1996

# ABSTRACT

The Sign test is employed to study the change-point problem with epidemic alternative. Discussions focus on the two different cases that under the null hypothesis the population median is known or unknown. The asymptotic distributions of the test statistic under the alternative hypothesis are proved. Numerical simulation is carried out to calculate the estimated change-points, test statistic values and their P-values.

# ACKNOWLEDGEMENTS

# Contents

# List of Tables

# Chapter 1

# INTRODUCTION

Change-point problem originally arose in the field of qualit control. When one monitors the output in a production line, one wants to keep the quality of the product within a required region and to detect the quality deviation across the threshold values as soon as possible. In Statistics, this problem can usually be modeled as follows. We have a sequence of observations of independent random variables $x_1, x_2, \cdots$ of identical distribution and want to detect whether a change at time $\tau$ could have occured in this sequence and that after time $\tau, x_\tau, x_{\tau+1}, \cdots$ have different distribution as that of $x_1, x_2, \cdots, x_{\tau-1}$. We call this problem the change-point problem and the $\tau$ the change-point.

In change-point problem, one usually needs to consider following :

[1]Testing the hypotheses:

$$H_0: \quad x_1, \cdots, x_n, \cdots \ i.i.d \sim F(x) \qquad No\ change.$$

$$H_1: \quad x_1, \cdots, x_{\tau-1} \ i.i.d \sim F(x) \qquad \qquad (1.1)$$

$$x_\tau, \cdots, x_n, \cdots \ i.i.d \sim G(x), \ F(x) \neq G(x) for\ some\ x.\ There\ is\ a\ change.$$

$\tau$ is the unknown change-point.

[2]Employing a suitable statistic $T_n$ for the test problem to obtain an estimator $\hat{\tau}(n)$ for the unknown change-point $\tau$.

[3]Discussing the properties of $T_n$ and $\hat{\tau}(n)$ and carrying out some numerical simulations to confirm the theoretical results.

For the test hypotheses of change-point problem, usually we assume that $x_1, \cdots, x_n, \cdots$ are independent continuous random variables. Beside (1.1), there are many special forms to $\circ$      the test hypotheses. For example, the test for change in the location par,     er can be written as

$$H_0: \quad x_1, \cdots, x_n \ i.i.d \sim F(x)$$

$$H_1: \quad x_1, \cdots, x_{\tau-1} \ i.i.d \sim F(x); \qquad\qquad (1.2)$$

$$x_\tau, \cdots, x_n \ i.i.d \sim F(x + \Delta), \qquad -\infty < \Delta < +\infty.$$

(1.2) is equivalent to:

$$H_0: \quad \Delta = 0$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad (1.3)$$

$$H_1: \quad \Delta \neq 0$$

Sometimes people look at

$$H_0: \quad \mu_1 = \cdots = \mu_m = \mu_0, \quad \mu_i = E(x_i)$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad (1.4)$$

$$H_1: \quad \mu_1 = \cdots = \mu_{\tau-1} \neq \mu_\tau = \cdots = \mu_n$$

In these above models, one just considers the alternative hypothesis of at most one change, the so called AMOC model. A slight generalization of the AMOC model, that is often very useful, is the more than one change point

model with the epidemic or square wave alternative.

$$H_0: \quad x_1, \cdots, x_n \; i.i.d \sim F(x)$$

$$H_1: \quad x_1, \cdots, x_{\tau_1-1}, x_{\tau_2}, \cdots, x_n \; i.i.d \sim F(x) \tag{1.5}$$

$$x_{\tau_1}, \cdots, x_{\tau_2-1} \; i.i.d \sim G(x), \quad F(x) \neq G(x),$$

where $\tau_1, \tau_2$ are unknown change-points.

There are also many different viewpoints used in change-point research. When one observes the output in a production line, one can use a sequential procedure where one observes the products sequentially and stops the line at a random time when one detects a change in quality, or fixed sample size procedure also called retroactive change-point detection procedure where one observes a large finite sequence of output such as the product produced in a day to determine possible change within the collection. People also use classical and Bayesian approaches, parametric and nonparametric models for change-point problem. So, there has been much research done for the change-point problem with the combination of different methods and models.

The basic AMOC problem was first considered by Page (1954, 1955) in the model (1.2). Assuming the initial value $\mu_0$ known, Page studied testing the null hypothesis of no change ($H_0 : \Delta = 0$) against either one or two sided alternatives ($H_1 : \Delta > 0$ or $H_1 : \Delta \neq 0$). Let $S_0 = 0$ and $S_k = \sum_{j=1}^{k} V_j \; k = 1, \cdots, n$

$$V_j = \begin{cases} a & \text{if } x_j \geq \mu_0, \\ b & \text{if } x_j < \mu_0 \end{cases}$$

where $a > 0$, $b > 0$ are constants, such that $E_{\mu_0}(V_j) = 0$, $j = 1, \cdots, n$. Page's decision rule rejects $H_0 : \Delta = 0$ in favour of the alternative of one change

3

$H_1 : \Delta > 0$, if

$$T_n = max_{0 \le k \le n}\{S_k - min_{0 \le i \le k}S_i\} \tag{1.6}$$

is too large.

S. Csörgö and Horváth (1983) calculated the limit distribution of $T_n$.

$$\lim_{n \to \infty}\{T_n/(nab)^{\frac{1}{2}} < x\} = P\{sup_{0 \le t \le 1}|w(t)| \le x\} \tag{1.7}$$

$$= 1 - 4\sum_{k=1}^{\infty}(-1)^{k+1}\Phi(-(2k-1)x), \qquad x > 0.$$

where w(t) is a Weiner process and $\Phi$ is the standard normal distribution function. A table for this limit distribution was also given.

G. K. Bhattacharyya and Johnson (1968) considered a general class of locally optimal rank tests for the change-point problem in the following two cases:

1 The initial distribution $F_0$ is known and symmetric around the origin. Testing of these hypothesis corresponds to the shift problem in model (2) with unknown $\Delta > 0$. Bhattacharyya and Johnson employed the criterion of maximizing the average local power $\overline{\beta(\Delta)} = \sum_{i=1}^{n} q_i\beta(\Delta|i)$ with respect to arbitrary set of weights $q_i$ that satisfies $q_1 = 0$, $q_i \ge 0$, $i = 2,\cdots,n$ and $\sum_{i=1}^{n} q_i = 1$ to get a nonparametric statistic

$$T_n = \sum_{i=1}^{n} Q_i sgn(x_i)E\{-f_0'(V^{(R_i)})/f_0(V^{(R_i)})\} \tag{1.8}$$

to reject $H_0$ at large value of $T_n$. Here $V^{(1)} \le \cdots \le V^{(n)}$ is an ordered statistic of n i.i.d random variables having a distribution $F_0$, $Q_i = \sum_{j=1}^{i} q_i$,

4

$(R_1, \cdots R_n)$ is the vector of the rank of $(x_1, \cdots x_n)$ and $\beta(\Delta|i)$ is the power at $\Delta$ when the change occurs at time i. From the Bayesian viewpoint, $q_i$ may be regarded as the prior probability of a change to occur at time i.

**2** When initial level is unknown, they proposed the

$$S_n = \sum_{i=1}^{n} Q_i E\{-f'(V^{(R_i)})/f(V^{(R_i)})\}, \tag{1.9}$$

and suggested to reject $H_0$ for larger value of $S_n$. In both cases, the tests are distribution free, they depend upon the weight function $\{q_i\}$ and are unbiased for general classes of shift alternatives. The asympototic distribution of the test statistic under the local translation alternative was also reached.

A. Sen and Srivastava (1975) proposed two nonlinear rank test for one-sided alternative with $x_i \sim N(\mu_i, \sigma^2)$, $i = 1, \cdots n$ and unknown initial $\mu_0$ and $\sigma^2$. They suggested rejecting $H_0 : \Delta = 0$ in favor of $H_1 : \Delta > 0$ for a large value of

$$I_1 = max_{1 \le k \le n-1}\{[M_{k,n-k} - E_0(M_{k,n-k})]/[Var_0(M_{k,n-k})]^{\frac{1}{2}}\}, \tag{1.10}$$

or

$$I_2 = max_{1 \le k \le n-1}\{[U_{k,n-k} - E_0(U_{k,n-k})]/[Var_0(U_{k,n-k})]^{\frac{1}{2}}\}, \tag{1.11}$$

where

$$M_{k,n-k} = \sum_{i=k+1}^{n} \Psi\{x_i - median_{1 \le j \le n}(x_j)\},$$

$$U_{k,n-k} = \sum_{i=k+1}^{n} \sum_{j=1}^{k} \Psi(x_i - x_j),$$

5

and

$$\Psi(t) = \begin{cases} 1 & t > 0, \\ 0 & t \leq 0, \end{cases}$$

$E_0(\cdot)$, $Var_0(\cdot)$ above denotes the mean and variance taken under null hypothesis respectively. Some Monte Carlo simulations for the estimated critical values were also provided.

A model that is similar to A.Sen and Srivastava's was studied by Hawkins(1977) for the two-sided alternative hypothesis. He provided the test statistic

$$U_n = max_{1 \leq k \leq n-1} \mid T_k \mid, \tag{1.12}$$

where

$$T_k = (\frac{n}{k(n-k)})^{\frac{1}{2}} \sum_{i=1}^{n} (x_i - \bar{x}_n), \qquad k = 1, \cdots, n-1.$$

The recursive formulae for the exact determination of the distribution of $U_n$ were also proved. With the normality of $T_1, \cdots, T_n$, he got the asympototic distribution of $U_n$ from the behavior of the maximum properties of a Gaussian process.

Pettitt (1979) proposed quite similar statistic to that of A. Sen and Srivastava for the one and two-sided tests. For the one-sided test: $H_0 : \Delta = 0$ vs $H_1 : \Delta > 0$, he suggested the statistic

$$\begin{aligned} J_1 &= min_{1 \leq k \leq n-1} \{ \sum_{i=1}^{k} \sum_{j=k+1}^{n} sgn(x_i - x_j) \} \\ &= min_{1 \leq k \leq n-1} \{ V_{k,n} \}. \end{aligned} \tag{1.13}$$

for the test and rejected $H_0$ for its large value. Here $V_{k,n} = \sum_{i=1}^{k} \sum_{j=k+1}^{n} sgn(x_i - x_j)$.

6

Pettitt proposed rejecting $H_0 : \Delta = 0$ in favour of the two-sided alternative $H_1 : \Delta \neq 0$ for large values of

$$J_2 = max_{1 \leq k \leq n-1} |\sum_{i=1}^{k} \sum_{j=k+1}^{n} sgn(x_i - x_j)|. \tag{1.14}$$

Pettitt proved that the limit distribution of

$$y_n(x) = n^{-1} \{ \frac{3}{n+1} \}^{\frac{1}{2}} V_{k,n} \tag{1.15}$$

is a Brownian bridge $y(x)$ and we know that

$$P\{sup \mid y(x) \mid \leq a\} = 1 + 2 \sum_{r=1}^{\infty} (-1)^r exp(-2r^2 a^2).$$

This is the limiting distribution of the Kolmogorov-Smirnov goodness of fit statistic and is extensively tabulated.

Comparing statistic $I_2$ and $J_2$, Schechtman and Wolfe (1981) proposed the following statistic

$$I_3 = max_{1 \leq k \leq n-1} \{ |U_{k,n-k} - E_0(U_{k,n-k})| / [Var_0(U_{k,n-k})]^{\frac{1}{2}} \} \tag{1.16}$$

for the two-sided test to reject $H_0 : \Delta = 0$ in favor of $H_1 : \Delta \neq 0$ for large value of $I_3$. The asymptotic properties of $I_3$ were also studied.

Lombard (1987) studied the smooth change model:

$$H_0 : \quad \mu_1 = \cdots = \mu_n = \xi_1 \qquad No\ change.$$

$$H_1 : \quad \mu_i = \theta_i \qquad i = 1, \cdots, n$$

$$\theta_i = \begin{cases} \xi_1 & i \leq \tau_1 \\ \xi_1 + (i - \tau_1)(\xi_2 - \xi_1)/(\tau_2 - \tau_1) & \tau_1 < i \leq \tau_2 \\ \xi_2 & i > \tau_2. \end{cases} \tag{1.17}$$

7

He considered the statistic

$$q_n = \sum_{t_1=1}^{n-1} \sum_{t_2=t_1+1}^{n} \{V^*_{t_1,t_2}\}^2$$

where

$$V^*_{t_1,t_2} = \sum_{j=t_1+1}^{t_2} \sum_{i=1}^{j} S(r_i),$$

$$S(r_i) = \{\Phi[i/(n+1)] - \overline{\Phi}\}/A,$$

$$\overline{\Phi} = n^{-1} \sum_{i=1}^{n} \phi[i/(n+1)],$$

$$(n-1)^{-1} \sum_{i=1}^{n} \{\phi[i/(n+1)] - \overline{\Phi}\}^2,$$

$\phi$ is an arbitrary score function satisfying $0 < \int_0^1 \phi^2 du < \infty$, and $r_1, \cdots, r_n$ are the ranks of $x_1, \cdots, x_n$, respectively. As $n \to \infty$, the null distribution of $n^{-5}q_n$ goes to the random variable $q = \sum_{n=1}^{\infty} (n\pi)^{-4} Z_n^2$ where $Z_1, \cdots, Z_n, \cdots$ are i.i.d $N(0,1)$ random variables. For the onset of trend model with $\tau_2 = n$, Lombard's statistics is $q_n^* = \sum_{t=1}^{n-1} \{V_{t,n}\}^2$. As $n \to \infty$, the null distribution of the $T^{-4}q_n^*$ approaches that of the random variable $q^* = \sum_{n=1}^{\infty} \lambda_n Z_n^2$ where $\lambda_1 > \lambda_2 > \cdots > 0$ is the positive real solution of the equation $tan\lambda^{-\frac{1}{4}} + tanh\lambda^{-\frac{1}{4}} = 0$. The AMOC model and multiple change point model were also discussed by Lombard.

For model (1.3) or (1.4), we have a general statistic to reject $H_0$ in favour of $H_1$ for large values of

$$max_{1 \leq k \leq n}\{| S_k - kS_n | /(k(1-k/n))^{\frac{1}{2}}\}. \tag{1.18}$$

Csörgö and Horváth (1986) studied the above statistic by considering

$$Z_n(t) = \begin{cases} (S_{[(n+1)t]} - [(n+1)t]S_n/n)/(n^{\frac{1}{2}}\sigma) & 0 \leq t < 1, \\ 0 & t = 1, \end{cases} \tag{1.19}$$

8

where $\sigma^2 = E(x_1 - E(x_1))^2$.

They proved that the process $Z_n(t)$, $(0 \le t \le 1)$ has the same asymptotic behaviour as the uniform quantile and empirical processes. Many asymptotic properties of nonparametric statistics were also given in their paper.

Gombay (1994) considered rank and sign stastistic for the epidemic alternative model of (1.5). She suggested the statistic

$$T_n = max_{k<l}|n^{\frac{1}{2}} \sum_{i=k}^{l-1} S(R_i)|.$$

for the rank test and proved the asymptotic distribution of $T_n$:

$$lim_{n\to\infty}P\{T_n \le c\} = 1 - \sum_{j=1}^{\infty} 2(4j^2c^2 - 1)e^{-2j^2c^2}. \qquad (1.20)$$

The asympototic consistency of $T_n$ was proved under some regularity conditions.

For the sign statistic, Gombay proposed the statisitc

$$U_n = max_{1 \le k<l \le n} \sum_{i=k}^{l-1} sgn(x_i - \xi_0). \qquad (1.21)$$

for the test:

$H_O : \quad x_i, \; i = 1, \cdots, n \; have \; known \; median \; \xi_0$

$H_1 : \quad x_i, \; i = 1, \cdots, \tau_1 - 1, \tau_2, \cdots, n \; have \; median \; \xi_0$

$\quad\quad x_i, \; i = \tau_1 - 1, \cdots, \tau_2 \; have \; median \; \xi_1, \; \xi_0 \ne \xi_1.$

Gombay suggested that $H_0$ be rejected for large values of $U_n$ when $\xi_1 > \xi_2$, and similarly for the $\xi_1 < \xi_0$ case. She also got the asymptotic distribution under the null hypothesis of the $U_n$

$$lim_{n\to\infty}P\{n^{-\frac{1}{2}}U_n \ge c\} = 1 - \frac{4}{\pi} \sum_{k=0}^{\infty} \frac{(-1)^k}{2k+1} exp\{-\pi^2 \frac{2k+1^2}{8c^2}\}, \quad c > 0. \quad (1.22)$$

9

Based on the research of Gombay, I consider the sign statistic for the epidemic alternative hypothesis. When $\xi_0$ is known, the asymptotic distribution under the alternative hypothesis is proved.

When the initial value $\xi_0$ is unknown, I suggest the use of the statistic

$$M(n) = max_{1 \leq k < l < n} n^{-\frac{1}{2}} | \sum_{k \leq i < l} sgn(x_i - \hat{\xi}_n)|.$$

where $\hat{\xi}_n = median(x_1, \cdots, x_n)$. The exact distribution of $n^{\frac{1}{2}} M(n)$ under $H_0$ is same as that of the maximum deviation in a simple symmetric random walk and it has been calculated exactly for each sample size.

To estimate the change-points, I use

$$(\hat{\tau}_1(n), \hat{\tau}_2(n)) = argmax_{k<l} |n^{-\frac{1}{2}} \sum_{i=k}^{l-1} sgn(x_i - \hat{\xi}_n)|.$$

as the estimatiors of $\tau_1, \tau_2$.

Finally, under $H_1$, I proved the asymptotic normality of $M(n)$ and

$$|\hat{\tau}_1(n) - \tau_1| + |\hat{\tau}_2(n) - \tau_2| = O_p(1)$$

In the third part, I calculated the power of sign test and do some numerical simulations.

# Chapter 2

# SIGN TEST FOR THE CHANGE-POINT PROBLEM WITH KNOWN INITIAL MEDIAN

Let $x_1, \cdots, x_n$ be a sequence of independent continuous random variables. Consider the following hypothesis test with two change-points :

$$H_O : \quad x_i, \; i = 1, \cdots, n \; have \; known \; median \; \xi_0,$$

$$H_1 : \quad x_i, \; i = 1, \cdots, \tau_1 - 1, \tau_2, \cdots, n \; have \; median \; \xi_0, \qquad (2.1)$$

$$x_i, \; i = \tau_1, \cdots, \tau_2 - 1 \; have \; median \; \xi_1, \quad \xi_1 \neq \xi_0,$$

where $\xi_0$ is known. The unknown integers $\tau_1$, $\tau_2$ are the change-points. We assume $\tau_1 = [n\lambda_1]$, $\tau_2 = [n\lambda_2]$ for some $0 < \lambda_1 < \lambda_2 < 1$. By $[a]$, we denote the integer part of a.

We employ the sign statistic for our test problem. Let

$$U_n = max_{k<l} \sum_{i=k}^{l-1} sgn(x_i - \xi_0)$$

$$= max_{k<l}(S_{l-1} - S_{k-1}). \qquad (2.2)$$

where $S_k = \sum_{i=1}^{k} sgn(x_i - \xi_0)$.

Under $H_0$, Gombay (1994) proved that the exact distribution of $U_n$ is

$$P\{U_n \geq N\} = 1 - \frac{2}{2N+1} \sum_{j=1}^{2N} (c(j))^n s(j(N+1)) \frac{1+c(j)}{s(j)} \frac{1-(-1)^j}{2}, \quad (2.3)$$

where $N$ is a positive integer and

$$c(j) = cos(\frac{j\pi}{2N+1}), \qquad s(j) = sin(\frac{j\pi}{2N+1}).$$

Also, the asymptotic distribution of the test statistic under the null hypothesis $H_0$ was shown to be

$$\lim_{n \to \infty} P\{n^{-\frac{1}{2}} U_n \geq c\} = 1 - \frac{4}{\pi} \sum_{k=0}^{\infty} \frac{(-1)^k}{2k+1} exp(-\pi^2 \frac{(2k+1)^2}{8c^2}), \quad (2.4)$$

for all $c > 0$.

Denote $\delta_n = P\{x_{\tau_1} > \xi_0\}$. We consider the following two cases of alternatives.

$$
\begin{aligned}
(i): \quad & \delta_n = \delta \qquad for\ all\ n \\
(ii): \quad & \delta_n \to \tfrac{1}{2} \ \ and \ \ \sqrt{n}|\delta_n - \tfrac{1}{2}| \to \infty.
\end{aligned}
\qquad (2.5)
$$

Case $(i)$ is the fixed alternative, while case $(ii)$ is the local but not contiguous alternative.

Assume $\delta_n > \tfrac{1}{2}$, as the other case is similar by symmetry.

**Lemma 2.1** *Under the hypothesis $H_1$:*

$$P\{max_{1 \leq l < \tau_2}(S_l - S_{\tau_2-1}) \geq max_{\tau_2 \leq l}(S_l - S_{\tau_2-1})\} \to 0, \quad as\ n \to \infty. \quad (2.6)$$

12

**Proof:** We want to prove the lemma by considering $l$ in different areas for the first term and keeping the second term unchanged.

(a) $l < \tau_1$

$$max_{l<\tau_1}(S_l - S_{\tau_2-1}) = max_{l<\tau_1}[(S_l - S_{\tau_1-1}) + (S_{\tau_1-1} - S_{\tau_2-1})].$$

Note that $S_l - S_{\tau_1-1}$, $l < \tau_1$ is a simple symmetic random walk.

In the sum $S_{\tau_1-1} - S_{\tau_2-1}$, the terms are

$$sgn(x_i - \xi_0) = \begin{cases} 1 & \text{w.p} & \delta_n \\ -1 & \text{w.p} & 1-\delta_n \end{cases} \quad \tau_1 \le i < \tau_2 \qquad (2.7)$$

so

$$\begin{aligned} E(sgn(x_i - \xi_0)) &= \delta_n + (-1)(1 - \delta_n) \\ &= 2\delta_n - 1 > 0 \qquad \tau_1 \le i < \tau_2 \end{aligned} \qquad (2.8)$$

$$\begin{aligned} Var(sgn(x_i - \xi_0)) &= E[(sgn(x_i - \xi_0))^2] - (Esgn(x_i - \xi_0))^2 \\ &= \delta_n + (1 - \delta_n) - (1 - 2\delta_n)^2 \\ &= 4\delta_n(1 - \delta_n) \qquad \tau_1 \le i < \tau_2 \end{aligned} \qquad (2.9)$$

from (2.8), (2.9) we have

$$\begin{aligned} E(S_{\tau_1-1} - S_{\tau_2-1}) &= -(\tau_2 - \tau_1)(2\delta_n - 1) \qquad (2.10) \\ &< 0 \end{aligned}$$

$$Var(S_{\tau_1-1} - S_{\tau_2-1}) = (\tau_2 - \tau_1)4\delta_n(1 - \delta_n) \qquad (2.11)$$

13

Employ $C.L.T$ to $S_{\tau_1-1} - S_{\tau_2-1}$

$$\frac{1}{\sqrt{(\tau_2-\tau_1)4\delta_n(1-\delta_n)}}[S_{\tau_1-1} - S_{\tau_2-1} + (\tau_2-\tau_1)(2\delta_n - 1)]$$

$$\xrightarrow{D} N(0,1) \qquad n \to \infty \qquad (2.12)$$

Note that $\tau_2 - \tau_1 = n(\lambda_2 - \lambda_1)$.

Under the condition of $(i)$ or $(ii)$ of $(2.5)$

$$-\frac{(\tau_2-\tau_1)(2\delta_n - 1)}{\sqrt{(\tau_2-\tau_1)4\delta_n(1-\delta_n)}} \to -\infty, \qquad n \to \infty, \qquad (2.13)$$

so

$$\frac{S_{\tau_1-1} - S_{\tau_2-1}}{\sqrt{(\tau_2-\tau_1)4\delta_n(1-\delta_n)}} \qquad (2.14)$$

has mean that converges to $-\infty$ and it has asymptotic variance 1.

According to Billingsley (1968), the properties of the simple random walk $S_l - S_{\tau_1-1}$ are

$$\frac{1}{\sqrt{\tau_1-1}}max_{l\leq\tau_1}(S_l - S_{\tau_1-1}) \sim |N(0,1)|, \qquad (2.15)$$

hence

$$\frac{1}{\sqrt{(\tau_2-\tau_1)4\delta_n(1-\delta_n)}}max_{l\leq\tau_1}(S_l - S_{\tau_1-1})$$

$$= \frac{\sqrt{\tau_1-1}}{\sqrt{(\tau_2-\tau_1)4\delta_n(1-\delta_n)}}\frac{1}{\sqrt{\tau_1-1}}max_{l\leq\tau_1}(S_l - S_{\tau_1-1}) \qquad (2.16)$$

$$= O_p(1).$$

Similarly, for the simple random variable $S_l - S_{\tau_2-1}, \tau_2 \leq l \leq n$,

$$\frac{1}{2\sqrt{(\tau_2-\tau_1)\delta_n(1-\delta_n)}}max_{\tau_2\leq l}(S_l - S_{\tau_2-1})$$

$$= \frac{\sqrt{n-\tau_2+1}}{2\sqrt{(\tau_2-\tau_1)\delta_n(1-\delta_n)}}\frac{1}{\sqrt{n-\tau_2+1}}max_{\tau_2\leq l}(S_l - S_{\tau_2-1}) \qquad (2.17)$$

$$= O_p(1).$$

14

Combine (2.14) (2.16) (2.17), we get

$$P\{max_{l<\tau_1}(S_l - S_{\tau_2-1}) \geq max_{\tau_2 \leq l}(S_l - S_{\tau_2-1})\} \to 0 \qquad as \; n \to \infty \qquad (2.18)$$

(b) $\tau_1 \leq l < \tau_2$.

First consider the case (i) of (2.5), that is when $\delta_n = \delta$ for all $n$. It is a well known fact (see, page 72 in Billingsley (1968)) that

$$\frac{1}{b_n} max_{\tau_2 \leq l \leq n}(S_l - S_{\tau_2-1}) \xrightarrow{D} |N(0,1)|, \qquad n \to \infty, \qquad (2.19)$$

where $b_n = \sqrt{n - \tau_2 + 1}$.

Let $\{\epsilon_n\}$ such that $\epsilon_n b_n \to \infty$ and $\epsilon_n \to 0$.

Denote $A = \{max_{\tau_1 \leq l < \tau_2}(S_l - S_{\tau_2-1}) \geq max_{\tau_2 \leq l \leq n}(S_l - S_{\tau_2-1})\}$.

Then we have

$$P\{\frac{1}{b_n} max_{\tau_2 \leq l \leq n}(S_l - S_{\tau_2-1}) < \epsilon_n\} = O(\epsilon_n). \qquad (2.20)$$

and

$$P\{max_{\tau_1 \leq l < \tau_2}(S_l - S_{\tau_2-1}) \geq max_{\tau_2 \leq l \leq n}(S_l - S_{\tau_2-1})\}$$

$$= P\{A \cap (\frac{1}{b_n} max_{\tau_2 \leq l \leq n}(S_l - S_{\tau_2-1}) \geq \epsilon_n)\}$$

$$+ P\{A \cap (\frac{1}{b_n} max_{\tau_2 \leq l \leq n}(S_l - S_{\tau_2-1}) < \epsilon_n)\} \qquad (2.21)$$

$$\leq P\{max_{\tau_1 \leq l < \tau_2}(S_l - S_{\tau_2-1}) \geq b_n \epsilon_n\}$$

$$+ P\{\frac{1}{b_n} max_{\tau_2 \leq l \leq n}(S_l - S_{\tau_2-1}) < \epsilon_n\}.$$

The first term in (2.21) can be written as

$$P\{\max_{1 \leq j \leq (\tau_2-\tau_1)} \sum_{i=1}^{j} \eta_i \geq b_n \epsilon_n\}, \qquad (2.22)$$

15

where

$$\eta_i = -sgn(x_{i+\tau_1-1} - \xi_0) \qquad 1 \leq i \leq (\tau_2 - \tau_1),$$

Denote $\tilde{S}_j(\omega) = \sum_{i=1}^{j} \eta_i(\omega)$.

By the Strong Law of Large Numbers

$$P\{\omega : \lim_{j \to \infty} \frac{\tilde{S}_j(\omega)}{j} = a\} = 1 \tag{2.23}$$

where $a = E\eta_i = 1 - 2\delta < 0$.

i.e. $\exists \Omega_0 \subset \Omega$, $P(\Omega_0) = 0$ s.t. for all $\omega \notin \Omega_0$

$$\frac{\tilde{S}_j(\omega)}{j} \to a < 0, \qquad j \to \infty, \tag{2.24}$$

i.e. for all $\omega \notin \Omega_0$, $\exists j_0 = j_0(\omega)$ s.t.

$$\tilde{S}_j(\omega) < 0, \qquad for\ all\ j \geq j_0. \tag{2.25}$$

Combine (2.25) and $b_n \epsilon_n \to \infty$, for every $\omega \notin \Omega_0$

$$max_{1 \leq j \leq j_0} \tilde{S}_j(\omega) < b_n \epsilon_n, \qquad when\ n\ is\ large. \tag{2.26}$$

so

$$P\{\omega : max_{1 \leq j \leq (\tau_2 - \tau_1 - 1)} \tilde{S}_j > b_n \epsilon_n\} \to 0, \qquad n \to \infty. \tag{2.27}$$

Using (2.20) and (2.27), we get

$$P\{max_{\tau_1 \leq l < \tau_2}(S_l - S_{\tau_2-1}) \geq max_{\tau_2 \leq l}(S_l - S_{\tau_2-1})\} \to 0, \qquad as\ n \to \infty. \tag{2.28}$$

The last case is $\tau_1 \leq l < \tau_2$ and $\sqrt{n}|\delta_n - \frac{1}{2}| \to 0$ as $n \to \infty$.

16

Let $r$ satisfy $\lambda_1 < r < \lambda_2$ and $rn$ is an integer.

$$max_{\tau_1 \leq l \leq rn}(S_l - S_{\tau_2-1}) = max_{\tau_1 \leq l \leq rn}(S_l - S_{rn}) + (S_{rn} - S_{\tau_2-1}). \qquad (2.29)$$

$$E(S_{rn} - S_{\tau_2-1}) = (\tau_2 - 1 - rn)(1 - 2\delta_n)$$

$$\approx (\lambda_2 - r)(1 - 2\delta_n)n.$$

$$Var(S_{rn} - S_{\tau_2-1}) = (\tau_2 - 1 - rn)4\delta_n(1 - \delta_n)$$

$$\approx (\lambda_2 - r)4\delta_n(1 - \delta_n)n$$

$$= dn.$$

where $d = (\lambda_2 - r)4\delta_n(1 - \delta_n)$.

The same way as in $(a)$ by C.L.T.

$$\frac{(S_{rn} - S_{\tau_2-1}) - (\tau_2 - 1 - rn)(1 - 2\delta_n)}{\sqrt{4\delta_n(1 - \delta_n)(\tau_2 - 1 - rn)}} \xrightarrow{D} N(0,1), \qquad n \to \infty. \qquad (2.30)$$

On the other hand, we have

$$\frac{1}{\sqrt{dn}} max_{\tau_1 \leq l \leq rn}(S_l - S_{rn})$$

$$= \frac{\sqrt{r - \lambda_1}}{\sqrt{d}} \frac{1}{\sqrt{n(r-\lambda_1)}} max_{\tau_1 \leq l \leq rn}(S_l - S_{rn}) \qquad (2.31)$$

$$= O_p(1),$$

and

$$\frac{1}{\sqrt{dn}} max_{\tau_2 \leq l \leq n}(S_l - S_{\tau_2-1})$$

$$= \frac{\sqrt{1 - \lambda_2}}{\sqrt{d}} \frac{1}{\sqrt{n - \tau_2}} max_{\tau_2 \leq l \leq n}(S_l - S_{\tau_2-1}) \qquad (2.32)$$

$$= O_p(1).$$

17

The mean of the dominating term in (2.29) is

$$E\left(\frac{S_{rn} - S_{\tau_2-1}}{\sqrt{dn}}\right) = \frac{(\tau_2 - 1 - rn)(1 - 2\delta_n)}{\sqrt{dn}} \qquad (2.33)$$

$$\to -\infty, \qquad as \ n \to \infty.$$

Join (2.30), (2.31), (2.32), (2.33), we have

$$P\{max_{\tau_1 \leq l \leq rn}(S_l - S_{\tau_2-1}) \geq max_{\tau_2 \leq l \leq n}(S_l - S_{\tau_2-1})\} \to 0, \qquad n \to \infty. \ (2.34)$$

Now consider a sequence $r_{(n)} \nearrow \lambda_2$, $r_{(n)}n$ is an integer. We want to show that

$$P\{max_{r_{(n)}n \leq l < \tau_2}(S_l - S_{\tau_2-1}) \geq max_{\tau_2 \leq l}(S_l - S_{\tau_2-1})\} \to 0, \qquad n \to \infty. \ (2.35)$$

Dividing by $\sqrt{n}$ on both sides of the inequality of (2.35) and employing (2.32), we can write (2.35) equivalently as

$$P\{\frac{1}{\sqrt{n}}max_{1 \leq j \leq \tilde{r}_{(n)}n}\tilde{S}_j \geq O_p(1)\} \to 0, \qquad n \to \infty. \qquad (2.36)$$

where $\tilde{S}_j$ has independent, identically distributed terms with mean $(1 - 2\delta_n) < 0$ and finite variance $4\delta_n(1 - \delta_n)$ and $\tilde{r}_{(n)} = \frac{[\lambda_2 n]}{n} - r_{(n)} \searrow 0$. But

$$max_{1 \leq j \leq \tilde{r}_{(n)}n}\tilde{S}_j \leq max_{1 \leq j \leq \tilde{r}_{(n)}n}(\tilde{S}_j + j(2\delta_n - 1))$$

$$= max_{1 \leq j \leq \tilde{r}_{(n)}n}\tilde{S}_j^*, \qquad (2.37)$$

where $\tilde{S}_j^*$ has terms with mean zero and

$$\frac{1}{\sqrt{n}}max_{1 \leq j \leq \tilde{r}_{(n)}n}\tilde{S}_j^* = \frac{\sqrt{\tilde{r}_{(n)}n}}{\sqrt{n}}\frac{1}{\sqrt{\tilde{r}_{(n)}n}}max_{1 \leq j \leq \tilde{r}_{(n)}n}S_j^*$$

$$= O_p(\sqrt{\tilde{r}_{(n)}}). \qquad (2.38)$$

18

From (2.36), (2.38), using that $\hat{r}_{(n)} \to 0$, we can conclude (2.35) is true.

It remains to prove that

$$P\{\frac{1}{\sqrt{n}}max_{rn<l<r_{(n)}n}(S_l - S_{\tau_2-1}) \geq O_p(1)\} \to 0, \qquad n \to \infty. \qquad (2.39)$$

We may write

$$max_{rn<l<r_{(n)}n}\frac{1}{\sqrt{n}}(S_l - S_{\tau_2-1})$$

$$= max_{rn<l<r_{(n)}n}\frac{1}{\sqrt{n}}\{(S_l - S_{r_{(n)}n}) + (S_{r_{(n)}n} - S_{\tau_2-1})\} \qquad (2.40)$$

$$\leq max_{rn<l<r_{(n)}n}\frac{1}{\sqrt{n}}\{(S_l - S_{r_{(n)}n} + (r_{(n)}n - l)(2\delta_n - 1)) + (S_{r_{(n)}n} - S_{[\lambda_2 n]-1})\},$$

Similarly as before

$$max_{rn<l<r_{(n)}n}\frac{1}{\sqrt{n}}\{(S_l - S_{r_{(n)}n}) + (r_{(n)}n - l)(2\delta_n - 1)\}$$

$$= O_p(1), \qquad n \to \infty. \qquad (2.41)$$

we have

$$E\{\frac{1}{\sqrt{n}}(S_{r_{(n)}n} - S_{[\lambda_2 n]-1})\}$$

$$\approx \sqrt{n}(\lambda_2 - r_{(n)})(1 - 2\delta_n) < 0, \qquad (2.42)$$

and

$$S.D.\{\frac{1}{\sqrt{n}}(S_{r_{(n)}n} - S_{[\lambda_2 n]-1})\}$$

$$\approx \sqrt{4\delta_n(1 - \delta_n)(\lambda_2 - r_{(n)})} \qquad (2.43)$$

Hence, if we choose a seguence $r_{(n)} \to \lambda_2$ such that we also have $\sqrt{n(\lambda_2 - r_{(n)})}(1 - 2\delta_n) \to -\infty$. We get that (2.39) holds.

Now consider

$$max_{\tau_1 \leq l \leq \tau_2}(S_l - S_{\tau_2-1})$$

$$= max\{max_{\tau_1 \leq l<rn}(S_l - S_{\tau_2-1}),$$

$$max_{rn\leq l<r_{(n)}n}(S_l - S_{\tau_2-1}), max_{r_{(n)}n\leq l<\tau_2}(S_l - S_{\tau_2-1})\} \qquad (2.44)$$

$$= max\{A_1, A_2, A_3\},$$

19

$$A_1 = max_{\tau_1 \le l < rn}(S_l - S_{\tau_2 - 1}),$$

$$A_2 = max_{rn \le l < r_{(n)}n}(S_l - S_{\tau_2 - 1}),$$

$$A_3 = max_{r_{(n)}n \le l < \tau_2}(S_l - S_{\tau_2 - 1}).$$

Also, denote $B = n \quad r_{\tau_2 \le l}(S_l - S_{\tau_2 - 1})$ and we have

$$P\{max(A_1, \ A_2, \ A_3) > B\}$$

$$= P\{A_1 > B \ or \ A_2 > B \ or \ A_3 > B\}$$

$$\le P\{A_1 > B\} + P\{A_2 > B\} + P\{A_3 > B\}$$

$$\to 0, \qquad\qquad n \to \infty.$$

(2.45)

Combining the above (2.34), (2.35) and (2.39), we get

$$P\{max_{\tau_1 \le l < \tau_2}(S_l - S_{\tau_2 - 1}) \ge max_{\tau_2 \le l}(S_l - S_{\tau_2 - 1})\} \to 0, \qquad n \to \infty, \quad (2.46)$$

as claimed, and the Lemma 2.1 is proved. □

**Lemma 2.2** *Under hypothesis* $H_1$:

$$P\{max_{k \ge \tau_1}(S_{\tau_1} - S_k) \ge max_{k < \tau_1}(S_{\tau_1} - S_k)\} \to 0 \qquad as \ n \to \infty \qquad (2.47)$$

**Proof:** It can be proved by using the method in the proof of Lemma 2.1 and observing the symmetry. □

To estimate $\tau_1$, $\tau_2$, it is customary to use

$$(\hat{\tau}_1(n), \hat{\tau}_2(n))$$

$$= ar\,gmax_{k < l}(S_{l-1} - S_{k-1})$$

(2.48)

$$= \{(min(k), max(l)) : \sum_{i=k}^{l-1} sgn(x_i - \xi_0) = max_{1 \le u < v \le n} \sum_{i=u}^{v-1} sgn(x_i - \xi_0)\}.$$

20

Combining Lemma 2.1 and Lemma 2.2 together, when $n \to \infty$ we get that

$$P\{\hat{\tau}_1 < \tau_1\} \to 1,$$

$$P\{\hat{\tau}_2 > \tau_2\} \to 1.$$

(2.49)

We can write

$$max_{k<l}(S_{l-1} - S_{k-1})$$

$$= max_{k<l}\{(S_{l-1} - S_{\tau_2-1}) + (S_{\tau_2-1} - S_{\tau_1}) + (S_{\tau_1} - S_{k-1})\}$$

(2.50)

$$= max_{k<l}\{max_l(S_{l-1} - S_{\tau_2-1}) + max_k(S_{\tau_1} - S_{k-1})\} + (S_{\tau_2-1} - S_{\tau_1}),$$

and

$$max_l(S_{l-1} - S_{\tau_2-1}) \quad = \quad max\{max_{l<\tau_2-1}(S_{l-1} - S_{\tau_2-1}), max_{\tau_2 \le l}(S_{l-1} - S_{\tau_2-1})\}$$

$$max_k(S_{\tau_1} - S_{k-1}) \quad = \quad max\{max_{k \ge \tau_1}(S_{\tau_1} - S_{k-1}), max_{k<\tau_1}(S_{\tau_1} - S_{k-1})\}$$

To get the asymptotic distribution of $\sqrt{n}max_{k<l}(S_{l-1} - S_{k-1})$ under the alternative hypothesis, we have by the properties of simple symmetric random walk that

$$\frac{1}{\sqrt{n\lambda_1}}max_{1 \le i < \tau_1} \sum_{j=1}^{i} sgn(x_j - \xi_0) \xrightarrow{D} |N_1(0,1)|, \qquad n \to \infty, \qquad (2.51)$$

and

$$\frac{1}{\sqrt{n(1-\lambda_1)}}max_{\tau_2 \le i \le n} \sum_{j=1}^{i} sgn(x_j - \xi_0) \xrightarrow{D} |N_2(0,1)|, \qquad n \to \infty, \qquad (2.52)$$

where $N_1$ and $N_2$ are indepedent standard normal random variables.

Furthermore (2.12) can also be written as

$$\frac{1}{\sqrt{4n\delta_n(1-\delta_n)(\lambda_2-\lambda_1)}}\{ \sum_{\tau_1 \le i < \tau_2} sgn(x_i - \xi_0) - \mu_n\} \to N_3(0,1), \qquad (2.53)$$

where $N_3$ is a standard normal random variable, independent of $N_1$ and $N_2$, and $\mu_n \approx n(\lambda_2 - \lambda_1)(2\delta_n - 1)$. So using Lemmas 2.1 and 2.2 and (2.49), we get that

21

the asymptotic distribution of $n^{-\frac{1}{2}}max_{k<l}(S_{l-1} - S_{k-1})$ under the alternative hypothesis is that of

$$\sqrt{\lambda_1}|N_1| + \sqrt{1 - \lambda_2}|N_2| + 2\sqrt{\delta_n(1 - \delta_n)(\lambda_2 - \lambda_1)}N_3 + M_n, \qquad (2.54)$$

where $M_n = (2\delta_n - 1)\sqrt{\frac{n(\lambda_2 - \lambda_1)}{4\delta_n(1-\delta_n)}}$. As $M_n > 0$, from (2.55) the consistency of our test follows.

Recall that

$$P\{\hat{\tau}_1(n) \leq \tau_1(n)\} \to 1,$$

$$P\{\hat{\tau}_2(n) \geq \tau_2(n)\} \to 1.$$

In the case of symmetric simple random walk $\{S_i\}$, given $\hat{\tau}_1(n) < \tau_1(n)$, $\hat{\tau}_2(n) > \tau_2(n)$, the distribution of $\hat{\tau}_1(n)$ and $\hat{\tau}_2(n)$ can be obtained. By symmetry, it is sufficient to consider

$$P\{\hat{\tau}_2(n) = j \mid \hat{\tau}_2(n) \geq \tau_2(n)\}$$

$$= P\{argmax_{\tau_2 \leq l < n}(S_l - S_{\tau_2 - 1}) = j\}$$

$$= P\{S_k < 0; \ k = 1, \cdots, j\}P\{S_k - S_j \leq 0; \ k = j + 1, \cdots, n\} \qquad (2.55)$$

$$= P\{S_k < 0; \ k = 1, \cdots, j\}P\{S_l \leq 0; \ l = 1, \cdots, n - j\}.$$

using time reversal and indepedent increments property, the expressions for the two factors are well known:

$$P\{S_k < 0; \ k = 1, \cdots, j\}$$

$$= \begin{cases} \frac{1}{2}P\{S_j = 0\} = \binom{j}{j/2}(\frac{1}{2})^{j+1} & j \text{ is even} \\ \\ \sum_{b=-1}^{-j} \frac{b}{j}P\{S_j = b\} & j \text{ is odd} \end{cases} \qquad (2.56)$$

and

$$P\{S_k \leq 0; \ k = 1, \cdots, L\}$$

$$= \begin{cases} P\{S_L = 0\} & j \text{ is even.} \\ P\{S_{L+1} = 0\} & j \text{ is odd.} \end{cases} \tag{2.57}$$

From (2.49), $P\{\hat{\tau}_1 > \tau_1\} \to 0$ and $P\{\hat{\tau}_2 < \tau_2\} \to 0$, we conclude that $(\hat{\tau}_1(n), \hat{\tau}_2(n)) = argmax_{k<l}(S_{l-1} - S_{k-1})$ is not a good estimation for the change-points $(\tau_1(n), \tau_2(n))$.

We will see in the next chapter that we can get a better estimator for the change-points $(\tau_1, \tau_2)$ even though we do not know the initial value of the median $\xi_0$, but have to use its estimator under $H_0$.

# Chapter 3

# SIGN TEST FOR THE CHANGE-POINT PROBLEM WITH UNKNOWN MEDIAN

Let $x_1, \cdots, x_n$ be a sequence of independent continuous random variables. We want to consider the following hypothesis test for the change-point problem.

$$H_0 : \quad x_i \stackrel{D}{=} Y; \quad i = 1, \cdots, n$$

$$H_1 : \quad x_i \stackrel{D}{=} Y; \quad i = 1, \cdots, \tau_1 - 1, \tau_2, \cdots, n \tag{3.1}$$

$$x_i \stackrel{D}{=} Z; \quad i = \tau_1, \cdots, \tau_2 - 1 \quad Y \neq Z.$$

$H_1$ is called the epidemic or square-wave alternative. We assume $Y$ and $Z$ have distribution functions $F(x)$ and $G(x)$ respectively where $F(x) \neq G(x)$ at least for some $x$ and

$$F^{-1}(\frac{1}{2}) = \xi_0, \qquad G^{-1}(\frac{1}{2}) = \xi_A,$$

where $\xi_0$ and $\xi_A$ are unknown. Also unknown are parameters $\tau_1$ and $\tau_2$ the change-points. We assume $\tau_1 = [n\lambda_1]$, $\tau_2 = [n\lambda_2]$, for some $0 < \lambda_1 < \lambda_2 < 1$.

First we consider the distribution of test statistic under $H_0$.

Let

$$\hat{\xi}_n = median\{x_1, \cdots, x_n\}, \tag{3.2}$$

and

$$V_n(u) = n^{-\frac{1}{2}} \sum_{1 \leq i \leq [nu]} sgn(x_i - \hat{\xi}_n), \qquad 0 \leq u \leq 1. \tag{3.3}$$

We employ the following statistic

$$
\begin{aligned}
M(n) &= max_{k<l}|V_n(\frac{l-1}{n}) - V_n(\frac{k-1}{n})| \\
&= max_{1 \leq k < l \leq n}|n^{-\frac{1}{2}} \sum_{k \leq i < l} sgn(x_i - \hat{\xi}_n)|. 
\end{aligned}
\tag{3.4}
$$

Under $H_0$, as $n \to \infty$

$$V_n(u) = n^{-\frac{1}{2}} \sum_{1 \leq i \leq [nu]} sgn(x_i - \hat{\xi}_n) \xrightarrow{D} B(u), \qquad 0 \leq u \leq 1. \tag{3.5}$$

where $\{B(u); 0 \leq u \leq 1\}$ is a Brownian bridge (see Billingsley 1968).

The asymptotic distribution of $M(n)$, under $H_0$, when $\lambda_2 = 1$, is the same as that of the Kolmogorov-Smirnov test of the equal sample sizes, $m = \frac{n}{2}$ in case $n$ is even and $\frac{n-1}{2}$ when $n$ is odd, (in this case $sgn(0) = 0$). This case is the at-most-one-change alternative

$$H_1^* : \quad x_1, \cdots, x_{\tau-1} \sim F(x),$$

$$x_\tau, \cdots, x_n \sim G(x), \quad F(x) \neq G(x).$$

Tables for this distribution can be used.

When $0 < \lambda_1 < \lambda_2 < 1$ is assumed, i.e. when we test $H_0$ against $H_1$, from Gnedenko (1954, $p.53$), we have the exact formula:

$$P\{M(n) < z\}$$

$$= 1 + \frac{2}{\binom{2m}{m}}\{[\alpha \sum_{s=1}^{[\frac{m}{\alpha+1}]} \binom{2m}{m-s(\alpha+1)}) - (\alpha - 1)\sum_{s=1}^{[\frac{m}{\alpha}]}\binom{2m}{m-s\alpha})] \tag{3.6}$$

$$- [\sum_{i=1}^{\alpha} \sum_{s=1}^{[\frac{m+i}{\alpha+1}]}\binom{2m}{m+i-s(\alpha+1)}) - \sum_{i=1}^{\alpha-1}\sum_{s=1}^{[\frac{m+i}{\alpha}]}\binom{2m}{m+i-s\alpha})]\},$$

where $m = [\frac{n}{2}]$, $\alpha = [z\sqrt{n}]$.

Now we consider the distribution of test statistic under the alternative hypothesis. Suppose

$$\delta = P\{x_{\tau_1} \leq \xi_0\} > \frac{1}{2} \tag{3.7}$$

where $\xi_0 = F^{-1}(\frac{1}{2})$. We use the notation

$$\hat{\xi}_n = median\{x_1, \cdots, x_n\}, \qquad under\ H_1. \tag{3.8}$$

$$H(x) = [1 - (\lambda_2 - \lambda_1)]F(x) + (\lambda_2 - \lambda_1)G(x) \tag{3.9}$$

Let $\eta_0 = H^{-1}(\frac{1}{2})$. Note that $\eta_0$ is an unknown number that does not depend on $n$. We have

$$G^{-1}(\frac{1}{2}) < H^{-1}(\frac{1}{2}) < F^{-1}(\frac{1}{2}). \tag{3.10}$$

Let

$$s = P\{x_i \leq \eta_0\} < \frac{1}{2} \qquad i = 1, \cdots, \tau_1 - 1, \tau_2, \cdots, n,$$
$$s' = P\{x_i \leq \eta_0\} > \frac{1}{2} \qquad i = \tau_1, \cdots, \tau_2 - 1. \tag{3.11}$$

By the well known property of quantile process (see e.g Csörgő-Révész 1981), we get

$$\hat{\xi}_n - \eta_0 = O_p(n^{-\frac{1}{2}}), \quad for \ large \ n. \tag{3.12}$$

From (3.11), we get that

$$
\begin{aligned}
r_n &= P\{x_i \le \hat{\xi}_n\} < \tfrac{1}{2}, \quad i = 1, \cdots, \tau_1 - 1, \tau_2, \cdots, n. \\
r'_n &= P\{x_i \le \hat{\xi}_n\} > \tfrac{1}{2}, \quad i = \tau_1, \cdots, \tau_2 - 1.
\end{aligned} \tag{3.13}
$$

**Lemma 3.1** *Under the alternative hypothesis $H_1$, for $0 < u < 1$, we have*

$$E\{V_n(u)\} \tag{3.14}$$

$$
= \begin{cases}
n^{\frac{1}{2}} u C_1 + O_p(n^{-\frac{1}{2}}), & u < \lambda_1, \\[2mm]
n^{\frac{1}{2}} [\lambda_1 C_1 + (u - \lambda_1) C_2] + O_p(n^{-\frac{1}{2}}), & \lambda_1 \le u < \lambda_2, \\[2mm]
n^{\frac{1}{2}} [(u - \lambda_2 + \lambda_1) C_1 + (\lambda_2 - \lambda_1) C_2] + O_p(n^{-\frac{1}{2}}), & \lambda_2 \le u.
\end{cases}
$$

$$Var\{n^{\frac{1}{2}} V_n(u)\} \tag{3.15}$$

$$
= \begin{cases}
n u D_1 + \frac{nu(nu-1)}{2} O_p(n^{-\frac{1}{2}}), & u < \lambda_1, \\[2mm]
n\lambda_1 D_1 + (u - \lambda_1) n D_2 + \frac{nu(nu-1)}{2} O_p(n^{-\frac{1}{2}}), & \lambda_1 \le u < \lambda_2, \\[2mm]
(u + \lambda_1 - \lambda_2) n D_1 + (\lambda_2 - \lambda_1) n D_2 \frac{nu(nu-1)}{2} O_p(n^{-\frac{1}{2}}), & \lambda_2 \le u.
\end{cases}
$$

*where $C_1 = 1 - 2r_n$, $C_2 = 1 - 2r'_n$, $D_1 = 4r_n(1 - r_n)$, $D_2 = 4r'_n(1 - r'_n)$*

**Proof:** For the sake of brevity, we give the proof for $u \ge \lambda_2$. The proof for other cases are quite similar. Hence they will be omitted. Note that $x_i \overset{D}{=} Y$, for $i < [n\lambda_1]$, or $i \ge [n\lambda_2]$, and $x_i \overset{D}{=} Z$ for $[n\lambda_1] \le i < [n\lambda_2]$. To calculate the mean,

we have

$$E\{V_n(u)\}$$

$$= E\{n^{-\frac{1}{2}} \sum_{1 \leq i \leq [nu]} sgn(x_i - \hat{\xi}_n\}$$

$$= n^{-\frac{1}{2}}[\sum_{1 \leq i < [n\lambda_1]} E\{sgn(x_i - \hat{\xi}_n)\} + \sum_{[n\lambda_1] \leq i < [n\lambda_2]} E\{sgn(x_i - \hat{\xi}_n)\}$$

$$+ \sum_{[n\lambda_2] \leq i < [nu]} E\{sgn(x_i - \hat{\xi}_n)\}]$$

$$= n^{-\frac{1}{2}}[(n\lambda_1 - 1)E\{sgn(Y - \hat{\xi}_n)\} + n(\lambda_2 - \lambda_1)E\{sgn(Z - \hat{\xi}_n)\}$$

$$+ (nu - n\lambda_2 + 1)E\{sgn(Y - \hat{\xi}_n)\}] + O_p(n^{-\frac{1}{2}}) \qquad (3.16)$$

$$= n^{-\frac{1}{2}}[n(\lambda_1 - 1)C_1 + n(\lambda_2 - \lambda_1)C_2 + (nu - n\lambda_2 + 1)C_1] + O_p(n^{-\frac{1}{2}})$$

$$= n^{\frac{1}{2}}[\lambda_1 C_1 + (\lambda_2 - \lambda_1)C_2 + (u - \lambda_2)C_1] + n^{-\frac{1}{2}}([nu] - nu + 1) + O_p(n^{-\frac{1}{2}})$$

$$= n^{\frac{1}{2}}[(n - \lambda_2 + \lambda_1)C_1 + (\lambda_2 - \lambda_1)C_2] + O_p(n^{-\frac{1}{2}}).$$

with

$$E\{sgn(Y - \hat{\xi}_n)\} = P\{Y > \hat{\xi}_n\} + (-1)P\{Y \leq \hat{\xi}_n\}$$

$$= 1 - 2P\{Y \leq \hat{\xi}_n\} \qquad (3.17)$$

$$= 1 - 2r_n.$$

and $E\{sgn(Z - \hat{\xi}_n)\} = 1 - 2r_n'$, we get (3.14)

For the variance calculation, we have

$$Var\{sgn(Y - \hat{\xi}_n)\} = E\{[sgn(Y > \hat{\xi}_n)]^2\} - [E\{sgn(Y - \hat{\xi}_n\}]^2$$

$$= 1 - (1 - 2r_n)^2 \qquad (3.18)$$

$$= 4r_n(1 - r_n).$$

28

$$Var\{sgn(Z - \hat{\xi}_n)\} = 4r'_n(1 - r'_n)$$

and

$$
\begin{aligned}
Var\{n^{\frac{1}{2}} V_n(u)\} &= Var\{ \sum_{1 \le i \le [nu]} sgn(x_i - \hat{\xi}_n)\} \\
&= \sum_{1 \le i \le [nu]} Var\{sgn(x_i - \hat{\xi}_n)\} \\
&\quad + \sum_{1 \le i < j \le [nu]} Cov\{sgn(x_i - \hat{\xi}_n), sgn(x_j - \hat{\xi}_n)\}
\end{aligned}
\tag{3.19}
$$

(1) We consider first $i$ and $j < [n\lambda_1]$ or $\ge [n\lambda_2]$

$$Cov\{sgn(x_i - \hat{\xi}_n), sgn(x_j - \hat{\xi}_n)\}$$

$$
\begin{aligned}
&= E\{sgn(x_i - \hat{\xi}_n)sgn(x_j - \hat{\xi}_n)\} - E\{sgn(x_i - \hat{\xi}_n)\}E\{sgn(x_j - \hat{\xi}_n)\} \\
&= E\{sgn(x_i - \eta_0)sgn(x_j - \eta_0)\} - (1 - 2r_n)^2 \\
&\quad + E\{sgn(x_i - \hat{\xi}_n)sgn(x_j - \hat{\xi}_n) - sgn(x_i - \eta_0)sgn(x_j - \eta_0)\} \\
&= E\{sgn(x_i - \eta_0)\}E\{sgn(x_j - \eta_0)\} - (1 - 2r_n)^2 \\
&\quad + E\{sgn(x_i - \hat{\xi}_n)sgn(x_j - \hat{\xi}_n) - sgn(x_i - \eta_0)sgn(x_j - \eta_0)\} \\
&= (1 - 2s)^2 - (1 - 2r_n)^2 \\
&\quad + E\{sgn(x_i - \hat{\xi}_n)sgn(x_j - \hat{\xi}_n) - sgn(x_i - \eta_0)sgn(x_j - \eta_0)\}.
\end{aligned}
\tag{3.20}
$$

Because when $x_i$ and $x_j \notin (\hat{\xi}_n \wedge \eta_0, \hat{\xi}_n \vee \eta_0)$

$$sgn(x_i - \hat{\xi}_n)sgn(x_j - \hat{\xi}_n) - sgn(x_i - \eta_0)sgn(x_j - \eta_0) = 0$$

we have

$$|E\{sgn(x_i - \hat{\xi}_n)sgn(x_j - \hat{\xi}_n) - sgn(x_i - \eta_0)sgn(x_j - \eta_0)\}|$$

$$\le \int_{(\hat{\xi}_n \wedge \eta_0, \hat{\xi}_n \vee \eta_0) \times (\hat{\xi}_n \wedge \eta_0, \hat{\xi}_n \vee \eta_0)} 2 dF_{i,j} \tag{3.21}$$

$$= O_p(n^{-\frac{1}{2}})$$

where $F_{i,j}$ is the joint d.f of $x_i$ and $x_j$.

Because $F_{i,j}$ is continuous and by (3.12)

$$(1 - 2s)^2 - (1 - 2r_n)^2 = 4(1 - r_n - s)(r_n - s) \tag{3.22}$$

As in (3.21), we get

$$\begin{aligned}
|r_n - s| &= |P\{x_i \le \hat{\xi}_n\} - P\{x_i \le \eta_0\}| \\
&\le P\{\hat{\xi}_n \wedge \eta_0 \le x_i < \hat{\xi}_n \vee \eta_0\} \\
&= \int_{(\hat{\xi}_n \wedge \eta_0, \hat{\xi}_n \vee \eta_0)} dF_i \\
&= O_p(n^{-\frac{1}{2}})
\end{aligned} \tag{3.23}$$

Combining (3.21), (3.22), (3.23), we have

$$Cov\{sgn(x_i - \hat{\xi}_n), sgn(x_j - \hat{\xi}_n)\} = O_p(n^{-\frac{1}{2}}), \tag{3.24}$$

for $i$ and $j < n\lambda_1$ or $\ge n\lambda_2$

(ii) For $i < n\lambda_1$ or $\ge n\lambda_2$ and $n\lambda_1 \le j < n\lambda_2$

$$\begin{aligned}
&Cov\{sgn(x_i - \hat{\xi}_n), sgn(x_j - \hat{\xi}_n)\} \\
&= E\{sgn(x_i - \hat{\xi}_n)sgn(x_j - \hat{\xi}_n)\} - E\{sgn(x_i - \hat{\xi}_n)\}E\{sgn(x_j - \hat{\xi}_n)\} \\
&= E\{sgn(x_i - \eta_0)sgn(x_j - \eta_0)\} - (1 - 2r_n)(1 - 2r_n') \\
&\quad + E\{sgn(x_i - \hat{\xi}_n)sgn(x_j - \hat{\xi}_n) - sgn(x_i - \eta_0)sgn(x_j - \eta_0)\} \\
&= (1 - 2s)(1 - 2s') - (1 - 2r_n)(1 - 2r_n') \\
&\quad + E\{sgn(x_i - \hat{\xi}_n)sgn(x_j - \hat{\xi}_n) - sgn(x_i - \eta_0)sgn(x_j - \eta_0)\}
\end{aligned} \tag{3.25}$$

Because

$$|ab - cd| \le |b||a - c| + |c||b - d|,$$

30

. have

$$|(1 - 2s)(1 - 2s') - (1 - 2r_n)(1 - 2r_n')|$$

$$\leq 2|1 - 2s||r_n - s| + 2|1 - 2r_n||r_n' - s'| \tag{3.26}$$

As in (3.23), we have

$$|r_n' - s'| = O_p(n^{-\frac{1}{2}}). \tag{3.27}$$

And

$$(1 - 2s)(1 - 2s') - (1 - 2r_n)(1 - 2r_n') = O_p(n^{-\frac{1}{2}}) \tag{3.28}$$

Using (3.21), (3.25), (3.28), we have proved

$$Cov\{sgn(x_i - \hat{\xi}_n), sgn(x_j - \hat{\xi}_n)\} = O_p(n^{-\frac{1}{2}}), \tag{3.29}$$

for $i < [n\lambda_1]$ or $\leq [n\lambda_2]$ and $[n\lambda_1] \leq j < [n\lambda_2]$. So, we have

$$Var\{n^{\frac{1}{2}}V_n(u)\}$$

$$= (n\lambda_1 - 1)D_1 + n(\lambda_2 - \lambda_2)D_2 + ([nu] - n\lambda_2 + 1)D_1 + \frac{[nu]([nu] - 1)}{2}O_p(n^{-\frac{1}{2}})$$

$$= n(\lambda_1 - \lambda_2 + u)D_1 + n(\lambda_2 - \lambda_2)D_2 + ([nu] - nu)D_1 + \frac{nu(nu - 1)}{2}O_p(n^{-\frac{1}{2}}).$$

$$= n(\lambda_1 - \lambda_2 + u)D_1 + n(\lambda_2 - \lambda_2)D_2 + \frac{nu(nu - 1)}{2}O_p(n^{-\frac{1}{2}})\}. \tag{3.30}$$

where

$$\frac{[nu]([nu] - 1)}{2}O_p(n^{-\frac{1}{2}}) = \frac{nu(nu - 1)}{2}O_p(n^{-\frac{1}{2}}) \tag{3.31}$$

Let $\hat{\tau}_1(n)$, $\hat{\tau}_2(n)$ be the estimators of unknown change-points $\tau_1$, $\tau_2$. Similarly

as in Chapter 2, we use the following estimators $\hat{\tau}_1(n)$, $\hat{\tau}_2(n)$ for $\tau_1$ and $\tau_2$,

$$(\hat{\tau}_1(n), \ \hat{\tau}_2(n))$$

$$= \quad argmax_{k<l}|S_{l-1} - S_{k-1}| \qquad\qquad (3.32)$$

$$= \quad \{(min\{k\}, max\{l\}: |\sum_{i=k}^{l-1} sgn(x_i - \hat{\xi}_n)| = max_{1\leq u<v\leq n}|\sum_{i=u}^{v-1} sgn(x_i - \hat{\xi}_n)|\}.$$

**Theorem 3.1** *Under the alternative hypothesis $H_1$, if (3.7) is true, then*

$$|\hat{\tau}_1(n) - \tau_1(n)| + |\hat{\tau}_2(n) - \tau_2(n)| = O_p(1), \qquad\qquad (3.33)$$

*where $\hat{\tau}_1(n)$, $\hat{\tau}_2(n)$ are the estimators of unknown change-points of $\tau_1$, $\tau_2$. Furthermore*

$$\frac{1}{\sigma\sqrt{(\lambda_2 - \lambda_1)n}}\{\sum_{j=\hat{\tau}_1}^{\hat{\tau}_2} sgn(x_j - \hat{\xi}_n) - (\lambda_2 - \lambda_1)nC_2\} \xrightarrow{D} N(0,1). \qquad (3.34)$$

*where $C_2$ is defined in (3.15) and $\sigma = 2\sqrt{s'(1-s')}$.*

**Proof:** For the sake of brevity, we give the proof for the alternative hypothesis $H_2 \subset H_1$, where $\lambda_2 = 1$. The more general claim will easily follow from this case. From now on we will drop the index of $\lambda$ and $\tau$.

First we prove

$$|\hat{\tau}(n) - \tau(n)| = O_p(1) \qquad\qquad (3.35)$$

This is equivalent to

$$\lim_{k\to\infty} limsup_{n\to\infty} P\{max_{i\leq\tau-k}V_n(\tfrac{i}{n}) \geq max_{\tau-k<i<\tau+k}V_n(\tfrac{i}{n})\}$$
$$+ \lim_{k\to\infty} limsup_{n\to\infty} P\{max_{\tau+k\leq i}V_n(\tfrac{i}{n}) \geq max_{\tau-k<i<\tau+k}V_n(\tfrac{i}{n})\} = 0 \qquad (3.36)$$

We show that the first term on the left hand side of (3.36) is zero. The claim for the second term can be proved the same way by symmetry. Hence it will be omitted.

Since $x_i$, $i < \tau$ are identically distributed,

$$P\{max_{i\le\tau-k}V_n(\frac{i}{n}) \ge max_{\tau-k<l<\tau+k}V_n(\frac{l}{n})\}$$

$$\le P\{max_{i\le\tau-k}\sum_{j=1}^{i} sgn(x_j - \hat{\xi}_n) \ge max_{\tau-k<l<\tau}\sum_{j=1}^{l} sgn(x_j - \hat{\xi}_n)\}$$

$$= P\{\exists i, i \le \tau - k : \sum_{j=1}^{i} sgn(x_j - \hat{\xi}_n) \ge max_{\tau-k<l<\tau}\sum_{j=1}^{l} sgn(x_j - \hat{\xi}_n)\}$$

$$= P\{\exists i, i \le \tau - k : 0 \ge max_{\tau-k<l<\tau}\sum_{j=i+1}^{l} sgn(x_j - \hat{\xi}_n)\} \qquad (3.37)$$

$$= P\{\exists i, i \le \tau - k : 0 \ge \sum_{j=i+1}^{\tau-k} sgn(x_j - \hat{\xi}_n)$$

$$+max_{\tau-k<l<\tau}\sum_{j=\tau-k+1}^{l} sgn(x_j - \hat{\xi}_n)\}$$

$$= 1 - P\{\forall i, i \le \tau - k : 0 < \sum_{j=i+1}^{\tau-k} sgn(x_j - \hat{\xi}_n)$$

$$+max_{\tau-k<l<\tau}\sum_{j=\tau-k+1}^{l} sgn(x_j - \hat{\xi}_n)\}$$

$$= 1 - P\{0 < min_{1\le i\le\tau-k} \sum_{j=i+1}^{\tau-k} sgn(x_j - \hat{\xi}_n)$$

$$+max_{\tau-k<l<\tau}\sum_{j=\tau-k+1}^{l} sgn(x_j - \hat{\xi}_n)\}$$

$$= P\{0 \ge \frac{1}{\sqrt{k}}min_{1\le i\le\tau-k} \sum_{j=i+1}^{\tau-k} sgn(x_j - \hat{\xi}_n)$$

$$+\frac{1}{\sqrt{k}}max_{\tau-k<l<\tau}\sum_{j=\tau-k+1}^{l} sgn(x_j - \hat{\xi}_n)\}$$

Let $k = 1, 2, \cdots$, be a sequence and $m(k)$ another sequence, such that

$$\frac{m(k)}{\sqrt{k}} \to \infty, \qquad \frac{m(k)}{k} \to 0. \tag{3.38}$$

$$\frac{1}{\sqrt{k}} \sum_{j=1}^{m(k)} sgn(x_j - \hat{\xi}_n)$$

$$\geq \frac{1}{\sqrt{k}} \sum_{j=1}^{m(k)} sgn(x_j - \eta_0) - \frac{1}{\sqrt{k}} \sum_{j=1}^{m(k)} |sgn(x_j - \hat{\xi}_n) - sgn(x_j - \eta_0)| \tag{3.39}$$

$$\overset{D}{=} cW(\tfrac{m(k)}{k}) + \delta\tfrac{m(k)}{\sqrt{k}} - \frac{2}{\sqrt{k}} \sum_{j=1}^{m(k)} I_j + o_p(1),$$

where $c$, $\delta > 0$ are constants, $W(\cdot)$ is a Wiener process, $I_j$ is the indicator of the

event $x_j \in \{\hat{\xi}_n \wedge \eta_0, \hat{\xi}_n \vee \eta_0\}$. As $|\hat{\xi}_n - \eta_0| = O_p(n^{-\frac{1}{2}})$, we get that $I_j = O_p(n^{-\frac{1}{2}})$.

As $\frac{m(k)}{k} \to 0$, by the continuity of the Wiener process

$$min_{1 \leq j \leq m(k)} W(\frac{j}{k}) \overset{P}{\to} 0 \tag{3.40}$$

and we also have $\delta = Esgn(x_j - \eta_0) = 1 - 2s > 0$ (see (3.11)).

Now,

$$\frac{1}{\sqrt{k}} \sum_{j=1}^{m(k)} I_j = \frac{m(k)}{\sqrt{k}} O_p(n^{-\frac{1}{2}})$$

$$= o_p(1),$$

So we get

$$-[min_{1 \leq l \leq m(k)} \frac{1}{\sqrt{k}} \sum_{j=1}^{l} sgn(x_j - \hat{\xi}_n)]^- = o_p(1). \tag{3.41}$$

where $[u]^- = min(u, 0)$. It is known that

$$min_{0 < t < 1} W(t) = O_p(1). \tag{3.42}$$

Also

$$-[\frac{1}{\sqrt{k}} min_{m(k) < l \leq k} sgn(x_j - \hat{\xi}_n)]^-$$

$$= -[min_{m(k)<l\le k}\{\frac{1}{\sqrt{k}}\sum_{j=1}^{l}sgn(x_j-\eta_0)+\frac{1}{\sqrt{k}}\sum_{j=1}^{l}[sgn(x_j-\hat{\xi}_n)-sgn(x_j-\eta_0)]\}]^-$$

$$\le -[min_{0<t<1}W(t)+o_p(1)-\frac{2}{\sqrt{k}}\sum_{j=1}^{k}I_j+\delta\frac{m(k)}{\sqrt{k}}]^-. \tag{3.43}$$

As $I_j=O_p(n^{-\frac{1}{2}})$, and $\frac{m(k)}{\sqrt{k}}\to\infty$, we get that

$$-[min_{m(k)<l\le k}\frac{1}{\sqrt{k}}sgn(x_j-\hat{\xi}_n)]^-=o_p(1). \tag{3.44}$$

To show that

$$-[min_{k<l<\tau-k}\frac{1}{\sqrt{k}}\sum_{j=1}^{l}sgn(x_j-\hat{\xi}_n)]^-=o_p(1). \tag{3.45}$$

we consider

$$\frac{1}{\sqrt{k}}\sum_{j=1}^{l}sgn(x_j-\hat{\xi}_n)$$

$$\ge\frac{\sqrt{l}}{\sqrt{k}}\frac{1}{\sqrt{l}}\sum_{j=1}^{l}sgn(x_j-\hat{\xi}_n)-\frac{2}{k}\sum_{j=1}^{l}I_j \tag{3.46}$$

$$=\sqrt{\frac{l}{k}}O_p((loglogl)^{\frac{1}{2}})+\frac{1}{\sqrt{k}}l\delta_n+\frac{1}{\sqrt{k}}lO_p(n^{-\frac{1}{2}}),\quad l\le\tau.$$

Here we used that by the law of iterated logarithm

$$liminf_{n\to\infty}\frac{S_n}{\sqrt{n}}\overset{a.s.}{=}O((loglogn)^{\frac{1}{2}}),$$

if $S_n$ is the sum of mean zero independent identically distributed random variables that have finite variance. From (3.46) we got that (3.45) is true. Putting (3.41), (3.44) and (3.45) we have

$$-[\frac{1}{\sqrt{k}}min_{1\le i\le\tau-k}\sum_{j=i+1}^{\tau-k}sgn(x_j-\hat{\xi}_n)]^-$$

$$\overset{D}{=}-[\frac{1}{\sqrt{k}}min_{1\le l\le\tau-k}\sum_{j=1}^{l}sgn(x_j-\hat{\xi}_n)]^- \tag{3.47}$$

35

$$= max\{-[\ min_{1\leq l\leq m(k)}\frac{1}{\sqrt{k}}\sum_{j=1}^{l} sgn(x_j - \hat{\xi}_n)]^-,$$

$$-[min_{m(k)<l\leq k}\frac{1}{\sqrt{k}}\sum_{j=1}^{l} sgn(x_j - \hat{\xi}_n)]^-,$$

$$-[min_{k<l<\tau-k}\frac{1}{\sqrt{k}}\sum_{j=1}^{l} sgn(x_j - \hat{\xi}_n)]^-\}$$

$$= o_p(1). \tag{3.48}$$

On the other hand,

$$\frac{1}{\sqrt{k}}max_{\tau-k<l<\tau}\sum_{\tau-k}^{l} sgn(x_j - \hat{\xi}_n)$$

$$\overset{D}{=} \frac{1}{\sqrt{k}}max_{1<l<k}\sum_{j=1}^{l} sgn(x_{j+\tau-k-1} - \hat{\xi}_n) \tag{3.49}$$

$$\geq \frac{1}{\sqrt{k}}max_{1<l<k}\sum_{j=1}^{l} sgn(x_{j+\tau-k-1} - \eta_0) - \frac{2}{\sqrt{k}}\sum_{j=1}^{l} I_j.$$

The error of the approximation is

$$\frac{2}{\sqrt{k}}\sum_{j=1}^{l} I_j \quad = \frac{l}{\sqrt{k}}O_p(n^{-\frac{1}{2}})$$

$$= \sqrt{k}O_p(n^{-\frac{1}{2}}).$$

Then we get by the strong law of large number

$$\frac{1}{\sqrt{k}}max_{1\leq l\leq k}\sum_{j=1}^{l} sgn(x_{j+\tau-k-1} - \eta_0)$$

$$\geq \sqrt{k}\frac{1}{k}\sum_{j=1}^{k} sgn(x_{j+\tau-k-1} - \eta_0) \tag{3.50}$$

$$\overset{a.s.}{\to} \infty, \qquad k \to \infty.$$

As $Esgn(x_j - \eta_0) > 0$, combining (3.47) and (3.49) we get

$$\lim_{k\to\infty} limsup_{n\to\infty}P\{max_{i\leq\tau-k}V(\frac{i}{n}) \geq max_{\tau-k<i<\tau+k}V(\frac{i}{n})\} = 0$$

As the sign statistic can be looked as a rank statistic with score function

$$\phi(u) = \begin{cases} -1, & 0 \leq u < \frac{1}{2} \\ 1, & \frac{1}{2} < u \leq 1 \\ 0 & otherwise. \end{cases} \tag{3.51}$$

we can use the Hajek (1968) result for the two sample sign statistic to get

$$\frac{V_n(\lambda) - \mu_n}{\sigma_n} \xrightarrow{D} N(0,1), \qquad (3.52)$$

where

$$\mu_n = E\{V_n(\lambda)\}, \quad \sigma_n = Var V_n(\lambda)$$

In the one change-point case, the test statistic

$$M(n) = max_{1 \leq k \leq n} V_n\left(\frac{k}{n}\right)$$
$$= \frac{1}{\sqrt{n}} \sum_{j=1}^{\hat{\tau}(n)} sgn(x_j - \hat{\xi}_n) \qquad (3.53)$$

If $\hat{\tau} < \tau$

$$M(n) - V_n(\lambda) = -\frac{1}{\sqrt{n}} \sum_{j=\hat{\tau}(n)+1}^{\tau(n)} sgn(x_j - \hat{\xi}_n)$$
$$\stackrel{D}{=} -\frac{1}{\sqrt{n}} \sum_{j=1}^{k} sgn(x_{j+\hat{\tau}} - \hat{\xi}_n) \qquad (3.54)$$
$$= o_p(1)$$

and similar statememt is true if $\hat{\tau} > \tau$.

It is easy to see that for $\sigma_n = 2\sqrt{r_n(1-r_n)}$ and $\sigma = 2\sqrt{s(1-s)}$

$$\frac{\sigma_n}{\sigma} \xrightarrow{P} 1 \qquad (3.55)$$

so by Slutsky's theorem, we can replace $\sigma_n$ by $\sigma$ in (3.51), and the proof of the theorem is concluded. $\square$

Besides showing the asymptotic distribution of the test statistic, the above Theorem implies the consistency of our test. Furthermore, the Theorem allows a comparision between statistic used for two-sample problems and those for change-point problems. It shows that asymptotically they have same limit. The

37

important implication of this is, that when we compare different change-point detection procedures, the results of asymptotic relation efficiencies of two-sample test are valid for the change-point tests as well. This statememnt is of course true only for at-most-one-change and for epidemic alternative cases.

# Chapter 4

# SIMULATION

In this Chapter we will consider the powers of our hypothesis tests in the previous chapters.

[1] The initial median is known case.

In Chapter 2, for the test (2.1),

$$H_0 : \quad x_i \text{ has median } \xi_0 \text{ for } i = 1, \cdots, n.$$

$$H_1 : \quad x_i \text{ has median } \xi_0 \text{ for } i = 1, \cdots, \tau_1 - 1 \ \tau_2, \cdots, n. \tag{4.1}$$

$$x_i \text{ has median } \xi_1 \text{ for } i = \tau_1, \cdots, \tau_2 - 1,$$

we consider the following test statistic

$$U_n = max_{k<l} \sum_{i=k}^{l-1} sgn(x_i - \xi_0). \tag{4.2}$$

To estimate the unknown change points $\tau_1$ and $\tau_2$, we use

$$(\hat{\tau_1}(n), \hat{\tau_2}(n)) = argmax_{k<l}(S_{l-1} - S_{k-1}) \tag{4.3}$$

Under $H_0$, Gombay proved that

$$P_{H_0}(U_n \geq N) = 1 - \frac{2}{2N+1} \sum_{j=1}^{2N} (c(j))^n s(j(N+1)) \frac{1 + c(j)}{s(j)} \frac{1 - (-1)^j}{2}. \tag{4.4}$$

where $N$ is a positive integer and

$$c(j) = cos(\frac{j\pi}{2N+1}), \qquad s(j) = sin(\frac{j\pi}{2N+1}).$$

Under the alternative hypothesis $H_1$, we have the distribution of the test statistic $U_n$

$$\sqrt{\lambda_1}|N_1| + \sqrt{1-\lambda_2}|N_2| + 2\sqrt{\delta_n(1-\delta_n)(\lambda_2-\lambda_1)}N_3 + M_n. \qquad (4.5)$$

All the notations here are the same as that in the Chapter 2.

For the hypothesis test (2.1), given significant level $\alpha$, we reject $H_0$ in favour of $H_1$, if $U_n \geq N_\alpha$, where $N_\alpha$ is a positive number and

$$P_{H_0}(U_n \geq N_\alpha) \leq \alpha.$$

The power of the test is

$$P_{H_1}(U_n \geq N_\alpha) \qquad (4.6)$$

$$= P\{\sqrt{\lambda_1}|N_1| + \sqrt{1-\lambda_2}|N_2| + 2\sqrt{\delta_n(1-\delta_n)(\lambda_2-\lambda_1)}N_3 + M_n > N_\alpha\}.$$

For the standard normal random variable $X \sim N(0,1)$, from

$$P\{|X| < x\} = \int_{-x}^{x} \frac{1}{\sqrt{2\pi}}e^{-x^2/2}dx,$$

we get the density function of $|N(0,1)|$

$$f_{|X|}(x) = \sqrt{\frac{2}{\pi}}e^{-x^2/2}.$$

40

Also denote

$$a = \sqrt{\lambda_1},$$

$$b = \sqrt{1 - \lambda_2},$$

$$c = 2\sqrt{\delta_n(1 - \delta_n)(\lambda_2 - \lambda_1)},$$

$$d = N_\alpha - M_n.$$

we have

$$P_{H_1}(U_n \geq N_\alpha)$$

$$= \int\int\int_{ax+by+cz>d,\ x>0,\ y>0} \frac{\sqrt{2}}{\pi^{\frac{3}{2}}} exp\{-\tfrac{1}{2}(x^2 + y^2 + z^2)\} dx dy dz. \tag{4.7}$$

When $n \to \infty$, $M \to \infty$, and $N_\alpha - M_n = d \to -\infty$, we can simply find

$$\lim_{n\to\infty} P_{H_1}(U_n \geq N_\alpha) = 1, \tag{4.8}$$

That means the power of the test converges to one as $n \to \infty$.

Unfortunately, we can not get the explicit expression for the integer (4.7). Numerical calculations are needed to carry out for different $a$, $b$, $c$, $d$.

To get some idea about the power for fixed sample size test, we do some simulation . We assume that the population distributions are normal and uniform. For hypothesis test (2.1), we consider seven different cases for test with epidemic alternatives. From (4.2), (4.3) and (4.4), we calculated the estimated change-points $\hat{\tau}_1$ and $\hat{\tau}_2$, test statistic values $u_n$ and their P-values list in Table 4.1. In the table, rnorm(20) denotes a set of twenty observations from a standard normal population, rnorm(20, 2, 1) denotes a set of twenty observations from a normal population with the mean of 2 and standard deviation of 1, runif(20) denotes a

41

Table 4.1: Initial median is known case

| $H_1$ | $\tau_1$ | $\tau_2$ | $\hat{\tau}_1$ | $\hat{\tau}_2$ | $u_n$ | P-value |
|---|---|---|---|---|---|---|
| rnorm(20), rnorm(20,2,1), rnorm(20) | 21 | 41 | 13 | 48 | 23 | 0.004539 |
| rnorm(20), rnorm(20,1.5,1), rnorm(20) | 21 | 41 | 8 | 40 | 24 | 0.002933 |
| rnorm(20), rnorm(20,1,1), rnorm(20) | 21 | 41 | 6 | 50 | 20 | 0.016071 |
| rnorm(20), rnorm(20,0.5,1), rnorm(20) | 21 | 41 | 8 | 18 | 8 | 0.548926 |
| runif(20), runif(20,0.75,1.75), runif(20) | 21 | 41 | 8 | 59 | 29 | 0.000223 |
| runif(20), runif(20,0.5,1.5), runif(20) | 21 | 41 | 10 | 56 | 32 | 0.000039 |
| runif(20), runif(20,0.25,1.25), runif(20) | 21 | 41 | 3 | 59 | 18 | 0.033904 |

set of twenty observations from a uniform population in $[0, 1]$ and runif(20, 0.5, 1.5) denotes a set of twenty observations from a uniform population in $[0.5, 1.5]$.

We consider only the case where the variance of the distribution population does not change. From the Table (4.1), we can see when the difference between $\xi_0$ and $\xi_1$ is getting large, the statistics value $u_n$ will likely get larger and there will be a more significant P-value. We confirm that the change-point estimators $\hat{\tau}_1$ and $\hat{\tau}_2$ are not good as our theory has predicted. They should be close to 21 and 41 but they are not. But we are able to detect the changes in all but one case, as the P-value is small.

Now we do the simulation on a real world data. We consider the sign test for Lon $\vdots$ $\mathbf{\ldots}$ d's (1987) data which give the radii of circular indentations cut by a milling machine. The sample size is 100. The data are time-ordered and to be read row by row. A constant, 3.9, has been subtracted from all the data. We

assume that they are independent random variables.

| 1.010 | 1.066 | 0.975 | 0.921 | 1.165 | 1.027 | 1.100 | 0.981 | 0.977 | 1.106 |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 0.932 | 0.990 | 0.940 | 0.877 | 0.987 | 0.958 | 1.112 | 0.878 | 1.029 | 0.971 |
| 1.004 | 1.087 | 1.038 | 1.119 | 0.768 | 1.096 | 1.114 | 1.007 | 0.978 | 0.957 |
| 0.884 | 1.004 | 1.032 | 1.130 | 0.961 | 1.066 | 1.029 | 1.107 | 1.150 | 1.190 |
| 1.152 | 1.049 | 1.183 | 0.993 | 1.161 | 0.988 | 1.087 | 1.034 | 0.889 | 1.109 |
| 1.196 | 1.098 | 0.954 | 0.986 | 0.943 | 1.058 | 0.960 | 1.073 | 0.904 | 1.171 |
| 1.060 | 1.189 | 1.019 | 1.213 | 1.204 | 1.148 | 1.033 | 1.023 | 1.145 | 0.994 |
| 1.147 | 1.054 | 1.059 | 0.972 | 1.141 | 1.082 | 0.931 | 0.848 | 1.039 | 1.043 |
| 1.016 | 1.027 | 0.932 | 0.879 | 0.754 | 0.911 | 0.971 | 1.180 | 0.849 | 0.870 |
| 1.003 | 0.843 | 1.018 | 1.145 | 0.995 | 0.895 | 1.085 | 1.055 | 0.992 | 1.141 |

To do the calculation, first we get an estimator of the initial median $\xi_0 = 0.987$ based on first 15 observations. Then we do calculation just like the known initial median case, we get the test statistic value is 34 and the two estimated change-points are $\hat{\tau}_1 = 16$ and $\hat{\tau}_2 = 82$. The P-value for the test is 0.00105002G and clearly indicates that there are changes along the sequence. The corresponding test of Pettitt (1979) of $H_0$ against one change alternative got the P-value of 0.1324 and did not detect signal change in this data.

For test (2.1), under the null hypothesis $H_0$, we use the relation (2.3) to calculate the exactly critical value $N_\alpha(n)$ for given n and $\alpha$ listed in Table 4.2.

[2] The initial median is unknown case.

Let's consider the hypothesis test (3.1), we use

$$M(n) = max_{1 \leq k < l \leq n} n^{-\frac{1}{2}} | \sum_{k \leq i < l} sgn(x_i - \hat{\xi}_n)|. \tag{4.9}$$

The distribution of the test statistic $M(n)$ under $H_0$ is

$$P\{M(n) \leq z\} \tag{4.10}$$

43

Table 4.2: $N_\alpha(n)$ for test (2.1) that $P_{H_0}\{U_n \geq N_\alpha(n)\} \leq \alpha$

| $n \backslash \alpha$ | 0.1 | 0.05 | 0.025 | 0.01 | 0.005 | 0.0025 | 0.001 |
|---|---|---|---|---|---|---|---|
| 4 | 4 | | | | | | |
| 5 | 4 | 5 | | | | | |
| 6 | | 5 | 6 | | | | |
| 7 | 5 | | 6 | 7 | | | |
| 8 | | 6 | 7 | | 8 | | |
| 9 | 6 | 7 | | 8 | | 9 | |
| 10 | 6 | 7 | 8 | | 9 | | 10 |
| 11 | 7 | | 8 | 9 | | 10 | 11 |
| 12 | 7 | 8 | | 9 | 10 | | 11 |
| 13 | 7 | 8 | 9 | 10 | | 11 | 12 |
| 14 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| 15 | 8 | 9 | 10 | 11 | | 12 | 13 |
| 16 | 8 | 9 | 10 | 11 | 12 | | 13 |
| 17 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| 18 | 8 | 10 | 11 | 12 | | 13 | 14 |
| 19 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 20 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 21 | 9 | 10 | 11 | 13 | | 14 | 16 |
| 22 | 9 | 11 | 12 | 13 | 14 | 15 | 16 |
| 23 | 9 | 11 | 12 | 13 | 14 | 15 | 16 |
| 24 | 10 | 11 | 12 | 14 | | 15 | 17 |
| 25 | 10 | 11 | 12 | 14 | 15 | 16 | 17 |
| 26 | 10 | 11 | 12 | 14 | 15 | 16 | 17 |
| 27 | 10 | 12 | 13 | 14 | | 16 | 18 |
| 28 | 10 | 12 | 13 | 15 | 16 | 17 | 18 |

$$= 1 + \frac{2}{\binom{2m}{m}}\{[\alpha \sum_{s=1}^{[\frac{m}{\alpha+1}]} \binom{2m}{m-s(\alpha+1)} - (\alpha - 1)\sum_{s=1}^{[\frac{m}{\alpha}]}\binom{2m}{m-s\alpha})]$$

$$-[\sum_{i=1}^{\alpha} \sum_{s=1}^{[\frac{m+i}{\alpha+1}]} \binom{2m}{m+i-s(\alpha+1)}) - \sum_{s=1}^{[\frac{m+i}{\alpha+1}]} \binom{2m}{m+i-s(\alpha+1)}) - \sum_{i=1}^{\alpha-1} \sum_{s=1}^{[\frac{n+i}{\alpha}]} \binom{2m}{m+i-s\alpha})]\}$$

Here all the notations are same as that in the Chapter 3.

Table 4.3: Initial median is unknown case

| $H_1$ | $\tau_1$ | $\tau_2$ | $\hat{\tau}_1$ | $\hat{\tau}_2$ | $n^{\frac{1}{2}}m(n)$ | P-value |
|---|---|---|---|---|---|---|
| rnorm(20), rnorm(20,2,1), rnorm(20) | 21 | 41 | 20 | 48 | 22 | 0.0000008 |
| rnorm(20), rnorm(20,1.5,1), rnorm(20) | 21 | 41 | 12 | 40 | 20 | 0.0000203 |
| rnorm(20), rnorm(20,1,1), rnorm(20) | 21 | 41 | 7 | 50 | 15 | 0.0088002 |
| rnorm(20), rnorm(20,0.5,1), rnorm(20) | 21 | 41 | 8 | 40 | 12 | 0.1025916 |
| runif(20), runif(20,0.75,1.75), runif(20) | 21 | 41 | 19 | 40 | 19 | 0.0000846 |
| runif(20), runif(20,0.50,1.50), runif(20) | 21 | 41 | 23 | 46 | 17 | 0.0010595 |
| runif(20), runif(20,0.25,1.25), runif(20) | 21 | 41 | 20 | 33 | 7 | 0.8687585 |

*Table* 4.2 : (*continued*)

| n \ $\alpha$ | 0.1 | 0.05 | 0.025 | 0.01 | 0.005 | 0.0025 | 0.001 |
|---|---|---|---|---|---|---|---|
| 29 | 11 | 12 | 13 | 15 | 16 | 17 | 18 |
| 30 | 11 | 12 | 14 | 15 | 16 | 17 | 19 |
| 35 | 12 | 13 | 15 | 16 | 18 | 19 | 20 |
| 40 | 12 | 14 | 16 | 18 | 19 | 20 | 22 |
| 45 | 13 | 15 | 17 | 19 | 20 | 21 | 23 |
| 50 | 14 | 16 | 18 | 20 | 21 | 23 | 24 |
| 55 | 15 | 17 | 18 | 21 | 22 | 24 | 26 |
| 60 | 15 | 17 | 19 | 22 | 23 | 26 | 27 |
| 70 | 16 | 19 | 21 | 23 | 25 | 27 | 29 |
| 80 | 18 | 20 | 22 | 25 | 27 | 29 | 31 |
| 90 | 19 | 21 | 24 | 27 | 29 | 30 | 33 |
| 100 | 20 | 22 | 25 | 28 | 30 | 32 | 35 |
| 110 | 21 | 24 | 26 | 29 | 32 | 34 | 36 |
| 120 | 21 | 25 | 27 | 31 | 33 | 35 | 38 |
| 130 | 22 | 26 | 28 | 32 | 34 | 37 | 40 |
| 140 | 23 | 27 | 30 | 33 | 36 | 38 | 41 |
| 150 | 24 | 27 | 31 | 34 | 37 | 39 | 42 |
| 200 | 28 | 32 | 35 | 40 | 43 | 46 | 49 |
| 250 | 31 | 35 | 39 | 44 | 48 | 51 | 55 |
| 300 | 34 | 39 | 43 | 49 | 52 | 56 | 60 |
| 350 | 37 | 42 | 47 | 52 | 57 | 60 | 65 |
| 400 | 39 | 45 | 50 | 56 | 60 | 64 | 70 |
| 500 | 44 | 50 | 56 | 63 | 68 | 72 | 78 |
| 1000 | 62 | 71 | 79 | 89 | 96 | 102 | 110 |
| 2000 | 88 | 100 | 112 | 126 | 135 | 144 | 156 |

Given significance level $\alpha$, we reject $H_0$ in favour of $H_1$ for large value of

$M(n)$, such that

$$P_{H_0}\{M(n) \geq M_\alpha(n)\} \leq \alpha. \tag{4.11}$$

where $M_\alpha(n)$ is a positive number.

From (3.34), we have the distribution of $M(n)$ under the alternative hypothesis $H_1$. Then, at the significance level $\alpha$, the power of test (3.1) is

$$P_{H_1}\{M(n) \geq M_\alpha(n)\}$$

$$= P_{H_1}\{\tfrac{M(n)-|\mu|}{\sigma} \geq \tfrac{M_\alpha(n)-|\mu|}{\sigma}\} \tag{4.12}$$

$$\approx 1 - \Phi(\tfrac{M_\alpha-|\mu|}{\sigma})$$

where $\Phi$ is the cumulative distribution function of the standard normal random variable.

For the simulation, first we consider the same data used in the simulations summ    ' in the Table 4.1. Assuming the intial medians of population under the n   hy  .thesis are unknown, we calculate the medians of the observations. Using relation (3.4), we get the estimated change-points $\hat{\tau}_1$ and $\hat{\tau}_2$, test statistic values $n^{\frac{1}{2}}m(n)$ and their P-values listed in Table 4.3.

From the Table 4.1 and Table 4.3, we find that estimated change-points calculated from statistic $M(n)$ are closer to the real change-points than that from statistic $U_n$.

For the Lombard's (1987) data, we use the meadian of the total data $\hat{\xi}_n = 1.027$ as the intial median of the population distribution is assumed to be unknown. Then we calculated the estimated change-points $\hat{\tau}_1 = 32$, $\hat{\tau}_2 = 76$. The

46

Table 4.4: $n^{\frac{1}{2}}M_\alpha(n)$  $P_{H_0}\{M(n) \geq M_\alpha(n)\} \leq \alpha$

| n \ α | 0.1 | 0.05 | 0.025 | 0.01 | 0.005 | 0.0025 | 0.001 |
|---|---|---|---|---|---|---|---|
| 14 | | 6 | | | | | 7 |
| 15 | | 6 | | | | | 7 |
| 16 | 6 | | 7 | | | | 8 |
| 17 | 6 | | 7 | | | | 8 |
| 18 | | 7 | | 8 | | | 9 |
| 19 | | 7 | | 8 | | | 9 |
| 20 | 7 | | 8 | | | 9 | 10 |
| 21 | 7 | | 8 | | | 9 | 10 |
| 22 | | 8 | | 9 | | | 10 |
| 23 | | 8 | | 9 | | | 10 |
| 24 | | 8 | 9 | | | 10 | 11 |
| 25 | | 8 | 9 | | | 10 | 11 |
| 26 | 8 | | 9 | | 10 | | 11 |
| 27 | 8 | | 9 | | 10 | | 11 |
| 28 | | 9 | | 10 | | 11 | 12 |
| 29 | | 9 | | 10 | | 11 | 12 |
| 30 | 9 | | 10 | | 11 | | 12 |
| 31 | 9 | | 10 | | 11 | | 12 |
| 32 | 9 | | 10 | 11 | | 12 | 13 |
| 33 | 9 | | 10 | 11 | | 12 | 13 |
| 34 | 9 | 10 | | 11 | | 12 | 13 |
| 35 | 9 | 10 | | 11 | | 12 | 13 |
| 36 | | 10 | 11 | | 12 | | 13 |
| 37 | | 10 | 11 | | 12 | | 13 |

test statistic value $n^{\frac{1}{2}}m(n))$ is 18 and its P-value is 0.0520979, indicating the presence of changes. These are in good agreement with the Gombay's (1994) results of rank test, with the suggestion of the cusum plot and with Lombard's (1987) conclusins.

In the general case for test (3.1), under the null hypothesis $H_0$, we use relation (3.6) to calculate the exact critical value $n^{\frac{1}{2}}M_\alpha(n)$ for given n and $\alpha$ listed in Table 4.4.

*Table* 4.4 : (*continued*)

| n \ $\alpha$ | 0.1 | 0.05 | 0.025 | 0.01 | 0.005 | 0.0025 | 0.001 |
|---|---|---|---|---|---|---|---|
| 38 | 10 | | 11 | 12 | | 13 | 14 |
| 39 | 10 | | 11 | 12 | | 13 | 14 |
| 40 | 10 | 11 | | 12 | 13 | | 14 |
| 45 | 11 | | 12 | 13 | | 14 | 15 |
| 50 | 11 | 12 | 13 | 14 | | 15 | 16 |
| 55 | 12 | 13 | 14 | | 15 | 16 | 17 |
| 60 | | 13 | 14 | 15 | 16 | 17 | 18 |
| 65 | 13 | 14 | 15 | 16 | | 17 | 18 |
| 70 | 14 | 15 | 16 | | 17 | 18 | 19 |
| 80 | 14 | 16 | 17 | 18 | | 19 | 21 |
| 90 | 15 | 17 | 18 | 19 | 20 | 21 | 22 |
| 100 | 16 | 17 | 19 | 20 | 21 | 22 | 23 |
| 110 | 17 | 18 | 19 | 21 | 22 | 23 | 24 |
| 120 | 18 | 19 | 20 | 22 | 23 | 24 | 25 |
| 130 | 18 | 20 | 21 | 23 | 24 | 25 | 26 |
| 140 | 19 | 21 | 22 | 24 | 25 | 26 | 27 |
| 150 | 20 | 21 | 23 | 24 | 26 | 27 | 28 |
| 160 | 20 | 22 | 24 | 25 | 26 | 28 | 29 |
| 170 | 21 | 23 | 24 | 26 | 27 | 28 | 30 |
| 180 | 22 | 23 | 25 | 27 | 28 | 29 | 31 |
| 190 | 23 | 24 | 26 | 28 | 29 | 30 | 32 |
| 200 | 23 | 25 | 26 | 28 | 30 | 31 | 32 |

# APPENDIX

[1] Program to calculate the test statistics value $u_n$ of (1.2) and the estimated change-points $\hat{\tau}_1$ and $\hat{\tau}_2$ when the mean $\xi_0$ is known.

```
v < -   Obeservations of H1
x < -   sign(v - ξ0)
z < -   x * 0
w < -   x * 0
for(i in 1 : length(x)) {
        y < -   x * 0
        for(j in i : length(x)) {
                s < -   0
                for(k in i : j) {
                        s < -   s + x[k] }
                y[j] < -   s
        }
        a < -   y[i]
        t < -   i
        for(l in i : length(x)) {
                b < -   y[l]
                if(a < b) {
                        a < -   b
                        t < -   l }
        }
        z[i] < -   a
        w[i] < -   t
}
print(z)
print(w)
a < -   z[1]
t < -   1
for(m in 2 : length(x))  {b < -   z[m]
        if(a < b)  {t < -   m
                a < -   b
        }
}
```

```
print("The Statistics value u_n is")
print(a)
print("The estimated τ̂_1 is")
print(t − 1)
print("The estimated τ̂_2 is")
print(w[t])
```

[2] Program to calculate the test statistics $n^{\frac{1}{2}}m(n)$ and the estimated change-points $\hat{\tau}_1$ and $\hat{\tau}_2$ when the mean is unknown.

```
v < −  Observations in H_1
e < −  median(v)
x < −  sign(v − e)
z < −  x * 0
w < −  x * 0
for(i in 1 : length(x))  {
        y < −  x * 0
        for(j in i : length(x))  {
                s < −  0
                for(k in i : j)  {
                        s < −  s + x[k] }
                y[j] < −  s
        }
        a < −  abs(y[i])
        t < −  i
        for(l in i : length(x))  {
                b < −  abs(y[l])
                if(a < b)  {
                        a < −  b
                        t < −  l
                }
        }
        z[i] < −  a
        w[i] < −  t
}
print(z)
print(w)
a < −  z[1]
t < −  1
```

```
for(l in 2 : length(x)) {
        b <-  z[l]
        if(a < b) {
                a <-  b
                t <-  l
        }
}
print("The Statistics value of n½m(n) is")
print(a)
print("The    l...ated τ̂₁ is")
print(t - 1
print("The (    ...ated τ̂₂ is")
print(w[t])
```

[3] Program to calculate the P-Value for the test statistics $n^{\frac{1}{2}}m(n)$ when the

mean is unknown.

```
n <-  Sample size
a <-  n½m(n)
m <-  floor(n/2)
d <-  2 * m
sum1 <-  0
sum2 <-  0
sum3 <-  0
sum4 <-  0
```

% calculate $\sum_{i=1}^{\alpha} \sum_{s=1}^{\frac{m+i}{a+1}} \binom{2m}{m+i-s(\alpha+1)}$ %

```
for(i in 1 : a) {b <-  floor((m + i)/(a + 1))
        sum11 <-  0
        for(s in 1 : b)  {c <-  m + i - s * (a + 1)
                if(c > 0 && c <= d)  {
                        e <-  1
                        f <-  c - 1
                        for(k in 0 : f)  {e <-  e * (d - k)/(c - k) }
                }
                else  if(c == 0)  {e <-  1}
                else  {e <-  0}
                sum11 <-  sum11 + e
        }
        sum1 <-  sum1 + sum11
}
```

51

% calculate $\sum_{i=1}^{a-1} \sum_{s=1}^{\frac{m+i}{a}} \binom{2m}{m+i-sa}$ %

$a1 <- a - 1$

$for(i\ in\ 1:a1)\ \{b <- floor((m+i)/a)$

      $sum22 <- 0$

      $for(s\ in\ 1:b)\ \{c <- m+i-s*a$

            $if(c > 0\ \&\&\ c <= d)\ \{e <- 1$

                $f <- c-1$

                $for(k\ in\ 0:f)\ \{e <- e*(d-k)/(c-k)\}$

            $\}$

            $else\ if(c == 0)\ \{e <- 1\ \}$

            $else\ \{e <- 0\ \}$

            $sum22 <- sum22 + e$

      $\}$

      $sum2 <- sum2 + sum22$

$\}$

% calculate $(\alpha - 1)\sum_{s=1}^{\frac{m}{a}} \binom{2m}{m-sa}$ %

$g <- floor(.n/a)$

$for(i\ in\ 1:g)\ \{c <- m - i*a$

      $if(c > 0\ \&\&\ c <= d)\ \{e <- 1$

          $f <- c-1$

          $for(k\ in\ 0:f)\ \{e <- e*(d-k)/(c-k)\}$

      $\}$

      $else\ if(c == 0)\ \{e <- 1\ \}$

      $else\ \{e <- 0\ \}$

      $sum3 <- sum3 + e$

$\}$

$sum3 <- (a-1)*sum3$

% calculate $\alpha\sum_{s=1}^{\frac{m}{a+1}} \binom{2m}{m-s(\alpha+1)}$ %

$h <- floor(m/(a+1))$

$for(i\ in\ 1:h)\ \{c <- m - i*(a+1)$

      $if(c > 0\ \&\&\ c <= d)\ \{e <- 1$

          $f <- c-1$

          $for(k\ in\ 0:f)\ \{e <- e*(d-k)/(c-k)\ \}$

      $\}$

      $else\ if(c == 0)\ \{e <- 1\ \}$

      $else\ \{e <- 0\ \}$

      $sum4 <- sum4 + e$

$\}$

$sum4 <- a*sum4$

```
% calculate (2m m) %
m1 < − m − 1
e < − 1
for(i in 0 : m1) { e < − e ∗ (d − i)/(m − i) }
p < − 2 ∗ (sum1 − sum2 + sum3 − sum4)/e
print("The P − Value for the test statistics value with unknow mean is")
print(p)
```

[4] Program to calculate the critical value for statistics $U_n$ of (2.2) when the

initial mean $\xi_0$ is known.

```
n < − Sample size
for(n1 in : n) {
        print(" ∗ ∗ ∗ ∗n = ∗ ∗ ∗")
        print(n1)
        for(k in 1 : n1) {s < − 0
                k1 < − 2 ∗ k
                for(j in 1 : k1) {
                        a < − cos(j ∗ pi/(2 ∗ k + 1))
                        b < − sin(j ∗ pi/(2 ∗ k + 1))
                        c < − sin(j ∗ pi ∗ (k + 1)/(2 ∗ k + 1))
                        s < − s + (a^{n1}) ∗ c ∗ (1 + a) ∗ (1 − (−1)^j)/(2 ∗ b)
                }
                p < − 1 − (2 ∗ s)/(2 ∗ k + 1)
                if(p > 0.05 && p <= 0.1) { print("M1 = ")
                        print(k) }
                if(p > 0.025 && p <= 0.05) { print("M2 = ")
                        print(k) }
                if(p > 0.01 && p <= 0.025) { print("M3 = ")
                        print(k) }
                if(p > 0.005 && p <= 0.01) { print("M4 = ")
                        print(k) }
                if(p > 0.0025 && p <= 0.005) { print("M5 = ")
                        print(k) }
                if(p > 0.001 && p <= 0.0025) { print("M6 = ")
                        print(k) }
                if(p <= 0.001) { print("M7 = ")
                        print(k) }
        }
}
```

[5] Program to calculate the P-Value for the test statistics value $u_n$ of (2.2)

when the initial mean $\xi_0$ is known.

```
n <- Sample size
N <- un
N1 <- 2 * N
s <- 0
for(i in 1 : N1) {
        a <- cos(i * pi/(2 * N + 1))
        b <- sin(i * pi/(2 * N + 1))
        c <- sin(i * pi * (N + 1)/(2 * N + 1))
        s <- s + (a^n) * c * (1 + a) * (1 - (-1)^i)/(2 * b)
}
p <- 1 - 2 * s/(2 * N + 1)
print("The P - Value for the test of known mean is")
print(p)
```

[6] Program to calculate the critical values for statistics $M(n)$ of (3.4) when

the initial mean is unknown

```
n <- sample size
m <- floor(n/2)
d <- 2 * m
for(a in 1 : n) {
        sum1 <- 0
        sum2 <- 0
        sum3 <- 0
        sum4 <- 0
        for(i in 1 : a) {
                b <- floor((m + i)/(a + 1))
                sum11 <- 0
                for(s in 1 : b) {
                        c <- m + i - s * (a + 1)
                        if(c > 0 && c <= d) {
                                e <- 1
                                f <- c - 1
                                for(k in 0 : f) {
                                        e <- e * (d - k)/(c - k) }
                        }
                }
```

54

```
            else  if(c == 0)  { e < −  1 }
            else  { c < −  0 }
            sum11 < −  sum11 + e
      }
      sum1 < −  sum1 + sum11
}
a1 < −  a − 1
for(i in 1 : a1)  {
      b < −  floor((m + i)/a)
      sum22 < −  0
      for(s in 1 : b)  {
            c < −  m + i − s * a
            if(c > 0 && c <= d) {
                  e < −  1
                  f < −  c − 1
                  for(k in 0 : f)  { e < −  e * (d − k)/(c − k) ¡
            }
            else  if(c == 0)  { e < −  1 }
            else  { e < −  0 }
            sum22 < −sum22 + e
      }
      sum2 < −sum2 + sum22
}
g < −  floor(m/a)
for(i in 1 : g)  {
      c < −  m − i * a
      if(c > 0 && c <= d)  {
            e < −  1
            f < −  c − 1
            for(k in 0 : f)  { e < −  e * (d − k)/(c − k) }
      }
      else  if(c == 0)  { e < −  1 }
      else  { e < −  0 }
      sum3 < −sum3 + e
}
sum3 < −  (m − 1) * sum3
```

```
h <- floor(m/(a + 1))
for(i in 1 : h) {
    c <- m - i * (a + 1)
    if(c > 0 && c <= d) {
        e <- 1
        f <- c - 1
        for(k in 0 : f) { e <- e * (d - k)/(c - k) }
    }
    else if(c == 0) { e <- 1 }
    else { e <- 0 }
    sum4 <- sum4 + e
}
sum4 <- a * sum4
m1 <- m - 1
e <- 1
for(i in 0 : m1) {
    e <- e * (d - i)/(m - i)}
p <- 2 * (sum1 - sum2 + sum3 - sum4)/e
if(p > 0.05 && p <= 0.1) { print("M1 = ")
    print(a) }
if(p > 0.025 && p <= 0.05) { print("M2 = ")
    print(a) }
if(p > 0.01 && p <= 0.025) { print("M3 = ")
    print(a) }
if(p > 0.005 && p <= 0.01) { print("M3 = ")
    print(a) }
if(p > 0.0025 && p <= 0.005) { print("M5 = ")
    print(a) }
if(p > 0.001 && p <= 0.0025) { print("M6 = ")
    print(a) }
if(p <= 0.001) { print("M7 = ")
    print(a) }

}
```

# Bibliography

[1] Bhattacharyya, G. K. (1984) *Tests of randomness against trend or serial correlations*, in: Handbook of Statistics, Vol.4 Ed. P.R.Krishnaiah and P. K. Sen, 89-111, Elsevier Science Publishers.

[2] Bhattacharyya, G. K. and Johnson, R. A. (1968) *Nonparametric tests for shift at unknown time point*, The Annals of Mathematical Statistics. **39**, 1731-1743.

[3] Billingsley, P. (1968) *Convergence of Probability Measure* (Wiley, New York).

[4] Chernoff, H. and Savage, R. I. (1958) *Asymptotic normality and efficiency of certain no-parametric test statistics*, The Annals of Mathematical Statistics, **29**, 972-994.

[5] Chow, Y. S. and Hsiung, A. C. (1976) *Limiting behavior of $max_{j \leq n} S_j j^{-\alpha}$ and the first passage times in a random walk with positive drift*, Bulletin of the Institute of the Mathematics Academic Sinica, **4**, 35-44.

[6] Csörgő, M. and Horváth, L. (1986) *Nonparametric methods for changepoint problems*, Technical Report Series of the Laboratory for Research in Statistics and Probability. Carleton University - University of Ottawa, Canada. No.**86**, 20-55.

[7] Csörgő, M. and Révész (1981) *Strong approximations in probability and statistics*, (Academic Press, New York).

[8] Gnedenko, B. V. (1954) *Kriterien für die Unveränderlichkeit der Wahrscheinlichkeitsverteilung von zwei unabhängigen Stichprobenreihen*, Mathematiche Nachrichten, **12** 29-66 (in Russian with German summary.)

[9] Gombay, E. (1994) *Testing for change-points with rank and sign statistics*, Statistics and Probability Letters, **20**, 49-55.

[10] Gombay, E. (1994) *A limit theorem for rank tests of the change-point problem*. preprint.

[11] Hájek, J. (1968) *Asymptotic normality of simple linear rank statistics under alternatives* , The Annals of Mathematical Statistics, **39**, 325-346.

[12] Hájek 1 and Šidák, Z. (1967) *Theory of rank tests* , (Academic Press, New York).

[13] Hawkins, D. M. (1977) *Testing a sequence of observations for a shift in location*, Journal of the American Statistical Association, **72**, 180-186.

[14] Lombard, F. (1987) *Rank test for changepoint problems*, Biometrika **74**, 615-624.

[15] Pag [?] S. (1954) *Continuous inspection schemes*, Biometrika **41**, 100-115.

[16] Page, E. S. (1955) *A test for a change in a parameter occurring at an unknown point*. Biometrika **42**, 523-526.

[17] Pettitt, A. N. (1979) *A non-parametric approach to the change-point problem* Applied St . **28**, 126-135.

[18] Pettitt, A. N. (1980) *Estimating a changepoint using nonparametric type statistics*, Journal of Statistical Computation and Simulation. **11**, 261-274.

[19] Sen, P. K. (1984) *Nonparametric procedures for some miscellaneous problems*, in: Handbook of Statistics, Vol.4 Ed. P.R.Krishnaiah and P. K. Sen, 669-739. Elsevier Science Publishers.

[20] Sen, A. and Srivastava, M. S. (1975) *On tests for detecting change in mean*, The Annals of Statistics, **3**, 90-108.

[21] Siegmund, D. O. (1986) *Bo ndary crossing probabilities with statistical applications*, The Annals of Statistics, **14**, 361-404.

[22] Wolfe, D. A. and Schechtman, E. (1984) *Nonparametric statistical procedures for the changepoint problem* Journal of Statistical Planning and Inference, **9**, 389-396.

[23] Zacks, S. (1983) *Survey of classical and bayesian approaches to the change-point problem: fixed sample and sequential procedures of testing and estimation*, in: Recent Advances in Statistics, eds. M. H. Rizvi, J. S. Rustagi and D. O. Siegmund, (Academic Press, New York) 245-267.