



[Journal Home](#) > [Volume 11, Issue - 1, March 30, 2016](#)

An Analytical Method for Detecting the Change-Point in Simple Linear Regression Model. Application at Weibull Distribution

 [Pdf Version](#)

[Dariush GHORBANZADEH](#)

[Philippe DURAND](#)

[Luan JAUPÍ](#)

Keywords

[Change-point](#), [simple linear regression mode](#), [Weibull distribution](#)

Abstract

In this paper we study an analytical method to detect the change-point in the model of simple linear regression. The study method is used to estimate the parameters of a Weibull model representative a change-point. The procedure proposed in this paper is illustrated through a classical change-point data. For the accuracy of the method a simulation study is performed.

[\(top\)](#)

[\(back to issue\)](#)

An analytical method for detecting the change-point in simple linear regression model. Application at Weibull distribution.

Dariush GHORBANZADEH⁽¹⁾ Philippe DURAND⁽²⁾ Luan JAUPI⁽³⁾
dariush.ghorbanzadeh@cnam.fr⁽¹⁾ *philippe.durand@cnam.fr*⁽²⁾ *jaupi@cnam.fr*⁽³⁾

CNAM- département IMATH, 292, Rue Saint-Martin, 75141-Paris CEDEX 03, France

Abstract

In this paper we study an analytical method to detect the change-point in the model of simple linear regression. The study method is used to estimate the parameters of a Weibull model representative of a change-point. The procedure proposed in this paper is illustrated through a classical change-point data. For the accuracy of the method a simulation study is performed.

Keywords: Change-point; simple linear regression model; Weibull distribution.

1 Introduction

Change-point models have originally been developed in connection with applications in quality control, where a change from the *in-control* to the *out-of-control* state has to be detected based on the available random observations. Up to now various change-point models have been suggested for a broad spectrum of applications like quality control, reliability, econometrics, medicine, signal processing, meteorology, etc.

The general change-point problem can be described as follows: A random process indexed by time is observed and we want to investigate whether a change in the distribution of the random elements occurs.

Formally, let X_1, \dots, X_n denote a sequence of independent random variables, where

the elements X_1, \dots, X_k have an identical distribution function f_1 and X_{k+1}, \dots, X_n are distributed according to f_2 and the change-point k is unknown.

The change point problem has been considered and studied by several authors. Change-point analysis concerns with the detection and estimation of the point at which the distribution changes. One change point problem or multiple change points problem have been studied in the literature, depending on whether one or more change points are observed in a sequence of random variables. Several methods, parametric or non-parametric, have been developed to approach the solution of this problem while the range of applications of change point analysis is broad.

There is an extensive bibliography on the subject and several methods to search for the change-point problem have appeared in the literature. *The CUSUM (cumulative sum) approach* : Basseville & Nikiforov [1], Lucas & Crosier [12], Ritov [15] and Yashchin [16]. *The maximum-logarithm of the likelihood ratio approach* : Guralnik & Srivastava [10], Gustafsson [11] and Ghorbanzadeh [7]. *The Bayesian approach* : Bradley & all [5], Barry & Hartigan [2] and Ghorbanzadeh & Lounes [9]. *The Non-Parametric approach* : Pettitt [13], Dehling & all [6] and Ghorbanzadeh & Picard [8].

In this work we consider the change-point model for a simple linear regression with one change point. Consider n pairs of observations (X_i, Y_i) and we suppose that the relationship between X and Y can be described by a simple linear regression, where the structure changes after a change point $k \in \{4, \dots, n-4\}$. This restriction on k is needed to ensure that the parameters in the model are estimable. Thus, the observations (X_i, Y_i) follow a linear model for $i \leq k$ and another linear model for $i > k$. Therefore, the model is given by

$$\begin{cases} Y_i = B_1 + A_1 h(X_i) + \varepsilon_{1,i} & \text{if } i = 1, \dots, k \\ Y_i = B_2 + A_2 h(X_i) + \varepsilon_{2,i} & \text{if } i = k+1, \dots, n \end{cases} \quad (1)$$

where A_j and B_j ($j = 1, 2$), are the unknown parameters, $\varepsilon_{j,i}$ are independent errors and h is a known function.

The methode proposed in this paper is illustrated through a classical change-point data from Quandt [14]. We use the model (1) to estimate the parameters of a Weibull model representative a change-point. For the accuracy of the method a simulation study is performed.

2 Analytical method for the change-point estimate

For $k_0 \in \{4, \dots, n-4\}$, we construct $n-7$ two subsamples as follows :

k_0	Sample ₁ (k_0)	Sample ₂ (k_0)
4	X_1, \dots, X_4	X_5, \dots, X_n
5	X_1, \dots, X_5	X_6, \dots, X_n
\vdots	\vdots	\vdots
$n-4$	X_1, \dots, X_{n-4}	X_{n-3}, \dots, X_n

Table 1: Distribution of data into two subsamples.

For each $k_0 \in \{4, \dots, n-4\}$, we consider the following models

$$\begin{cases} Y_i = B_1(k_0) + A_1(k_0) h(X_i) + \varepsilon_{1,i}(k_0) & \text{if } i = 1, \dots, k_0 \\ Y_i = B_2(k_0) + A_2(k_0) h(X_i) + \varepsilon_{2,i}(k_0) & \text{if } i = k_0 + 1, \dots, n \end{cases} \quad (2)$$

For each $k_0 \in \{4, \dots, n-4\}$, $A_1(k_0)$, $B_1(k_0)$, $A_2(k_0)$ and $B_2(k_0)$ solve the following minimization problem:

$$\min_{(A_1(k_0), B_1(k_0), A_2(k_0), B_2(k_0))} D(k_0, A_1(k_0), B_1(k_0), A_2(k_0), B_2(k_0))$$

where

$$D(k_0, A_1(k_0), B_1(k_0), A_2(k_0), B_2(k_0)) = \sum_{i=1}^{k_0} \varepsilon_{1,i}^2(k_0) + \sum_{i=k_0+1}^n \varepsilon_{2,i}^2(k_0) \quad (3)$$

By classics calculations, we obtain the estimator of $A_1(k_0)$, $B_1(k_0)$, $A_2(k_0)$ and $B_2(k_0)$

$$\begin{cases} \hat{A}_1(k_0) = \frac{\sum_{i=1}^{k_0} (h(X_i) - \bar{h}_{k_0}(X)) (Y_i - \bar{Y}_{k_0})}{\sum_{i=1}^{k_0} (h(X_i) - \bar{h}_{k_0}(X))^2} & , \hat{B}_1(k_0) = \bar{Y}_{k_0} - \hat{A}_1(k_0) \bar{h}_{k_0}(X) \\ \hat{A}_2(k_0) = \frac{\sum_{i=k_0+1}^n (h(X_i) - \bar{h}_{k_0}^*(X)) (Y_i - \bar{Y}_{k_0}^*)}{\sum_{i=k_0+1}^n (h(X_i) - \bar{h}_{k_0}^*(X))^2} & , \hat{B}_2(k_0) = \bar{Y}_{k_0}^* - \hat{A}_2(k_0) \bar{h}_{k_0}^*(X) \end{cases} \quad (4)$$

where

$$\begin{cases} \bar{h}_{k_0}(X) = \frac{1}{k_0} \sum_{i=1}^{k_0} h(X_i) & , \bar{h}_{k_0}^*(X) = \frac{1}{n-k_0} \sum_{i=k_0+1}^n h(X_i) \\ \bar{Y}_{k_0} = \frac{1}{k_0} \sum_{i=1}^{k_0} Y_i & , \bar{Y}_{k_0}^* = \frac{1}{n-k_0} \sum_{i=k_0+1}^n Y_i \end{cases} \quad (5)$$

Let $D(k_0) = D(k_0, \hat{A}_1(k_0), \hat{B}_1(k_0), \hat{A}_2(k_0), \hat{B}_2(k_0))$ and

$$k^* = \underset{k_0}{\operatorname{argmin}} D(k_0) \quad (6)$$

By equation (4), we deduce the estimators of A_1 , B_1 , A_2 and B_2

$$\hat{A}_1 = \hat{A}_1(k^*) , \hat{B}_1 = \hat{B}_1(k^*) , \hat{A}_2 = \hat{A}_2(k^*) , \hat{B}_2 = \hat{B}_2(k^*) \quad (7)$$

and the change-point time is estimated by

$$\hat{k} = \text{length of Sample}_1(k^*) \quad (8)$$

3 Application to Quandt's data

This data was illustrated by Quandt [14]. He considered a simple linear regression model with one point-change. The data, listed in Table 2.

i	1	2	3	4	5	6	7	8	9	10
X_i	4	13	5	2	6	8	1	12	17	20
Y_i	3.473	11.555	5.714	5.710	6.046	7.650	3.140	10.312	13.353	17.197
i	11	12	13	14	15	16	17	18	19	20
X_i	15	11	3	14	16	10	7	19	18	9
Y_i	13.036	8.264	7.612	11.802	12.551	10.296	10.014	15.472	15.650	9.871

Table 2: Quandt's data.

The results obtained by the model (1) show a change after the first $\hat{k} = 12$ observations, giving

$$\begin{cases} Y_i = 2.2215 + 0.6912 X_i & \text{if } i = 1, \dots, 12 \\ Y_i = 5.9141 + 0.4787 X_i & \text{if } i = 13, \dots, 20 \end{cases}$$

The following graph shows the estimation results.

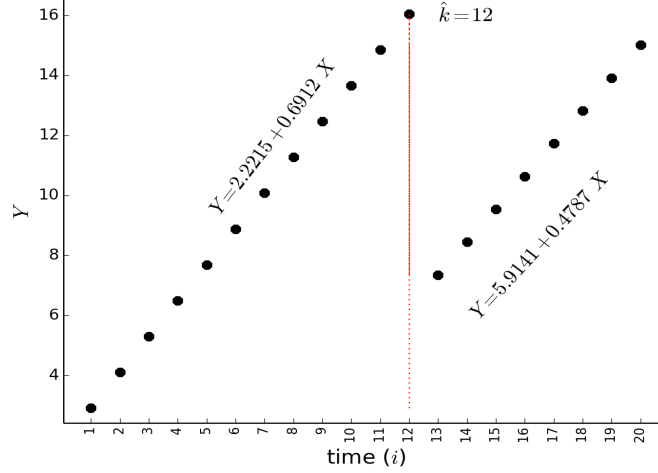


Figure 1: Estimation results for Quandt's data.

4 The change-point detection model for the Weibull distribution

In the following, we note $W(a, b)$ the Weibull distribution with the cumulative distribution function $F(x) = 1 - \exp\left(-\left(\frac{x}{a}\right)^b\right)$.

In this section we assume that a sequence of observations X_1, \dots, X_n represents a change point with :

$$\begin{cases} X_i \sim W(a_1, b_1) & \text{if } i = 1, \dots, k \\ X_i \sim W(a_2, b_2) & \text{if } i = k + 1, \dots, n \end{cases} \quad (9)$$

For $k_0 \in \{4, \dots, n - 4\}$, we build $n - 7$ two subsamples and we order them as the Table 3.

For each subsample, we use the Benard's approximation (Bernard & Bosi-Levenbach [4]) for median ranks, given by:

$$\begin{cases} MR_1(i) = \frac{i - 0.3}{k_0 + 0.4} & i \in \{1, \dots, k_0\} \\ MR_2(i) = \frac{i - 0.3}{n - k_0 + 0.4} & i \in \{k_0 + 1, \dots, n\} \end{cases} \quad (10)$$

The cumulative distribution function of Weibull distribution will be transformed to a

linear function:

$$\ln\left(-\ln(1-F(x))\right) = b \ln x - b \ln a$$

Let $Y = \ln\left(-\ln(1-F(x))\right)$, $A = b$ and $B = -b \ln a$.

To estimate the values of the cumulative distribution function, we use the median rank. For each subsample, we have

$$\begin{cases} Y_i = \ln\left(-\ln(1-MR_1(i))\right) & \text{if } i = 1, \dots, k_0 \\ Y_i = \ln\left(-\ln(1-MR_2(i))\right) & \text{if } i = k_0 + 1, \dots, n \end{cases} \quad (11)$$

Then the model (9) is written :

$$\begin{cases} Y_i = B_1(k_0) + A_1(k_0) \ln(X_i) & \text{if } i = 1, \dots, k_0 \\ Y_i = B_2(k_0) + A_2(k_0) \ln(X_i) & \text{if } i = k_0 + 1, \dots, n \end{cases} \quad (12)$$

By the equations (6), (7) and (8), we deduce the estimators of a_1 , b_1 , a_2 and b_2 :

$$\hat{a}_j = \exp\left(-\frac{\hat{B}_j(k^*)}{\hat{b}_j}\right) \quad , \quad \hat{b}_j = \hat{A}_j(k^*) \quad (j = 1, 2) \quad (13)$$

k_0	Sample ₁ (k_0) ordered	Sample ₂ (k_0) ordered
4	$X_1^{(1)}, \dots, X_1^{(4)}$	$X_2^{(1)}, \dots, X_2^{(n-4)}$
5	$X_1^{(1)}, \dots, X_1^{(5)}$	$X_2^{(1)}, \dots, X_2^{(n-5)}$
\vdots	\vdots	\vdots
$n-4$	$X_1^{(1)}, \dots, X_1^{(n-4)}$	$X_2^{(1)}, \dots, X_2^{(4)}$

Table 3: Distribution of data into two subsamples ordered. $X_j^{(i)}$ denotes the i -th order statistic of sample _{j} (k_0) ($j = 1, 2$).

5 Illustrative data, simulations and application

5.1 Illustrative data.

To illustrate all the steps of the method studied in this paper, we propose the data presented in the Table 4. These data have been simulated from Python 3.3. This is a sample with size 30 representing a point-change. The first 13 data are simulated

according to the Weibull distribution $W(6, 3)$ and the remains are simulated according to the Weibull distribution $W(10, 9)$.

i	1	2	3	4	5	6	7	8	9	10
X_i	5.66	4.78	5.49	6.30	4.69	7.29	4.02	5.01	5.59	3.79
i	11	12	13	14	15	16	17	18	19	20
X_i	5.48	5.48	6.37	8.94	8.81	11.09	8.17	9.86	10.31	9.72
i	21	22	23	24	25	26	27	28	29	30
X_i	10.12	9.66	9.89	10.40	10.01	8.47	7.14	10.30	11.20	10.44

Table 4: Illustrative data.

The steps of the study method are illustrated in the Table 5.

k_0	Sample ₁ (k_0) ordered	Y	Sample ₂ (k_0) ordered	Y	D
4	4.78, 5.49, 5.66, 6.3	-1.75, -0.72, -0.05, 0.61	3.79, 4.02, 4.69, 5.01, 5.48, 5.48, 5.59, 6.37, 7.14, 7.29, 8.17, 8.47, 8.81, 8.94, 9.66, 9.72, 9.86, 9.89, 10.01, 10.12, 10.3, 10.31, 10.4, 10.44, 11.09, 11.2	-3.62, -2.71, -2.23, -1.89, -1.63, -1.41, -1.23, -1.06, -0.92, -0.78, -0.65, -0.54, -0.42, -0.31, -0.21, -0.1, 0.0, 0.1, 0.21, 0.32, 0.43, 0.55, 0.68, 0.82, 1.01, 1.29	2.1898
5	4.69, 4.78, 5.49, 5.66, 6.3	-1.97, -0.97, -0.37, 0.14, 0.71	3.79, 4.02, 5.01, 5.48, 5.48, 5.59, 6.37, 7.14, 7.29, 8.17, 8.47, 8.81, 8.94, 9.66, 9.72, 9.86, 9.89, 10.01, 10.12, 10.3, 10.31, 10.4, 10.44, 11.09, 11.2	-3.58, -2.67, -2.19, -1.85, -1.59, -1.37, -1.18, -1.02, -0.87, -0.73, -0.6, -0.48, -0.37, -0.25, -0.15, -0.04, 0.07, 0.18, 0.29, 0.4, 0.52, 0.66, 0.81, 0.99, 1.28	2.3498
⋮					
13	3.79, 4.02, 4.69, 4.78, 5.01, 5.48, 5.48, 5.49, 5.59, 5.66, 6.3, 6.37, 7.29	-2.93, -2.0, -1.49, -1.13, -0.84, -0.59, -0.37, -0.16, 0.05, 0.25, 0.47, 0.72, 1.08	7.14, 8.17, 8.47, 8.81, 8.94, 9.66, 9.72, 9.86, 9.89, 10.01, 10.12, 10.3, 10.31, 10.4, 10.44, 11.09, 11.2	-3.19, -2.27, -1.78, -1.43, -1.16, -0.92, -0.72, -0.54, -0.37, -0.2, -0.05, 0.11, 0.27, 0.44, 0.62, 0.84, 1.17	1.3247
⋮					

25	3.79, 4.02, 4.69, 4.78, 5.01, 5.48, 5.48, 5.49, 5.59, 5.66, 6.3, 6.37, 7.29, 8.17, 8.81, 8.94, 9.66, 9.72, 9.86, 9.89, 10.01, 10.12, 10.31, 10.4, 11.09	-3.58, -2.67, -2.19, -1.85, - 1.59, -1.37, -1.18, -1.02, - 0.87, -0.73, -0.6, -0.48, -0.37, -0.25, -0.15, -0.04, 0.07, 0.18, 0.29, 0.4, 0.52, 0.66, 0.81, 0.99, 1.28	7.14, 8.47, 10.3, 10.44, 11.2	-1.97, -0.97, -0.37, 0.14, 0.71	3.0564
26	3.79, 4.02, 4.69, 4.78, 5.01, 5.48, 5.48, 5.49, 5.59, 5.66, 6.3, 6.37, 7.29, 8.17, 8.47, 8.81, 8.94, 9.66, 9.72, 9.86, 9.89, 10.01, 10.12, 10.31, 10.4, 11.09	-3.62, -2.71, -2.23, -1.89, - 1.63, -1.41, -1.23, -1.06, - 0.92, -0.78, -0.65, -0.54, - 0.42, -0.31, -0.21, -0.1, 0.0, 0.1, 0.21, 0.32, 0.43, 0.55, 0.68, 0.82, 1.01, 1.29	7.14, 10.3, 10.44, 11.2	-1.75, -0.72, -0.05, 0.61	3.1428

Table 5: The steps of the calculations for the data in the Table 4.

The following figure represents the sum of squared errors defined in equation (3).

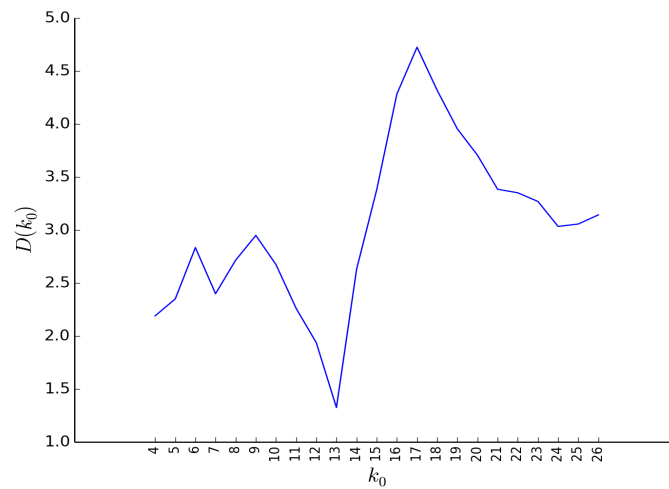


Figure 2: Sum of squared errors defined in equation (3).

The following figure shows the weibull probability plot.

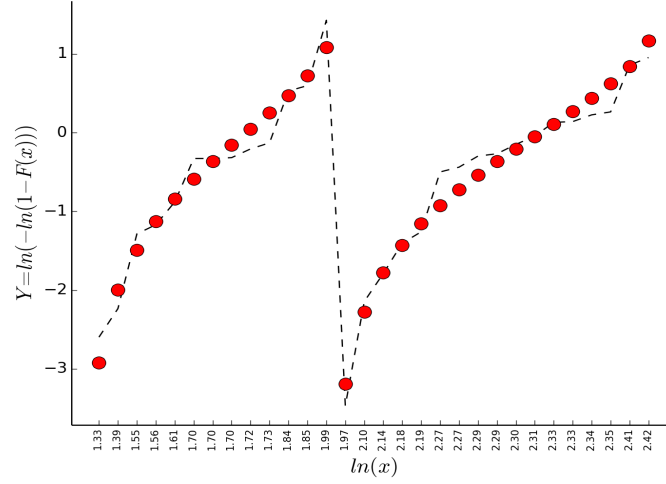


Figure 3: The weibull probability plot for the illustrative data. The estimators are the values : $\hat{a}_1 = 5.78$, $\hat{b}_1 = 6.15$, $\hat{a}_2 = 10.16$, $\hat{b}_2 = 9.83$ and $\hat{k} = 13$.

5.2 Simulations

In order to study the performance of the method, we simulated 1000-samples of sizes $n = 30$ and 100 with a change-points $k = 13$ and 41 . We considered three cases: the first relates to the change in the second parameter of the Weibull distribution, the second case, the change in the first parameter and the third case, the change in both parameters. For each sample we calculated the parameter estimators, the following table summarizes the results obtained for different values of a_1, a_2, b_1 and b_2 .

$a_1 = a_2 = 6, b_1 = 2, b_2 = 5$							
mean of \hat{a}_1	6.1633	mean of \hat{a}_2	6.0181	mean of \hat{b}_1	2.0747	mean of \hat{b}_2	4.8576
std of \hat{a}_1	1.0793	std of \hat{a}_2	0.4221	std of \hat{b}_1	0.7989	std of \hat{b}_2	2.2538
$a_1 = 6, a_2 = 10, b_1 = b_2 = 4$							
mean of \hat{a}_1	6.5194	mean of \hat{a}_2	9.3962	mean of \hat{b}_1	4.1095	mean of \hat{b}_2	3.4012
std of \hat{a}_1	0.9955	std of \hat{a}_2	0.8481	std of \hat{b}_1	2.1639	std of \hat{b}_2	1.2270
$a_1 = 6, a_2 = 10, b_1 = 3, b_2 = 9$							
mean of \hat{a}_1	7.0964	mean of \hat{a}_2	9.9036	mean of \hat{b}_1	2.7622	mean of \hat{b}_2	9.0453
std of \hat{a}_1	1.0870	std of \hat{a}_2	0.4775	std of \hat{b}_1	1.2196	std of \hat{b}_2	4.8527

Table 6: Statistics of estimators of a_1, a_2, b_1, b_2 for size $n = 30$ and change-point $k = 13$.

$a_1 = a_2 = 6, b_1 = 2, b_2 = 5$							
mean of \hat{a}_1	6.1116	mean of \hat{a}_2	6.0298	mean of \hat{b}_1	1.9759	mean of \hat{b}_2	4.8201
std of \hat{a}_1	0.6597	std of \hat{a}_2	0.2304	std of \hat{b}_1	0.5492	std of \hat{b}_2	1.4703
$a_1 = 6, a_2 = 10, b_1 = b_2 = 4$							
mean of \hat{a}_1	6.2783	mean of \hat{a}_2	9.3905	mean of \hat{b}_1	4.2397	mean of \hat{b}_2	3.4163
std of \hat{a}_1	0.8557	std of \hat{a}_2	0.7253	std of \hat{b}_1	1.5346	std of \hat{b}_2	0.9326
$a_1 = 6, a_2 = 10, b_1 = 3, b_2 = 9$							
mean of \hat{a}_1	6.8058	mean of \hat{a}_2	9.9860	mean of \hat{b}_1	2.7517	mean of \hat{b}_2	8.7966
std of \hat{a}_1	0.8215	std of \hat{a}_2	0.2530	std of \hat{b}_1	0.4662	std of \hat{b}_2	1.7604

Table 7: Statistics of estimators of a_1, a_2, b_1, b_2 for size $n = 100$ and change-point $k = 41$.

5.3 Application

We apply the method to study data used by Bhattacharya & Bhattacharjee [3] which represents the Average Monthly Wind Speed (m/s) at kolkata (from 1st March, 2009 to 31st March, 2009).

The following figure represents the sum of squared errors for the Average Monthly Wind Speed (m/s) at kolkata data.

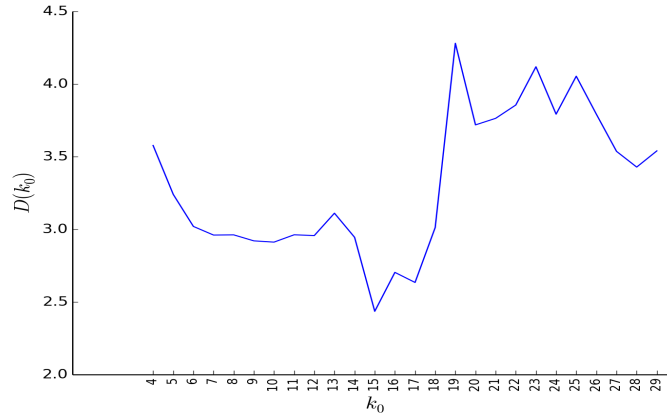


Figure 4: Sum of squared errors defined in equation (3).

The following figure shows the weibull probability plot for the Average Monthly Wind Speed (m/s) at kolkata data.

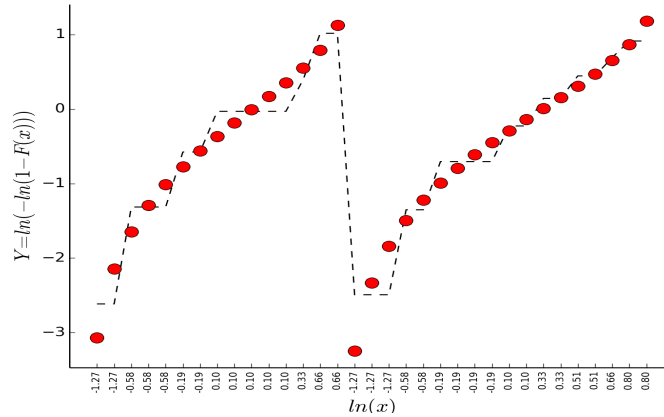


Figure 5: The weibull probability plot. The estimators are the values : $\hat{a}_1 = 1.13$, $\hat{b}_1 = 1.88$, $\hat{a}_2 = 1.27$, $\hat{b}_2 = 1.65$ and $\hat{k} = 15$.

6 conclusion

In this paper, we presented an analytical method of estimating change-point parameters. The results obtained for the Weibull distribution are satisfying. The proposed method is very simple to program that could be easily adapted to other distributions.

References

- [1] Basseville, M., Nikiforov, V. (1993). Detection of Abrupt Changes: Theory and Application. Prentice-Hall, Inc., Englewood Cliffs, N. J.
- [2] Barry, D. and Hartigan, J. A. (1993). A bayesian analysis for change point problems. Journal of the American Statistical Association, 88(421), 309-319.
- [3] Bhattacharya, P. and Bhattacharjee, R. (2010). A Study on Weibull Distribution for Estimating the Parameters. Journal of Applied Quantitative Methods. Vol 5, no:2 (summer 2010), 234-241.
- [4] Bernard, A., Bosi-Levenbach, E.C. (1953). The plotting of observations on probability paper. Stat. Neerlandica, 7, 163-173.
- [5] Bradley, P. Carlin, Alan E. Gelfand, and Adrian F. M. Smith (1992) . Hierarchical bayesian analysis of changepoints problems. Applied Statistics, 41(2), 389-405.
- [6] Dehling, H., Roach, A., and Taqqu, M. S. (2013). Non-parametric change-point tests for long-range dependent data. Scandinavian Journal of Statistics , 40(1), 153-173.
- [7] Ghorbanzadeh, D. (1995). Un test de détection de rupture de la moyenne dans un modèle gaussien. Revue de Statistique Appliquée. XLIII, 67-76.
- [8] Ghorbanzadeh, D. and Picard, D. (1992). Étude Asymptotique et Pratique du comportement de deux tests de détection de Rupture. Statistique et Analyse des données. Vol 16, no. 3, 63-84.
- [9] Ghorbanzadeh, D. and Lounes, R. (2001). Bayesian Analysis for detecting a change in exponential family. Applied Mathematics and Computation . Vol 124, 1-15.
- [10] Guralnik, V.J. and Srivastava (1999). Event detection from time series data. In Proceedings of the 5th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'99), 33-42.
- [11] Gustafsson, F. (1996). The marginalized likelihood ratio test for detecting abrupt changes. IEEE Transactions on Automatic Control, 41(1), 66-78.
- [12] Lucas, M. and Crosier, R. B. (1982). Fast Initial Response for CUSUM Quality-Control Schemes: Give your CUSUM a Head Start. Technometrics, vol. 24, no. 3, 199-205.

- [13] Pettitt, A.N. (1979). A non-parametric approach to the change-point problem. Appl. Statist., 28, 126-135.
- [14] Quandt, R. (1958). The estimation of the parameters of a linear regression system obeying two separate regimes. Journal of the American Statistical Association, 53, 873-880.
- [15] Ritov, Y. (1990). Decision Theoretic Optimality of the Cusum Procedure. The Annals of Statistics, vol. 18, no. 3, 1464-1469.
- [16] Yashchin, E. (1985). On the analysis and design of CUSUM-Shewhart control schemes. IBM Journal Research and Development, vol. 29, no. 4, 377-391.