

# Introduction to change point detection

## Computer lab worksheet

Euan McGonigle

### Instructions

- Answers are available to you at: if you feel you're spending too much time on any questions particularly Question 1, you can copy-paste the code to move onto questions 2 – 4.
- Bonus questions: Q 1(e) is a bonus, as is Question 5. Feel free to skip these if you want.
- For any functions you're not sure how to use, use `?function_name` in R to get help (or ask me!).
- In the time we have, I don't expect you to finish all questions. If you can, try to finish at least Q2.

### Question 1: Creating a test statistic to estimate change points

- (a) Without using pre-existing change point detection libraries, write a function to calculate **either**:
- the CUSUM statistic vector  $\{T_{0,k,n}\}_{k=1}^{n-1}$  for a change point in the mean,
  - the MOSUM statistic vector  $\{T_G(k)\}_{k=G}^{n-G}$  for a change point in the mean, for a given bandwidth  $G$ .
- (b) To set a threshold for declaring changes, we need to know the noise level  $\sigma$ . In reality, this is unknown; how could we estimate it?

*Hint:* the sample standard deviation would be positively biased by the change points. Medians are more robust than means: how can we transform the data to remove most of the effect of the change points, before using a median-like analogue of the standard deviation? If in doubt, Google is your friend.

- (c) Implement your chosen method of estimating  $\sigma$ , and add the functionality to your CUSUM or MOSUM function from part (a) so that the test statistic calculation is scaled by your estimate  $\hat{\sigma}$ .
- (d) What is the computational complexity of the code you wrote in part (a)? How fast could it be?
- (e) (Bonus) For simultaneous estimation of multiple change points with the MOSUM approach, we can calculate all  $\hat{\theta}$  that satisfy:

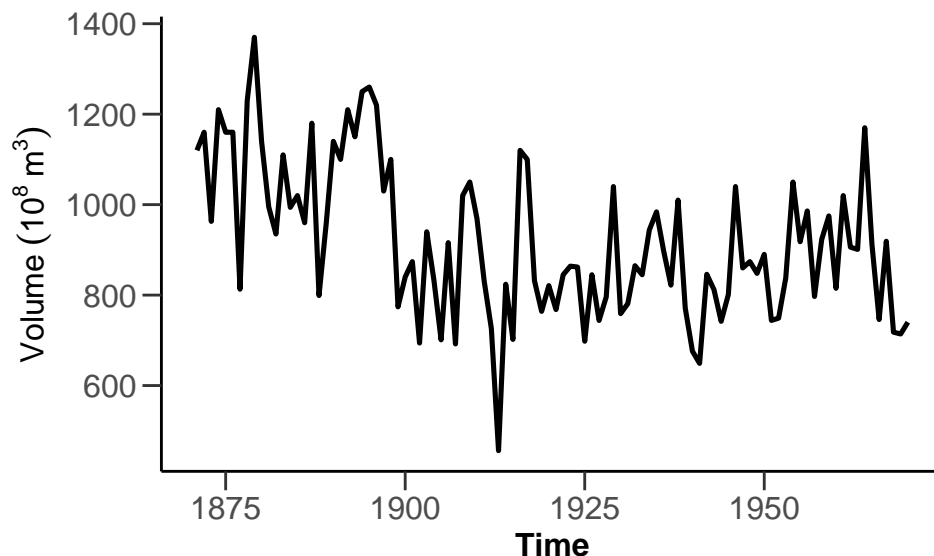
$$T_G(\hat{\theta}) > D \quad \text{and} \quad \hat{\theta} = \operatorname{argmax}_{k: |k-\hat{\theta}| \leq \eta G} T_G(k).$$

for some window fraction  $\eta \in (0, 1)$  and a threshold  $D$ . That is,  $\hat{\theta}$  is declared a change point if it is a local maximiser of  $T_G(k)$  over a sufficiently large interval of size  $\eta G$ , at which the threshold  $D$  is exceeded.

Extend the MOSUM function to return as output the change point estimators satisfying this criterion.

## Question 2: Nile annual river flow

In this question you will use the code you've written so far to analyse data collected on the river Nile. Data in the file `nile_volume.txt` records measurements of the annual volume (in units  $10^8 m^3$ ) of discharge from the Nile River at Aswan for the years 1871 to 1970. The measurements are of meteorological importance as evidence of a possible abrupt change in the rainfall levels around the turn of the 20th century.



For whichever method (CUSUM or MOSUM) you did not code up in Q1, you can use the code in the answers for this question.

- (a) Load the data into R from the `nile_volume.txt` file, and plot it. By eye, where does it look like there could be a change in mean?
- (b) Calculate the CUSUM statistic for the Nile volume data. Using the threshold  $D = \sqrt{2 \log(n)}$ , perform a test to decide if there is a change point. If there is one, where is it?
- (c) Calculate the MOSUM statistic for the Nile volume data using a bandwidth  $G = 25$ . Using the threshold  $D = 3.4$ , compute the change point estimators:
  - either using the code from Q1 (e) with  $\eta = 0.4$ , or
  - by eye by plotting the MOSUM test statistic.

Does this agree with your answer from part (b)?

- (d) Now use the MOSUM method with bandwidth  $G = 40$ . How does your answer differ from part (c)?
- (e) The estimated mean signal  $\hat{f}_t$  can be calculated using sample means of the segments defined by the estimated change points. For the CUSUM method, add the estimated  $\hat{f}_t$  to your plot from part (a).

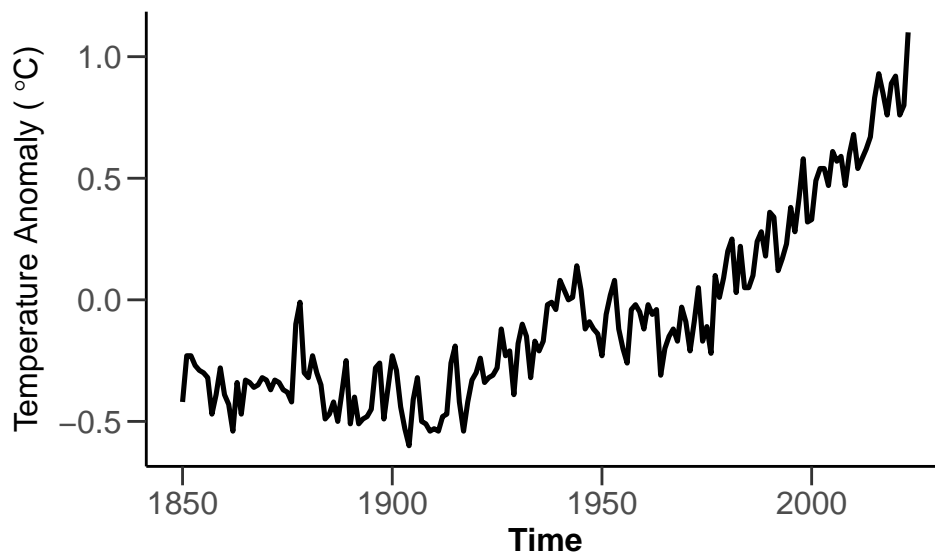
### Question 3: Using the `mosum` and `changepoint` packages

- (a) Load the `changepoint` and `mosum` R packages into your R workspace. Using the help files, get acquainted with the `cpt.mean()` and `mosum()` functions.
- (b) Setting the argument `method = "AMOC"`, use the `cpt.mean()` function on the Nile data set. Do the results agree with your answer from question 1(a)?
- (c) Use the `cpt.mean()` function to apply the PELT algorithm on the Nile data set. Is the answer the same as part (b)? **Note:** you will need to standardise the data first (subtract the mean and divide by the standard deviation).
- (d) By investigating the output returned by the `cpt.mean()` function, what was the value of the penalty function used?
- (d) Using the bandwidth  $G = 25$ , use the `mosum()` function on the Nile data set. Do the results agree with your answer from question 1(b)?
- (e) By investigating the output returned by the `mosum()` function, what are the p-value(s) of the detected change points?

*Python users:* you can use the `ruptures` function `Pelt()` and the `mosum` function `mosum()` for this question. See <https://centre-borelli.github.io/ruptures-docs/> and <https://pypi.org/project/mosum/> for help.

#### Question 4: Global yearly mean sea temperature anomalies

Global mean surface temperature series help enable the monitoring of global warming. The warming can be quantified by, for example, a change from a base period used as reference. The Hadley Centre/Climatic Research Unit (HadCRUT) surface temperature data set provides the annual global temperature anomalies from 1850-2023, where anomalies are calculated relative to the 1961–1990 period. The data are shown below, which can be seen to exhibit a gradual upward trend.



- (a) Load the data into R from the `temp_anomalies.txt` file, and plot it. Using e.g. the `lm()` function, fit two linear trends to the data: one from 1850 to 1963, and one from 1964 to 2023, and add these lines to the plot, as well as a vertical line through year 1963.
- (b) Using your own functions, and functions from the `changepoint` and `mosum` packages, fit 4 models to the data set;
  - 1. a constant mean model with no change points,
  - 2. a mean change point model,
  - 3. a linear trend model with no change points.
  - 4. a linear trend change point model.

Looking at the data, are there any other models that might be appropriate?

*Hint:* if you are struggling, the `EnvCpt` R package has everything you need.

*Python users:* you can use the `Pelt` function and change the `model` parameter.

- (c) Pick the “best” change point model from your candidates computed in part (b), and justify your choice. Add the estimated mean/trend function onto a plot of the data.

*Hint:* the Akaike information criterion (AIC) is given by  $AIC = 2k - 2\log(\hat{L})$ , where  $k$  is the number of parameters of the model, and  $\hat{L}$  is the maximised value of the likelihood function for the model. The smaller the AIC, the “better” the model. For normally distributed data, you can use `dnorm()` in R to compute the likelihood.

### Question 5 (Bonus): Multivariate data

- (a) Suppose we want to find mean change points in a time series  $\{X_t\}_{t=1}^n$ , where  $X_t = (X_{1t}, \dots, X_{pt})^\top$  is a  $p$ -dimensional vector. One approach is to calculate a test statistic for each variable, and then combine the results across variables to give a single test statistic. For example, for the  $i$ -th variable of  $\{X_t\}_{t=1}^n$ , denoted  $\{X_{it}\}_{t=1}^n$ , the MOSUM test statistic is given by

$$T_G(k, i) = \frac{1}{\sqrt{2G}} \left( \sum_{t=k+1}^{k+G} X_{it} - \sum_{t=k-G+1}^k X_{it} \right).$$

Then, the MOSUM statistic  $T_G(k)$  for a change point in  $\{X_t\}_{t=1}^n$  can be computed by aggregating the  $\{T_G(k, i)\}_{i=1}^p$  using some aggregating function  $f$ , so that

$$T_G(k) = f(T_G(k, 1), \dots, T_G(k, p)).$$

What possibilities could we use for the aggregating function  $f$ ?

- (b) Using your aggregating function  $f$ , write a function to compute the MOSUM or CUSUM statistic for a change in the mean vector of a multivariate time series.
- (c) Simulate a data set of length  $n = 500$  and dimension  $p = 40$  as follows. Generate  $\{X_t\}_{t=1}^{250}$  from the standard  $p$ -dimensional normal distribution, and generate  $\{X_t\}_{t=251}^{500}$  from the  $p$ -dimensional normal distribution with identity covariance matrix, and mean vector given by  $\mu = 1.4 \times (1/\sqrt{p}, \dots, 1/\sqrt{p})^\top$ . Use the method from part (b), with  $G = 80$ , to compute the test statistic for a change in mean. Is the change point easy to see?
- (d) Re-do part (c), but generate  $\{X_t\}_{t=251}^{500}$  from the  $p = 200$ -dimensional normal distribution with identity covariance matrix, and mean vector given  $\mu$  with first 2 entries equal to 0.8, and last  $p - 2$  entries equal to 0. Is the change point easy to see?