

Tutorial: Detecting a breakpoint a time series

Jan Verbesselt, Achim Zeileis

November 24, 2014

Abstract

This tutorial explain the basics of structural change test using with the BFAST concept. Basic principles are illustrated using the *Nile* data set (add reference).

1 Introduction

First load the BFAST package, which also loads the required packages for this small tutorial (i.e. strucchange, zoo).

```
> library("bfast")
> library("zoo")
> library("strucchange")
```

We will illustrate basic principles using the following time series of Annual Flow data of the river Nile. For more information also check the papers of Zeileis et al. (2005).

```
> plot(Nile, ylab="Annual Flow of the river Nile")
> abline(h= mean(Nile), col='blue')
```

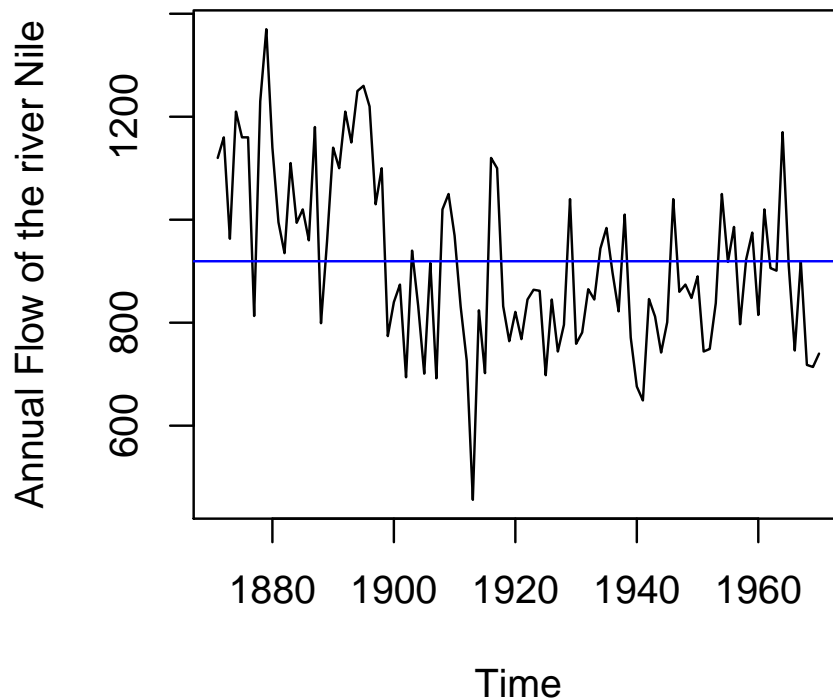


Figure 1: Annual river flow of the river Nile. The mean river flow (m³/sec) is plotted as a blue horizontal line.

2 CUSUM and MOSUM illustrated

Here the CUSUM and MOSUM of the residuals of the Nile data using a constant (intercept) as explanatory variable in the model.

```
> plot(merge(
+   Nile = as.zoo(Nile),
+   zoo(mean(Nile), time(Nile)),
+   CUSUM = cumsum(Nile - mean(Nile)),
+   zoo(0, time(Nile)),
+   MOSUM = rollapply(Nile - mean(Nile), 15, sum),
+   zoo(0, time(Nile))
+ ), screen = c(1, 1, 2, 2, 3, 3), main = "", xlab = "Time",
+   col = c(1, 4, 1, 4, 1, 4)
+ )
```

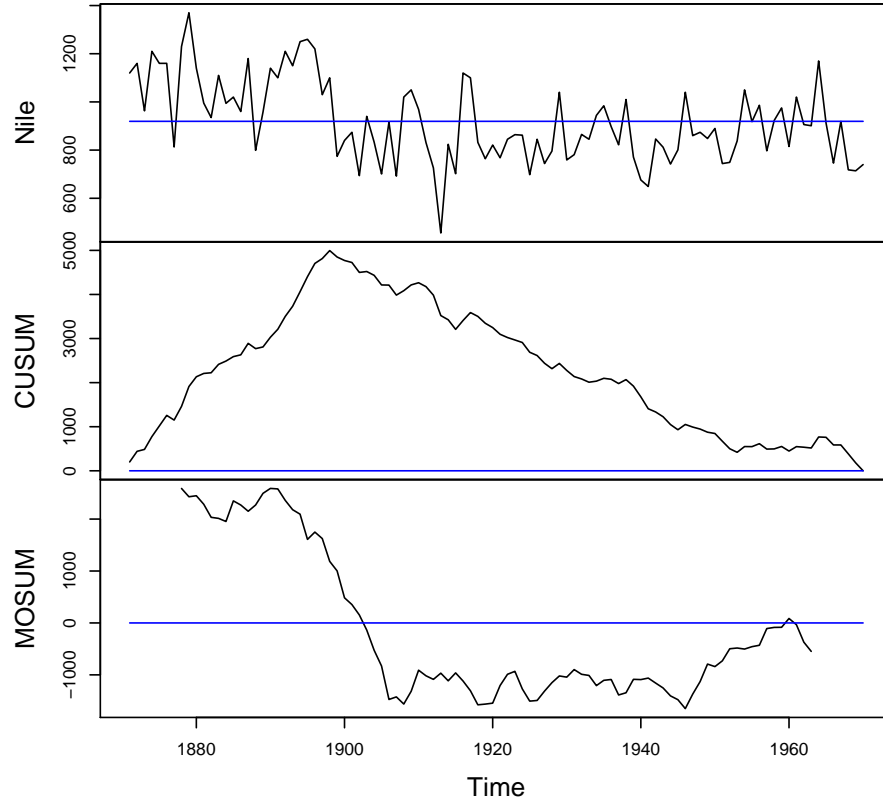
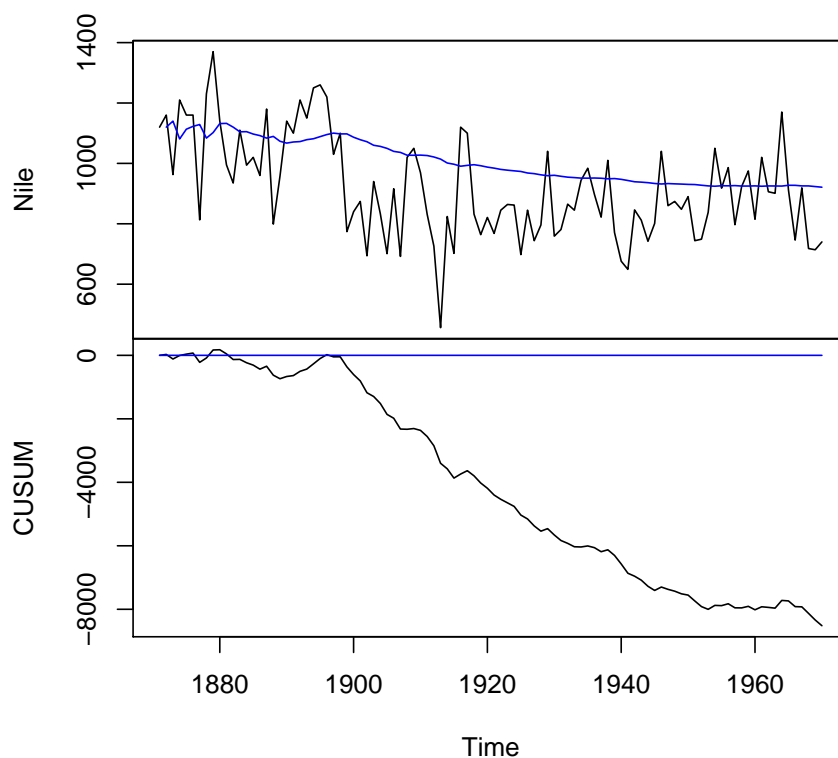


Figure 2: Nile data. MOSUM test. CUSUM test.

Fig. 2 illustrates how the CUSUM and MOSUM test evolve over time where CUSUM is a cumulative sum of the residuals and MOSUM is a moving sum (using a window of 15 years) of the residuals. The cumulative and moving sum of the residuals should fluctuate around zero (blue line) however significant deviation occurs from 1918 (?) onwards.

The recursive CUSUM illustrated.

```
> plot(merge(
+   Nile = as.zoo(Nile),
+   zoo(c(NA, cumsum(head(Nile, -1))/1:99), time(Nile)),
+   CUSUM = cumsum(c(0, recresid(lm(Nile ~ 1)))),
+   zoo(0, time(Nile))
+ ), screen = c(1, 1, 2, 2), main = "", xlab = "Time",
+ col = c(1, 4, 1, 4)
+ )
```



Rec-CUMSUM. The blue line in the first panel simply shows the mean of all observations prior to it, i.e., $\text{prediction}[t] = \text{mean}(\text{Nile}[1:(t-1)])$. The black line in the second panel shows the cumulated recursive residuals (= standardized one-step-prediction errors). You can see that up to the building of the dam, the residuals are approximately zero and then they start deviating from zero.

The following R lines reproduce the OLS-CUSUM, -MOSUM, and -Rec-CUSUM plots which are suitably scaled (see Zeileis et al. 2005)

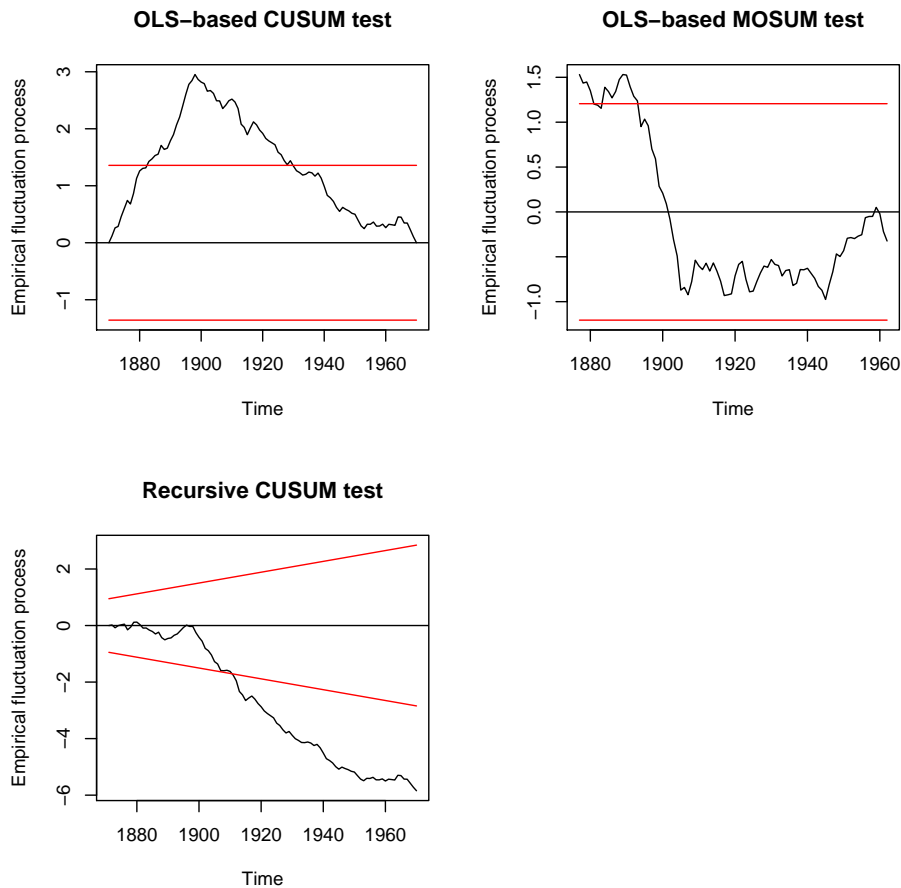
```
> ocus.nile <- efp(Nile ~ 1, type = "OLS-CUSUM")
> omus.nile <- efp(Nile ~ 1, type = "OLS-MOSUM")
> rocus.nile <- efp(Nile ~ 1, type = "Rec-CUSUM")
```

These plots illustrate that a significant structural change is detected by the CUSUM and the MOSUM test. Especially the CUSUM illustrate that around 1910 the empirical fluctuation process based on the CUSUM residuals goes outside the boundaries. The MOSUM based plot less clearly illustrates the potential structural change that occurs (i.e. efp does not go outside the significance boundaries expect around early 1900. There is however a clear drop in the scaled MOSUM residuals visual which indicates a significant change in the data variation of the residuals.

```

> opar <- par(mfrow=c(2,2))
> plot(ocus.nile)
> plot(omus.nile)
> plot(ocus.nile)
> par(opar)

```



3 Fitting a piecewise linear model

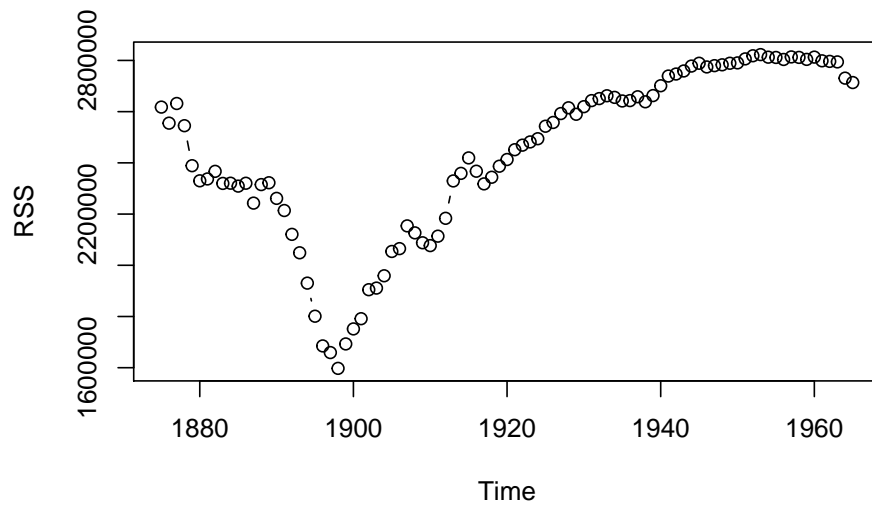
The following section illustrate how the location of the structural change in the Nile time series is determined.

First, determine the the minimum Residual Sum of Squares RSS to determine the position of the breakpoint:

```

> plot(1870 + 5:95, sapply(5:95, function(i) {
+   before <- 1:i
+   after <- (i+1):100
+   res <- c(Nile[before] - mean(Nile[before]), Nile[after] - mean(Nile[after]))
+   sum(res^2)
+ })), type = "b", xlab = "Time", ylab = "RSS")

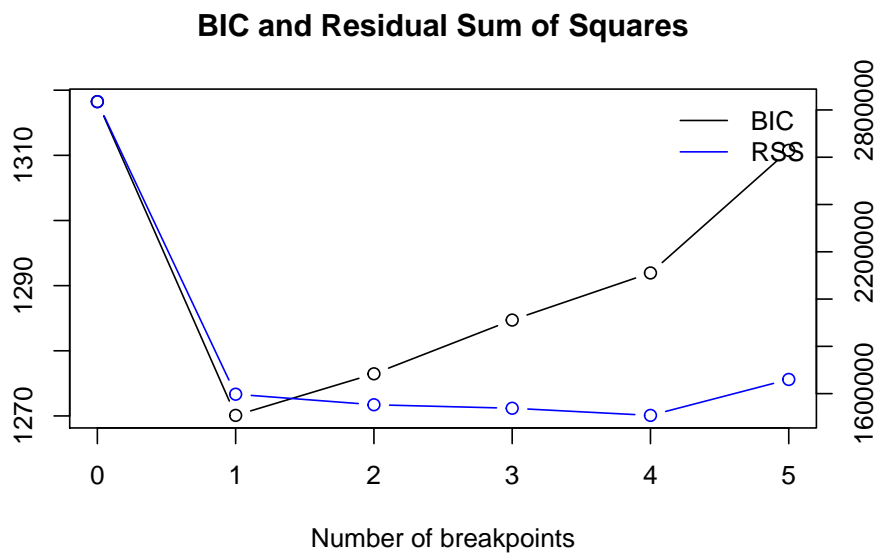
```



Second, use the strucchange functionality to determine the *Date* of change.

```
> bp.nile <- breakpoints(Nile ~ 1)
> nile.fac <- breakfactor(bp.nile, breaks = 1 )
> fm1.nile <- lm(Nile ~ nile.fac - 1)

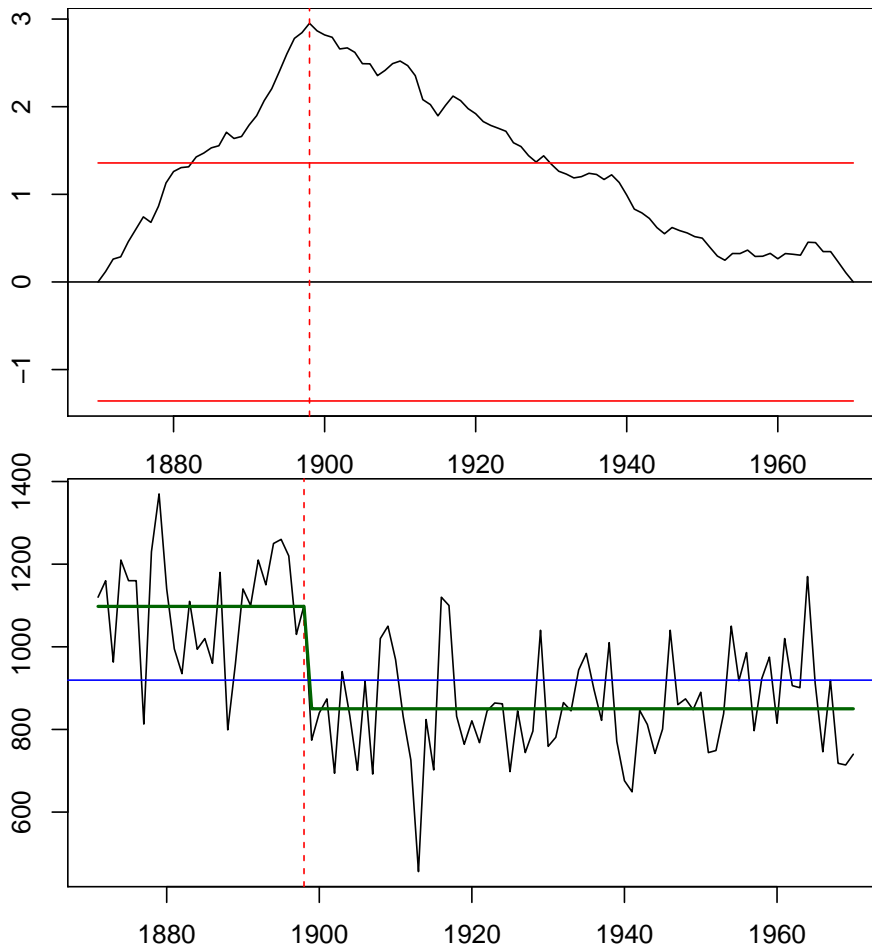
> plot(bp.nile)
```



```

> opar <- par(mfrow=c(2,1), mar=c(2,2,0,2))
> plot(ocus.nile, alt.boundary=F,main="")
> abline(v= 1898, lty=2, col='red')
> plot(Nile, ylab="Annual Flow of the river Nile")
> abline(h= mean(Nile),col='blue')
> abline(v= 1898, lty=2, col='red')
> lines(ts(predict(fm1.nile),start=1871,freq=1), col='darkgreen',lwd=2)
> par(opar)

```



4 Advanced: model selection

Next, critical when detecting breakpoints in a time series is the selection of the right model. The model should accommodate for normal data variation within a time series. The selection of the model should be done based on a time series without any significant breaks (if known).

```

> ## set up time series
> ndvi <- as.ts(zoo(cbind(a = som$NDVI.a, b = som$NDVI.b), som$Time))
> ndvi <- window(ndvi, start = c(2003, 1), end = c(2009, 23))
> ## prepare the time series and related variables needed within the linear
> ## regression

```



```

>
> d1 <- bfastpp(ndvi, order = 3)
> ## lm fit
> f1 <- lm(response ~ xreg, data = d1)
> f2 <- lm(response ~ trend + xreg, data = d1)
> f3 <- lm(response ~ trend + xreg + harmon, data = d1)
> ## check the fit and fit the best model
> anova(f1,f2,f3)

```

Analysis of Variance Table

```

Model 1: response ~ xreg
Model 2: response ~ trend + xreg
Model 3: response ~ trend + xreg + harmon
  Res.Df    RSS Df Sum of Sq    F Pr(>F)
1     159 1.33090
2     158 1.31244  1    0.01847  3.7497 0.05467 .
3     152 0.74858  6    0.56386 19.0820 < 2e-16 ***

```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

> ## here RSS of f3 is significantly smaller than f2 although 6 extra variables
> ## are used.
> z0 <- lm(response ~ 1, data = d1)
> z1 <- lm(response ~ harmon, data = d1)
> z2 <- lm(response ~ harmon + xreg, data = d1)
> z3 <- lm(response ~ harmon + xreg + trend, data = d1)
> anova(z0, z1, z2, z3)

```

Analysis of Variance Table

```

Model 1: response ~ 1
Model 2: response ~ harmon
Model 3: response ~ harmon + xreg
Model 4: response ~ harmon + xreg + trend
  Res.Df    RSS Df Sum of Sq    F    Pr(>F)
1     160 2.26204
2     154 1.05878  6    1.20326 40.7206 < 2.2e-16 ***
3     153 0.76929  1    0.28948 58.7802 1.948e-12 ***
4     152 0.74858  1    0.02071  4.2058    0.042 *

```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

> ## This also shows that the trend effect is just significant at the 5\% level.
> ## AIC also improves but BIC doesn't:
>
> AIC(z0, z1, z2, z3)

```

	df	AIC
z0	2	-225.7891
z1	8	-336.0126
z2	9	-385.4358
z3	10	-387.8301

```
> BIC(z0, z1, z2, z3)
```

	df	BIC
z0	2	-219.6263
z1	8	-311.3613
z2	9	-357.7031
z3	10	-357.0160

```
>
>
```

5 More information

More information can be found on the following website <http://bfast.r-forge.r-project.org/> and in the BFAST papers mentioned on that website.