

# ACR2Full

Martin Cavarga

02 February, 2020

## Advanced Methods of Time Series Analysis Applied to Quarterly Estimates of Unemployment Rate

### Introduction

The chosen source of data is the Labour Force Survey (LFS) quarterly estimates of unemployment rate in the UK since March 1971, up to March 2018.

---

### 1. Elementary Modeling by an AR Process

We begin by extracting the data from a downloaded file

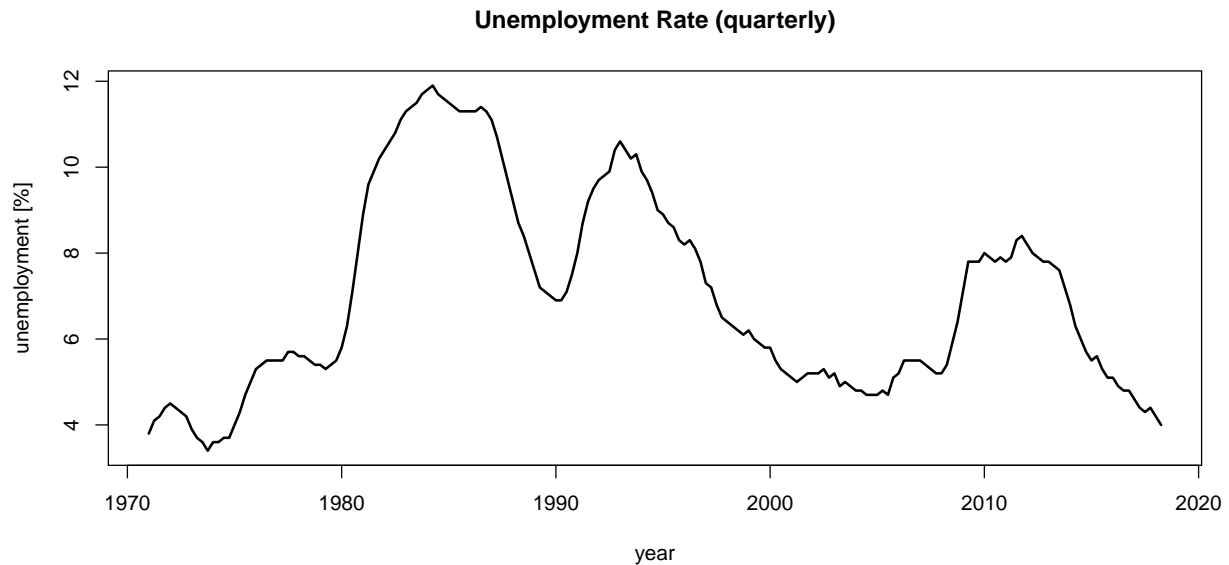
```
## head.dat.time.
## 1      1971 Q1
## 2      1971 Q2
## 3      1971 Q3
## 4      1971 Q4
## 5      1972 Q1
## 6      1972 Q2

## head.months.
## 1      0
## 2      3
## 3      6
## 4      9
## 5      0
## 6      3

## head.years.
## 1      1971
## 2      1971
## 3      1971
## 4      1971
## 5      1972
## 6      1972

## head.years.
## 1      1971.00
## 2      1971.25
## 3      1971.50
## 4      1971.75
## 5      1972.00
## 6      1972.25
```

## 1.1: Data Plot

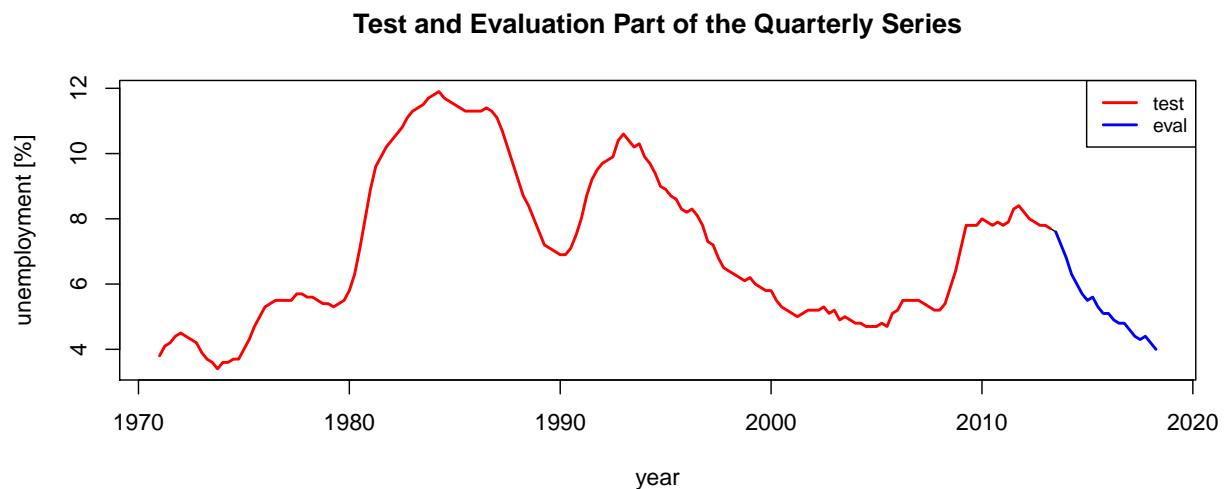


Now even though we are working with annual data there should not be any seasonal components or trend since the results do not depend on periodic observable phenomena, but rather the complex economic situation over multiple decades. Also the data may include an exponentially decaying decrease in unemployment, but only after year 1980, which would suggest a regime-switching stochastic process.

## 1.2: Test Part and Evaluation Part of the Time Series

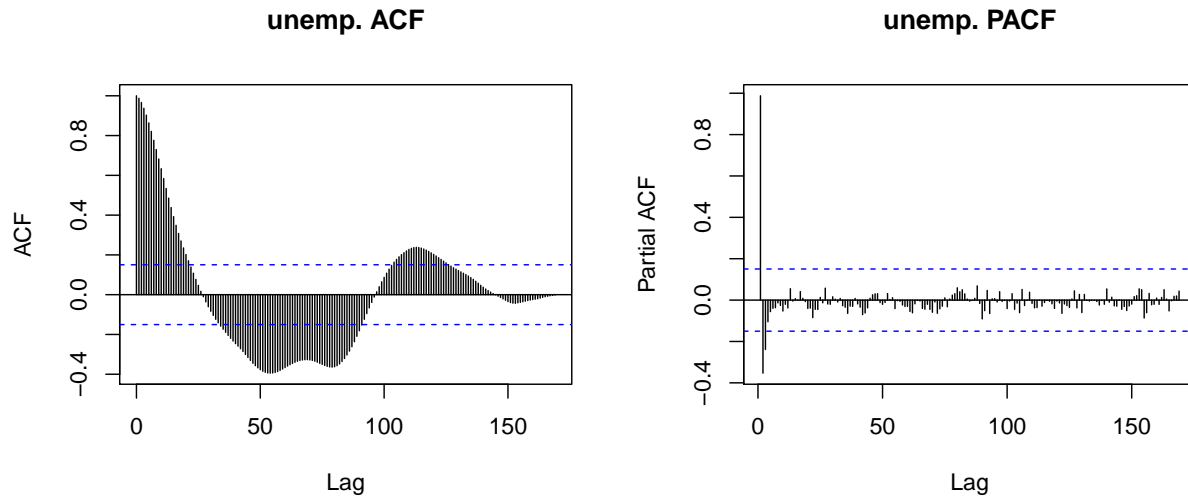
Now we separate the time series into test part, where a suitable model of a stochastic process will be found, and the evaluation part, where predictions given by such model are evaluated. Since our dataset contains quarterly data, we choose the length of the evaluation part of the time series as  $L = k * 4$  where  $k$  is an arbitrary (and sufficiently small) positive integer.

```
## [1] 20
```



### 1.3: Mean, Variance, ACF, and PACF of the Test Part

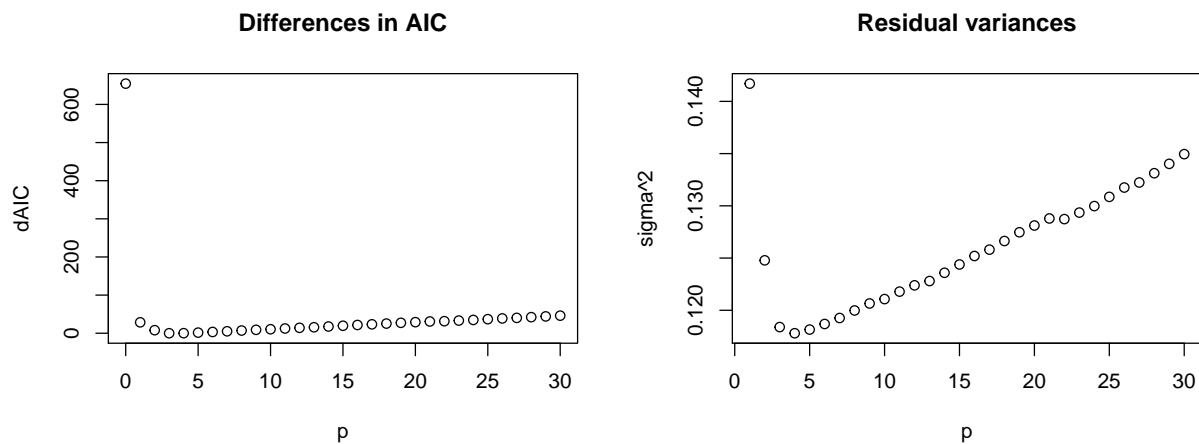
```
##      Min.      Max.      Mean      Median      Variance
## 3.400000 11.900000  7.188824  6.900000  5.661827
```



As we mentioned in section 1.1, the underlying process which gave rise to the observed results is aperiodic, yet it is undoubtedly a process with memory. Unemployment rate strongly depends (aside from other important aspects) on its own history which might extend generations into the past. The results are, however, significantly influenced by external phenomena, such as the global economic crisis in late 2000's.

### 1.4: Finding a Suitable AR Model

Since the economic situation and the job market remembers its past, we choose a simple  $AR(p)$  process with parameter  $p$  corresponding to the number of steps after which the process still “remembers” its past. The inbuilt `ar()` function automatically finds the model with the lowest AIC (Akaike's Information Criterion). And by plotting the `$aic` parameter we obtain differences  $AIC_{min} - AIC_k$  for all models.



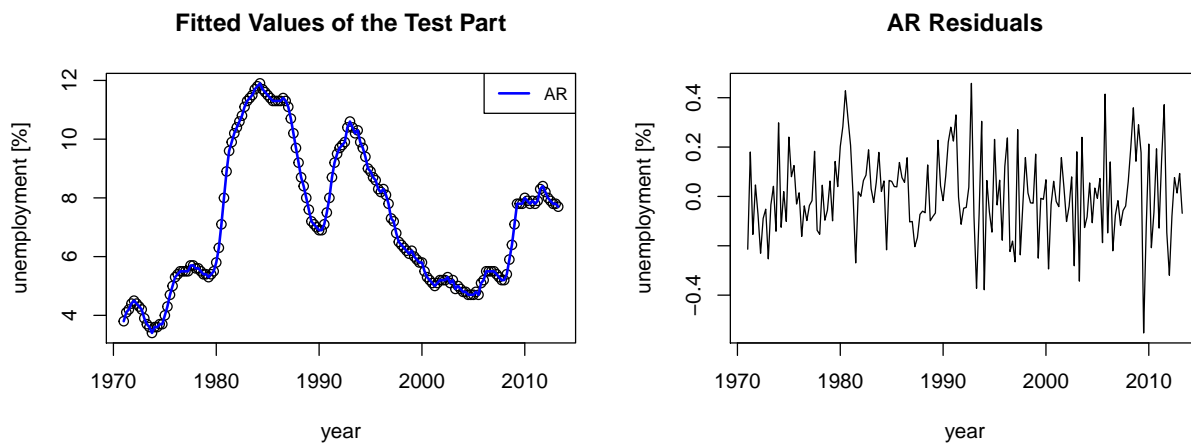
As we can see in the figures, the lowest variance of residues corresponds to an  $AR(3)$  process with coefficients:

```
##      [,1]      [,2]      [,3]
## coef 1.2519124 -0.03440575 -0.2384831
## se   0.0753756  0.12294642  0.0753756
```

```
## [1] 3
```

Unfortunately, the `ar` function does not return fitted values, thus we need to model the time series via the `arima` function using the AR order `p` from the previous result.

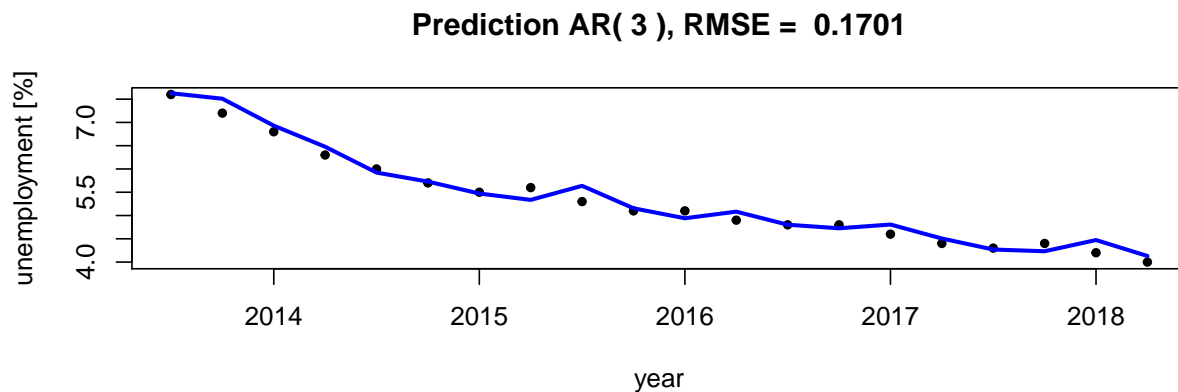
```
##  
## Call:  
## arima(x = as.numeric(unempseries$test), order = c(p, 0, 0))  
##  
## Coefficients:  
##          ar1          ar2          ar3  intercept  
##          1.6000    -0.4176    -0.1951     6.8458  
## s.e.    0.0752     0.1412     0.0754     0.9580  
##  
## sigma^2 estimated as 0.0289:  log likelihood = 56.88,  aic = -103.75
```



The given AR model seems to fit the time series very well, which may be due to its low oscillation rate.

### 1.5: 1-Step Predictions Over the Evaluation Part

```
## [1] 0.1701227
```



## 1.6.: Conclusion

Since it has very low rate of local oscillation, but does not have an easily predictable systematic pattern, the analyzed unemployment rate time series seems to be well-estimated by an  $AR(3)$  process with low prediction errors. However, as we mentioned earlier we might be dealing with a ‘regime-switching’ process. Further analysis will be carried out in the next chapter.

---

## 2. Finding the Parameters of a SETAR Model

As we mentioned, the unemployment time series might be a result of a regime-switching process. Naturally, the behavior of the unemployment rate in a given country should depend on the current economic situation. The change in the local economy can be described via a set of “thresholds” which determine whether the stochastic process changes its regime. The regime of a stochastic process is defined as a unique ARMA or any other linear stochastic process with unique parameters. We begin by finding the parameters of a Self-Exciting Threshold Autoregressive (SETAR) process, that is: a process whose regime is described by a random variable determined by the very process itself, more specifically its history of up to  $d$  steps behind, which in an essence means that the process “influences its regime” up to  $d$  time steps into the future.

For the purposes of this analysis we consider only 2 regimes, namely the regime of “job crisis” when the unemployment rate may fluctuate or drop more wildly compared to the regime of “job stability” when the unemployment rate stabilizes or grows.

### 2.1: Implementation of Useful Functions

First, we define a, so called, “indicator function” which essentially returns a boolean value from a given input process value  $x$  and threshold value  $c$ :

```
Indicator <- function(x, c) ifelse(x > c, 1, 0)
```

Afterwards, we define the basis function for a single regime

```
Yt <- function(x, t, p) c(1, x[(t - 1):(t - p)])
```

which can then be used in the “basis” for two regimes:

```
Xt <- function(x, t, p, d, c, z = x) {  
  # z is the threshold variable  
  I <- Indicator(z[t - d], c)  
  Y <- Yt(x, t, p)  
  c((1 - I) * Y, I * Y)  
}
```

We can test the function on a given subset of the unemployment time series. Due to the fact that the examined time series is rather ‘smooth’, for further use, we will examine its differences:

```
xt <- na.omit(diff(as.numeric(dat$unemp))) # take differences in data  
x_train <- xt[seq_along(unempseries$test)] # extract test part  
nt <- length(x_train)  
x_eval <- xt[(nt + 1):length(xt)] # extract eval part  
xt <- x_train # set the source data to test part values
```

```
## [1] 0.1 0.2
```

```
## [1] 0.0 0.0 0.0 1.0 0.1 0.3
```

```
## [1] -0.3 -0.1
```

```
## [1] 1.0 -0.3 -0.1 -0.2 0.0 0.0 0.0 0.0
```

As we can see, in the first case, with time series values crossing zero from above, the latter half of the coefficient vector gets expressed, corresponding to the series assuming the second regime.

Then we need a function defining a deterministic skeleton of the model:

```
SkeletonSETAR <- function(x, t, p, d, c, theta, z = x) theta %*% Xt(x, t, p, d, c, z)
```

where `theta` corresponds to the parameter vector, for example:

```
##      [,1]
## [1,]  0.4
```

and the last group of functions we need for the upcoming procedure are functions for the information criteria of a SETAR model:

```
# Akaike
AIC_SETAR <- function(orders, regimeDataCount, resVariances) {
  sum(regimeDataCount * log(resVariances) + 2 * (orders + 1))
}

# Bayesian
BIC_SETAR <- function(orders, regimeDataCount, resVariances) {
  sum(regimeDataCount * log(resVariances) + log(regimeDataCount) * (orders + 1))
}

#' and test it out:
AIC_SETAR(c(2, 2), c(10, 10), c(0.5, 0.7))

## [1] 1.501779
BIC_SETAR(c(2, 2), c(10, 10), c(0.5, 0.7))

## [1] 3.317289
```

## 2.2: The Estimation of Parameters of a SETAR Model

Given a dataset `x` and parameters `p` (AR order), `d` (SETAR delay), and the threshold `c` we find the coefficients of a SETAR model with these parameters by performing a multivariate linear regression. The coefficient vector `PhiParams` is the vector of unknowns of a linear system with matrix  $\mathbf{X}$  and a right-hand-side vector  $\mathbf{y}$  given by the time series. Although for higher values of `p` the inversion of matrix  $\mathbf{X}^T \mathbf{X}$  (with dimensions  $(2p + 2) \times (2p + 2)$ ) might be computationally demanding, we will determine the covariance matrix, i.e.:  $(\mathbf{X}^T \mathbf{X})^{-1}$  using a function `inv` from the `matlib` package:

```
suppressMessages(pkgTest("zeallot"))
suppressMessages(pkgTest("matlib"))

EstimSETAR <- function(x, p, d, c) {

  resultModel <- list()
  resultModel$p = p; resultModel$d = d; resultModel$c = c;
  resultModel$data = x; n = length(x); resultModel$n = n;
  k <- max(p, d)

  X <- as.matrix(apply(as.matrix((k + 1):n), MARGIN=1, function(t) Xt(x, t, p, d, c) ))
  y <- as.matrix(x[(k + 1):n])

  A = crossprod(t(X), t(X)); b = crossprod(t(X), y)

  if (abs(det(A)) > 0.000001) {
    inv <- inv(A)
    sol_phi <- as.numeric(t(inv %*% b)); sol_se <- sqrt(diag(inv)/n);
```

```

eps <- 0.01;

# filter out those coeffs that are of the same order of magnitude as their errors
filter <- sapply(1:(2*(p + 1)), function(i) ifelse(
  abs(sol_phi[i]) <= 2 * abs(sol_se[i]), 0, 1)
)

sol_phi <- sol_phi * filter
sol_se <- sol_se * filter

solution <- cbind(phi = sol_phi, se = sol_se)

resultModel$PhiParams <- solution[,1] # solving (X'X)*phi = X'y
resultModel$PhiStErrors <- solution[,2] # standard errors
skel <- crossprod(X, resultModel$PhiParams); resultModel$skel <- skel;
resultModel$residuals <- (y - skel)
resultModel$resSigmaSq <- 1 / (n - k) * sum(resultModel$residuals ^ 2)

return(resultModel)
} else {
  return(NA)
}
}

```

After performing this procedure for multiple parameters, i.e.: searching the discrete parameter space, we further process the model with minimum residual square sum. For that we'll use:

```

EstimSETAR_postproc <- function(model) {
  x <- model$data; k <- max(model$p, model$d); c <- model$c; n <- model$n;
  y <- as.matrix(x[(k + 1):n])
  skel <- model$skel; model$skel <- NULL; #skel attribute no longer needed

  model$n1 <- sum(apply(as.matrix(x), MARGIN = 1, function(xt) (1 - Indicator(xt, c))))
  model$n2 <- sum(apply(as.matrix(x), MARGIN = 1, function(xt) Indicator(xt, c)))

  model$resSigmaSq1 <- sum(
    apply(as.matrix(seq_along(y)), MARGIN = 1,
      function(t) ifelse((1 - Indicator(y[t], c)), (y[t] - skel[t])^2, 0))) / (model$n1 - k)
  model$resSigmaSq2 <- sum(
    apply(as.matrix(seq_along(y)), MARGIN = 1,
      function(t) ifelse(Indicator(y[t], c), (y[t] - skel[t])^2, 0))) / (model$n2 - k)

  model$AIC <- AIC_SETAR(c(p, p), c(model$n1, model$n2), c(model$resSigmaSq1, model$resSigmaSq2))
  model$BIC <- BIC_SETAR(c(p, p), c(model$n1, model$n2), c(model$resSigmaSq1, model$resSigmaSq2))

  return(model)
}

```

and now we test the function for suitable parameters:

```
str( model <- EstimSETAR_postproc(model) )
```

```

## List of 15
## $ p      : num 2
## $ d      : num 2
## $ c      : num 0
## $ data   : num [1:170] 0.3 0.1 0.2 0.1 -0.1 ...
## $ n      : int 170
## $ PhiParams : num [1:6] 0 0.4124 0.4661 0.0692 0.8179 ...
## $ PhiStErrors: num [1:6] 0 0.0534 0.0694 0.0155 0.0448 ...

```

```
## $ residuals : num [1:168, 1] 0.1002 -0.1157 -0.2168 -0.0703 -0.0121 ...
## $ resSigmaSq : num 0.028
## $ n1 : num 102
## $ n2 : num 68
## $ resSigmaSq1: num 0.0213
## $ resSigmaSq2: num 0.039
## $ AIC : num -597
## $ BIC : num -578
```

It should be noted that for some values of  $p$  and  $d$  the indices of arrays in the algorithms might get out of the range of regularity for the linear system. For that reason we implement exceptions for the outputs of `EstimSETAR` in the following algorithm.

## 2.3: SETAR Parameter Estimation Procedure

To answer the question: ‘how does one find the right parameters  $p$ ,  $d$  and  $c$  for their desired SETAR model?’, we implement the following procedure:

```
pmax <- 7 # set maximum order p
# limit the c parameter by the 7.5-th and 92.5 percentile
cmin <- as.numeric(quantile(xt, 0.075)); cmax <- as.numeric(quantile(xt, 0.925));
h = (cmax - cmin) / 100 # determine the step by which c should be iterated
models <- list()
modelColumns <- list()
for (p in 1:pmax) {
  for (d in 1:p) {
    pdModels <- list()
    for (c in seq(cmin, cmax, h)) {
      tmp <- EstimSETAR(xt, p, d, c) # try to run the function
      # then test whether it returns `NA` as a result
      if (!as.logical(sum(is.na(tmp)))) {
        pdModels[[length(pdModels) + 1]] <- tmp
      }
    }
    sigmas <- as.numeric(lapply(pdModels, function(m) m$resSigmaSq))
    orders <- order(sigmas)
    # only the model whose parameter c gives the lowest residual square sum is chosen for postprocessing
    min_sigma_model <- EstimSETAR_postproc(pdModels[[ orders[1] ]])
    models[[length(models) + 1]] <- min_sigma_model
    modelColumns[[length(modelColumns) + 1]] <- c(
      p, d, min_sigma_model$c,
      min_sigma_model$n1, min_sigma_model$n2,
      min_sigma_model$AIC, min_sigma_model$BIC,
      min_sigma_model$resSigmaSq)
  }
}
```

```
##      p d      c  n1  n2      AIC      BIC resSigmaSq
## 1 1 1 0.000325 102 68 -595.3267 -585.6377 0.02955918
## 2 2 1 0.205425 142 28 -596.8388 -583.9747 0.02908718
## 3 2 2 0.102875 129 41 -616.5615 -602.8413 0.02624694
## 4 3 1 0.300650 153 17 -601.5380 -586.0834 0.02777303
## 5 3 2 0.102875 129 41 -606.5288 -588.2352 0.02668915
## 6 3 3 -0.197450 34 136 -592.0333 -574.2772 0.02778604
## 7 4 1 0.300650 153 17 -595.8224 -576.5041 0.02800968
## 8 4 2 0.102875 129 41 -600.5090 -577.6421 0.02681811
## 9 4 3 0.000325 102 68 -587.9837 -563.7613 0.02750077
## 10 4 4 0.205425 142 28 -597.1284 -575.6882 0.02698219
## 11 5 1 0.205425 142 28 -589.0005 -563.2723 0.02758865
```



```
## 12 5 2 0.102875 129 41 -598.8347 -571.3944 0.02590032
```

Now we have a set of models in their original order. To find the best suitable model, we choose 12 models with the lowest BIC (Bayesian Information Criterion):

```
##      p d      c n1 n2      AIC      BIC resSigmaSq
## 3  2 2 0.102875 129 41 -616.5615 -602.8413 0.02624694
## 5  3 2 0.102875 129 41 -606.5288 -588.2352 0.02668915
## 4  3 1 0.300650 153 17 -601.5380 -586.0834 0.02777303
## 1  1 1 0.000325 102 68 -595.3267 -585.6377 0.02955918
## 2  2 1 0.205425 142 28 -596.8388 -583.9747 0.02908718
## 8  4 2 0.102875 129 41 -600.5090 -577.6421 0.02681811
## 7  4 1 0.300650 153 17 -595.8224 -576.5041 0.02800968
## 10 4 4 0.205425 142 28 -597.1284 -575.6882 0.02698219
## 6  3 3 -0.197450 34 136 -592.0333 -574.2772 0.02778604
## 15 5 5 0.102875 129 41 -600.5783 -573.1380 0.02542239
## 14 5 4 0.205425 142 28 -597.2758 -571.5476 0.02618782
## 12 5 2 0.102875 129 41 -598.8347 -571.3944 0.02590032
```

and we can also include errors of the estimated regression coefficients:

```
## $`2/2/0.1029`
##      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]
## Phi      0 0.43170625 0.42834406 0.09138806 1.08026643 -0.42454326
## stdError  0 0.04164731 0.05038936 0.02581381 0.05901015 0.07802637
##
## $`3/2/0.1029`
##      [,1]      [,2]      [,3] [,4]      [,5]      [,6]      [,7] [,8]
## Phi      0 0.41811454 0.40547459 0 0.08019403 1.07869738 -0.31056937 0
## stdError  0 0.04340129 0.05441896 0 0.02619414 0.05963855 0.09969198 0
##
## $`3/1/0.3007`
##      [,1]      [,2]      [,3] [,4]      [,5]      [,6]      [,7]
## Phi      0 0.53829523 0.20456086 0 -0.1620930 1.1391822 0.6438414
## stdError  0 0.04368799 0.04233384 0 0.0792829 0.1808141 0.1619481
##
##      [,8]
## Phi     -0.9934059
## stdError 0.1556738
##
## $`1/1/3e-04`
##      [,1]      [,2]      [,3]      [,4]
## Phi      0.02224840 0.78551645 -0.06795137 0.93659667
## stdError 0.01108097 0.05520027 0.01518854 0.04389936
##
## $`2/1/0.2054`
##      [,1]      [,2]      [,3] [,4]      [,5] [,6]
## Phi      0 0.50086623 0.22030155 0 0.8652436 0
## stdError  0 0.04539197 0.03816144 0 0.1019782 0
##
## $`4/2/0.1029`
##      [,1]      [,2]      [,3] [,4] [,5]      [,6]      [,7]      [,8]
## Phi      0 0.4229181 0.4253314 0 0 0.08769237 1.06176402 -0.31828032
## stdError  0 0.0434935 0.0555784 0 0 0.02700914 0.06146456 0.09992171
##
##      [,9] [,10]
## Phi      0 0
## stdError  0 0
##
## $`4/1/0.3007`
##      [,1]      [,2]      [,3] [,4] [,5] [,6]      [,7]      [,8] [,9]
## Phi      0 0.54084131 0.20899552 0 0 0 0.9452945 0.7317165 0
## stdError  0 0.04390061 0.04344423 0 0 0 0.1823845 0.1623094 0
```

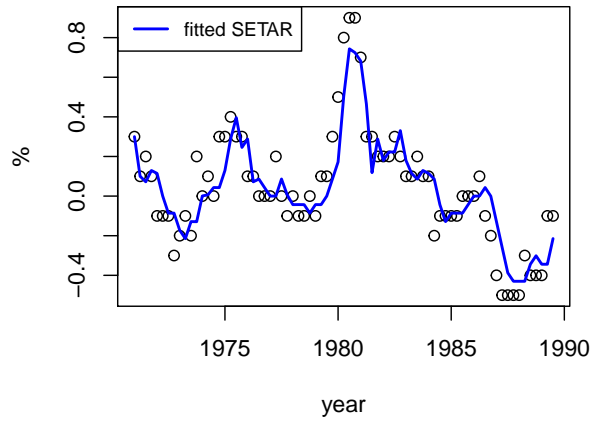
```

##          [,10]
## Phi      -1.2916701
## stdError  0.1591013
##
## $`4/4/0.2054`
##          [,1]      [,2]      [,3] [,4]      [,5]      [,6]      [,7]
## Phi        0 0.65727359 0.38978649 0 -0.16773802 0.14173560 0.6582130
## stdError    0 0.03774249 0.04650459 0 0.04752073 0.04020795 0.0967757
##          [,8]      [,9]      [,10]
## Phi      -0.3314201 0.4747251 -0.4104174
## stdError  0.1119202 0.1083704 0.1026022
##
## $`3/3/-0.1974`
##          [,1] [,2]      [,3]      [,4] [,5]      [,6]      [,7]      [,8]
## Phi        0 0 0.46444897 0.3403872 0 0.75043610 0.12079526 -0.1255942
## stdError    0 0 0.08494846 0.1407401 0 0.03879669 0.04629148 0.0410894
##
## $`5/5/0.1029`
##          [,1]      [,2]      [,3] [,4] [,5]      [,6] [,7]      [,8]
## Phi      0.016299311 0.64365718 0.38921338 0 0 -0.1625067 0 0.47118620
## stdError 0.008034404 0.04103325 0.05004459 0 0 0.0539052 0 0.07585229
##          [,9]      [,10]      [,11] [,12]
## Phi      -0.19079304 0.53307667 -0.3303657 0
## stdError  0.08097692 0.09450724 0.1058457 0
##
## $`5/4/0.2054`
##          [,1]      [,2]      [,3] [,4] [,5]      [,6]      [,7]      [,8]
## Phi        0 0.63242444 0.40442194 0 0 -0.20691700 0.16324920 0.63510151
## stdError    0 0.03804026 0.04658868 0 0 0.03955612 0.04111555 0.09720526
##          [,9]      [,10]      [,11] [,12]
## Phi      -0.2803961 0.4214599 -0.4884326 0
## stdError  0.1138933 0.1103896 0.1218861 0
##
## $`5/2/0.1029`
##          [,1]      [,2]      [,3]      [,4] [,5]      [,6]
## Phi      0.018185178 0.42527088 0.43763948 0.11801416 0 -0.14016981
## stdError 0.008024606 0.04349921 0.05570058 0.04784635 0 0.04199251
##          [,7]      [,8]      [,9] [,10] [,11] [,12]
## Phi      0.11594786 1.03987845 -0.3667765 0 0 0
## stdError 0.02955887 0.06321611 0.1031897 0 0 0

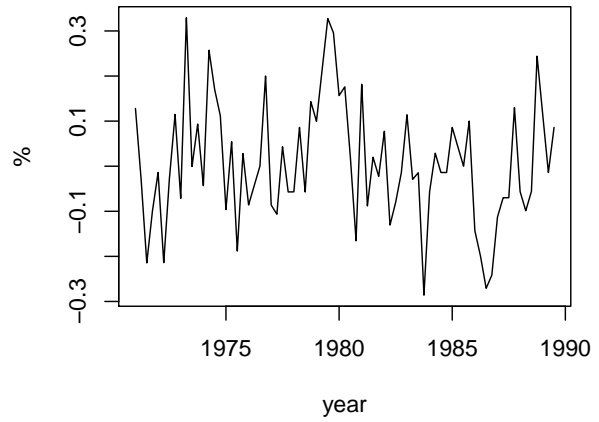
```

We can now visualize the results of the top 3 models:

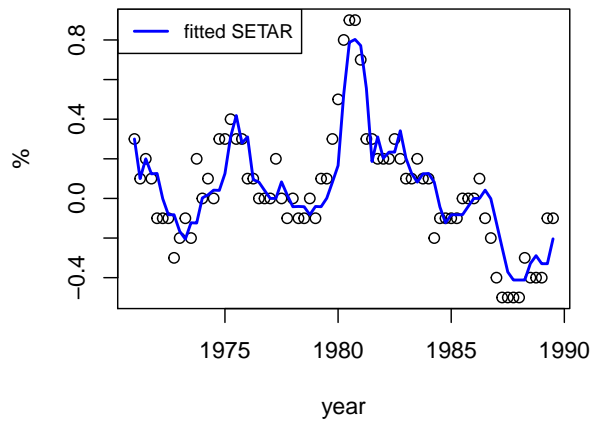
**SETAR( 2 , 2 , 0.1029 )**



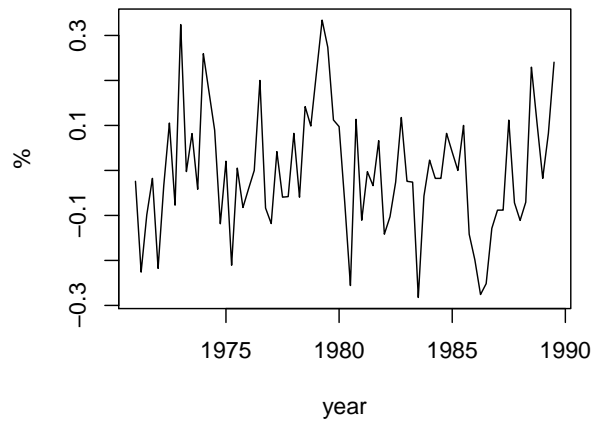
**SETAR( 2 , 2 , 0.1029 ) Residuals**



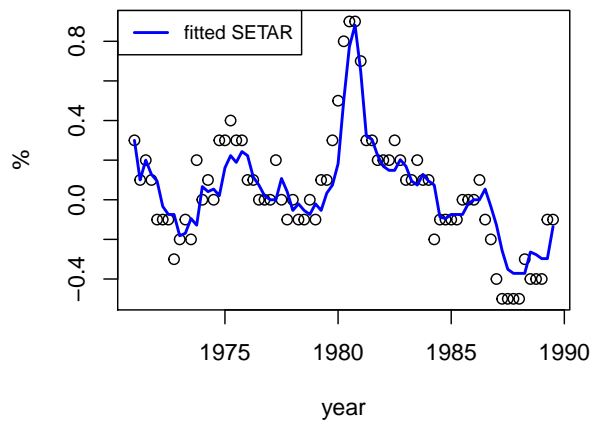
**SETAR( 3 , 2 , 0.1029 )**



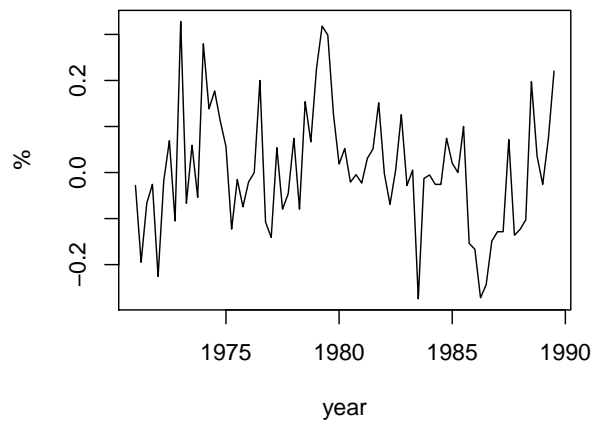
**SETAR( 3 , 2 , 0.1029 ) Residuals**



**SETAR( 3 , 1 , 0.3007 )**



**SETAR( 3 , 1 , 0.3007 ) Residuals**

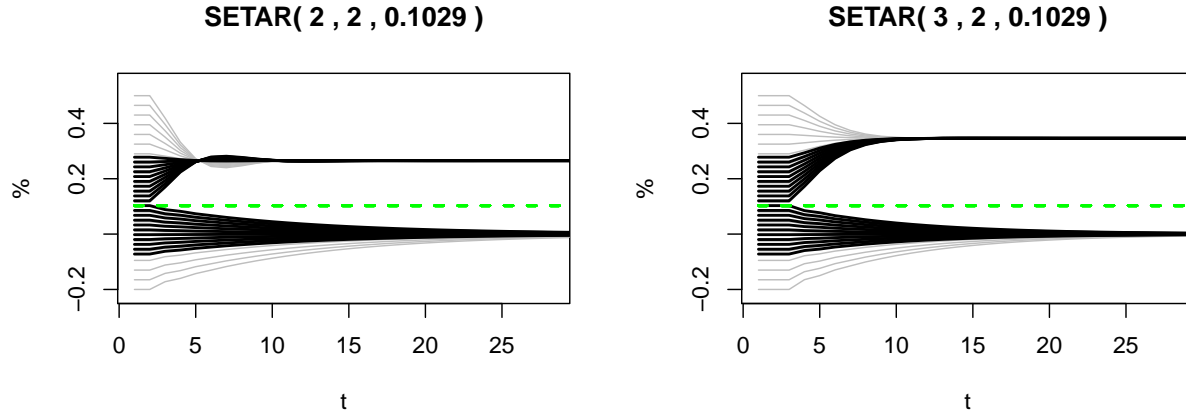


The results suggest that the best SETAR models have a threshold  $c$  quite close to zero and more-or-less the same

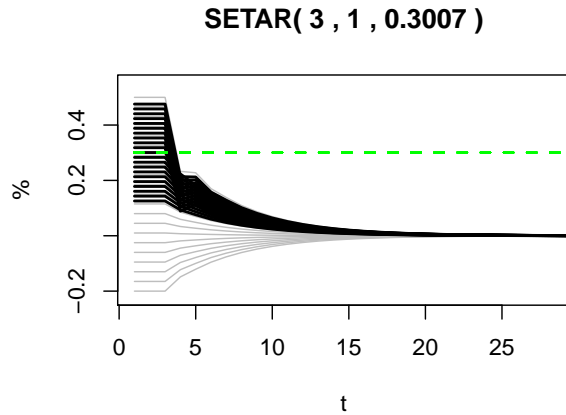
RSS. The very first with a lower BIC (Bayesian Information Criterion), has slightly larger RSS than the models that come after it. To verify the correctness of our procedure we will need to compare it with inbuilt functions from a verified library.

## 2.4: SETAR Equilibria and Equilibrium Simulations

It is also essential to find out whether the skeletons of the selected SETAR models have some equilibria. The estimation of the exact equilibria of the piecewise-linear skeletons with  $p = 1$  is straightforward: We find the fixed points of the skeletons by finding the intersections between their graphs and the identity line  $id x = x$ , given the model parameters (coefficients). However, the results of our search have mostly higher AR degrees, thus we will need to determine the models' equilibria using a more general method, namely letting the model skeletons evolve with multiple input initial conditions.

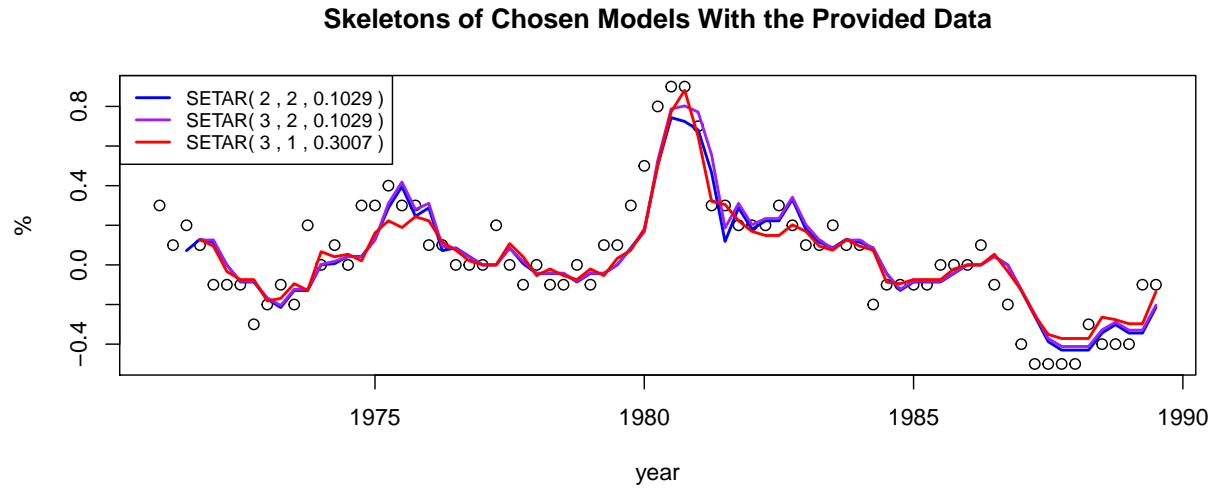


```
## $`SETAR( 3 , 1 , 0.301 ) equilibria`
## [1] 0.0000 0.2654
##
## $`SETAR( 3 , 1 , 0.301 ) equilibria`
## [1] 0.0000 0.3459
##
## $`SETAR( 3 , 1 , 0.301 ) equilibria`
## [1] 0
```

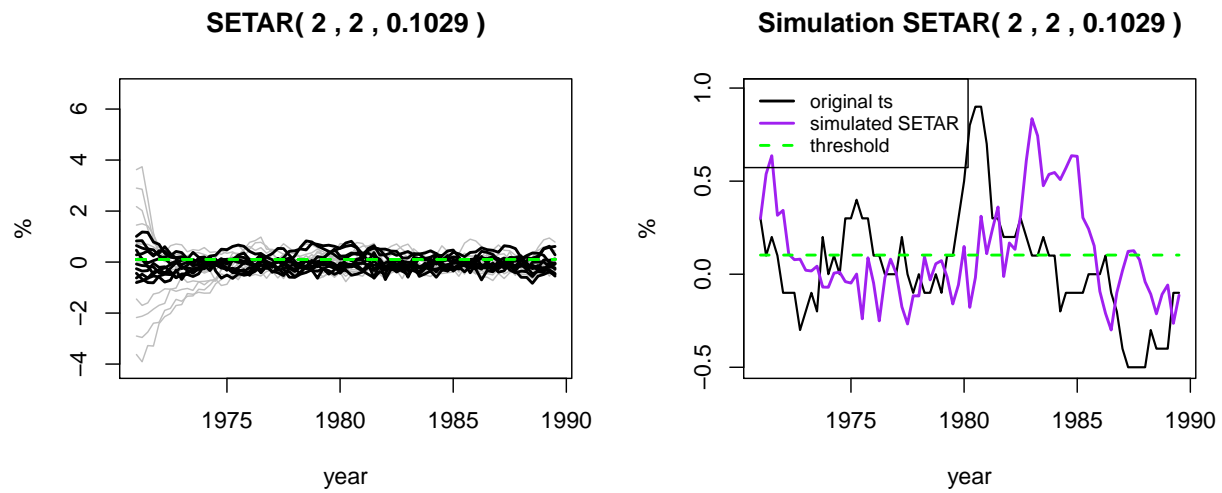


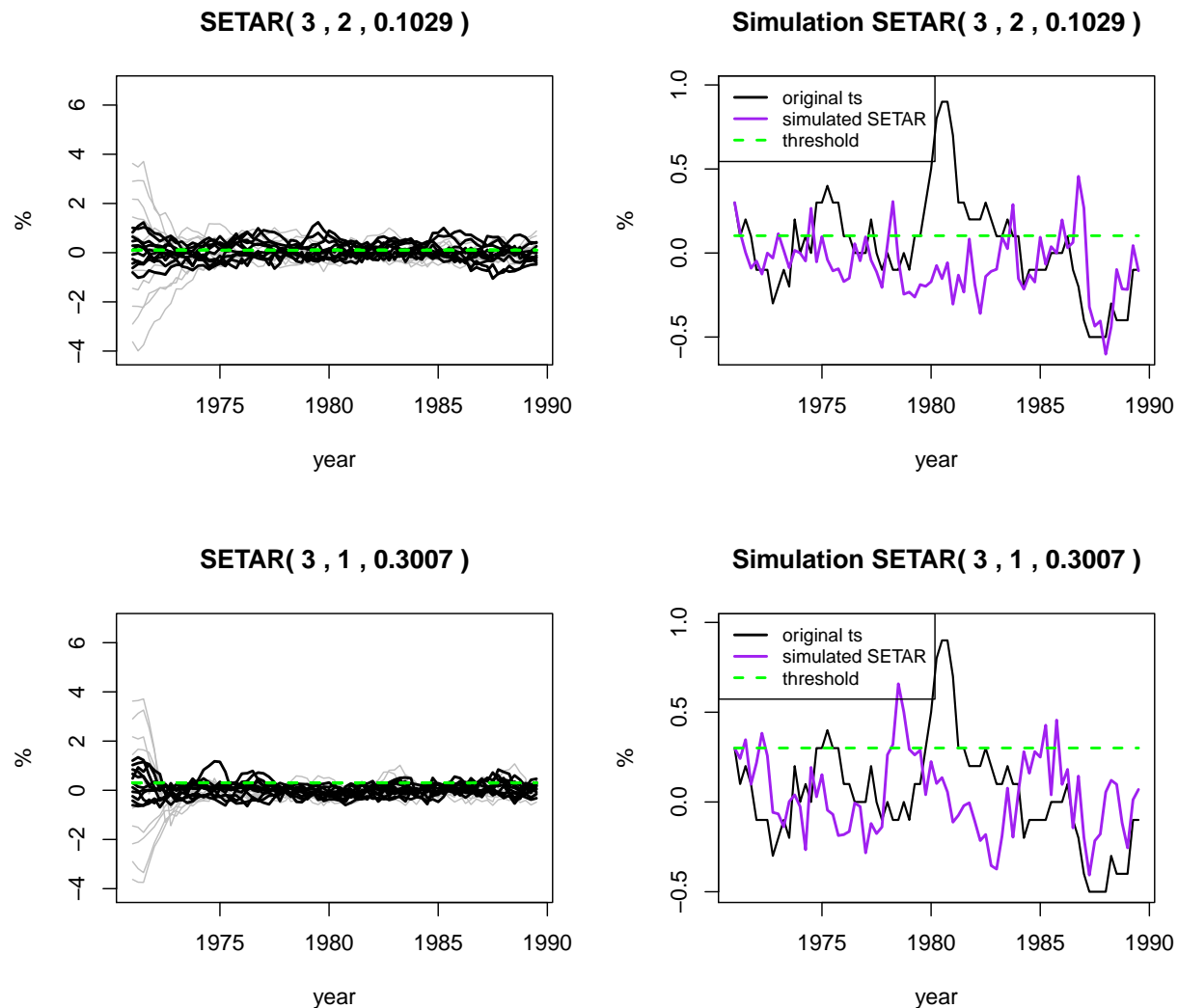
As we see, the trajectories of the top 3 models gravitate towards 0 in all models, but in the first and second model they can end up in one more position, close to zero. It might also be interesting to see how the trajectories evolve

when we add an iid noise on top of the model skeleton. First we observe the skeleton behavior in our data:



And then we carry out multiple simulations with initial conditions close to the threshold. The added noise will have the same deviance as the residual square sum.





The trajectories of all of the first three models seem to gravitate toward 0 significantly fast (or alternatively: towards their threshold values which are close to zero as well). The relatively low oscillation rate of the original time series suggests that the differences of this time series will, at most, fluctuate around 0. The change between the 'high' and 'low' regimes does not seem very significant, at least on the larger scale. The validity of the model will be tested in chapter 3.

## 2.5: Comparison Of the Results With Inbuilt Functions

To verify the correctness of our methods we proceed to construct the top 3 SETAR models by plugging their parameters into inbuilt functions:

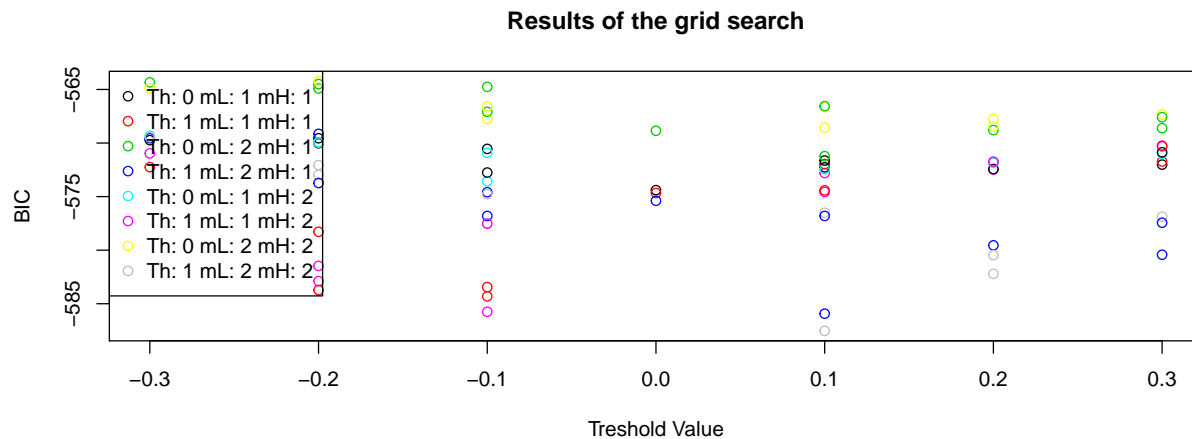
```
suppressMessages(pkgTest("tsDyn"))

#Testing a function which selects an orders automatically:
mmax <- 2

par(mfrow=c(1,1))
( result1 <- selectSETAR(xt, m=mmax, thDelay=0:(mmax-1), criterion="BIC", same.lags=T, trim=0.1) )

## Using maximum autoregressive order for low regime: mL = 2
```

```
## Using maximum autoregressive order for high regime: mH = 2
## Searching on 21 possible threshold values within regimes with sufficient ( 10% ) number of observations
## Searching on 84 combinations of thresholds ( 21 ), thDelay ( 2 ) and m ( 2 )
```



```
## Results of the grid search for 1 threshold
```

```
##   thDelay m   th      BIC
## 1      1 2  0.1 -587.5300
## 2      1 2  0.1 -585.9265
## 3      1 1 -0.1 -585.7518
## 4      1 1 -0.1 -584.3112
## 5      1 1 -0.2 -583.7230
## 6      1 1 -0.1 -583.4471
## 7      1 1 -0.2 -582.8902
## 8      1 2  0.2 -582.2097
## 9      1 1 -0.2 -581.4681
## 10     1 2  0.2 -580.4653
```

the estimated thDelay corresponds to d-1.

```
## List of 15
```

```
## $ p      : num 2
## $ d      : num 1
## $ c      : num -0.1
## $ data    : num [1:170] 0.3 0.1 0.2 0.1 -0.1 ...
## $ n      : int 170
## $ PhiParams : num [1:6] 0.0711 0.7998 0.1755 0 0.6911 ...
## $ PhiStErrors: num [1:6] 0.023 0.0946 0.0666 0 0.0444 ...
## $ residuals : num [1:168, 1] 0.0838 -0.0539 -0.2005 -0.0466 -0.0736 ...
## $ resSigmaSq : num 0.0293
## $ n1       : num 50
## $ n2       : num 120
## $ resSigmaSq1: num 0.0357
## $ resSigmaSq2: num 0.0272
## $ AIC      : num -583
## $ BIC      : num -564
```

Note that we set `thDelay=0:(mmax-1)` instead of `1:mmax`. `selectSETAR` uses `thDelay = 0` for step `d=1` delay correspondence:  $x_{t-d} < c$  or  $x_{t-d} > c$ . the resulting BIC's are different, possibly due to the package using a different formula

```
## Results of the grid search for 1 threshold
```

```
##   thDelay m   th      BIC
## 1      1 2  0.1 -590.7913
```

```
## 2      1 2  0.1 -589.0388
## 3      1 1 -0.1 -588.2120
## 4      1 1 -0.1 -586.5338
## 5      1 1 -0.2 -585.9216
## 6      1 1 -0.1 -585.7998
## 7      1 2  0.2 -585.1638
## 8      1 1 -0.2 -585.0533
## 9      4 2  0.1 -585.0007
## 10     4 2  0.1 -584.5197
```

Setting higher `mmax`, the function returns a list of models similar to the one given by our procedure in section 2.3. We can also compare the accuracy of the computation of the regression coefficients in our `EstimSETAR` method, with for example: `setar()` function (from `tsDyn` library as well):

```
setars <- list()
coeffComparison <- list()
resSigmaComparison <- list()
n <- length(xt)
for (i in 1:3) {
  p <- models[[ orders[i] ]]$p
  d <- models[[ orders[i] ]]$d
  c <- models[[ orders[i] ]]$c
  setars[[i]] <- setar(xt, m=p)
  k <- max(p, d)
  resSigmaComparison[[i]] <- c(
    (1 / (n - k) * sum(setars[[i]]$residuals ^ 2)),
    models[[ orders[i] ]]$resSigmaSq
  )
  inbuiltParams <- t(setars[[i]]$coefficients)
  key <- paste(p, d, round(c, digits=4), sep=" / ")
  coeffComparison[[key]] <- rbind(
    inbuiltParams,
    t(append(models[[orders[i]]]$PhiParams, models[[orders[i]]]$c))
  )
  row.names(coeffComparison[[key]]) <- t(c("inbuilt", "custom"))
}
```

```
##
## 1 T: Trim not respected:  0.8511905 0.1488095 from th: 0.3
## 1 T: Trim not respected:  0.8502994 0.1497006 from th: 0.3
## 1 T: Trim not respected:  0.8502994 0.1497006 from th: 0.3
## 1 T: Trim not respected:  0.8502994 0.1497006 from th: 0.3
## 1 T: Trim not respected:  0.8502994 0.1497006 from th: 0.3
```

```
coeffComparison
```

```
## $`2 / 2 / 0.1029`
##      const.L   phiL.1   phiL.2   const.H   phiH.1   phiH.2
## inbuilt -0.01014873 0.4562865 0.2714053 -0.03164651 0.9628571 -0.1110418
## custom  0.00000000 0.4317063 0.4283441  0.09138806 1.0802664 -0.4245433
##
##      th
## inbuilt 0.100000
## custom  0.102875
##
## $`3 / 2 / 0.1029`
##      const.L   phiL.1   phiL.2   phiL.3   const.H   phiH.1
## inbuilt -0.009656784 0.4908912 0.1958345 0.03467407 0.09657510 0.5947604
## custom  0.000000000 0.4181145 0.4054746 0.00000000 0.08019403 1.0786974
##
##      phiH.2   phiH.3   th
## inbuilt  0.5578468 -0.6112305 0.300000
## custom   -0.3105694  0.0000000 0.102875
```



```
##
## $`3 / 1 / 0.3007`
##          const.L   phiL.1   phiL.2   phiL.3   const.H   phiH.1
## inbuilt -0.009656784 0.4908912 0.1958345 0.03467407 0.0965751 0.5947604
## custom  0.000000000 0.5382952 0.2045609 0.000000000 -0.1620930 1.1391822
##          phiH.2   phiH.3   th
## inbuilt 0.5578468 -0.6112305 0.30000
## custom  0.6438414 -0.9934059 0.30065

# comparing RSS

resSigmaComparison <- data.frame(matrix(unlist(resSigmaComparison), nrow=3, byrow=T))
colnames(resSigmaComparison) <- c("inbuilt", "custom")
row.names(resSigmaComparison) <- t(paste("rss",1:3))
resSigmaComparison

##          inbuilt      custom
## rss 1 0.02844427 0.02624694
## rss 2 0.02765805 0.02668915
## rss 3 0.02765805 0.02777303
```

Without specifying the threshold value, the inbuilt `setar` function finds threshold values quite close to those of our custom procedures. The AR order `p` for both regimes, however, has to be specified in advance. The comparison of the model coefficients suggests that our custom method was more-or-less accurate.

## 2.6: Conclusion

The results of the SETAR Parameter Estimation Procedure in section 2.3 show that the 3 best 2-regime SETAR models are:

```
##          unlist.results.
## 1 SETAR( 2 , 2 , 0.1029 )
## 2 SETAR( 3 , 2 , 0.1029 )
## 3 SETAR( 3 , 1 , 0.3007 )
```

The first model with the lowest BIC (Bayesian Information Criterion) has the most accurate estimation of its 4 regression parameters, with the highest residual square sum. The first model seems to have a stable equilibrium at their threshold values.

## 3: Tests of Linearity/Nonlinearity of SETAR models

We need to make sure a non-linear model (SETAR, for example) is really suitable for describing the process. In order to find out, we test the null hypothesis that a linear model is more suitable than a non-linear one. In the case of a 2-regime model we are looking for, so called, nuisance parameters, i.e.:  $H_0 : \Phi_1 = \Phi_2$  where  $\Phi_1$  and  $\Phi_2$  are the parameters of the low and the high regime respectively.

### 3.1: Hansen's Conditions

Hansen proposed three conditions to test whether a SETAR model can be tested for linearity using the so called Likelihood-Ratio (LR) test:

```
Hansen <- function(d, c, Phi) {
  p <- (length(Phi)/2) - 1
  #separate regimes into rows
  Phi <- do.call(rbind, split(Phi,rep(1:2,each=(p + 1))))
  # (p10-p20)+(p1d-p2d)*c <= 0
  c1 <- !isTRUE(all.equal( 0, apply(Phi[,c(1,1 + d),drop=F],2, diff) %*% c(1,c) ))
```

```

# p1j neq p2j, j notin {0,d}
c2 <- all(apply(Phi[, -c(1,1 + d), drop=F], 2, function(x) !identical(0, diff(x))))
# sum_j|pij| < 1 forall i=1,2
c3 <- all(apply(Phi[, -1, drop=F], 1, function(x) sum(abs(x))) < 1)
c(cond1=c3, cond2=c2, cond3=c3)
}

```

If all three are satisfied the model can be tested using the LR test:

```

##          cond1 cond2 cond3
## 2 / 2 / 0.103 FALSE TRUE FALSE
## 3 / 2 / 0.103 FALSE TRUE FALSE
## 3 / 1 / 0.301 FALSE TRUE FALSE
## 1 / 1 / 0      TRUE TRUE TRUE
## 2 / 1 / 0.205 TRUE TRUE TRUE
## 4 / 2 / 0.103 FALSE TRUE FALSE
## 4 / 1 / 0.301 FALSE TRUE FALSE
## 4 / 4 / 0.205 FALSE TRUE FALSE
## 3 / 3 / -0.197 TRUE TRUE TRUE
## 5 / 5 / 0.103 FALSE TRUE FALSE
## 5 / 4 / 0.205 FALSE TRUE FALSE
## 5 / 2 / 0.103 FALSE TRUE FALSE

```

It appears that only the first and the fifth model can be tested using the LR test. The rest will have to be assessed using the Lagrange Multiplier (LM) test.

### 3.2: LR and LM Tests

In this section we formulate the basic procedures for the LR (Likelihood Ratio), and LM (Lagrange Multiplier) tests:

```

LRtest <- function(x, p, var, alpha=0.05) {
  tmp <- ar(x, aic=F, order.max=pmax, method = "ols")
  tmp <- tmp$var.pred # linear model residual variance
  teststat <- length(x)*(tmp-var)/tmp # test statistic
  CDF <- Vectorize( function(t) { # test statistic CDF
    fun <- function(t) 1 + sqrt(t/(2*pi))*exp(-t/8) + 1.5*exp(t)*pnorm(-1.5*sqrt(t)) -
      (t+5)*pnorm(-sqrt(t)/2)/2
    if(abs(t)>300 || is.infinite(t)) return(sign(t))
    if(t >= 0 ) fun(t) else 1-fun(-t)
  })
  # for alpha=2.5%: CV=11.03329250
  if(alpha==0.05) critval <- 7.68727553
  else critval <- uniroot(function(x) CDF(x) - (1-alpha), c(-1000,1000))$root
  # (test statistics, critical value, p-value)
  c(TS=teststat, CV=critval, p_value=1-CDF(teststat))
}

```

```

LRtest(xt, models[[ orders[1] ]]$p, models[[ orders[1] ]]$resSigmaSq)

```

```

##          TS          CV          p_value
## 17.152950411  7.687275530  0.007928351

```

```

suppressMessages(pkgTest("dynlm"))

```

```

LMtest <- function(x, p, d, alpha = 0.05) {
  # prevent from passing (accidental and needless) name to result
  names(p) <- NULL
  # if x is not a ts object, by chance
  x <- as.ts(x)
  # requires dynlm package (it can be implemented withou dynlm, see model2)

```

```

model1 <- dynlm(x ~ L(x,1:p))
y <- model1$residuals
# a list of shifted time series
tmp <- c(
  list(y),
  lapply(1:p, function(i) stats::lag(x, -i)),
  lapply(1:p, function(i) stats::lag(x, -i)*stats::lag(x,-d)),
  list(stats::lag(x,-d)^3)
)
tmp <- do.call(function(...) ts.intersect(..., dframe=T), tmp)
names(tmp) <- c("y", paste0("x",1:p), paste0("xd",1:p), "xd^3")
# cannot be done with the dynlm package
model2 <- lm(y ~ ., data = tmp)
z <- model2$residuals
teststat <- (length(x)-p) * (sum(y^2)/sum(z^2) - 1)
c(TS=teststat, CV=qchisq(1-alpha, df=p+1), p_value=1-pchisq(teststat, df=p+1))
}

LMtest(xt, models[[ orders[1] ]]$p, models[[ orders[1] ]]$d)

```

```

##          TS          CV      p_value
## 1.696816e+01 7.814728e+00 7.174773e-04

```

We can easily automate the testing procedure with the following results:

```

##      p d      c Hansen Cond.      TS      CV      p-value
## 3  2 2  0.102875      FALSE 16.968165  7.814728 0.0007174773
## 5  3 2  0.102875      FALSE 16.343839  9.487729 0.0025908411
## 4  3 1  0.300650      FALSE  8.351652  9.487729 0.0795136152
## 1  1 1  0.000325       TRUE -2.135613  7.687276 0.7970372758
## 2  2 1  0.205425       TRUE  0.613035  7.687276 0.3539835179
## 8  4 2  0.102875      FALSE 16.214712 11.070498 0.0062570609
## 7  4 1  0.300650      FALSE 12.687710 11.070498 0.0264878050
## 10 4 4  0.205425      FALSE 10.173680 11.070498 0.0704610459
## 6  3 3 -0.197450       TRUE  8.190160  7.687276 0.0448575191
## 15 5 5  0.102875      FALSE 15.615335 12.591587 0.0159745090
## 14 5 4  0.205425      FALSE 10.919148 12.591587 0.0909076284
## 12 5 2  0.102875      FALSE 17.915958 12.591587 0.0064456881

```

For significance level  $\alpha = 0.05$  the linearity hypothesis is not rejected only for the following models:

```

##          unlist.res.
## 1 SETAR( 2 , 2 , 0.1028750000000001 )
## 2 SETAR( 3 , 2 , 0.1028750000000001 )
## 3 SETAR( 4 , 2 , 0.1028750000000001 )

```

The remaining models can be considered non-linear.

### 3.3 Modified LR Test Via Bootstrapping

The proposed LR test has a significant drawback in the fact that it can only be done when Hansen's conditions are satisfied. This is due to the fact that we do not know the distribution of the resulting F-statistic. According to Hansen (1996), however, the distribution of a bootstrapped statistic  $F^*$  converges weakly in probability to the distribution of  $F$ , so that repeated bootstrap draws from  $F^*$  can be used to approximate the asymptotic distribution of  $F$ . A parallelized implementation can be seen in the following snippet:

```

...

if (FALSE %in% Hansen(p, d, model$PhiParams)) {
  suppressMessages(pkgTest("tsDyn"))
}

```

```

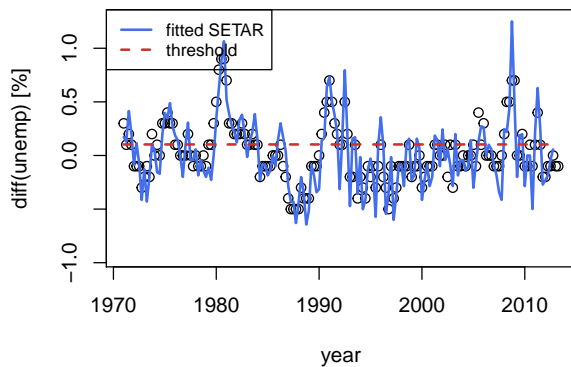
suppressMessages(pkgTest("parallel"))
suppressMessages(pkgTest("doSNOW")) # using doSNOW package for parallel computing
n_cores <- detectCores() - 1
cl <- makeCluster(n_cores, type="SOCK")
registerDoSNOW(cl)
log <- capture.output({
  testResults <- suppressWarnings(
    setarTest(x, m=p, thDelay=0:(d - 1), nboot=nboot ,trim=0.1, test="lvs", hpc="foreach")
  )
})
stopCluster(cl)
...
}
...

```

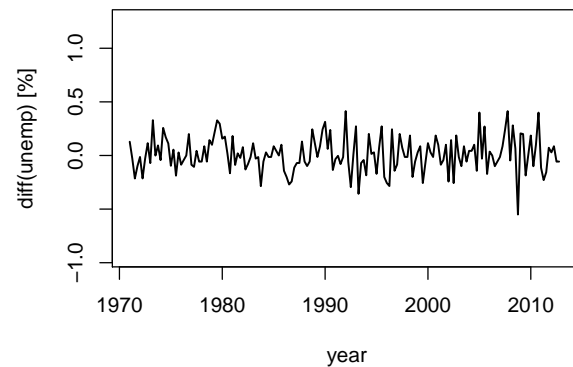
### 3.4 Visualisation of Non-Linear Models

From the results of the previous procedure, we will visualize the models for which the linearity null-hypothesis was rejected based on the LR and LM tests:

**SETAR( 2 , 2 , 0.1029 )**

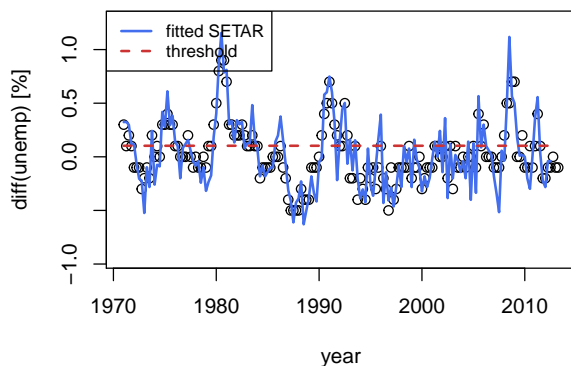


**SETAR( 2 , 2 , 0.1029 ) residuals**

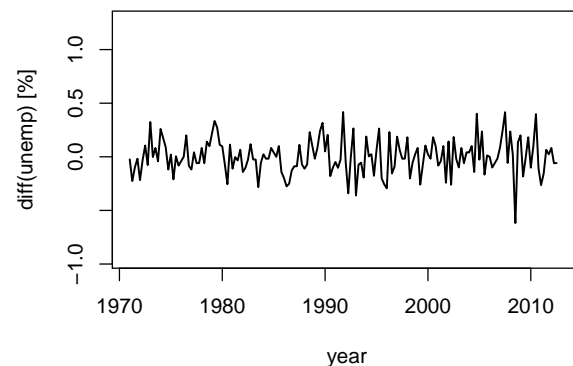


```
## resSigmaSq
## 0.02624694
```

**SETAR( 3 , 2 , 0.1029 )**

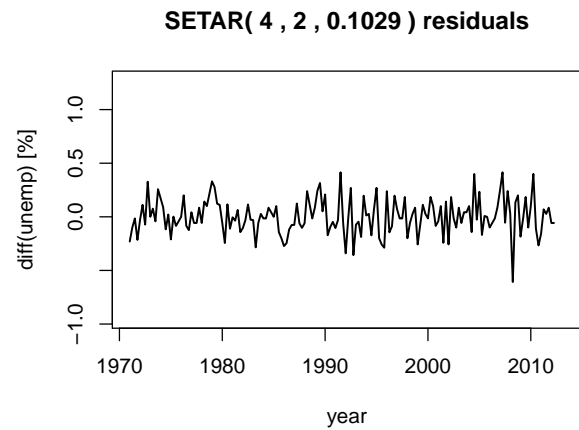
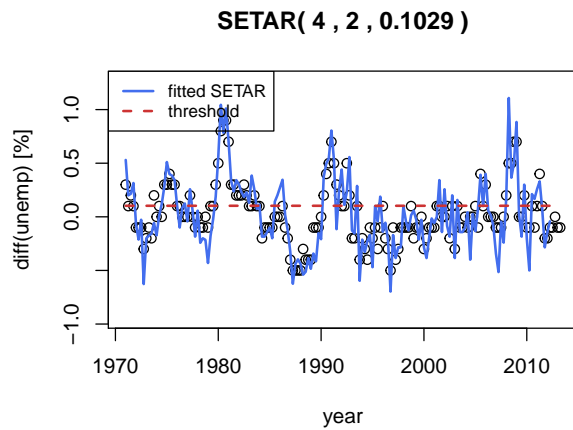


**SETAR( 3 , 2 , 0.1029 ) residuals**

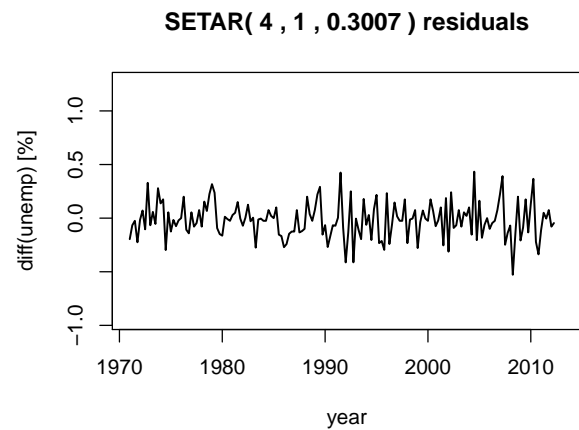
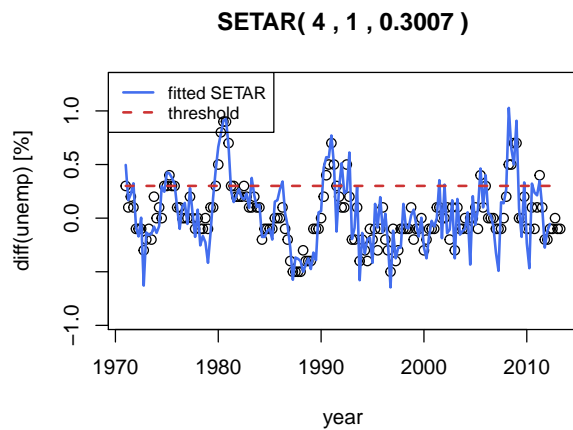


```
## resSigmaSq
```

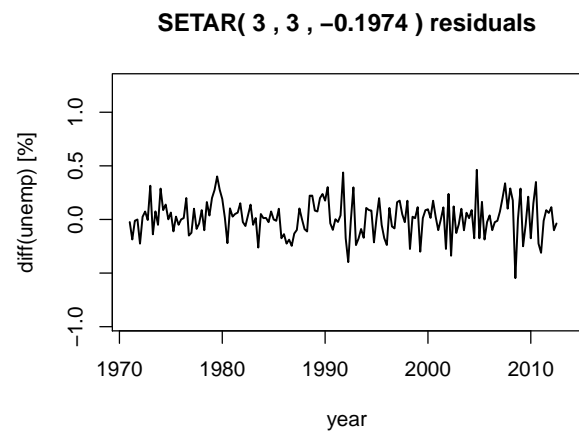
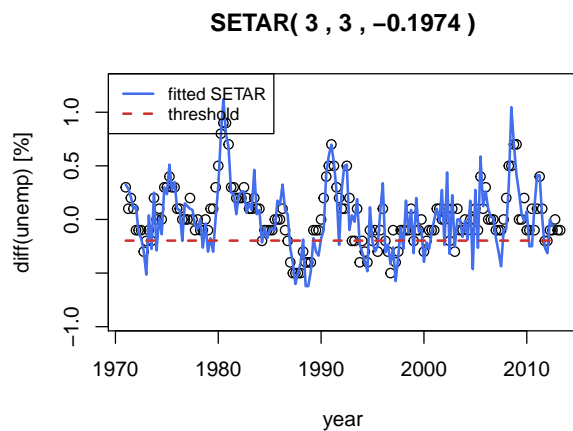
```
## 0.02668915
```



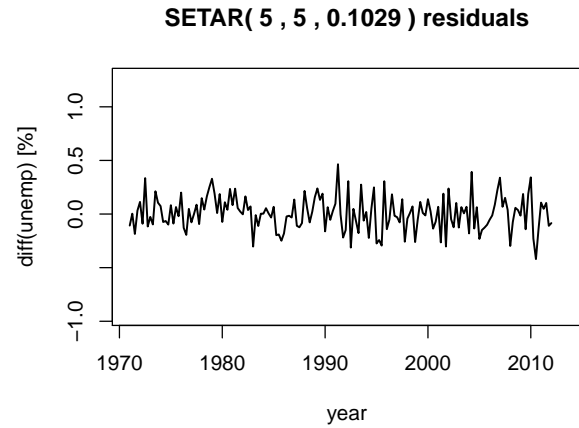
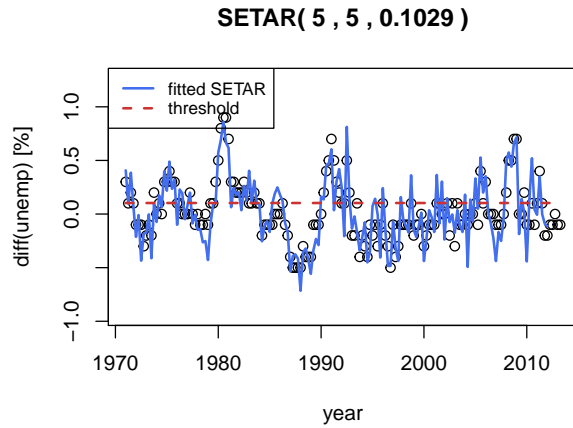
```
## resSigmaSq  
## 0.02681811
```



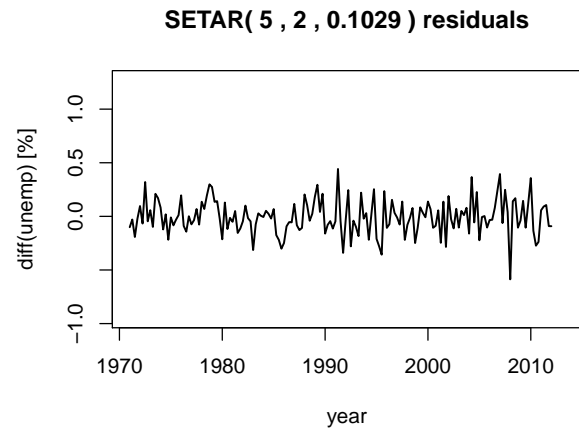
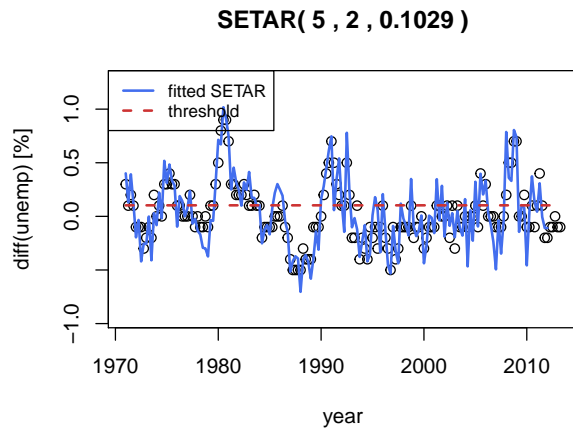
```
## resSigmaSq  
## 0.02800968
```



```
## resSigmaSq
## 0.02778604
```



```
## resSigmaSq
## 0.02542239
```



```
## resSigmaSq
## 0.02590032
```

### 3.5 Conclusion

Since the differences in the unemployment rate have been used, we show the threshold value as well as the fitted values of the models in the same plot. The switching between the high and the low regimes can be clearly visible for all the selected models (perhaps, except the second one, with its threshold value quite close to zero). It is not yet clear whether another regime should be present in the stochastic process. This will be assessed in the following chapter.

## 4. 3-Regime SETARs and Diagnostic Tests of SETAR Models

The next step in the analysis using SETAR models is verifying whether 2 regimes suffice. If they do not, we will have to consider the possibility that a third regime needs to be added. In that case, we need to write methods for such model.

## 4.1 Useful Functions

```
# the indicator function for 3 regimes:
Indicator3 <- function(x, c) {
  tmp <- rep(F,3)
  tmp[findInterval(x, c, left.open = T) + 1] <- T
  tmp
}

Indicator3(-4, c(-1,1))

## [1] TRUE FALSE FALSE
Indicator3(0, c(-1,1))

## [1] FALSE TRUE FALSE
Indicator3(4, c(-1,1))

## [1] FALSE FALSE TRUE

# SETAR3 basis vector
Yt3 <- function(x, t, p) c(1, x[(t - 1):(t - p)])

# SETAR3 skeleton
Xt3 <- function(x, t, p, d, c, z = x) {
  # z is the threshold variable
  I <- Indicator3(z[t - d], c)
  Y <- Yt(x, t, p)
  c(I[1] * Y, I[2] * Y, I[3] * Y)
}

# covariance matrix of the 3 regime SETAR
CovMat3 <- function(x, p, d, c) {
  n <- length(x)
  # this will become the covariance matrix
  Yc <- matrix(0., ncol = (3 * p + 3), nrow = (3 * p + 3))
  k <- max(p, d)
  for (t in (k + 1):n) {
    XT <- Xt3(x, t, p, d, c)
    Yc <- Yc + (XT %o% XT)
  }
  det <- det(Yc)
  if (det > -0.00001 && det < 0.00001) {
    return(NA)
  } else {
    return(inv(Yc))
  }
}

CovMat3(xt, p=2, d=1, c=c(-0.1, 0.2))

##           [,1]      [,2]      [,3]      [,4]      [,5]      [,6]
## [1,] 0.09009639 0.2885639 -0.00608143 0.00000000 0.00000000 0.00000000
## [2,] 0.28856387 1.5217885 -0.52175095 0.00000000 0.00000000 0.00000000
## [3,] -0.00608143 -0.5217510 0.75481283 0.00000000 0.00000000 0.00000000
## [4,] 0.00000000 0.0000000 0.00000000 0.01209957 -0.01656548 0.00025222
## [5,] 0.00000000 0.0000000 0.00000000 -0.01656548 1.66822853 -0.3459105
## [6,] 0.00000000 0.0000000 0.00000000 0.00025220 -0.34591053 0.4059073
## [7,] 0.00000000 0.0000000 0.00000000 0.00000000 0.00000000 0.00000000
## [8,] 0.00000000 0.0000000 0.00000000 0.00000000 0.00000000 0.00000000
```

```
## [9,] 0.00000000 0.0000000 0.00000000 0.00000000 0.00000000 0.00000000
##      [,7]      [,8]      [,9]
## [1,] 0.00000000 0.0000000 0.00000000
## [2,] 0.00000000 0.0000000 0.00000000
## [3,] 0.00000000 0.0000000 0.00000000
## [4,] 0.00000000 0.0000000 0.00000000
## [5,] 0.00000000 0.0000000 0.00000000
## [6,] 0.00000000 0.0000000 0.00000000
## [7,] 0.15071243 -0.3415514 0.06411848
## [8,] -0.34155138 1.3811044 -0.73962398
## [9,] 0.06411848 -0.7396240 0.77066336
```

To find out, whether the third regime should be added, we need to test for the independence of residuals:

## 4.2 The Brock-Dechert-Scheinkman (BDS) Test

Regarded as the most successful tests for nonlinearity due to its universality, the BDS test relies on evaluating a correlation integral  $C(q, r)$  as a measure of repeated occurrence of patterns in the time series. It is the estimate of the probability of two arbitrary  $q$ -dimensional points in  $\mathbb{R}^q$  being no further than  $r\hat{\sigma}_\varepsilon$  apart ( $0.5 \leq r \leq 1.5$ ). If the data is generated by an iid process, the correlation integral should approach  $C(q, r) \rightarrow C(1, r)^q$ .

```
## possible remaining nonlinearity for with mean p-val = 0.875
## possible remaining nonlinearity for with mean p-val = 1
## possible remaining nonlinearity for with mean p-val = 0.875
## remaining nonlinearity rejected with mean p-val = 0.25
## remaining nonlinearity rejected with mean p-val = 0.25
## possible remaining nonlinearity for with mean p-val = 0.625
## possible remaining nonlinearity for with mean p-val = 0.75

## ==== Remaining nonlinearities detected for: =====
## NULL
```

After filtering out the models with remaining nonlinearity, and since we get 3 models with the same  $d$  parameter, it will be enough to build just one SETAR3.

## 4.3 SETAR3 Parameter Estimation

Similarly to section 2.3, we construct an estimation procedure with two distinct threshold parameters  $c_1$  and  $c_2$ . Our helper functions are ready from section 4.1. Using them we implement:

```
suppressMessages(pkgTest("zeallot"))
suppressMessages(pkgTest("matlib"))

EstimSETAR3 <- function(x, p, d, c) {
  resultModel <- list()
  resultModel$p = p; resultModel$d = d; resultModel$c = c;
  resultModel$data = x; n = length(x); resultModel$n = n;

  k <- max(p, d)

  X <- as.matrix(apply(as.matrix((k + 1):n), MARGIN=1, function(t) Xt3(x, t, p, d, c) ))
  y <- as.matrix(x[(k + 1):n])
  K <- CovMat3(x, p, d, c); b <- crossprod(t(X), y);

  if (as.logical(sum(is.na(K)))) {
    return(NA)
  } else {
    sol_phi <- as.numeric(t(K %*% b)); sol_se <- sqrt(diag(K)/n);
    eps <- 0.01;
  }
}
```



```

# filter out those coeffs that are of the same order of magnitude as their errors
filter <- sapply(1:(3*(p + 1)), function(i) ifelse(
  abs(sol_phi[i]) <= 2 * abs(sol_se[i]), 0, 1
))

sol_phi <- sol_phi * filter
sol_se <- sol_se * filter

solution <- cbind(phi = sol_phi, se = sol_se)
resultModel$PhiParams <- solution[,1] # solving (X'X)*phi = X'y
resultModel$PhiStErrors <- solution[,2] # standard errors
skel <- crossprod(X, resultModel$PhiParams); resultModel$skel <- skel;
resultModel$residuals <- (y - skel)
resultModel$resSigmaSq <- 1 / (n - k) * sum(resultModel$residuals ^ 2)

return(resultModel)
}
}

str( test_model <- EstimSETAR3(xt, p=2, d=1, c=c(-0.1, 0.2)) )

## List of 10
## $ p : num 2
## $ d : num 1
## $ c : num [1:2] -0.1 0.2
## $ data : num [1:170] 0.3 0.1 0.2 0.1 -0.1 ...
## $ n : int 170
## $ PhiParams : num [1:9] 0.0711 0.7998 0.1755 -0.0171 0.3404 ...
## $ PhiStErrors : num [1:9] 0.02302 0.09461 0.06663 0.00844 0.09906 ...
## $ skel : num [1:168, 1] 0.0908 0.1688 0.0662 -0.0265 -0.0264 ...
## $ residuals : num [1:168, 1] 0.1092 -0.0688 -0.1662 -0.0735 -0.0736 ...
## $ resSigmaSq : num 0.0285

After we find a suitable SETAR3 model, we implement a postprocessing method:
EstimSETAR3_postproc <- function(model) {
  x <- model$data; k <- max(model$p, model$d); p <- model$p; c <- model$c; n <- model$n;
  y <- as.matrix(x[(k + 1):n])
  skel <- model$skel; # model$skel <- NULL; #skel attribute no longer needed

  # regime counts
  nRegCounts <- rowSums( matrix(as.numeric((sapply(x, function(xi) Indicator3(xi, c)))), nrow=3) )
  model$nRegCounts <- nRegCounts
  # names(model$nRegCounts) <- c("n1", "n2", "n3")

  # regime sigma sq:
  regSigmaSq <- rowSums(
    matrix(
      as.numeric((sapply(y, function(xi) Indicator3(xi, c)))), nrow=3)
      %*% (y - skel)^2 ) / (nRegCounts - k)
  )
  model$regSigmaSq <- regSigmaSq
  # names(model$regSigmaSq) <- c("rss1", "rss2", "rss3")

  # count valid p orders
  pOrders <- rowSums(
    matrix(as.numeric((
      sapply(1:(3 * (p + 1)),
        function(i) ifelse( (i % (p + 1)) != 1 && model$PhiParams[i] != 0, 1, 0) )
    ))
  )

```

```

    )
  ), nrow=3, byrow=T)
)
model$pOrders <- pOrders
# names(model$pOrders) <- c("p1", "p2", "p3")

model$AIC <- AIC_SETAR(pOrders, model$nRegCounts, model$resSigmaSq)
model$BIC <- BIC_SETAR(pOrders, model$nRegCounts, model$resSigmaSq)

return(model)
}

str( test_model <- EstimSETAR3_postproc(test_model) )

```

```

## List of 15
## $ p      : num 2
## $ d      : num 1
## $ c      : num [1:2] -0.1 0.2
## $ data    : num [1:170] 0.3 0.1 0.2 0.1 -0.1 ...
## $ n      : int 170
## $ PhiParams : num [1:9] 0.0711 0.7998 0.1755 -0.0171 0.3404 ...
## $ PhiStErrors: num [1:9] 0.02302 0.09461 0.06663 0.00844 0.09906 ...
## $ skel     : num [1:168, 1] 0.0908 0.1688 0.0662 -0.0265 -0.0264 ...
## $ residuals : num [1:168, 1] 0.1092 -0.0688 -0.1662 -0.0735 -0.0736 ...
## $ resSigmaSq : num 0.0285
## $ nRegCounts : num [1:3] 50 84 36
## $ regSigmaSq : num [1:3] 0.0302 0.0178 0.0553
## $ pOrders    : num [1:3] 2 2 1
## $ AIC        : num -589
## $ BIC        : num -573

```

Since the number of regimes  $m$  appears as a parameter in the implementation above, we can compose a generalized method for estimation and postprocessing for any number of regimes:

```

Indicator_m <- function(x, c, m) {
  tmp <- rep(F, m)
  tmp[findInterval(x, c, left.open = T) + 1] <- T
  tmp
}

Xt_m <- function(x, t, p, d, c, m, z = x) {
  I <- Indicator_m(z[t - d], c, m)
  Y <- Yt(x, t, p)
  sapply(1:m, function(j) I[j] * Y)
}

# test for packages
suppressMessages(pkgTest("zeallot"))
suppressMessages(pkgTest("matlib"))

EstimSETAR_m <- function(x, p, d, c, m) {
  m = as.integer(m)
  if (m <= 0) {
    message("Error: regime count m has to be a positive integer")
    return(NA)
  }
  if (length(c) != m - 1) {
    message("Error: Incompatible dimensions of threshold vector and regime count.");
  }
}

```

```

    return(NA)
  }

  resultModel <- list()
  resultModel$nReg <- m
  resultModel$p = p; resultModel$d = d; resultModel$c = c;
  resultModel$data = x; n = length(x); resultModel$n = n;

  k <- max(p, d)

  X <- as.matrix(apply(as.matrix((k + 1):n), MARGIN=1, function(t) Xt_m(x, t, p, d, c, m) ))
  y <- as.matrix(x[(k + 1):n])
  K <- crossprod(t(X), t(X)); b <- crossprod(t(X), y);

  detK <- abs(det(K))

  if (detK < 0.000001) {
    return(NA)
  } else {
    K <- inv(K)
    sol_phi <- as.numeric(t(K %*% b)); sol_se <- sqrt(diag(K)/n);
    eps <- 0.01;

    # filter out those coeffs that are of the same order of magnitude as their errors
    filter <- sapply(1:(m*(p + 1)), function(i) ifelse(
      abs(sol_phi[i]) <= 2 * abs(sol_se[i]), 0, 1
    ))
  }

  sol_phi <- sol_phi * filter
  sol_se <- sol_se * filter

  solution <- cbind(phi = sol_phi, se = sol_se)
  resultModel$PhiParams <- solution[,1] # solving (X'X)*phi = X'y
  resultModel$PhiStErrors <- solution[,2] # standard errors
  skel <- crossprod(X, resultModel$PhiParams); resultModel$skel <- skel;
  resultModel$residuals <- (y - skel)
  resultModel$resSigmaSq <- 1 / (n - k) * sum(resultModel$residuals ^ 2)

  return(resultModel)
}

}

EstimSETAR_m_postproc <- function(model) {
  m <- model$nReg; x <- model$data; n <- model$n;
  k <- max(model$p, model$d); p <- model$p; c <- model$c;
  y <- as.matrix(x[(k + 1):n])
  skel <- model$skel;

  # regime counts
  nRegCounts <- rowSums( matrix(as.numeric( sapply(x, function(xi) Indicator_m(xi, c, m)) ), nrow=m) )
  model$nRegCounts <- nRegCounts

  # count valid p orders
  pOrders <- rowSums(
    matrix(as.numeric((
      sapply(1:(m * (p + 1)),
        function(i) ifelse( (i % (p + 1)) != 1 && model$PhiParams[i] != 0, 1, 0) )
    ))
  )

```

```

    )
  ), nrow=m, byrow=T)
)
model$pOrders <- pOrders

k_m <- pmax(pOrders, model$d) # k-offset for different regimes

# regime sigma sq:
regSigmaSq <- rowSums(
  matrix(as.numeric(
    sapply(y, function(xi) Indicator_m(xi, c, m)) ), nrow=m) %*% (y - skel)^2 ) / (nRegCounts - k_m)
model$regSigmaSq <- regSigmaSq
# names(model$regSigmaSq) <- sapply(1:m, function(j) paste0("rss", j))

model$AIC <- AIC_SETAR(pOrders, model$nRegCounts, model$resSigmaSq)
model$BIC <- BIC_SETAR(pOrders, model$nRegCounts, model$resSigmaSq)

c <- round(c, digits=3) # 3 dec. places seems enough
model$name <- paste0("SETAR(", p, ",", d, ",", paste(na.omit(c), collapse=','), ")")

return(model)
}

getModelName <- function(model) {
  p <- model$p; d <- model$d; c <- round(model$c, digits=3) # 3 dec. places seems enough
  paste0("SETAR(", p, ",", d, ",", paste(na.omit(c), collapse=','), ")")
}

str( EstimSETAR_m_postproc( EstimSETAR_m(xt, p=2, d=1, c=-0.1, m=2) ) ) # 2 regimes

## List of 17
## $ nReg      : int 2
## $ p         : num 2
## $ d         : num 1
## $ c         : num -0.1
## $ data      : num [1:170] 0.3 0.1 0.2 0.1 -0.1 ...
## $ n         : int 170
## $ PhiParams : num [1:6] 0.0711 0.7998 0.1755 0 0.6911 ...
## $ PhiStErrors: num [1:6] 0.023 0.0946 0.0666 0 0.0444 ...
## $ skel      : num [1:168, 1] 0.1162 0.1539 0.1005 -0.0534 -0.0264 ...
## $ residuals : num [1:168, 1] 0.0838 -0.0539 -0.2005 -0.0466 -0.0736 ...
## $ resSigmaSq : num 0.0293
## $ nRegCounts : num [1:2] 50 120
## $ pOrders    : num [1:2] 2 2
## $ regSigmaSq : num [1:2] 0.0357 0.0272
## $ AIC        : num -588
## $ BIC        : num -574
## $ name       : chr "SETAR(2,2,-0.1)"

str( EstimSETAR_m_postproc( EstimSETAR_m(xt, p=2, d=1, c=c(-0.1, 0.2), m=3) ) ) # 3 regimes

## List of 17
## $ nReg      : int 3
## $ p         : num 2
## $ d         : num 1
## $ c         : num [1:2] -0.1 0.2
## $ data      : num [1:170] 0.3 0.1 0.2 0.1 -0.1 ...
## $ n         : int 170

```

```
## $ PhiParams : num [1:9] 0.0711 0.7998 0.1755 -0.0171 0.3404 ...
## $ PhiStErrors: num [1:9] 0.02302 0.09461 0.06663 0.00844 0.09906 ...
## $ skel : num [1:168, 1] 0.0908 0.1688 0.0662 -0.0265 -0.0264 ...
## $ residuals : num [1:168, 1] 0.1092 -0.0688 -0.1662 -0.0735 -0.0736 ...
## $ resSigmaSq : num 0.0285
## $ nRegCounts : num [1:3] 50 84 36
## $ pOrders : num [1:3] 2 2 1
## $ regSigmaSq : num [1:3] 0.0302 0.0178 0.0538
## $ AIC : num -589
## $ BIC : num -573
## $ name : chr "SETAR(2,2,-0.1,0.2)"

str( EstimSETAR_m_postproc( EstimSETAR_m(xt, p=2, d=1, c=c(-0.1, 0.1, 0.2), m=4) ) ) # 4 regimes

## List of 17
## $ nReg : int 4
## $ p : num 2
## $ d : num 1
## $ c : num [1:3] -0.1 0.1 0.2
## $ data : num [1:170] 0.3 0.1 0.2 0.1 -0.1 ...
## $ n : int 170
## $ PhiParams : num [1:12] 0.0711 0.7998 0.1755 0 0.4321 ...
## $ PhiStErrors: num [1:12] 0.023 0.0946 0.0666 0 0.1364 ...
## $ skel : num [1:168, 1] 0.1028 0.1688 0.083 -0.0233 -0.0264 ...
## $ residuals : num [1:168, 1] 0.0972 -0.0688 -0.183 -0.0767 -0.0736 ...
## $ resSigmaSq : num 0.0285
## $ nRegCounts : num [1:4] 50 65 19 36
## $ pOrders : num [1:4] 2 2 1 1
## $ regSigmaSq : num [1:4] 0.0305 0.0172 0.0206 0.0536
## $ AIC : num -585
## $ BIC : num -567
## $ name : chr "SETAR(2,2,-0.1,0.1,0.2)"
```

#### 4.4 SETAR Estimation procedure

Now that we prepared all necessary functions we may proceed to search for 3-regime SETAR's in a suitable search space. This time we will construct our outer loop through delays  $d$  which will be reduced to only the delays that are contained within the models with detected remaining nonlinearity:

```
## unique delays:
## [1] 2 5
```

Still, even after this alleviation, the search might be computationally demanding. To obtain our results within reasonable time we use `foreach` and `doParallel` packages to compute search through  $c_1$  and  $c_2$  thresholds, and then process the results:

```
m <- 3 # 3-regime setars
pmax <- 7 # set maximum order p
# limit the c parameter by the 7.5-th and 92.5 percentile
cmin <- as.numeric(quantile(xt, 0.075)); cmax <- as.numeric(quantile(xt, 0.925));
h = (cmax - cmin) / 50 # determine the step by which c should be iterated

models3 <- list()
model3Columns <- list()

suppressMessages(pkgTest("foreach"))
suppressMessages(pkgTest("doParallel"))

pkgs <- c("zeallot", "matlib")
```

```

n_cores <- (detectCores() - 1)

for (d in delays) {
  for (p in d:pmax) {
    # PARALLEL LOOP

    cl <- makeCluster(n_cores)
    registerDoParallel(cl)
    pdModels <- foreach(c1 = seq(from = cmin, to = cmax - h, by = h), .packages = pkgs) %:%
      foreach(c2 = seq(c1 + h, cmax, by = h), .packages = pkgs) %dopar% {
        # skip models with slim regime
        if(sum(xt > c1 & xt < c2) < length(xt) * 0.15) {
          NA
        } else {
          tmp <- EstimSETAR_m(xt, p, d, c(c1, c2), m) # try to run the function
          # then test whether it returns `NA` as a result
          if (!as.logical(sum(is.na(tmp)))) {
            list(tmp)
          }
        }
      }
    }

    stopCluster(cl)

    # OLD LOOP:

    #pdModels <- list()
    #for(c1 in seq(from = cmin, to = cmax - h, by = h)) {
    #  for(c2 in seq(c1 + h, cmax, by = h)) {
    #    if(sum(xt > c1 & xt < c2) < length(xt) * 0.15) next # skip models with slim regime
    #    tmp <- EstimSETAR_m(xt, p, d, c(c1, c2), m) # try to run the function
    #    # then test whether it returns `NA` as a result
    #    if (!as.logical(sum(is.na(tmp)))) {
    #      pdModels[[length(pdModels) + 1]] <- tmp
    #    }
    #  }
    #}

    # frankly, this is a mess, but I can only get this nested list from the parallel loop
    pdOmitted <- lapply(unlist(pdModels, recursive=F),
      function(m) if(!is.logical(m) && !is.null(m)) m else list(list(resSigmaSq = Inf)))
    sigmas <- as.numeric(lapply(pdOmitted, function(m) m[[1]]$resSigmaSq))
    s_orders <- order(sigmas)
    # only the model whose parameter c gives the lowest residual square sum is chosen for postprocessing
    min_sigma_model <- EstimSETAR_m_postproc(pdOmitted[[ s_orders[1] ]][[1]])

    models3[[length(models3) + 1]] <- min_sigma_model
    model3Columns[[length(model3Columns) + 1]] <- c(
      min_sigma_model$p, min_sigma_model$pOrders, d, round(min_sigma_model$c, digits=4),
      min_sigma_model$nRegCounts,
      min_sigma_model$AIC, min_sigma_model$BIC,
      min_sigma_model$resSigmaSq)
  }
}

```

```

##   p p1 p2 p3 d      c1      c2  n1  n2 n3      AIC      BIC resSigmaSq
## 1 2  2  2  2 2 -0.2853 0.1102  20 109 41 -606.8036 -590.6016 0.02534142
## 2 3  2  2  3 2  0.1102 0.4032 129  28 13 -603.6662 -588.8303 0.02551154

```

```
## 7 5 2 2 4 5 -0.0949 0.1102 77 52 41 -606.4192 -584.9662 0.02480813
## 4 5 2 4 2 2 -0.2853 0.0077 20 82 68 -603.1484 -581.4691 0.02529007
## 3 4 2 4 2 2 -0.2853 0.0077 20 82 68 -599.7517 -578.0724 0.02580045
## 9 7 3 4 3 5 -0.1974 0.3007 34 119 17 -595.8356 -572.5017 0.02578773
## 8 6 3 3 5 5 -0.1974 0.1102 34 95 41 -597.4841 -570.8817 0.02524018
## 6 7 4 4 6 2 0.0077 0.3007 102 51 17 -593.9537 -565.3372 0.02487616
## 5 6 4 5 4 2 -0.1974 0.0077 34 68 68 -586.9812 -554.9348 0.02622437
```

These are the SETAR3 models that can replace the SETAR2's with remaining nonlinearity, ordered by *BIC* with estimated coefficients:

```
## $`2/2/-0.2853/0.1102`
##           [,1]      [,2]      [,3] [,4]      [,5]      [,6]      [,7]
## Phi      0.34687231 0.3175267 1.2763044 0 0.4718739 0.56917316 0.09138806
## stdError 0.08510169 0.1031526 0.2302884 0 0.0460411 0.07670688 0.02581381
##           [,8]      [,9]
## Phi      1.08026643 -0.42454326
## stdError 0.05901015 0.07802637
##
## $`3/2/0.1102/0.4032`
##           [,1]      [,2]      [,3] [,4] [,5]      [,6] [,7]      [,8]
## Phi      0 0.41811454 0.40547459 0 0 1.19010389 0 -0.23475893
## stdError 0 0.04340129 0.05441896 0 0 0.08788666 0 0.07578564
##           [,9]      [,10]      [,11]      [,12]
## Phi      0.20246225 1.06478118 -0.8092012 0.2569796
## stdError 0.09725574 0.08660258 0.2156003 0.1267073
##
## $`5/5/-0.0949/0.1102`
##           [,1]      [,2]      [,3] [,4] [,5] [,6] [,7]      [,8]
## Phi      0.06573613 0.72504540 0.31344645 0 0 0 0 0.56641566
## stdError 0.01696259 0.05804913 0.07220894 0 0 0 0 0.05909938
##           [,9] [,10] [,11] [,12] [,13]      [,14]      [,15]      [,16]
## Phi      0.46198906 0 0 0 0 0.47118620 -0.19079304 0.53307667
## stdError 0.07050245 0 0 0 0 0.07585229 0.08097692 0.09450724
##           [,17] [,18]
## Phi      -0.3303657 0
## stdError 0.1058457 0
##
## $`5/2/-0.2853/0.0077`
##           [,1]      [,2]      [,3] [,4] [,5] [,6]      [,7]      [,8]
## Phi      0.31942101 0.2938117 1.2154507 0 0 0 0.04951122 0.45458870
## stdError 0.09212305 0.1123581 0.2723039 0 0 0 0.01396373 0.06405413
##           [,9]      [,10] [,11]      [,12]      [,13]      [,14] [,15]
## Phi      0.8435796 0.21470044 0 -0.18897376 0.07606997 0.83058402 0
## stdError 0.1370189 0.06106055 0 0.05445145 0.01665680 0.04628859 0
##           [,16] [,17] [,18]
## Phi      -0.2019964 0 0
## stdError 0.0686019 0 0
##
## $`4/2/-0.2853/0.0077`
##           [,1]      [,2]      [,3] [,4] [,5]      [,6]      [,7]
## Phi      0.34409573 0.3118261 1.2875116 0 0 0.04355466 0.43838144
## stdError 0.08873654 0.1108958 0.2625368 0 0 0.01385785 0.06388366
##           [,8]      [,9]      [,10]      [,11]      [,12] [,13]
## Phi      0.7760593 0.19839730 -0.14724371 0.06390110 0.84271491 0
## stdError 0.1356306 0.06087958 0.04825214 0.01599119 0.04567036 0
##           [,14] [,15]
## Phi      -0.21031590 0
## stdError 0.06618311 0
##
```

```

## $`7/5/-0.1974/0.3007`
##      [,1]      [,2]      [,3] [,4] [,5] [,6]      [,7] [,8]      [,9]
## Phi      0 0.5415693 0.4210005 0 0 0 -0.2367908 0 0.025819168
## stdError 0 0.1009274 0.1190252 0 0 0 0.0841457 0 0.007674703
##      [,10]      [,11] [,12] [,13]      [,14]      [,15] [,16] [,17]
## Phi      0.66976648 0.26856721 0 0 -0.46931684 0.12332324 0 0
## stdError 0.04260034 0.05210529 0 0 0.07377245 0.05252777 0 0
##      [,18] [,19]      [,20] [,21]      [,22]      [,23] [,24]
## Phi      0 0 0.4055609 0 -0.5926826 0.5607706 0
## stdError 0 0 0.1345050 0 0.2416369 0.1814230 0
##
## $`6/5/-0.1974/0.1102`
##      [,1]      [,2]      [,3] [,4] [,5] [,6]      [,7]      [,8]
## Phi      0 0.51519188 0.3614104 0 0 0 -0.21550315 0.024084657
## stdError 0 0.09922349 0.1114737 0 0 0 0.08281528 0.008288851
##      [,9]      [,10] [,11] [,12]      [,13] [,14] [,15]      [,16]
## Phi      0.65482079 0.39266114 0 0 -0.5626092 0 0 0.5048728
## stdError 0.04672132 0.05853316 0 0 0.1059246 0 0 0.0771052
##      [,17]      [,18]      [,19] [,20]      [,21]
## Phi      -0.26200260 0.64842165 -0.3995987 0 0.32971940
## stdError 0.08209169 0.09731003 0.1069363 0 0.06466259
##
## $`7/2/0.0077/0.3007`
##      [,1]      [,2]      [,3]      [,4] [,5]      [,6] [,7] [,8]
## Phi      0.02341172 0.38968193 0.46653491 0.23801400 0 -0.16028150 0 0
## stdError 0.01130769 0.05573268 0.07445277 0.05535027 0 0.05038605 0 0
##      [,9]      [,10] [,11]      [,12] [,13] [,14]      [,15]      [,16]
## Phi      0 0.81219656 0 -0.27758165 0 0 0.3806256 -0.27796262
## stdError 0 0.06152922 0 0.07875932 0 0 0.1019897 0.07170216
##      [,17]      [,18]      [,19] [,20]      [,21] [,22]      [,23]
## Phi      0.2487305 1.2487327 -0.7890808 0.4129725 -1.022866 0 1.9675142
## stdError 0.0817996 0.1439504 0.2204510 0.2029055 0.255567 0 0.4123954
##      [,24]
## Phi      -1.2986218
## stdError 0.2743571
##
## $`6/2/-0.1974/0.0077`
##      [,1]      [,2] [,3]      [,4]      [,5] [,6]      [,7]
## Phi      -0.18973908 0.39250631 0 0.32459817 0.22062185 0 -0.45527112
## stdError 0.04825684 0.08481379 0 0.08746111 0.08813265 0 0.09514038
##      [,8]      [,9]      [,10]      [,11]      [,12]      [,13]
## Phi      0.03992378 0.51451199 0.6313824 0.28209233 -0.18357269 -0.23118701
## stdError 0.01604076 0.07843829 0.2303664 0.07434585 0.07418622 0.06738637
##      [,14]      [,15]      [,16] [,17]      [,18] [,19]      [,20]
## Phi      0 0.06337701 0.82804564 0 -0.15770506 0 -0.25743535
## stdError 0 0.01776836 0.04643295 0 0.07122721 0 0.08731039
##      [,21]
## Phi      0.17899366
## stdError 0.07599395

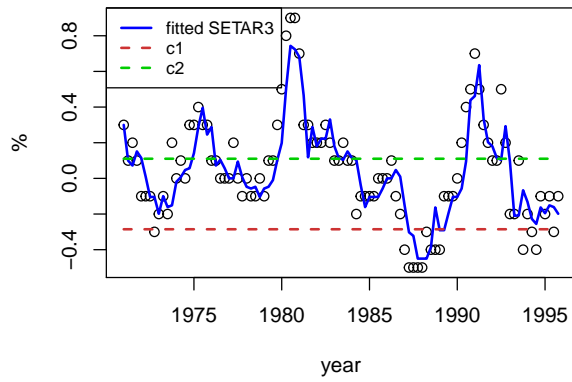
```

As we might notice, some models differ only in their information criteria and coefficients, since they were estimated from a different maximum order  $p$ .

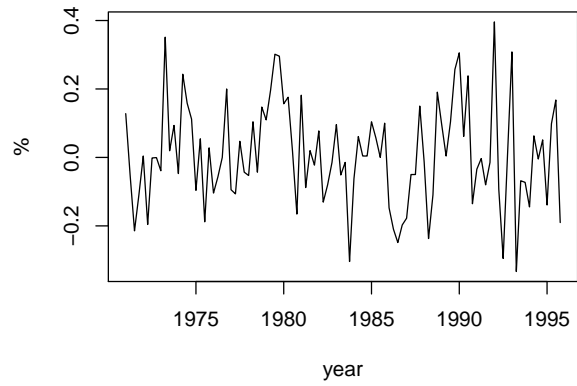


#### 4.5 SETAR3 Visualisation

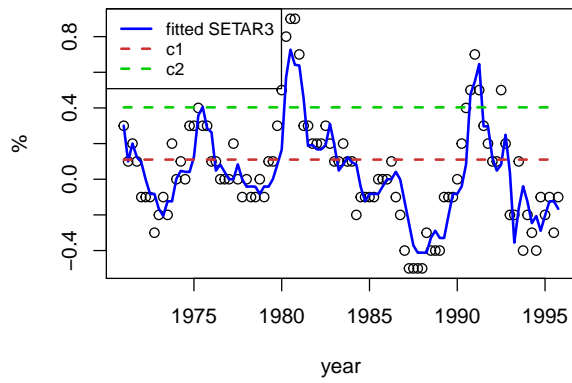
**SETAR(2,2,-0.285,0.11)**



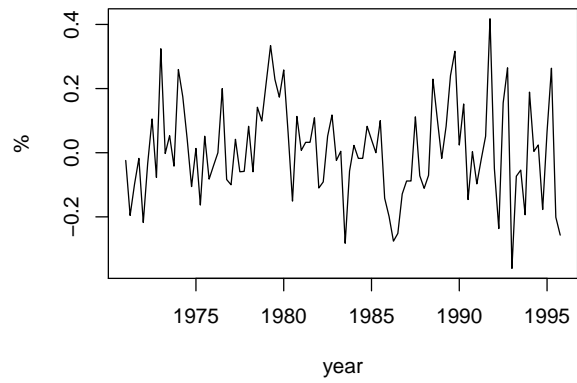
**SETAR(2,2,-0.285,0.11) Residuals**



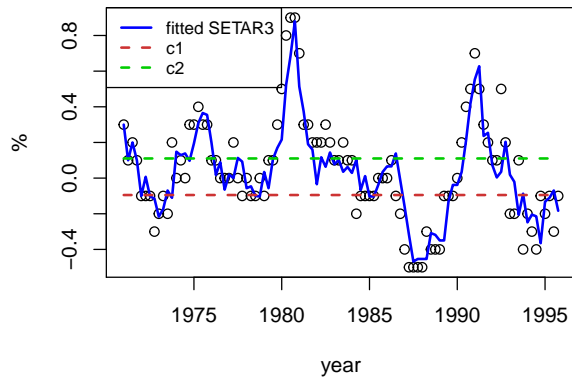
**SETAR(3,2,0.11,0.403)**



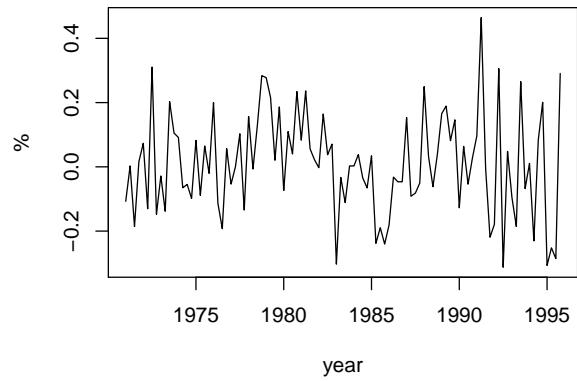
**SETAR(3,2,0.11,0.403) Residuals**



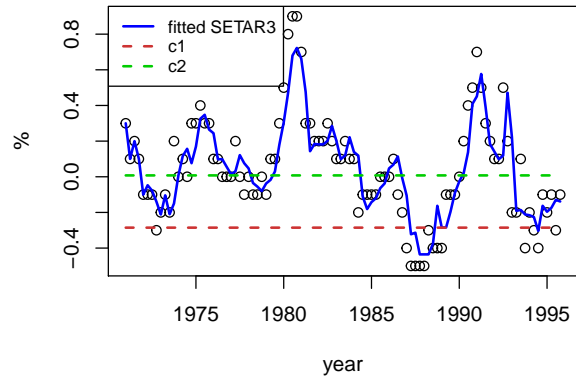
**SETAR(5,5,-0.095,0.11)**



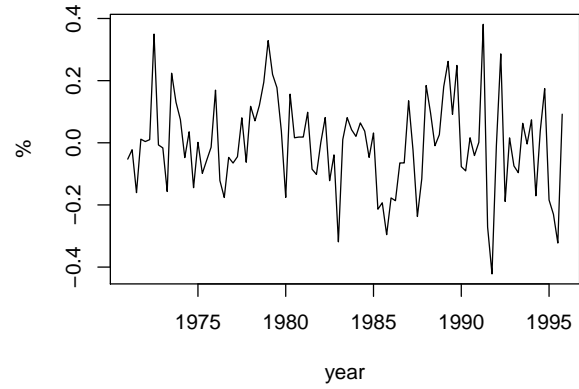
**SETAR(5,5,-0.095,0.11) Residuals**



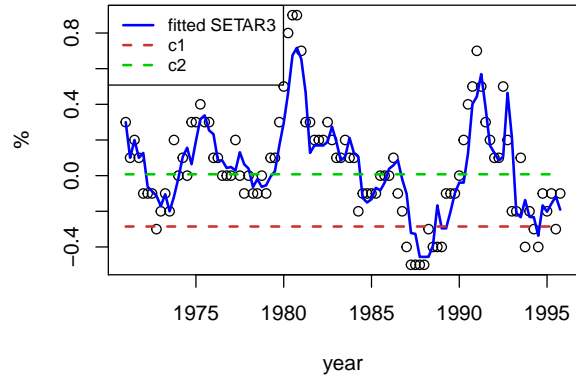
**SETAR(5,2,-0.285,0.008)**



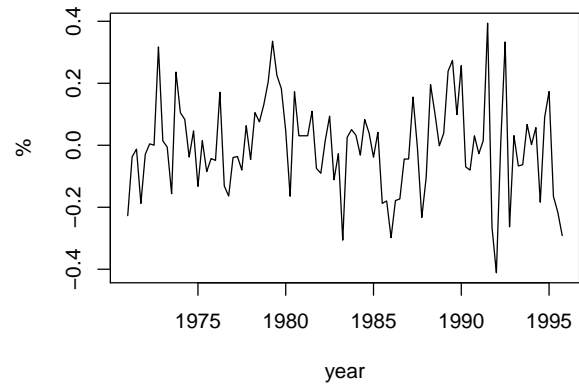
**SETAR(5,2,-0.285,0.008) Residuals**



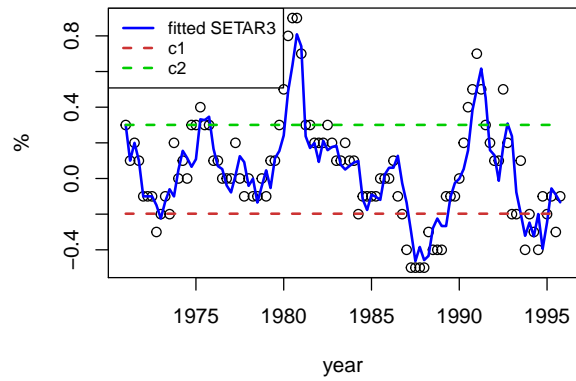
**SETAR(4,2,-0.285,0.008)**



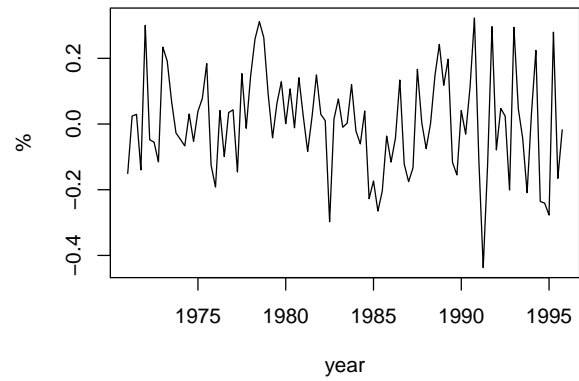
**SETAR(4,2,-0.285,0.008) Residuals**

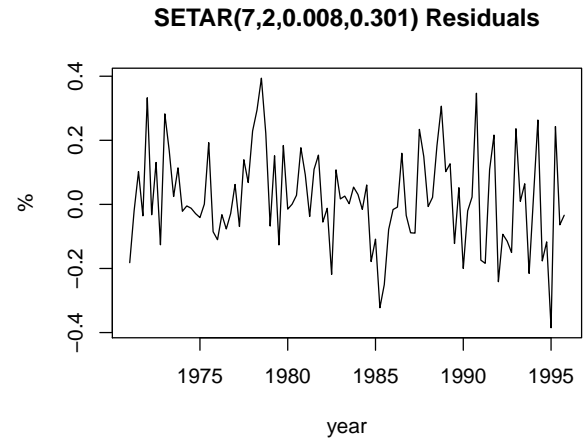
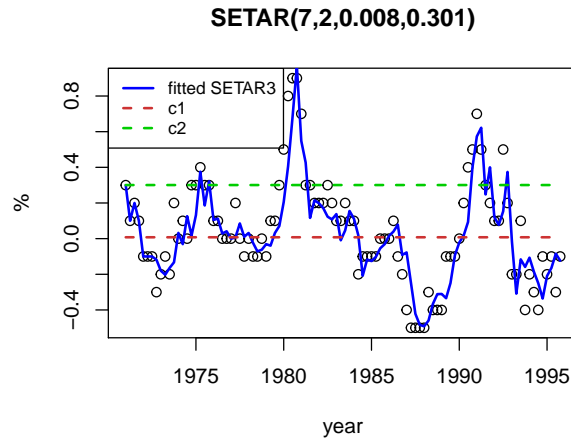
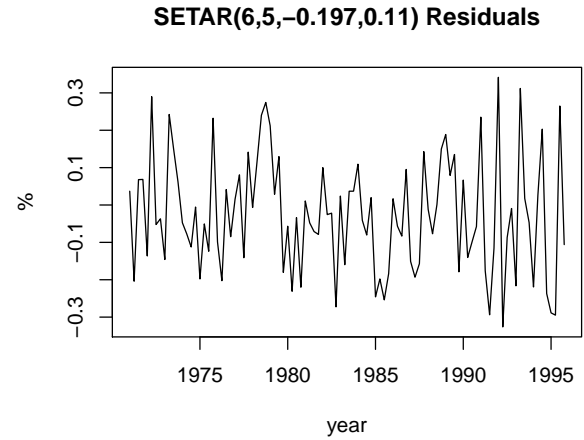
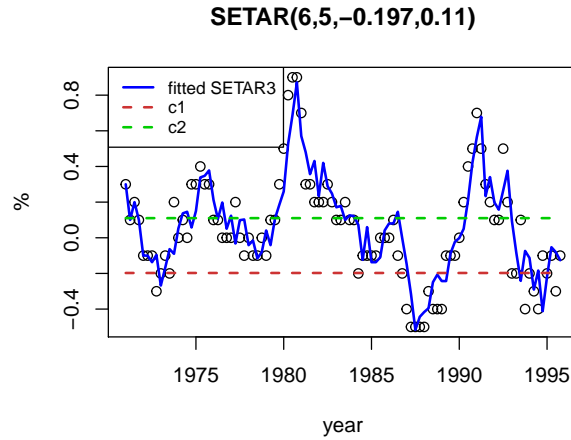


**SETAR(7,5,-0.197,0.301)**



**SETAR(7,5,-0.197,0.301) Residuals**





## 4.6 Conclusion

Due to our limited sampling space of `delays` we obtained 9 possible  $SETAR3(p, d, c_1, c_2)$  models, and since remaining SETAR3 nonlinearity was detected in 5 of the original SETAR2 models, we replace them by the top 5 newly found SETAR3's:

##	model	BIC
## 1	SETAR(2,2,-0.285,0.11)	-590.602
## 2	SETAR(3,2,0.11,0.403)	-588.83
## 3	SETAR(3,1,0.301)	-586.083
## 4	SETAR(1,1,0)	-585.638
## 5	SETAR(5,5,-0.095,0.11)	-584.966
## 6	SETAR(2,1,0.205)	-583.975
## 7	SETAR(5,2,-0.285,0.008)	-581.469
## 8	SETAR(4,2,-0.285,0.008)	-578.072
## 9	SETAR(4,1,0.301)	-576.504
## 10	SETAR(4,4,0.205)	-575.688
## 11	SETAR(3,3,-0.197)	-574.277
## 12	SETAR(5,4,0.205)	-571.548

## 5. Predictions via SETAR Models and Their Evaluation

### 5.1 Helper functions

Since we have not yet defined a skeleton function, i.e. one that would continue plugging in the time series values even after the end of testing data. For that purpose we implement a step-forward function for an  $m$ -regime SETAR:

```
SETAR_m_singleStep <- function(model, x, t) {
  m <- model$nReg; n <- model$n;
  p <- model$p; d <- model$d; c <- model$c;
  # parameter matrix with regime coefficients by row
  Phi <- matrix(model$PhiParams, nrow=m)

  # extract regime vector
  X <- Xt_m(x, t, p, d, c, m)
  reg_id <- which(colSums(X!=0)!=0)
  x_reg <- X[,reg_id]
  Phi[reg_id,] %*% x_reg
}
```

We can test it on a particular SETAR model:

```
n_ahead <- 1
model <- models[[ orders[3] ]]
x_out <- c(x_train, rep(0, n_ahead)); nt <- length(x_train)

for (i in 1:n_ahead) {
  x_out[nt + i] <- SETAR_m_singleStep(model, x_out, nt + i)
}
xt <- c(xt, x_eval) # fuse test and eval data
```

```
## x_out:
```

```
## [1] -0.100000000 -0.004246784
```

```
## data:
```

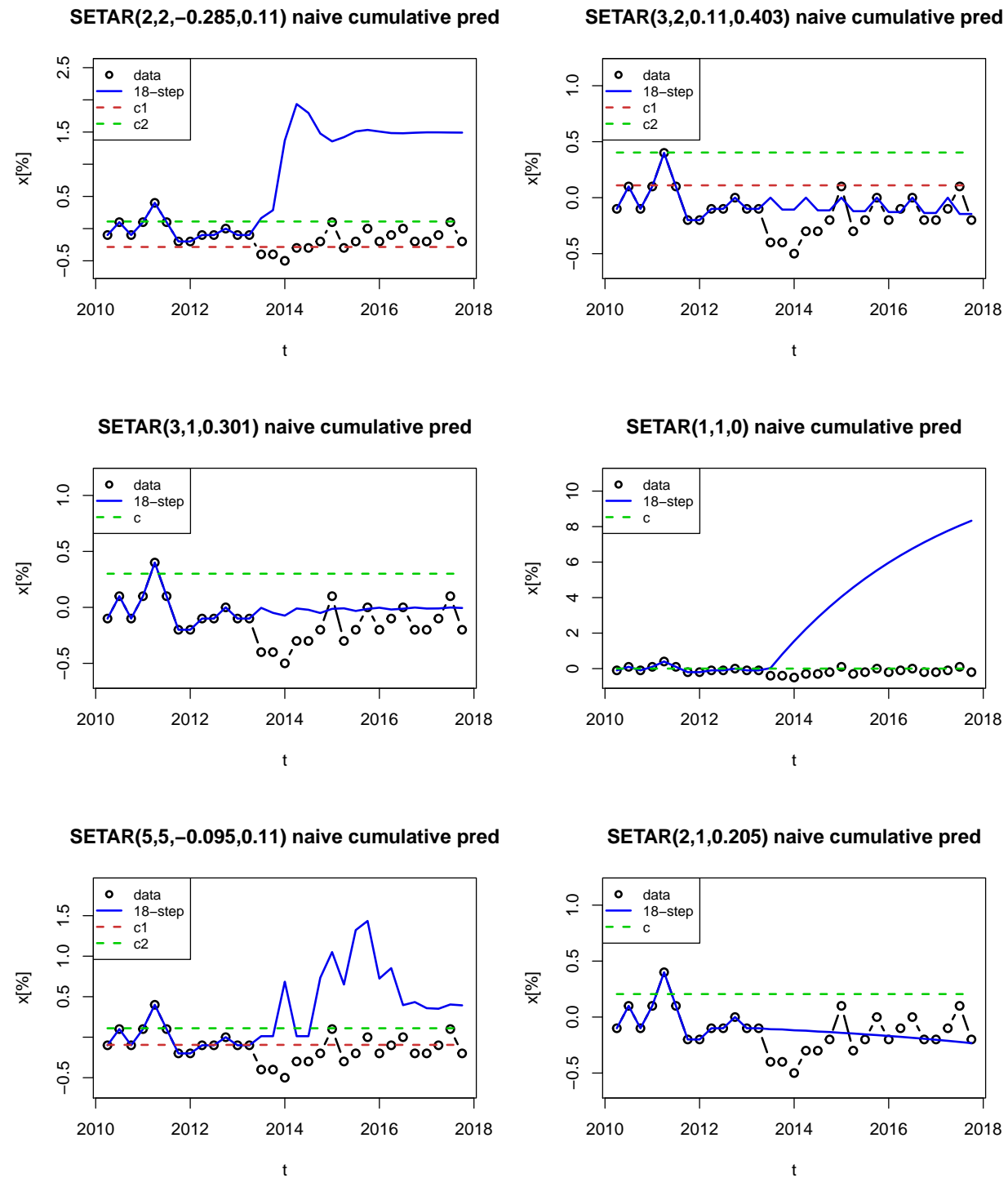
```
## [1] -0.1 -0.4
```

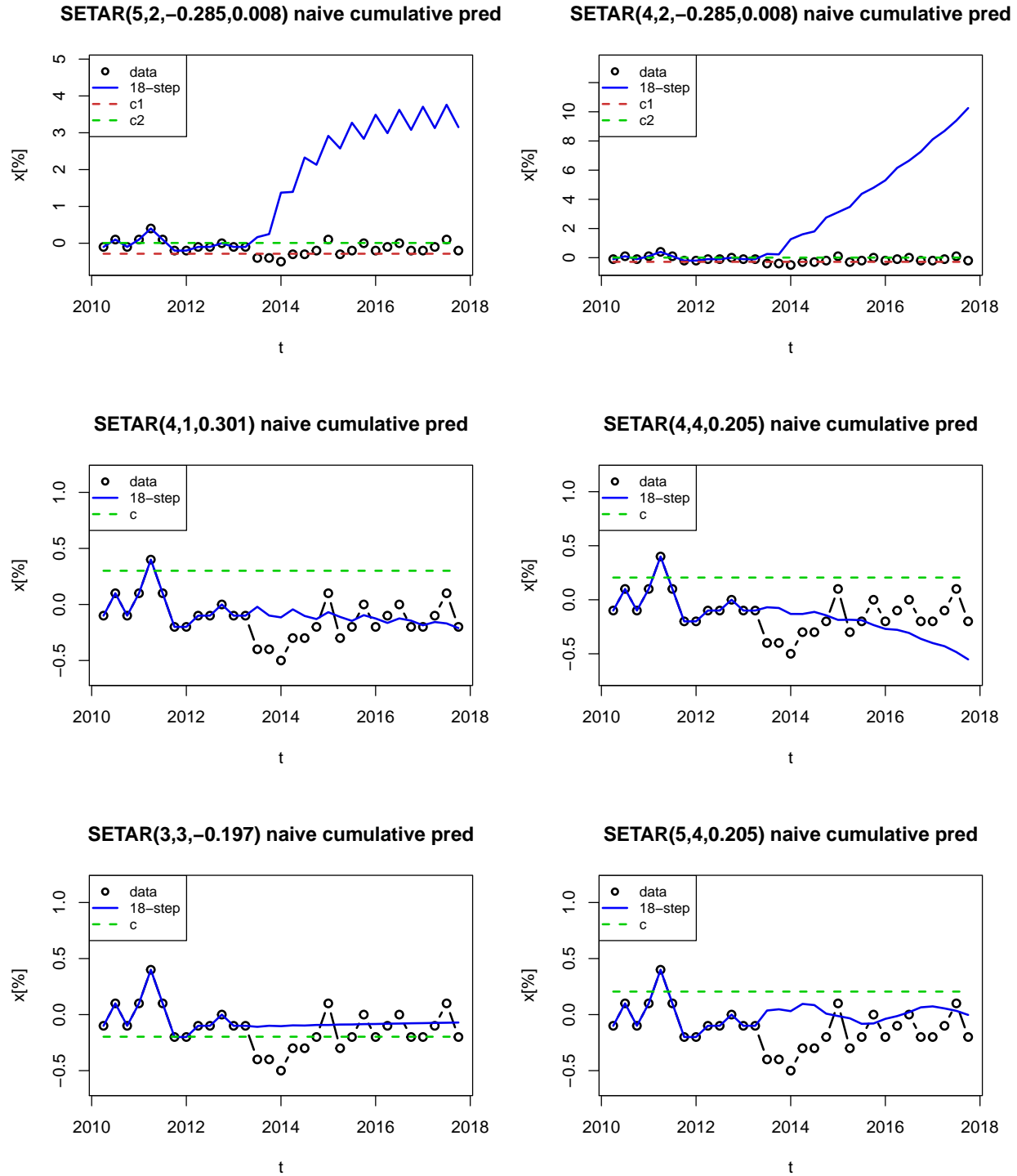
This is a single-step prediction of the data using the first model. When we set `n_ahead > 1`, the step function builds up upon previous predicted values and the skeleton converges to a model's particular equilibrium:

```
##          x_out data
## 1 -0.100000000 -0.1
## 2 -0.004246784 -0.4
## 3 -0.049043561 -0.4
## 4 -0.073728157 -0.5
## 5 -0.009866531 -0.3
## 6 -0.021643760 -0.3
## 7 -0.050297408 -0.2
## 8 -0.013133059  0.1
## 9 -0.008468800 -0.3
## 10 -0.031987160 -0.2
## 11 -0.013626195  0.0
## 12 -0.003055054 -0.2
## 13 -0.019010891 -0.1
## 14 -0.012166789  0.0
## 15 -0.001374286 -0.2
## 16 -0.010548972 -0.2
## 17 -0.009768627 -0.1
## 18 -0.001173186  0.1
## 19 -0.005448426 -0.2
```

```
##
## equilibria:
## [1] 0
```

If the model is explosive, the predictions will diverge:





The examples above tested 18 steps of a naive approach to prediction, by assuming the process evolves via its skeleton. More convenient approaches are **Monte Carlo** ("MC") and **Bootstrap**. Both rely on adding noise to series predictions. Monte Carlo approach simulates normally distributed noise  $\epsilon \sim N(0, \hat{\sigma}_\epsilon)$  from residual standard error  $\hat{\sigma}_\epsilon$ , and Bootstrap, on the other hand, does not assume the normality of model residuals and instead randomly samples the residuals themselves.

We will use these two approaches in the following implementation:

```

PredictSETAR <- function(
  model, x_train, x_eval, horizon=(length(x_eval) + 1), n_ahead=1,
  type=c("naive", "MC", "bootstrap"), Nboot=100, alpha=0.2, refit=F,
  single.step=F, return.paths=F) {

  type <- match.arg(type)

  p <- model$p; d <- model$d; # model dims

  if(missing(x_train)) {
    x_train = model$data # training sample
  }

  # result series
  x_res <- x_train
  # training sample size
  nt <- length(x_res)
  sd_res <- sqrt(model$resSigmaSq)

  # extract model residuals
  resid <- as.numeric(model$residuals)
  resid <- resid[!is.na(resid)]

  # fill the prediction part of the array with zeros
  x_res <- c(x_res, rep(0, horizon))
  xrange <- p - ((p - 1):0)

  if(type=="naive") Nboot <- 1

  predictions <- function(x_res, eval.model=model, tmin=nt) {
    noise <- switch(
      type,
      "naive"= rep(0, n_ahead),
      "MC"= rnorm(n_ahead, mean = 0, sd=sd_res),
      "bootstrap" = sample(resid, size=n_ahead, replace=T)
    )

    for(t in (tmin + (1:n_ahead))) {
      x_res[t] <- SETAR_m_singleStep(eval.model, x_res, t)
      x_res[t] <- x_res[t] + noise[t - tmin]
    }

    if (!single.step) {
      return(x_res)
    }
    return(x_res[tmin + n_ahead])
  }

  if(single.step) {
    n_ahead <- 1
    x_data <- c(x_train, x_eval)
    if (nt + horizon > length(x_data)) horizon <- (length(x_eval) + 1)

    c <- model$c; m <- model$nReg
    if (refit) {
      fit_model <- EstimSETAR_m_postproc( EstimSETAR_m(x_data, p, d, c, m) )
    } else {
      fit_model <- model
    }
  }
}

```

```

}

x_simulations <- matrix(rep(x_data[nt], Nboot), ncol=Nboot)
for (t in (nt + 1:horizon)) {
  x_source <- x_data[1:(t - 1)]
  x_simulations <- rbind(x_simulations, replicate(Nboot,
    predictions(x_source, eval.model=fit_model, tmin=(t - 1))
  ))
  x_res[t] <- mean(x_simulations[t - nt,])
}
x_pred <- x_res[nt + 1:horizon]
} else {
  # === MULTISTEP ===
  if (n_ahead == 1) n_ahead <- horizon

  x_simulations <- replicate(Nboot, predictions(x_res))
  x_sim_means <- rowMeans(x_simulations, na.rm=T)
  x_pred <- x_sim_means[(nt - 1) + 1:n_ahead]
}

# if not naive compute conf. intervals:
x_errors <- x_pred
if(type != "naive") {
  x_errors <- t(apply(
    x_simulations[(nt - 1) * (!single.step) + 1:horizon, ,drop=F], MARGIN=1, quantile,
    prob=sort(c(alpha, 1 - alpha)), na.rm=T))
}

# compute prediction errors
MSE <- sum((x_eval[1:horizon] - x_pred)^2, na.rm=T) / horizon
# RMSE <- sqrt(MSE)

if(type == "naive"){
  result <- list(pred=x_pred, MSE=MSE)
} else {
  if (return.paths) result <- list(pred=x_pred, se=x_errors, MSE=MSE, alpha=alpha,
    paths=x_simulations[(nt - 1) * (!single.step) + 1:horizon,])
  else result <- list(pred=x_pred, se=x_errors, MSE=MSE, alpha=alpha)
}

return(result)
}

```

Now we test it on our data:

## naive (10-step):

```
## [1] -0.100 -0.100  0.036  0.047  0.031  0.095  0.085  0.008 -0.013 -0.033
## [11] -0.082
```

## Monte Carlo (10-step):

```
## [1] -0.100 -0.100  0.032  0.029  0.003  0.087  0.182  0.207  0.213  0.292
## [11]  0.375
```

## Bootstrap (10-step):

```
## [1] -0.100 -0.100  0.064  0.050  0.063  0.124  0.189  0.254  0.274  0.321
## [11]  0.416
```

## data:

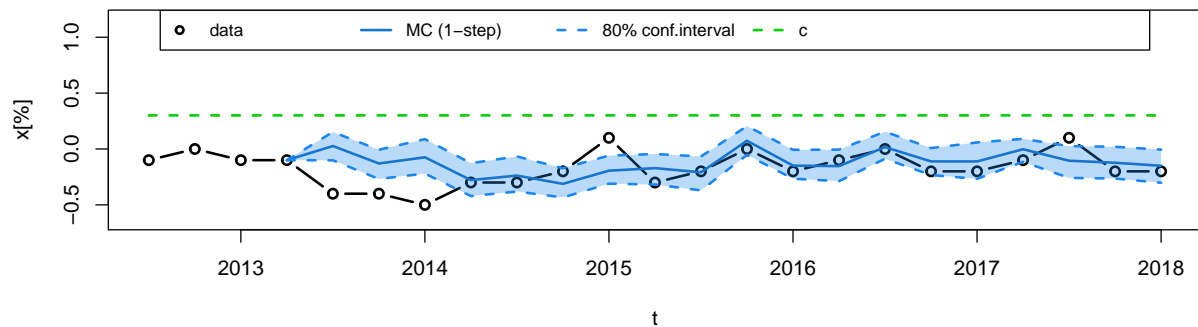


```
## [1] -0.1 -0.4 -0.4 -0.5 -0.3 -0.3 -0.2  0.1 -0.3 -0.2  0.0
## naive (1-step):
## [1] -0.100 -0.100  0.036 -0.129 -0.150 -0.191 -0.026  0.105  0.206  0.320
## [11]  0.077
## Monte Carlo (1-step):
## [1] -0.100 -0.100  0.052 -0.105 -0.154 -0.172 -0.021  0.116  0.225  0.315
## [11]  0.094
## Bootstrap (1-step):
## [1] -0.100 -0.100  0.075 -0.126 -0.147 -0.170 -0.021  0.157  0.224  0.401
## [11]  0.070
## data:
## [1] -0.1 -0.4 -0.4 -0.5 -0.3 -0.3 -0.2  0.1 -0.3 -0.2  0.0
```

To test other prediction results, such as confidence intervals and simulation paths, we implement a plot procedure:

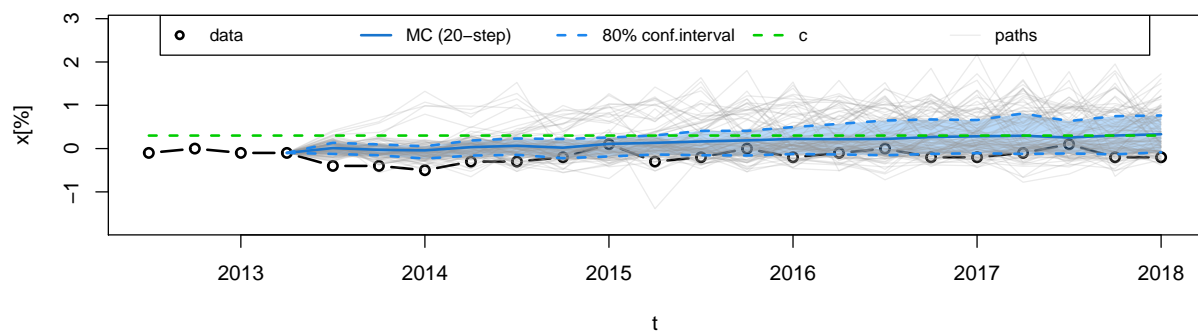
```
par(mfrow=c(1, 1))
predictSETAR_andPlot(
  models[[ orders[3] ]], time=dat$time, x_train=x_train, x_eval=x_eval,
  pred_type="MC", single.step=T, plot.paths=F)
```

**SETAR(3,1,0.301) Monte Carlo (1-step) , RMSE = 0.2084**



```
predictSETAR_andPlot(
  models[[ orders[3] ]], time=dat$time, x_train=x_train, x_eval=x_eval,
  pred_type="MC", single.step=F, plot.paths=T)
```

**SETAR(3,1,0.301) Monte Carlo (20-step) sim, RMSE = 0.3454**

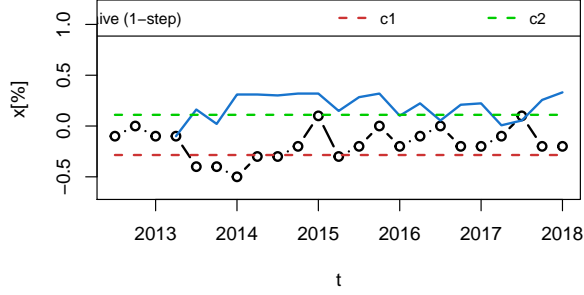


Now we possess all necessary tools to proceed evaluating our SETAR models according to their predictive abilities.

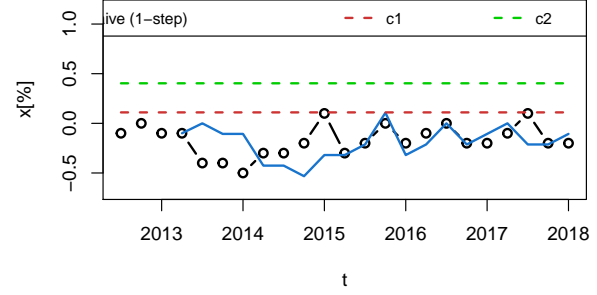
## 5.2 Single-Step Predictions of SETAR Models

Naive:

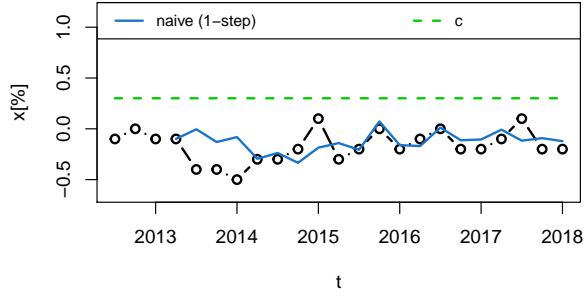
**SETAR(2,2,-0.285,0.11) Naive (1-step) , RMSE = 0.4054**



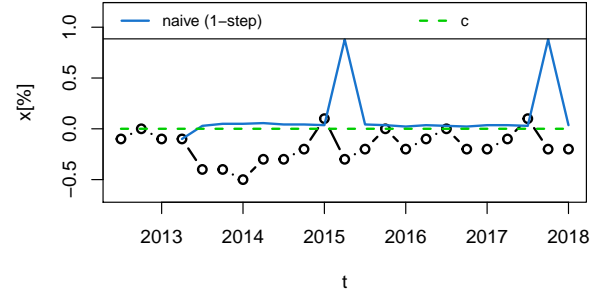
**SETAR(3,2,0.11,0.403) Naive (1-step) , RMSE = 0.2457**



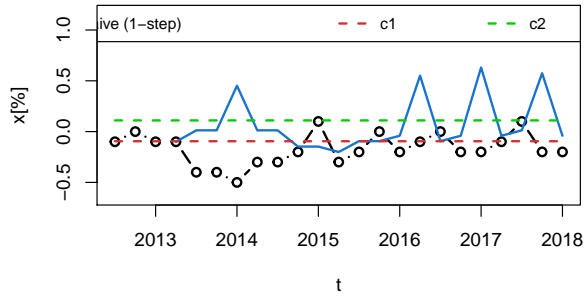
**SETAR(3,1,0.301) Naive (1-step) , RMSE = 0.2086**



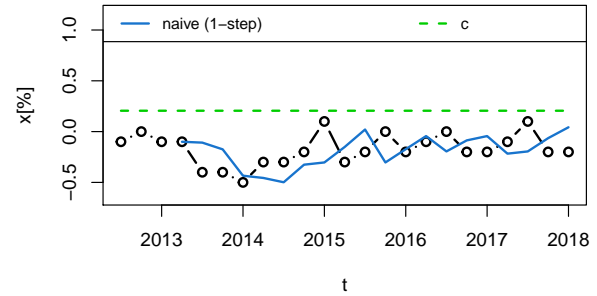
**SETAR(1,1,0) Naive (1-step) , RMSE = 0.4231**



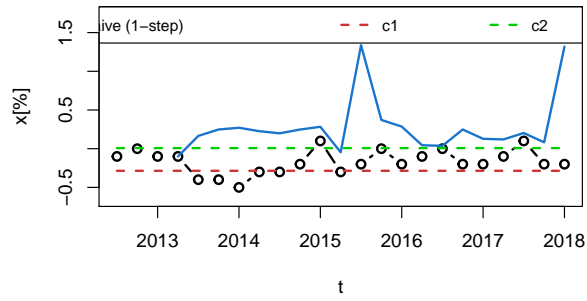
**SETAR(5,5,-0.095,0.11) Naive (1-step) , RMSE = 0.3792**



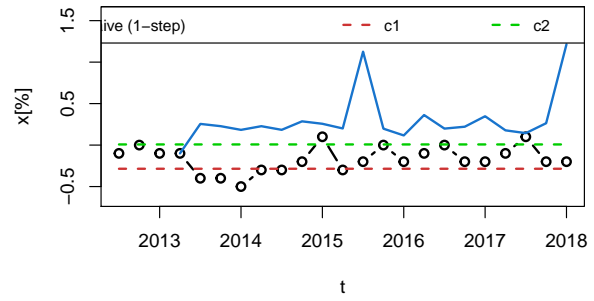
**SETAR(2,1,0.205) Naive (1-step) , RMSE = 0.1936**



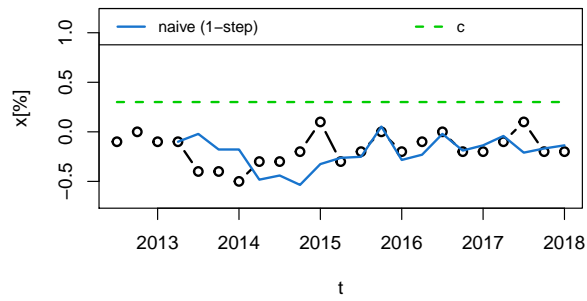
**SETAR(5,2,-0.285,0.008) Naive (1-step) , RMSE = 0.4967**



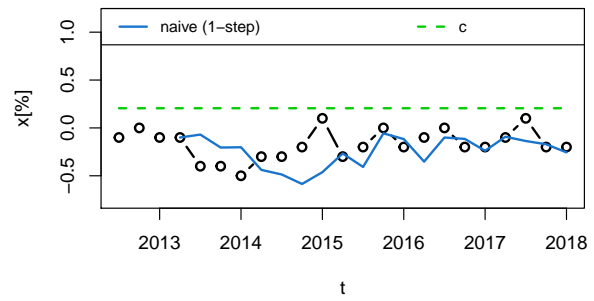
**SETAR(4,2,-0.285,0.008) Naive (1-step) , RMSE = 0.4844**



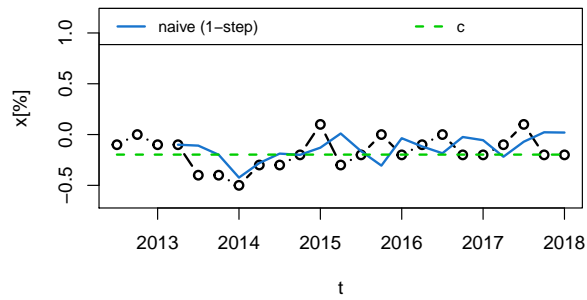
**SETAR(4,1,0.301) Naive (1-step) , RMSE = 0.2367**



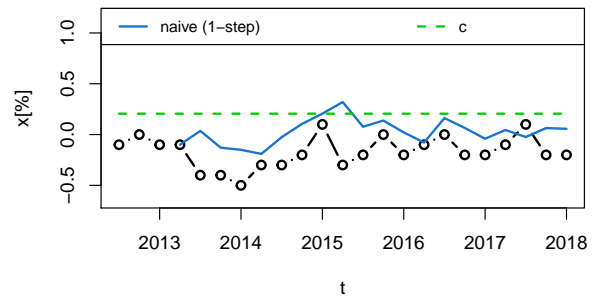
**SETAR(4,4,0.205) Naive (1-step) , RMSE = 0.2537**



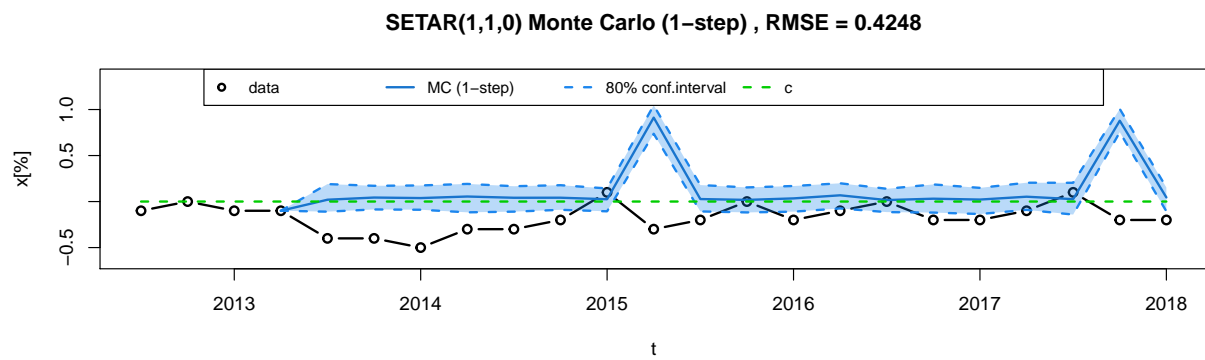
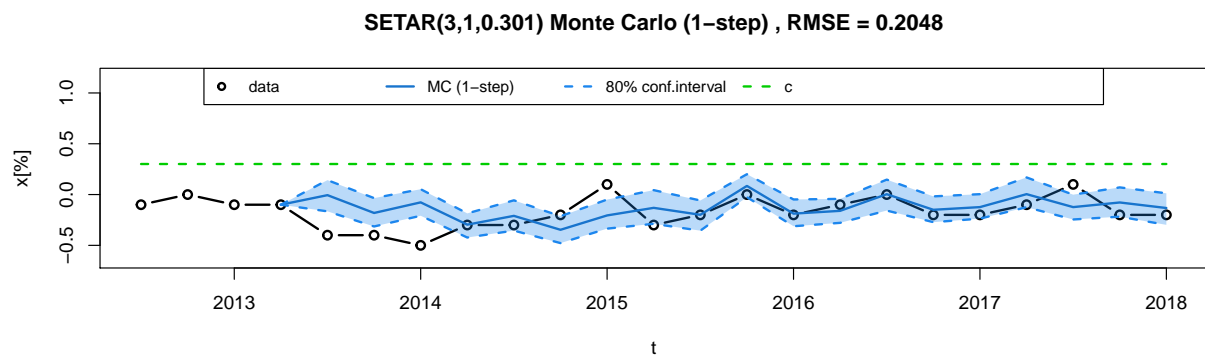
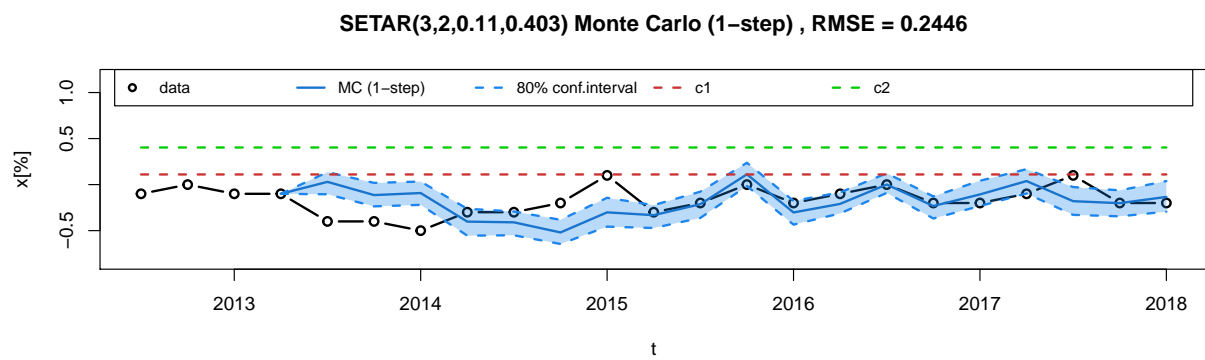
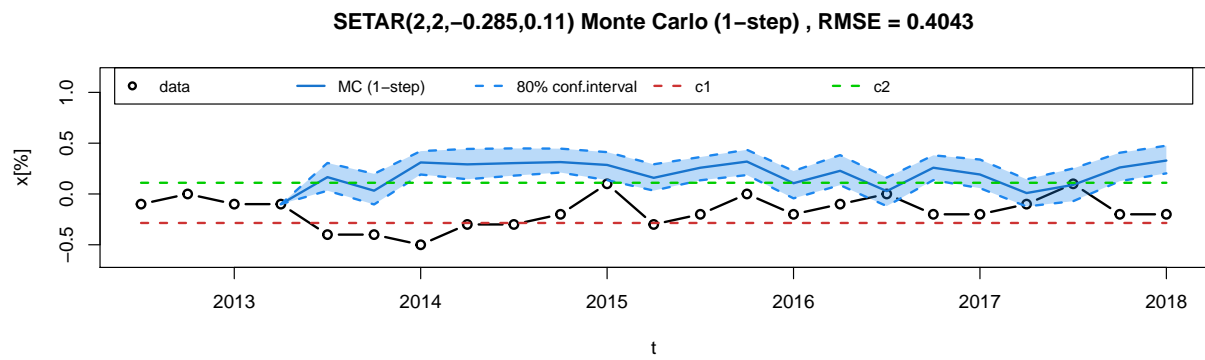
**SETAR(3,3,-0.197) Naive (1-step) , RMSE = 0.1872**

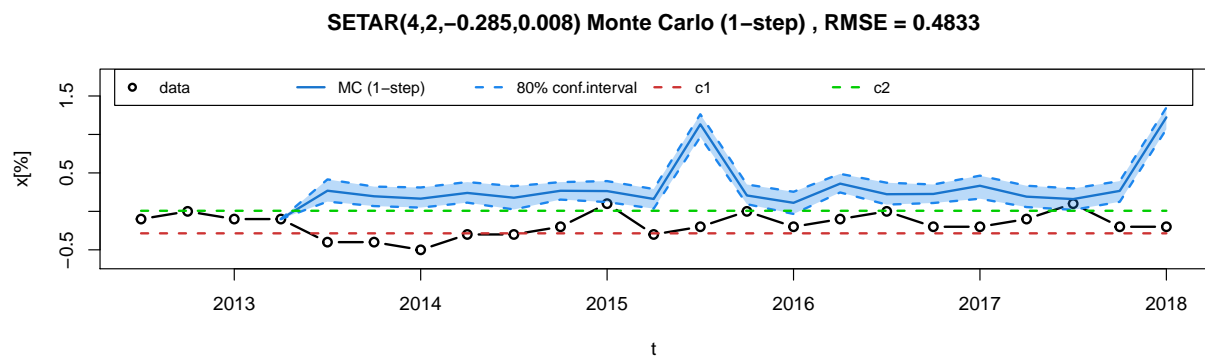
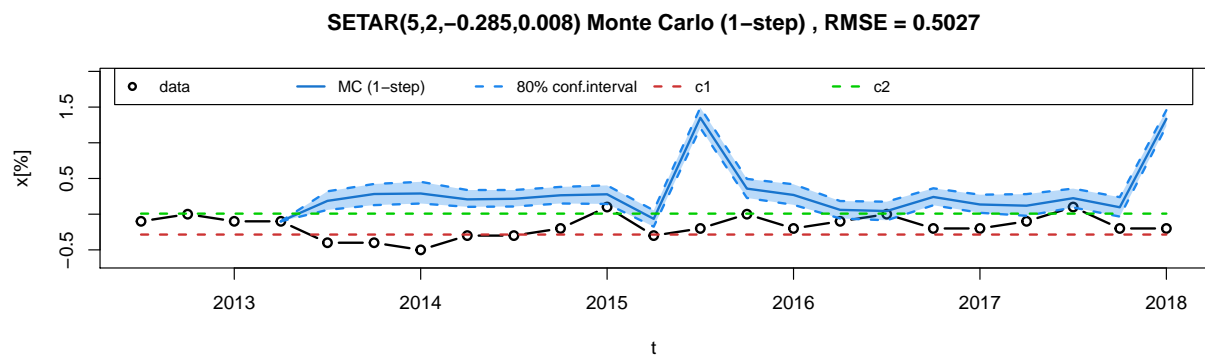
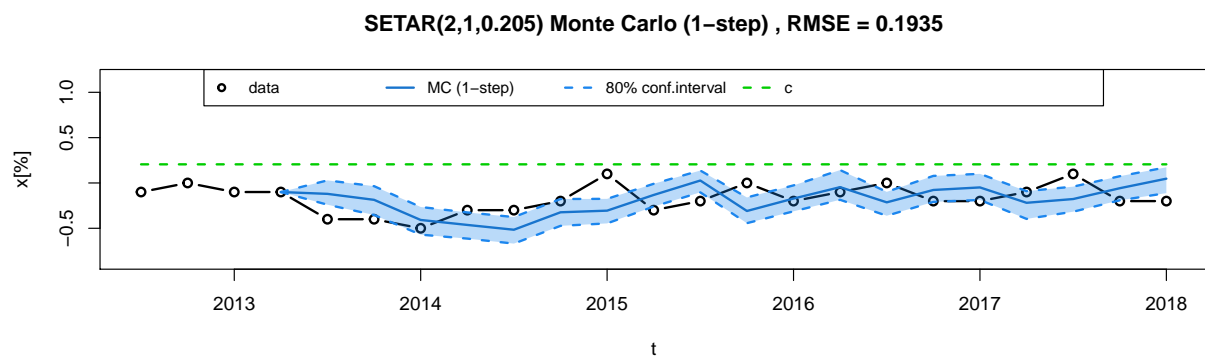
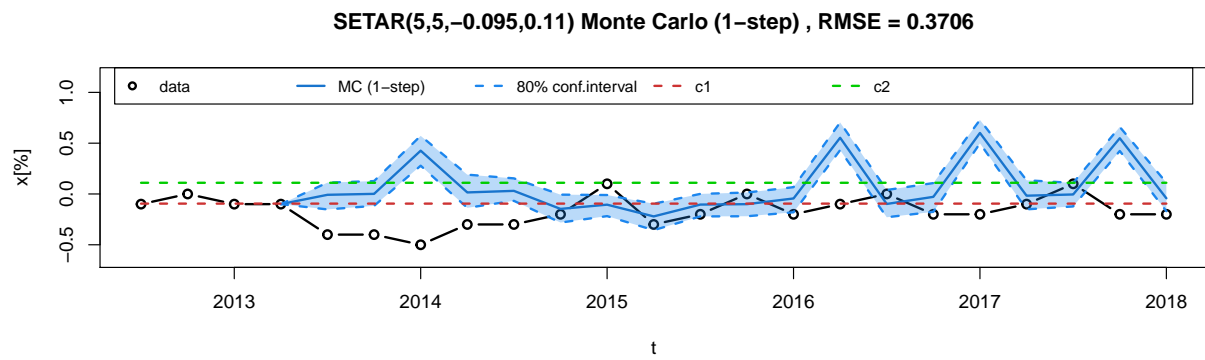


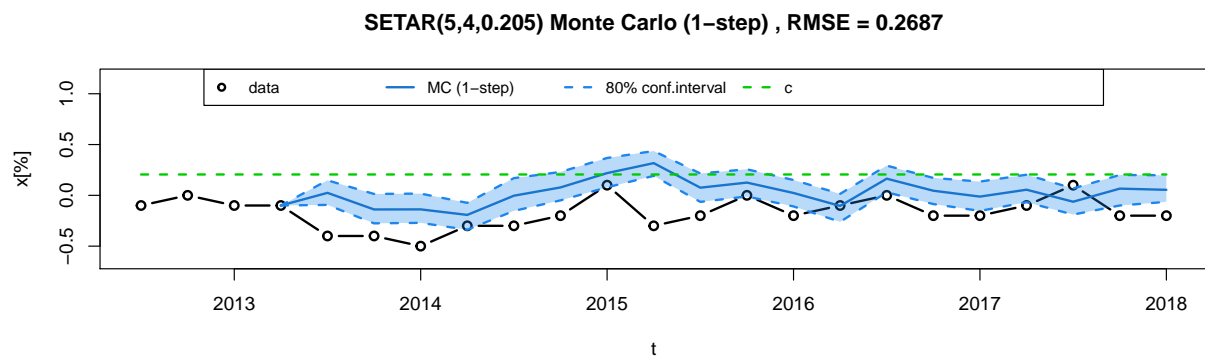
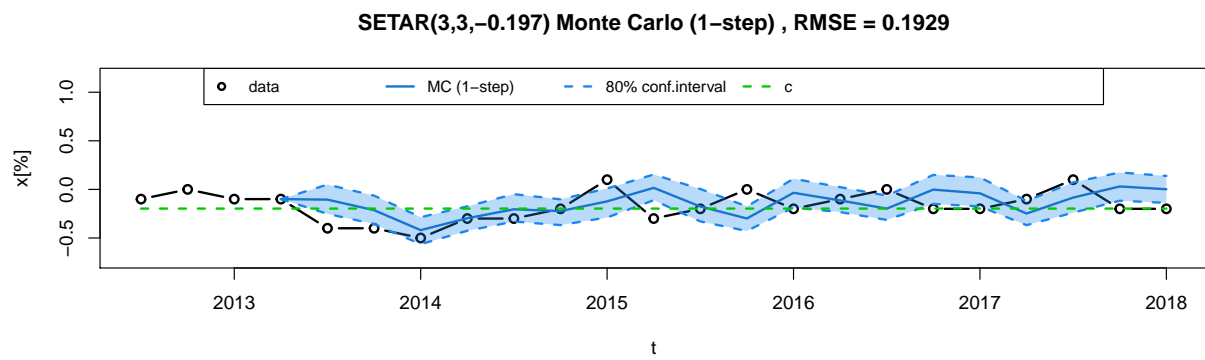
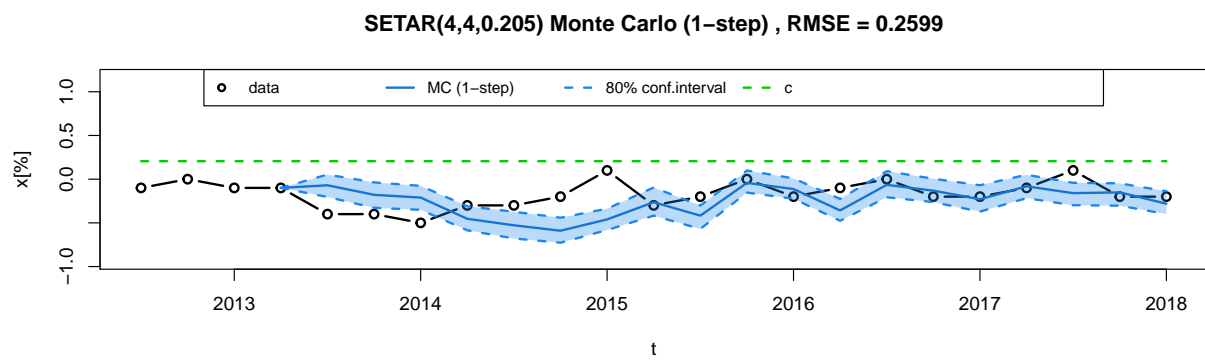
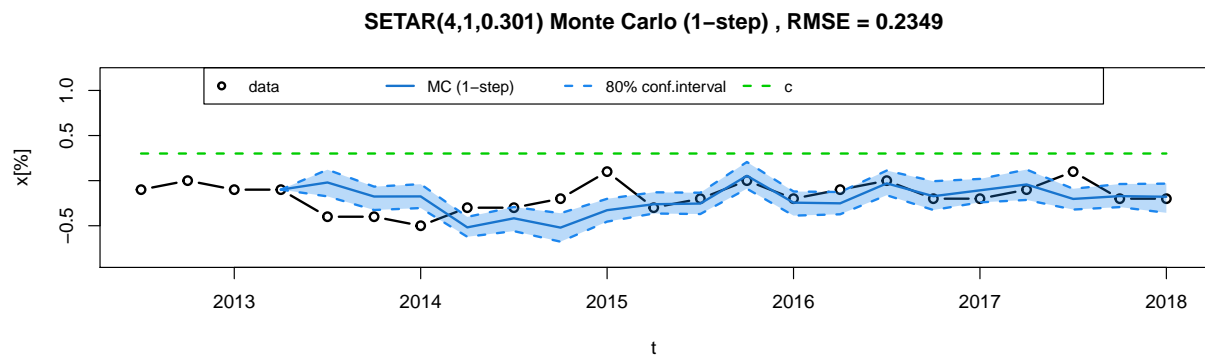
**SETAR(5,4,0.205) Naive (1-step) , RMSE = 0.27**



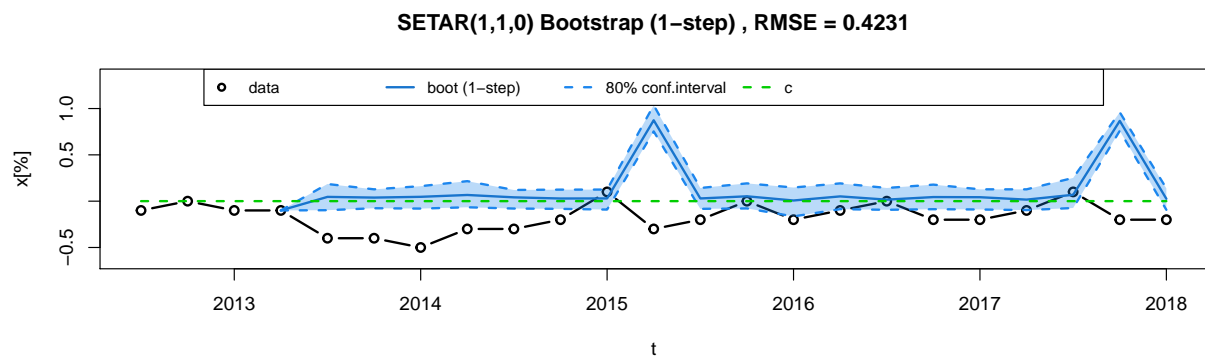
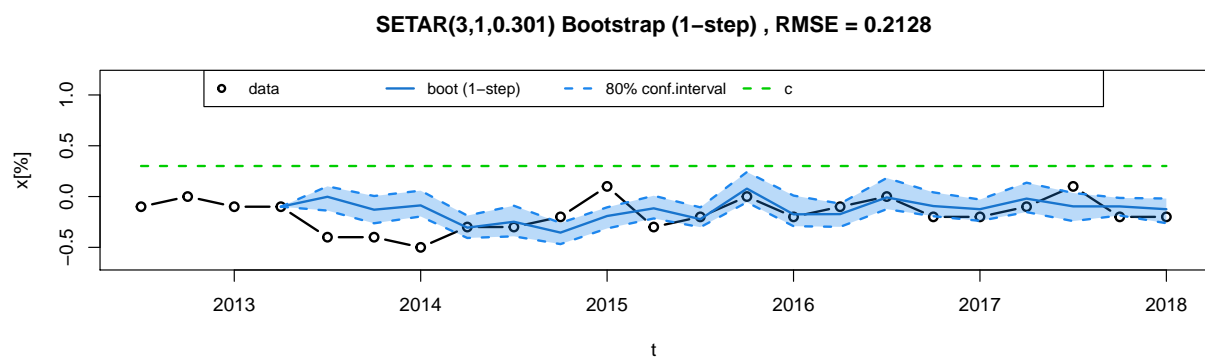
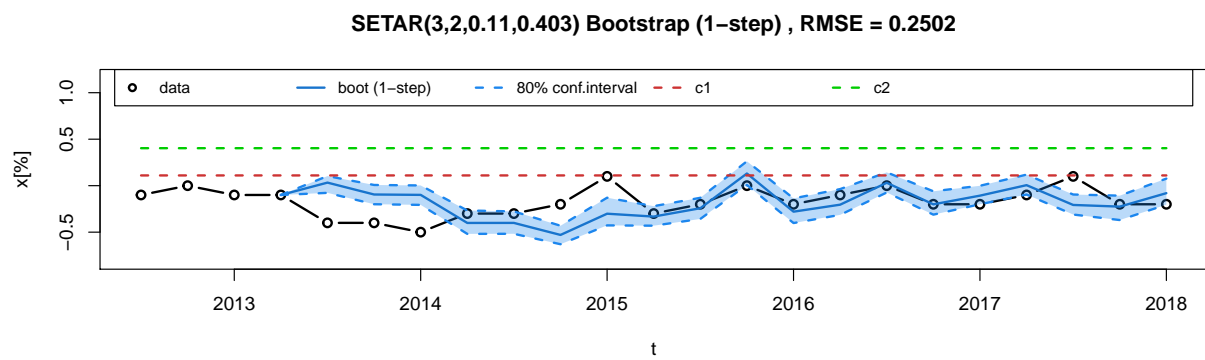
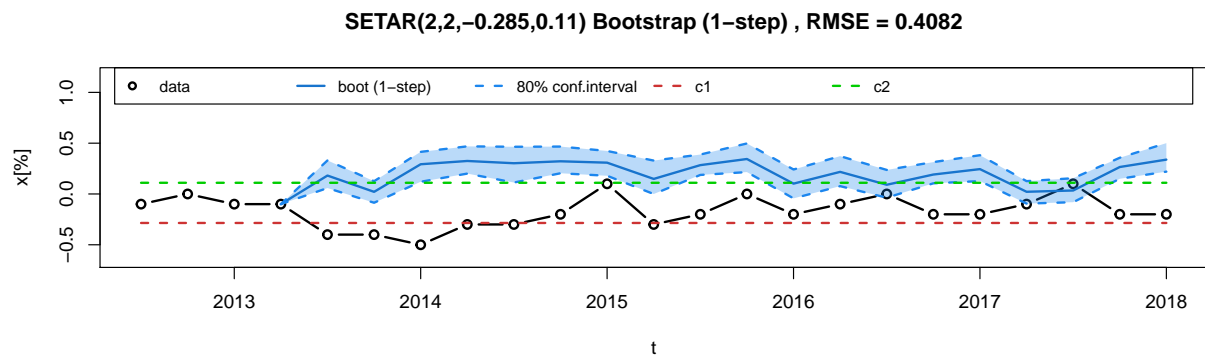
Monte Carlo:

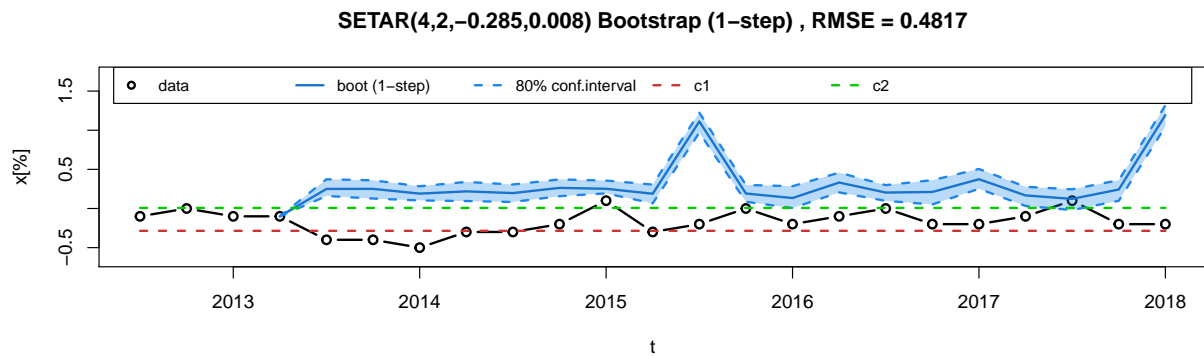
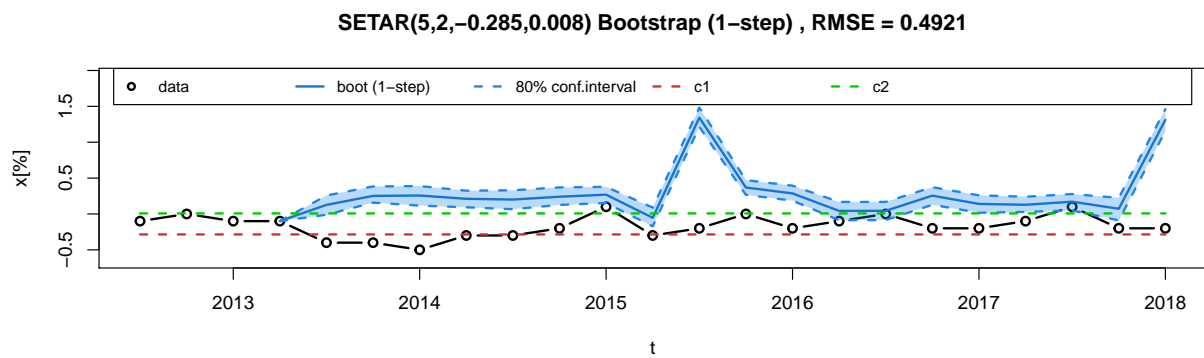
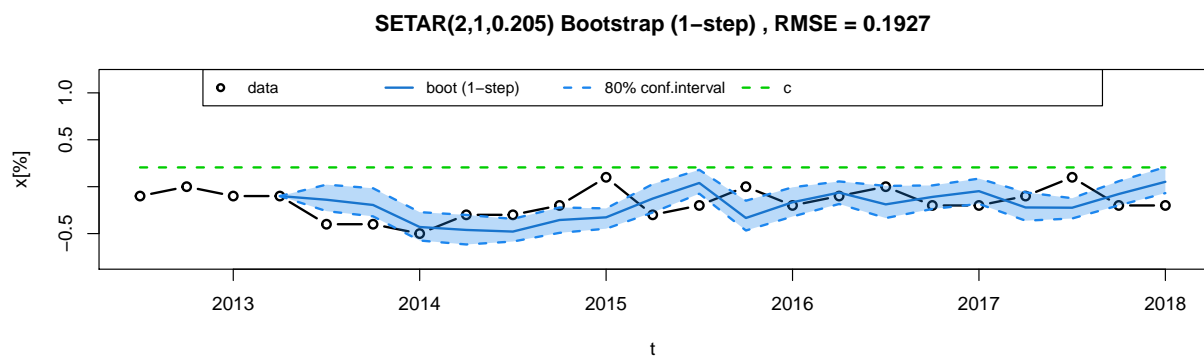
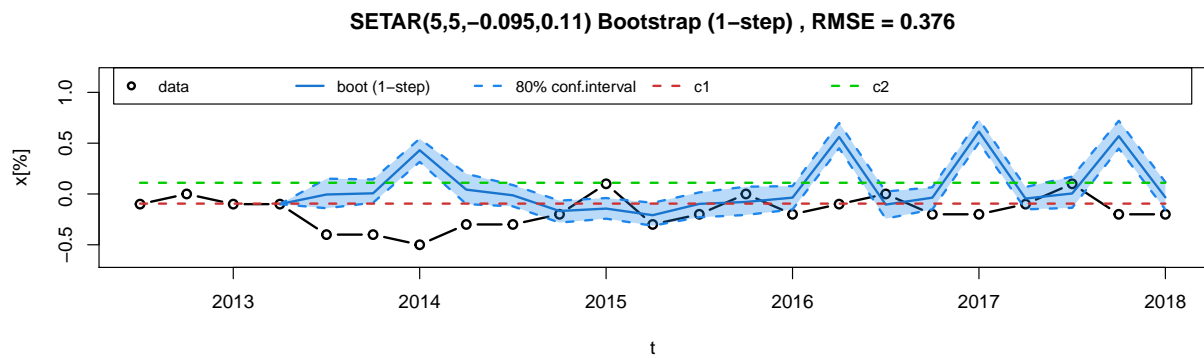




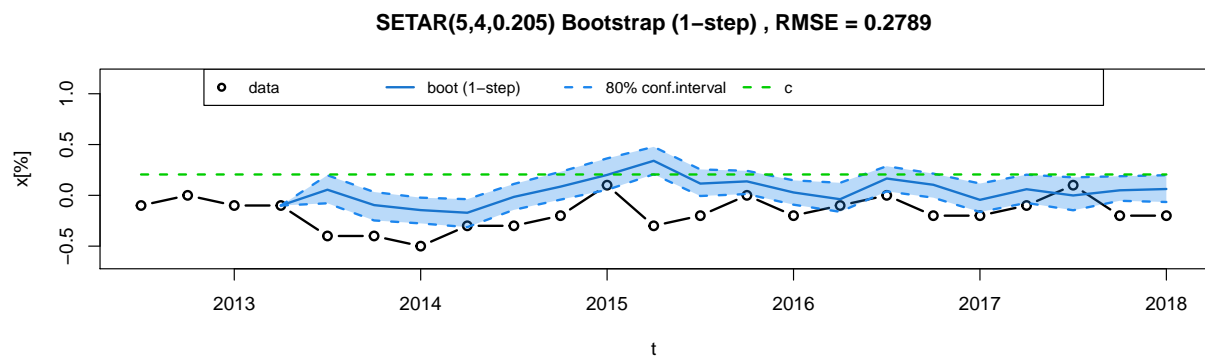
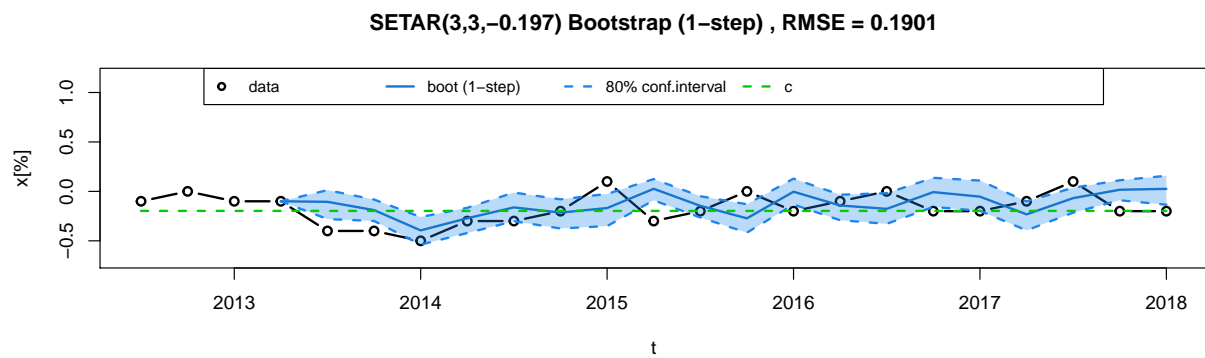
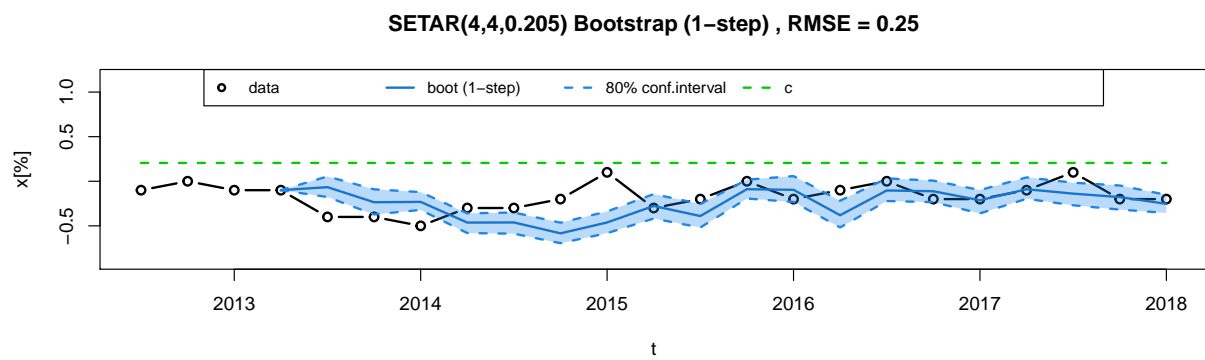
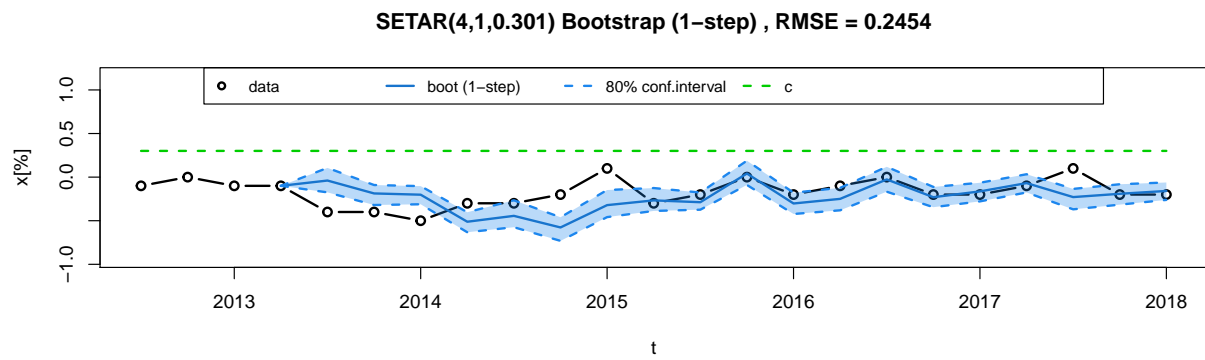


**Bootstrap:**



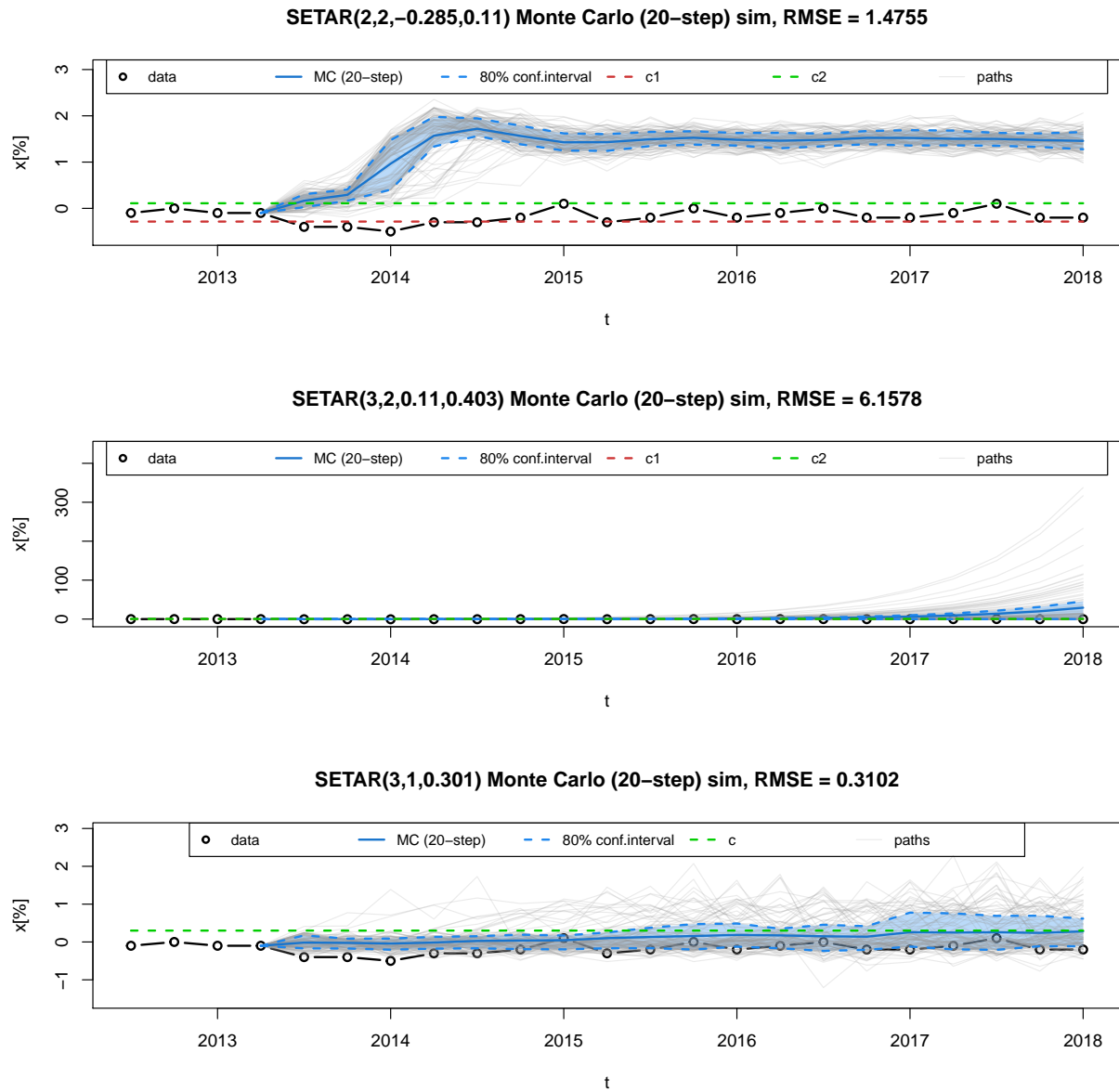




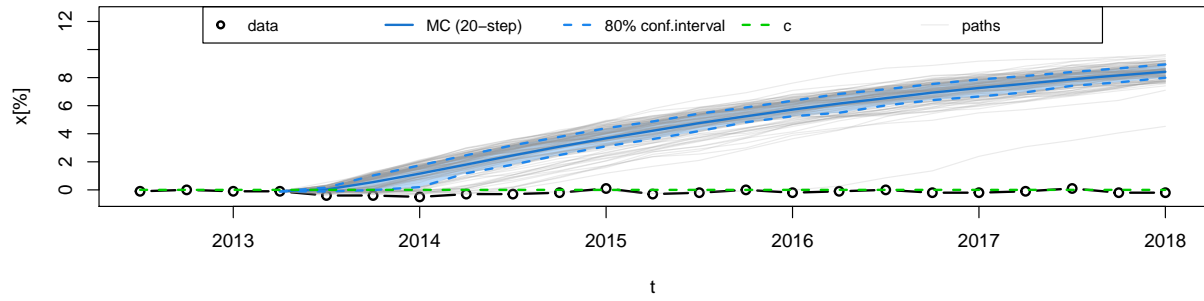


### 5.3 Multi-Step Predictions of SETAR Models

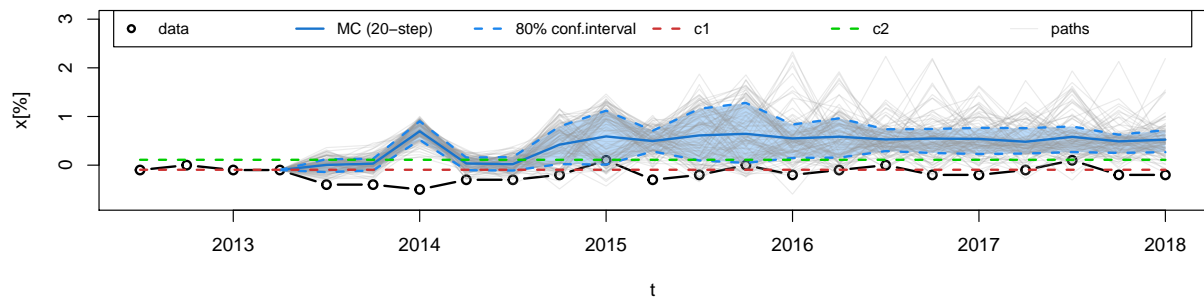
Monte Carlo:



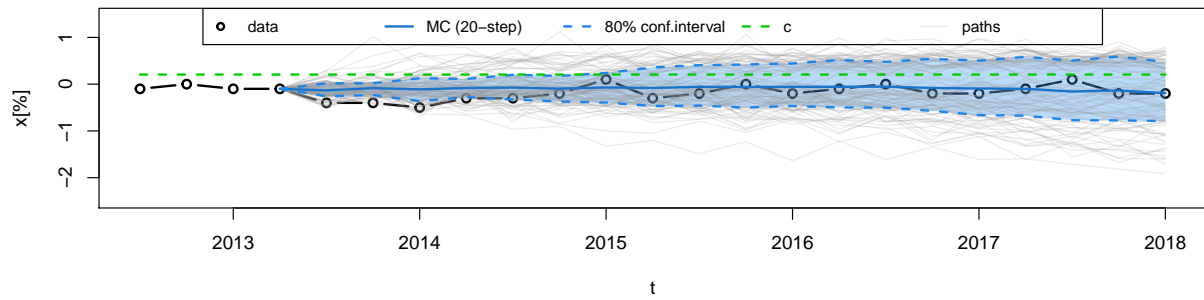
**SETAR(1,1,0) Monte Carlo (20-step) sim, RMSE = 5.1209**



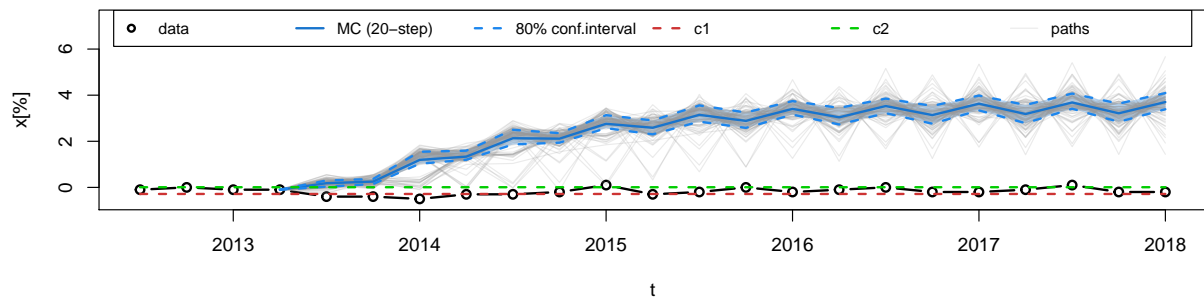
**SETAR(5,5,-0.095,0.11) Monte Carlo (20-step) sim, RMSE = 0.6184**



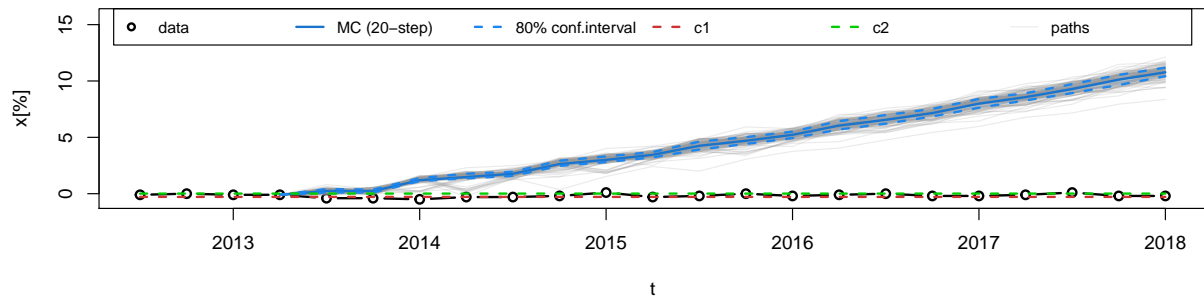
**SETAR(2,1,0.205) Monte Carlo (20-step) sim, RMSE = 0.1797**



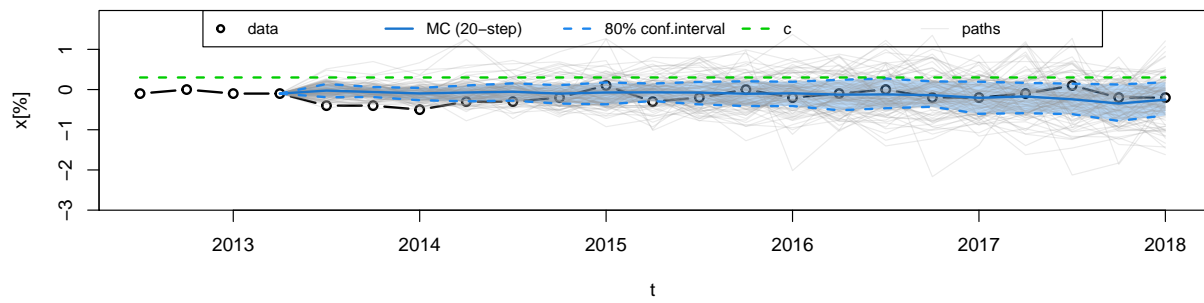
**SETAR(5,2,-0.285,0.008) Monte Carlo (20-step) sim, RMSE = 2.7261**



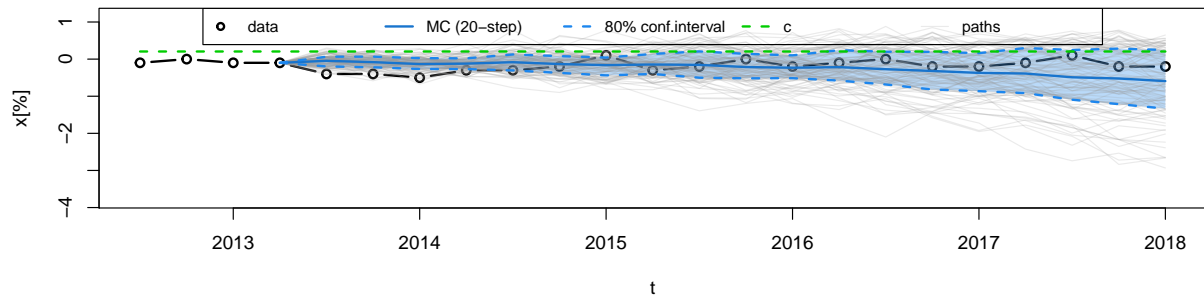
**SETAR(4,2,-0.285,0.008) Monte Carlo (20-step) sim, RMSE = 5.3906**



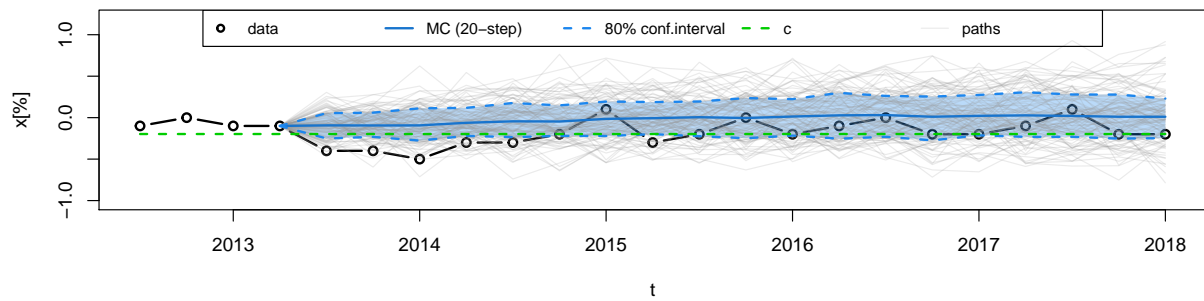
**SETAR(4,1,0.301) Monte Carlo (20-step) sim, RMSE = 0.1995**

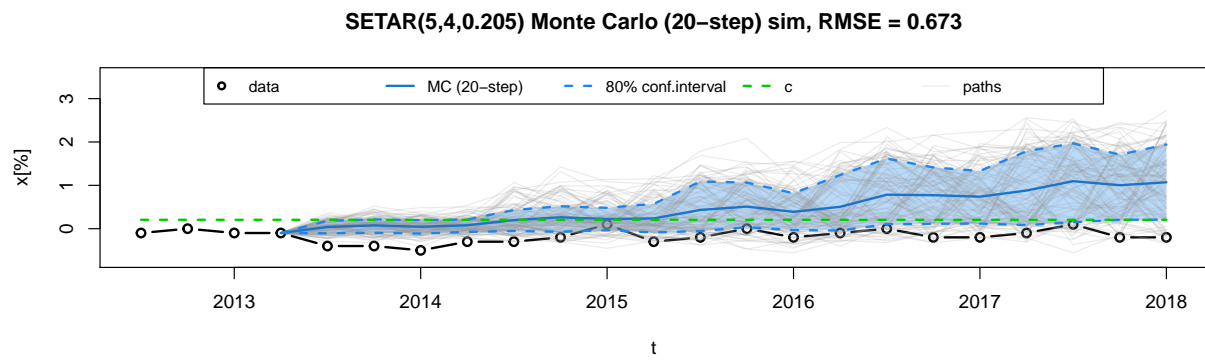


**SETAR(4,4,0.205) Monte Carlo (20-step) sim, RMSE = 0.2396**

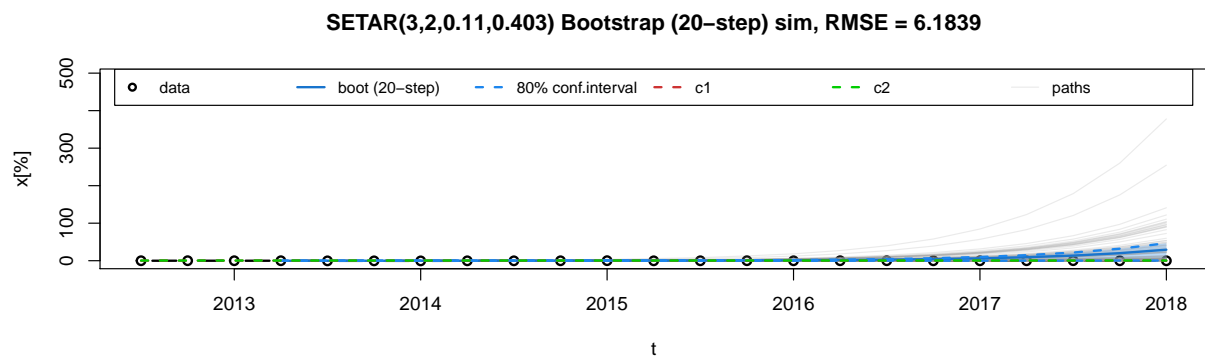
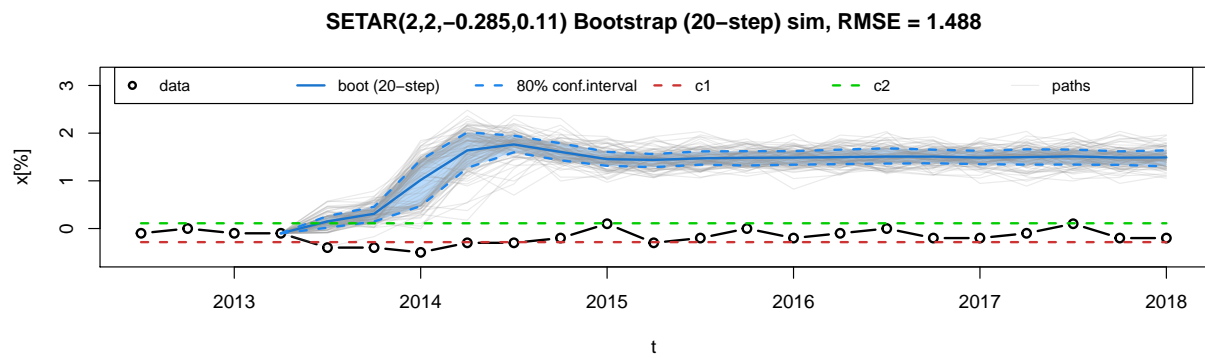


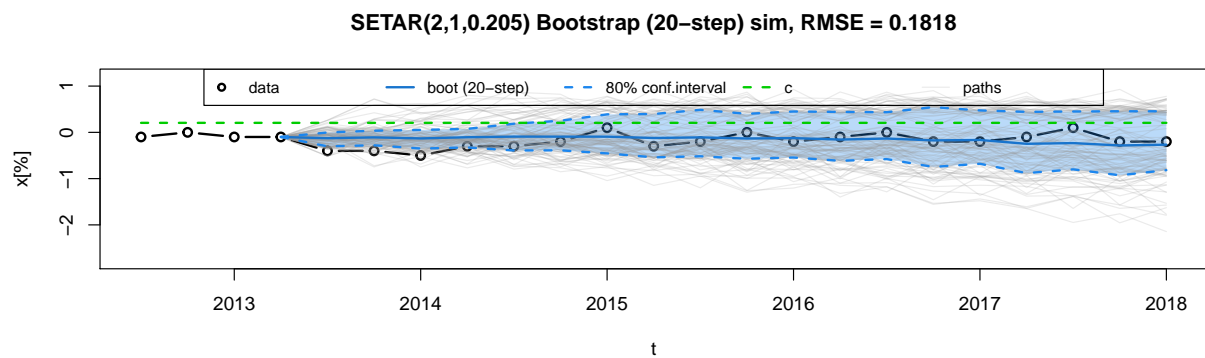
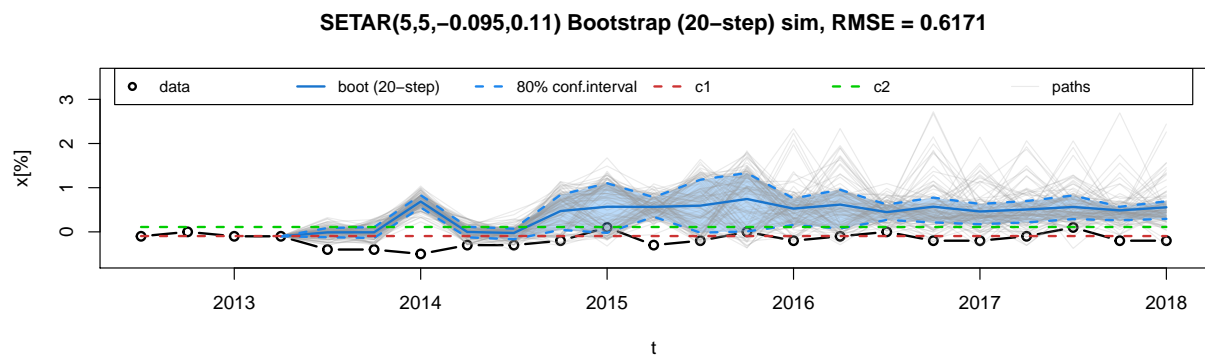
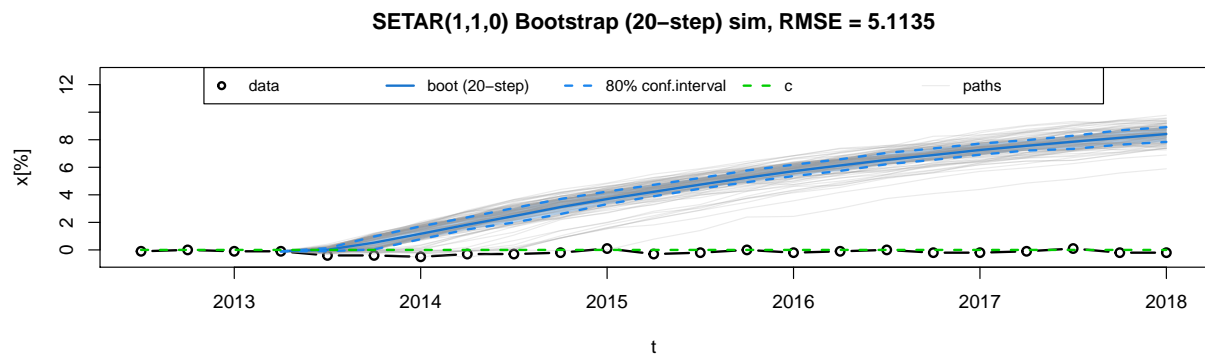
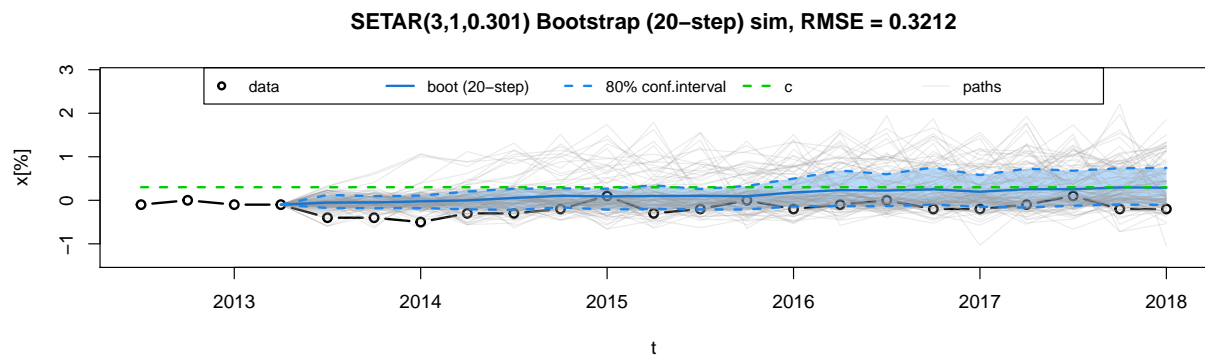
**SETAR(3,3,-0.197) Monte Carlo (20-step) sim, RMSE = 0.2092**



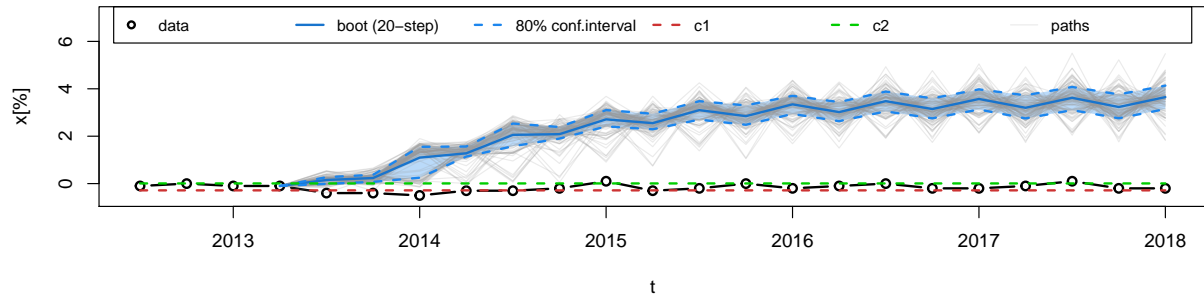


**Bootstrap:**

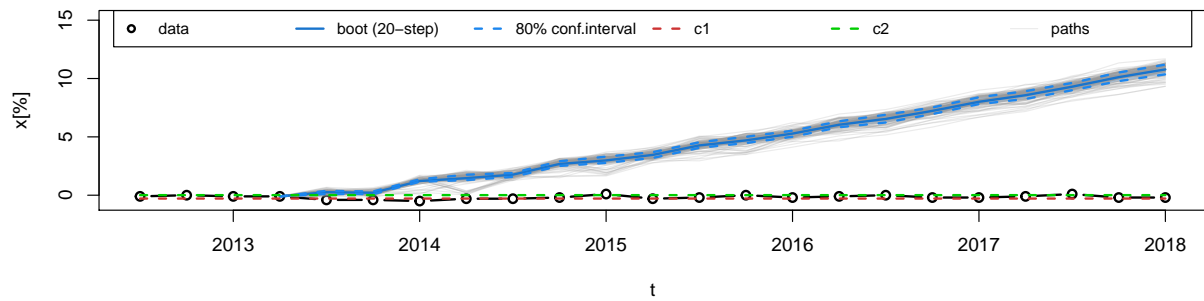




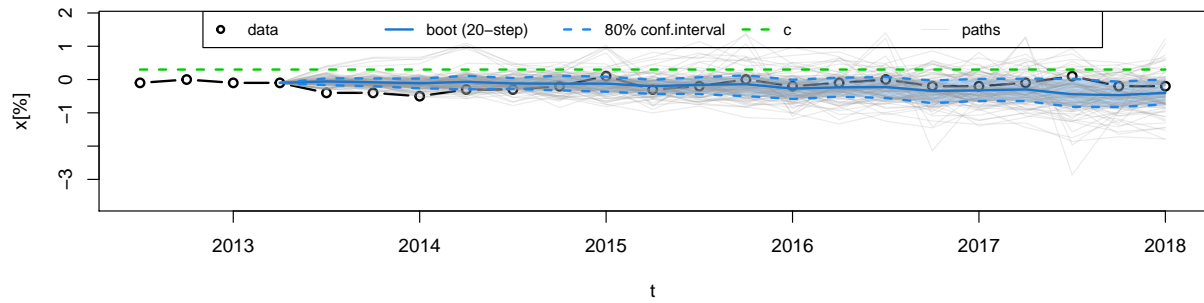
**SETAR(5,2,-0.285,0.008) Bootstrap (20-step) sim, RMSE = 2.693**



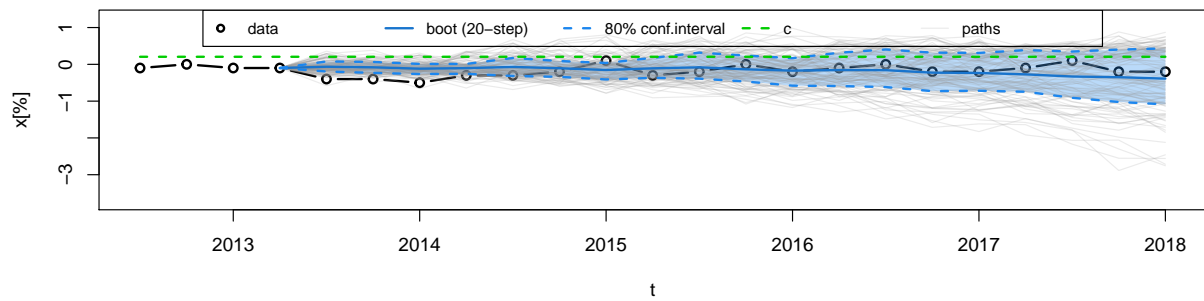
**SETAR(4,2,-0.285,0.008) Bootstrap (20-step) sim, RMSE = 5.4009**

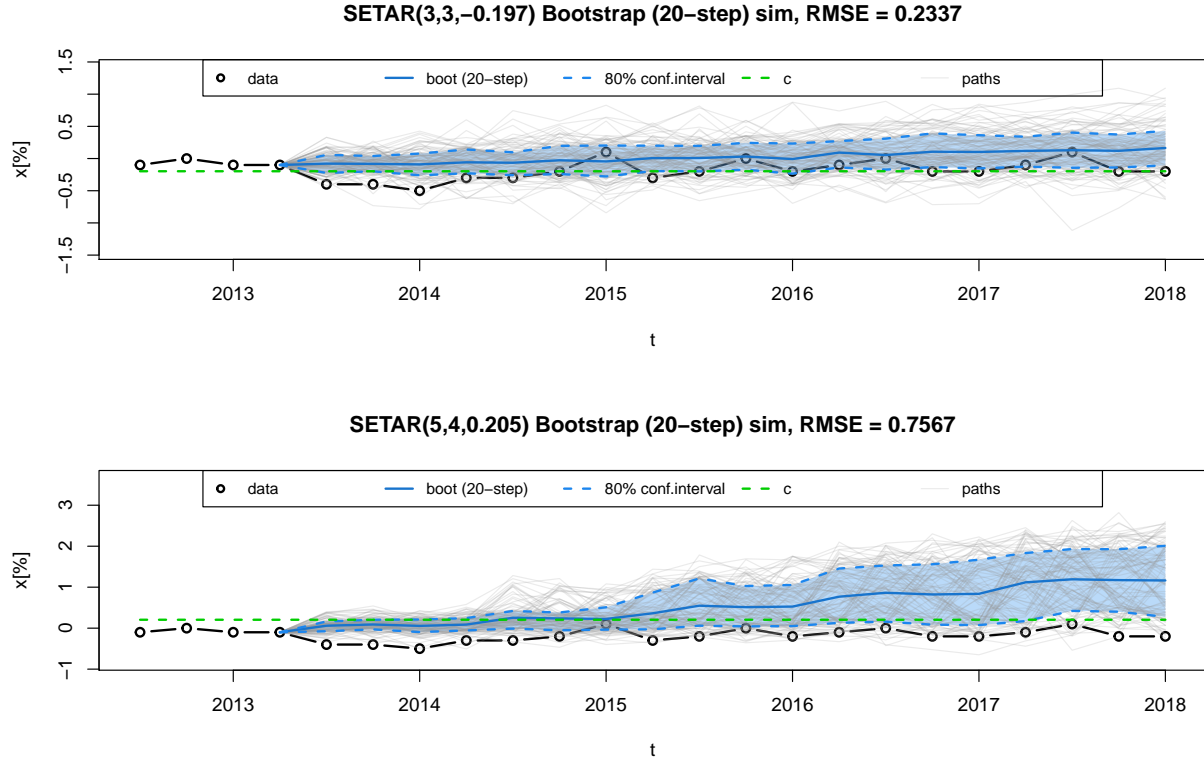


**SETAR(4,1,0.301) Bootstrap (20-step) sim, RMSE = 0.2274**



**SETAR(4,4,0.205) Bootstrap (20-step) sim, RMSE = 0.1972**





## 5.4 Evaluating Models

From what we observe, some models exhibit explosive properties, and have equilibria significantly distant from the evaluated data set, which, of course, impacts the resulting prediction error. Now we evaluate models according to their predictive properties, and filter out those that produce distant and/or explosive predictions.

##	model	BIC	rss	MSE_naive	MSE_mc	MSE_boot	divergence
## 1	SETAR(2,2,-0.285,0.11)	-590.602	0.0253	0.1643	0.1581	0.1687	! ! !
## 2	SETAR(3,2,0.11,0.403)	-588.83	0.0255	0.0604	0.0630	0.0619	
## 3	SETAR(3,1,0.301)	-586.083	0.0278	0.0435	0.0460	0.0419	
## 4	SETAR(1,1,0)	-585.638	0.0296	0.1790	0.1773	0.1808	! ! !
## 5	SETAR(5,5,-0.095,0.11)	-584.966	0.0248	0.1438	0.1444	0.1427	! ! !
## 6	SETAR(2,1,0.205)	-583.975	0.0291	0.0375	0.0366	0.0431	
## 7	SETAR(5,2,-0.285,0.008)	-581.469	0.0253	0.2467	0.2433	0.2448	! ! !
## 8	SETAR(4,2,-0.285,0.008)	-578.072	0.0258	0.2347	0.2364	0.2357	! ! !
## 9	SETAR(4,1,0.301)	-576.504	0.0280	0.0560	0.0566	0.0561	
## 10	SETAR(4,4,0.205)	-575.688	0.0270	0.0643	0.0665	0.0619	
## 11	SETAR(3,3,-0.197)	-574.277	0.0278	0.0350	0.0356	0.0327	
## 12	SETAR(5,4,0.205)	-571.548	0.0262	0.0729	0.0762	0.0804	

We notice that models that quickly diverge have significantly higher  $MSE$  than their  $\hat{\sigma}_\varepsilon^2$  (**rss**). One could then formulate a condition that model is “too divergent” if its training data  $MSE$  (**rss**) is significantly lower than prediction  $MSE$ . The significance factor could be taken, for example, as 4, i.e.: if the prediction  $MSE$  exceeds 4-multiple of training data **rss**, the prediction diverges. In the above table, we mark divergent predictions with !.

## 5.5 Conclusion and SETAR Evaluation

After considering all possible configurations of SETAR models, testing them for remaining nonlinearity, and evaluating their predictive properties, we conclude the following:



Some models placed within the top 3 best in their fit onto the training data (according to their *BIC*), contained SETAR3-type remaining nonlinearity, and thus had their thresholds estimated, turned out to have radically different equilibria than their 2-regime predecessors. Some were even explosive when we examined their predictive properties. Hence their training data fit quality cannot justify their mismatch between data and their properties as stochastic dynamical systems.

With regards to their predictive properties, we pick the best models according to their given *MSE* depending on the prediction method ("naive", "mc", "boot"):

```
## Models sorted by 1-step naive MSE:

##           model      BIC    rss MSE_naive
## 11  SETAR(3,3,-0.197) -574.277 0.0278    0.0350
## 6   SETAR(2,1,0.205) -583.975 0.0291    0.0375
## 3   SETAR(3,1,0.301) -586.083 0.0278    0.0435
## 9   SETAR(4,1,0.301) -576.504 0.0280    0.0560
## 2   SETAR(3,2,0.11,0.403) -588.83 0.0255    0.0604
## 10  SETAR(4,4,0.205) -575.688 0.0270    0.0643
## 12  SETAR(5,4,0.205) -571.548 0.0262    0.0729
```

```
## Models sorted by 1-step Monte Carlo MSE:

##           model      BIC    rss MSE_mc
## 11  SETAR(3,3,-0.197) -574.277 0.0278 0.0356
## 6   SETAR(2,1,0.205) -583.975 0.0291 0.0366
## 3   SETAR(3,1,0.301) -586.083 0.0278 0.0460
## 9   SETAR(4,1,0.301) -576.504 0.0280 0.0566
## 2   SETAR(3,2,0.11,0.403) -588.83 0.0255 0.0630
## 10  SETAR(4,4,0.205) -575.688 0.0270 0.0665
## 12  SETAR(5,4,0.205) -571.548 0.0262 0.0762
```

```
## Models sorted by 1-step Bootstrap MSE:

##           model      BIC    rss MSE_boot
## 11  SETAR(3,3,-0.197) -574.277 0.0278    0.0327
## 6   SETAR(2,1,0.205) -583.975 0.0291    0.0431
## 3   SETAR(3,1,0.301) -586.083 0.0278    0.0419
## 9   SETAR(4,1,0.301) -576.504 0.0280    0.0561
## 2   SETAR(3,2,0.11,0.403) -588.83 0.0255    0.0619
## 10  SETAR(4,4,0.205) -575.688 0.0270    0.0619
## 12  SETAR(5,4,0.205) -571.548 0.0262    0.0804
```

and as we observe, the order remains the same regardless of the prediction method.

## 6. Tests for Non-Linearity of STAR Models