

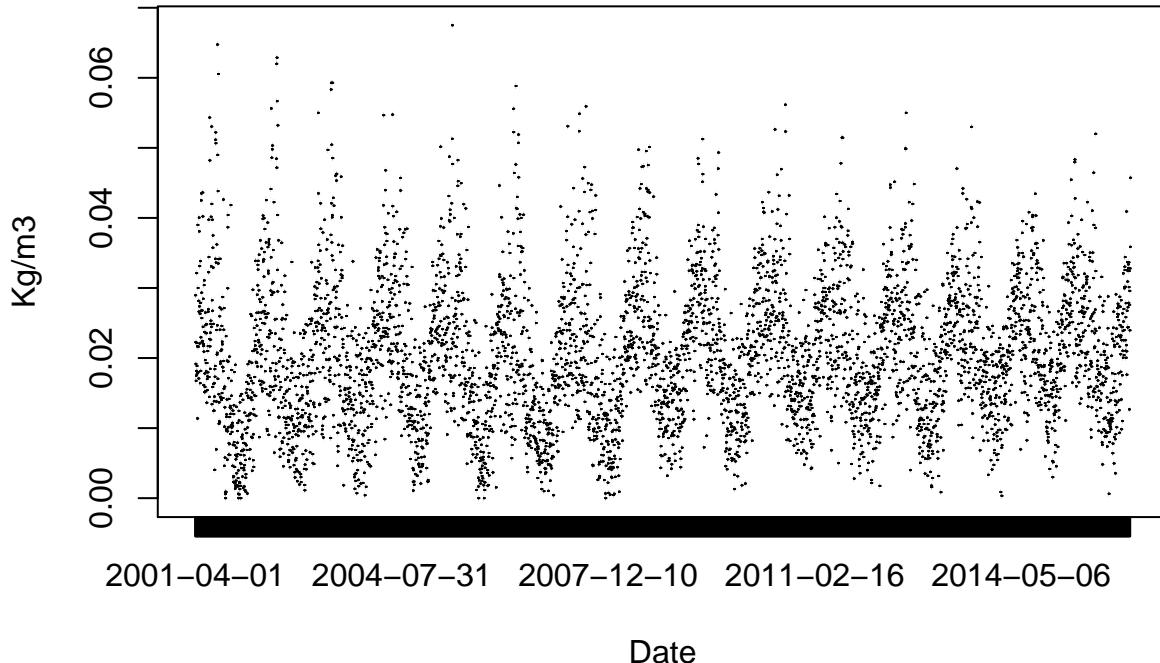
# Prediction of ozone level in Boston

Lucas Emanuel Resck Domingues\*

Lucas Machado Moschen†

Load and visualize

**Daily average level of O<sub>3</sub> in Boston**



## Methodology

1. Data treatment: analyse missing and duplicated data.
2. Compare the daily models: we use rollapply using 2 years + 6 days to predict the day, that is, we use  $\{X_t\}_{t=2 \cdot 365}^t$  to predict  $X_{t+7}$ . We calcute, for each  $t > 2 \cdot 365 + 6$ ,  $|\hat{X}_t - X_t|$  and after calculate  $\sum_t |\hat{X}_t - X_t|$ . We used Decomposition, Regression, Holt-Winters and ARMA.
3. Compare the weekly models: we use rollapply using 2 years + 1 week to predict the week, that is, we use  $\{X_t\}_{t=2 \cdot 52}^t$  to predict  $X_{t+1}$ . We calcute, for each  $t > 2 \cdot 52 + 1$ ,  $|\hat{X}_t - X_t|$  and after calculate  $\sum_t |\hat{X}_t - X_t|$ . We used Decomposition, Regression, Holt-Winters and ARMA.
4. Compare the best models for the 2 cases in the test data using the same approach.

## Data treatment

We noticed that some days do not exist in the dataset, for example, the day August 31, 2001 does not have information in the dataset.

```
##      X   City           State Site.Num Date.Local  O3.Mean
## 148 148 Boston Massachusetts        42 2001-08-28 0.024583
## 149 149 Boston Massachusetts        42 2001-08-29 0.015000
```

\*Escola de Matemática Aplicada

†Escola de Matemática Aplicada

```

## 150 150 Boston Massachusetts      42 2001-08-30 0.022333
## 151 151 Boston Massachusetts      42 2001-09-01 0.021958
## 152 152 Boston Massachusetts      42 2001-09-02 0.018750
## 153 153 Boston Massachusetts      42 2001-09-03 0.028708

```

Also, there is duplicated days, as June 9, 2002:

```

##      X   City       State Site.Num Date.Local 03.Mean
## 412 412 Boston Massachusetts      42 2002-06-08 0.022917
## 413 413 Boston Massachusetts      42 2002-06-09 0.036190
## 414 414 Boston Massachusetts      42 2002-06-09 0.037000
## 415 415 Boston Massachusetts      42 2002-06-10 0.023389

```

The duplicated one is easier to deal, but the missing values are harder. First we calculate the mean value between the duplicated.

The rate of missing values is almost 5% of the dataset.

```
## [1] 0.04453367
```

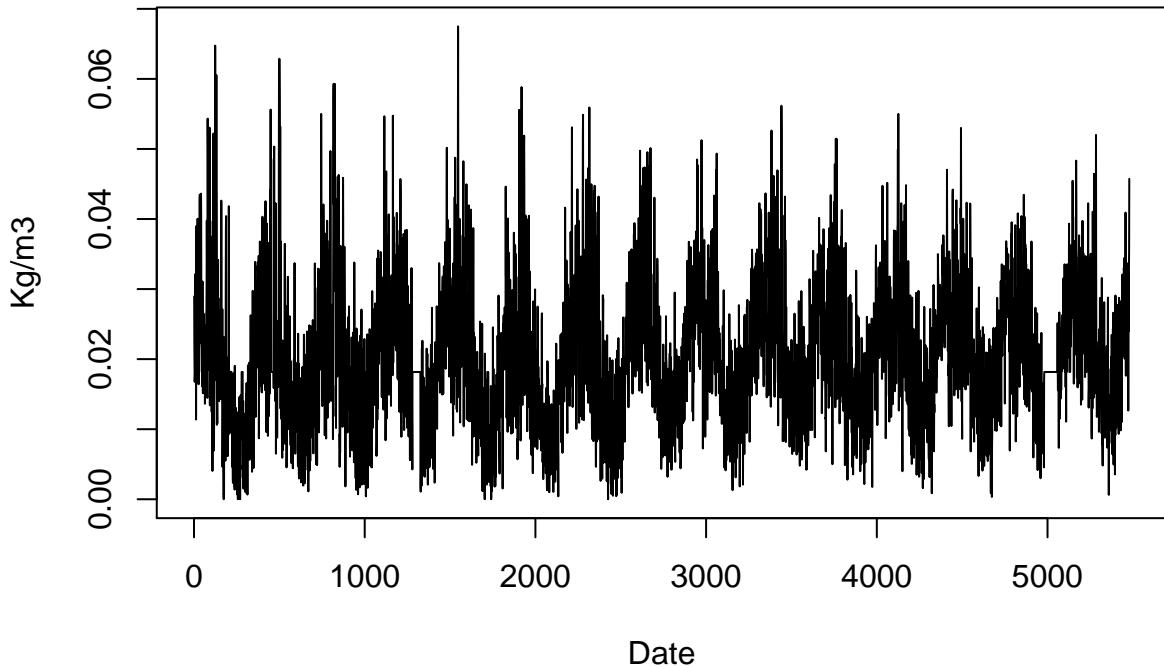
So as to solve that problem, we make a knn imputation using the month ( $k = 30$ )

```

o3.clean <- knn.impute(as.matrix(o3.ts), k = 30)
o3.clean <- as.ts(o3.clean)
plot(o3.clean, main = 'Daily average level of O3 in Boston (after imputation)',
     xlab = 'Date', ylab = 'Kg/m3')

```

## Daily average level of O3 in Boston (after imputation)



We also separate our data between training and test because of modelling best practices.

```

o3_train = o3.clean[1:(length(o3.clean)[1] - 365),]
o3_test = o3.clean[(length(o3.clean)[1] - 365 + 1):length(o3.clean)[1],]

```

## Models: case 1

Now we develop some models using the train data.

The metric to compare is the Mean Absolute Error (MAE) in the predictions:

```
mae <- function(ytrue, ypred)
{
  return(mean(abs(ytrue - ypred)))
}
```

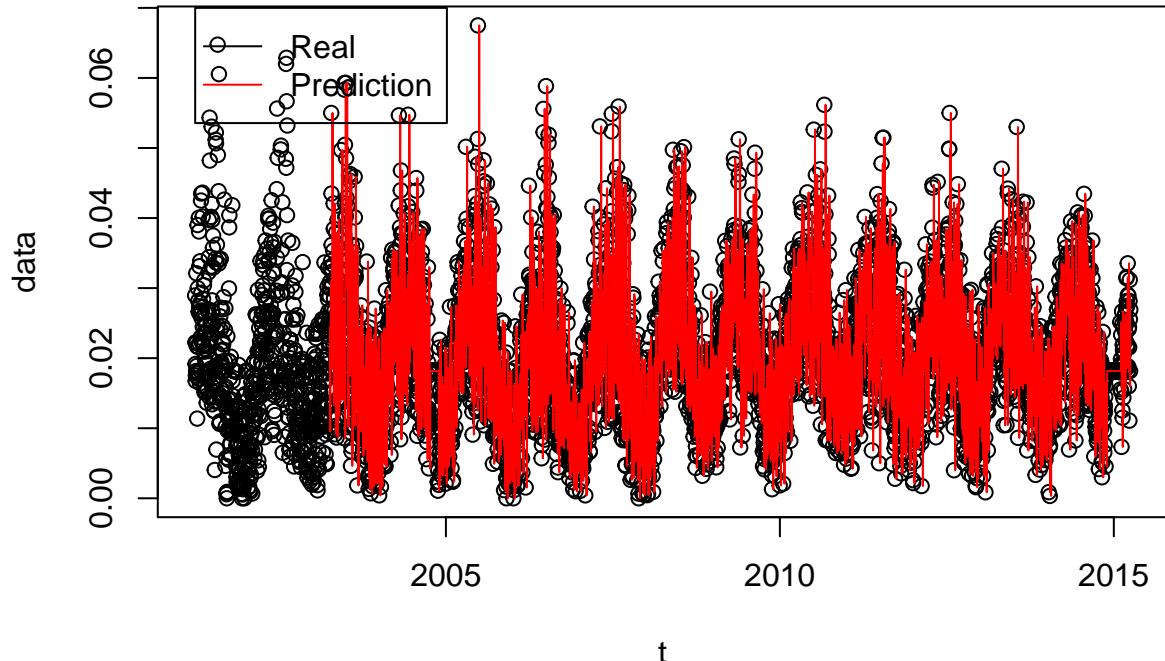
We will use `rollapply` in order to calculate the error, considering the last two years to predict one week forward.

### Baseline Model

We will do the naive forecast to the baseline model.

```
## [1] 0.007844892
```

### Baseline model prediction



### Decompose

First of all we make a seasonality test using Kruskal-Wallis. Actually it tests whether samples originate from the same distribution. We can organize it to be samples for each corresponding day. We compare two different frequencies: monthly and yearly. The second one showed the smallest p-value, in particular less than 0.05. For that reason, we will use 365 as the seasonality.

```
##
## Kruskal-Wallis rank sum test
##
## data: o3_train and g
## Kruskal-Wallis chi-squared = 32.114, df = 30, p-value = 0.3622
```

```

## Kruskal-Wallis rank sum test
##
## data: o3_train and g
## Kruskal-Wallis chi-squared = 2179.6, df = 364, p-value < 2.2e-16

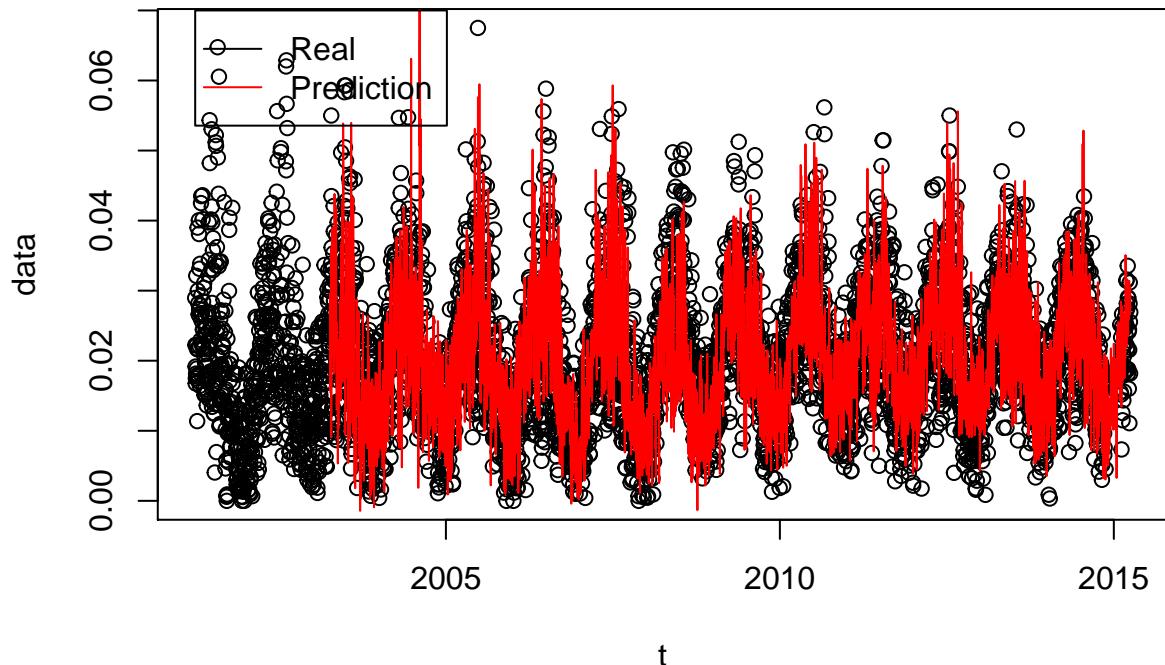
```

### Additive model

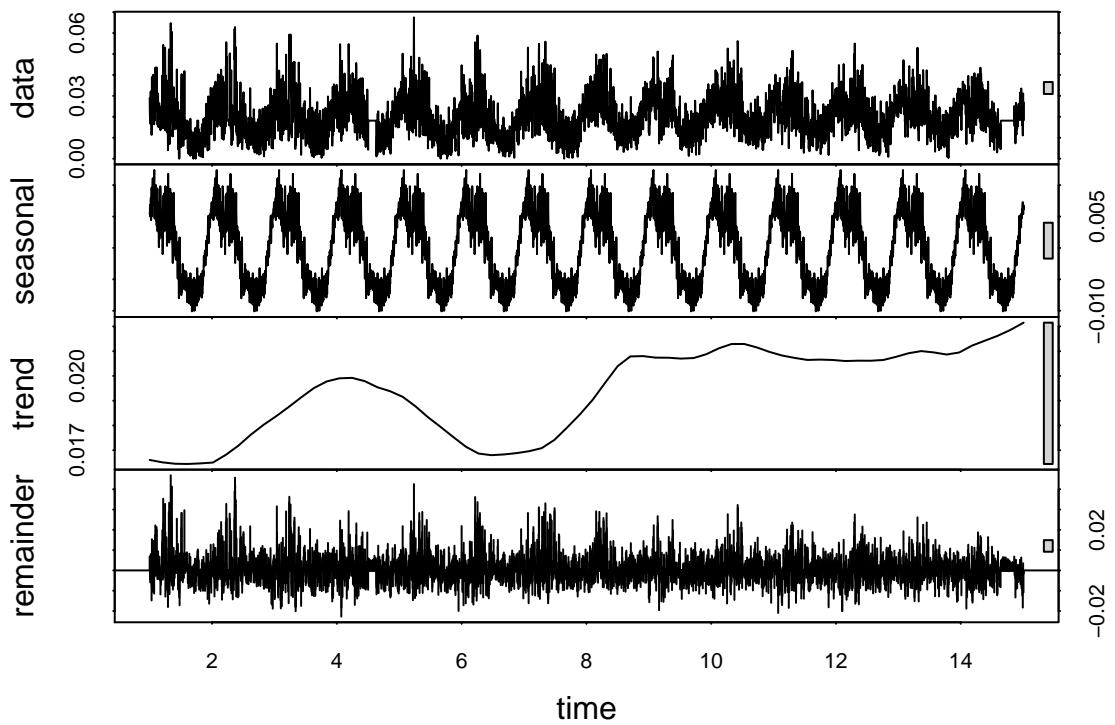
First we analyse the MAE.

```
## [1] 0.007832648
```

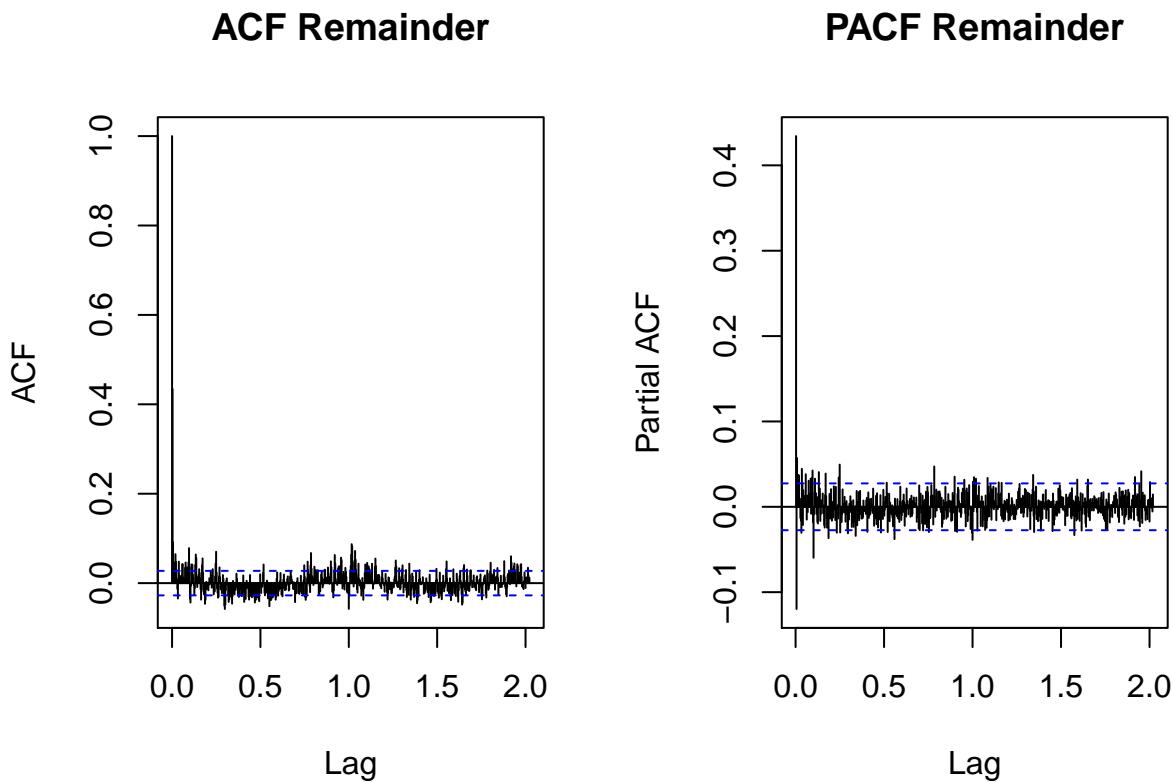
### Additive decompose prediction



We also can fit the model using `t.window` and analyse the reminder of the method.



The ACF and the PACF of the reminder:



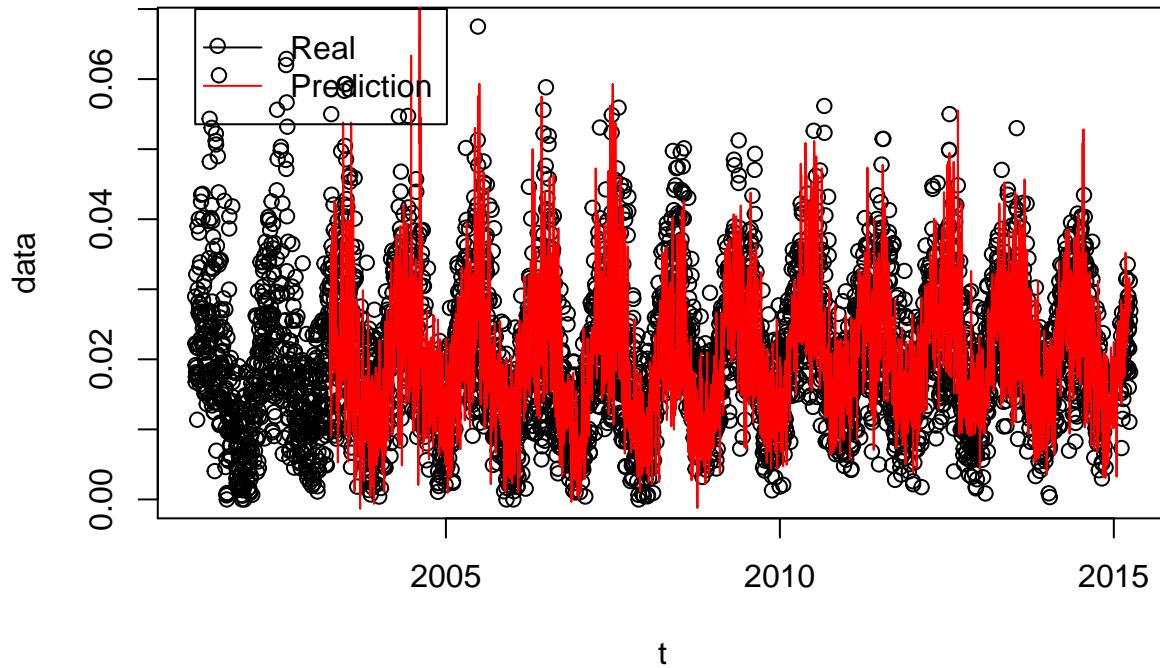
We see that there are a big spike when lag = 365. It seems not so good for a reminder. We could fit an ARMA model in this reminder yet.

**Multiplicative model**

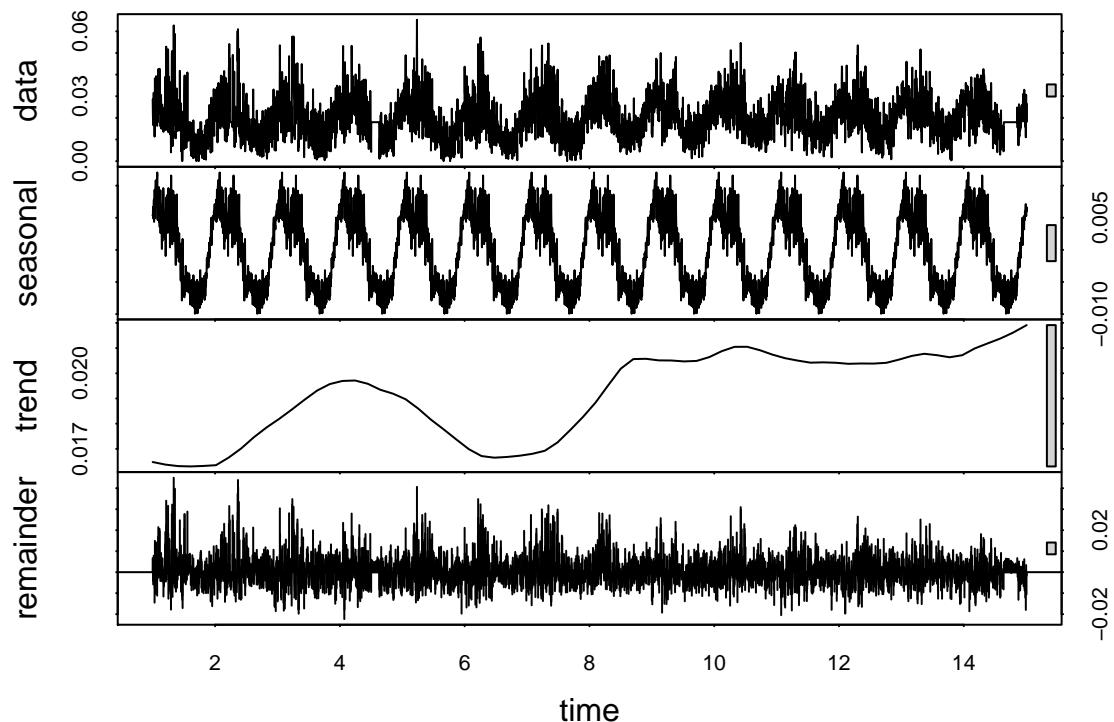
First we analyse the MAE.

```
## [1] 0.007822273
```

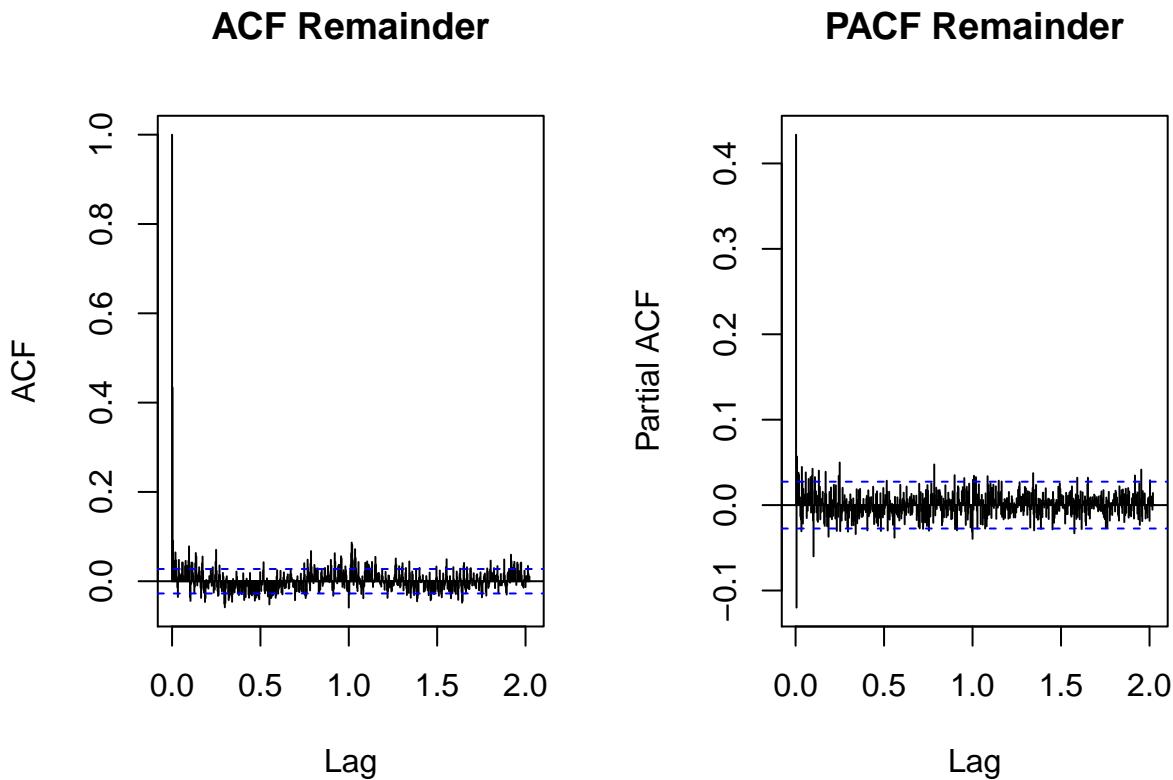
### Multiplicative decompose prediction



We also can fit the model using `t.window` and analyse the reminder of the method.



The ACF and the PACF of the remainder:



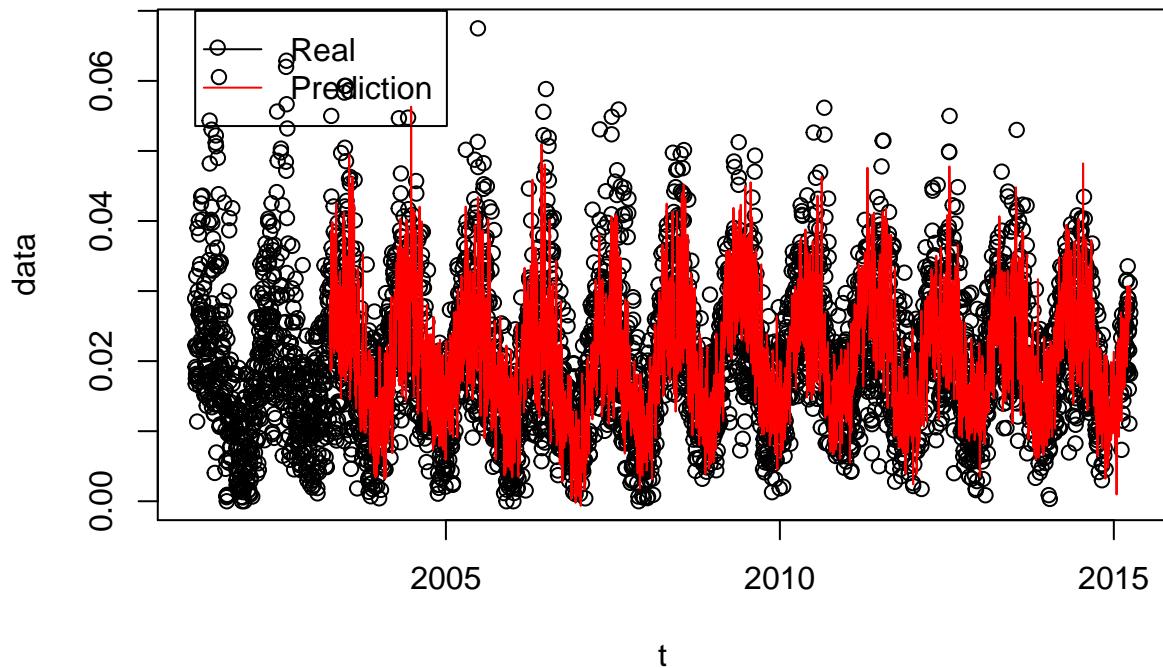
We see that there are a big spike when lag = 365. It seems not so good for a reminder. We could fit an ARMA model in this reminder yet. The same problem as before.

### Regression

We tested for seasonalities, and we settled with 365. So we will fit a regression model with the seasonality dummies. We see the MAE:

```
## [1] 0.007423855
```

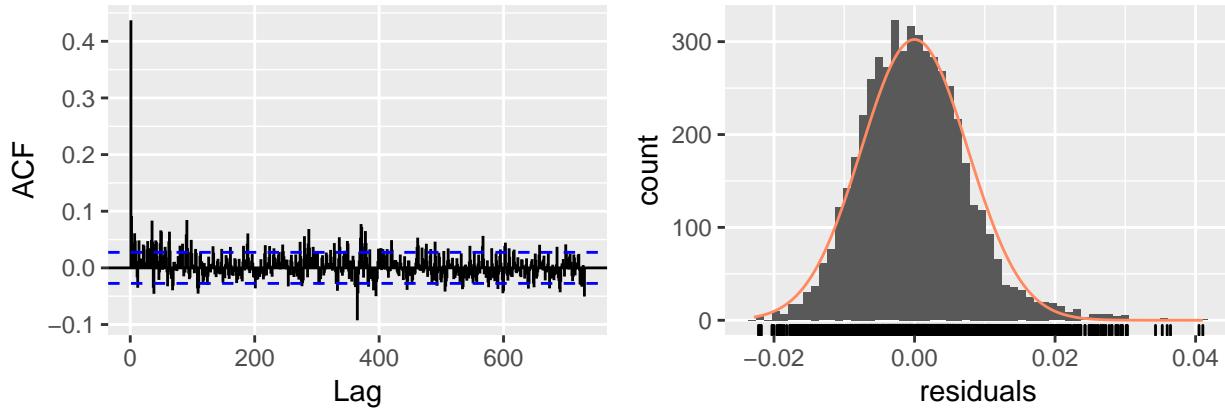
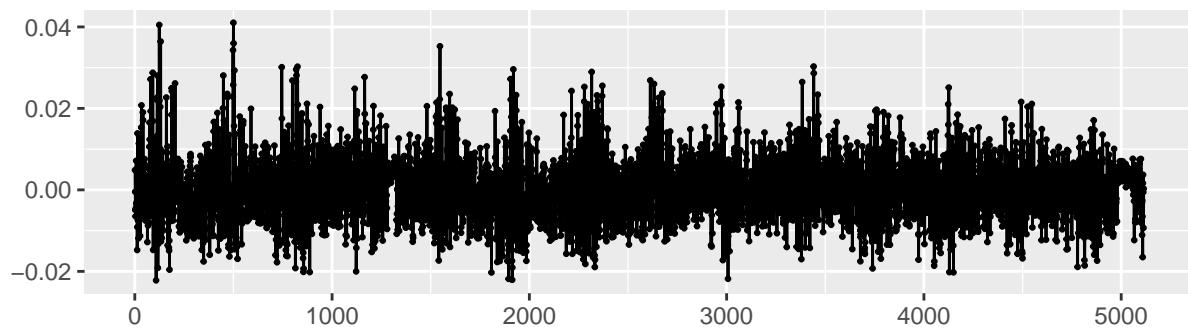
## Regression prediction



Now we analyse the residuals. We fit a LM model in all training data in order to analyse it.

```
train = data.frame(  
  t = t,  
  o3_train = o3_train,  
  Q = Q  
)  
mod = lm(o3_train~t+Q, data = train)  
checkresiduals(mod, lag = 2*freq, lag.max = 2*freq)
```

## Residuals



```
##  
## Breusch-Godfrey test for serial correlation of order up to 730  
##  
## data: Residuals  
## LM test = 1747.3, df = 730, p-value < 2.2e-16
```

As before, we see spikes in lag = 365, 730. We expect a WN to not have this.

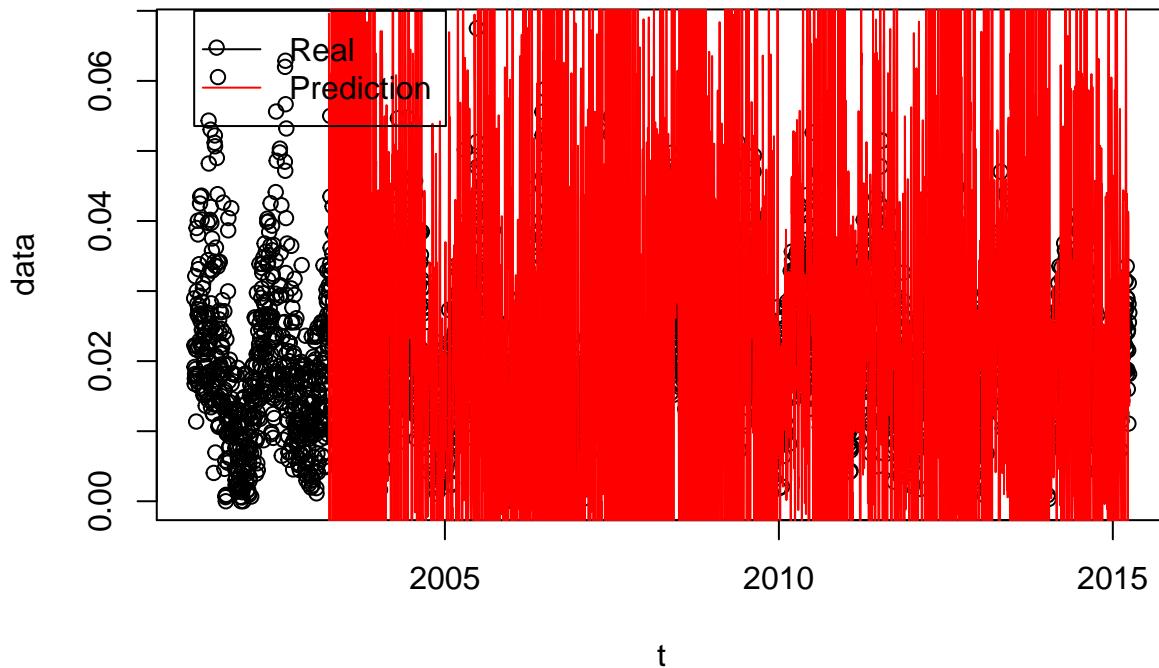
## Holt-Winters

Now we will try Holt-Winters models. In fact, because of apparently seasonality, we will consider complete Holt-Winters models, both additive and multiplicative.

### Additive

```
## [1] 0.03434545
```

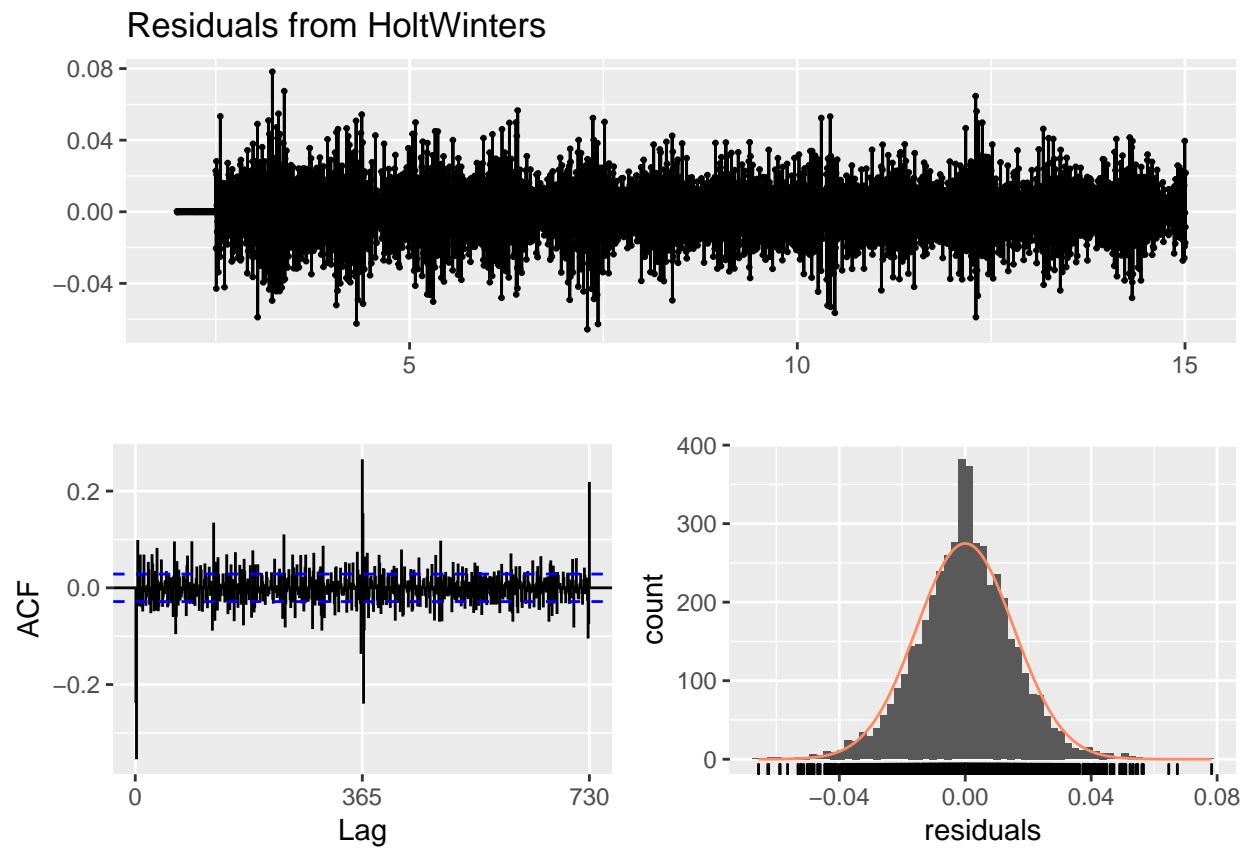
## Additive Holt–Winters prediction



MAE is not so good, nor is the resulting graph. HW tends to “explode”. We have seem better results. Let's analyse the residuals:

```
mod = HoltWinters(ts(o3_train, frequency = 365), beta = T, gamma = T, seasonal = "additive")
checkresiduals(mod, lag = 2*freq, lag.max = 2*freq)

## Warning in modeldf.default(object): Could not find appropriate degrees of
## freedom for this model.
```

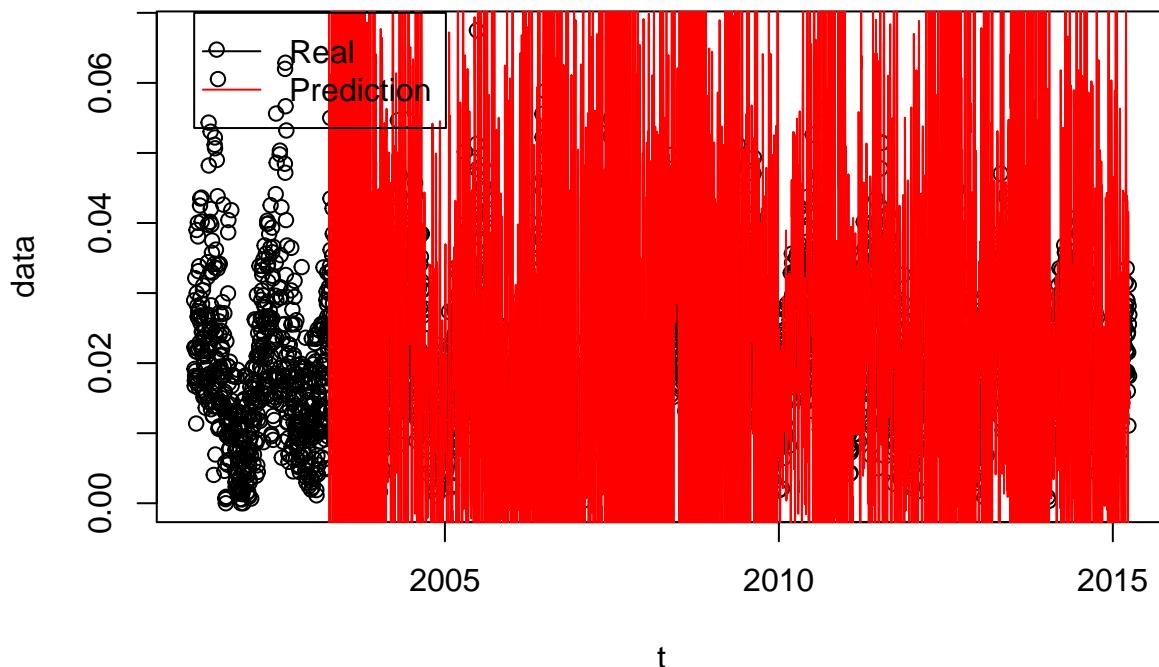


We see the same problems as before: high correlated lag = 365, 730, evidence of this not being a WN.

#### Multiplicative

```
## [1] 0.03427301
```

## Multiplicative Holt–Winters prediction

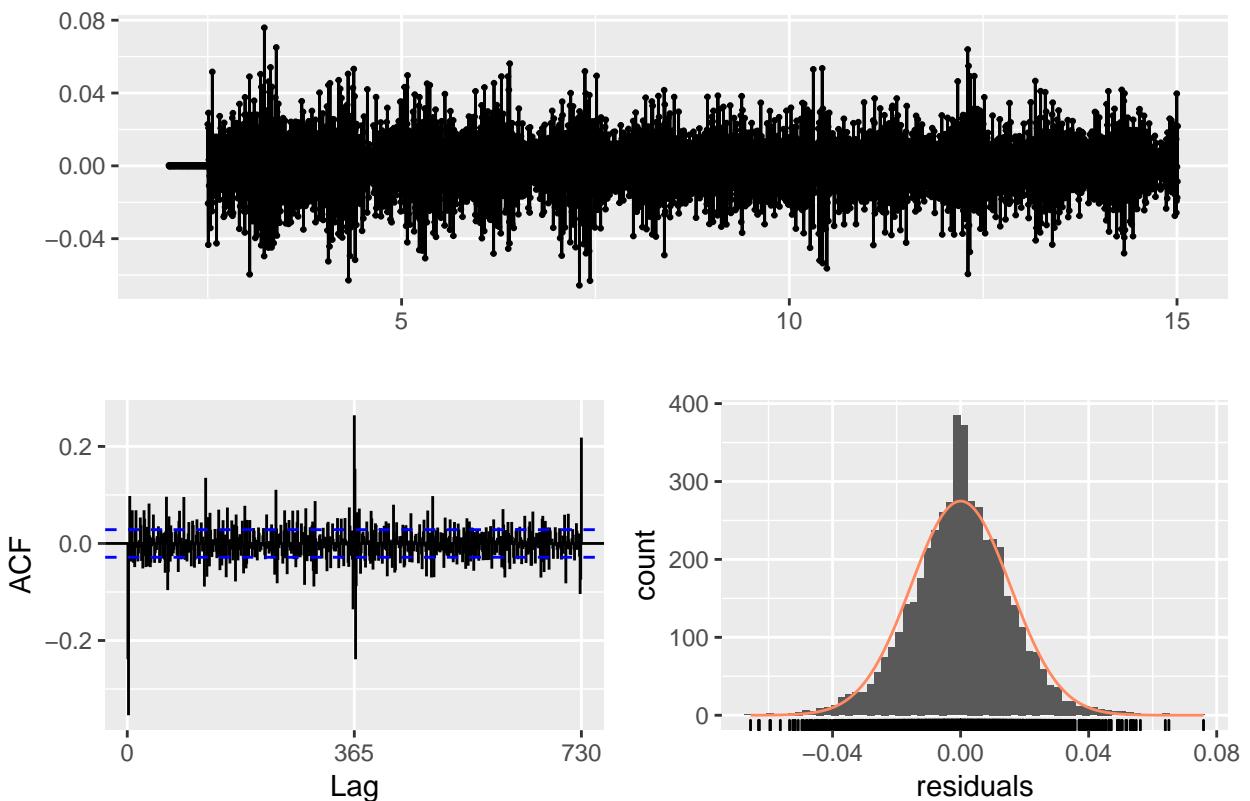


We see also a not so good MAE and graph. Let's analyse the residuals of an model fitted in all training data:

```
mod = HoltWinters(ts(o3_train+1, frequency = 365), beta = T, gamma = T, seasonal = "multiplicative")
checkresiduals(mod, lag = 2*freq, lag.max = 2*freq)

## Warning in modelfdf.default(object): Could not find appropriate degrees of
## freedom for this model.
```

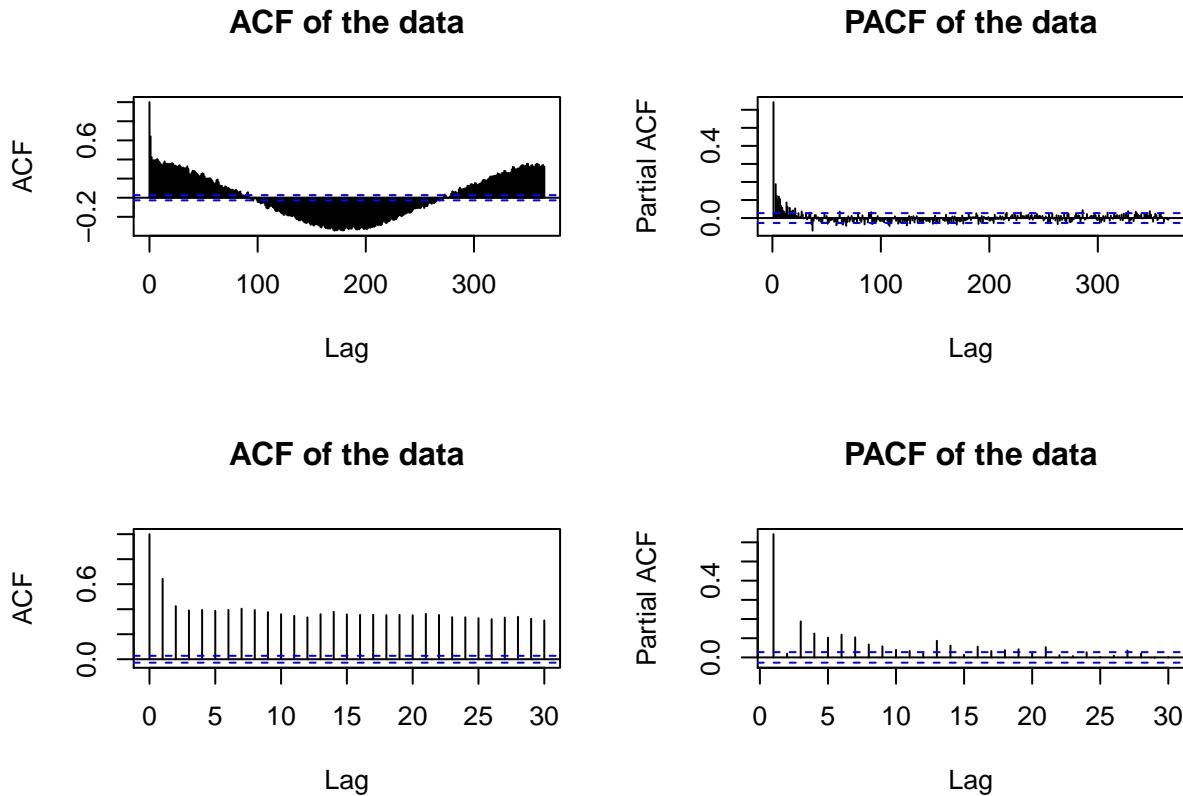
### Residuals from HoltWinters



Same problems as before. Holt-Winters doesn't appear to work well in this case.

### ARMA

We can see the ACF and PACF:



Based on these graphs, we see both graphs has a exponentially decay, the first after the  $q - p = -1$  or  $q - p = 0$ . In order to identify the model, we will compare the adjusted ARMA models with different  $p$  and  $q$ . First we simply fit it to look at the Akaike Information Criteria (AIC) and the significance of the parameters estimated.

The AIC measures the goodness of fit and the simplicity of the model into a single statistic. Generally we aim to reduce the AIC.

$$AIC = 2k - 2 \ln(\hat{L}),$$

where  $k = p + q + 2$  and  $\hat{L}$  is the maximum value of the likelihood for the model.

```
##
## Call:
## arma(x = o3_train, order = c(1, 2))
##
## Model:
## ARMA(1,2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -0.0277211 -0.0048904 -0.0002275  0.0042125  0.0397965
## 
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)    
## ar1     9.889e-01  2.568e-03 385.12 < 2e-16 ***
## ma1    -4.774e-01  1.271e-02 -37.58 < 2e-16 ***
## ma2    -3.884e-01  1.235e-02 -31.45 < 2e-16 ***
## intercept 2.230e-04  5.335e-05    4.18 2.92e-05 ***
## ---
```

```

## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Fit:
## sigma^2 estimated as 5.143e-05, Conditional Sum-of-Squares = 0.26, AIC = -35981.74
##
## Call:
## arma(x = o3_train, order = c(2, 1))
##
## Model:
## ARMA(2,1)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.0256785 -0.0049644 -0.0002157  0.0042057  0.0398909
##
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)
## ar1        1.401e+00  1.469e-02   95.40 < 2e-16 ***
## ar2       -4.069e-01  1.432e-02  -28.41 < 2e-16 ***
## ma1       -9.288e-01  5.849e-03 -158.79 < 2e-16 ***
## intercept 1.146e-04  2.878e-05    3.98 6.89e-05 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Fit:
## sigma^2 estimated as 5.241e-05, Conditional Sum-of-Squares = 0.27, AIC = -35884.66
##
## Call:
## arma(x = o3_train, order = c(3, 2))
##
## Model:
## ARMA(3,2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.0277453 -0.0048477 -0.0001787  0.0042406  0.0393626
##
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)
## ar1        1.232e+00  8.610e-02   14.306 < 2e-16 ***
## ar2       -3.087e-01  1.244e-01   -2.482  0.0131 *
## ar3        6.808e-02  4.091e-02    1.664  0.0961 .
## ma1       -6.989e-01  8.560e-02   -8.164 2.22e-16 ***
## ma2       -1.928e-01  8.026e-02   -2.402  0.0163 *
## intercept 1.776e-04  4.444e-05    3.997 6.41e-05 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Fit:
## sigma^2 estimated as 5.127e-05, Conditional Sum-of-Squares = 0.26, AIC = -35993.35
##
## Call:

```

```

## arma(x = o3_train, order = c(2, 3))
##
## Model:
## ARMA(2,3)
##
## Residuals:
##      Min       1Q    Median       3Q      Max
## -0.0275175 -0.0048786 -0.0001909  0.0042262  0.0393501
##
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)
## ar1        7.861e-01  2.123e-01   3.702 0.000214 ***
## ar2        2.012e-01  2.099e-01   0.959 0.337659
## ma1       -2.510e-01  2.120e-01  -1.184 0.236467
## ma2       -4.576e-01  1.011e-01  -4.528 5.96e-06 ***
## ma3       -1.349e-01  8.342e-02  -1.617 0.105835
## intercept 2.564e-04  7.909e-05   3.242 0.001186 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Fit:
## sigma^2 estimated as 5.127e-05, Conditional Sum-of-Squares = 0.26, AIC = -35993.36
##
## Call:
## arma(x = o3_train, order = c(4, 3))
##
## Model:
## ARMA(4,3)
##
## Residuals:
##      Min       1Q    Median       3Q      Max
## -0.0276405 -0.0048667 -0.0001821  0.0042292  0.0391527
##
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)
## ar1        9.671e-01  4.143e-01   2.335 0.0196 *
## ar2       -1.015e-01  4.961e-01  -0.205 0.8379
## ar3        1.572e-01  1.670e-01   0.941 0.3466
## ar4       -3.447e-02  5.153e-02  -0.669 0.5036
## ma1       -4.328e-01  4.136e-01  -1.046 0.2955
## ma2       -2.585e-01  2.821e-01  -0.916 0.3596
## ma3       -1.636e-01  1.365e-01  -1.198 0.2307
## intercept 2.353e-04  9.353e-05   2.516 0.0119 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Fit:
## sigma^2 estimated as 5.129e-05, Conditional Sum-of-Squares = 0.26, AIC = -35987.24

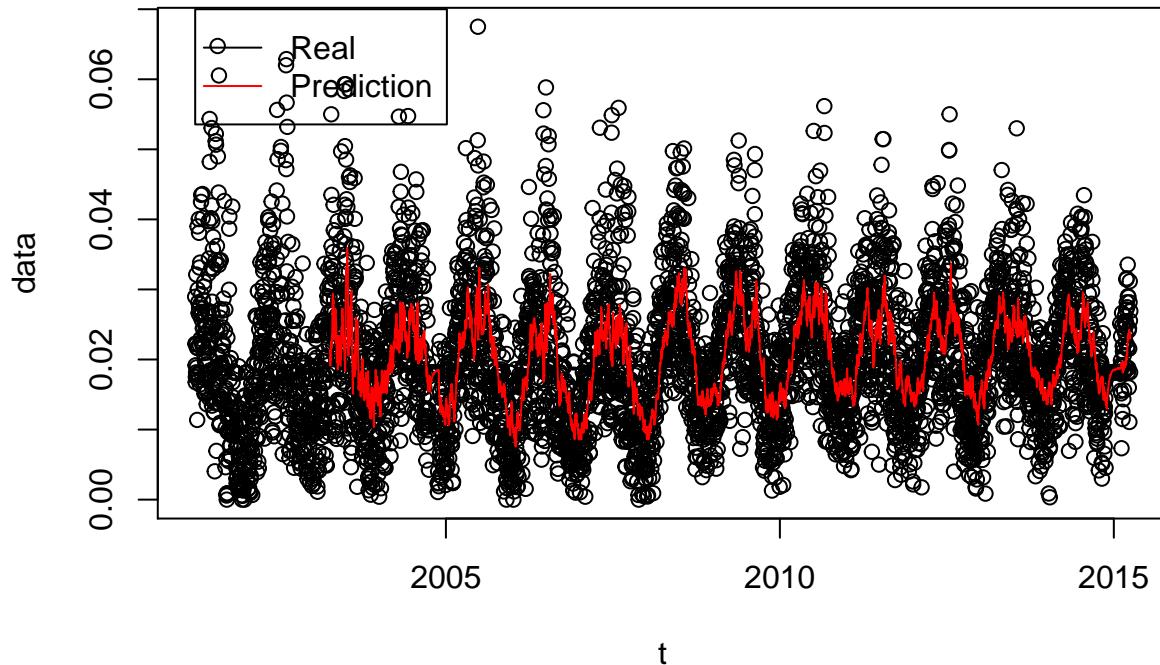
```

First we see that in the last two model, there is no statistical significance in the major part of the parameters, so we'll no consider it. The first, second and third models have significant parameters and similar AIC. In especial the 2° and 3° has the smallest AIC. So we will compare them both.

We see the MAE of the ARIMA(1, 0, 2):

```
## [1] 0.006330668
```

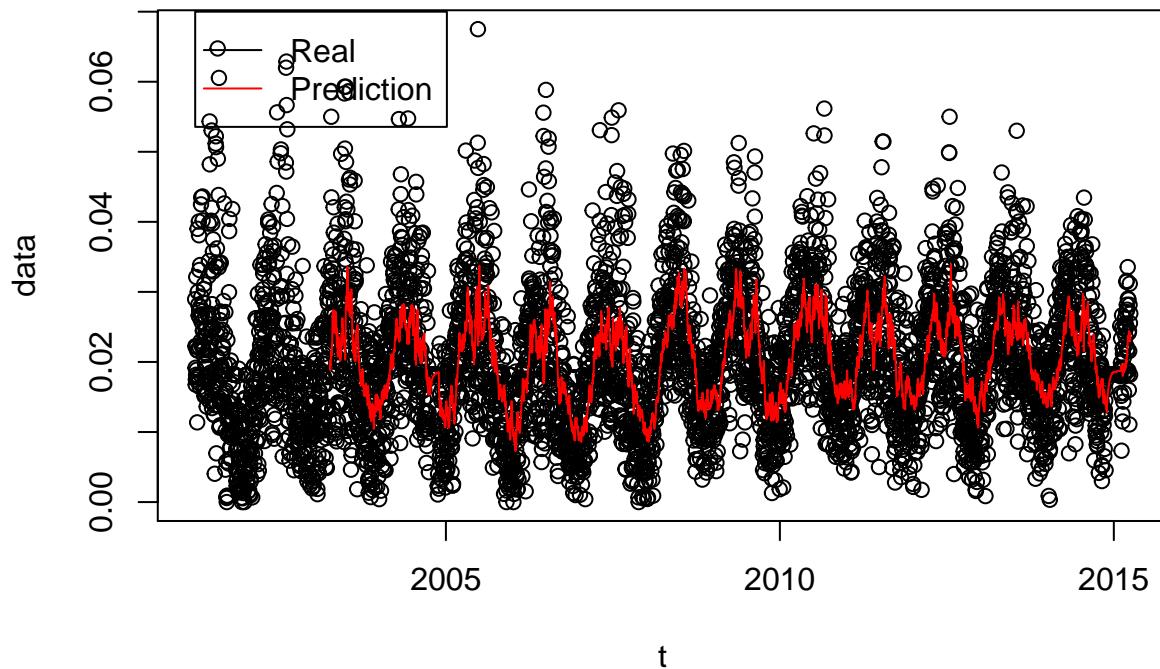
### ARIMA(1,0,2) prediction



And the MAE of the ARIMA(3, 0, 1):

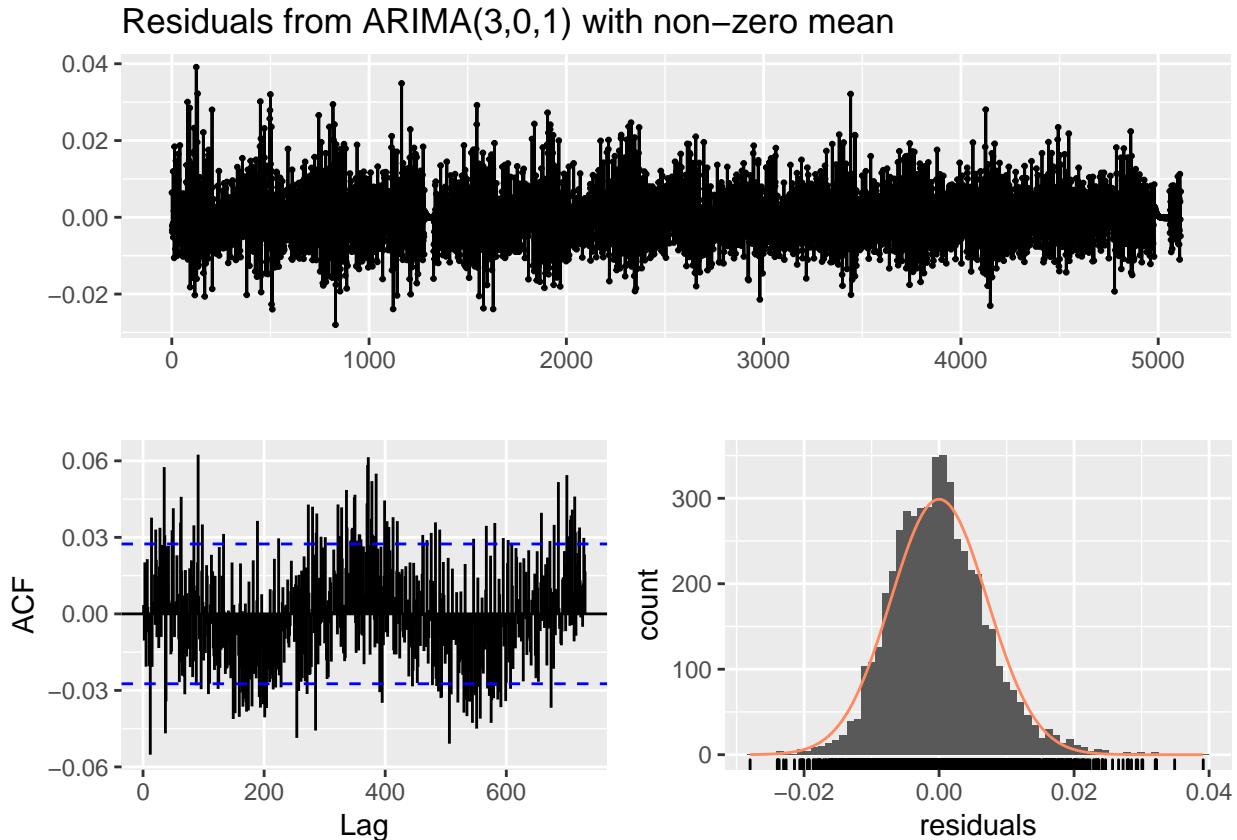
```
## [1] 0.006327397
```

### ARIMA(3,0,1) prediction



The second model seems a little better. So, let's check the residuals to observe it's problems.

```
model <- arima(o3_train, order = c(3,0,1))
checkresiduals(model, lag = 2*freq, lag.max = 2*freq)
```



```
##  
## Ljung-Box test  
##  
## data: Residuals from ARIMA(3,0,1) with non-zero mean  
## Q* = 1704.5, df = 725, p-value < 2.2e-16  
##  
## Model df: 5. Total lags used: 730
```

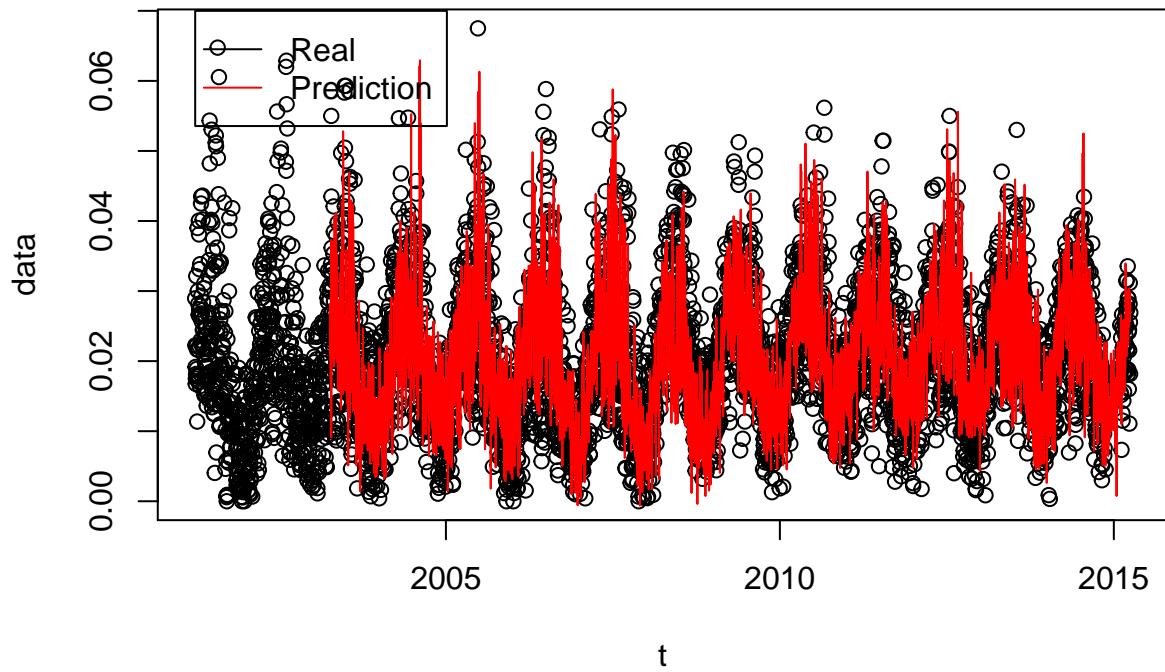
It's interesting to note the ACF has a peak around the 365, so the ARMA model did not seem to capture the seasonality. It would be better to fit an Seasonal ARIMA further. The histogram is pretty similar to normal distribution. The test made analyses if the mean is 0. It may be because the p-value is really small. That's the reason to adapt ARMA with STL, that is, fit an ARMA model in the residuals.

### Adapting ARMA

The ARMA seems to fit well as we can see so far. However, it's not capturing other characteristics on the data, as seasonality. For that reason, we will combine the stl and arma model and extract the best of each one. We will decompose the series in trend and seasonality and in the reminder, we fit an arima model with `auto.arima()`.

```
## [1] 0.007726959
```

## STL + ARIMA prediction



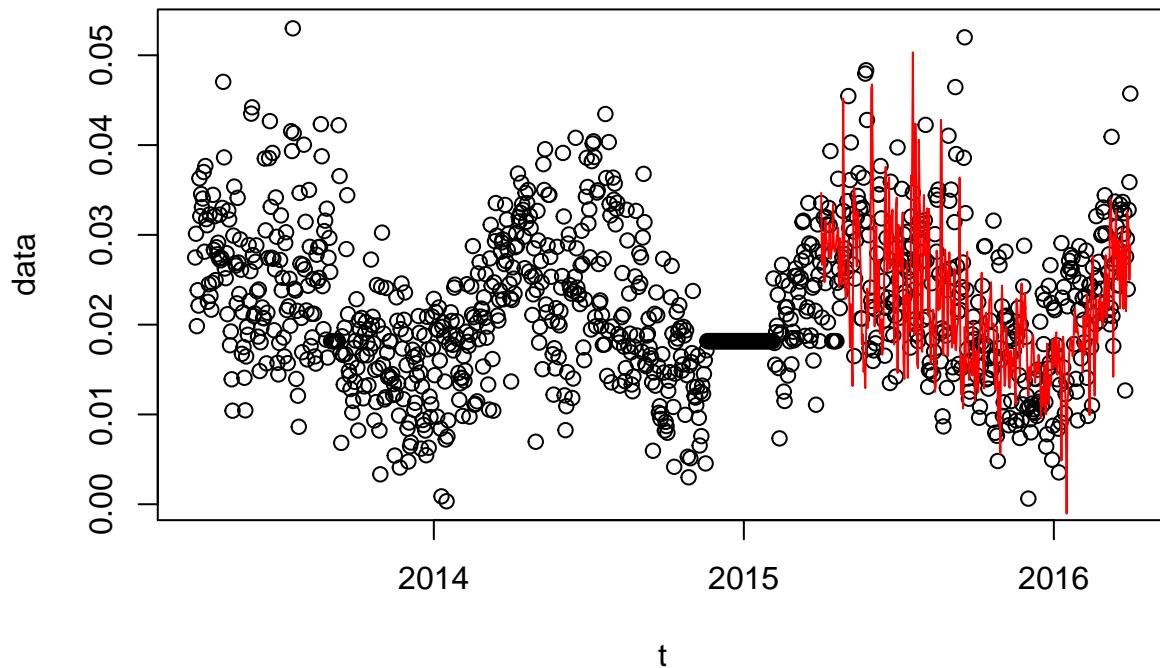
However the MAE isn't improved by this model. For that reason, this model was disregarded.

### Comparing models in the test data.

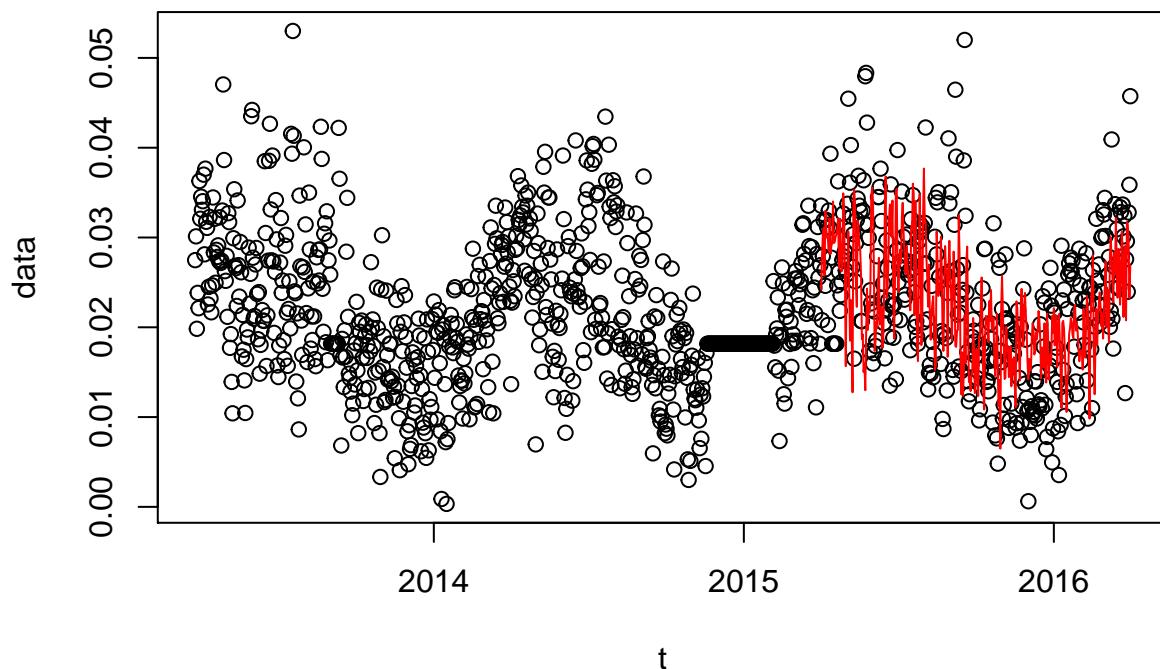
We chose some of the best models to compare with the testing data.

```
## [1] 0.007542956
```

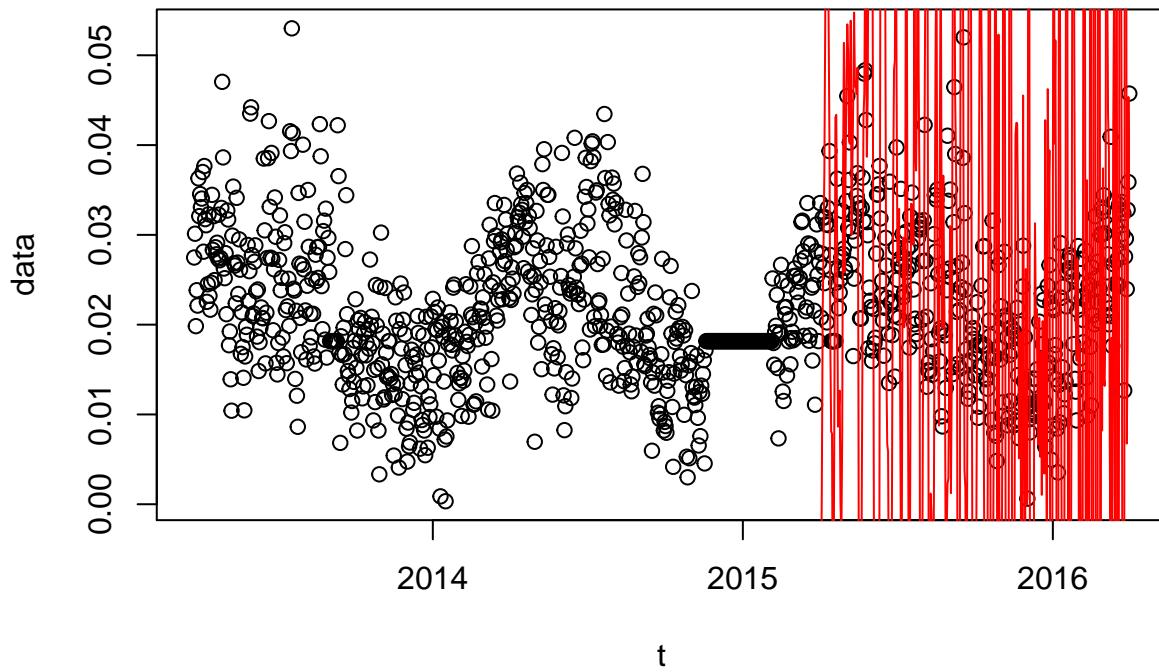
### Multiplicative decompose prediction



### Regression prediction

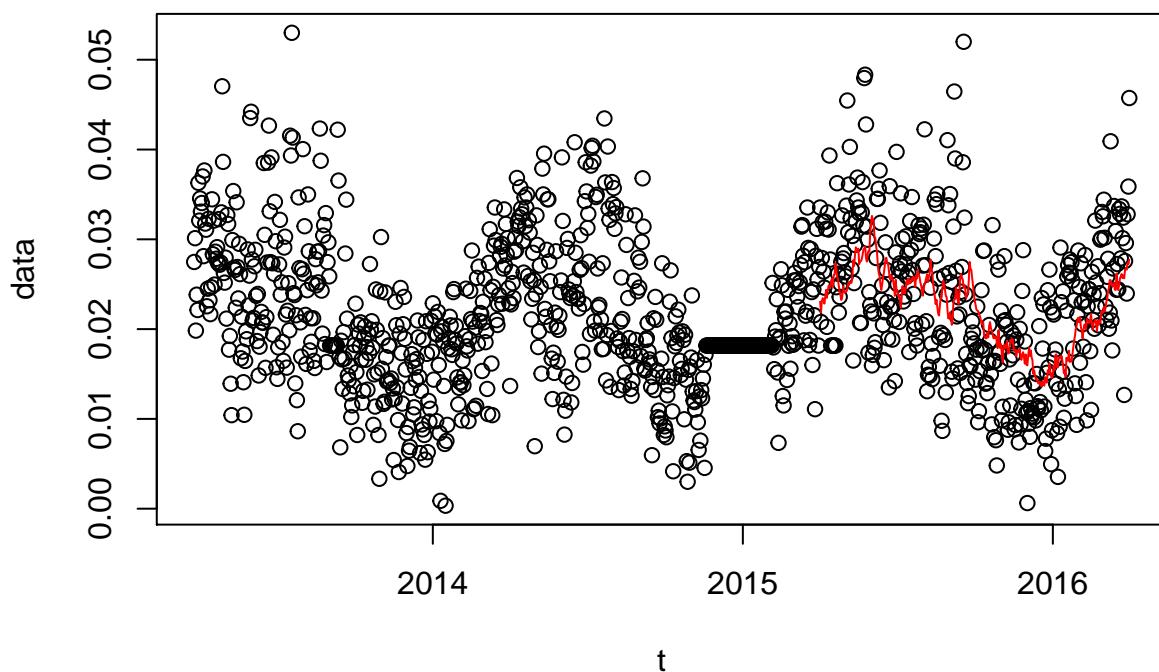


## Multiplicative Holt–Winters prediction



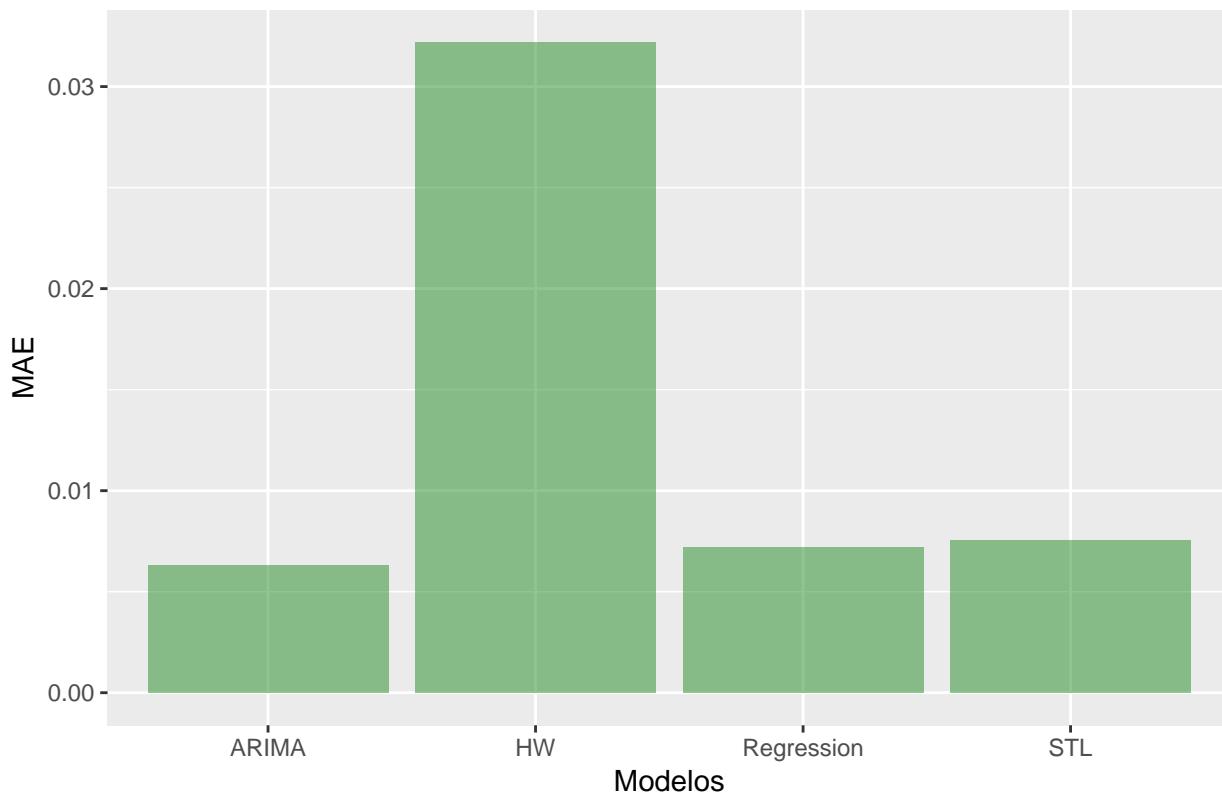
```
## [1] "MAE ARIMA(3,0,1)"  
## [1] 0.006289096
```

## ARIMA(3,0,1) prediction



Finally, we have this bar graphic to compare the values:

## Comparing model's MAE in test data



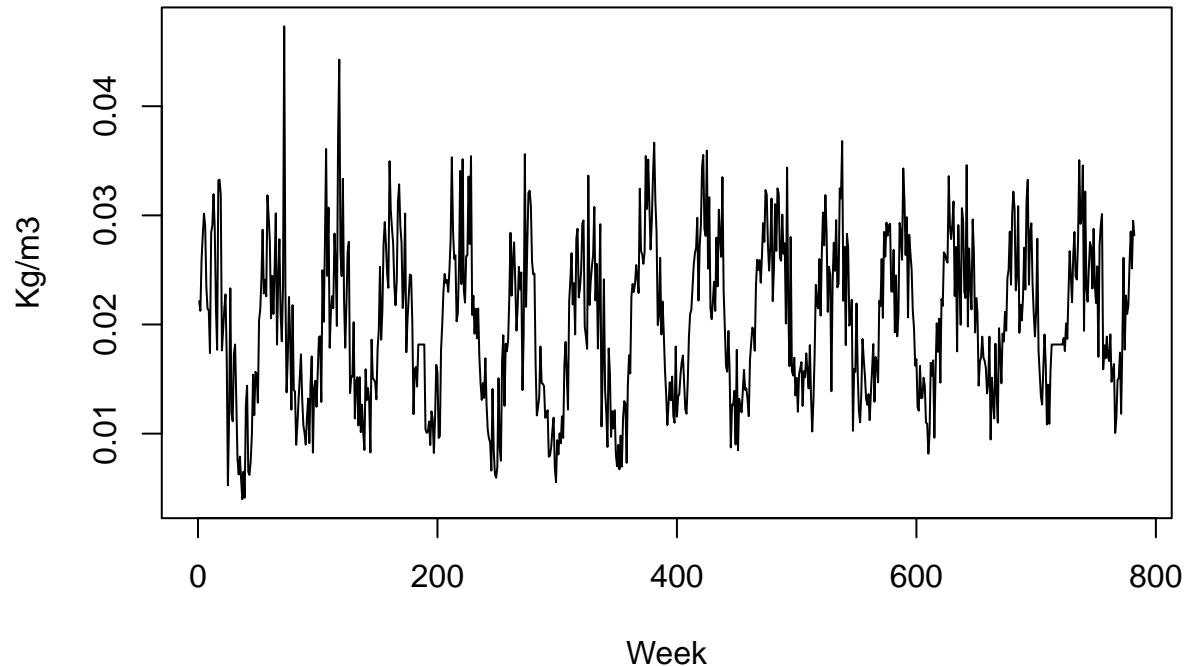
### Models: case 2

In case two, we have to aggregate the diary days in a week, starting from the sunday, as requested. So we calculate the mean value in the week to be its representant. The models may be very similar to the previous. We may see less outliers. We also will separate train and test data. The first day in the data is April 1, 2001, a Sunday. So we do not worry about that.

```
o3_week <- c(1:floor(length(o3.clean)/7))
for(i in seq(1,length(o3.clean)-7, 7)){
  o3_week[ceiling(i/7)] = mean(o3.clean[i:(i+6)])
}
dates = seq(from = as.Date(time(o3[1])), to = as.Date(time(o3[length(o3]))), by = "weeks")[1:782]

plot(ts(o3_week), main = 'Weekly average level of O3 in Boston (after imputation)',
      xlab = 'Week', ylab = 'Kg/m3')
```

## Weekly average level of O<sub>3</sub> in Boston (after imputation)



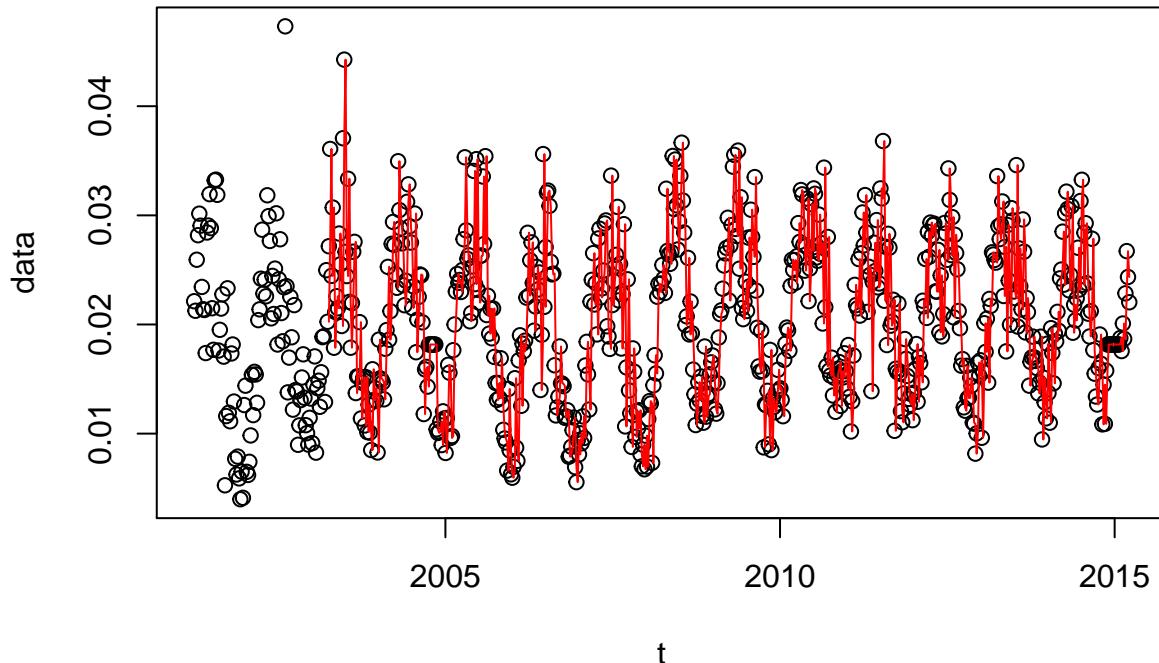
```
o3_train_week = o3_week[1:(length(o3_week)[1] - 52)]  
o3_test_week = o3_week[-c(1:(length(o3_week)[1] - 52))]
```

### Baseline Model

We will do the naive forecast to the baseline model.

```
## [1] 0.01039995
```

## Baseline model prediction



### Decompose

We now check for seasonality considering monthly (4 records) and yearly (52 records). For seasonality of 52, we see better p-value. In fact, less than 0.05. So we will use 52.

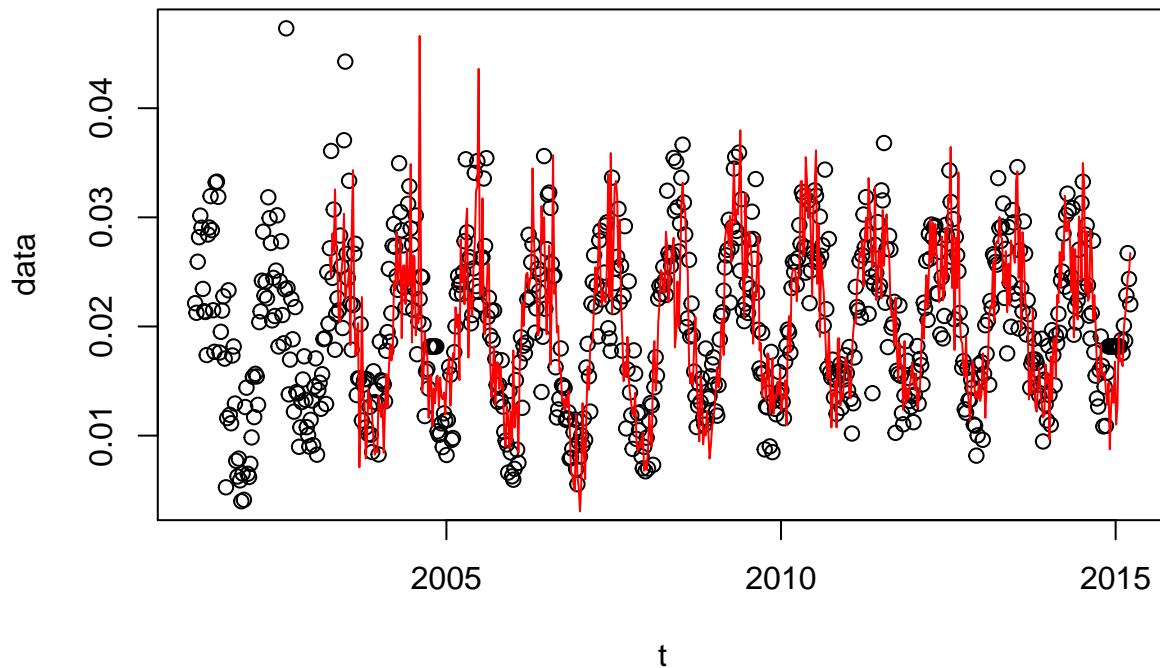
```
##  
## Kruskal-Wallis rank sum test  
##  
## data: o3_train_week and g  
## Kruskal-Wallis chi-squared = 0.12693, df = 3, p-value = 0.9884  
##  
## Kruskal-Wallis rank sum test  
##  
## data: o3_train_week and g  
## Kruskal-Wallis chi-squared = 527.25, df = 51, p-value < 2.2e-16
```

### Additive model

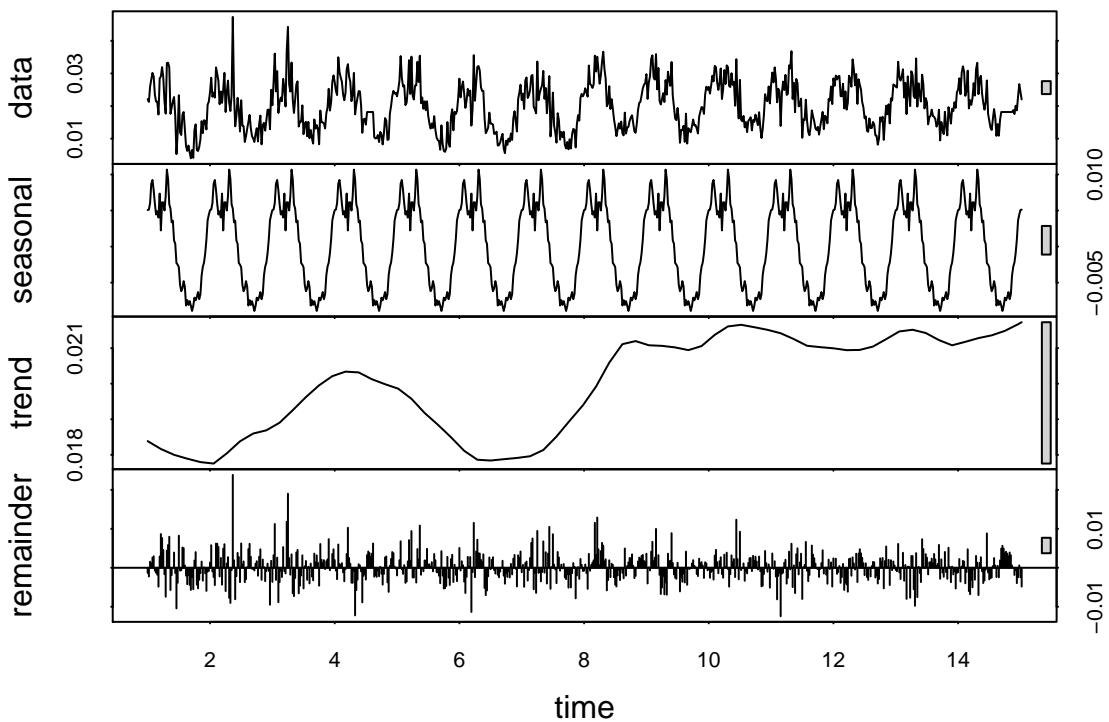
We see the MAE of the additive decompose model:

```
## [1] 0.003993143
```

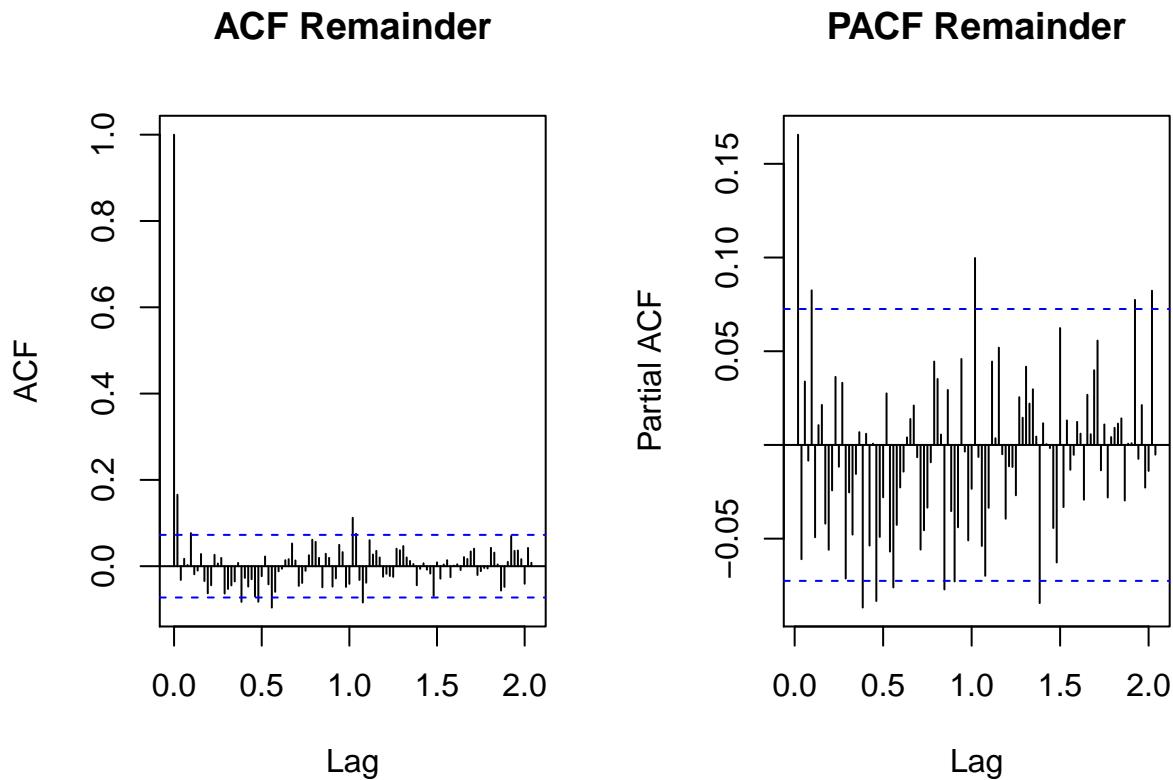
## Additive decompose prediction



We fit the model using `t.window` and analyse the reminder:



The ACF and the PACF of the reminder:



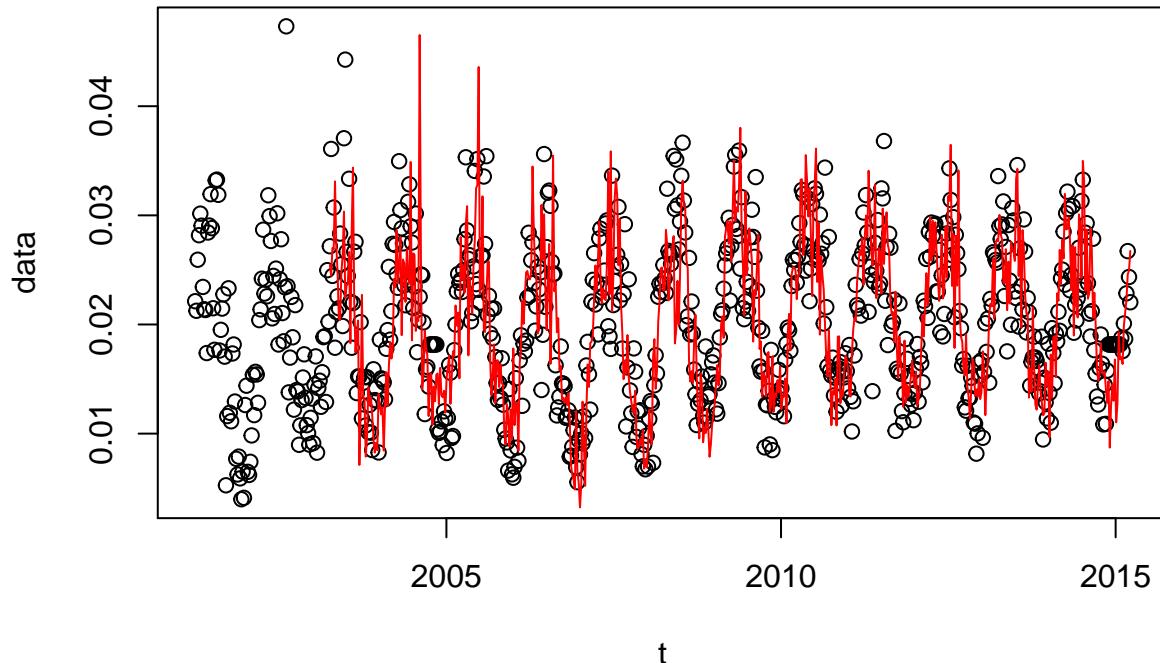
We see a big spike when lag = 52. It seems not so good for a reminder

### Multiplicative model

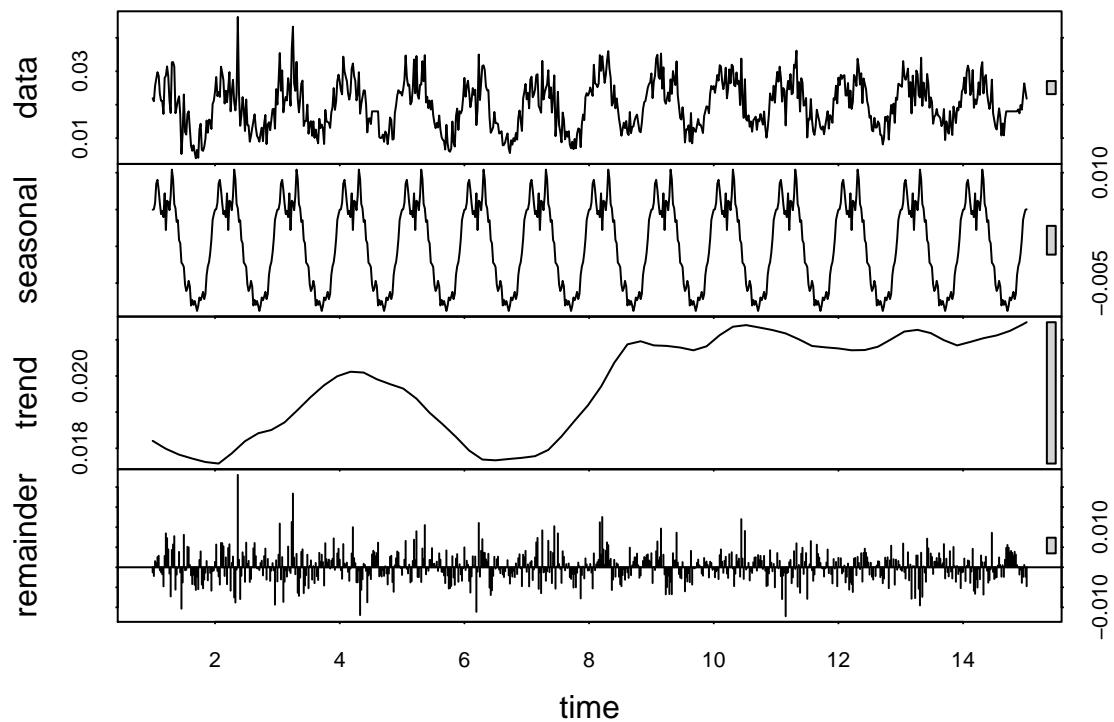
We analyse the MAE:

```
## [1] 0.003968754
```

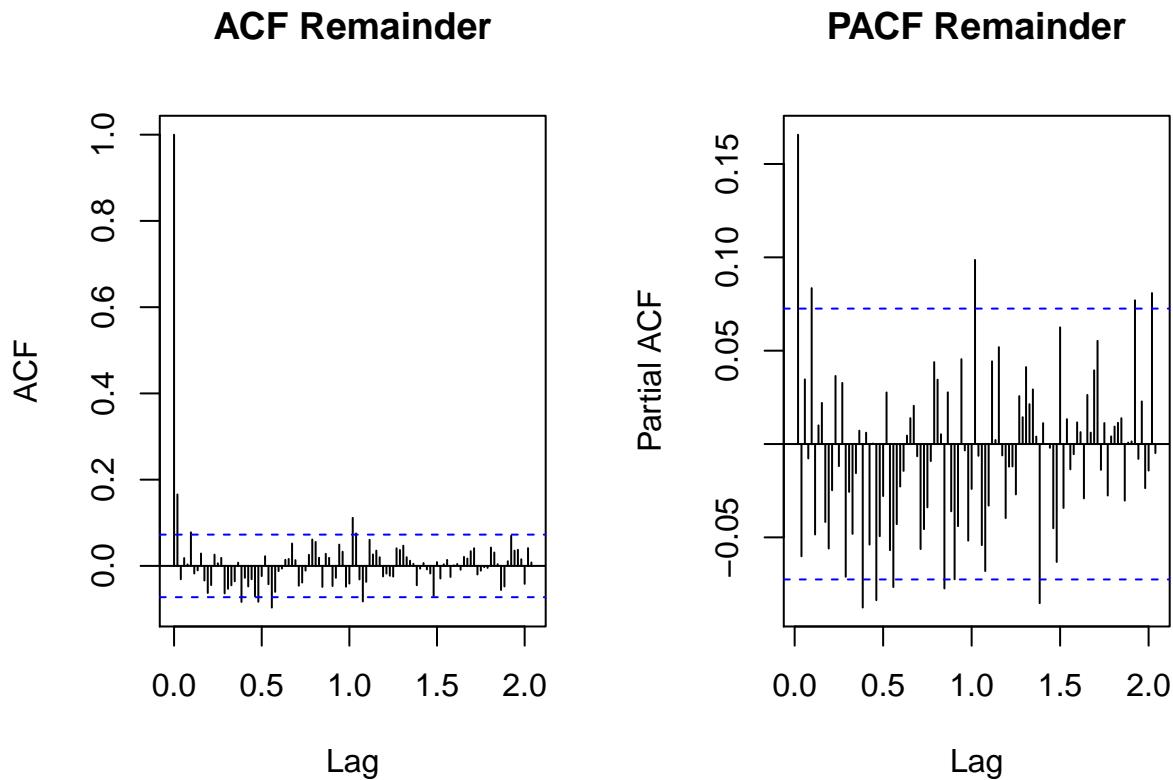
## Multiplicative decompose prediction



Using `t.window` and analysing the reminder:



The ACF and the PACF of the reminder:



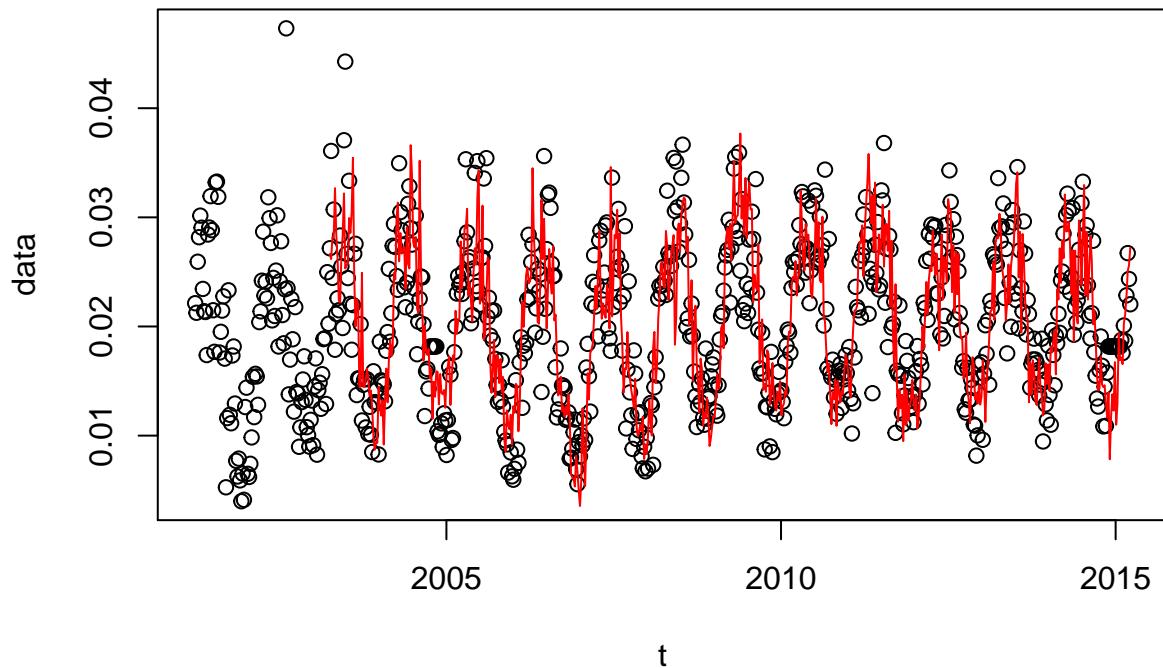
We also see the spike when lag = 52. It seems not so good for a reminder. The same as before.

### Regression

We tested for seasonalities, and we settled with 52. So we fit a regression model with the seasonality dummies. We see the MAE:

```
## [1] 0.003803593
```

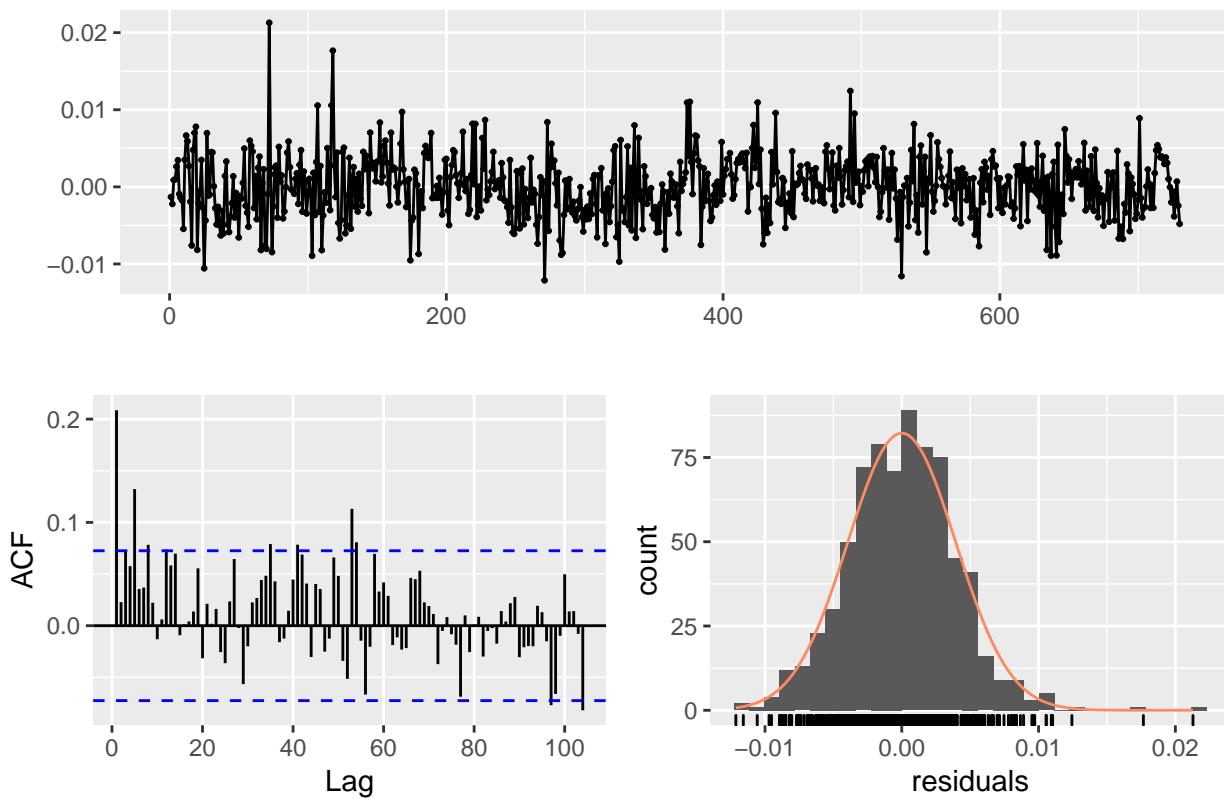
## Regression prediction



Now we analyse the residuals.

```
train = data.frame(  
  t = t,  
  o3_train_week = o3_train_week,  
  Q = Q  
)  
mod = lm(o3_train_week~t+Q, data = train)  
checkresiduals(mod, lag = 2*freq, lag.max = 2*freq)
```

## Residuals



```
##  
## Breusch-Godfrey test for serial correlation of order up to 104  
##  
## data: Residuals  
## LM test = 142.53, df = 104, p-value = 0.007267
```

We see spikes in lag = 52, 104. Same as before.

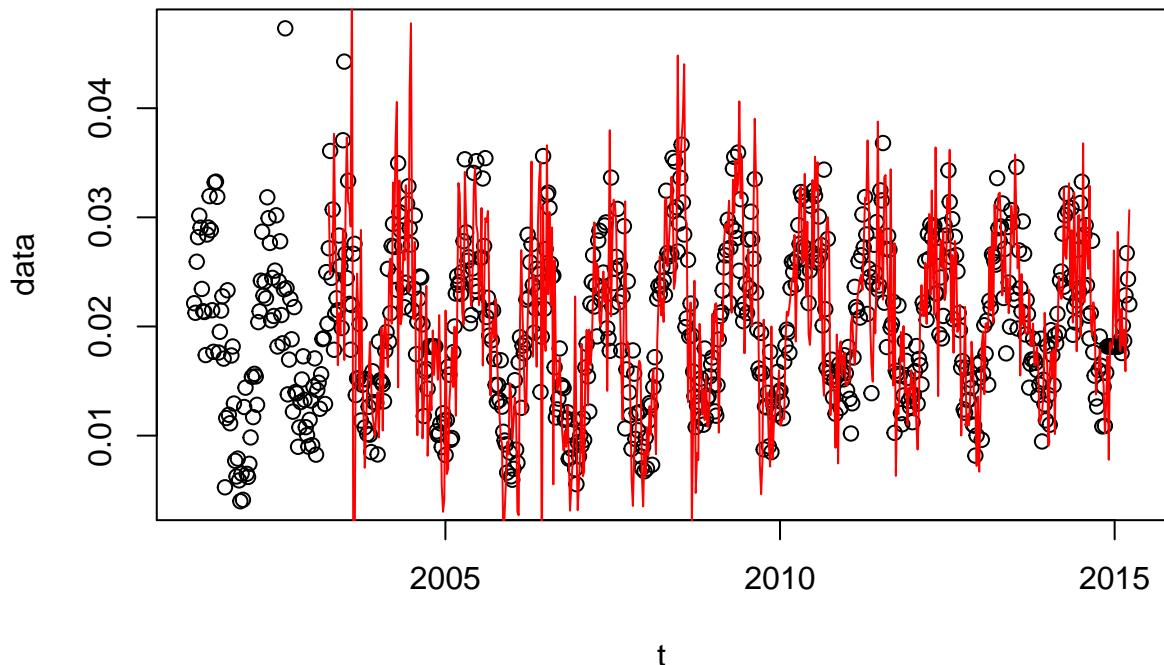
## Holt-Winters

Now we will try Holt-Winters models. We will consider complete Holt-Winters models, those with seasonality.

### Additive

```
## [1] 0.005158324
```

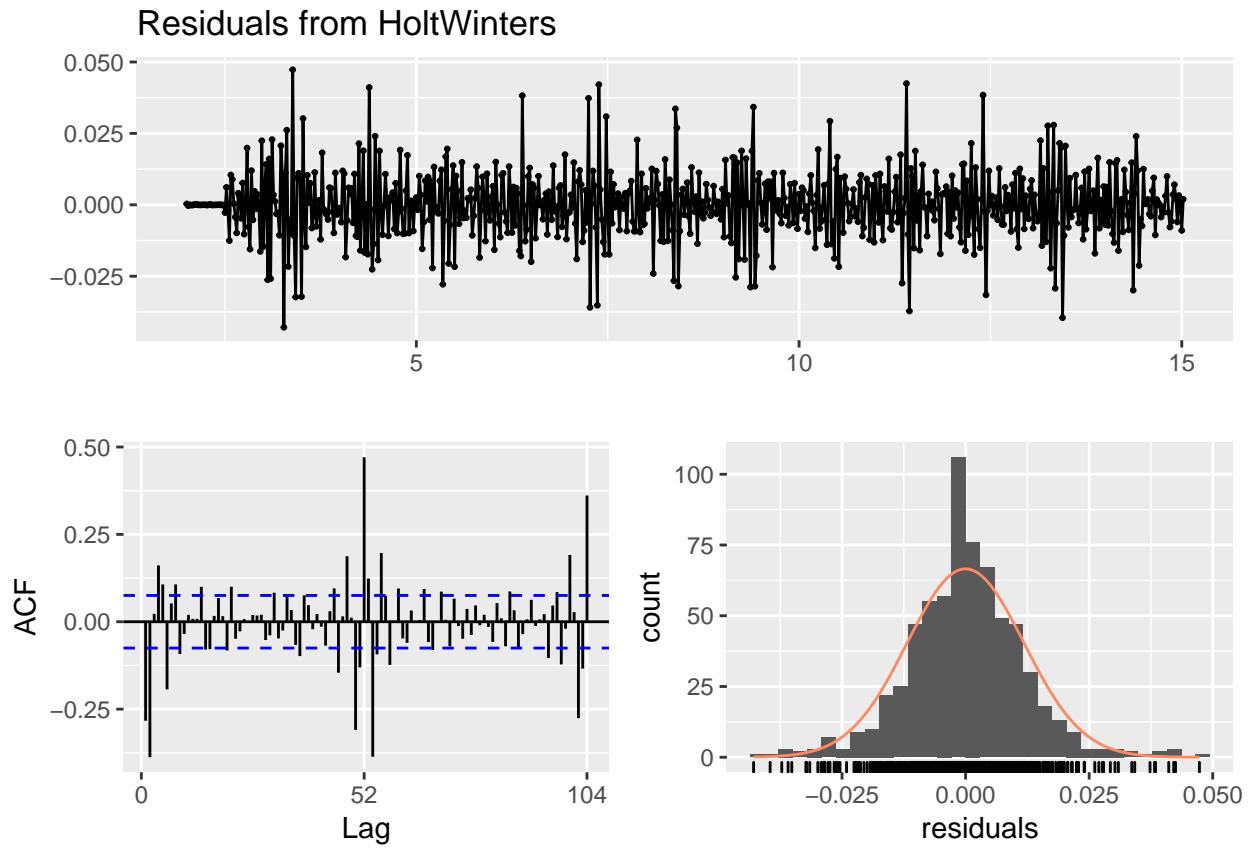
## Additive Holt–Winters prediction



MAE is not so good as the other models seen so far. We see some kind of explosion in the resulting graph. We know that HW tends to “explode”. We have seem better results. Let's analyse the residuals:

```
mod = HoltWinters(ts(o3_train_week, frequency = 52), beta = T, gamma = T, seasonal = "additive")
checkresiduals(mod, lag = 2*freq, lag.max = 2*freq)

## Warning in modeldf.default(object): Could not find appropriate degrees of
## freedom for this model.
```

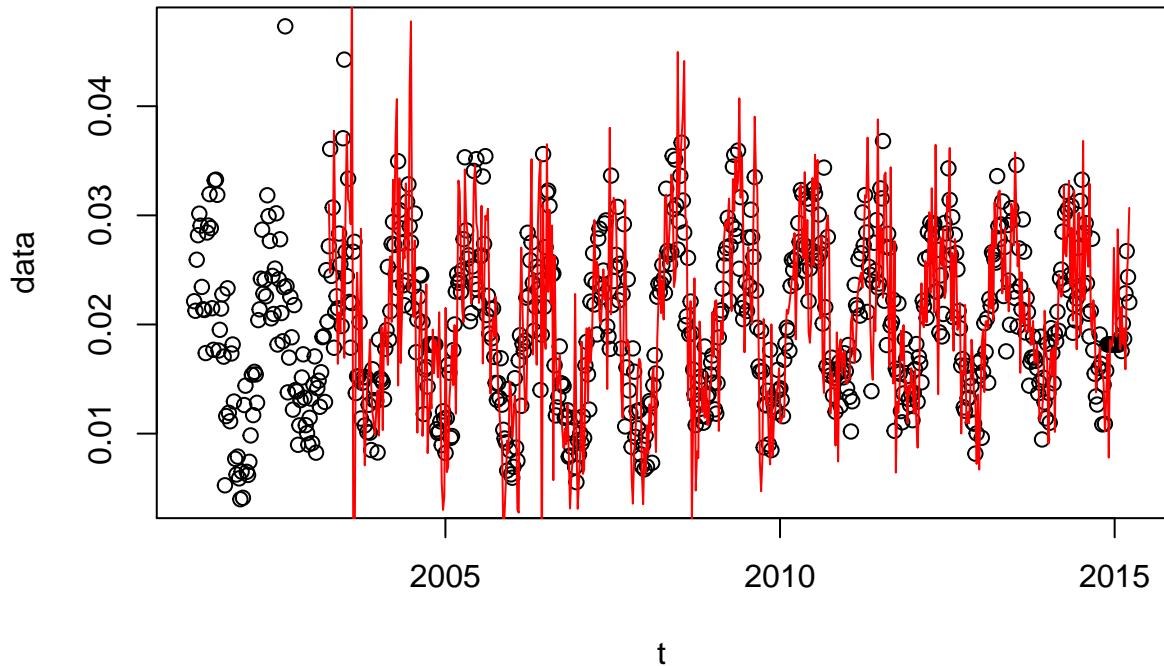


We see the same problems as before: high correlated lag = 52, 104, evidence of this not being a WN.

#### Multiplicative

```
## [1] 0.005168137
```

## Multiplicative Holt–Winters prediction

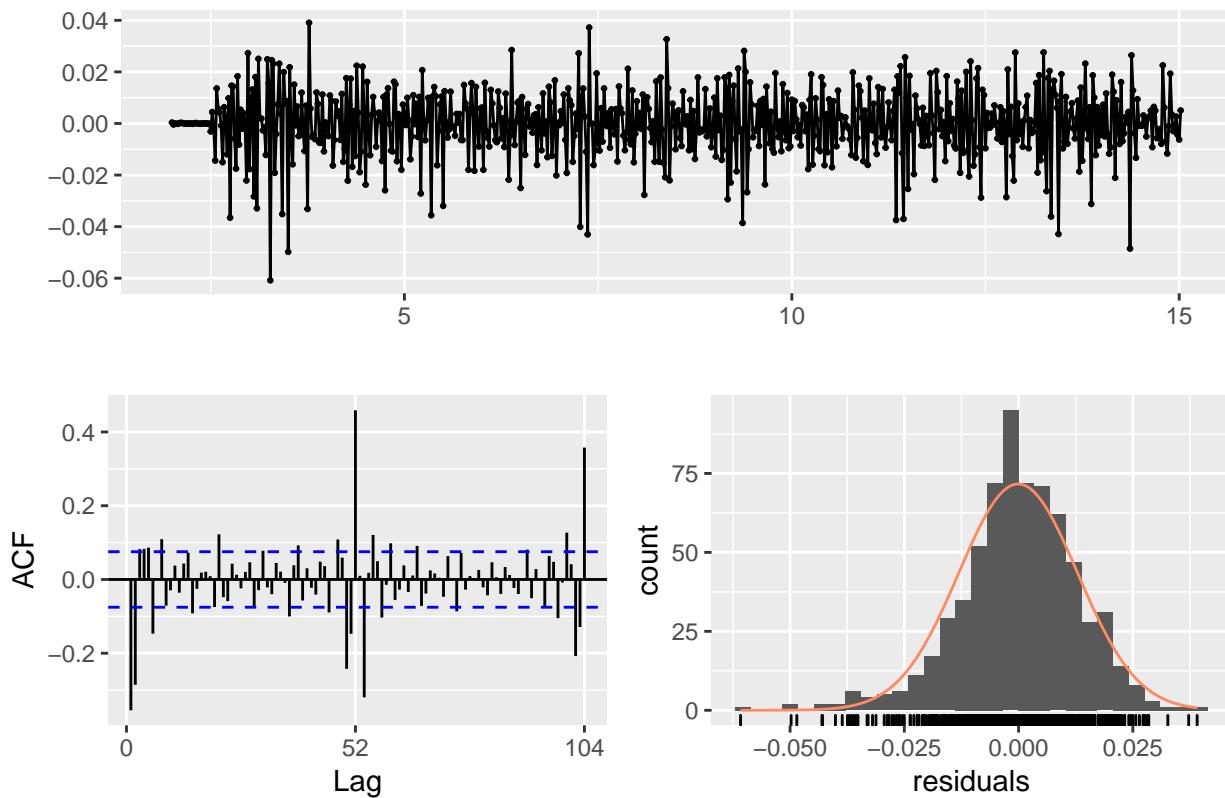


We see also a not so good MAE and graph. Let's analyse the residuals:

```
mod = HoltWinters(ts(o3_train_week, frequency = 52), beta = T, gamma = T, seasonal = "multiplicative")
checkresiduals(mod, lag = 2*freq, lag.max = 2*freq)

## Warning in modelfdf.default(object): Could not find appropriate degrees of
## freedom for this model.
```

### Residuals from HoltWinters

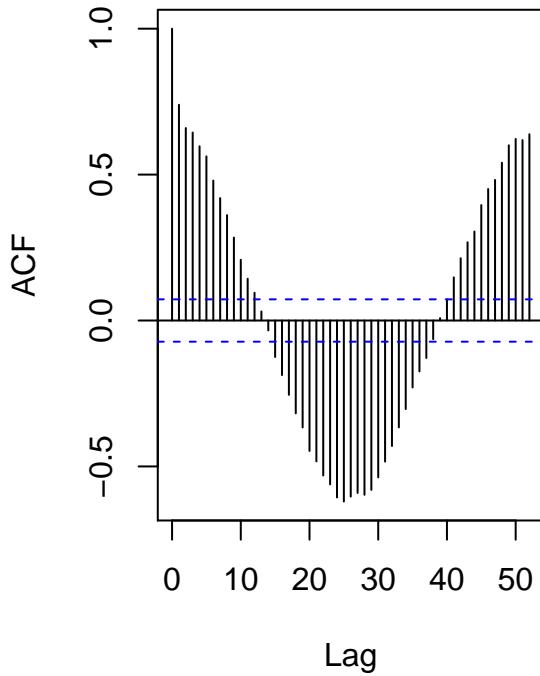


Same problems as before. Holt-Winters doesn't appear to work well in this case.

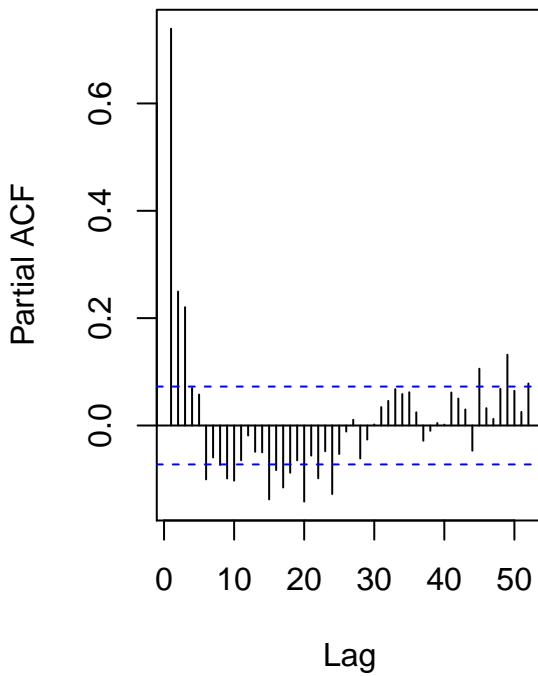
-### ARMA

We can see the ACF and PACF:

**ACF of the data**



**PACF of the data**



Based on these graphs, we see both graphs has a exponentially decay, the first after the  $q - p = -1$  or  $q - p = 0$ . In order to identify the model, we will compare the adjusted ARMA models with different  $p$  and  $q$ . First we simply fit it to look at the Akaike Information Criteria (AIC) and the significance of the parameters estimated.

```
##
## Call:
## arma(x = o3_train_week, order = c(2, 0))
##
## Model:
## ARMA(2,0)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.0151516 -0.0032246 -0.0002512  0.0028414  0.0257926
##
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)
## ar1        0.5546040  0.0358395 15.475 < 2e-16 ***
## ar2        0.2498525  0.0358459  6.970 3.17e-12 ***
## intercept  0.0039316  0.0005489  7.162 7.93e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Fit:
## sigma^2 estimated as 2.257e-05,  Conditional Sum-of-Squares = 0.02,  AIC = -5732.5
##
## Call:
## arma(x = o3_train_week, order = c(2, 1))
```

```

## 
## Model:
## ARMA(2,1)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.014583 -0.003300 -0.000134  0.002970  0.025345
## 
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)    
## ar1        0.9751555  0.0700310 13.925 < 2e-16 ***
## ar2       -0.0463356  0.0607827 -0.762   0.446    
## ma1       -0.4952763  0.0585446 -8.460 < 2e-16 ***
## intercept  0.0014296  0.0003658  3.908 9.29e-05 *** 
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Fit:
## sigma^2 estimated as 2.162e-05, Conditional Sum-of-Squares = 0.02, AIC = -5761.97
## 
## Call:
## arma(x = o3_train_week, order = c(3, 1))
## 
## Model:
## ARMA(3,1)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.0140809 -0.0032711 -0.0001538  0.0029207  0.0251467
## 
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)    
## ar1        0.7724162  0.1110272  6.957 3.48e-12 ***
## ar2       -0.0260143  0.0701409 -0.371  0.71072    
## ar3        0.1551277  0.0497065  3.121  0.00180 **  
## ma1       -0.2899243  0.1088996 -2.662  0.00776 **  
## intercept  0.0019742  0.0005461  3.615  0.00030 *** 
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Fit:
## sigma^2 estimated as 2.133e-05, Conditional Sum-of-Squares = 0.02, AIC = -5769.92
## Warning in arma(o3_train_week, order = c(3, 2)): Hessian negative-
## semidefinite
## Warning in sqrt(diag(object$vcov)): NaNs produced
## Warning in sqrt(diag(object$vcov)): NaNs produced
## 
## Call:
## arma(x = o3_train_week, order = c(3, 2))
## 
## Model:

```

```

## ARMA(3,2)
##
## Residuals:
##      Min       1Q    Median       3Q      Max
## -1.401e-02 -3.242e-03 -9.053e-05  2.926e-03  2.527e-02
##
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)
## ar1          0.6314855     NA        NA        NA
## ar2          0.1161892     NA        NA        NA
## ar3          0.1431114    0.0393657   3.635 0.000278 ***
## ma1         -0.1477659     NA        NA        NA
## ma2         -0.0693593     NA        NA        NA
## intercept   0.0021847    0.0004249   5.142 2.71e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Fit:
## sigma^2 estimated as 2.135e-05, Conditional Sum-of-Squares = 0.02, AIC = -5767.21
## Warning in arma(o3_train_week, order = c(4, 2)): Hessian negative-
## semidefinite

## Warning in arma(o3_train_week, order = c(4, 2)): NaNs produced

## Warning in arma(o3_train_week, order = c(4, 2)): NaNs produced
##
## Call:
## arma(x = o3_train_week, order = c(4, 2))
##
## Model:
## ARMA(4,2)
##
## Residuals:
##      Min       1Q    Median       3Q      Max
## -0.0140735 -0.0032840 -0.0001482  0.0029149  0.0251517
##
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)
## ar1          0.5552339     NA        NA        NA
## ar2          0.1374963     NA        NA        NA
## ar3          0.1507390    0.0504035   2.991 0.00278 **
## ar4          0.0362919    0.0564229   0.643 0.52009
## ma1         -0.0750043     NA        NA        NA
## ma2         -0.0568046     NA        NA        NA
## intercept   0.0024065    0.0005688   4.231 2.33e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Fit:
## sigma^2 estimated as 2.134e-05, Conditional Sum-of-Squares = 0.02, AIC = -5765.57
## Warning in arma(o3_train_week, order = c(4, 3)): Hessian negative-
## semidefinite

```

```

## Warning in arma(o3_train_week, order = c(4, 3)): NaNs produced
## Warning in arma(o3_train_week, order = c(4, 3)): NaNs produced
##
## Call:
## arma(x = o3_train_week, order = c(4, 3))
##
## Model:
## ARMA(4,3)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.0134441 -0.0031889 -0.0001146  0.0029098  0.0249980
##
## Coefficient(s):
##             Estimate Std. Error t value Pr(>|t|)
## ar1        0.7282249     NA       NA       NA
## ar2        0.4180584     NA       NA       NA
## ar3       -0.4544918    0.0666640  -6.818 9.25e-12 ***
## ar4        0.2009348     NA       NA       NA
## ma1       -0.2669488     NA       NA       NA
## ma2       -0.4451386     NA       NA       NA
## ma3        0.4472807    0.1095129   4.084 4.42e-05 ***
## intercept  0.0021247    0.0004557   4.662 3.13e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Fit:
## sigma^2 estimated as 2.09e-05, Conditional Sum-of-Squares = 0.02, AIC = -5778.63

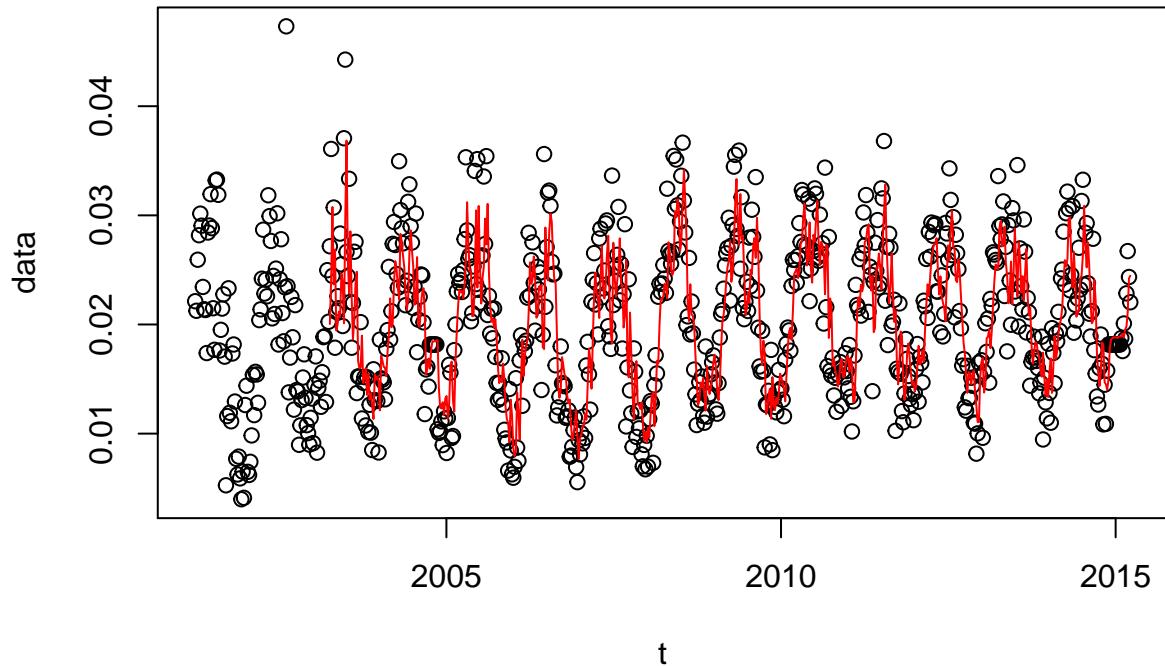
```

We see that from the second onwards, we have parameters without statistical significance. However the second and third models have good AIC either. The last three has problems with the Hessian, so the p-values could not be calculated. For that reason, we shall disregard them and compare the others. In special we'll test ARMA(3,1), AR(2) and ARMA(2,1).

We see the MAE of the AR(2):

```
## [1] 0.003525595
```

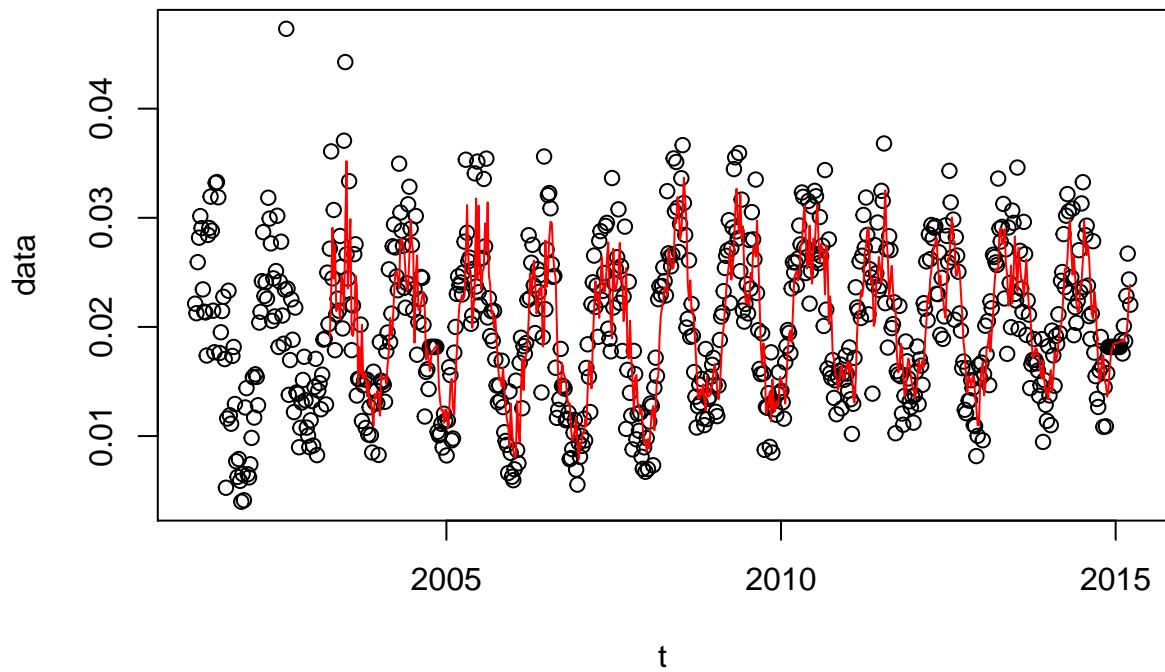
### ARIMA(3,0,1) prediction



And the MAE of the ARMA(2, 1):

```
## [1] 0.003507216
```

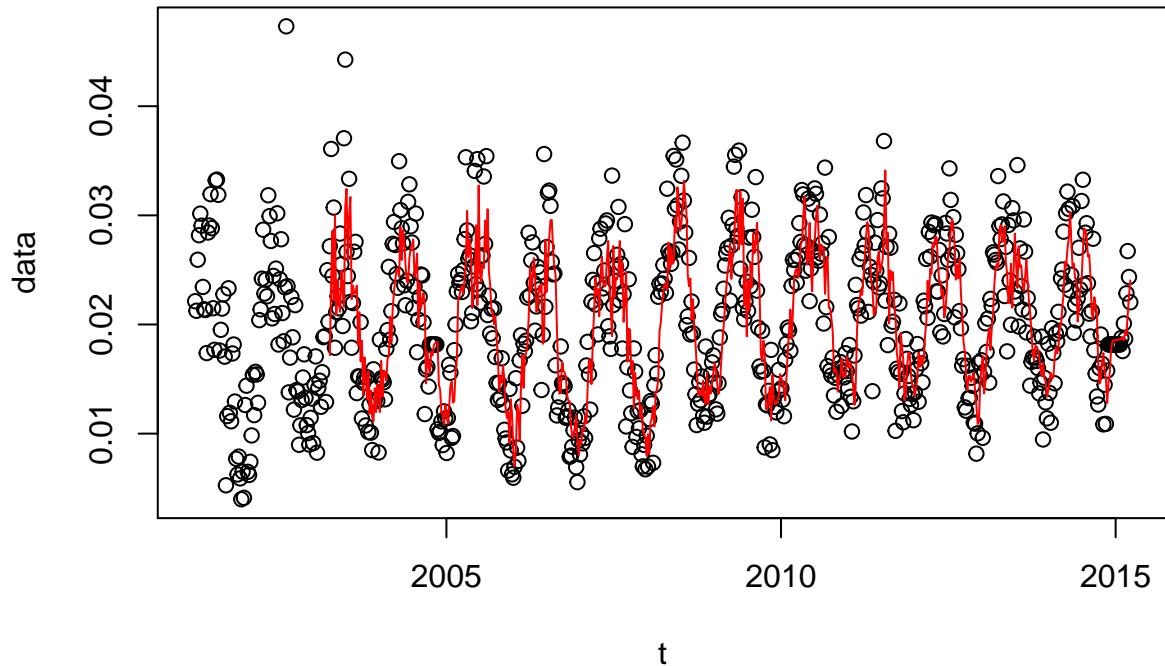
### ARIMA(2,0,1) prediction



And the MAE of the ARMA(3, 1):

```
## [1] 0.003558913
```

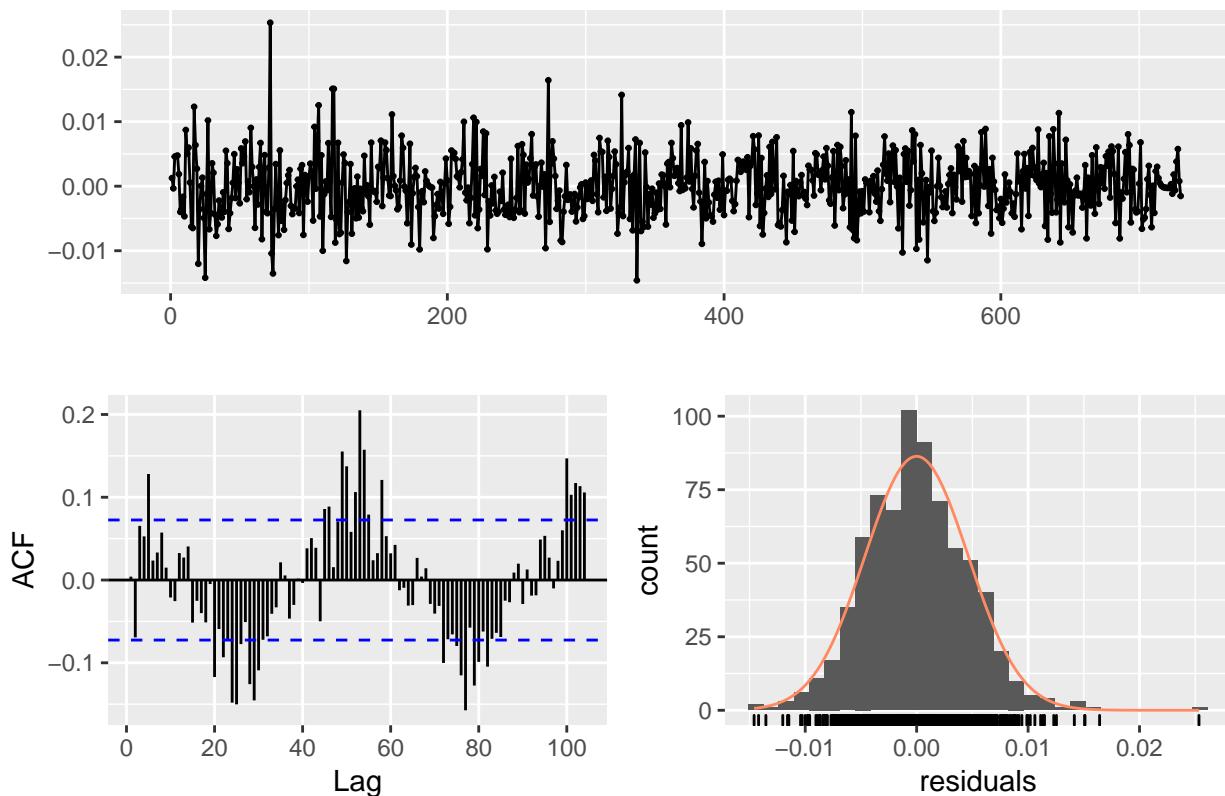
## ARIMA(3,0,1) prediction



The second model seems a little better. So, let's check the residuals to observe it's problems.

```
model <- arima(o3_train_week, order = c(2,0,1))
checkresiduals(model, lag = 2*freq, lag.max = 2*freq)
```

## Residuals from ARIMA(2,0,1) with non-zero mean



```
##  
## Ljung-Box test  
##  
## data: Residuals from ARIMA(2,0,1) with non-zero mean  
## Q* = 465.72, df = 100, p-value < 2.2e-16  
##  
## Model df: 4. Total lags used: 104
```

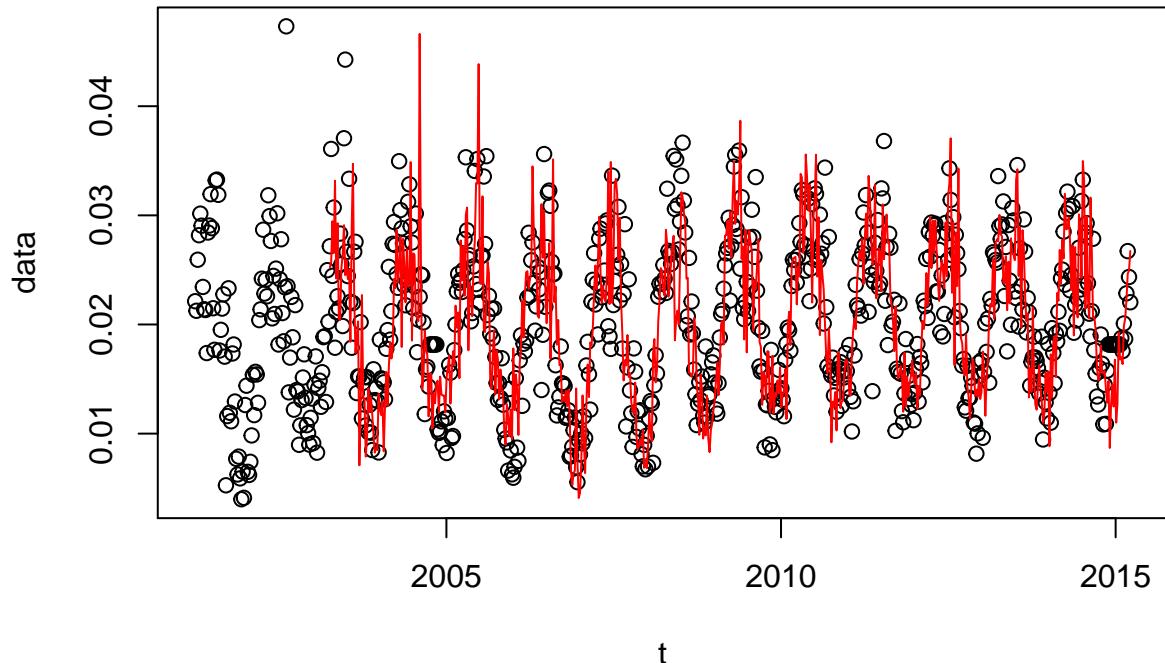
It's interesting to note the ACF has a peak around the 52, so the ARMA model did not seem to capture the seasonality. It would be better to fit an Seasonal ARIMA further. The histogram is pretty similar to normal distribution. The test made analyses if the mean is 0. It may be because the p-value is really small. The ACF also suggests that we include a high order parameter in MA and AR. However, in trying to fit an ARMA(3,2), we struggle with stationarity in AR part. For that reason, the model has it fails and alone it seems that it cannot be improved.

### Adapting ARMA

The ARMA seems to fit well as we can see so far. However, it's not capturing other characteristics on the data, as seasonality as mentioned. For that reason, we will combine the STL and ARMA model and extract the best of each one. We will decompose the series in trend and seasonality and in the reminder, we fit an arima model with `auto.arima()`.

```
## [1] 0.003987221
```

## STL + ARIMA prediction



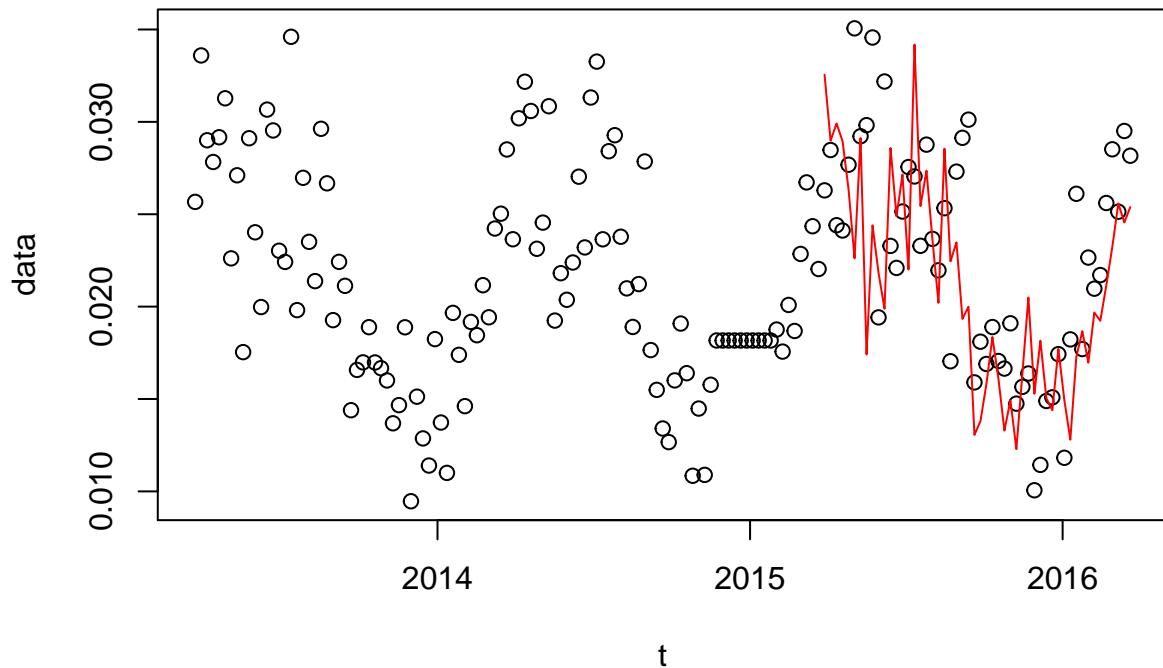
In spite of the improvement in the STL model, the MAE isn't improved by this model when compared to ARMA best model. For that reason, this model was disregarded.

### Comparing models in the test data.

We chose some of the best models to compare with the testing data.

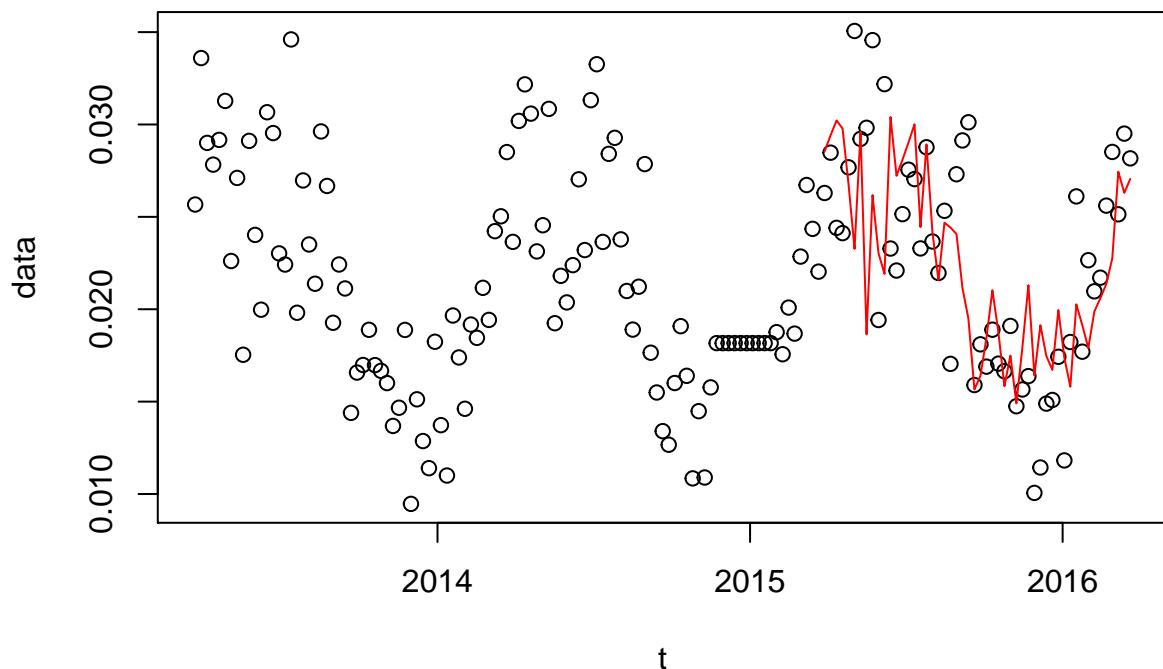
```
## [1] 0.004080176
```

### Multiplicative decompose prediction



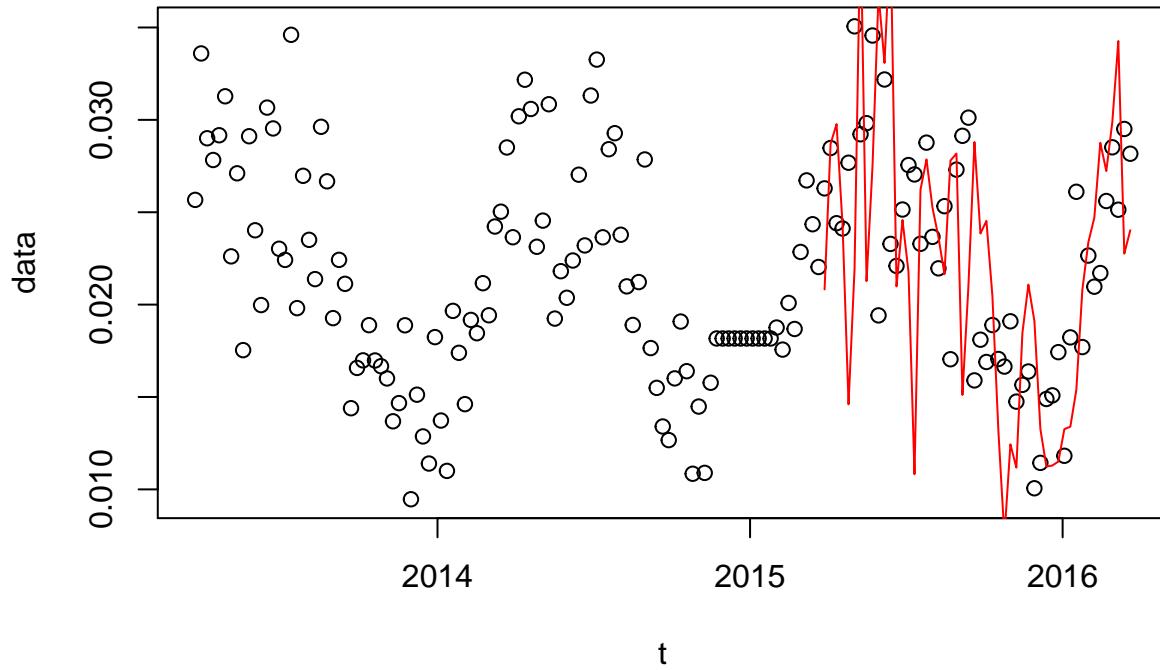
```
## [1] 0.003646437
```

### Regression prediction

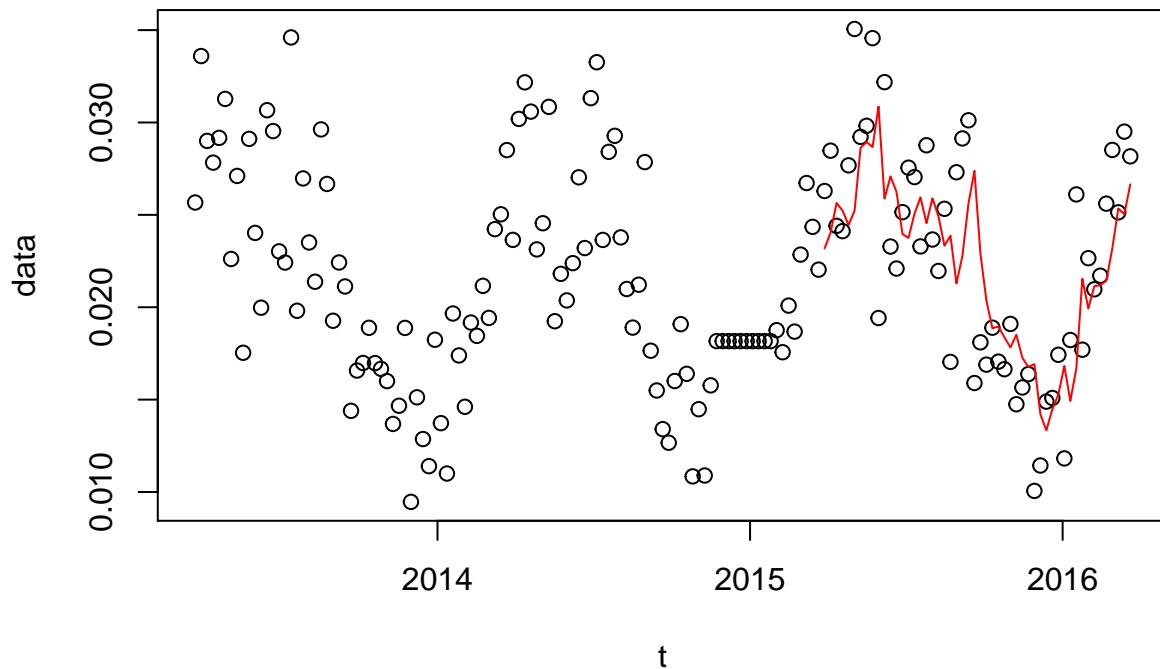


```
## [1] 0.00589603
```

### Multiplicative Holt–Winters prediction



### ARIMA(2,0,1) prediction



Finally, we have this bar graphic to compare the values:

Comparing model's MAE in test data

