

Danmarks
Tekniske
Universitet



Time Series Analysis Assignment 4

AUTHOR

Ghassen Lassoued - s196609

May 15, 2020

Contents

1	Question 4.1	1
2	Question 4.2	2
3	Question 4.3	3
4	Question 4.4	8
5	Question 4.5	10
6	Question 4.6	10
7	Code	11
	List of Figures	I

1 Question 4.1

The assignment will only look into the salinity and the dissolved oxygen. Two time series are defined then :

S_t : salinity

DO_t : dissolved oxygen

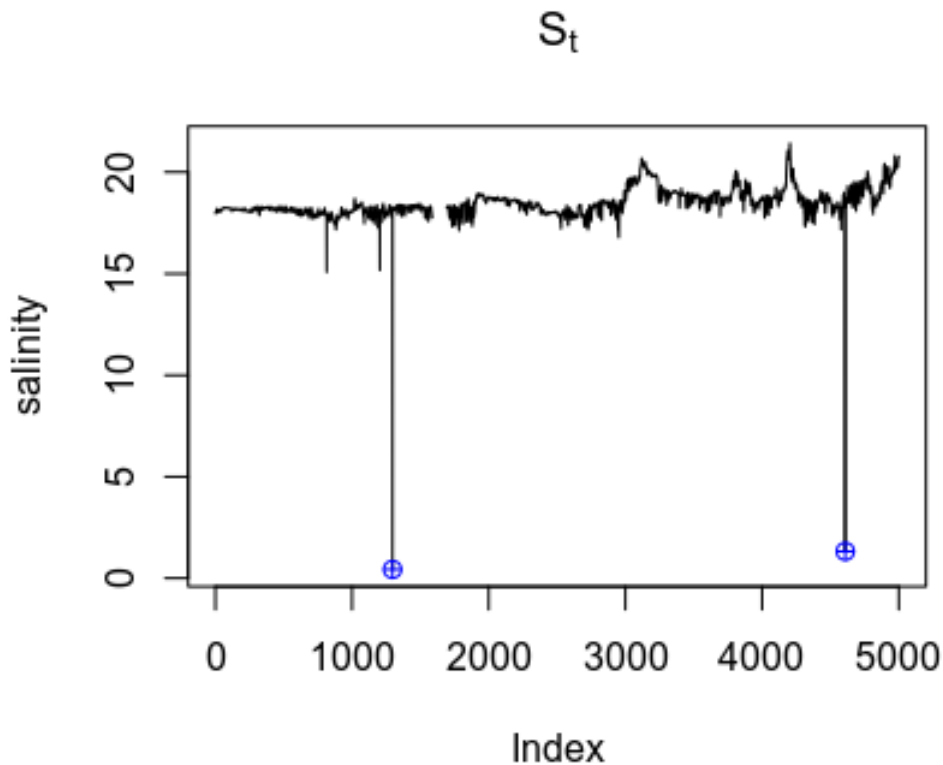


Figure 1: plotting data

It can be seen that the salinity data has some outliers. One can ask whether the first 2 points that are not following the same trend around index 1000 are outliers or no. If one considers only the first part of the plot (till index 2500), one will say that they are outliers. Considering the second part of plot, one might consider them normal because there are several points having more salinity than the rest majority. Still more likely at this point, these two first points are most likely outliers.

It can also be seen the missing data around index 1500

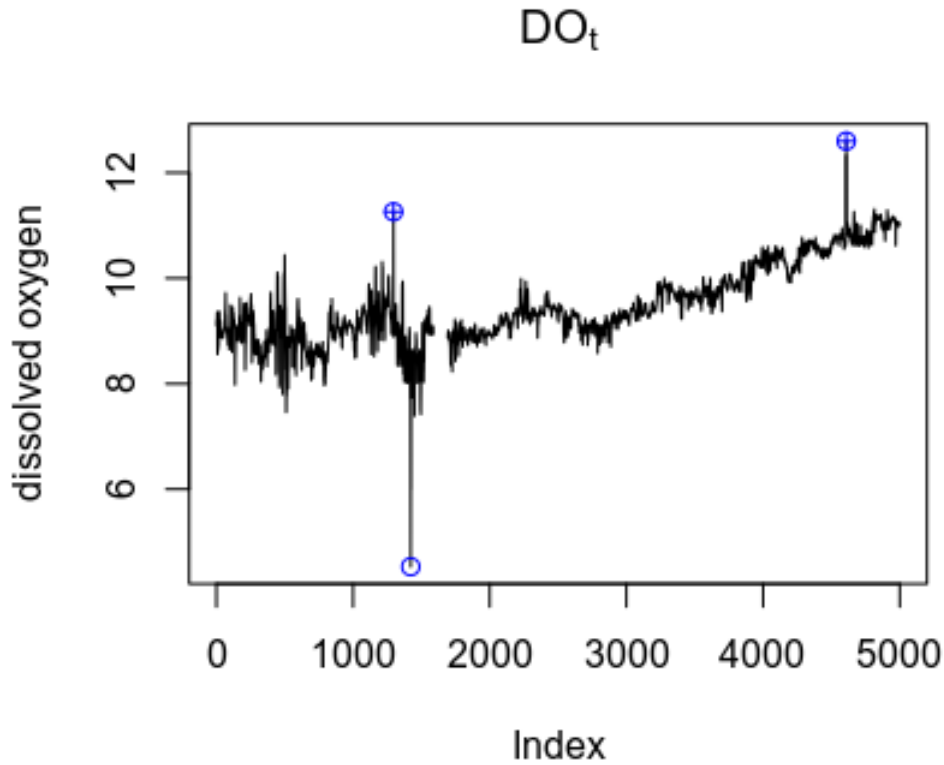


Figure 2: plotting data

There are some outliers also in the second plot. There are also missing data around index 1500. Going back to the data in csv form, it can be said that the two times series have the gap at exact same moments. In fact, from index 1588 to 1692 S_t and DO_t are NAs. What is new in the second plot is a linear increase from around index 3000.

2 Question 4.2

According to 10.1 from the book p284:

$$X_t = \mathbf{A}_t X_{t-1} + \mathbf{B}_t u_{t-1} + e_t \quad (1)$$

$$S_t = \mathbf{C}_t X_t + m_t \quad (2)$$

such that X_t : state of salinity

$A_t=1$, $B_t=1$ because there is no input, the latent noise is

$$e_t \sim \mathcal{N}(0, \Sigma_1) \quad (3)$$

S_t observed (i.e. measured) salinity (denotes the observations)
 $C_t=1$ and the noise of the measurable output S_t is

$$m_t \sim \mathcal{N}(0, \Sigma_1) \quad (4)$$

e_t is white noise. m_t is white noise. They are uncorrelated.

It should be mentioned that in the model, the state space and observations are univariate.

3 Question 4.3

To plot the one step prediction error, the missing data is substituted by the corresponding predictions.

To calculate the 95% prediction interval, the following formula is used

$$X_{t+1|t} \pm 1.96 \cdot \sqrt{\Sigma_{t+1|t}^{ss}} \quad (5)$$

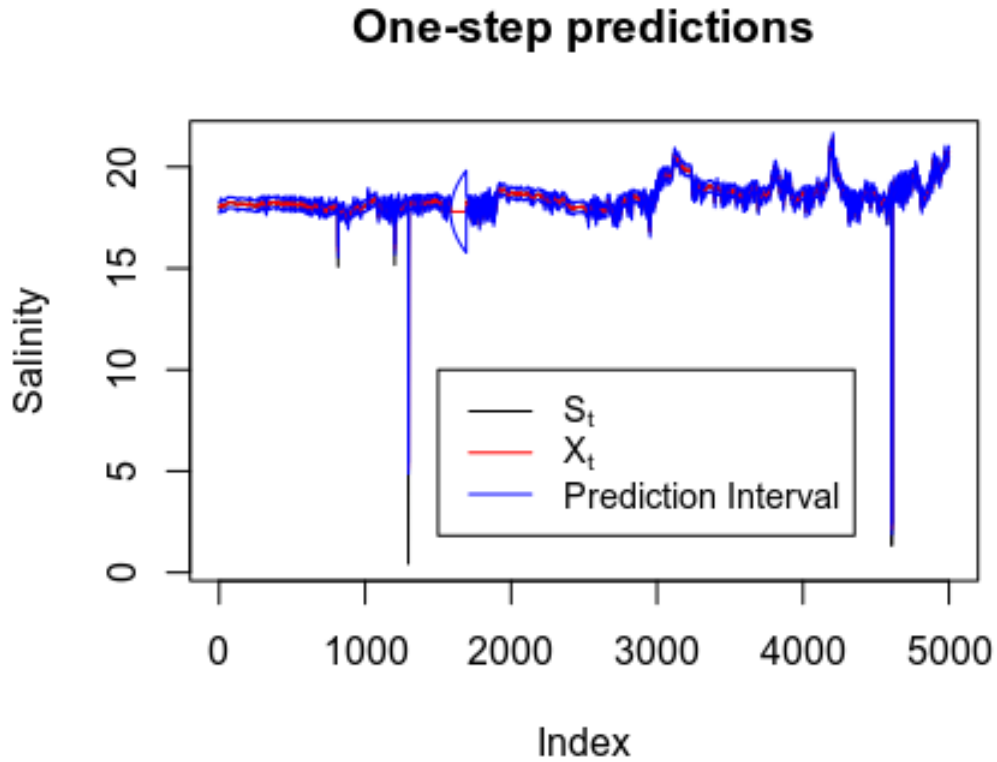


Figure 3

1)

The predicted values are within the confidence interval. It is hard too see it from this plot. It is much easier to refer to the zoom plot next to see this.

One can observe that the variance of the prediction interval increased around index 1500. This is due to the missing data there.

It is seen that the prediction intervals goes with the outliers i.e they follow them.

2)

$\hat{X}_{t+1|t} = E[X_{t+1}]$ is the one-step prediction of X_t

The definition of one step prediction error is:

$$\tilde{S}_{t+1|t} = S_{t+1} - \hat{X}_{t+1|t} \quad (6)$$

its variance:

$$\text{Var}[S_{t+1}|S_0, \dots, S_t] = \text{Var}[\tilde{S}_{t+1}] = \Sigma_{t+1|t}^{ss} \quad (7)$$

So the standardized one step prediction errors are defined as:

$$\tilde{S}_{t+1|t} = \frac{S_{t+1} - \hat{X}_{t+1|t}}{\sqrt{\Sigma_{t+1|t}^{ss}}} \quad (8)$$

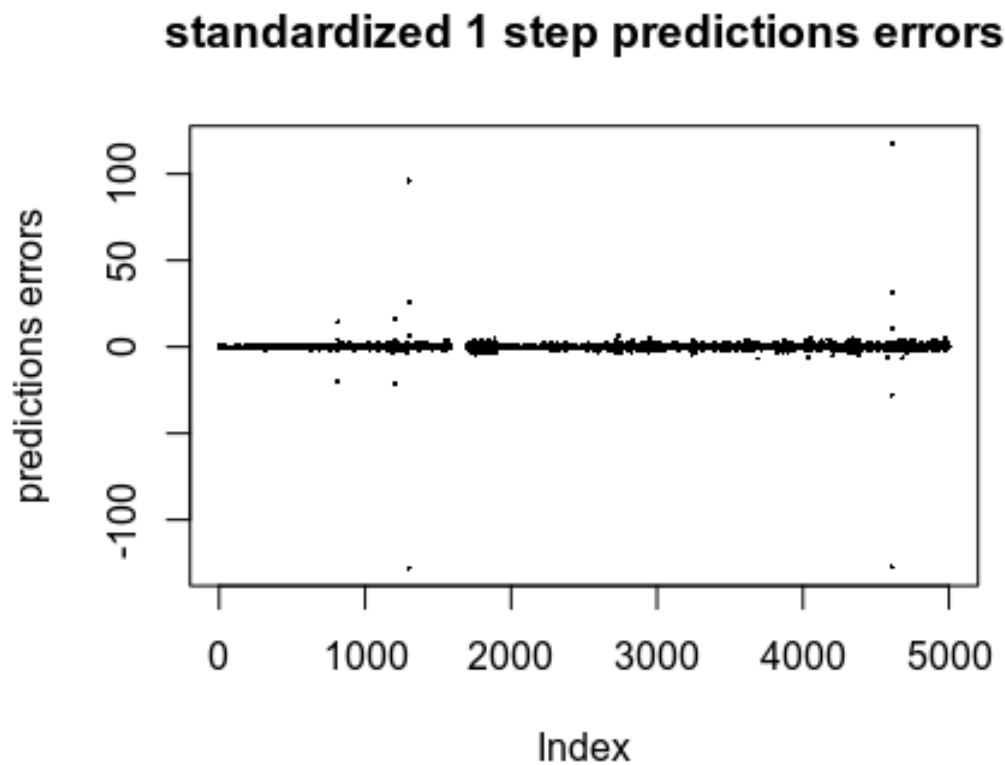


Figure 4

Again, the gap in data caused a gap in the one step prediction errors plot. One can also notice that each time there is an outlier in the salinity plot, there are 2 points in predictions errors that do not follow the rest of points which are in 0. It seems that the absolute value of the positive point is smaller than the negative one.

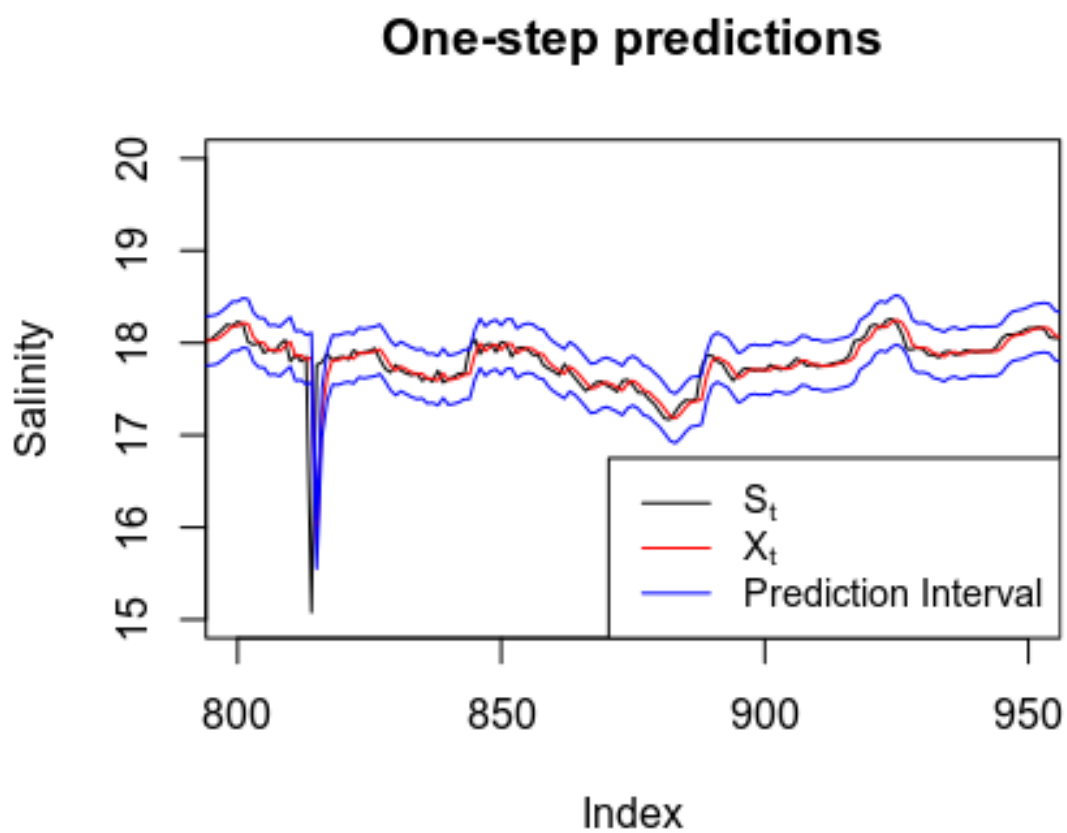


Figure 5

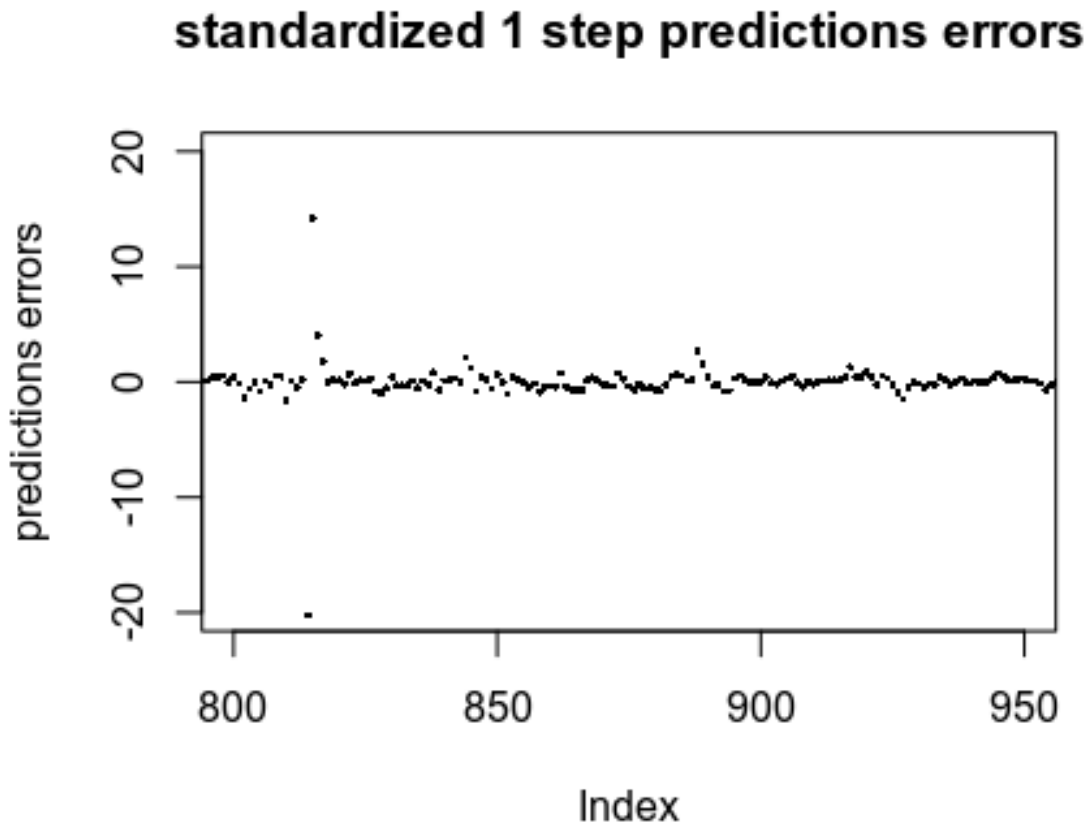


Figure 6

3)

With the zoom, one can see the good fit of confidence interval i.e inside it there is the predictions except for outliers. In fact it can be seen how the prediction interval follows the outlier down.

The outlier in black in S_t cause the next prediction and prediction interval to go down. As a result, the prediction error goes down also (negative point), then it goes up (positive point). The absolute negative value is greater than the positive one.

4)

The final state of filter is given by:

$$\hat{X}_{5000} = 20.75815 \quad (9)$$

and

$$\Sigma_{5000|5000}^{xx} = 0.003660254 \quad (10)$$

4 Question 4.4

Please refers to the function "Kalman_filter" in appendix in code.

1)

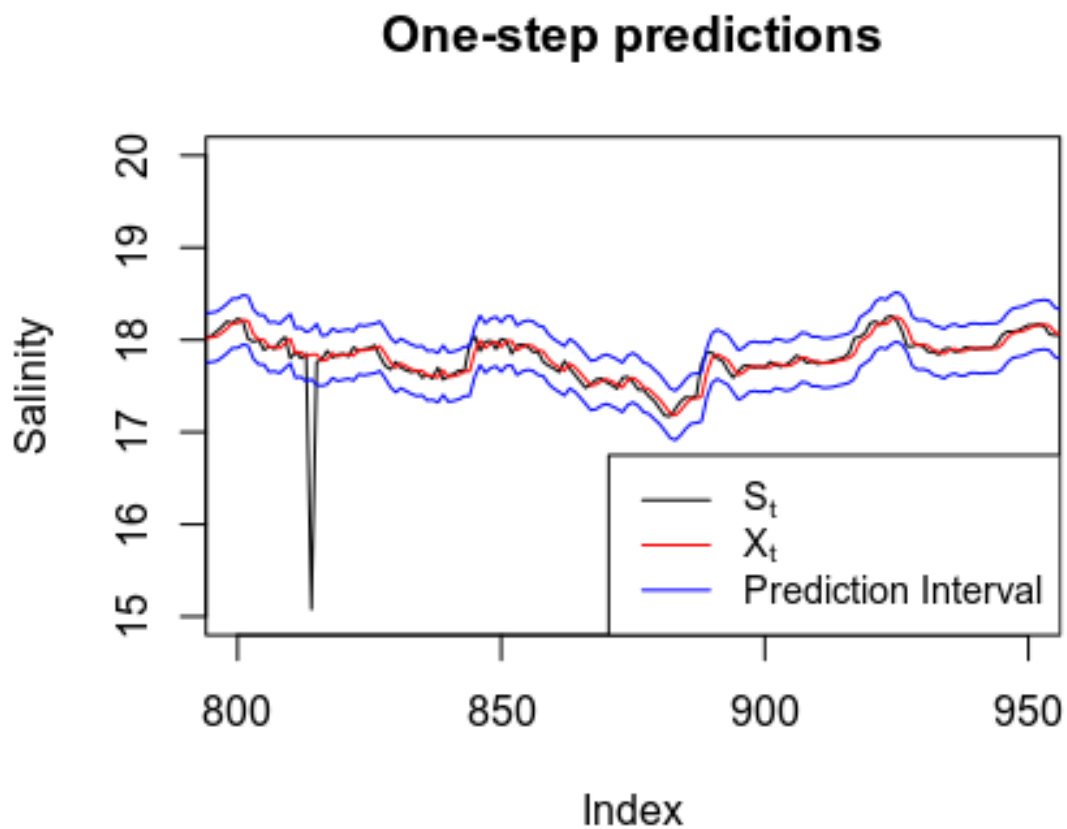


Figure 7

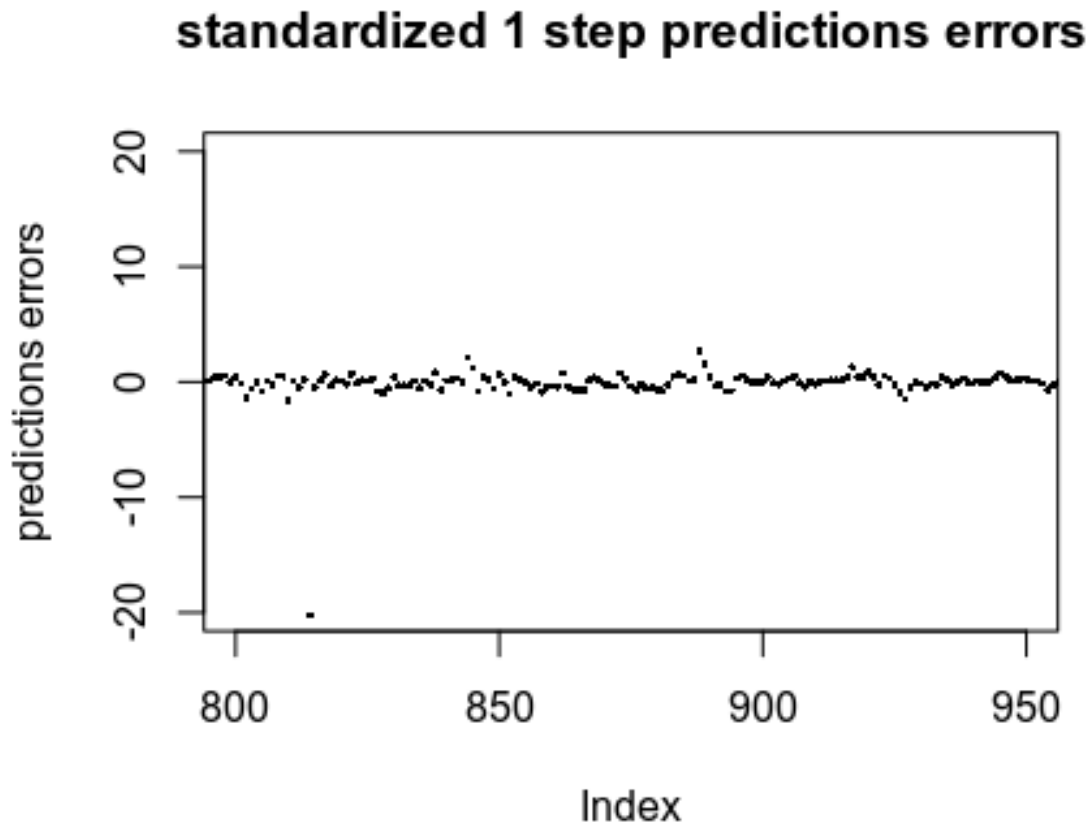


Figure 8

1)

It can be seen that there is no more down-up : the negative point followed by positive point in predictions errors. This confirms that that was due to outliers. 2)

The first five detected outliers have the following indexes:

814, 1203, 1296, 2733, 3692

3)

The number of skipped observation are 13. Note that this number does not include the NAs in data which are 111. Hence, if it is wanted to consider those too, the answer will be $13+111=124$

4)

The final state of filter is given by:

$$\hat{X}_{5000} = 20375815 \quad (11)$$

and

$$\Sigma_{5000|5000}^{xx} = 0.003660254 \quad (12)$$

It is the same final state as found in question 4.3) This question is remained without answer, in other words does the outliers not affect the final state of filter ?

5 Question 4.5

6 Question 4.6

The physical models where there are production, consumption of a quantity are numerous. A famous example is the model of temperature in a house. In this section, the results are basically obtained by analogy to that model.

1 and 2)

$$DO_t = DO_{t-1} + \alpha_{prod}I_{t-1} + \alpha_{cons}R_{t-1} + \alpha_{exch}(DO_{sat,t-1} - DO_{t-1}) + \epsilon_{DO,t} \quad (13)$$

$$R_t = R_{t-1} + \epsilon_{R,t} \quad (14)$$

3)

The state space model for dissolved oxygen defined by :
System equation:

$$X_t = \mathbf{A}_t X_{t-1} + \mathbf{B}_t u_{t-1} + e_t \quad (15)$$

where the state is $X_t = \begin{bmatrix} DO_t \\ R_t \end{bmatrix}$, this gives:

$$\mathbf{X}_t = \begin{bmatrix} DO_t \\ R_t \end{bmatrix} = \begin{bmatrix} 1 - \alpha_{exch} & \alpha_{cons} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} DO_{t-1} \\ R_{t-1} \end{bmatrix} + \begin{bmatrix} \alpha_{prod} & \alpha_{exch} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} I_{t-1} \\ DO_{sat,t-1} \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} e_{DO,t} \\ e_{R,t} \end{bmatrix} \quad (16)$$

Observation equation:

$$Y_t = \mathbf{C}_t X_t + \epsilon_{Y,t} \quad (17)$$

$$(18)$$

$$Y_t = [1, 0] \mathbf{X}_t + \epsilon_{Y,t} = [1, 0] \begin{bmatrix} DO_t \\ R_t \end{bmatrix} + \epsilon_{Y,t} = DO_t + \epsilon_{Y,t} \quad (19)$$

$$\Sigma_1 = \Sigma_{DO,R} = \begin{bmatrix} \sigma_{DO}^2 & 0 \\ 0 & \sigma_R^2 \end{bmatrix} \quad (20)$$

$$\Sigma_2 = \Sigma_Y = \sigma_Y^2 \quad (21)$$

The assumptions made here are :

$\epsilon_{DO,t}$ is white noise. $\epsilon_{R,t}$ is white noise. They are uncorrelated.

$\epsilon_{Y,t}$ is white noise. It is uncorrelated with $\epsilon_{DO,t}$ and $\epsilon_{R,t}$

Explanation of the procedure by analogy in eq(13):

a) DO_t is the dissolved oxygen at time t, equivalent to the temperature inside the house $T_{i,t}$

b) $\alpha_{prod}I_{t-1}$ this is the production of oxygen from chlorophyll and sunlight i.e photosynthesis. This is equivalent to the production of temperature via radiator in house. As a result I_{t-1} is considered as the first input to the model and α_{prod} denotes the amount of production.

c) $\alpha_{cons}R_{t-1}$ is the consumption of oxygen i.e respiration. α denotes its quantity and should be negative to denote that it is consumption i.e a quantity going out from the system. This is similar to any cooling event in the house, for instance if someone opens the door for seconds. R_{t-1} is modelled as random walk.

d) $\alpha_{exch}(DO_{sat,t-1} - DO_{t-1})$ is the exchange of oxygen with the atmosphere so that the dissolved oxygen concentration approaches the saturation concentration. ($DO_{sat,t-1}$ is the second input. This is similar to $\alpha_{exch}(T_{atm,t-1} - T_{t-1})$ where the atmosphere temperature is the second input and denoted $T_{atm,t-1}$ equivalent to ($DO_{sat,t-1}$ here.

4)

It is believed that the answer to this question is probably obtained when coding this question. One should verify the assumption made above about errors . Sadly and due to lack of time, I did not do it.

7 Code

```
1
2 rm(list=ls())
3
4 library(data.table)
5 library("dlm")
6 library("FKF")
7 library(scales)
8 library("plot3D")
9 library(plotly)
10 library(numDeriv)
11 library(MASS)
12 library(mvtnorm)
13 require(graphics)
14
15 #Question1
#####
```

```
16 data = read.csv(file = "/home/ghassen97/Desktop/S8/time series analysis/
    assignment/assignment 4/A4_Kulhuse.csv", fill = TRUE, header = TRUE)
17
18 S = data$Sal
19 DO = data$ODO
20 DT = data$DateTime
21 i = 1:length(S) # Time
22
23 plot(S, type = 'l', main = expression(S[t]), ylab='salinity')
24 points(c(1296, 4609),c(0.43,1.32), col = 'blue', pch = 10)
25 plot(DO, type = 'l', main = expression(DO[t]), col = 'black',ylab='dissolved
    oxygen')
26
27 points(c(1422), c(4.53), col = 'blue', type = 'p')
28 points(c(4609, 1296),c(12.60, 11.26), col = 'blue', pch = 10)
29
30 #Question2 see report
    #####
31
32 #Question3
    #####
33
34 Kalman_filter = function(obs_condition, outlier_condition) {
35     S = data$Sal # S_t
36     N = length(S)
37     X = vector(length = N) # State of salinity (state vector/number)
38     X_pred = vector(length = N) # One-step predictions of X
39     S_xx = vector(length = N) # This is Conditional covariance of X given all
        previous S
40     S_yy = vector(length = N) # This is Variance of S given all previous S
41     S_xx_pred = vector(length = N) # one step prediction variance
42     S_yy_pred = vector(length = N) # one step prediction variance
43     K = vector(length = N) # Kalman gain
44     outlier_ind = vector(length = N) # Outlier index
45     S_t = vector(length = N)
46     error1 = vector(length = N)
47     i = 1:N # Time
48     V_e = 0.01 #system variance
49     V_m = 0.005 #observation variance
50
51     # Initial Condition
```

```
52 X[1] = S[1]
53 X_pred[1] = S[1]
54 S_xx[1] = V_e
55 S_yy[1] = V_e + V_m
56 S_xx_pred[1] = V_e
57 S_yy_pred[1] = V_e + V_m
58
59 # Kalman Filter
60 for (t in i) {
61
62     error1[t] = S[t] - X_pred[t]
63
64     # fix missing observations (if obs_condition = 1) fix outliers (if
        outlier_condition = 1). This if statement skips the reconstruction
        step
65     # if obs_condition == TRUE then fix missing observations, if
        outlier_condition == TRUE then skip outliers.
66     # There is no reconstruction in this if below
67     if((S[t] %in% NA && obs_condition == TRUE ) | (abs(error1[t]) > 6*sqrt(
        S_xx[t]) && outlier_condition == TRUE)){
68         outlier_ind[t] = t
69         X[t+1] = X[t]
70         X_pred[t+1] = X[t+1]
71         S_xx[t+1] = S_xx[t] + V_e
72         S_yy[t+1] = S_xx[t+1] + V_m
73         S_xx_pred[t+1] = S_xx[t]
74         S_yy_pred[t+1] = S_xx[t+1]
75     }else{
76
77
78         K[t] = S_xx[t]/S_yy[t] #Kalman gain
79
80         #Reconstruction , updating
81         X[t] = X_pred[t] + K[t]*(S[t] - X_pred[t]) # X_t|t
82         S_xx[t] = S_xx[t] - K[t]*S_xx[t]
83
84         #Prediction
85         X[t+1] = X[t]
86         X_pred[t+1] = X[t+1]
87         S_xx[t+1] = S_xx[t] + V_e
88         S_yy[t+1] = S_xx[t+1] + V_m
89         S_xx_pred[t+1] = S_xx[t]
90         S_yy_pred[t+1] = S_xx[t+1]
91     }
```

```
92 }
93
94 CI_upper = X_pred + 1.96*sqrt((S_yy)) # 95% confidence interval upper
95 CI_lower = X_pred - 1.96*sqrt((S_yy)) # 95% Confidence interval lower
96
97 outlier_ind = outlier_ind[!is.na(S)]
98 outlier_ind = outlier_ind[outlier_ind != 0]
99 # with fkf implementation, check if Kalman_filter is ok
100 mod = fkf(a0 = c(S[1]), P0 = matrix(V_e), dt = matrix(0), ct = matrix(0),
101           Tt = matrix(1), Zt = matrix(1),
102           HHt = matrix(V_e), GGt = matrix(V_m), yt = matrix(t(S), nrow =
103               1, ncol = length(S)), check.input = TRUE)
104 check_fkf_1 = sum(X_pred - mod$at) #check_fkf_1=5,329071e-15
105 #print(check_fkf_1)
106 check_fkf_2 = sum(X[i] - mod$att)#check_fkf_2=5,329071e-15
107 print(check_fkf_2)
108 error2 = S - X_pred[i]
109
110 # Return Values
111 newlist = list("X" = X, "X_pred" = X_pred, "S_xx" = S_xx, "S_yy" = S_yy,
112               "CI_upper" = CI_upper, "CI_lower" = CI_lower, "fkf" = mod,
113               "check_fkf_1" = check_fkf_1,
114               "check_fkf_2" = check_fkf_2, "S_xx_pred" = S_xx_pred, "
115               S_yy_pred" = S_yy_pred,
116               "error1" = error1, "error2" = error2, 'outlier_ind' =
117               outlier_ind, "K" = K)
118
119 return(newlist)
120
121 }
122 #1)#####
123 K_filter_1 = Kalman_filter(TRUE, FALSE)
124 plot(S, type = 'l', col = 'black', main = "One-step predictions", ylab = "
125     Salinity")
126 lines(K_filter_1$X_pred, type = 'l', col = 'red')
127 lines(K_filter_1$CI_upper, type = 'l', col = 'blue')
128 lines(K_filter_1$CI_lower, type = 'l', col = 'blue')
129 legend(1500, 10, legend = c(expression(S[t]), expression(X[t]), "Prediction
130     Interval"), col = c("black", "red", "blue"), lty=1:1, cex = 0.9)
131 #2)#####
132
133 i = 1:length(S)
134
135 S_pred_error = (S - K_filter_1$X_pred[i])/sqrt(K_filter_1$S_yy[i])
```



```
128 plot(S_pred_error, type = 'p', ylab = 'predictions errors', main = "
      standardized 1 step predictions errors",pch=16,cex=0.3)
129
130 #3)#####
131
132 K_filter_1 = Kalman_filter(TRUE, FALSE)
133
134 plot(S, type = 'l', col = 'black', main = "One-step predictions", ylab = "
      Salinity",xlim = c(800,950), ylim = c(15,20))
135 lines(K_filter_1$X_pred, type = 'l', col = 'red')
136 lines(K_filter_1$CI_upper, type = 'l', col = 'blue')
137 lines(K_filter_1$CI_lower, type = 'l', col = 'blue')
138 legend("bottomright", legend = c(expression(S[t]), expression(X[t]), "
      Prediction Interval"), col = c("black", "red", "blue"), lty=1:1, cex =
      0.9)
139
140
141 S_pred_error = (S - K_filter_1$X_pred[i])/sqrt(K_filter_1$S_yy[i])
142 plot(S_pred_error, type = 'p', col = 'black', xlim = c(800,950), ylim = c
      (-20,20), ylab = 'predictions errors', main = "standardized 1 step
      predictions errors",pch=16,cex=0.4)
143
144 #4)#####
145 #using function to get X
146 K_filter_1$X[5000]
147 K_filter_1$S_xx[5000]
148
149 # using library FKF
150 K_filter_1$fkf$att[5000] #same values
151 K_filter_1$fkf$Ptt[5000]
152
153 #Question4
      #####
154
155 #1)#####
156 K_filter_no_outlier = Kalman_filter(TRUE, TRUE)
157
158 plot(S, type = 'l', col = 'black', main = "One-step predictions", ylab = "
      Salinity",xlim = c(800,950), ylim = c(15,20))
159 lines(K_filter_no_outlier$X_pred, type = 'l', col = 'red')
160 lines(K_filter_no_outlier$CI_upper, type = 'l', col = 'blue')
161 lines(K_filter_no_outlier$CI_lower, type = 'l', col = 'blue')
162 legend("bottomright", legend = c(expression(S[t]), expression(X[t]), "
      Prediction Interval"), col = c("black", "red", "blue"), lty=1:1, cex =
```

```

0.9)
162
163 S_pred_error = (S - K_filter_no_outlier$X_pred[i])/sqrt(K_filter_no_outlier$
    S_yy[i])
164 plot(S_pred_error, type = 'p', col = 'black', xlim = c(800,950), ylim = c
    (-20,20), ylab = 'predictions errors', main = "standardized 1 step
    predictions errors",pch=16,cex=0.4)
165
166 #2)#####
167 K_filter_no_outlier$outlier_ind[1:5] #Outliers from the output of tje kalman
    filter K_filter_no_outlier = Kalman_filter(1,0,1)
168 #3)#####
169 length(K_filter_no_outlier$outlier_ind) #Number of skipped observations
    excluding the missing observations (NA)
170 #4)#####
171 K_filter_no_outlier$X[5000]
172 K_filter_no_outlier$S_xx[5000]
173
174 #Question5
    #####
175
176 #Question6
    #####
177 #see report

```

List of Figures

1	plotting data	1
2	plotting data	2
3	4
4	5
5	6
6	7
7	8
8	9