

Exercise 3

TMA4300

Erling Fause Steen, Christian Oppegård Moen

Spring 2022

Contents

Introduction	2
Problem A	2
1	2
2	5
Problem B	5
1	6
2	7
3	7

Introduction

Problem A

In this problem we will analyze a dataset which contain a sequence of length $T = 100$ of a non-Gaussian time-series for which we will compare two different parameter estimators. Consider the AR(2) specified by the relation

$$x_t = \beta_1 x_{t-1} + \beta_2 x_{t-2} + e_t, \quad (1)$$

where e_t are independent and identically distributed (iid) random variables with zero mean and constant variance. Also, consider the loss functions with respect to $\beta = [\beta_1, \beta_2]^T$ given by

$$Q_{LS}(\mathbf{x}) = \sum_{t=3}^T (x_t - \beta_1 x_{t-1} - \beta_2 x_{t-2})^2$$
$$Q_{LA}(\mathbf{x}) = \sum_{t=3}^T |x_t - \beta_1 x_{t-1} - \beta_2 x_{t-2}|.$$

Then, the least sum residuals (LS) and least sum of absolute residuals (LA) are obtained by minimizing $Q_{LS}(\mathbf{x})$ and $Q_{LA}(\mathbf{x})$ respectively. We denote the minimisers by $\hat{\beta}_{LS}$ and $\hat{\beta}_{LA}$, and define the estimated residuals by $\hat{e}_t = x_t - \hat{\beta}_1 x_{t-1} - \hat{\beta}_2 x_{t-2}$ for $t = 3, \dots, T$ with mean \bar{e} . Then, $\hat{\varepsilon} = \hat{e} - \bar{e}$ is re-centered to have mean zero.

1

Now we will use the residual resampling bootstrap method to evaluate the relative performance of the two parameter estimators, $\hat{\beta}_{LS}$ and $\hat{\beta}_{LA}$, which are calculated by the given function `ARp.beta.est`. We consider the variance and mean of β_{LS}^* and β_{LA}^* obtained by minimizing $Q_{LS}(\mathbf{x}^*)$ and $Q_{LA}(\mathbf{x}^*)$ for bootstrap sample \mathbf{x}^* , respectively. The full method is shown in Algorithm 1. In Table 1 the variance of β_{LA} is slightly lower than that of β_{LS} , and the absolute value of the bias is slightly lower as well. This suggest that the LS estimator is not optimal.

Algorithm 1

```
1: for LS and LA estimator of  $\beta$  do
2:   for  $t = 2, \dots, n$  do
3:     Define estimated residuals.
4:   end for
5:   Resample  $n + 1$  pseudo residuals from the estimated residuals with replacement.
6:   Generate pseudo data.
7:   for each bootstrap sample of pseudo data do
8:     estimate  $\hat{\beta}^*$ .
9:   end for
10: end for
```

```
x = data3A$x
plot(x, type = "l", xlab = "t", ylab = "x")
```

```
detach("package:dplyr", unload = T)
rsBoot = function(B, x, p = 2) {
  T = length(x)
```

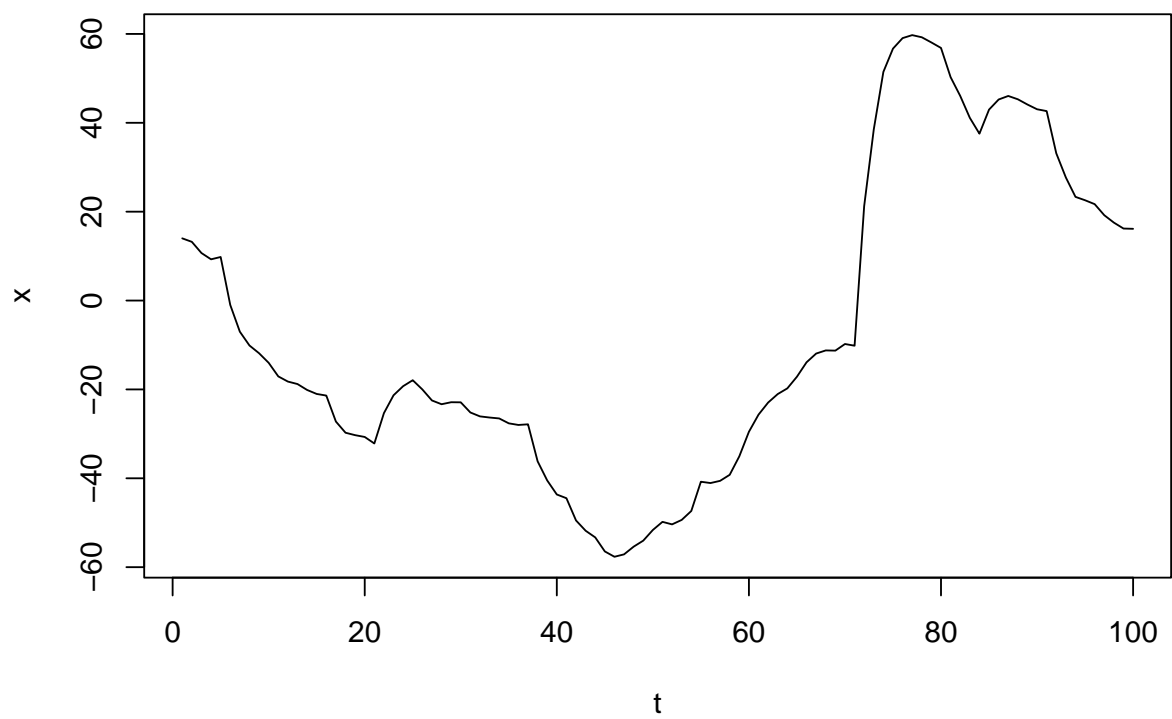


Figure 1: Given dataset. write more

```

# Estimate beta
beta.hat = ARp.beta.est(x, p)
# calculate observed residuals of AR(p) seq
e.LS.observed = ARp.resid(x, beta.hat$LS)
e.LA.observed = ARp.resid(x, beta.hat$LA)
# Allocate memory
beta.LS.star = matrix(nrow = B, ncol = 2)
beta.LA.star = matrix(nrow = B, ncol = 2)
# e.LS.star = vector(mode = 'double', length = B) x.LS.star = matrix(nrow =
# B, ncol = 2) x.LA.star = matrix(nrow = B, ncol = 2)

# Bootstrap
for (b in 1:B) {
  # Resample from observed residuals to yield pseudo innovations
  e.LS.star = sample(e.LS.observed, size = T, replace = TRUE)
  e.LA.star = sample(e.LA.observed, size = T, replace = TRUE)
  # Generate pseudo data (Timeseries based on sampled residuals and
  # beta.hat)
  i = sample(T - 1, 1)
  x.LS.star = ARp.filter(x[c(i, i + 1)], beta.hat$LS, e.LS.star)
  x.LA.star = ARp.filter(x[c(i, i + 1)], beta.hat$LA, e.LA.star)
  # compute beta star
  beta.LS.star[b, ] = ARp.beta.est(x.LS.star, 2)$LS
  beta.LA.star[b, ] = ARp.beta.est(x.LA.star, 2)$LA
}
return(list(beta.hat = beta.hat, beta.LS = beta.LS.star, beta.LA = beta.LA.star,
  e.LS = e.LS.observed, e.LA = e.LA.observed, x.LS = x.LS.star, x.LA = x.LA.star))
}

B = 1500
set.seed(420)
boot = rsBoot(B, x)

# Compute variance and bias

beta.var = rbind(apply(boot$beta.LS, 2, var), apply(boot$beta.LA, 2, var))
beta.bias = rbind(apply(boot$beta.LS, 2, mean) - boot$beta.hat$LS, apply(boot$beta.LA,
  2, mean) - boot$beta.hat$LA)
beta.vb = cbind(beta.var, beta.bias)
rownames(beta.vb) = c("LS", "LA")
# colnames(beta.vb) = c('var_1', 'var_2', 'bias_1', 'bias_2')

kable(round(beta.vb, 6), caption = "Variance and mean for  $\beta_1$  and  $\beta_2$ ",
  format = "latex") %>%
  kable_styling() %>%
  add_header_above(c(" ", "Var_{1}", "Var_{2}", "bias_{1}", "bias_{2}"),
    escape = FALSE)

```

Table 1: Variance and mean for β_1^* and β_2^* .

	Var ₁	Var ₂	bias ₁	bias ₂
LS	0.005451	0.005240	-0.012071	0.005988
LA	0.000362	0.000358	-0.001311	0.000829

Table 2: Prediction interval of x_{101} .

	2.5%	97.5%
LS	10.01	23.34
LA	7.21	23.67

2

Next, we will compute a 95% prediction interval for x_{101} for both the LS and the LA estimator. That is, we bootstrap sample B times

$$x_{101} = \beta_1^* x_{100} + \beta_2^* x_{99} + \hat{e}_i \quad (2)$$

where β_1^* and β_2^* are sampled from all 1500 bootstrap samples found in section A1, and \hat{e}_i is sampled from the estimated residuals \hat{e}_t for $t = 2, \dots, 100$. The prediction intervals are found in Table 2.

```
set.seed(420)
# Bootstrap x_101 using beta.star and e.observed
x_101.LS = vector(mode = "double", length = B)
x_101.LA = vector(mode = "double", length = B)
for (b in 1:B) {
  # i.e = sample(98, 1)
  i.beta = sample(1500, 1)
  # x_101.LS[b] = boot$beta.LS[i.beta,] %*% x[c(100,99)] + boot$e.LS[i.e]
  # x_101.LA[b] = boot$beta.LS[i.beta,] %*% x[c(100,99)] + boot$e.LA[i.e]
  x_101.LS[b] = boot$beta.LS[i.beta, ] %*% x[c(100, 99)] + sample(boot$e.LS, 1)
  x_101.LA[b] = boot$beta.LS[i.beta, ] %*% x[c(100, 99)] + sample(boot$e.LA, 1)
}

pi.boot = rbind(LS = quantile(x_101.LS, probs = c(0.025, 0.975)), LA = quantile(x_101.LA,
  probs = c(0.025, 0.975)))
kable(round(pi.boot, 2), caption = "Prediction interval of $x_{101}$.")
```

Problem B

We will investigate the concentration of bilirubin (mg/dL), which is a breakdown of haemoglobin, in blood samples taken from three young men shown in Table 3.

Individual	Concentration (mg/dL)										
1	0.14	0.20	0.23	0.27	0.27	0.34	0.41	0.41	0.55	0.61	0.66
2	0.20	0.27	0.32	0.34	0.34	0.38	0.41	0.41	0.48	0.55	
3	0.32	0.41	0.41	0.55	0.55	0.62	0.71	0.91			

Table 3: Concentration of bilirubin in three different individuals.

1

The logarithms of the concentration for each individual are displayed by the boxplot in Figure 2. We see that the mean logarithm of bilirubin is smallest for individual 1 and highest for individual 2. The standard deviation is smallest for individual 2 while individual 1 have slightly larger standard deviation than individual 3.

```
ylim = c(-2.1, 0)
bilirubin <- read.table("./additionalFiles/bilirubin.txt", header = T)
boxplot(log(meas) ~ pers, data = as.data.frame(bilirubin), xlab = "Person", ylab = "log(meas)",
        ylim = ylim)
```

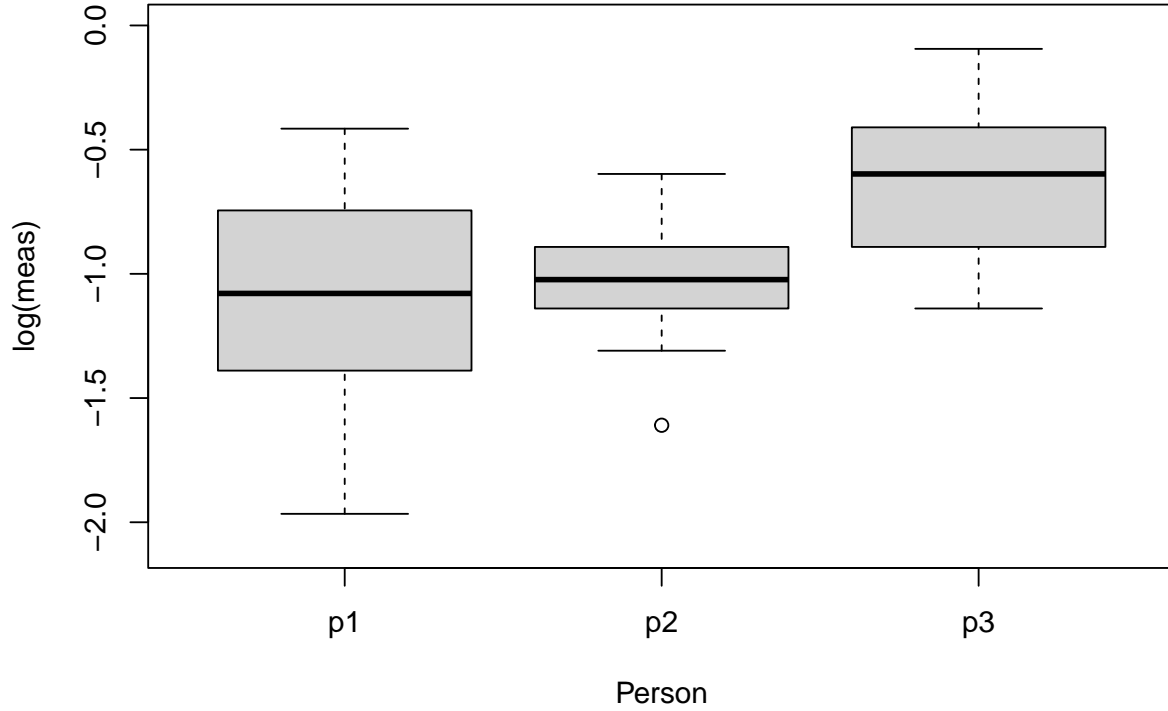


Figure 2: The logarithm of `meas` for each individual.

```
fitB1 = lm(log(meas) ~ pers, data = bilirubin)
fitB1.Fstat = summary(fitB1)$fstatistic
Fval = fitB1.Fstat["value"]
Pval = pf(Fval, fitB1.Fstat[2], fitB1.Fstat[3], lower.tail = F)
```

Let $\log Y_{ij} = \log(\text{meas}_{ij})$ for individual $i = 1, 2, 3$ and observation $j = 1, \dots, n_i$, where $n_1 = 11$, $n_2 = 10$ and $n_3 = 8$. Consider the regression model

$$\log Y_{ij} = \beta_i + \varepsilon_{ij}, \quad (3)$$

where $\varepsilon_{ij} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$. Then, the F-statistic on $3 - 1$ and $29 - 3$ degrees of freedom of the hypothesis that $\beta_1 = \beta_2 = \beta_3$ is 4 with p-value 0.03946. With the F-statistic we are investigating whether the bilirubin of

each individual share significantly equally distributed. Judging from the p-value we reject the hypothesis on a $\alpha = 0.05$ -level. That is, on a $\alpha = 0.05$ -level, the individuals do not share the same bilirubin.

2

To further investigate the individuals bilirubin we will consider a permutation test where the idea is that shuffling the labels will not change the joint distribution of the data. That is, we assume that the distributions are equal so that shuffling the order of the data to generate bootstrap samples should be valid. This is done by first randomly assigning bilirubin values, `meas`, to the three individuals, then generating $B = 999$ bootstrap samples. Implementation of the permutation and computation of the F-statistic is seen in the chunk below.

```
permTest = function(data = bilirubin) {  
  bilirubin$perm = bilirubin$meas[sample(29, 29)]  
  return(summary(lm(log(perm) ~ pers, data = bilirubin))$fstatistic["value"])  
}
```

3

We generate $B = 999$ samples of the F-statistic, `Fvals`, using the function `permTest`. Then we compute the p-value defined by $Pval = (Fvals \geq Fval) / B$, where `Fval` is the F-statistic found in section B1.

```
set.seed(420)  
B = 999  
Fvals = vector(mode = "double", length = 999)  
for (b in 1:B) {  
  Fvals[b] = permTest()  
}  
Pval.perm = sum(Fvals >= Fval) / B
```

The p-value of the 999 F-statistics for this run was 0.03704 which also reject the hypothesis on a $\alpha = 0.05$ -level. Comparing it to the p-value from section B1 we see that they are almost equal, with the value from the permutation test being slightly lower.