

Time Series Analysis

UC Berkeley, School of Information: MIDS w271

2023-07-24

Contents

Problem Set 11	5
1 (10 points) ARIMA model	7
1.1 Time series plot	7
1.2 Fit a model	8
1.3 Is this model appropriate?	9
1.4 Forecasts, by hand!	10
1.5 Interpret roots	10
2 (10 points) Seasonal ARIMA model	13
2.1 Time series plot	14
2.2 Check for Stationary	14
2.3 Model identification and estimation	15
2.4 Model diagnostic	18
2.5 Forecasting	20
3 (10 points) Time Series Linear Model and Cointegration	21
3.1 Plot electricity	21
3.2 Cointegration test	22
3.3 Fit Model	22
3.4 Residuals Plot	22
3.5 Forecasting model	22
4 (12 points) Vector autoregression	23
4.1 Time series plot	23
4.2 Check for the unit root	23
4.3 Determine VAR model	24
4.4 Estimation	24
4.5 Model diagnostic	24
4.6 Forecasting	24

Problem Set 11

This is the final problem set that you will have to work on for this class. Congratulations! (Although there is still a group lab that will be the final assignment in the course.)

You will start with some guided work, and then proceed into less structured work that will let you stretch and demonstrate what you have learned to date.

Notice, in particular, that the last few questions are asking you essentially to “produce a model” using a method. At this point in the course, you should be familiar with many of the model forms that you *might* fit; and, you are familiar with methods that you can use to evaluate models’ performances. In these questions, we are asking you to, essentially, fit a good model with a method and then to evaluate how a good model with “this” method is doing compared to another good model with “that” method.

In several of these questions, there isn’t a correct answer, *per se*. Instead, there is the process that you will undertake and record as you are producing your argument for the model that you think is best meeting your objectives. This is a **very** applied task that we anticipate you will see many times in your work.

```
knitr::opts_chunk$set(echo=TRUE)
```

We are providing you with an additional challenge, but one that is also very evocative of work that you’re likely to come across. This is a well-built repository, that uses a well-documented framework to produce reports, namely **bookdown**.

Once you have done your work, you can render the entire book using the following call in your console:

```
> bookdown::render_book()
```

This will ingest each of the files `01-time_series_...`, `02-cross_validation.Rmd`, `03-ARIMA_model.Rmd` and so on ... and produce a PDF that is stored in `./_book/_main.pdf`. If you would like to read more about this framework, you can do so at the following website: <https://bookdown.org/yihui/bookdown/>.

Chapter 1

(10 points) ARIMA model

Consider `fma::sheep`, the sheep population of England and Wales from 1867–1939. `:sheep`:

```
#install.packages('fma')
library(fma)
sheep_ts <- as_tsibble(fma::sheep)
head(sheep_ts)
```

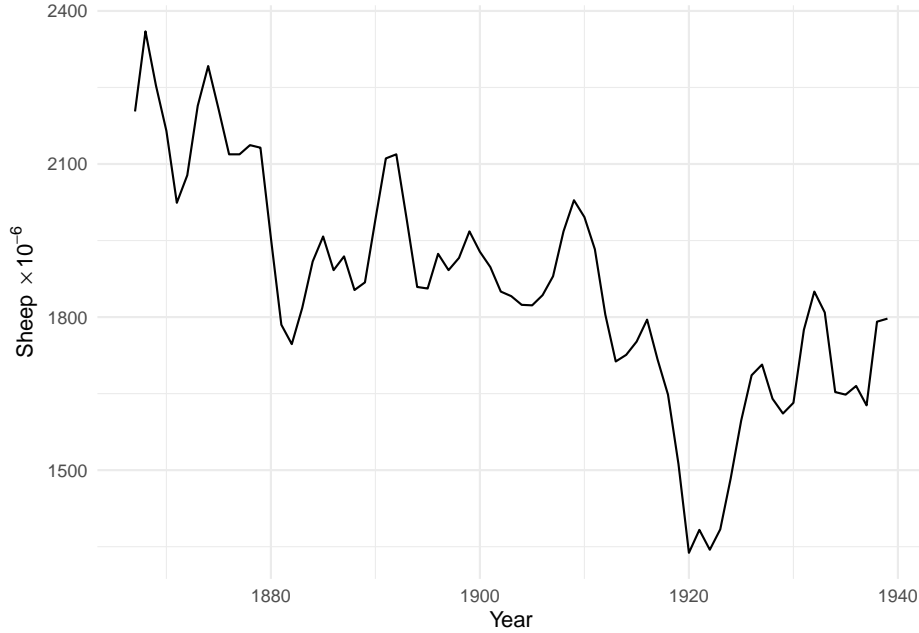
```
## # A tsibble: 6 x 2 [1Y]
##   index value
##   <dbl> <dbl>
## 1  1867  2203
## 2  1868  2360
## 3  1869  2254
## 4  1870  2165
## 5  1871  2024
## 6  1872  2078
```

1.1 Time series plot

Produce a time plot of the time series, and comment on what you observe.

```
sheep_plot <- autoplot(sheep_ts) + labs(x = "Year", y = TeX(r'(Sheep $\times 10^{-6}$'))))

## Plot variable not specified, automatically selected `.vars = value`
sheep_plot
```



Sheep population is gradually decreasing over this period and appears to stabilize in the 1930's at around 1700×10^6 . There is a significant local dip that coincides with the WWI and subsequent period of economic turmoil in Europe.

1.2 Fit a model

Assume you decide to fit the following model:

$$y_t = y_{t-1} + \phi_1(y_{t-1} - y_{t-2}) + \phi_2(y_{t-2} - y_{t-3}) + \phi_3(y_{t-3} - y_{t-4}) + \epsilon_t,$$

where ϵ_t is a white noise series.

1.2.1 Model type

What sort of ARIMA model is this (i.e., what are p , d , and q)?

This is a representation of ARIMA(3.1.0). There are no error terms from the past lags, therefore MA term $q = 0$; there is only one differing step, so $d = 1$, and the last lag taken into account is y_{t-4} which means $p = 3$ (given $d = 1$).

1.2.2 Back to the future

Express this ARIMA model using backshift operator notation.

$$(1 - B)y_t = \phi_1(1 - B)y_{t-1} + \phi_2(1 - B)y_{t-2} + \phi_3(1 - B)y_{t-3} + \epsilon_t$$

$$(1 - B)y_t = \phi_1(1 - B)By_t + \phi_2(1 - B)B^2y_t + \phi_3(1 - B)B^3y_t + \epsilon_t$$

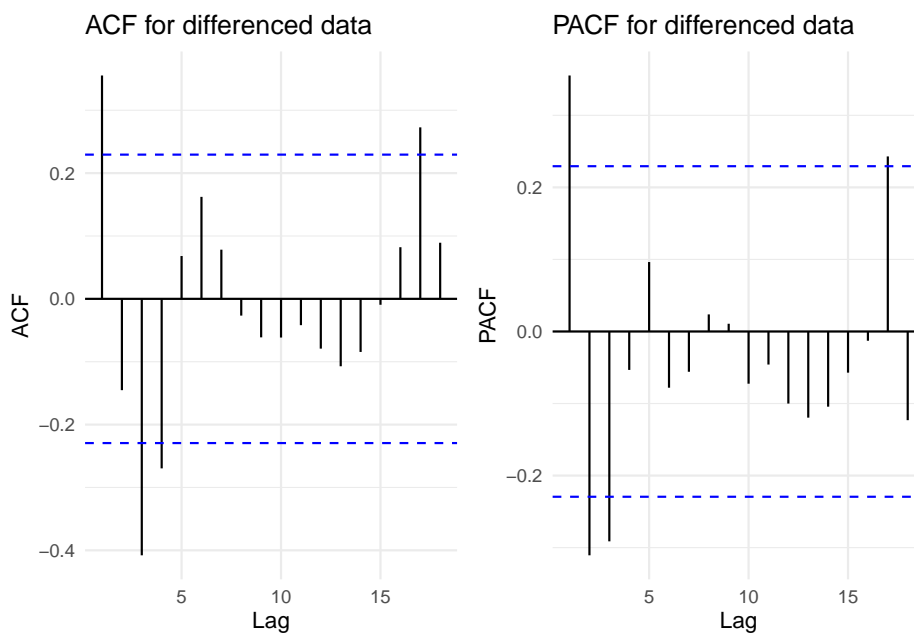
$$(1 - B)y_t = (1 - B)(\phi_1By_t + \phi_2B^2y_t + \phi_3B^3y_t) + \epsilon_t$$

$$(1 - B)(1 - \phi_1B - \phi_2B^2 - \phi_3B^3)y_t = \epsilon_t$$

1.3 Is this model appropriate?

Examine the ACF and PACF of the differenced data. Evaluate whether this model is appropriate.

```
# Add a column with differenced data to the original tsibble
sheep_ts <- sheep_ts %>% mutate(diff = difference(value, lag = 1))
# Generate acf plot
acf_plot <- acf(sheep_ts$diff, na.action = na.pass, plot = FALSE) %>%
  autoplot() + labs(title = "ACF for differenced data")
# Generate pacf plot
pacf_plot <- pacf(sheep_ts$diff, na.action = na.pass, plot = FALSE) %>%
  autoplot() + labs(title = "PACF for differenced data")
acf_plot | pacf_plot # Display plots
```



The model does not capture all the aspects of the underlying data. On one hand, PACF plot shows only 3 significant lags, indicating that this is indeed an AR(3) process. On the other hand, ACF plot shows an abrupt decrease of significance after lag 4 and some oscillating

pattern, indicating MA(4) process with some periodicity. Therefore ARIMA(3.1.4) would be more appropriate model.

1.4 Forecasts, by hand!

The last five values of the series are given below:

Year	1935	1936	1937	1938	1939
Millions of sheep	1648	1665	1627	1791	1797

The estimated parameters are:

- $\phi_1 = 0.42$;
- $\phi_2 = -0.20$; and,
- $\phi_3 = -0.30$.

Without using the forecast function, calculate forecasts for the next three years (1940–1942).

$$y_{1940} = y_{1939} + \phi_1(y_{1939} - y_{1938}) + \phi_2(y_{1938} - y_{1937}) + \phi_3(y_{1937} - y_{1936}) + 0$$

$$y_{1941} = y_{1940} + \phi_1(y_{1940} - y_{1939}) + \phi_2(y_{1939} - y_{1938}) + \phi_3(y_{1938} - y_{1937}) + 0$$

$$y_{1942} = y_{1941} + \phi_1(y_{1941} - y_{1940}) + \phi_2(y_{1940} - y_{1939}) + \phi_3(y_{1939} - y_{1938}) + 0$$

```
y1936<-1665
```

```
y1937<-1627
```

```
y1938<-1791
```

```
y1939<-1797
```

```
phi_1 <- 0.42
```

```
phi_2 <- -0.20
```

```
phi_3 <- -0.30
```

```
y1940 <- y1939 + phi_1*(y1939-y1938) + phi_2*(y1938-y1937) + phi_3*(y1937-y1936)
```

```
y1941 <- y1940 + phi_1*(y1940-y1939) + phi_2*(y1939-y1938) + phi_3*(y1938-y1937)
```

```
y1942 <- y1941 + phi_1*(y1941-y1940) + phi_2*(y1940-y1939) + phi_3*(y1939-y1938)
```

Forecasted values for 1940, 1941 and 1942 are 1778, 1721 and 1699, respectively.

1.5 Interpret roots

Find the roots of your model's characteristic equation. Is this process stationary?.

```
roots <- polyroot(c(1, -phi_1, -phi_2, -phi_3))
```

```
mroot1 <- Mod(roots[1])
```

```
mroot2 <- Mod(roots[2])  
mroot3 <- Mod(roots[3])
```

Characteristic polynomial for the ARMA model is $1 - \phi_1 B - \phi_2 B^2 - \phi_3 B^3 = 0$. For the process to be stationary, all of the roots should lie strictly outside of a unit circle on the complex plane, i.e. modules of all roots should be more than 1. We have modules 1.26, 2.09, 1.26, all above 1. That means that the process is stationary.

Chapter 2

(10 points) Seasonal ARIMA model

```
library(fredr)
if (fredr_has_key()){
  ecom_df <- fredr(series_id = "ECOMPCTNSA",
    observation_start = as.Date("1990-01-01"))
  ecom_df <- cbind.data.frame(ecom_df[c(1,3)], index = 1:nrow(ecom_df))
  ecom_ts <- as_tsibble(ecom_df, index = "date")
} else {
  print("Expect FREDR API key as an environment variable")
  quit(save="ask")
}

split <- as.Date("2020-12-31")
train_ts <- ecom_ts %>% filter(date < split)
test_ts <- ecom_ts %>% filter(date >= split)
```

Download the series of E-Commerce Retail Sales as a Percent of Total Sales [here](#).

(Feel free to explore the `fredr` package and API if interested.)

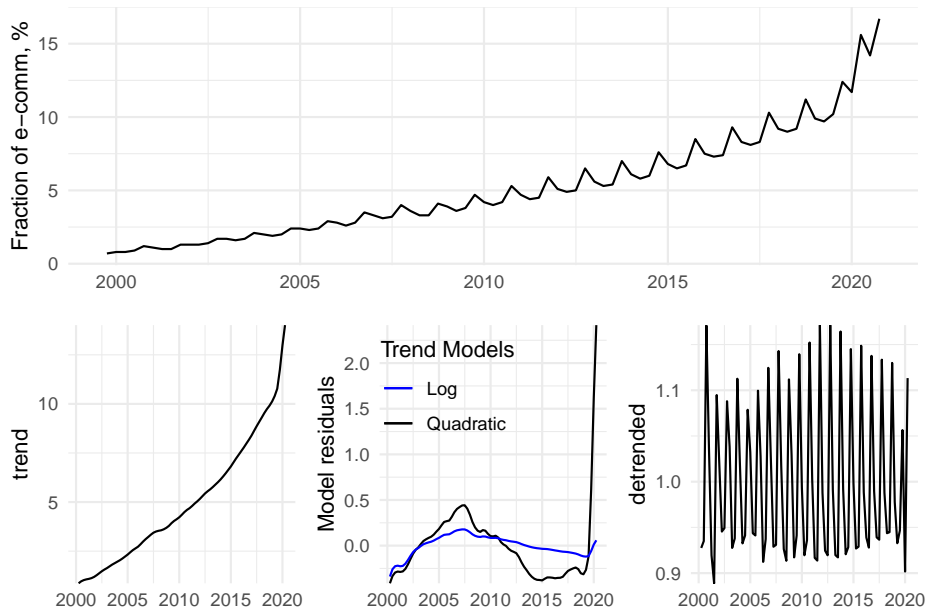
Our goal is to Build a Seasonal ARIMA model, following all appropriate steps for a univariate time series model.

Separate the data set into training and test data. The training data is used to estimate model parameters, and it is for 10/1999-12/2020. The test data is used to evaluate its accuracy, and it is for 01/2021-01/22.

2.1 Time series plot

Plot training data set of Retail Sales. What do you notice? Is there any transformation necessary?

```
## Plot variable not specified, automatically selected `.vars = value`
```



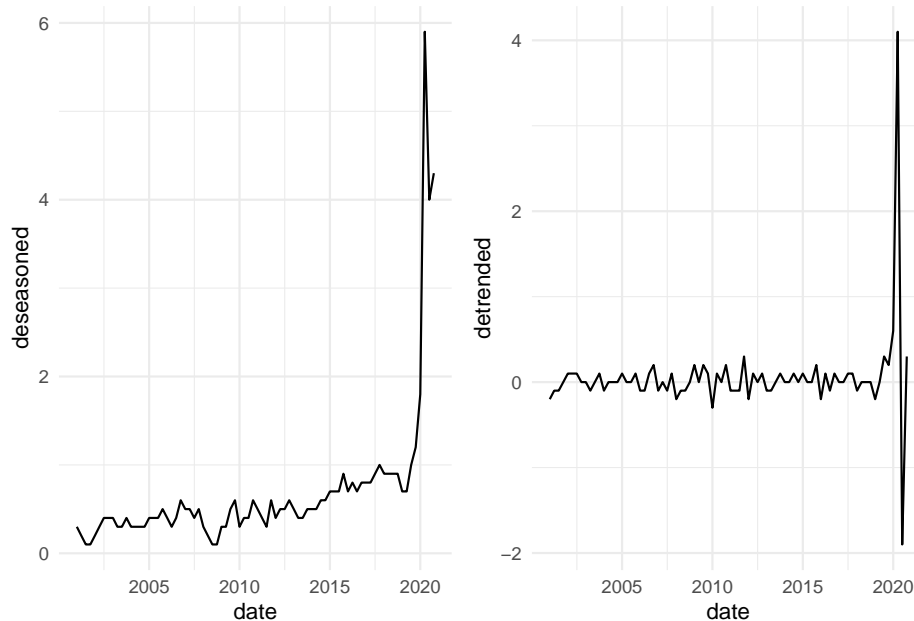
Fraction of e-commerce in the overall retail sales is growing, and its growth is accelerating. It appears that at around the start of the COVID-19 pandemic the acceleration increased. This is evident from abrupt change in the overall trend curvature at the beginning of 2020. The series also has marked seasonal component, and the magnitude of the seasonal component is increasing, suggesting a multiplicative seasoning. The trend has obvious curvature, but neither quadratic, nor logarithmic transformation result in white-noise residuals.

2.2 Check for Stationary

Use ACF/PACF and a unit root test to check if Retail Sales is stationary. If data is not stationary, difference the data, and apply the test again until it becomes stationary? How many differences are needed to make data stationary?

```
my_lag <- 4
train_ts <- mutate(train_ts,
  deseasoned = difference(value, lag = my_lag),
  detrended = difference(deseasoned, lag = 1))
train_ts <- na.omit(train_ts)
```

```
deseasoned_plot <- ggplot(data = train_ts) + aes(x = date, y = deseasoned) +
  geom_line()
detrended_plot <- ggplot(data = train_ts) + aes(x = date, y = detrended) +
  geom_line()
deseasoned_plot | detrended_plot
```



```
value_p <- PP.test(train_ts$value)
deseasoned_p <- PP.test(train_ts$deseasoned)
detrended_p <- PP.test(train_ts$detrended)
```

Given strong trend and seasonality in the data it is obvious that the original series is not stationary. Deseasoning via differencing with lag 4 (1 year), does not eliminate the trend. Further de-trending with lag = 1 eliminates the trend, but the resulting trend is clearly heteroschedastic. Results of Phillips-Perron test reflect these observations. The test fails to reject the hypothesis of non-stationarity for the original series (p-value = 0.983) and de-seasoned series (p-value = 0.535), but strongly rejects this hypothesis for the de-trended series with p-value = 0.01)

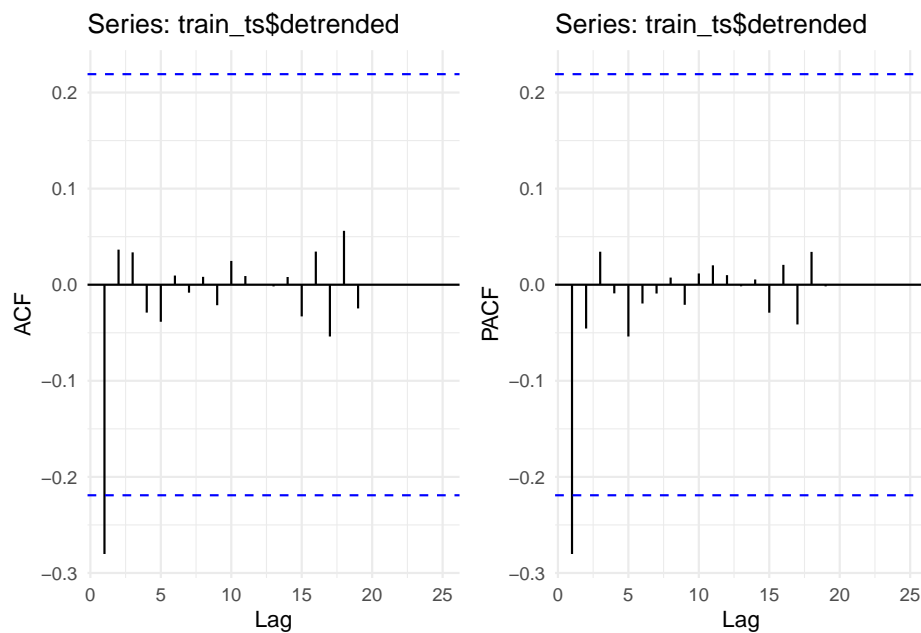
2.3 Model identification and estimation

Use ACF/PACF to identify an appropriate SARIMA model. Estimate both select model and model chosen by ARIMA()

```
acf_plot <- acf(train_ts$detrended, plot = F) %>% autoplot() + xlim(1,25)

## Scale for x is already present.
## Adding another scale for x, which will replace the existing scale.
pacf_plot <- pacf(train_ts$detrended, plot = F) %>% autoplot() + xlim(1,25)

## Scale for x is already present.
## Adding another scale for x, which will replace the existing scale.
acf_plot | pacf_plot
```



```
model.manual <- arima(train_ts$value, order = c(1, 1, 1),
  seasonal = list(order = c(0, 1, 0), period = 4))
model.manual

##
## Call:
## arima(x = train_ts$value, order = c(1, 1, 1), seasonal = list(order = c(0, 1,
## 0), period = 4))
##
## Coefficients:
##          ar1          ma1
##      -0.2077   -0.0673
## s.e.    0.3814    0.3856
##
## sigma^2 estimated as 0.2697:  log likelihood = -57.32,  aic = 120.64
```



```

model.auto <- train_ts$value %>% ts(frequency = 4) %>% auto.arima(d = 1, D = 1,
                             max.p = 5, max.q = 5, max.P = 2, max.Q = 2, max.d = 2, max.D = 2,
                             max.order = 10,
                             start.p = 0, start.q = 0, start.P = 0, start.Q = 0,
                             ic="aic", seasonal = TRUE, stepwise = FALSE, approximation = FALSE, trace = FALSE)
model.auto

## Series: .
## ARIMA(1,1,0)(0,1,0)[4]
##
## Coefficients:
##          ar1
##        -0.2702
## s.e.    0.1105
##
## sigma^2 = 0.2735: log likelihood = -57.34
## AIC=118.67   AICc=118.84   BIC=123.31

train_short <- train_ts %>% filter(date < as.Date("2020-01-01"))
model.short <- train_short$value %>% ts(frequency = 4) %>% auto.arima(d = 1, D = 1,
                             max.p = 5, max.q = 5, max.P = 2, max.Q = 2, max.d = 2, max.D = 2,
                             max.order = 10,
                             start.p = 0, start.q = 0, start.P = 0, start.Q = 0,
                             ic="aic", seasonal = TRUE, stepwise = FALSE, approximation = FALSE, trace = FALSE)
model.short

## Series: .
## ARIMA(0,1,3)(2,1,0)[4]
##
## Coefficients:
##          ma1          ma2          ma3          sar1          sar2
##        -0.2811   -0.0158   -0.4032    0.0990    0.3872
## s.e.    0.1236    0.1477    0.1312    0.1379    0.1206
##
## sigma^2 = 0.0114: log likelihood = 59.74
## AIC=-107.47   AICc=-106.16   BIC=-93.9

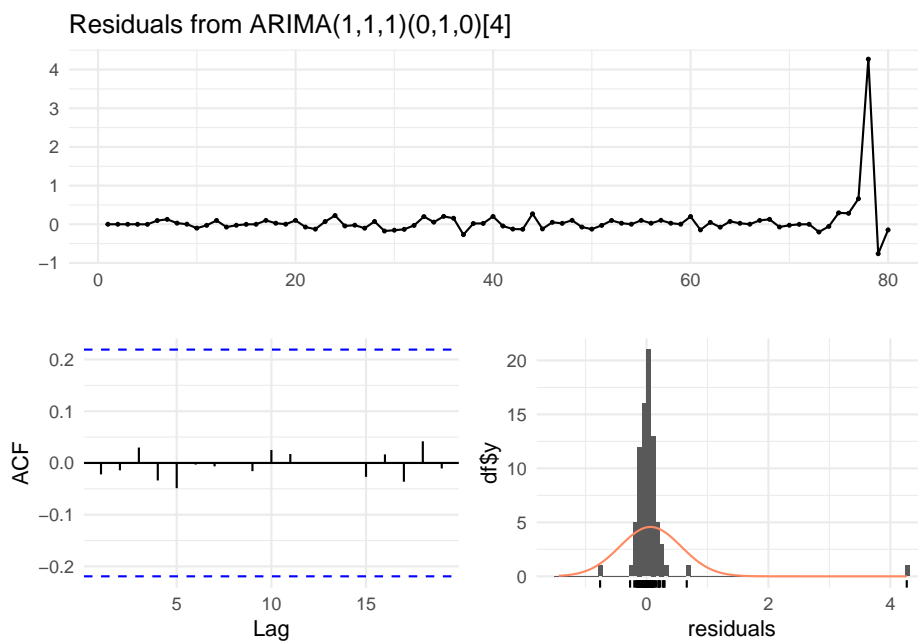
```

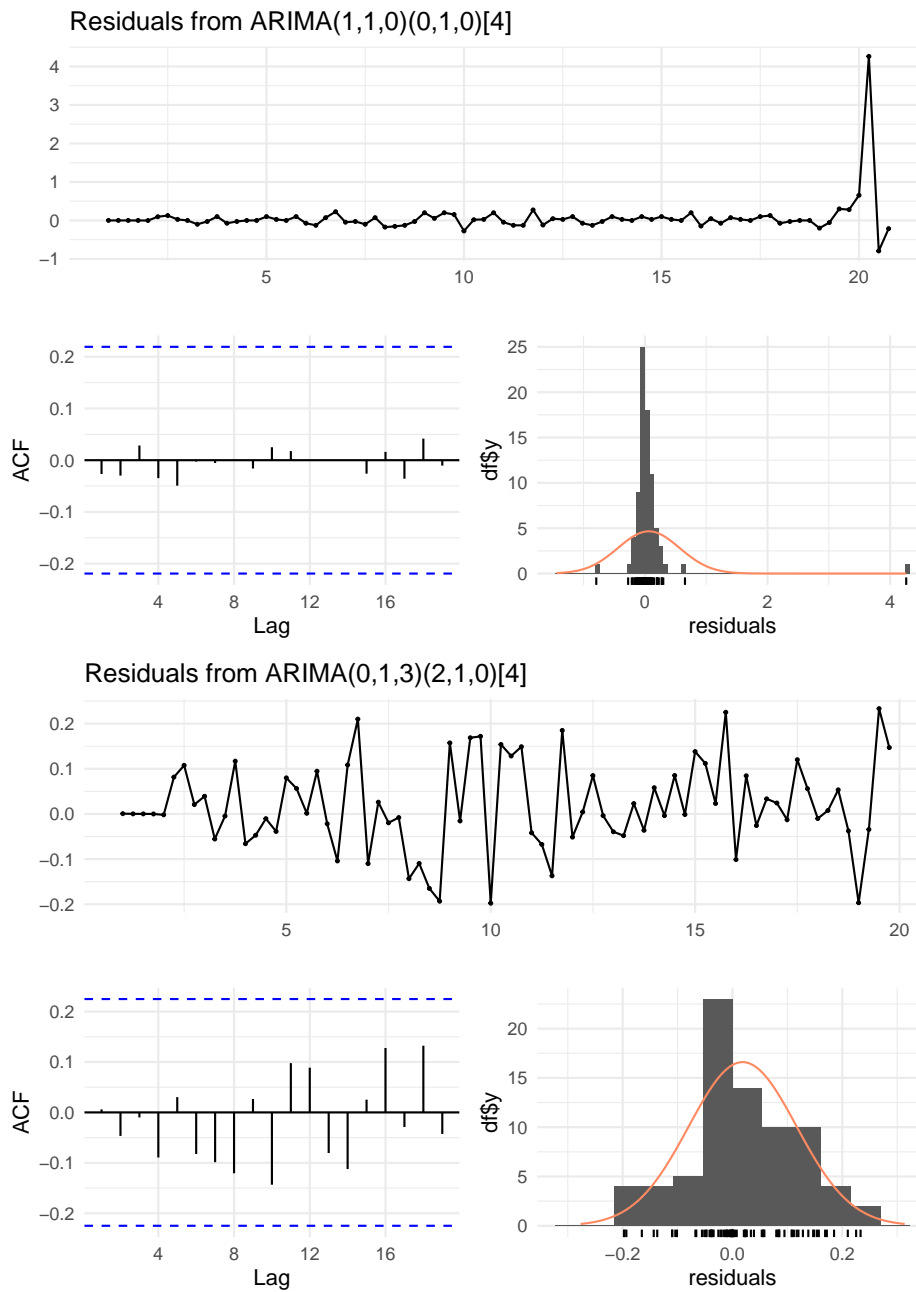
ACF plot for de-seasoned and de-trended series contains one strongly significant lag = 1, same as PACF plot. Formally, this is typical for ARMA(1,1) process, however, there is also a chance that this is an artifact caused by the few lags at the end of the series that are affected by the pandemic abnormality. Having abnormally large values at the tail of the series might make short lags a lot more significant. This considerations aside, ACF and PACF plot suggest SARIMA(1,1,1)(0,1,0)[4] model. Estimating this model results in AIC = 120.6. Grid search with `auto.arima` function yields very similar result. The function finds the lowest AIC 118.7 for SARIMA(1,1,0)(0,1,0)[4]. The fact that MA component of ARIMA does not improve the model, despite the fact that PACF

plot goes down abruptly after lag 1, supports the hypothesis that observed plots are just artifacts. Repeating this grid search on the shortened training set that excludes the pandemic data results in a much lower AIC of -107.5.

2.4 Model diagnostic

Do residual diagnostic checking of both models. Are the residuals white noise? Use the Ljung-box test to check if the residuals are white noise.





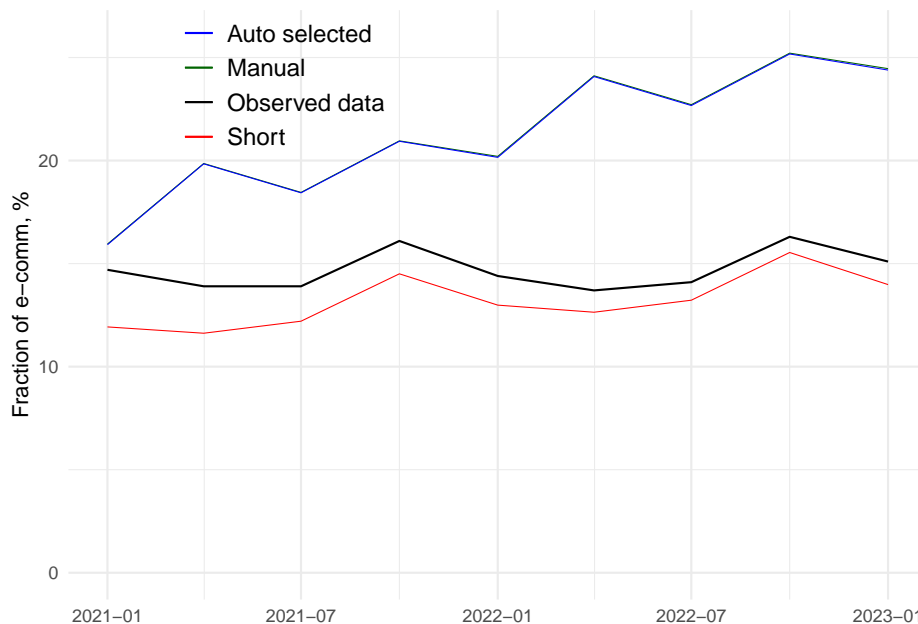
Both manually estimated model as well as `auto.arima` result fail visual tests for model quality: in both cases residual distribution has a significant outlier and residual plots contain strong spikes. However, Ljung-Box Test fails to reject the null hypothesis with p-values 0.84 and 0.81 respectively. The model estimated on the shorter data does have all the characteristics of good fit: normally

distributed residuals with no pattern in acf plot and p-value 0.96 on Ljung-Box Test, indicating no autocorrelation in residuals.

2.5 Forecasting

Use the both models to forecast the next 12 months and evaluate the forecast accuracy of these models.

```
frcst_length = length(test_ts$value)
frcst_manual <- forecast(model.manual, h=frcst_length)
frcst_auto <- forecast(model.auto, h=frcst_length)
frcst_short <- forecast(model.short, h=frcst_length+4)
```



Both manually selected and auto-selected models make essentially the same predictions, both significantly overestimate the growth rate for the e-commerce sales. This is because the models were trained on the data that contained a large shock and they implicitly expect more of the similar shocks to come. The model that was trained on truncated data, on the other hand, captures the underlying market forces without compounding effect of once-in-a-lifetime anomaly. As a result, once the anomaly passes, this model predicts the actual observed data much better than the previous two.

Chapter 3

(10 points) Time Series Linear Model and Cointegration

Daily electricity demand and temperature (in degrees Celsius) is recorded in `./data/temperature_demand.csv`. Please work through the following questions to build a time series linear model against this data.

```
library(tidyverse)
```

```
temperature <- read_csv('./data/temperature_demand.csv') %>%  
  rename(  
    'index'      = '...1',  
    'demand'     = 'Demand',  
    'work_day'   = 'WorkDay',  
    'temperature' = 'Temperature'  
  )  
#glimpse(temperature)
```

3.1 Plot electricity

Plot electricity demand and temperature as time series. Is there any correlation between these two variables? If yes, Do you think is it a spurious correlation?

```
#'fill this in'
```

'Fill this in

3.2 Cointegration test

Use the Engle-Granger test to check for cointegration. What do you conclude?

#'fill this in'

3.3 Fit Model

Based on cointegration test, fit a regression model for demand with temperature as an explanatory variable (or their first difference).

#'fill this in'

3.4 Residuals Plot

Produce a residual plot of the estimated model in previous part. Is the model adequate? Describe any outliers or influential observations, and discuss how the model could be improved.

#Fill this in

'Fill this in:

3.5 Forecasting model

Use a model to forecast the electricity demand (with **prediction** intervals) that you would expect for the next day if the maximum temperature was 15° . Compare this with the forecast if the with maximum temperature was 35° . Do you believe these forecasts? Why or why not?

#Fill this in

'Fill this in:'

Chapter 4

(12 points) Vector autoregression

```
library(tidyverse)
```

Annual values for real mortgage credit (RMC), real consumer credit (RCC) and real disposable personal income (RDPI) for the period 1946-2006 are recorded in `./data/mortgage_credit.csv`.

All of the observations are measured in billions of dollars, after adjustment by the Consumer Price Index (CPI).

Our goal is to develop a VAR model for these data for the period 1946-2003, and then forecast the last three years, 2004-2006.

```
credit <- read_csv('./data/mortgage_credit.csv')  
#glimpse(credit)
```

4.1 Time series plot

Plot the time-series of real mortgage credit (RMC), real consumer credit (RCC) and real disposable personal income (RDPI)? Do they look stationary?

```
#Fill this in
```

'Fill this in

4.2 Check for the unit root

Plot ACF/PACF and Perform the unit root test on these variables and report the results. Do you reject the null of unit root for them? Is the first differencing

necessary?

#Fill this in

'Fill this in

4.3 Determine VAR model

Based on the unit root results transform the variables and determine the lag length of the VAR using the information criteria.

#Fill this in

'Fill this in

4.4 Estimation

Estimate the selected VAR in previous part and comment on the results.

#Fill this in

'Fill this in

4.5 Model diagnostic

Do diagnostic checking of the VAR model.

#Fill this in

'Fill this in

4.6 Forecasting

forecast the last three years, 2004-2006.

#Fill this in

'Fill this in