

Computer exercises 5

Demo exercises

The goal is to build a regression model that explains alcohol consumption expenditures per capita **Q1CPC** with the real price index of alcohol **R1C** and total consumption expenditures per capita **QTOTALPC**. The data we use in these demo exercises is real data from Finland and the data consists of yearly time series between 1950 and 1981. The prices of different years have been inflation adjusted to be comparable with the fixed prices of the year 1975.

An elementary model in economics is the log-linear regression model, defined as,

$$\log(\mathbf{Q1CPC}) = \beta_0 + \beta_1 \log(\mathbf{R1C}) + \beta_2 \log(\mathbf{QTOTALPC}) + \varepsilon, \quad (1)$$

where $\log(\cdot)$ denotes natural logarithm. Model (1) will be estimated in Exercise 5.1 for years 1950-1981.

A problem with Model (1) is that the alcohol legislation changed in the beginning of the year 1969. This motivates us to add an indicator (or dummy) variable **LAW** to the Model (1). The variable **LAW** represents the change in the legislation and it is defined as follows,

$$\mathbf{LAW}_t = \begin{cases} 0, & 1950 \leq t \leq 1968 \\ 1, & 1969 \leq t \leq 1981, \end{cases}$$

where $t \in \{1950, \dots, 1981\}$ denotes the year. Hence, we obtain the model

$$\log(\mathbf{Q1CPC}) = \beta_0 + \beta_1 \log(\mathbf{R1C}) + \beta_2 \log(\mathbf{QTOTALPC}) + \beta_3 \mathbf{LAW} + \varepsilon. \quad (2)$$

The indicator variable **LAW** tries to take into account the jump in the level of the alcohol expenditures. Namely, in Model (2) the constant term is of the form,

$$\begin{aligned} &\beta_0, && 1950 \leq t \leq 1968 \\ &\beta_0 + \beta_3, && 1969 \leq t \leq 1981. \end{aligned}$$

The presumption is that the regression coefficient β_3 is statistically significant and positive. Model (2) will be estimated in Exercise 5.2 for the years 1950-1981.

However, by performing some regression diagnostics, we find that Model (2) is not satisfactory when trying to describe the behavior of the variable $\log(\mathbf{Q1CPC})$. The problem with model (2) is that the residuals of the estimated model are heavily correlated. We can sometimes get rid of the autocorrelation issue by utilizing so-called difference models. Therefore, we try the following difference model to describe alcohol expenditures,

$$\mathbf{D} \log(\mathbf{Q1CPC}) = \beta_0 + \beta_1 \mathbf{D} \log(\mathbf{R1C}) + \beta_2 \mathbf{D} \log(\mathbf{QTOTALPC}) + \beta_3 \mathbf{D} \mathbf{LAW} + \varepsilon. \quad (3)$$

The difference operation on the dummy variable **LAW** produces a so called impulse dummy. Model (3) is estimated in Exercise 5.3 for the years 1950-1981.

Note that, in Model (3) we have a different response variable than in Models (1) and (2). Thus, different models are not directly comparable.

In Exercise 5.4, we modify Model (2) by adding dynamic components. We try the following regression model for the alcohol expenditures,

$$\begin{aligned} \log(\mathbf{Q1CPC}_t) = & \beta_0 + \beta_1 \log(\mathbf{Q1CPC}_{t-1}) \\ & + \beta_2 \log(\mathbf{R1C}_t) + \beta_3 \log(\mathbf{R1C}_{t-1}) \\ & + \beta_4 \log(\mathbf{QTOTALPC}_t) + \beta_5 \log(\mathbf{QTOTALPC}_{t-1}) \\ & + \beta_6 \mathbf{LAW}_t + \beta_7 \mathbf{LAW}_{t-1} + \varepsilon_t, \end{aligned} \quad (4)$$

where X_{t-1} is the variable X_t with lag one. Note that the explanatory variable $\log(Q1CPC_{t-1})$ is not independent of the error term ε_t . Therefore, the standard assumptions are not satisfied and it is not possible to draw conclusions from the coefficient of determination directly.

In Exercises 5.3 and 5.4, we study the autocorrelation of the residuals by using Breusch–Godfrey test, which is similar to Ljung–Box test. In general, Ljung–Box can be applied to test autocorrelation of the residuals of fitted SARIMA models. However, it is not justified to use Ljung–Box test in regression diagnostics, if the model involves endogenous explanatory variables, that is, variables that are not independent of the residuals. On the other hand, Breusch–Godfrey test is applicable in these situations. In the Breusch–Godfrey test, the null hypothesis is that there is no autocorrelation up to the lag p . The test can be conducted in R with the function `bgtest`, which is implemented in the package `lmtest`.

First, we attach relevant packages, read the data and create `ts` objects. Notice that variables `LQ1CPC`, `LR1C` and `LQTOTALPC` are logarithms of variables `Q1CPC`, `R1C` and `QTOTALPC`, respectively.

```
library(car) # Calculate VIF
library(lmtest) # Breusch-Godfrey test

alko <- read.table("data/alkohol.txt", header = TRUE, sep = "\t")
alko <- alko[, 1:5]
head(alko)
```

```
##   YEAR   Q1CPC   LQ1CPC   LR1C LQTOTALPC
## 1 1950 148.0456 4.997520 4.614780 8.449157
## 2 1951 152.4226 5.026657 4.652774 8.510380
## 3 1952 168.9036 5.129328 4.605994 8.561820
## 4 1953 168.7201 5.128241 4.617447 8.546479
## 5 1954 169.6960 5.134008 4.630082 8.598005
## 6 1955 180.9960 5.198475 4.648004 8.659339
```

```
consumption <- ts(alko$LQ1CPC, start = 1950)
price <- ts(alko$LR1C, start = 1950)
total <- ts(alko$LQTOTALPC, start = 1950)
```

To ease our analysis, we make R functions for plotting diagnostics and performing Breusch–Godfrey test. Many of the defined functions are just wrappers of familiar R functions.

```
#' Plot original time series and fit
#'
#' @param y Response variable.
#' @param model Linear regression model object of class lm.
#' @param name Name of the response variable.
plot_fit <- function(y, model, name) {
  fit <- ts(model$fitted.values, start = start(y)[1])
  plot(y, col = "red", xlab = "Time", ylab = "")
  lines(fit, col = "blue")
  legend("topleft", legend = c(name, "Fit"), col = c("red", "blue"),
        lty = c(1, 1))
}

#' Plot Cook's distances
#'
#' @param y Response variable.
#' @param model Linear regression model object of class lm.
plot_cook <- function(y, model) {
  cooksds <- cooks.distance(model)
  plot(cooksds, xaxt = "n", type = "h", lwd = 3, xlab = NA,
```

```
      ylab = "Cook's distances")
axis(side = 1, at = seq(1, length(y), 5), cex.axis = 0.9,
      labels = seq(start(y)[1], (start(y)[1] + length(y) - 1), 5))
}

#' Plot residuals versus time
#'
#' @param y Response variable.
#' @param model Linear regression model object of class lm.
plot_res <- function(y, model) {
  res <- ts(model$residuals, start = start(y)[1])
  plot(res, type = "p", pch = 16)
}

#' ACF plot of residuals
#'
#' @param model Linear regression model object of class lm.
#' @param lag_max Maximum lag at which to calculate the acf.
plot_acf <- function(model, lag_max = NULL){
  acf(model$residuals, main = "", lag.max = lag_max)
}

#' Histogram of residuals (8 bins)
#'
#' @param model Linear regression model object of class lm.
plot_hist <- function(model) {
  res <- model$residuals
  breaks <- seq(min(res), max(res), length.out = 9)
  hist(model$residuals, xlab = "Residuals", ylab = "Frequency", main = "",
        breaks = breaks)
}

#' QQ plot
#'
#' @param model Linear regression model object of class lm.
plot_qq <- function(model) {
  res <- model$residuals
  qqnorm(res, pch = 16, main = "")
  qqline(res, col = "red", lwd = 2)
}

#' Perform Breusch-Godfrey test
#'
#' Breusch-Godfrey can be performed up to order
#' 'sample size' - 'number of estimated parameters'.
#'
#' @param model Linear regression model object of class lm.
#' @param m Number of estimated parameters.
#'
#' @return Vector of p-values.
res_test <- function(model, m) {
  n <- length(model$residuals)
  pvalues <- rep(NA, n - m)
```

```
for (i in 1:(n - m)) {  
  pvalues[i] <- bgtest(model, order = i)$p.value  
}  
pvalues  
}
```

5.1

Estimate Model (1) and study the goodness of fit.

Solution

Next, estimate Model (1).

```
model_log <- lm(consumption ~ price + total)  
summary(model_log)  
  
##  
## Call:  
## lm(formula = consumption ~ price + total)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -0.146202 -0.083305 -0.009638  0.082679  0.161945   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept) -2.89170     1.98506  -1.457   0.1559      
## price       -1.00346     0.37255  -2.693   0.0116 *      
## total        1.46489     0.05904   24.813 <2e-16 ***    
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.09986 on 29 degrees of freedom  
## Multiple R-squared:  0.9642, Adjusted R-squared:  0.9617   
## F-statistic: 390.1 on 2 and 29 DF,  p-value: < 2.2e-16
```

Comments about `summary(model_log)`:

- The regression coefficients corresponding to the variable `price` and variable `total` are statistically significant with 5% level of significance.
- The signs of the regression coefficients of the price and total expenditures variables are as expected: the coefficient of the price variable (`price`) is negative and the coefficient of the total expenditures variable (`total`) is positive.
- Interpretations of the regression coefficients as elasticities:
 - If the price goes up by 1 %, then the alcohol expenditures are reduced by 1.003%.
 - If the total expenditures are increased by 1 %, then the alcohol expenditures are increased by 1.465%.
- The coefficient of determination of the model is 96.42%.

Next, we study the normality of the residuals, Cook's distances and compare the fitted model with the original time series.

```
plot_fit(consumption, model_log, name = "Consumption")  
plot_cook(consumption, model_log)  
plot_res(consumption, model_log)  
plot_acf(model_log)  
plot_hist(model_log)  
plot_qq(model_log)
```

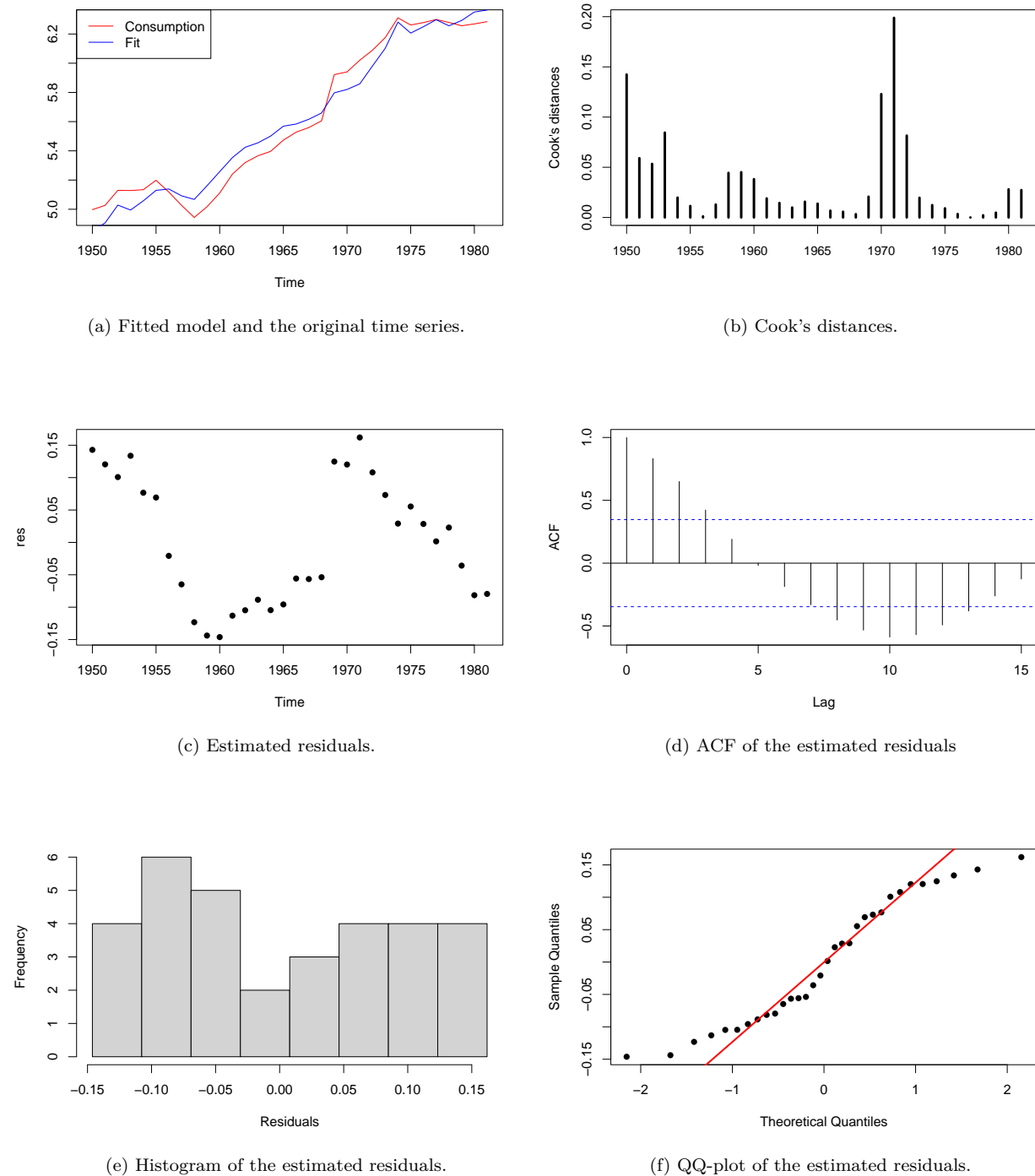


Figure 1: Diagnostic plots for Model (1).

```
vif(model_log)

## price total
## 1.16005 1.16005
```

Comments about diagnostics of Model (1):

- By Figures 1(e) and 1(f), the residuals do not look normally distributed.
- By Figures 1(c) and 1(d), the residuals seem to be heavily correlated.
- By the variance inflation factor (VIF), multicollinearity is not a problem here.
- The reason for the correlatedness of the residuals can be seen from Figure 1(a), where the fitted curve stays above and below the response variable **consumption** for long time periods.
- The model does not take account of the change in the legislation (the beginning of the year 1969). This is also visible in the Cook's distances (Figure 1(b)).

All in all, Model (1) cannot be considered to be sufficient.

5.2

Estimate and study Model (2).

Solution

First, let us fit Model (2).

```
law <- ts(c(rep(0, 19), rep(1, 13)), start = 1950)
model_law <- lm(consumption ~ price + total + law)

summary(model_law)

##
## Call:
## lm(formula = consumption ~ price + total + law)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.144576 -0.031179  0.008463  0.048923  0.086176
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   0.21798    1.40182   0.155  0.87754
## price        -0.88570    0.24650  -3.593  0.00124 **
## total         1.04355    0.07818  13.349 1.16e-13 ***
## law           0.31738    0.05106   6.216 1.03e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.06588 on 28 degrees of freedom
## Multiple R-squared:  0.9849, Adjusted R-squared:  0.9833
## F-statistic: 610.5 on 3 and 28 DF,  p-value: < 2.2e-16
```

Comments about `summary(model_law)`:

- The regression coefficients corresponding to the price variable (**price**) and the total expenditures variable (**total**) are statistically significant with 5% level of significance.
- The estimates for the regression coefficients differ from the estimates of Model (1).
- The signs of the regression coefficients for the price and total expenditures are as expected: the coefficient of the price variable is negative and the coefficient of the total expenditures variable is positive.

- Interpretations of the regression coefficients as elasticities:
 - If the price goes up by 1%, then the alcohol expenditures are reduced by 0.886%.
 - If the total expenditures are increased by 1%, then the alcohol expenditures are increased by 1.044%.
- The regression coefficient 0.317 corresponding to **law** is statistically significant with 5% level of significance.
- The sign of the regression coefficient of **law** is as expected.
- The coefficient of determination has increased to 98.49%.

Next let us study the goodness of the fit.

```
plot_fit(consumption, model_law, name = "Consumption")
plot_cook(consumption, model_law)
plot_res(consumption, model_law)
plot_acf(model_law)
plot_hist(model_law)
plot_qq(model_law)
```

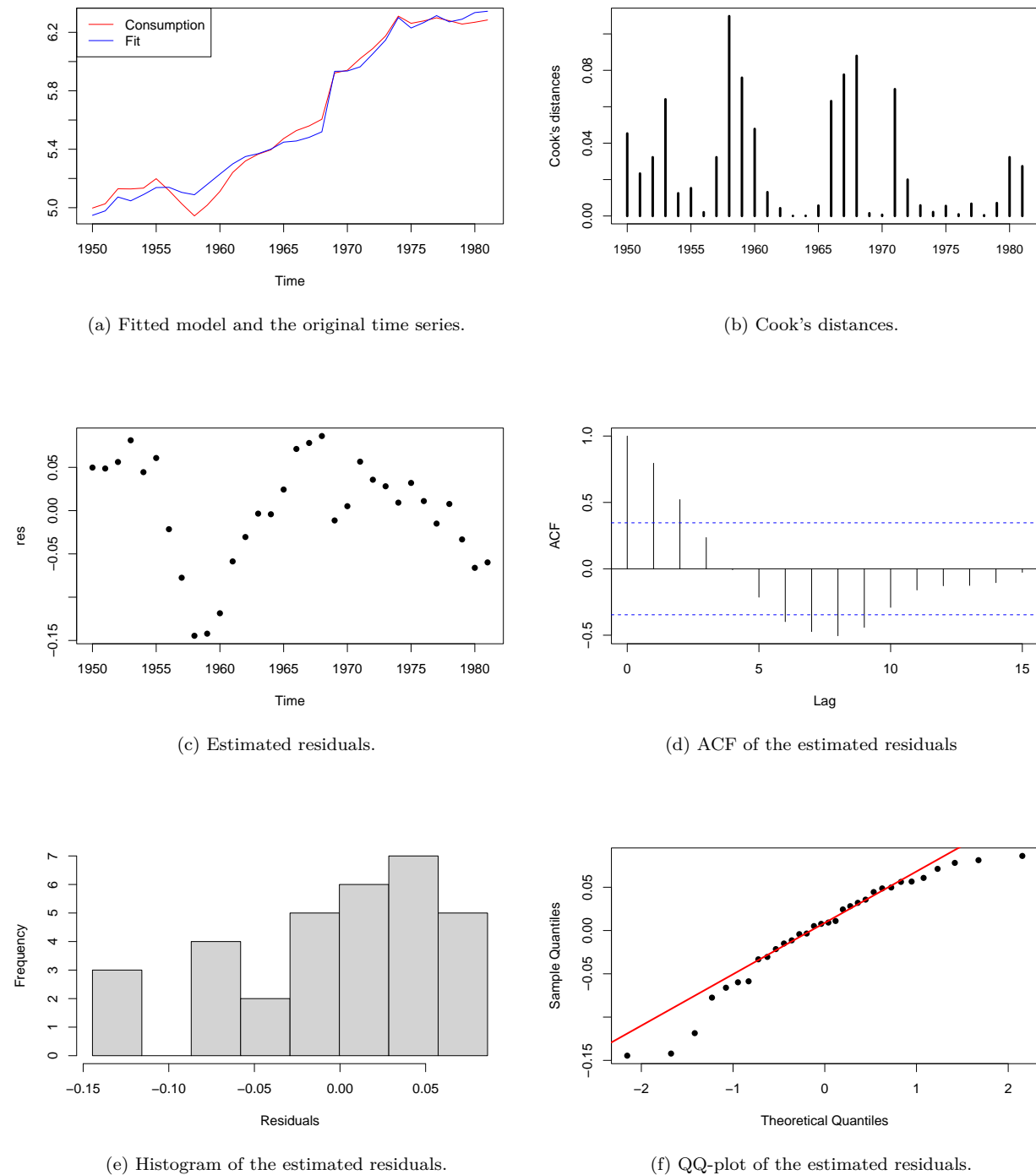



Figure 2: Diagnostic plots for Model (2).

```
vif(model_low)
```

```
## price total law
## 1.166943 4.674154 4.636612
```

Comments about diagnostics of Model (2):

- By Figures 2e) and 2f), the distribution of the residuals does not seem to be normal. The histogram of the residuals is skewed, which is evidence against normality.
- By Figures 2c) and 2d), the residuals are strongly correlated.
- By VIF, there is no problem with multicollinearity.
- The reason for the correlation of the residuals can be seen from Figure 2a). The fitted curve stays long time periods above and below the values of the response variable **consumption**.
- The model takes into account the change in the legislation (beginning of 1969).

By the regression diagnostics, this model is not satisfactory in explaining the alcohol expenditures.

5.3

Estimate and study Model (3).

Solution

First, we compute differenced variables and fit Model (3).

```
consumption_d <- diff(consumption)
price_d <- diff(price)
total_d <- diff(total)
law_d <- diff(law)

model_diff <- lm(consumption_d ~ price_d + total_d + law_d)
summary(model_diff)

##
## Call:
## lm(formula = consumption_d ~ price_d + total_d + law_d)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.08272 -0.01250  0.00000  0.02027  0.05541
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.010343   0.009273  -1.115   0.275
## price_d      -0.816697   0.133193  -6.132 1.50e-06 ***
## total_d       1.372390   0.225936   6.074 1.74e-06 ***
## law_d         0.196386   0.037765   5.200 1.78e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.03383 on 27 degrees of freedom
## Multiple R-squared:  0.8264, Adjusted R-squared:  0.8071
## F-statistic: 42.84 on 3 and 27 DF,  p-value: 2.113e-10
```

Comments about `summary(model_diff)`:

- All the coefficients of Model (3) are statistically significant (with the exception of the constant term) with 5% level of significance.
- The estimates for Model (3) are clearly different when compared to the estimates of Model (2).

- The signs of the regression coefficients of the price variable (**price_d**) and total expenditures variable (**total_d**) are as expected: the coefficient of the price variable is negative and the coefficient of the total expenditures variable is positive.
- Interpretations of the regression coefficients as elasticities:
 - If the price goes up by 1%, then the alcohol expenditures are reduced by 0.817%.
 - If the total expenditures are increased by 1%, then the alcohol expenditures are increased by 1.372%.
- The coefficient of the instant effect of the dummy variable **law_d** is 0.196.
- The coefficient of determination is 82.64%.
- The coefficient of determination of the Model (3) is not comparable with the coefficients of determinations corresponding to Models (1) and (2), since the response variable is not the same.

Next, let us study goodness of the fit. Here we also use Breusch–Godfrey test to study autocorrelations of the residuals.

```
plot_fit(consumption_d, model_diff, name = "D(Consumption)")
plot_cook(consumption_d, model_diff)
plot_res(consumption_d, model_diff)
plot_acf(model_diff)
plot_hist(model_diff)
plot_qq(model_diff)
```

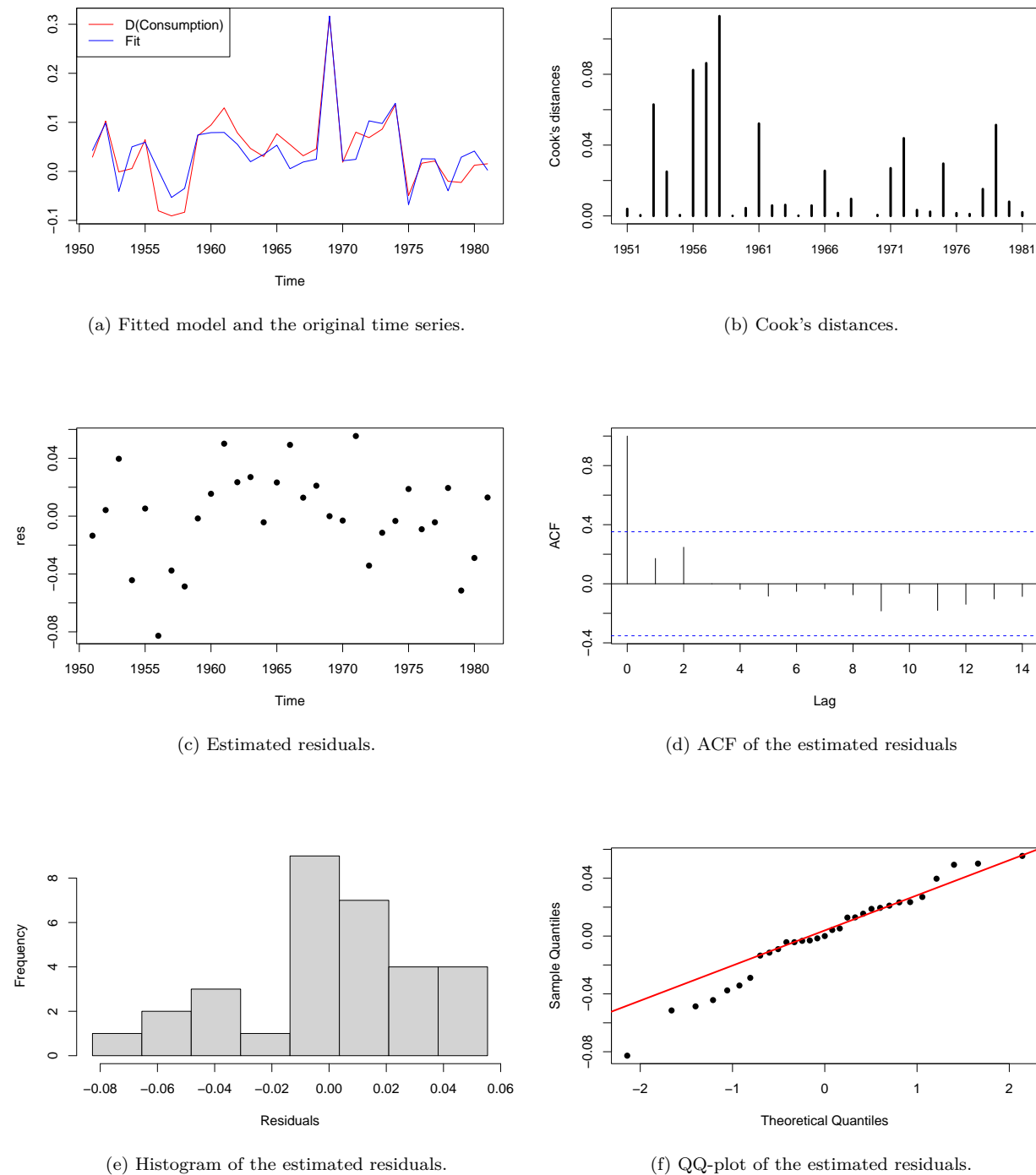


Figure 3: Diagnostic plots for Model (3).

```
vif(model_diff)
```

```
## price_d total_d law_d
## 1.061327 1.273221 1.206251
```

```
res_test(model_diff, m = 4) > 0.05
```

```
## [1] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [16] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
```

Comments about diagnostics of Model (3):

- By Figures 3e) and 3f), the residuals could be normally distributed.
- By Figures 3c) and 3d), the residuals are not correlated.
- By the Breusch-Godfrey test, the residuals are not correlated, since the null hypothesis is accepted with 5% level of significance for all lags.
- By VIF, multicollinearity is not a problem.
- By plotting residuals against time (Figure 3c)), there does not seem to be evidence of heteroscedasticity.

All in all, by diagnostics Model (3) seems satisfactory.

5.4

Estimate and study Model (4).

Solution

First, we fit Model (4). Note that, when variables of the form X_{t-1} and X_t are considered, the last and the first observation are omitted, respectively, when Model (4) is estimated.

```
n <- length(consumption)
model_lag <- lm(consumption[-1] ~ consumption[-n] + price[-1] + price[-n] +
                total[-1] + total[-n] + law[-1] + law[-n])
summary(model_lag)
```

```
##
## Call:
## lm(formula = consumption[-1] ~ consumption[-n] + price[-1] +
##     price[-n] + total[-1] + total[-n] + law[-1] + law[-n])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.073939 -0.013803  0.002322  0.013555  0.071924
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -0.36137    0.86190  -0.419  0.678911
## consumption[-n]  0.91197    0.10461   8.718 9.52e-09 ***
## price[-1]      -0.82927    0.16234  -5.108 3.57e-05 ***
## price[-n]       0.71352    0.17579   4.059 0.000486 ***
## total[-1]       1.46946    0.24887   5.905 5.10e-06 ***
## total[-n]      -1.31522    0.27235  -4.829 7.13e-05 ***
## law[-1]         0.17751    0.04468   3.973 0.000602 ***
## law[-n]        -0.19188    0.04707  -4.077 0.000465 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.03442 on 23 degrees of freedom
## Multiple R-squared:  0.9964, Adjusted R-squared:  0.9954
```

F-statistic: 922.3 on 7 and 23 DF, p-value: < 2.2e-16

Comments about `summary(model_lag)`:

- It is not possible to draw direct conclusions regarding the significance of the regression coefficients based on the t -tests. However, the results give some general direction, and the results indicate that all regression coefficients would be statistically significant (with the exception of the constant).
- The coefficient of the variable `consumption` with lag 1 is 0.912, which implies that the adjustment to changes in prices and total expenditures is rather fast.
- The signs of the coefficients of variables `price` and `total` with lag 0 are as expected: the coefficient -0.829 of the price variable is negative and the coefficient 1.469 of the total expenditures variable is positive. These coefficients describe the instant effects of changes in prices and total expenditures.
- The signs of the coefficients of the price and total expenditures variables with lag 1 are also as expected.
- Long term elasticities are:

$$\begin{array}{ll} \text{Price:} & \frac{\beta_2 + \beta_3}{1 - \beta_1} \approx -1.31, \\ \text{Total expenditures:} & \frac{\beta_4 + \beta_5}{1 - \beta_1} \approx 1.75. \end{array}$$

- Interpretations of the regression coefficients of price and total expenditures variables with lag 0:
 - If the price goes up by 1%, then the alcohol expenditures are instantly reduced by (without a lag) 0.829%.
 - If the total expenditures are increased by 1%, then the alcohol expenditures are increased by 1.469%.

Interpretations of the long term elasticities of price and total expenditures variables:

- If the price goes up by 1%, then the alcohol expenditures are reduced by 1.31% in the long term.
- If the total expenditures are increased by 1%, then the alcohol expenditures are increased by 1.75% in the long term.
- The coefficient of the instant effect of the dummy variable `law` is 0.177 and the long term coefficient is small ($(\beta_6 + \beta_7)/(1 - \beta_1) \approx -0.16$). Hence, the change in the legislation has rather minor effect on the behaviour of the consumers in a long term, which seems plausible.
- Additionally, it is not possible to draw conclusions from the coefficient of the determination.

Next, let us study goodness of the fit. Here we also use Breusch–Godfrey test to study autocorrelations of the residuals.

```
y <- ts(consumption[-1], start = 1951)
plot_fit(y, model_lag, name = "Consumption")
plot_cook(y, model_lag)
plot_res(y, model_lag)
plot_acf(model_lag)
plot_hist(model_lag)
plot_qq(model_lag)
```

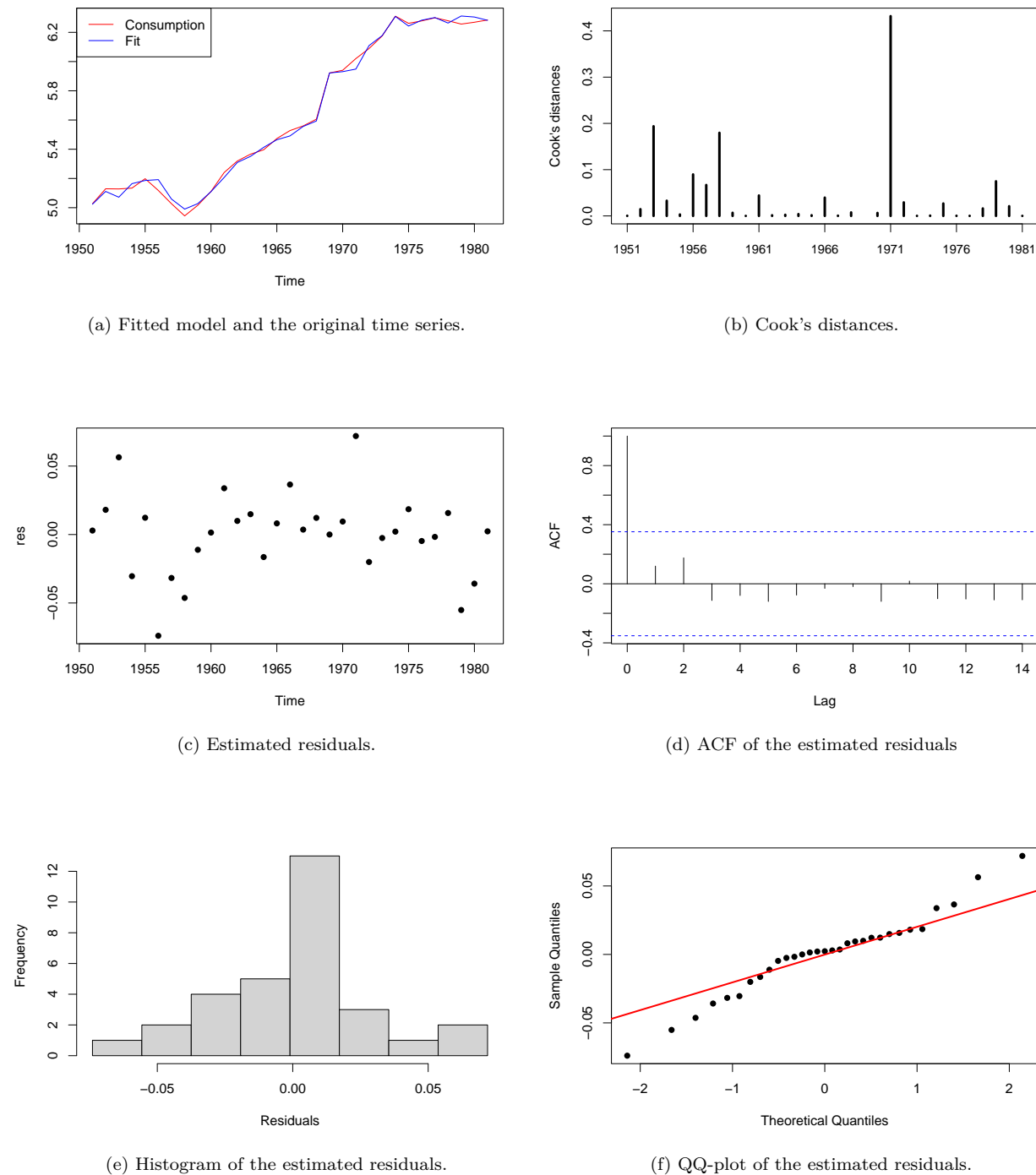


Figure 4: Diagnostic plots for Model (3).

```
vif(model_lag)
```

```
## consumption[-n]    price[-1]    price[-n]    total[-1]    total[-n]
##          70.276140      1.846537      2.128280    158.246107    192.764336
```

```
##          law[-1]          law[-n]
##          12.720381         13.752088
```

```
res_test(model_lag, m = 8) > 0.05
```

```
## [1] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [16] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
```

Comments about diagnostics of Model (3):

- By Figures 4e) and 4f), the residuals could be normally distributed. However, tails of the distribution seem to be heavier than tails of the normal distribution.
- By Figures 4c) and 4d), the residuals are not correlated.
- By the Breusch-Godfrey test, the residuals are not correlated. The null hypothesis is accepted with 5% level of significance for all lags.
- By VIF, there is strong multicollinearity in the model. This is unsurprising, as the model involves same variables with different lags.
- By Figure 4c), there is no evidence of heteroscedasticity.
- The model takes into account the change in the legislation.
- By Figure 4a), the fitted model coincides better with the original time series than the fits of Models (1) and (2).

We consider this model to be sufficient in explaining the alcohol expenditures.

Homework

5.5

The file `t38.txt` contains three quarterly time series. The time series start from the first quarter of the year 1953 and the corresponding time series are,

CONS = total consumption (billions)
INC = income (billions)
INFLAT = inflation (%)

The time series **CONS** and **INC** represent the observed total consumption and income in an imaginary country. The time series **INFLAT** represents inflation. The goal is to estimate a so-called consumption function that explains the time series **CONS** with the time series **INC** and **INFLAT**.

- a) Estimate the static linear regression model

$$\text{CONS}_t = \beta_0 + \beta_1 \text{INC}_t + \beta_2 \text{INFLAT}_t + \varepsilon_t.$$

and study the goodness of fit.

- b) Estimate the difference model

$$D(\text{CONS}_t) = \beta_0 + \beta_1 D(\text{INC}_t) + \beta_2 D(\text{INFLAT}_t) + \varepsilon_t.$$

and study the goodness of fit.

- c) Estimate the dynamic linear regression model

$$\text{CONS}_t = \beta_0 + \beta_1 \text{CONS}_{t-1} + \beta_2 \text{INC}_t + \beta_3 \text{INC}_{t-1} + \beta_4 \text{INFLAT}_t + \beta_5 \text{INFLAT}_{t-1} + \varepsilon_t$$

and study the goodness of fit.