

1 **Advanced Data Analysis with R - Part Time**
2 **Series Analysis**

3 **Summer Term 2025**

4 Johannes, Sebastian, and Kai (held by Kai)

Table of contents

6	Preface	3
7	Welcome	3
8	Who I am	3
9	The aims of the time series part are	3
10	Where are we in the course	4
11	Motivation for Time Series Analysis	5
12	Properties of Time Series Data - What is a Time Series?	9
13	Exemplary Time Series and Components of Time Series	13
14	Components of Time Series in More Detail	15
15	Stationarity	19
16	Descriptive Statistics and statistical modeling	28
17	Classical Decomposition	28
18	Autoregressive Models (AR)	30
19	Regression with time dynamics - Temporal regression	33
20	References	34

21 Preface

22 Welcome

23 to the time series part of the course advanced data analysis with R. In the **three** time series
24 lessons, we will

- 25 • **understand** why time series are an exciting type of data for us and where we usually
26 come in touch with them,
- 27 • get familiar with the **properties** of time series data and with their most relevant differ-
28 ences to other types of data that we already know,
- 29 • learn how to analyse time series data **descriptively** and with simple **time series re-**
30 **gression models**,
- 31 • and we will learn how to account for/ **correct for** time-dynamic covariates in regression
32 models.

33 All you need is this document and the respective data. You find both on GitHub. However,
34 this document will probably develop within the next few weeks. I let you know, once it is
35 finalised and stable.

36 Who I am

- 37 • [Kai Husmann](#)
- 38 • [Department Forest Economics and Sustainable Land-use Planning](#) (Prof. Carola Paul)
- 39 • [Projects and topics of Forest Econometrics](#)
- 40 • Contact: kai.husmann@uni-goettingen.de and StudIP

41 The aims of the time series part are

- 42 • to **make you aware** of time series as a specific kind of data,
- 43 • to understand ways and methods for **detection** time dynamic pattern in data,
- 44 • to introduce you to the exciting world of **time series regression** (even if we only scratch
45 the surface of the simplest models),
- 46 • and to account for (correct for) time-dynamic covariates in **regression models**.

Advanced Data Analysis with R: Outline

1. Statistical Modelling of spatio-temporal Data

Working with data in R & Research Data Management

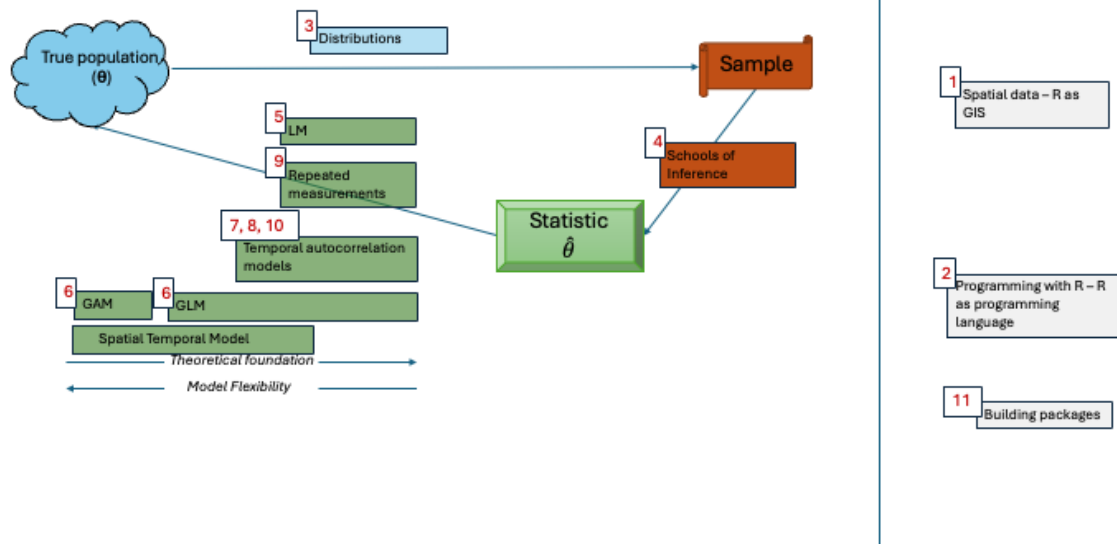


Figure 1: Overview

Motivation for Time Series Analysis

Time series data are ubiquitous in many fields, including economics and finance (where most of today's methods originate from), biology, and environmental sciences. As a result of the discussion about the resilience and resilience of ecosystems, time series data and time series methods have also become increasingly important in the context of ecology, agriculture, and also **forestry**. Time series methods are particularly promising in this area, as the time pattern (direct response, delayed response, ...) and the time horizon (how long is the recovery period, is there a recovery, ...) of the responses of ecosystem variables to disturbances are usually the main interest. Furthermore, many ecosystem variables themselves show a time trend. As time series models have evolved from the field of economics, they are also in a forestry context often used to describe the dynamics of economic variables, such as marked reactions in the sense of e.g. *how does the (timber) price react on supply and demand changes?* and does this relationship persist sudden and extreme supply changes (e.g. due to storms) (e.g. Fuchs et al. 2022)? Is it resistant and resilient to calamities? Time series models are more and more used to describe the dynamics of ecological variables as well, such as the relationship between tree growth and climate variables, or the relationship between tree mortality and tree health variables (e.g. Lemoine 2021).



Time series exercise 1

Consider Figure Figure 2.

1. Do you think, there is a relation between harvested volume and revenue or between share of damaged wood and revenue?
2. How would you analyse these relationships? Suggest a statistical model that you are already familiar with.

Following Lütkepohl and Krätzig (2004, 1), a time series is a sequence of observations of one variable over a period in. The observations are thus ordered in time and usually have equal observation frequency. Most economic measures, like the gross national product, wood prices, or wood material flows, are often provided at an annual base. In contrast, meteorological data, like temperature or precipitation, are often provided at a daily base, or even more frequent. Ecosystem data, like tree dimension's measurements or tree health data, are seldom found in a frequency higher than annually. The forest health survey in Germany (Waldzustandserhebung) e.g. takes place every year, while the national forest inventory (Bundeswaldinventur) is conducted every 10 years only.

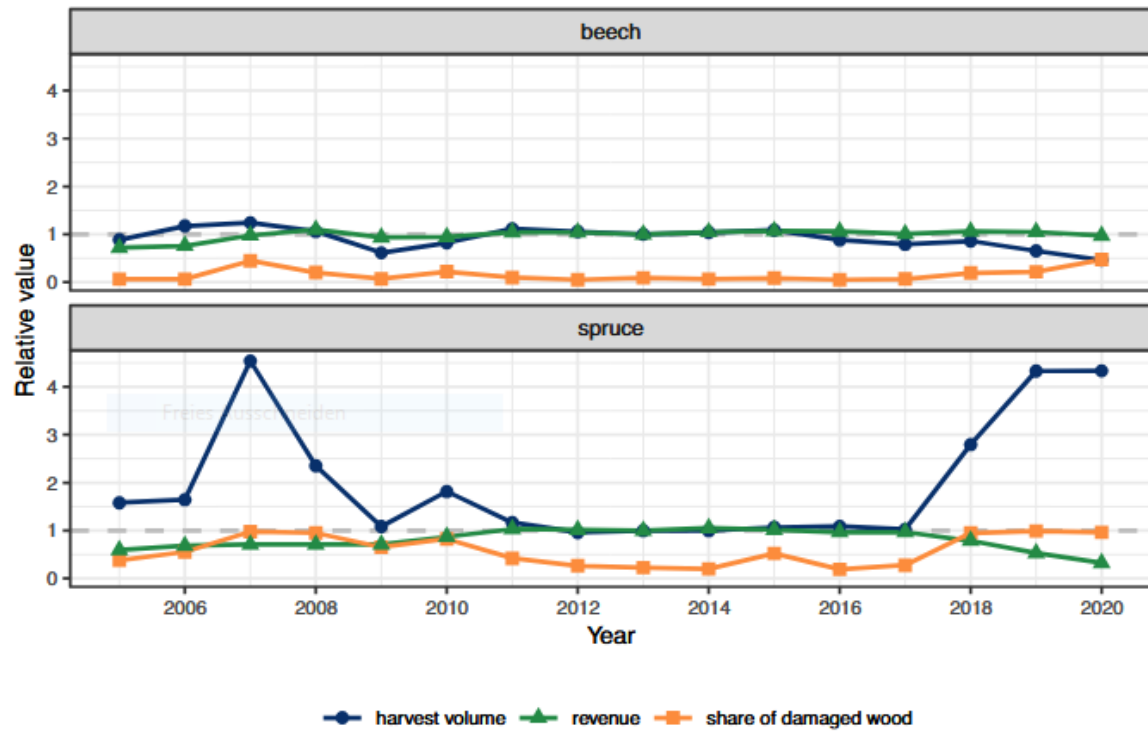


Figure 2: Fuchs et al. (2022): *Is there a relation connection between harvest volume, wood revenue and share of damaged wood? What is your guess?*

As with other types of data, time series data can be used in a regression context to describe correlations of the past, to forecast the future, or to estimate parameters for further use in e.g. causal simulation models. In addition, time series are used to integrate or correct for temporal dynamics (autocorrelation) in regression models, particularly in ecosystem sciences. In dynamic ecosystems, the relationships between the variables of interest are often confounded by temporal dynamics. If we want to infer the relationship between crown defoliation and precipitation, for example, we need to consider the state of crown defoliation in the past (diseased trees with high defoliation in the last year will never be 100% healthy in the current year, even if precipitation is currently sufficient). This may also explain why the data availability of economic variables is better than that of ecological variables. The challenge to measure variables in an even frequency without changing the measuring or estimation principles is a higher challenge in environmental sciences than it is in economics.

Typical time series projects start with a descriptive analysis of the time-dynamic patterns (Lütkepohl and Krätzig 2004, 5), whereby, in contrast to the previous descriptive analysis, an important aspect is whether the data are actually time series data.

The typical issues of interest in time series analysis are to do

- descriptive statistics of time-dynamic patterns,
- filtering (however, we won't do this in this course),
- hypotheses testing/ statistical inference,
- forecasting, and
- accounting for time-dynamics of covariates (autocorrelation) in regression models

In the three time series lessons, we will introduce some of the most common methods of univariate time series analysis and provide practised examples in R.

Topic 1: What is a time series?

- Examples of time series
- Relevant **properties** and **assumptions**
- **Differences** (and similarities) to other data types
- Concept of **stationarity**
- Detection of autocorrelation
- Practised programming features for time series in R

Topic 2: Analysis of time-dynamic patterns - Detection of autocorrelation - **Descriptive** statistics (classical decomposition) - Statistical **modeling** (exponential smoothing) - Hypotheses **testing** and causality

Topic 3: Accounting for autocorrelation in linear mixed models

- **Detection** of autocorrelated residuals in ordinary models
- Most common **procedures**

i Note

With this in mind. What are your wishes and expectations on the course. Let me know by next week.



<https://flinga.fi/s/FQ3KSVC>

Properties of Time Series Data - What is a Time Series?

We start with an example. The formal properties of time series are illustrated using the example of the wood price of oak (*Quercus robur* and *Quercus petraea*) in Germany. Data taken from <https://www-genesis.destatis.de/> (Code: 61231-0001). You find it as `stemwood_prices_annually.csv` in the data folder. We also use weather data of the weather station Göttingen (`month_mean_temp_goe.csv`, https://opendata.dwd.de/climate_environment/CDC/observations_germany/climate/monthly/kl/historical/).

This chapter introduces the specific properties of time series data. We will compare the time series as a random variable with other types that you are already familiar with. We introduce the most famous descriptive methods for time series data, thereby introducing the concept of autocorrelation. By doing so, we will also introduce and discuss some practical tools for data handling and visualization of time series data in R. We will already discuss, which attributes models need to bring to analyse univariate time series data and in which situations the time series properties have to be considered in regression models. Remembering Figure 3, we have three stages in the process of statistical inference. The population, which we usually want to make estimations for, the sample, which we actually have, and the estimated population, which allows us to create simulation based confidence intervals and which we used to illustrate the concept of unbiasedness (see also Chapter 4 Random Variables).

So far, a central assumption was that all observations from the population come from an arbitrary but common distribution and can be independent sampled. This assumption is not valid for time series data, as the observations are ordered in time and thus not independent from each other. In the case of normal distribution, e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$. This independence enabled unbiased estimates of the unconditioned average, e.g. the arithmetic mean, the unconditioned deviation, e.g. the standard deviation (Chapter 5 Random Variables), and to regress the data with other random variables (Chapter 5 Statistical Inference and 5 The Linear Model). In time series data, however, the observations are not independent, and the assumption of independence is violated. This prohibits to calculate averages and deviation that do not consider this dependency. Regressions with other covariates would be confounded by that dependency. All these methods would lead to biased estimates. However, we will not provide the proof of biases estimates in this course. Nor will we cover the simulation of time series. However, it is possible to simulate a time series by a *Brownian Motion* (e.g. Hamilton 2020, chap. 17.1 and 17.2) if you are interested.

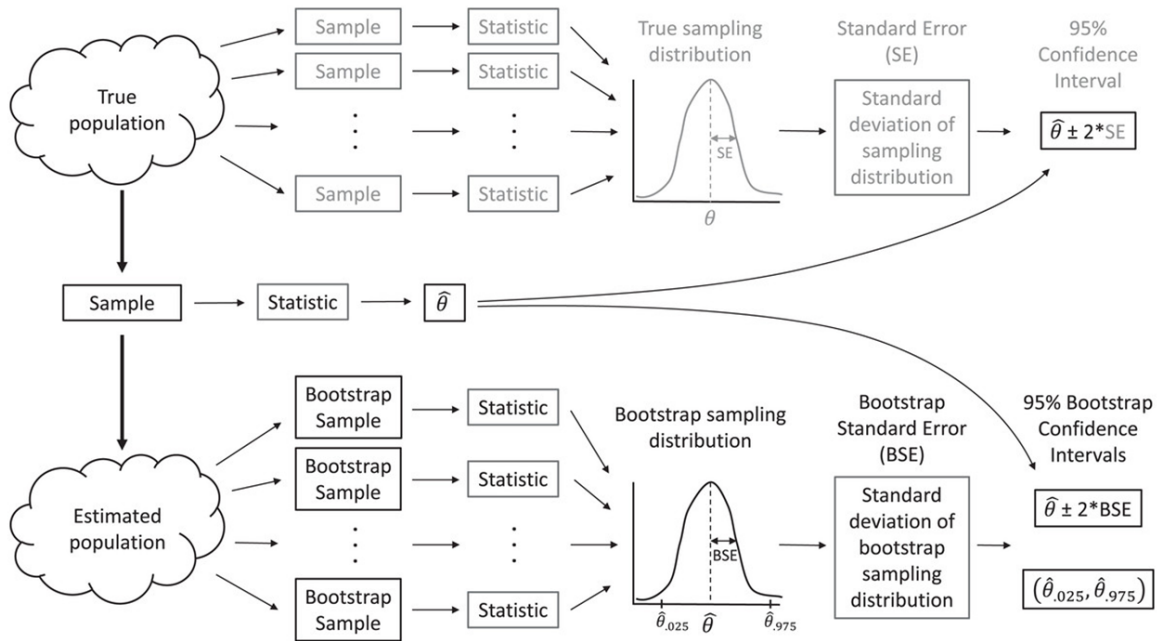


Figure 3: Fieberg, Vitense, and Johnson (2020) Resampling-based methods for biologists. See Chapter Resampling-based methods.

In practice, the arithmetic mean will estimate the center point of a time series (remember the Central Limit Theorem) but ignore the dependent part of the data, i.e. the autocorrelation. The same applies for the standard deviation, which will be constant over time. Consider e.g. the stem timber price index of oak (Figure 4). The arithmetic mean is usually not a suitable descriptive statistic for time series data. Instead of *what is the average of the data*, questions like *is there a seasonal trend?* or *To what extent does the data from the past describe my current situation in terms of time horizon and relevance?*

💡 Ask yourself

- Do I expect autocorrelation in the data?
- Does it possibly confound my statistic of interest?
- Am I interested in analysing the autocorrelation?

Another typical question could be *is there a linear trend?*, which brings us back to the ordinary linear regression (Chapter 6). The linear regression is suitable to describe the global linear trend of a time series (Figure 5). It will thus detect a long-term development of a series. However, the autocorrelation is ignored. Thus, typical question like *how is my recent observation related to last observations of my series?* or *Is there a seasonal pattern?* cannot be analysed

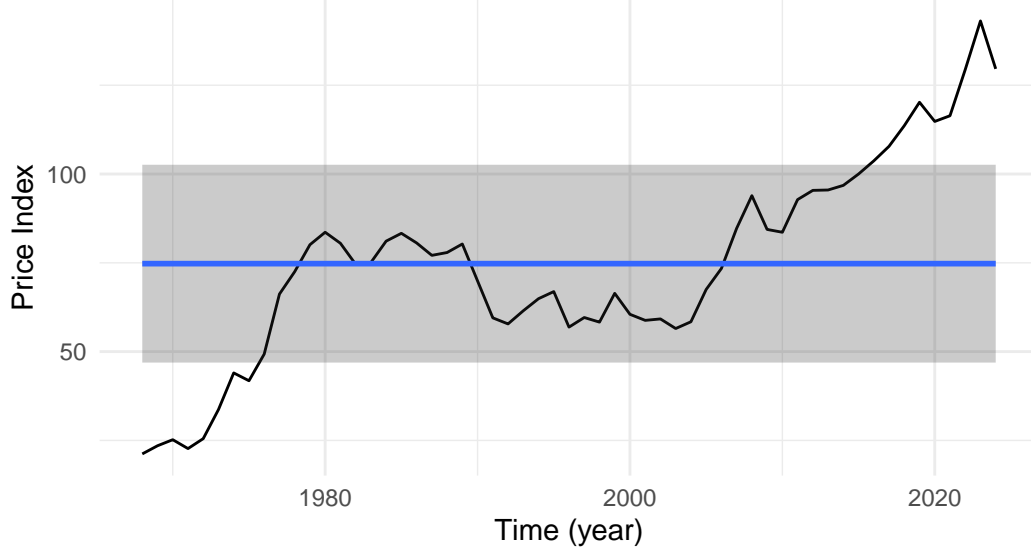


Figure 4: The black line shows the stem timber price index of oak (*Quercus robur* and *Quercus petraea*) in Germany from 1968 - 2024. Data taken from <https://www-genesis.destatis.de/> (Code: 61231-0001). You find it as `stemwood_prices_annually.csv` in the data folder. The blue line shows the arithmetic mean, and the grey band shows the standard deviation.

by linear regression¹.

More formally, a time series is a sequence of T observations $y_t, t = 1, 2, \dots, T$ that are ordered (dependent) in time and which emerge from one random variable (Lütkepohl and Krätzig 2004, 11). Considering this ordering and some heterogeneity assumptions that we will come back to later, in a time series, any observation at any time t is a (so far unknown) function of its history as

$$y_t = f_t(t, y_{t-1}, y_{t-2}, \dots).$$

If we consider this time dependent function as a common function over the entire series, the discrepancy between this function and the actual observation is a stochastic component u_t , which is usually assumed to be an iid error process with mean zero and constant variance σ^2 . Thus, the function can be rewritten as

$$y_t = f(t, y_{t-1}, y_{t-2}, \dots) + u_t,$$

¹Note that many time series methods are actually specific variants of linear regression. We use the term *linear regression* here to mean *ordinary linear regression* without any correction, generalized term, mixed term, etc.

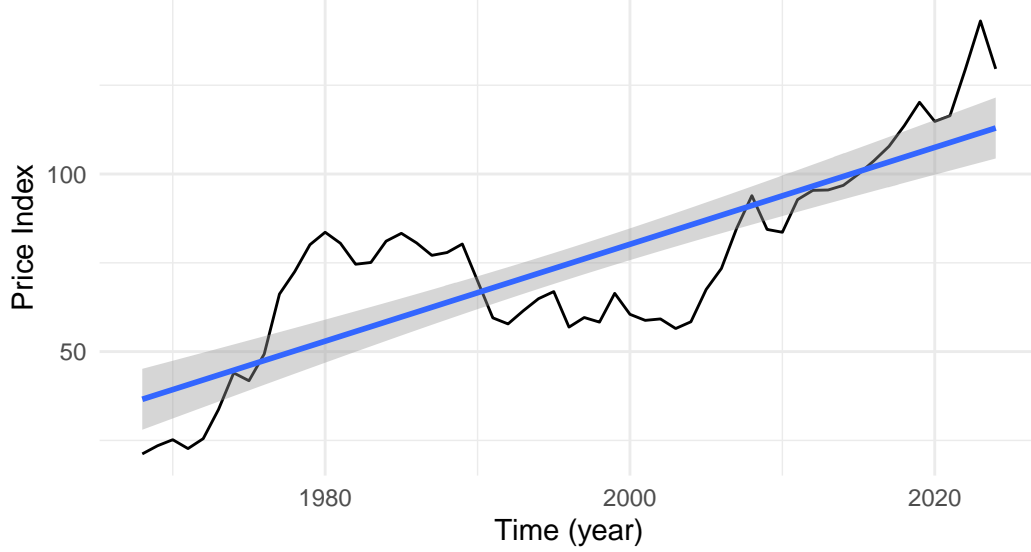


Figure 5: The black line shows the stem timber price index of oak (*Quercus robur* and *Quercus petraea*) in Germany from 1968 - 2024. Data taken from <https://www-genesis.destatis.de/> (Code: 61231-0001). You find it as `stemwood_prices_annually.csv` in the data folder. The blue line shows the linear regression and its standard error (grey).

which means that the entire time series can also be described by a function f and a stochastic component u_t , just as we can do it for any regression. In practice, the function f is limited to a significant lag order P , thus

$$y_t \approx f(t, y_{t-1}, y_{t-2}, \dots, y_{t-P}) + u_t.$$

Theorem 0.1. This representation allows to further distinguish f into a *deterministic* part $g(t)$ and an *autocorrelative* past, as

$$y_t \approx g(t), \alpha_1 y_{t-1} + \alpha_2 y_{t-2}, \dots, \alpha_P y_{t-P} + u_t.$$

$g(t)$ is able to capture e.g. seasonality and/ or a common linear trend and/ or a constant. $g(t)$ captures all components that commonly apply independent from recent historic observations. A time series model with a linear trend (and optionally a constant) and without autocorrelation is thus an ordinary linear model (optionally with intercept) (see Figure 5). The linear regression in the example was able to describe this common trend, but the temporal dynamics remained in the residuals u_t .

Exemplary Time Series and Components of Time Series

When creating time series models, it is particularly important to analyse the characteristics of the series and also to take into account the theoretically assumed characteristics, as different models exist for different data-generating processes in time series statistics (Lütkepohl and Krätzig 2004, 8). The most relevant components that are to be investigated or hypothesised prior modelling are the

- constant components (intercept and/or slope),
- the seasonal component, and
- the autocorrelation component.

We use another example to illustrate the three components and thereby also learn some features for practised programming in R. R accounts for the properties of time series and provides functions for practical programming with time series data, which enables an effective working flow, beginning with a time series data type. The native function `ts` converts a data frame into a time series data type. `ts` requires the data to be ordered and a regular time pattern.

```
stemwood_prices <- read_csv2("data/stemwood_prices_annually.csv")
stemwood_prices <- stemwood_prices |> select(-time) |>
  # Time is not required any more as a column as it is included in the ts object
  ts(start = min(stemwood_prices$time), frequency = 1)
# frequency = 1 as we have annual data
```

The `autoplot` function is a wrapper for the `ggplot2` package, which provides complete plots for particular data types. The class of the object transmitted to the function determines the type of the plot, which can then be further modified using the well-known `ggplot2` syntax. To include the time series feature in `autoplot`, also the `forecast` package is required. `forecast` is a package that contains numerous tools for time series analysis. To get an nice overview over the time series of the prices of stem wood for oak, beech, and spruce, for example, we can use the `autoplot` as follows.

```
library(forecast)
plot_stemwood_prices <- stemwood_prices |>
  autoplot(facets = TRUE, colour = TRUE)
```

Adding elements that might help interpreting the time series data, such as vertical lines, can be done straightforwardly using `geom_vline`. The `annotate` function can be used to add labels to the plot. In the following example, we add the most severe storm events after 2000.

```
plot_stemwood_prices <- plot_stemwood_prices +
  ylab("Price Index") +
  guides(colour = "none") + # legend not necessary as the facets are annotated
  geom_smooth(method = "lm") + # Add linear trends
  theme_minimal() +
  geom_vline(xintercept = c(2000, 2007, 2018)) + # Add storm events
  annotate(x = 2000, y = +Inf, label = "Lothar", vjust = 1, geom = "label") +
  annotate(x = 2007, y = +Inf, label = "Kyrill", vjust = 1, geom = "label") +
  annotate(x = 2018, y = +Inf, label = "Friederike", vjust = 1, geom = "label")

plot_stemwood_prices
```

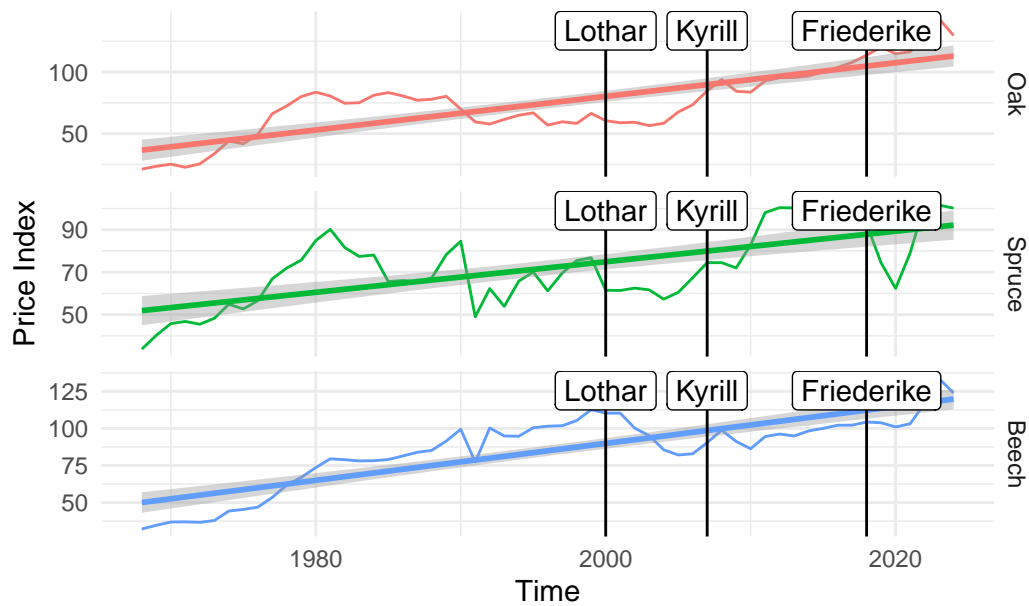


Figure 6: Oak, spruce and Beech stem wood price indices in Germany from 1968 - 2024. Data taken from <https://www-genesis.destatis.de/> (Code: 61231-0001). An example of three series with similar global trend but differing autocorrelations.

🔥 Time series exercise 2

Now it's your turn. Please organize yourself in small groups of 2 - 3 students and choose one species (Norway spruce, Scots pine or European beech) from the forest defoliation data set per group (see Chapter 2_datasets.pdf on Github for more details of the data). It would be great if we would cover all species. Please coordinate with the other groups to ensure that all three types are analysed.

1. Data preparation

- Load the data set into the variable `dat`.
- Filter it to your species.
- Create a univariate time series with the mean loss for each year. Call your `ts` object `dat_*[your species]`.

2. Visualization

- Plot your time series using `autoplot`.

According to the UBA (Umweltbundesamt) (<https://www.umweltbundesamt.de/themen/wasser/extremereignisseklimawandel/trockenheit-in-deutschland-fragen-antworten#trockenheit-aktuelle-situation>), the years 2018, 2019, 2020 and 2022 have been severely dry within the time horizon from 1990 to 2023.

3. Interpretation

- Emphasize the drought years 2018, 2019, 2020 and 2023 in your plot.
- Please present your plot to the colleagues. What are the 3 main findings of your plot?

4. Save your workspace.

205

206 Components of Time Series in More Detail

207 It can be seen that all three series increase by trend (*trend* component), which is emphasised
208 by the 3 linear trend lines. It can also be seen that the series differ fundamentally in terms of
209 their short-time dynamic patterns (autocorrelation). Additionally to the trend, there appears
210 to be a correlation between the observations, a time-dynamic which on a time horizon shorter
211 than the trend. This short term dynamics seem to be different among the three species, in
212 contrast to the trend. While the price index for oak stemwood is relatively stable in terms
213 of short-term time-dynamic pattern (see also Figure 5) and does not react on the events
214 displayed, spruce is more sensible to dynamic pattern including a very severe storm reaction
215 (Friederike) that led to a price decline to the index of 1975. Visual inspection shows that there
216 are obvious trend components and that there might be *autocorrelative* components as well.
217 Seasonal components cannot be followed from this figure. However, the data is annual, and
218 thus, the seasonal component is not expected to be visible in the plot.

219 Industrial wood could be hypothesised to have a *seasonal* trend, as the demand for wood is
220 prospectively higher in winter than in summer, since it is often used as firewood. In forestry,
221 the timber sales prices are usually negotiated on a long-term basis, meaning that short-term

222 demand and supply rarely have a direct impact and that possible seasonal trends are therefore
 223 masked (Fuchs et al. 2022). However, the higher the quality of the wood, the more this
 224 applies. Among all wood assortments, industrial prices are thus most likely to have a seasonal
 225 component. The time series also shows a linear trend, but evolves more slowly and appears to
 226 have a much stronger autocorrelative component. There may also be a seasonal trend, which
 227 appears to be masked by autocorrelative trends in periods with higher fluctuation.

```
ind_prices <- read_csv2("data/industrialwood_prices_monthly.csv")
ind_prices <- ind_prices |> select(Spruce) |>
  ts(start = min(ind_prices$year), frequency = 12) # Monthly data

ind_prices |>
  autoplot() +
  guides(colour = "none") + # A legend is not necessary.
  theme_minimal() + geom_smooth(method = "lm") + ylab("Price Index")
```

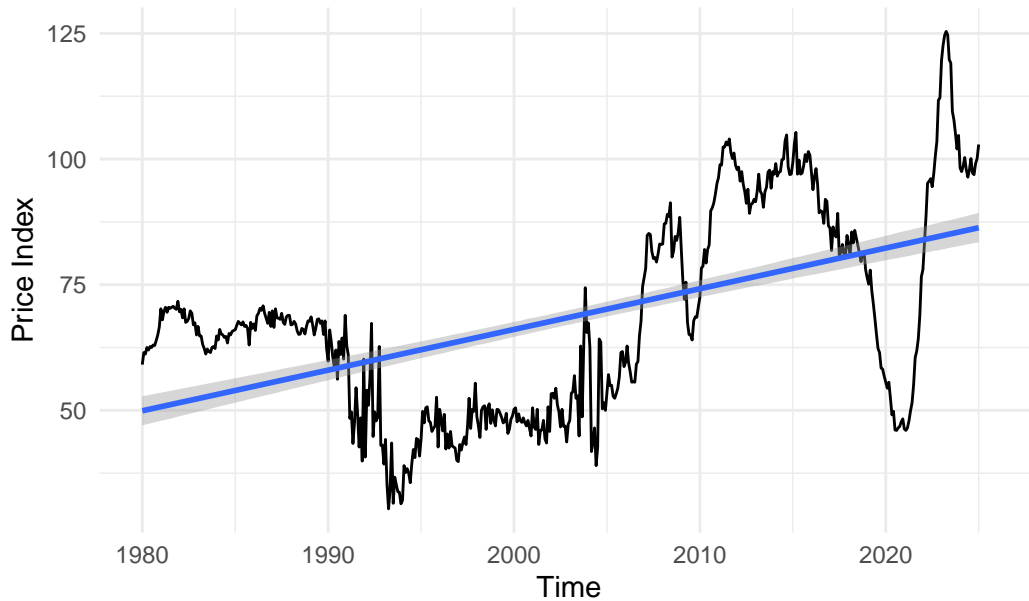


Figure 7: Monthly spruce industrial wood price index in Germany from 1980 - 2024. Data taken from <https://www-genesis.destatis.de/> (Code: 61231-0002).

228 The seasonal trend is (among the linear trend and constant trend) another typical fixed term.
 229 Fixed term is defined as common trend that appears for the time series in general (i.e. every-
 230 thing that is not directly correlated with the historic observations). A clear *saisonal* component
 231 can e.g. be found in the air temperature data of the German Weather Service (DWD). Temper-
 232 ature data with a resolution finer than a year is a common example of seasonality that needs

233 to be taken into account in the inference model. The data is available at a daily base and
234 can be used to illustrate the concept of seasonality. The following code visualises the mean
235 air temperature at the weather station of Göttingen. The data contains a strong seasonal
236 (weather pattern in the annual seasons) and a strong trend component (climate change).

```
temp_goe <- read_csv2("data/month_mean_temp_goe.csv")
temp_goe <- temp_goe |> select(mean_daymean_temp) |>
  ts(start = c(min(temp_goe$year), 1), # Starting year = min year
     # Starting month = Jan
     frequency = 12) # monthly data

temp_goe |>
  autoplot() +
  guides(colour = "none") + # legend not necessary as the facets are annotated
  theme_minimal() + geom_smooth(method = "lm") + ylab("Mean Temperature (°C)")
```

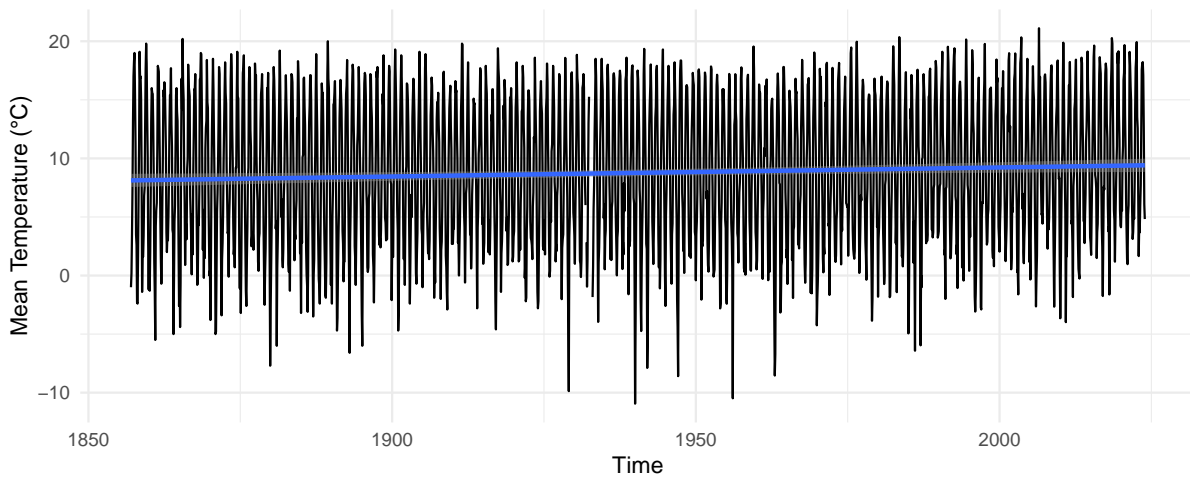


Figure 8: Monthly mean air temperature (mean of day means) at 2m height at the weather station of Göttingen from January 1857 till December 2023. Taken from (https://opendata.dwd.de/climate_environment/CDC/observations_germany/climate/monthly/kl/historical/).

237 We can use the native function `window` to extract a part of the time series.

```
temp_goe |> window(start = c(2000, 1), end = c(2020, 12)) |>
  autoplot() +
  guides(colour = "none") + # legend not necessary as the facets are annotated
  theme_minimal() + geom_smooth(method = "lm") + ylab("Mean Temperature (°C)")
```

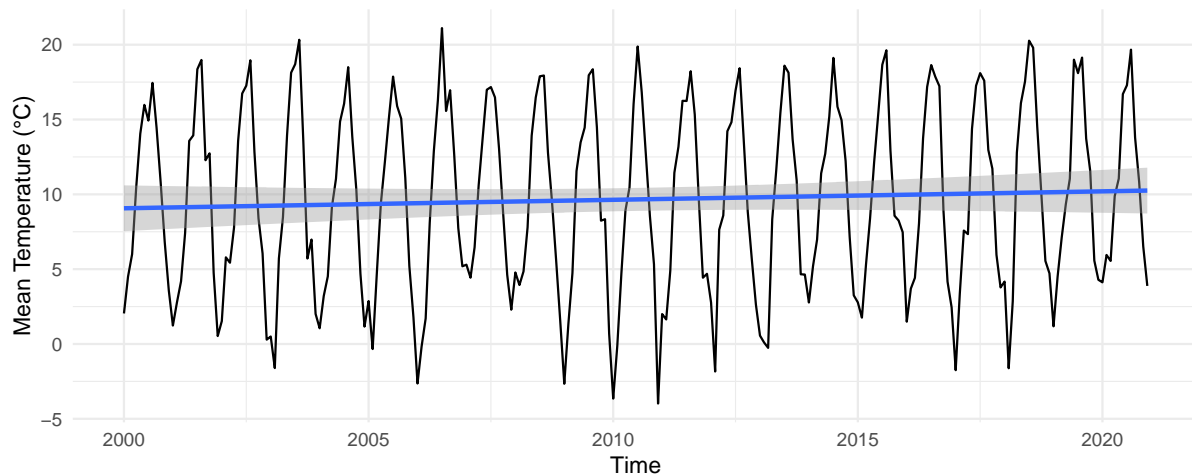


Figure 9

🔥 Time series exercise 3



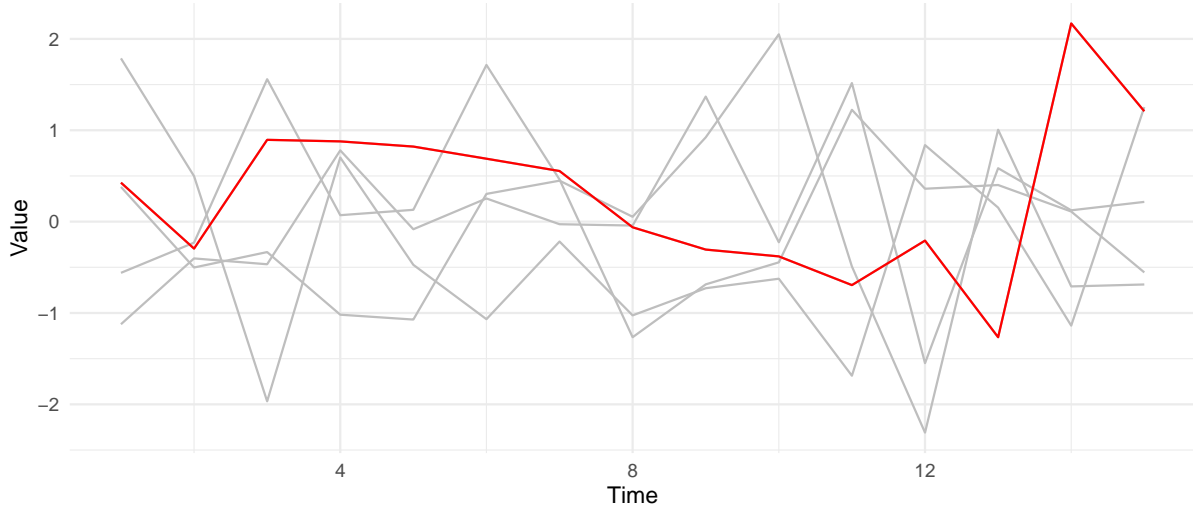
Go to <https://flinga.fi/s/FCZ78B4>

1. Imagine one time series. Decide which of the components (trend, seasonality, autocorrelation) apply.
2. Describe your time series briefly (heading + ~ 1-2 lines) and write the description into a purple sticky note. Arrange your sticky notes in a horizontal line.
3. Consider whether your data can be analysed using methods that are non-time-series, or if you need to make use of time series methods. Don't write down your answer - just think about it.
4. If possible, suggest a non-time-series method for analysing your data series on a blue sticky note. Stick this note somewhere, but not directly under your purple sticky note.
5. Now let's go through all the notes together in class.
 - Which data series is analysable using non-time-series-methods?
 - If so, which method would suit?
 - What would be the additional information/ the advantage of a time series method instead?

Stationarity

How is the concept of stationarity related to the three components of time series?

To understand time series data and the assumptions of the time series model, a time series can be thought of as a stochastic process with a latent data generating process (the population) and a realisation at the level of the random sample (Figure 3), just as we did for the other data types. In contrast to the populations so far, an observed time series y_t , $t = 1, \dots, T$ is regarded as one realization of a finite part of a stochastic process $y_t(\omega)$ (Lütkepohl and Krätzig 2004, 10, 11). We do only describe the stochastic theory very superficially. You can find a deeper insight and references to the textbooks from which the time series theory originate in Lütkepohl and Krätzig (2004, 11). A stochastic time series process is **stationary** if all of its members are mutually independent, which in particular for time series processes means that all members are time invariant. Stationary observations are independent. Such a stationary process, also called white noise, would generate observations that fluctuate around a stable mean and have a constant variance. Such a process would meet all assumptions (iid, common variance) that we have talked about so far (Chapter 2, see also Figure 3) and would not require time series methods. The ordinary (unconditional) arithmetic mean, calculated from any realised series, would be an unbiased estimator of the population mean. The same would appear for the standard deviation. In reality, of course, we never know whether our apparently observed non-stationary time series has arisen from a stationary process, or whether it really has arisen from a non-stationary process (see also Chapter 2). Consider the following 5 simulated observations emerging from a stationary process with mean 0 and a standard deviation of 1. All of these 5 series have a mean close to 0, of course. Indeed, the mean and standard deviation would thus provide unbiased estimates for the population but for the red line, as an example, it is difficult to recognise visually that it has emerged from a stationary process. The line could also be interpreted as a increasing trend or autocorrelation. This problem occurs to any observed time series. While it sometimes seems obvious that there is an autocorrelation component (e.g. Figure 7), a seasonal component (e.g. Figure 9), or a trend component (e.g. Figure 5), in fact this is no clear advice that a series does not evolve from a stationary process. When it comes to testing for stationarity, we must remember that we are only testing the realisation, never the population, in the sense of *how likely is it that this realisation could arise from a stationary process?*



270

271 More formally², stationarity means that each member of a series in the population has the
 272 same expectation and expected variance (homoscedasticity).

$$E[y_1] = E[y_2] = \dots = E[y_T] = \mu$$

$$Var[y_1] = E[(y_1 - \mu)(y_1 - \mu)] = Var[y_2] = \dots = Var[y_T] = \gamma_0$$

273 From which follows that the covariance between two arbitrary members y_t and y_{t+h} is a
 274 function of the lag h only. h is the difference between two time points within one series. The
 275 covariance is called the autocovariance and is denoted as γ_h .

$$Cov[y_{1+h}, y_1] = E[(y_{1+h} - \mu)(y_1 - \mu)] = Cov[y_{2+h}, y_2] = \dots = Cov[y_T, y_{T-h}] = \gamma_h$$

276 Under the assumption of stationarity, the ordinary descriptive statistics for iid sampled statis-
 277 tics are therefore appropriate to describe time series as well.

- 278 • Mean: $\hat{\mu} = \bar{y} = \frac{1}{T} \sum_{t=1}^T y_t$
- 279 • Variance: $\hat{\gamma}_0 = \frac{1}{T-1} \sum_{t=1}^T (y_t - \bar{y})^2$
- 280 • Covariance, also called autocovariance: $\hat{\gamma}_h = \frac{1}{T} \sum_{t=1}^{T-h} (y_{t+h} - \bar{y})(y_t - \bar{y}), h = 1, 2, \dots$

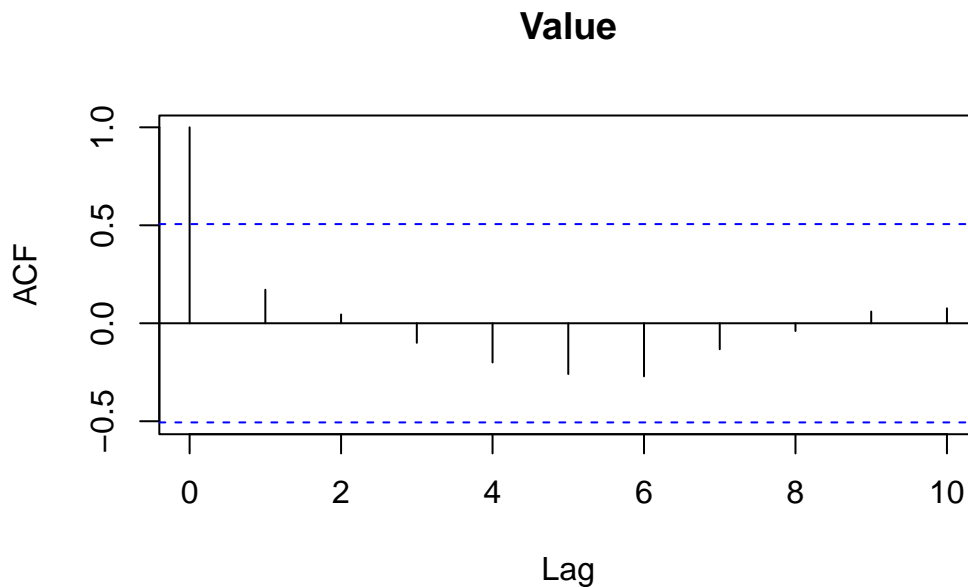
281 By convention, it is not common to do finite correction in time series analyses. The autocor-
 282 relation is calculated as the relation between the covariance and variance.

- 283 • Autocorrelation (AC): $\hat{\rho}_h = \frac{\hat{\gamma}_h}{\hat{\gamma}_0}, h = 1, 2, \dots$

²the formulations are mainly taken from Lütkepohl and Krätzig (2004, 12 ff.) but aligned to the nomenclature of the course

284 A white noise process would lead to a stationary time series with common mean $\hat{\mu}$ for all ob-
 285 servations and time-invariant variance $\hat{\gamma}_0$, as mentioned earlier, and also to an autocorrelation
 286 of 0. Consider for example the autocorrelation of the red time series. R comes with a native
 287 function `acf` that calculates the autocorrelation for h from 0 to `lag.max`. It can be seen that
 288 except for $h = 0$, which is of course always 1, the autocorrelation is close to 0 for all h . The
 289 set of ordered autocorrelations with increasing h , is also called autocorrelation function. The
 290 autocorrelation functions helps in identifying the time-dynamic component of a time series.
 291 The red series appears to be *stationary* indeed.

```
example_whitenoise_red_ts |> acf(lag.max = 10)
```



292

293 Note that `urca` package (among others) provides a unit root test to perform a statistical test
 294 for stationarity. However, we will not capture this or any other testing procedure in this
 295 course.

Time series exercise 4

Join your group from Exercise 2 again. Load your workspace from Exercise 2.

1. Create an autocorrelation function (ACF) for your species with a lag order of 24 months.
2. Considering both, the plot created in exercise 2 and the ACF. Which components do you expect to have in your series?
3. Save your workspace.

296

297 A white noise process with linear trend but without autocorrelation, is called *trend stationary*
 298 in time series statistics. Such a trend would be sufficiently described by means of an ordinary

linear regression. Or in other words: A stationary process shows no autocorrelation after correcting for the linear trend. Followingly, the residuals of a linear regression would be stationary. The function `tslm` can be used to wrap an `lm` for `ts` objects. However, putting the `ts` object in `lm` directly would also work. You then need to define the years as the only covariate. Correcting for the linear trend of the stem wood prices (Figure 6), for example, leads to the following time series and autocorrelation functions.

```
# Calculate lm and save the residuals
detrended_stemwood_prices <- ts.union(
  tslm(stemwood_prices[, "Oak"] ~ trend) |> residuals(),
  tslm(stemwood_prices[, "Spruce"] ~ trend) |> residuals(),
  tslm(stemwood_prices[, "Beech"] ~ trend) |> residuals()

# Keep the original names
colnames(detrended_stemwood_prices) <- colnames(stemwood_prices)

detrended_stemwood_prices |>
  autoplot(facets = TRUE, colour = TRUE) + facet_wrap(~ series) +
  ylab("Price Index") +
  guides(colour = "none") + # legend not necessary as the facets are annotated
  theme_minimal()
```

- ① `ts.union` is the `cbind`-pendant for `ts` objects.
- ② Calculate a regression model with `trend = Detrending`.
- ③ Store (only) the residuals of that model.

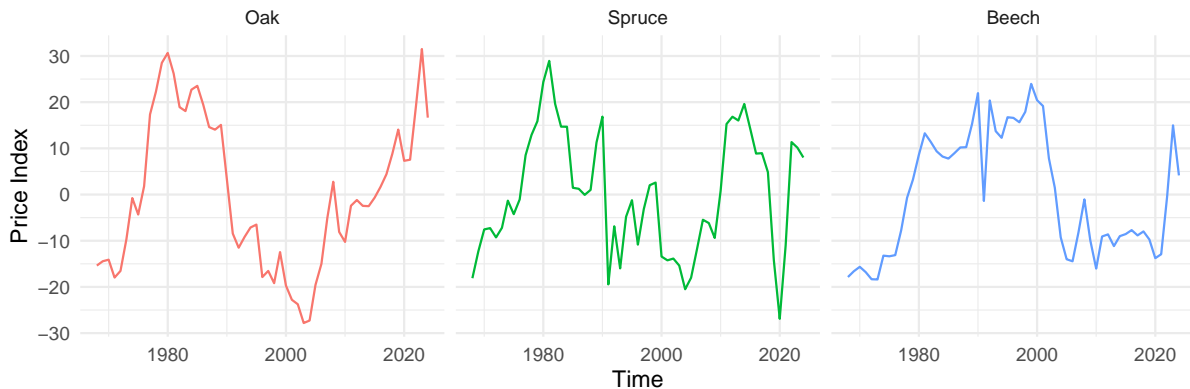


Figure 10: Detrended Oak, spruce and Beech stem wood price indices in Germany (the original series are shown in Figure 6).

Here we see that there is indeed still evidence of autocorrelation after detrending. All three

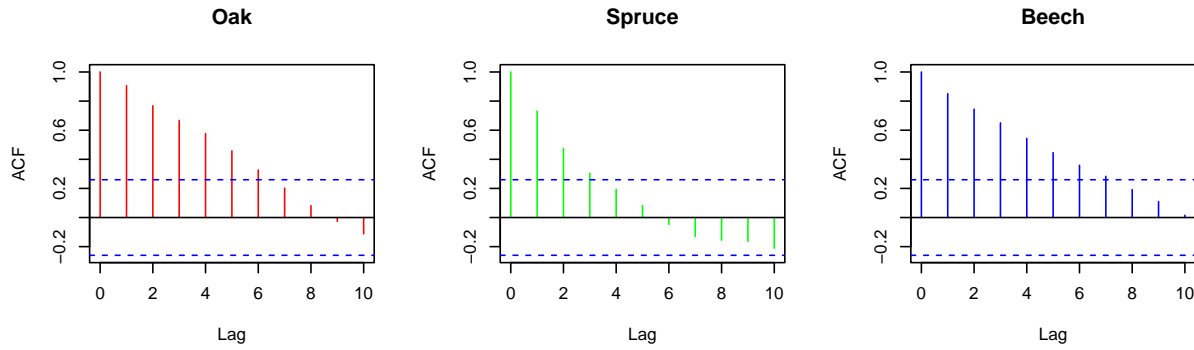
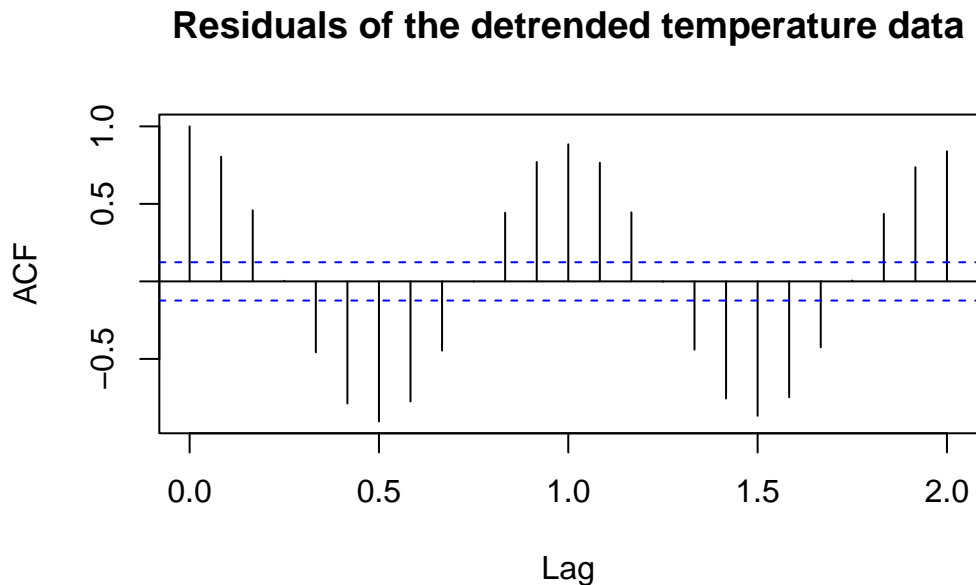


Figure 11: The respective autocorrelation functions of the detrended series from Figure Figure 10 using the acf function.

species tend to have decreasing autocorrelations with increasing lag (h). Note that the dashed lines cannot be interpreted as confidence intervals in the sense of *significant correlation must be above a critical value* even if this is sometimes proclaimed in scientific literature. It is the autocorrelation under perfect noise. Nevertheless, the autocorrelation functions indicate that none of the time series are stationary or trend stationary. Doing the same (detrending and then calculating an autocorrelation function using `acf`) for the windowed temperature data (Figure 9) leads to the following autocorrelation function.

```
tslm(temp_goe |> window(start = c(2000, 1), end = c(2020, 12)) ~ trend) |>
  residuals() |> acf(lag.max = 24, main =
    "Residuals of the detrended temperature data")
```



316

317 Note that a `lag.max` of 24 means that h is set to 2 years (= 24 months). As expected, the
 318 autocorrelation is close to 1 after a full cycle (year) and highly negative correlated after the
 319 half cycle. This series is thus also neither stationary nor trend stationary. Additionally to the
 320 autocorrelation function of the timber wood prices (Figure 6), the autocorrelation function
 321 reveals a seasonal component. It can be followed that the time series is not trend stationary.
 322 Yet, it remains unclear whether this strong remaining autocorrelation after detrending is only
 323 due to the seasonal component (*the temperature in one month is autocorrelated with the same*
 324 *month of the previous year*) or whether the series is also autocorrelated in the sense of *the*
 325 *temperature in one month is autocorrelated with the temperature of the previous months*. The
 326 seasonal component can be removed in the same way as the trend component. A linear
 327 regression using only dummy variables, one dummy for each point in the cycle, 12 months in
 328 our example. The `tslm` function saves some programming effort here, as it automatically uses
 329 the frequency information of the `ts` object to create the number of dummy variables. A model
 330 containing this seasonal component and also the trend can be parameterised as follows:

```

tslm(temp_goe |>
  window(start = c(2000, 1), end = c(2020, 12)) ~ trend + season) |>
  residuals() |> acf(lag.max = 24, main =
    "Resid. of the detrended & seasons-corrected temp.")
  
```

Resid. of the detrended & seasons-corrected temp.

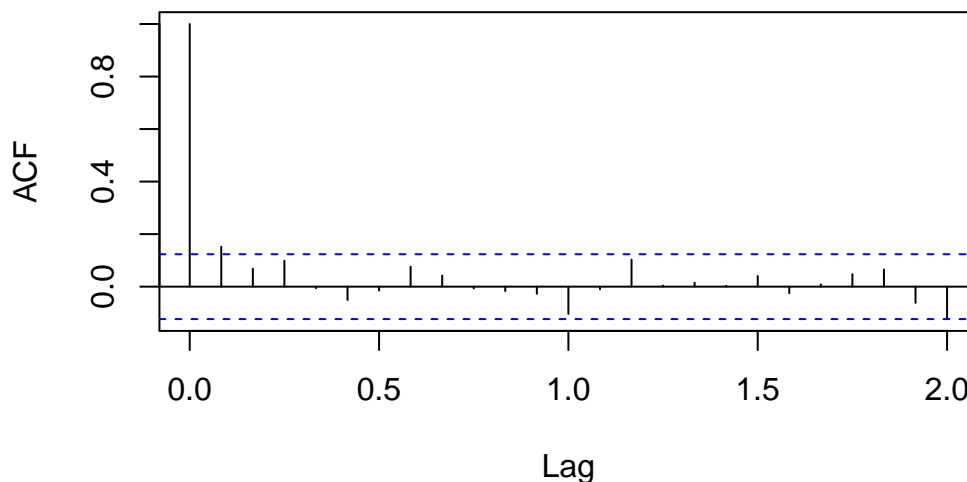


Figure 12: Autocorrelation function of the residuals of the detrended and seasons-corrected temperature data in Göttingen.

331 A look at the autocorrelation function of the residuals (Figure 12) reveals that the autocor-
 332 relation is now close to 0 for all h . There is evidence for stationary of the residuals. The

temperature data thus mainly consists of the season and the trend component. There is no further time dynamic in the data then the trend and the season. We can access the parameters of the trend and the season component just as we do it in lms. `summary` of the model gives us:

```
tslm(temp_goe ~ trend + season) |> summary()
```

```
337 Call:
338 tslm(formula = temp_goe ~ trend + season)
339
340 Residuals:
341      Min       1Q   Median       3Q      Max
342 -11.7326  -1.0655   0.0466   1.1868   5.4121
343
344 Coefficients:
345             Estimate Std. Error t value Pr(>|t|)
346 (Intercept) -2.878e-01  1.676e-01  -1.717 0.086152 .
347 trend        6.284e-04  7.489e-05   8.390 < 2e-16 ***
348 season2      7.927e-01  2.122e-01   3.735 0.000193 ***
349 season3      3.793e+00  2.122e-01  17.872 < 2e-16 ***
350 season4      7.847e+00  2.122e-01  36.972 < 2e-16 ***
351 season5      1.235e+01  2.122e-01  58.201 < 2e-16 ***
352 season6      1.553e+01  2.122e-01  73.191 < 2e-16 ***
353 season7      1.707e+01  2.126e-01  80.324 < 2e-16 ***
354 season8      1.649e+01  2.126e-01  77.571 < 2e-16 ***
355 season9      1.317e+01  2.122e-01  62.042 < 2e-16 ***
356 season10     8.715e+00  2.126e-01  40.999 < 2e-16 ***
357 season11     4.147e+00  2.126e-01  19.512 < 2e-16 ***
358 season12     1.147e+00  2.129e-01   5.387 8.02e-08 ***
359 ---
360 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
361
362 Residual standard error: 1.939 on 1985 degrees of freedom
363 (6 Beobachtungen als fehlend gelöscht)
364 Multiple R-squared:  0.9098,    Adjusted R-squared:  0.9092
365 F-statistic: 1668 on 12 and 1985 DF,  p-value: < 2.2e-16
```

Plotting these parameters and the residuals provides us graphical evidence of the relevance of the 3 components. In our example we see that there is a slight but significant trend component (climate change) and a strong and significant seasonal component. The remainder appears to be white noise only at first sight and by consideration of the autocorrelation function (Figure 12).

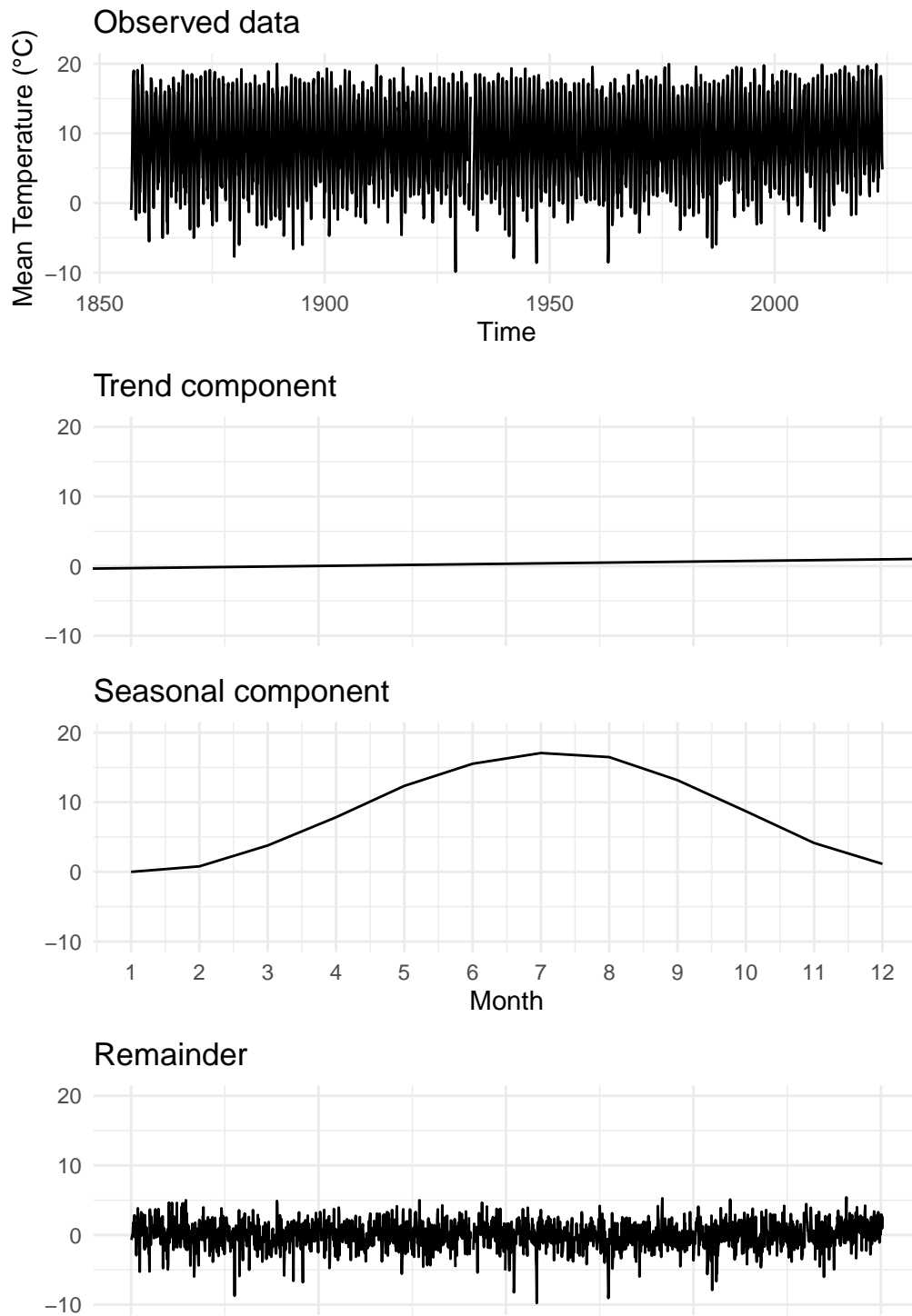


Figure 13: Raw data, trend component, season component and remainder to visualise all components of time series data. Note that the seasonal component has an x-axis different to the other diagrams in order to better visualise the annual development.

Time series exercise 5

Join your group from Exercise 4 again. Load your workspace from Exercise 4.

1. Which of the 3 components (trend, season, autocorrelation) do you expect in your time series (see Exercise 4)?
2. Use `ts.lm` and `acf` to test your expectations.
3. Save your workspace.

371

Descriptive Statistics and statistical modeling

Classical Decomposition

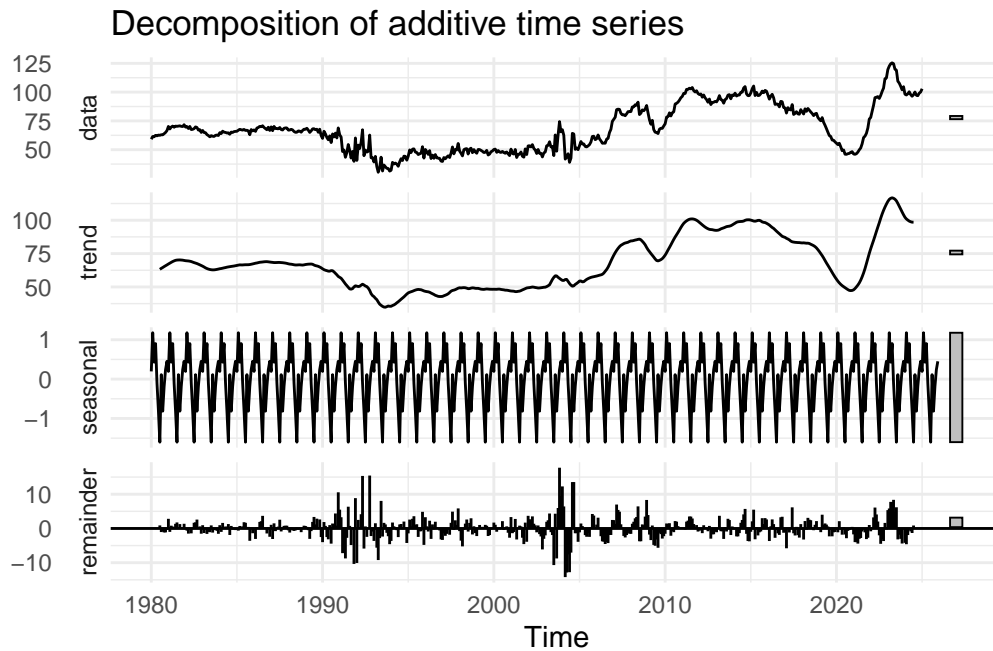
To decompose a time series into its components *trend*, *season* and *autocorrelation*, as we have done it in the last subchapter, is a commonly used technique to check whether time series statistics need to be applied to a series or whether ordinary models are sufficient. The aim is to determine whether there is an autoregressive time pattern or if the series has deterministic components (trend, season, and also constant) only. The so called classical decomposition is a set of descriptive statistics in time series statistics. In the previous subchapter, we developed a simple additive decomposition with a linear trend and a linear seasonal component. Others, such as polynomial trends or trigonometric seasonal components, are also commonly used. The native R functions `decompose` or `stl`, among others, provide numerous methods for decomposing a time series and for visualisation. The ordinary linear model that we parameterised above can be used if the following assumption of additivity holds:

$$y_t = m_t + s_t + u_t$$

where y_i is the observed time series, m_t is the trend component, s_t is the seasonal component, and u_t is the remainder. In general, any regression model can be used to decompose a time series into deterministic components and the possibly autocorrelated remainder (residuals). The most simple and straightforward model is a linear model with one parameter for the trend, as we have already performed. The native R function `decompose` used a symmetric moving average approach to estimate the trend. Advantage of the moving average approach is that autocorrelation (firstly in this course) is considered as the trend is calculated by means of last P observations. Per default, the last 6 observations are used with equal weights. Disadvantage is that we do not get a parameter for the trend component. `decompose` does not deliver any parameter information. The seasonal trend is then estimated by means of a linear model, just as we did in the last subchapter.

Decomposition of our industrial wood price of spruce (Figure 7) leads to the following picture:

```
decompose(ind_prices, type = "additive") |> autoplot() + theme_minimal()
```

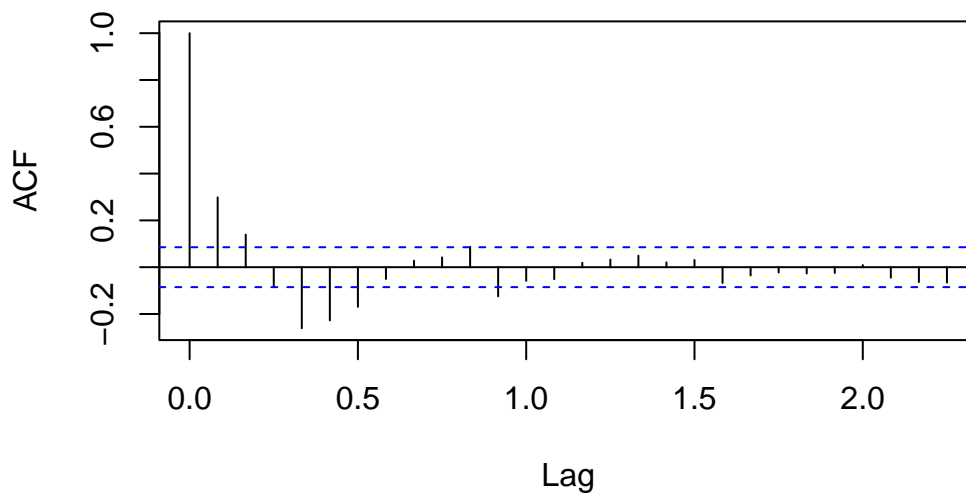


397

398 and to the following autocorrelation function of the remainder:

```
d <- decompose(ind_prices, type = "additive")
d$random |> na.omit() |> acf()
```

Series `na.omit(d$random)`



399

400 Which is already pretty good in describing this relatively complex data set.

Time series exercise 6

Join your group from Exercise 5 again. Load your workspace from Exercise 5.

1. Perform a classical decomposition of your series.
2. Plot a autocorrelation function of the remainder.
3. Compare the results to the results of Exercises 4 and 5. What is different? What is similar? Do you come to the same conclusions?
4. Save your workspace.

Autoregressive Models (AR)

In case there is autocorrelation in the data, even after correcting for the deterministic components, ARMA models are capable to estimate those autocorrelative terms via coefficients. Autoregressive (AR) and moving average (MA) models are commonly used to estimate coefficients for autocorrelative terms in univariate time series. According to Verbeek (2004, 279), a population that follows a ARMA process can be estimated by means of ordinary or nonlinear least squares, or by maximum likelihood. The `Arima` function from the package `forecast` uses the maximum likelihood approach as a default (`R` comes with the native function `arima` which has the same functionality as `Arima` - however, `Arima` provides it is streamlined with further functionalities of the `forecast` package, such as visualisation and post processing). AR models of order lag order p can be expressed as

$$Y_t = \delta + \theta_1 Y_{t-1} + \theta_2 Y_{t-2} + \dots + \theta_p Y_{t-p} + \varepsilon_t,$$

which is a consistent estimator for our population. ε_t are white the noise remainders, δ is the intercept, and θ_i are the coefficients to be estimated. Estimation of an autoregressive model is thus no different than that of a linear regression model with a lagged dependent variable. Consider e.g. Verbeek (2004, 280) for more details on the estimation. Even though MA models are almost as common as AR models, we will not cover them in this course. Take a look at Verbeek (2004, 106, 281) if you want to learn more about MA data structures in comparison to AR data and about MA modelling. We thus only use the `ar` part of the `Arima` function. We thus do not have to choose the type of model but only the deterministic components and the lag order for the autocorrelation. Both can be done visually by the time series plot and the ACF.

Let's come back to the annual oak prices. We have not yet found a satisfactory model. There was an obvious overarching linear trend towards higher prices, but also a certain temporal remainder that is obviously not a seasonal trend (Figure 6, Figure 10, Figure 11).

detrended_stemwood_prices

```
426 Time Series:
427 Start = 1968
428 End = 2024
429 Frequency = 1
430      Oak      Spruce      Beech
431 1968 -15.3764065 -18.10828796 -17.8569873
432 1969 -14.4402191 -12.32741552 -16.6015686
433 1970 -14.1040316 -7.54654308 -15.6461499
434 1971 -17.9678442 -7.26567064 -16.7907311
435 1972 -16.5316567 -9.28479820 -18.3353124
436 1973 -9.6954693 -7.20392576 -18.3798937
437 1974 -0.7592818 -1.32305332 -13.2244750
438 1975 -4.3230944 -4.24218088 -13.3690563
439 1976 1.8130931 -1.06130844 -13.1136375
440 1977 17.3492805 8.51956400 -7.7582188
441 1978 22.2854680 12.90043644 -0.7028001
442 1979 28.5216554 15.88130888 3.2526186
443 1980 30.6578429 24.36218132 8.6080373
444 1981 26.1940303 28.94305376 13.2634561
445 1982 18.9302178 19.62392619 11.4188748
446 1983 18.0664052 14.70479863 9.3742935
447 1984 22.7025927 14.68567107 8.2297122
448 1985 23.5387801 1.46654351 7.7851309
449 1986 19.4749676 1.24741595 8.9405496
450 1987 14.6111550 -0.07171161 10.1959684
451 1988 14.0473425 1.00916083 10.2513871
452 1989 15.0835299 11.29003327 15.3068058
453 1990 3.3197174 16.87090571 21.9622245
454 1991 -8.4440952 -19.44822185 -1.3823568
455 1992 -11.5079077 -6.86734941 20.3730620
456 1993 -9.1717202 -15.98647697 13.7284807
457 1994 -7.1355328 -4.80560453 12.2838994
458 1995 -6.4993453 -1.22473209 16.7393181
459 1996 -17.8631579 -10.84385965 16.5947368
460 1997 -16.5269704 -3.16298721 15.6501556
461 1998 -19.1907830 2.01788523 17.9055743
462 1999 -12.4545955 2.59875767 23.9609930
463 2000 -19.7184081 -13.42036989 20.4164117
464 2001 -22.7822206 -14.23949745 19.1718304
465 2002 -23.7460332 -13.85862501 7.8272492
```

466	2003	-27.8098457	-15.37775257	1.4826679
467	2004	-27.2736583	-20.49688013	-9.2619134
468	2005	-19.5374708	-18.01600769	-14.0064947
469	2006	-15.0012834	-11.83513525	-14.4510760
470	2007	-5.0650959	-5.45426281	-8.1956572
471	2008	2.7710915	-6.17339037	-1.0402385
472	2009	-8.0927210	-9.39251793	-9.8848198
473	2010	-10.2565336	0.78835451	-16.0294011
474	2011	-2.4203461	15.26922695	-9.0739824
475	2012	-1.1841587	16.85009939	-8.6185637
476	2013	-2.4479712	16.03097183	-11.1631449
477	2014	-2.5117838	19.61184427	-9.0077262
478	2015	-0.6755963	14.29271671	-8.5523075
479	2016	1.6605911	8.87358915	-7.6968888
480	2017	4.3967786	8.95446158	-8.8414701
481	2018	8.8329660	4.83533402	-7.9860513
482	2019	14.0691535	-14.18379354	-9.7306326
483	2020	7.3053409	-26.90292110	-13.7752139
484	2021	7.5415284	-11.02204866	-12.9197952
485	2022	19.1777158	11.35882378	-0.4643765
486	2023	31.5139033	10.13969622	14.9910423
487	2024	16.6500907	8.02056866	4.1464610

488 **Regression with time dynamics - Temporal**
489 **regression**

490 Upcoming

References

- Fieberg, John R, Kelsey Vitense, and Douglas H Johnson. 2020. "Resampling-Based Methods for Biologists." *PeerJ* 8: e9089.
- Fuchs, Jasper M, Hilmar v Bodelschwingh, Alexander Lange, Carola Paul, and Kai Husmann. 2022. "Quantifying the Consequences of Disturbances on Wood Revenues with Impulse Response Functions." *Forest Policy and Economics* 140: 102738.
- Hamilton, James D. 2020. *Time Series Analysis*. Princeton university press.
- Lemoine, Nathan P. 2021. "Unifying Ecosystem Responses to Disturbance into a Single Statistical Framework." *Oikos* 130 (3): 408–21.
- Lütkepohl, Helmut, and Markus Krätzig. 2004. *Applied Time Series Econometrics*. Cambridge university press.
- Verbeek, Marno. 2004. "A Guide to Modern Econometrics. Erasmus University Rotterdam." *John Wiley & Sons Ltd., Hoboken, in: Organizational Behavior and Human Decision Processes* 67: 326–44.