

Lecture 13:

Time Series Analysis

Olexandr Isayev






Department of Chemistry, CMU

olexandr@cmu.edu

Applications

The fields of application of Time series Analysis are numerous:

Demand Planning is one of the most common application, however, from industry to industry there are other possible uses. For instance:

	Logistics & Transportation	▪ Forecasting of shipped packages : <u>workforce planning</u>
	Retail grocery	▪ Forecasting of sales during promotions : <u>optimizing warehouses</u>
	Insurance	▪ Claims prediction : <u>determining insurance policies</u>
	Manufacturing	▪ Predictive Maintenance : <u>improving operational efficiency</u>
	Energy & Utilities	▪ Energy load forecasting : <u>better planning and trading strategies</u>

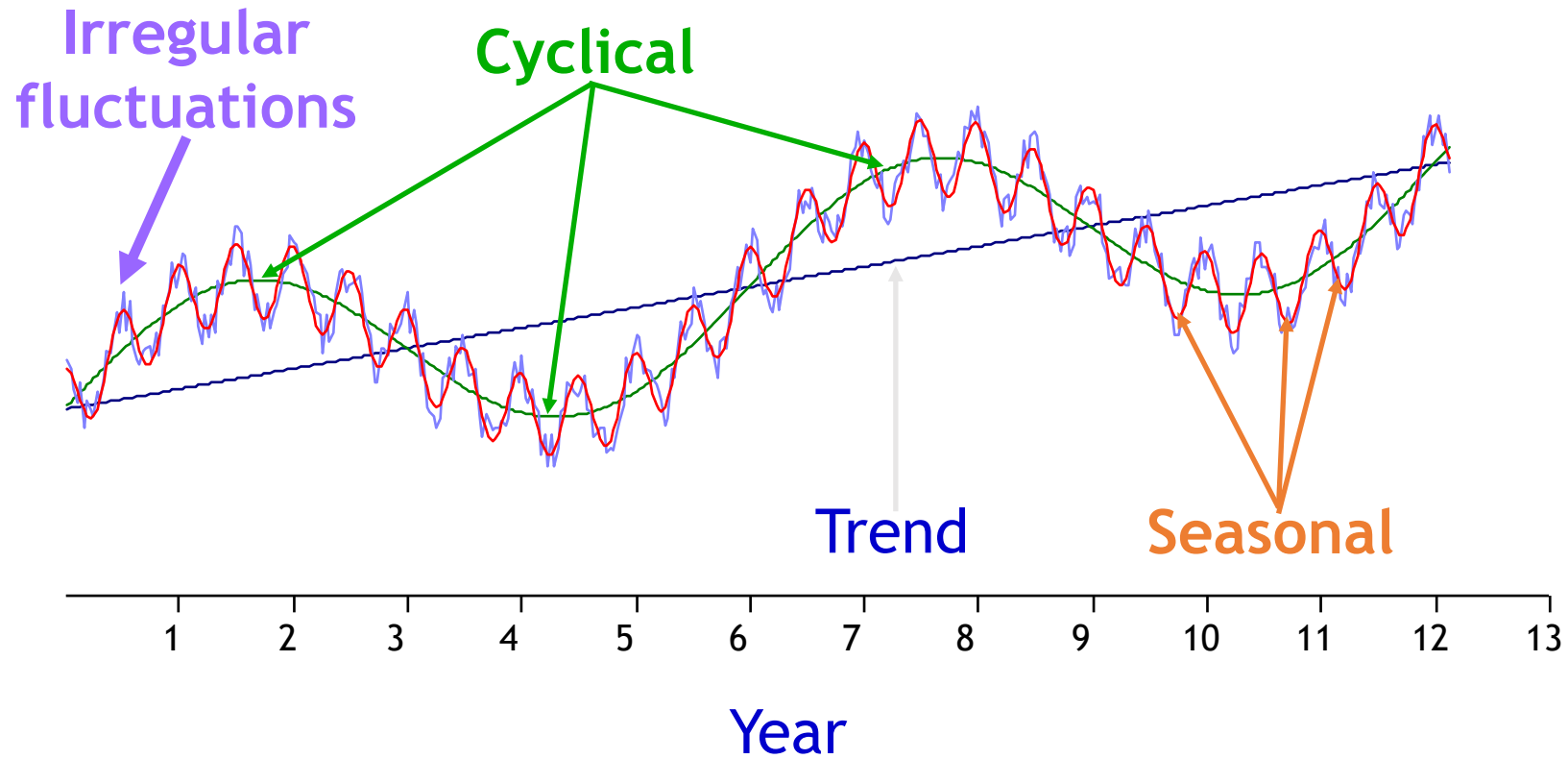
Time series

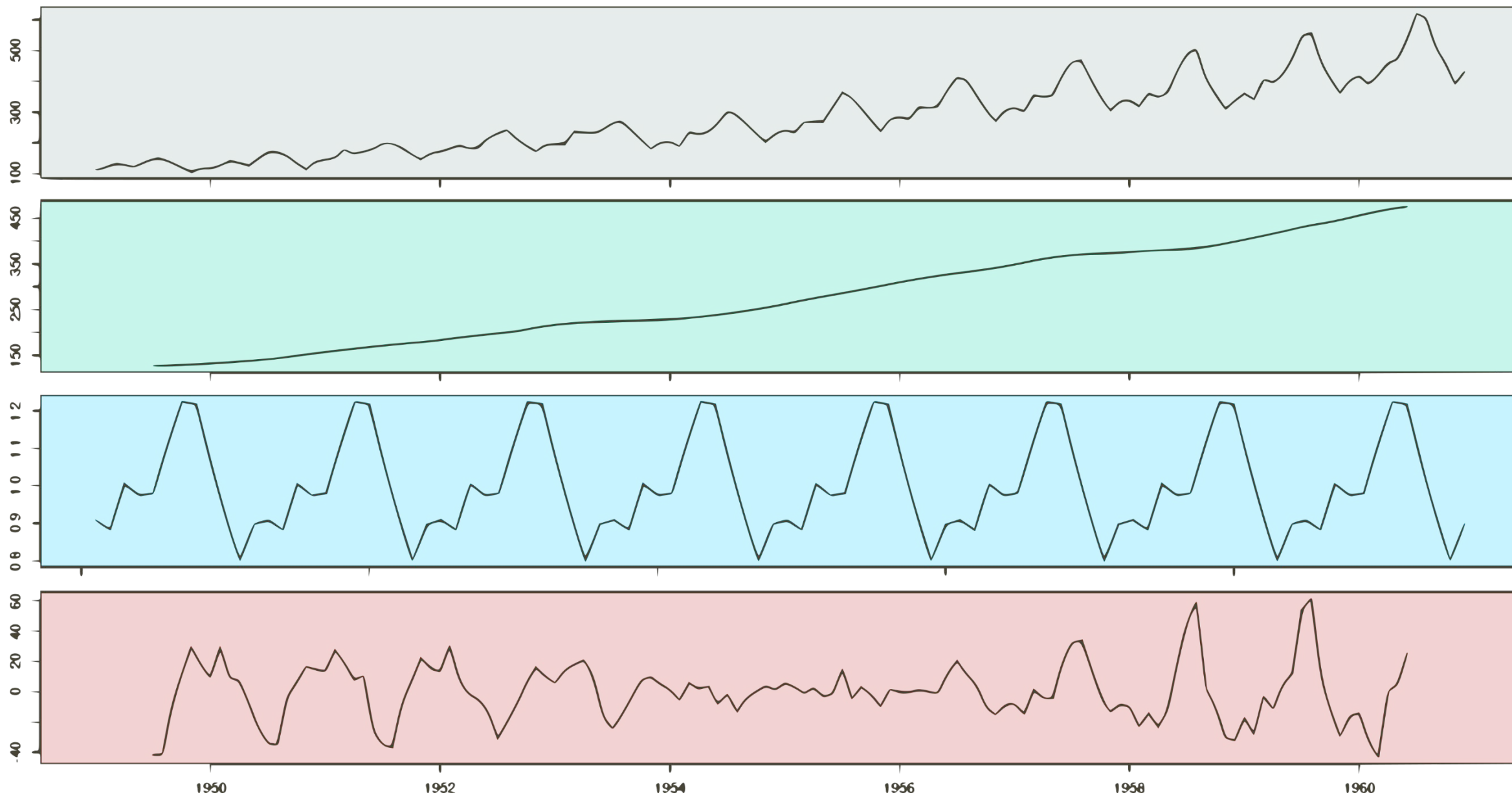
A time series is a series of data points indexed in time order.

Most commonly, a time series is a sequence taken at successive **equally spaced points** in time.

In a time series, time is often the independent variable and the goal is usually to make a forecast for the future.

Components of Time-Series Data





Typical questions

Is it **stationary**?

Is there a **seasonality**?

Is the target variable **autocorrelated**?

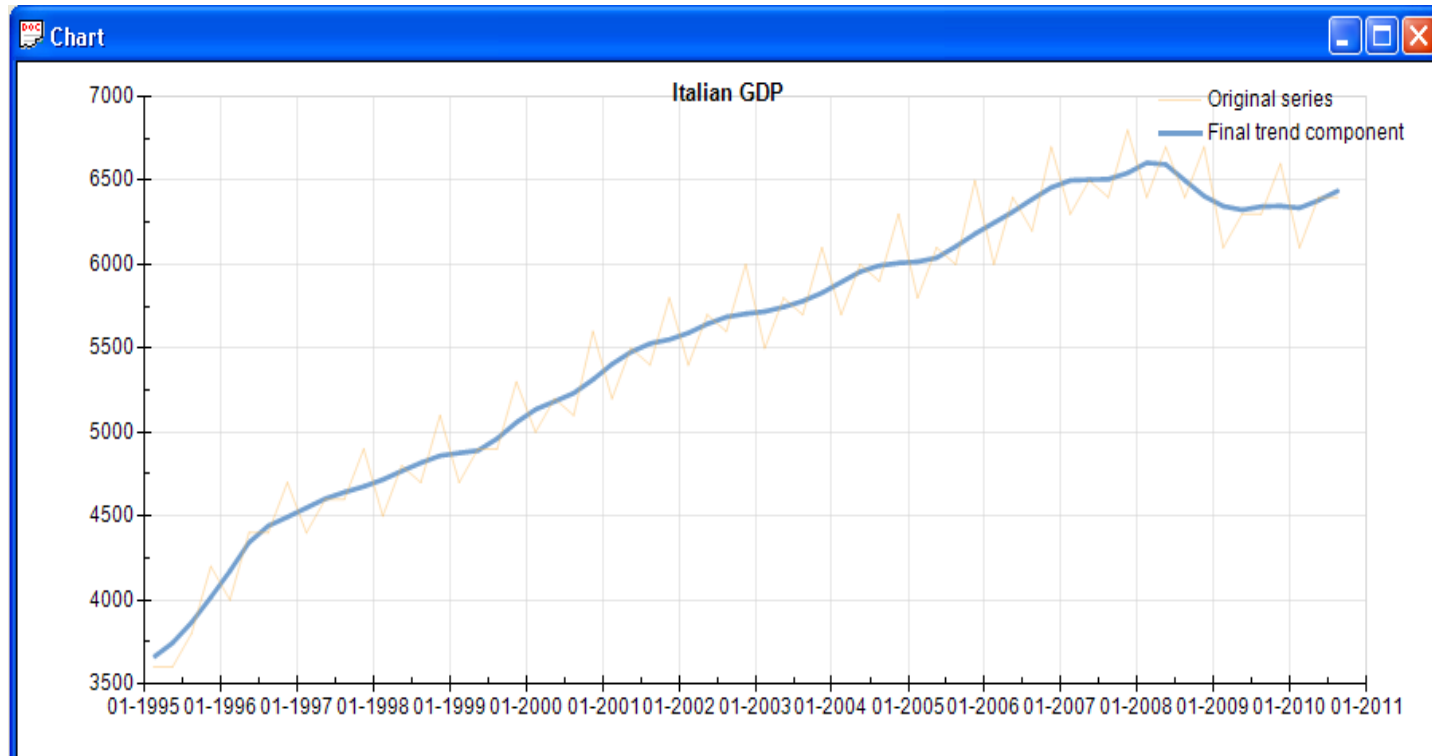
Usual Components

- The Trend Component
 - The Trend is the long term evolution of the series that can be observed on several decades
- The Cycle Component
 - The Cycle is the smooth and quasi-periodic movement of the series that can usually be observed around the long term trend
- The Seasonal Component (Seasonality)
 - Fluctuations observed during the year (each month, each quarter) and which appear to repeat themselves on a more or less regular basis from one year to other

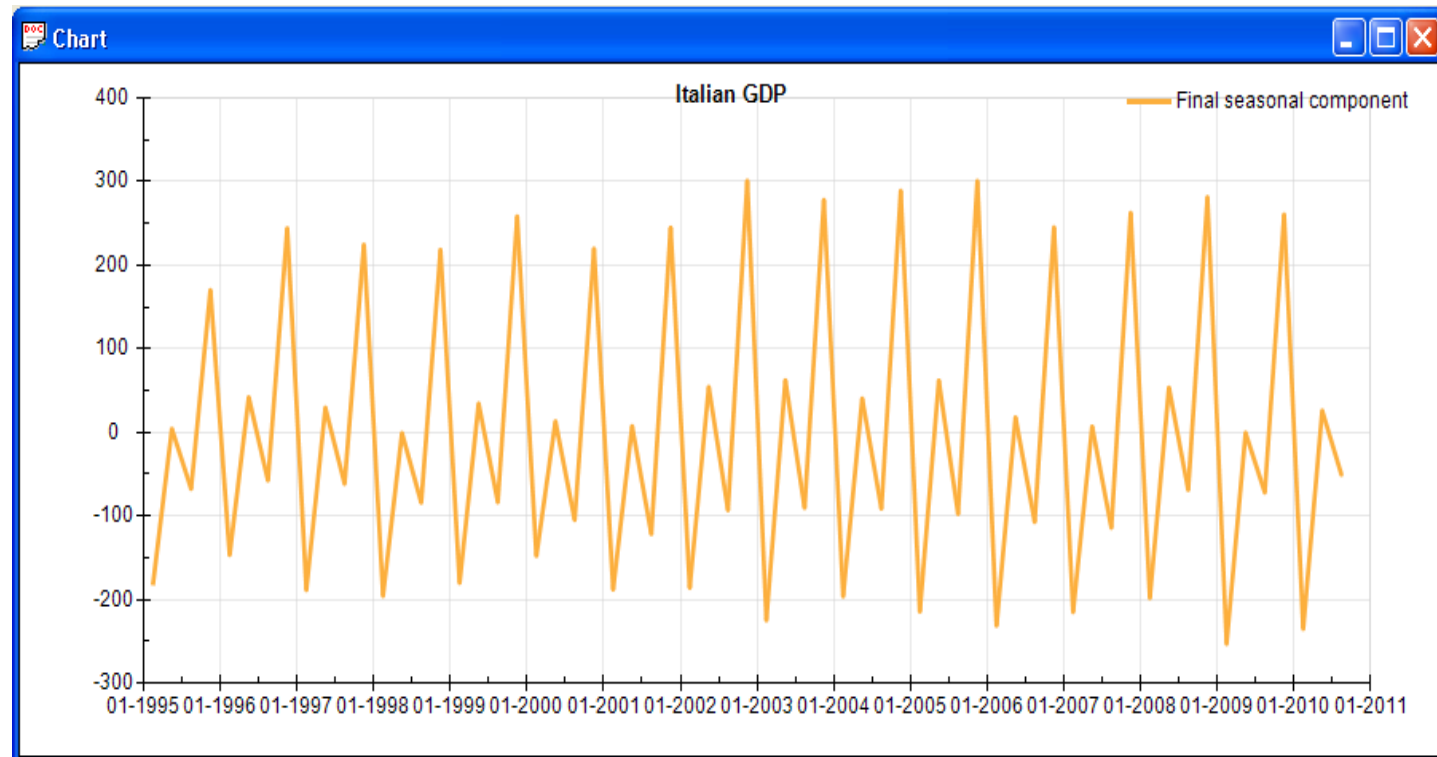
Usual Components

- The Calendar Effect
 - Any economic effect which appears to be related to the calendar (e.g. one more Sunday in the month can affect the production)
- The Irregular Component
 - The Irregular Component is composed of residual and random fluctuations that cannot be attributed to the other “systematic” components
- Outliers
 - Different kinds of Outliers can be defined
 - year to other

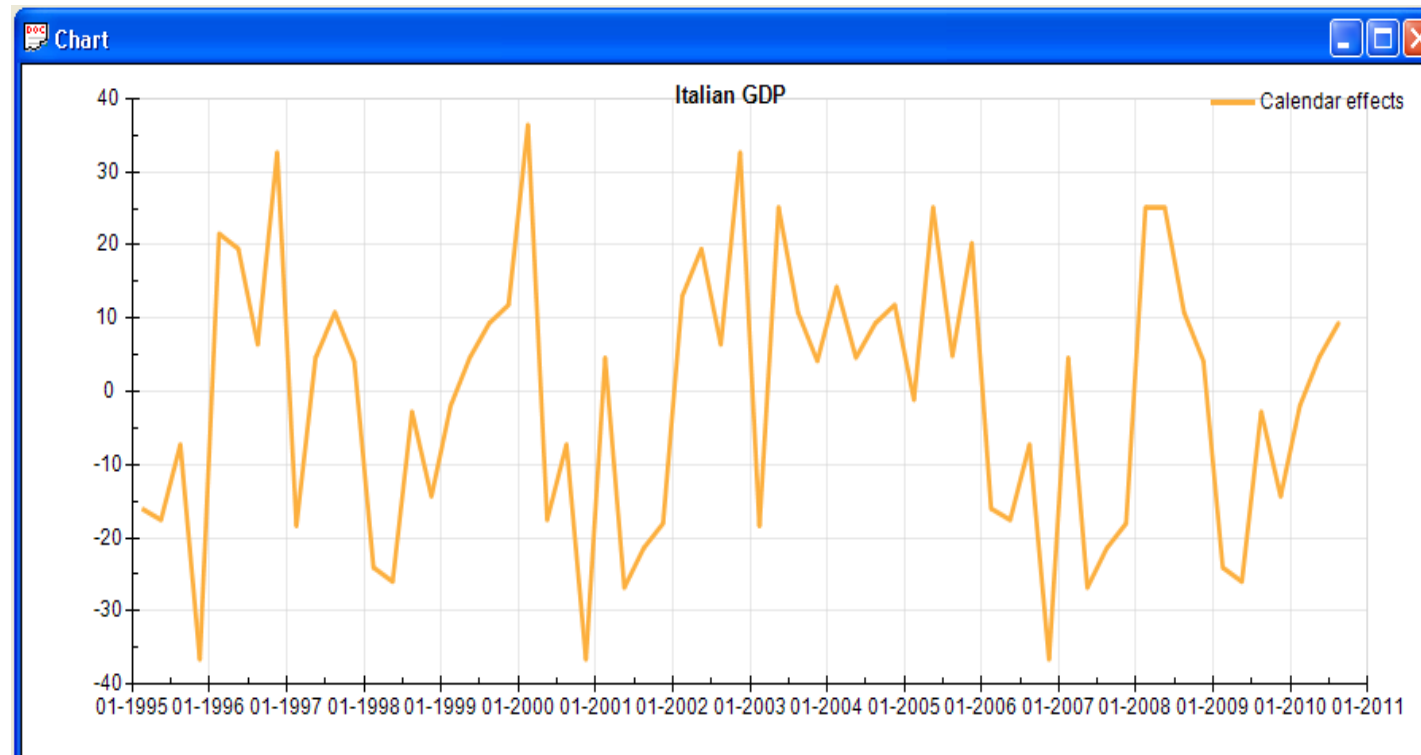
Italian GDP – Trend component



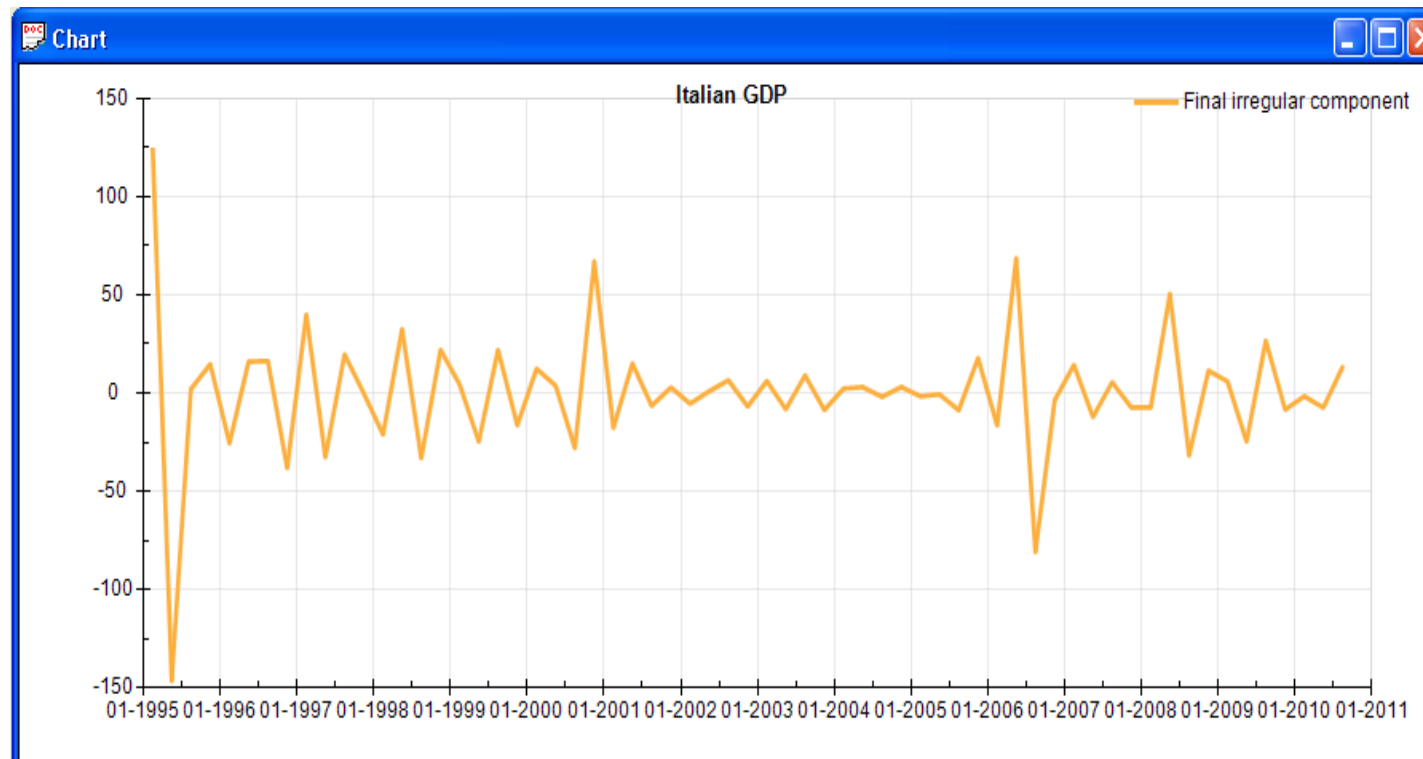
Italian GDP – Seasonal component



Italian GDP – Calendar component



Italian GDP – Irregular component



Italian GDP – Irregular, Seasonal and Calendar



Trend

The Trend Component is defined as the **long-term movement** in a series

The Trend is a reflection of the underlying level of the series. This is typically due to influences such as population growth, price inflation and general economic development

The Trend Component is sometimes referred to as the Trend-Cycle (see Cycle Component)

Cause of Seasonality

Seasonality and Climate: due to the variations of the weather and of the climate (seasons!)

- Examples: agriculture, consumption of electricity (heating)

Seasonality and Institutions: due to the social habits and practices or to the administrative rules

- Examples: effect of Christmas on the retail trade, of the fiscal year on some financial variables, of the academic calendar

Indirect Seasonality: due to the Seasonality that affects other sectors

- Examples: toy industry is affected a long time before Christmas. A Seasonal increase in the retail trade has an impact on manufacturing, deliveries, etc..

Seasonal Adjustment

Seasonal Adjustment is the process of estimating and removing the Seasonal Effects from a Time Series, and by Seasonal we mean an effect that happens at the same time and with the same magnitude and direction every year

The basic goal of Seasonal Adjustment is to decompose a Time Series into several different components, including a Seasonal Component and an Irregular Component

Because the Seasonal effects are an unwanted feature of the Time Series, Seasonal Adjustment can be thought of as focused noise reduction

Seasonal Adjustment

Since Seasonal effects are annual effects, the data must be collected at a frequency less than annually, usually monthly or quarterly

For the data to be useful for Time Series analysis, the data should be comparable over time. This means:

- The measurements should be taken over discrete (nonoverlapping) consecutive periods, i.e., every month or every quarter
- The definition of the concept and the way it is measured should be consistent over time

Seasonal Adjustment

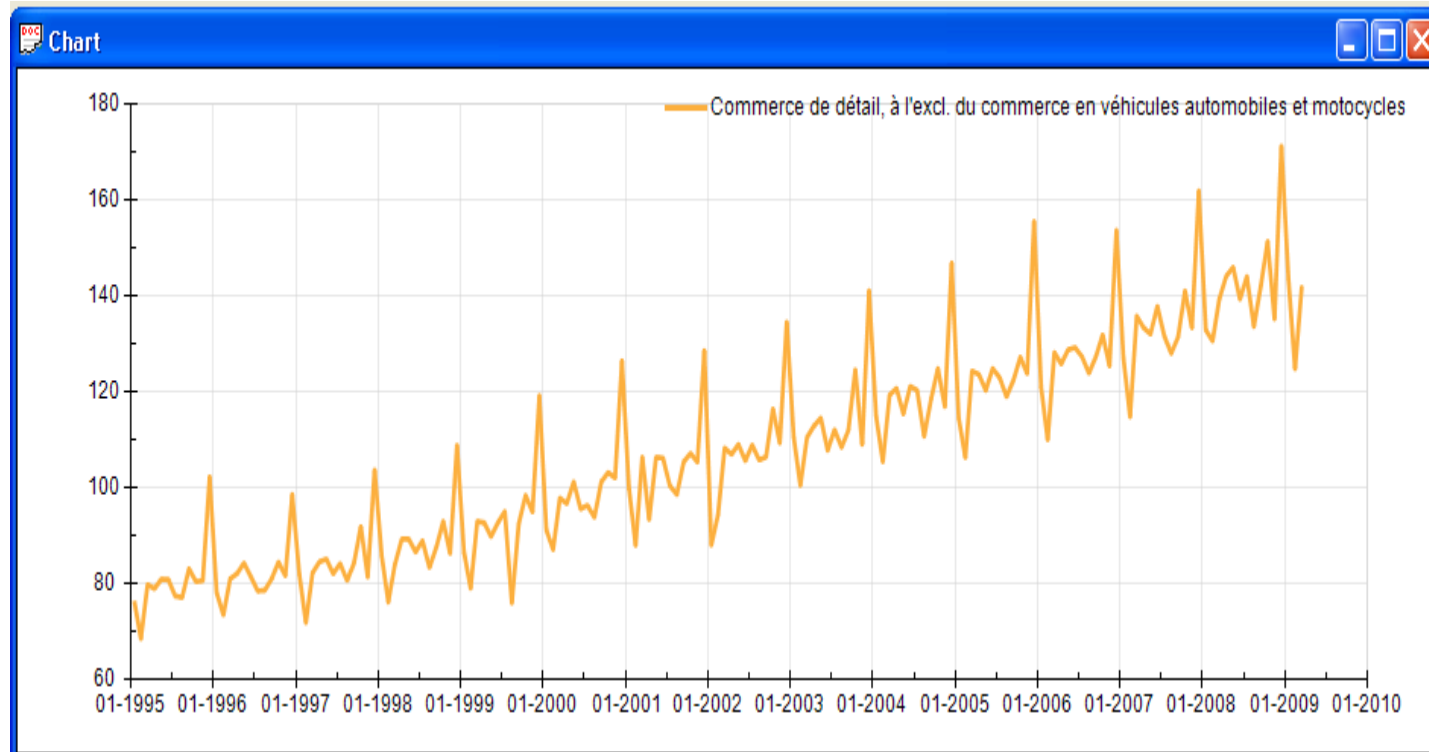
- Keep in mind that longer series are NOT necessarily better. If the series has changed the way the data is measured or defined, it might be better to cut off the early part of the series to keep the series as homogeneous as possible
- The best way to decide if your series needs to be shortened is to investigate the data collection methods and the economic factors associated with your series and choose a length that gives you the most homogeneous series possible

Types of Stationarity

- **Strict stationarity** - This means that the unconditional joint distribution of any moments (e.g. expected values, variances, third-order and higher moments) remains constant over time. This type of series is rarely seen in real-life practice.
- **First-order stationarity** - These series have a mean constant over time. Any other statistics (like variance) can change at the different points in time.
- **Second-order stationarity (also called weak stationarity)** - These time series have a constant mean and variance over time. Other statistics in the system are free to change over time. This is one of the most commonly observed series in real-life practice.
- **Trend stationarity** - These series are the series with a trend. This trend when removed from the series leaves a stationary series.
- **Difference stationarity** - These series are the series that need one or more differencing to become stationary.

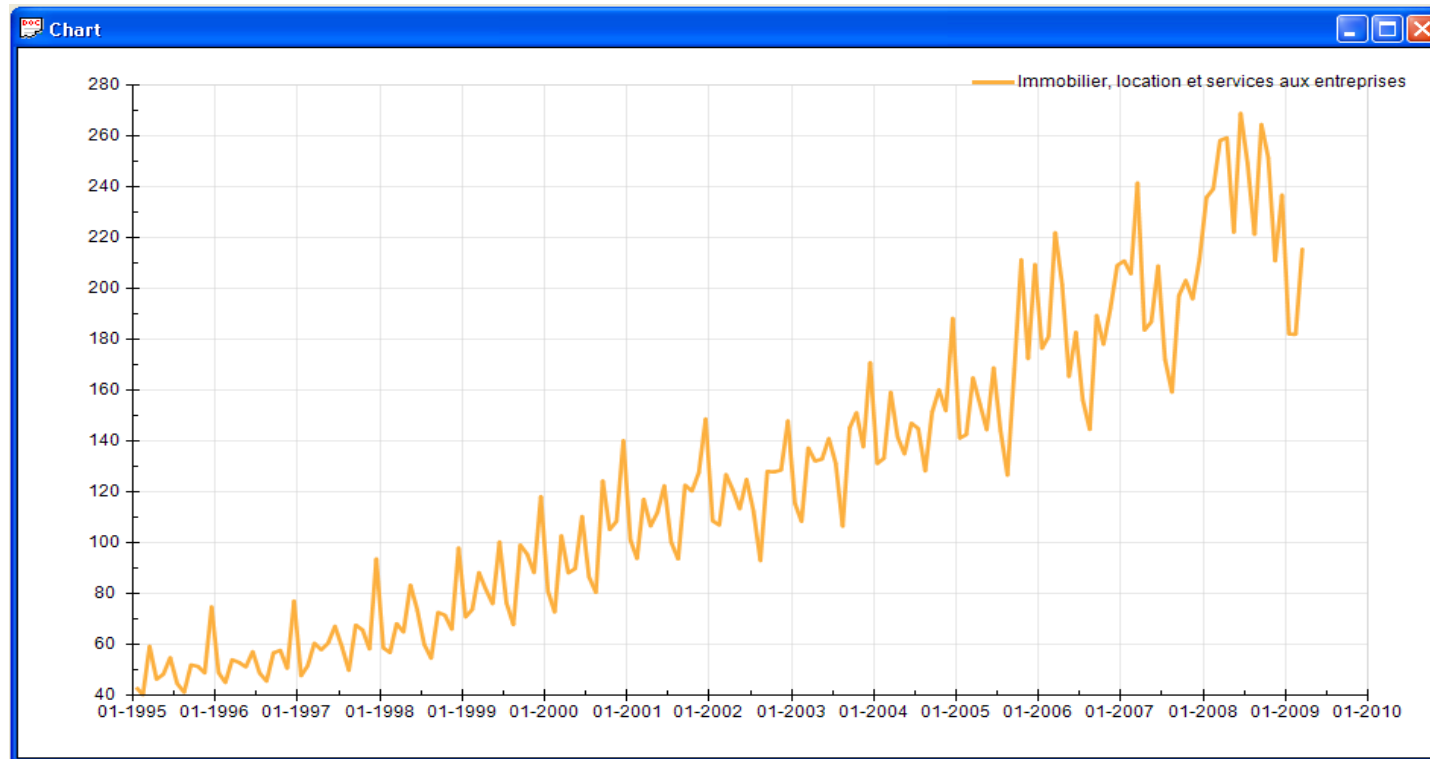
Seasonal Adjustment

A first overview – No Stationary in mean (example)



Seasonal Adjustment

A first overview – No Stationary in variance (example)



Differencing

One way to make a time series stationary is to compute the differences between consecutive observations.

This is known as **DIFFERENCING**

Differencing can help **stabilize the mean** of a time series by removing changes in the level of a time series, and so eliminating trend (and also seasonality, using a specific differencing order)

The **Order of Integration** for a Time Series, denoted $I(d)$, reports the minimum number of differences (d) required to obtain a stationary series (*note: $I(0) \rightarrow$ it means the series is stationary!*)

**Differenced
Time Series (first order)**

y_t $y'_t = y_t - y_{t-1}$

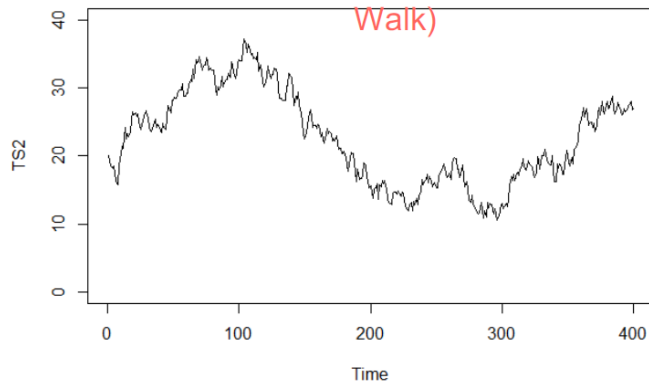
↓ ↓

CYCLE_	WEEK_	DATE_	COLLI_ARR	DIFF_1
1	1	1 1	983	.
1	2	1 2	1478	495
1	3	1 3	1822	345
1	4	1 4	1883	61
1	5	1 5	1913	30
1	6	1 6	2001	88
1	7	1 7	2077	76

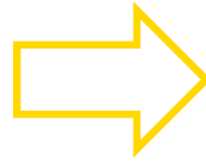
← 1478 - 983 = 495

Differencing

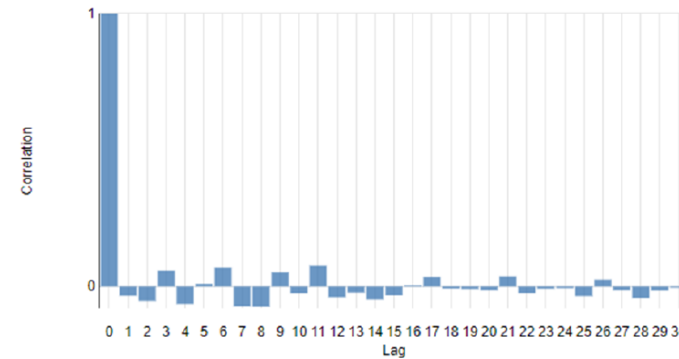
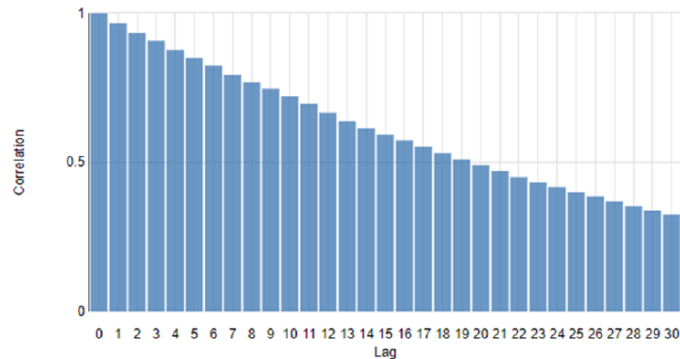
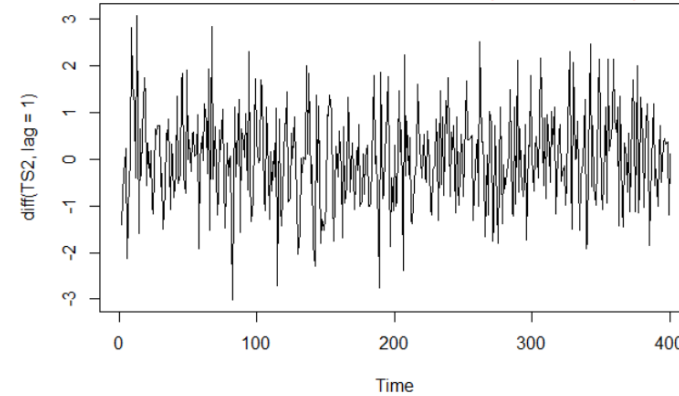
Non-Stationary Time Series example 1 (Random Walk)



$$TS2_t - TS2_{t-1}$$



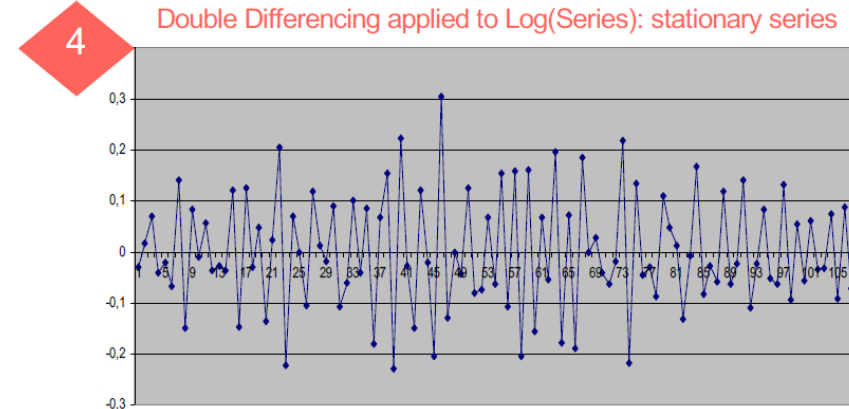
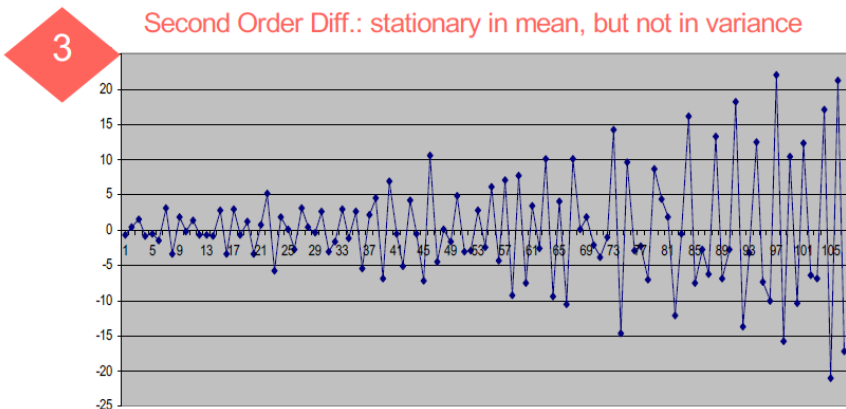
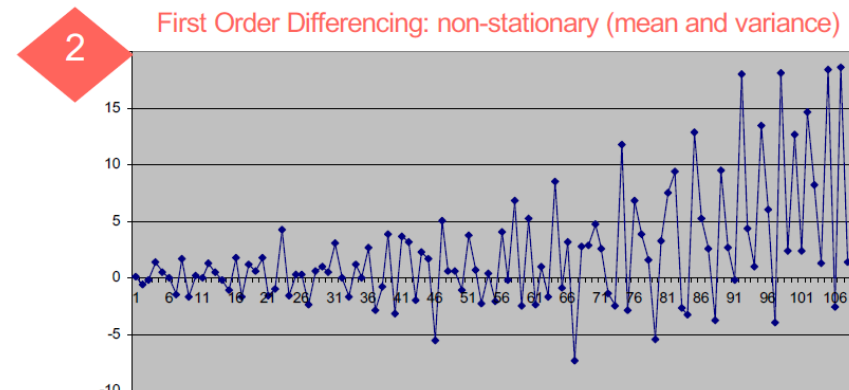
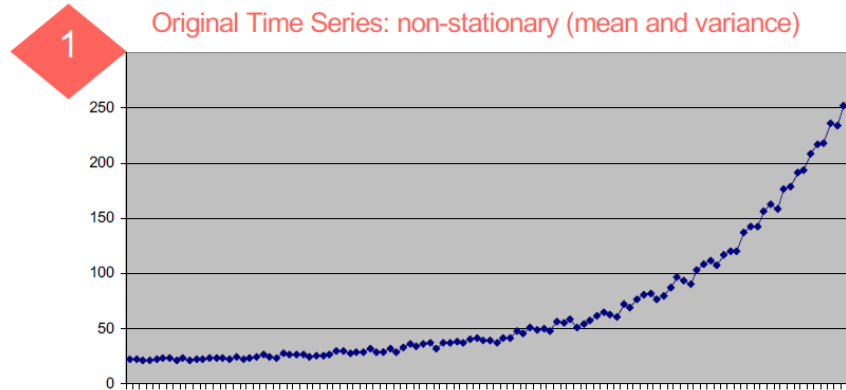
Differenced Time Series (first order)



No significant autocorrelation exists → applying first differences to a random walk generates a white noise

Example: use differencing to make stationary a non-stationary series

Differencing



* it's almost never necessary to go beyond second-order differences

Occasionally the differenced data will not appear stationary and it may be necessary to difference the data a second time to obtain a stationary series

Decomposition Models

Usual Additive and Multiplicative Models

$$X_t = T_t + C_t + S_t + I_t$$

$$X_t = T_t * C_t * S_t * I_t$$

$$X_t = T_t * (1 + C_t) * (1 + S_t) * (1 + I_t)$$

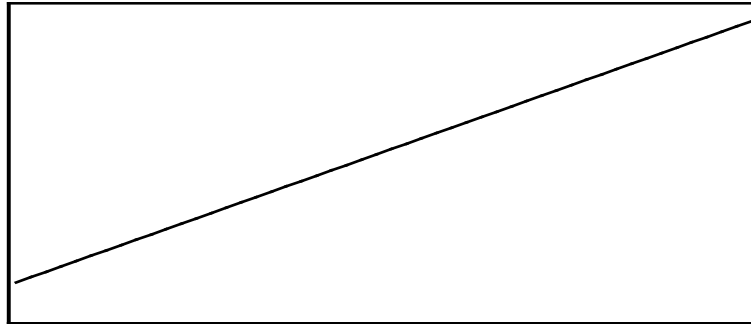
More components: Outliers, Calendar Effects

$$X_t = T_t + C_t + S_t + O_t + TD_t + MH_t + I_t$$

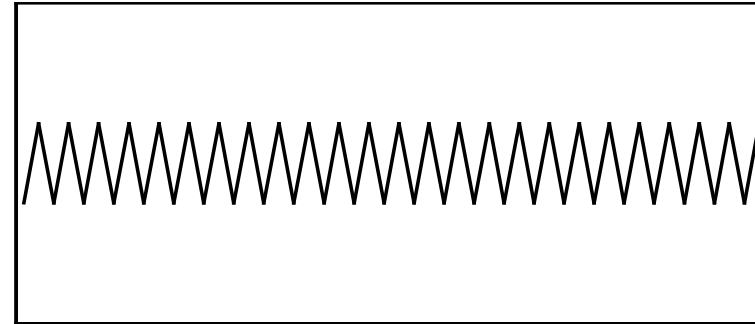
Decomposition Models

Some usual shapes

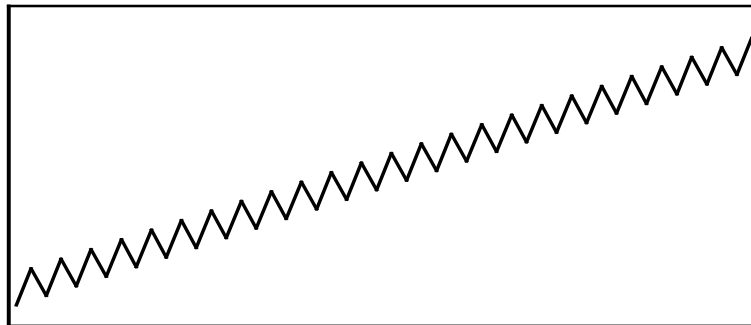
Trend



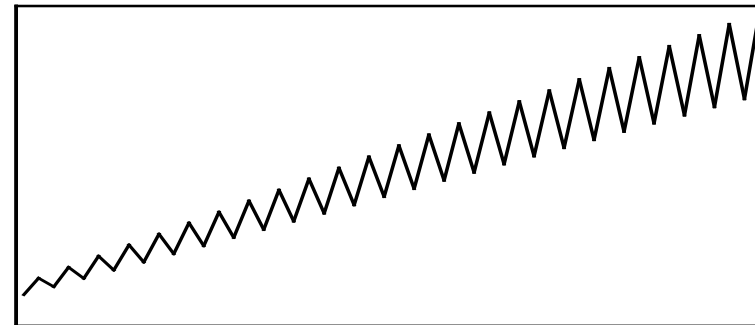
Seasonality



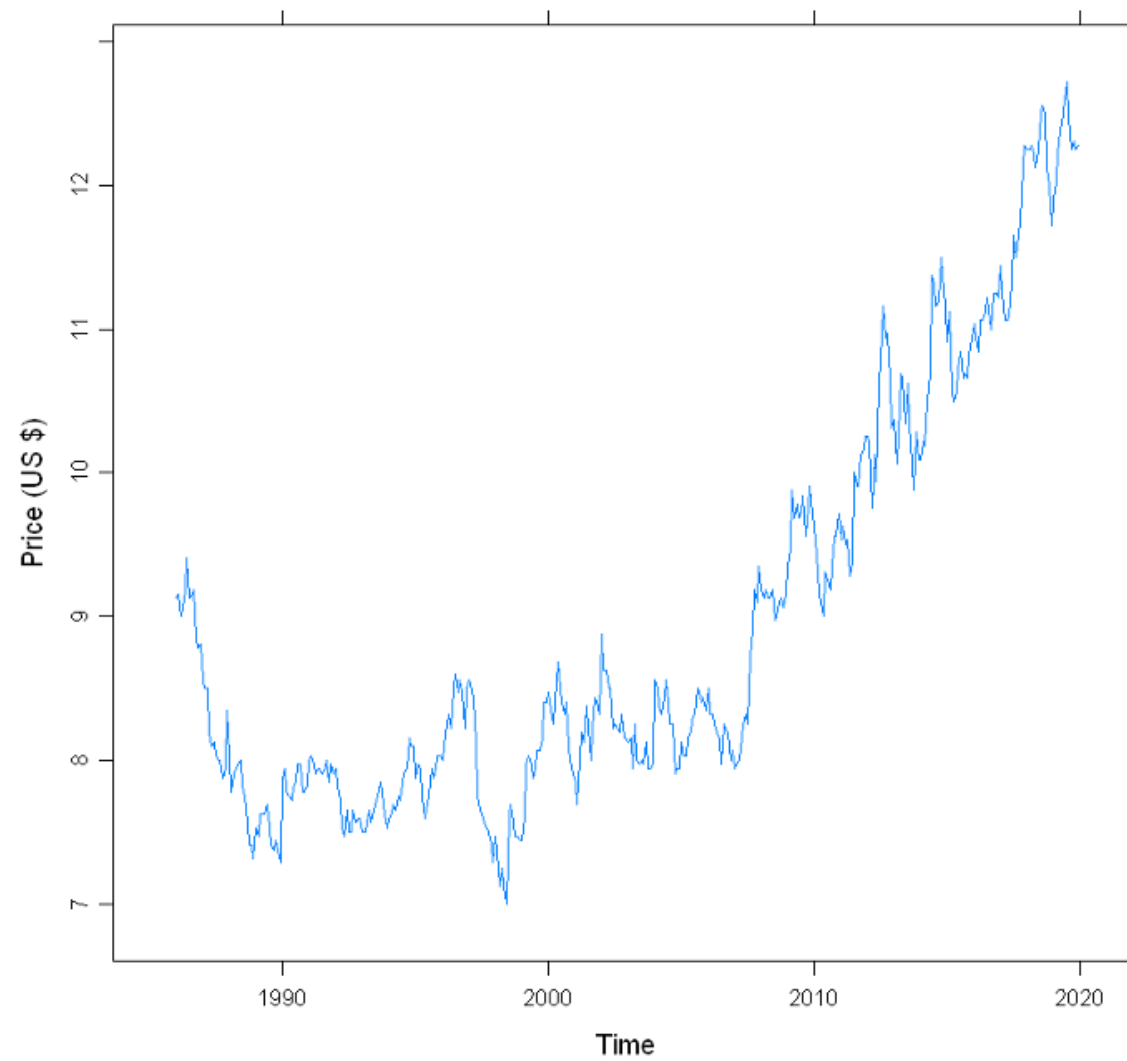
Additive Model



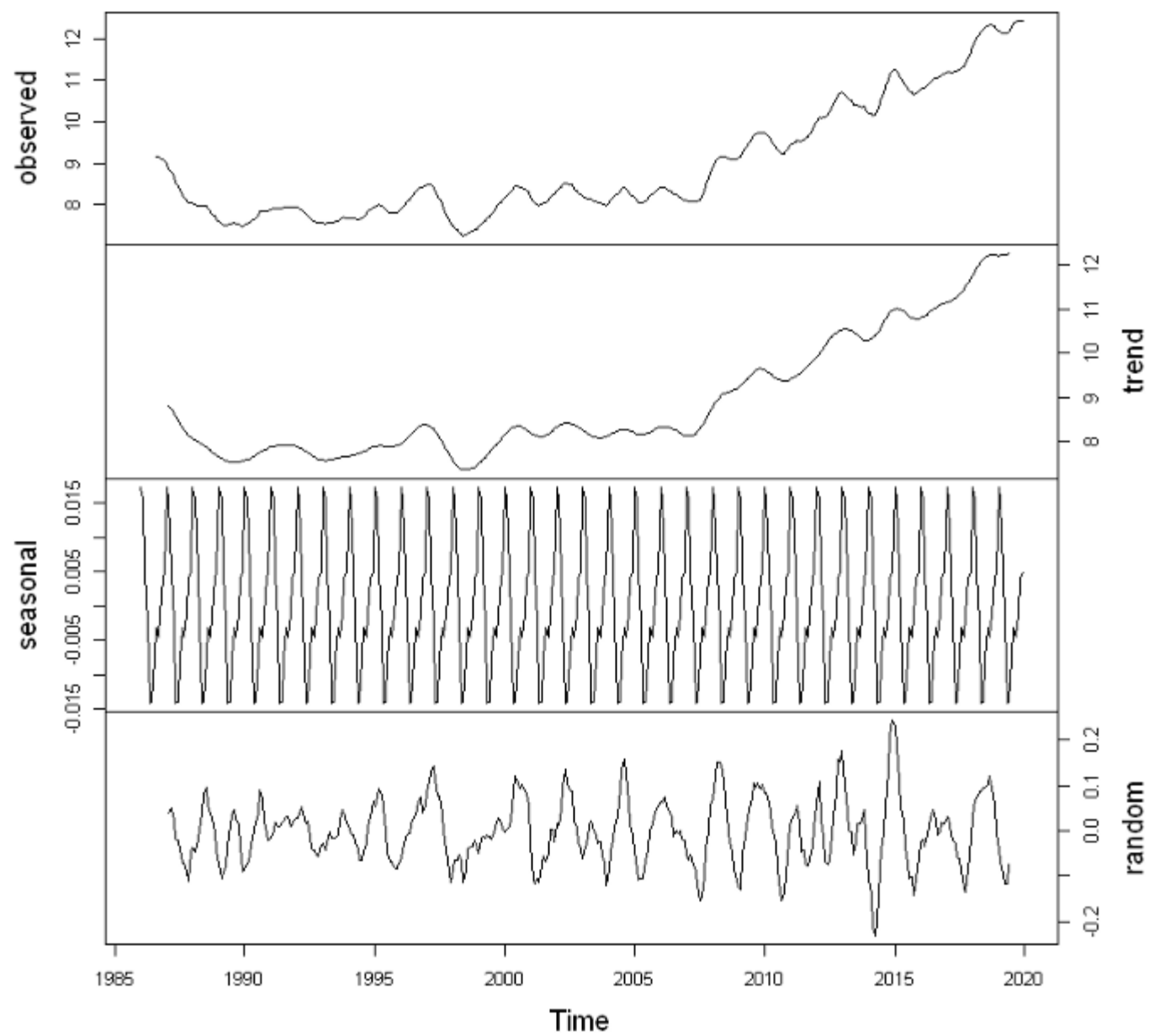
Multiplicative Model



Time series plot for Schlumberger price



Decomposition of additive time series



Calendar Adjustment

Calendar Effects typically include:

- Different number of Working Days in a specific period
- Composition of Working Days
- Leap Year effect
- Moving Holidays (Easter, Ramadan, etc.)

Calendar Adjustment - Trading Day Effect

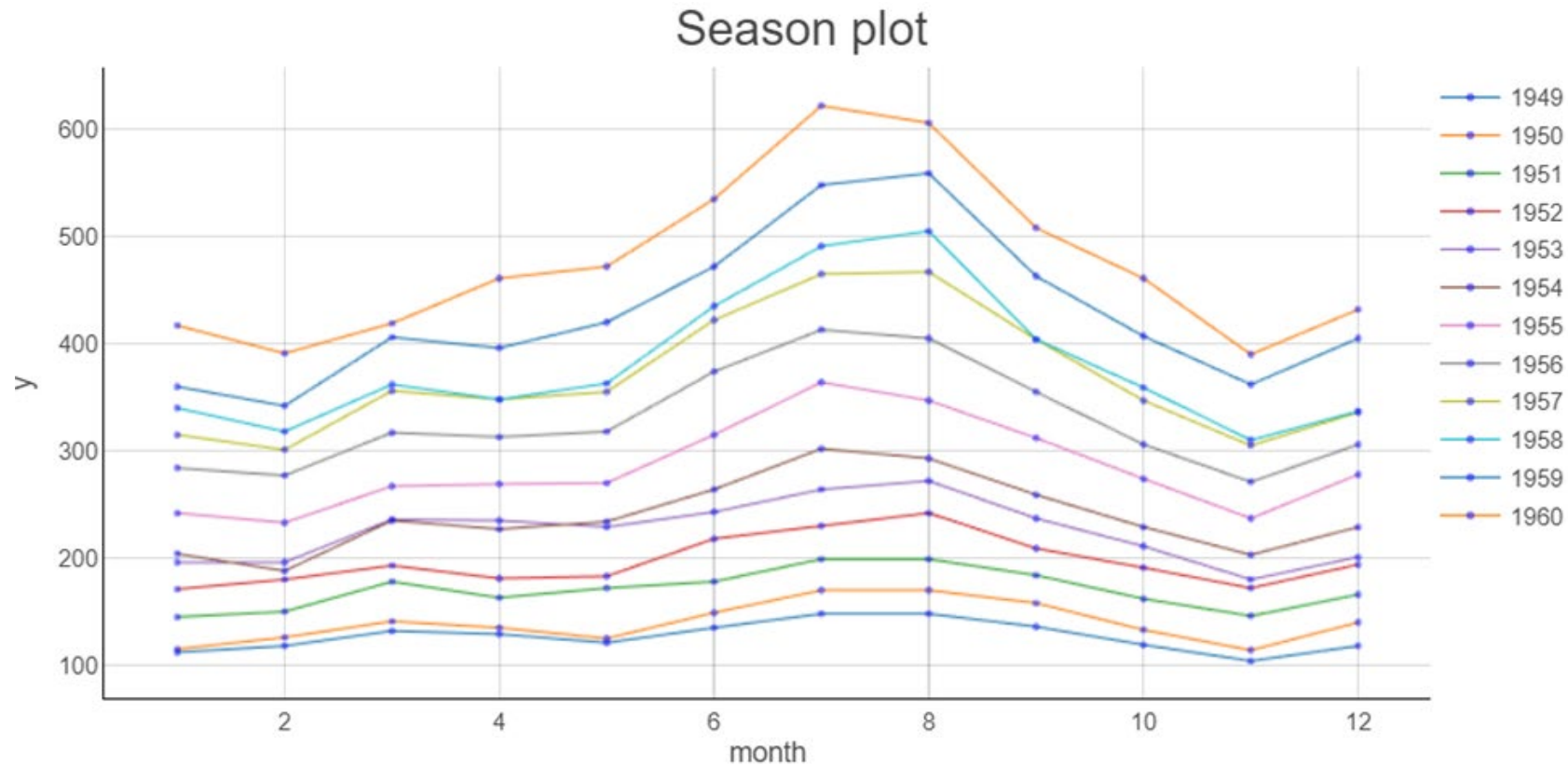
- Recurring effects associated with individual days of the week. This occurs because only non-leap-year Februarys have four of each type of day: four Mondays, four Tuesdays, etc.
- All other **months have an excess of some types of days**. If an activity is higher on some days compared to others, then the series can have a Trading Day effect. For example, building permit offices are usually closed on Saturday and Sunday
- Thus, the number of building permits issued in a given month is likely to be higher if the month contains a surplus of weekdays and lower if the month contains a surplus of weekend days

Calendar Adjustment - Moving Holiday Effect

Effects from **holidays that are not always on the same day of a month**, such as Labor Day or Thanksgiving.

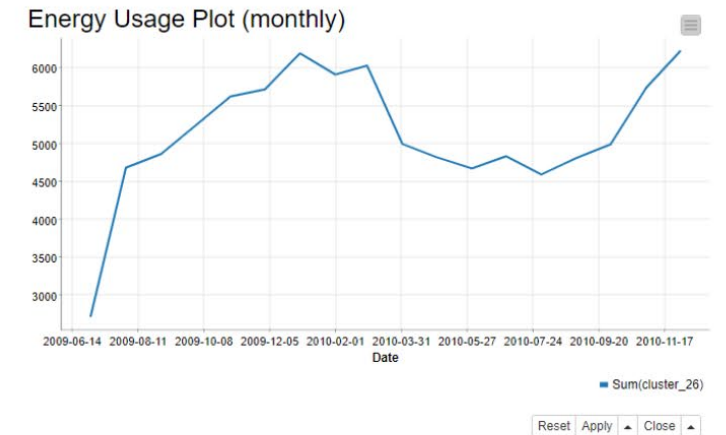
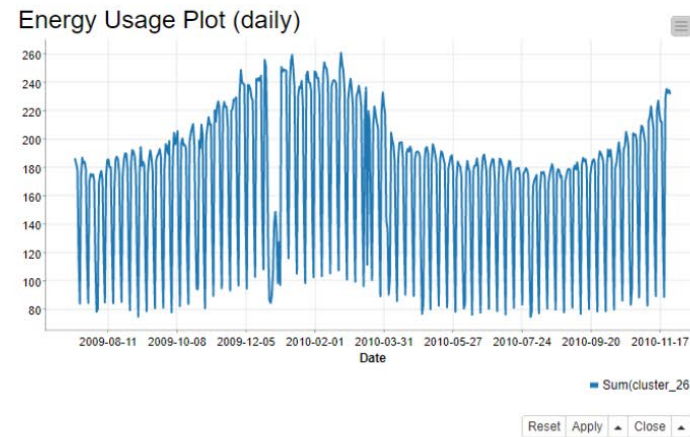
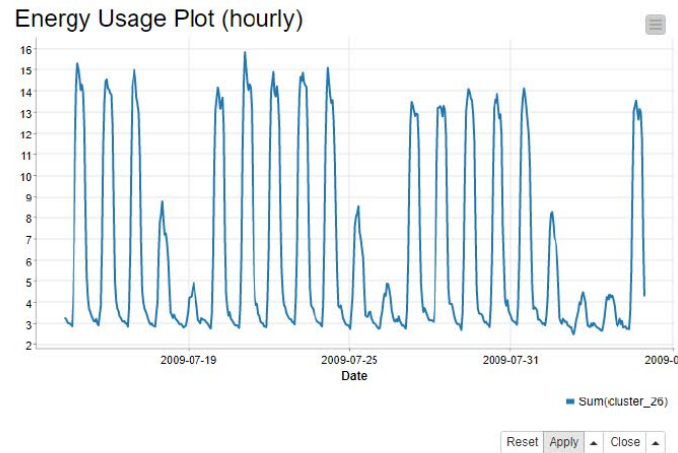
The most important Moving Holiday in the US and European countries is Easter, not only because it moves between days, but also because it moves between months since it can occur in March or April

Graphical Analysis: Seasonal Plot



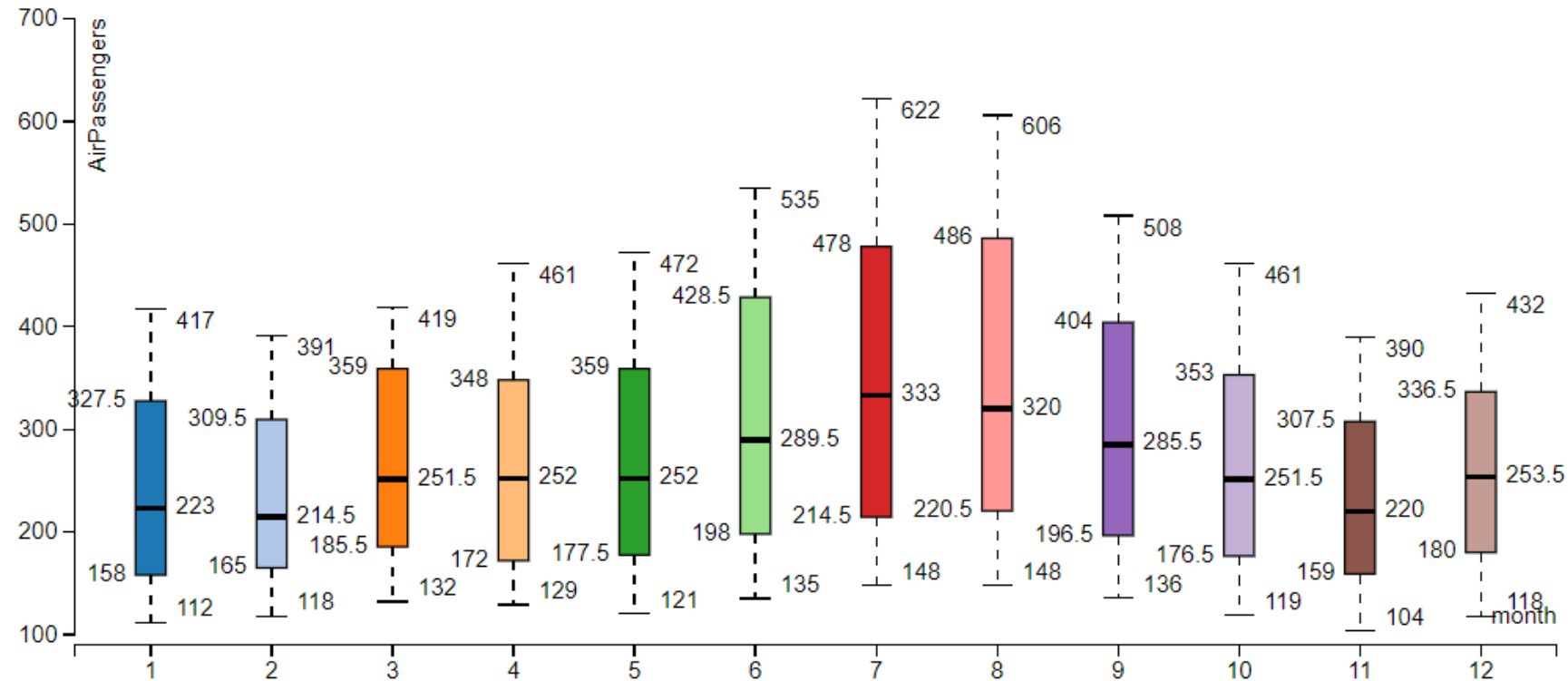
Produce the **Seasonal plot** of the Time series in order to analyze more in detail the seasonal component (and possible changes in seasonality over time)

Granularity and Line Plots



Show time series by hour, day, and month in line plots
Identify daily, weekly, and yearly seasonality

Graphical Analysis: Box Plot



Create the **conditional Box plot** of the Time series in order to deeply understand the distribution of data in the same period of each seasons and focusing on specific aspects such as outliers, skewness, variability,...

Irregular Component

The Irregular Component is the remaining component of the series **after the Seasonal and Trend Components have been removed** from the original data

For this reason, it is also sometimes referred to as the Residual Component. It attempts to capture the remaining short term fluctuations in the series which are neither systematic nor predictable

Exploration – Basic tools

Exploration is a very essential step when analyzing a Time Series

Looking for “structures” in the series

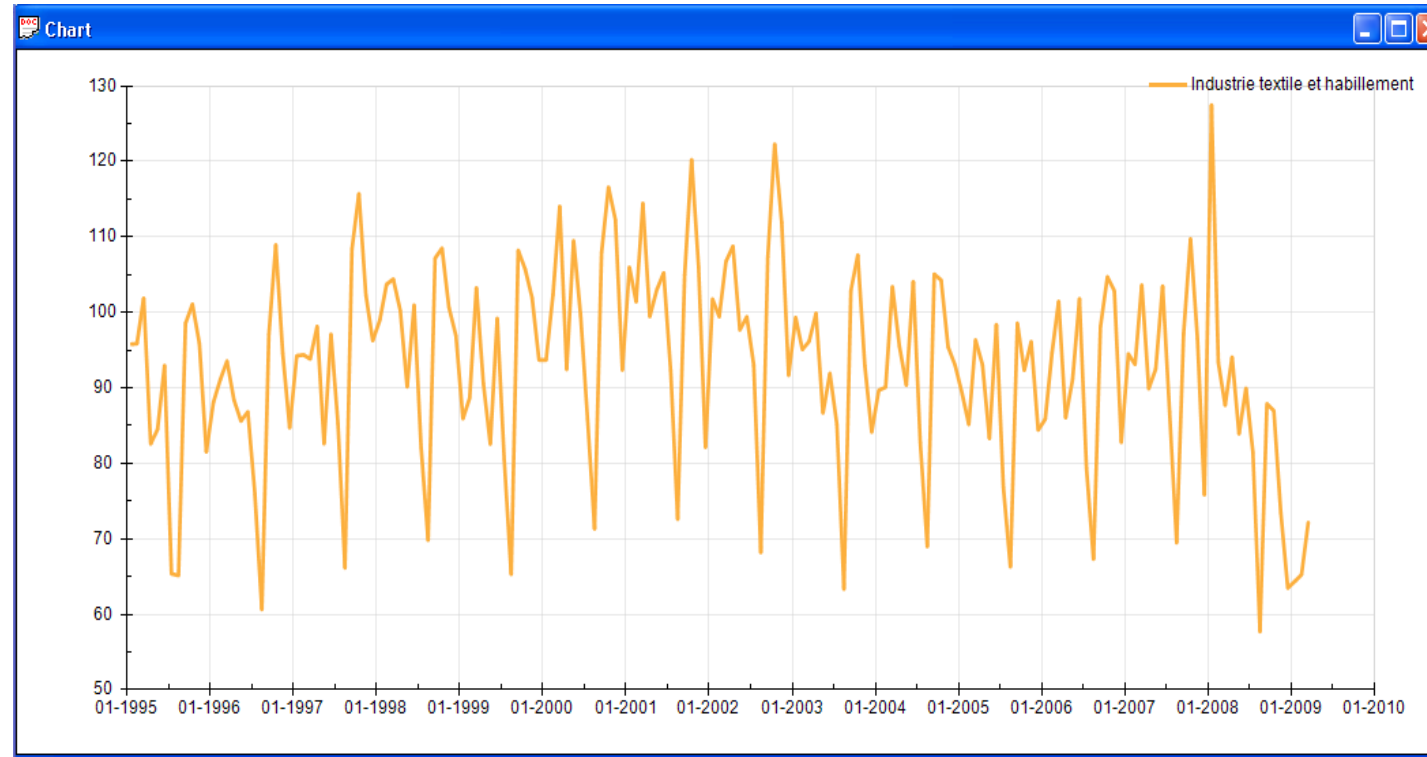
- Trend, Seasonality, “strange” points or behavior etc.

Helps to formulate a global or decomposition model for the series

Graphics are a key player in this exploration

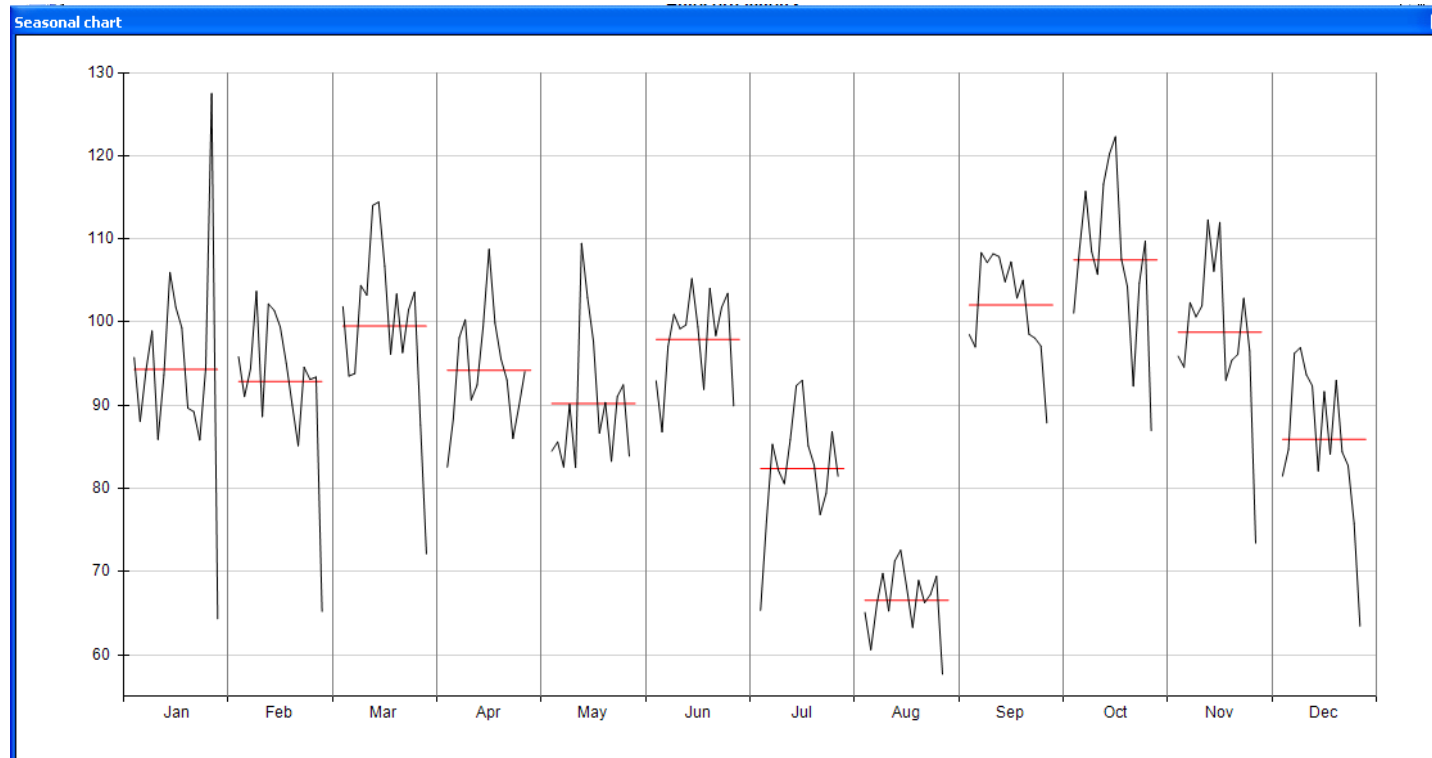
Exploration – Usual representation

Textile industry and clothing



Exploration – Seasonal Chart

Textile industry and clothing





Why Seasonal Adjustment?

Business cycle analysis

To improve comparability:

Over time:

- Example: how to compare the first quarter (with February) to the fourth quarter (with Christmas)?

Across space:

- Never forget that while we are freezing at work, Australians are burning on the beach!
- Very important to compare national economies (convergence of business cycles) or sectors

Why Seasonal Adjustment?

The aim of Seasonal Adjustment is to eliminate Seasonal and Calendar Effects. Hence **there are no Seasonal and Calendar Effects in a perfectly Seasonally Adjusted series**

In other words: Seasonal Adjustment transforms the world we live in into a world where no Seasonal and Calendar Effects occur. In a Seasonally Adjusted world the temperature is exactly the same in winter as in the summer, there are no holidays, Christmas is abolished, people work every day in the week with the same intensity (no break over the weekend), etc.

Outliers

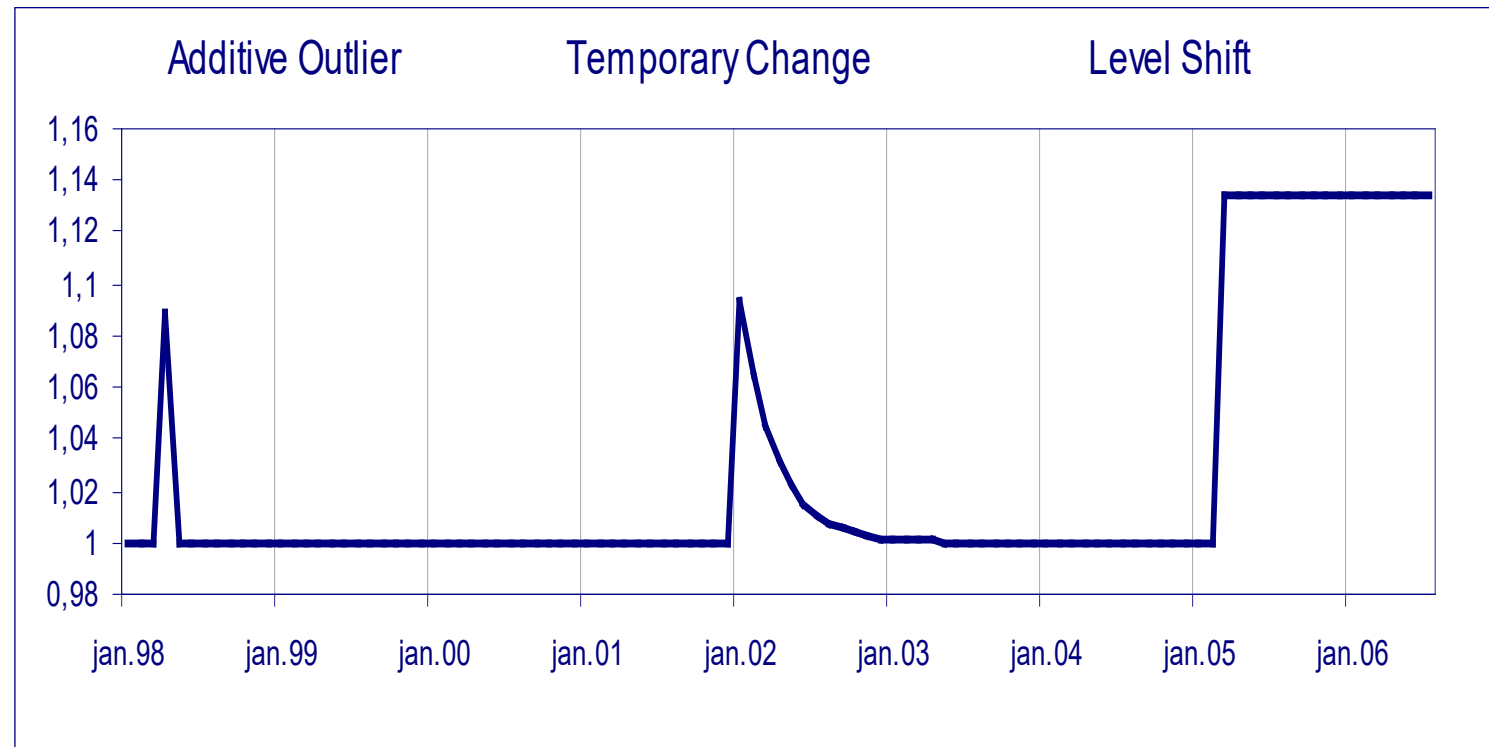
Outliers are data which do not fit in the tendency of the Time Series observed, which fall outside the range expected on the basis of the typical pattern of the Trend and Seasonal Components

Additive Outlier (AO): the value of only one observation is affected. AO may either be caused by random effects or due to an identifiable cause as a strike, bad weather or war

Temporary Change (TC): the value of one observation is extremely high or low, then the size of the deviation reduces gradually (exponentially) in the course of the subsequent observations until the Time Series returns to the initial level

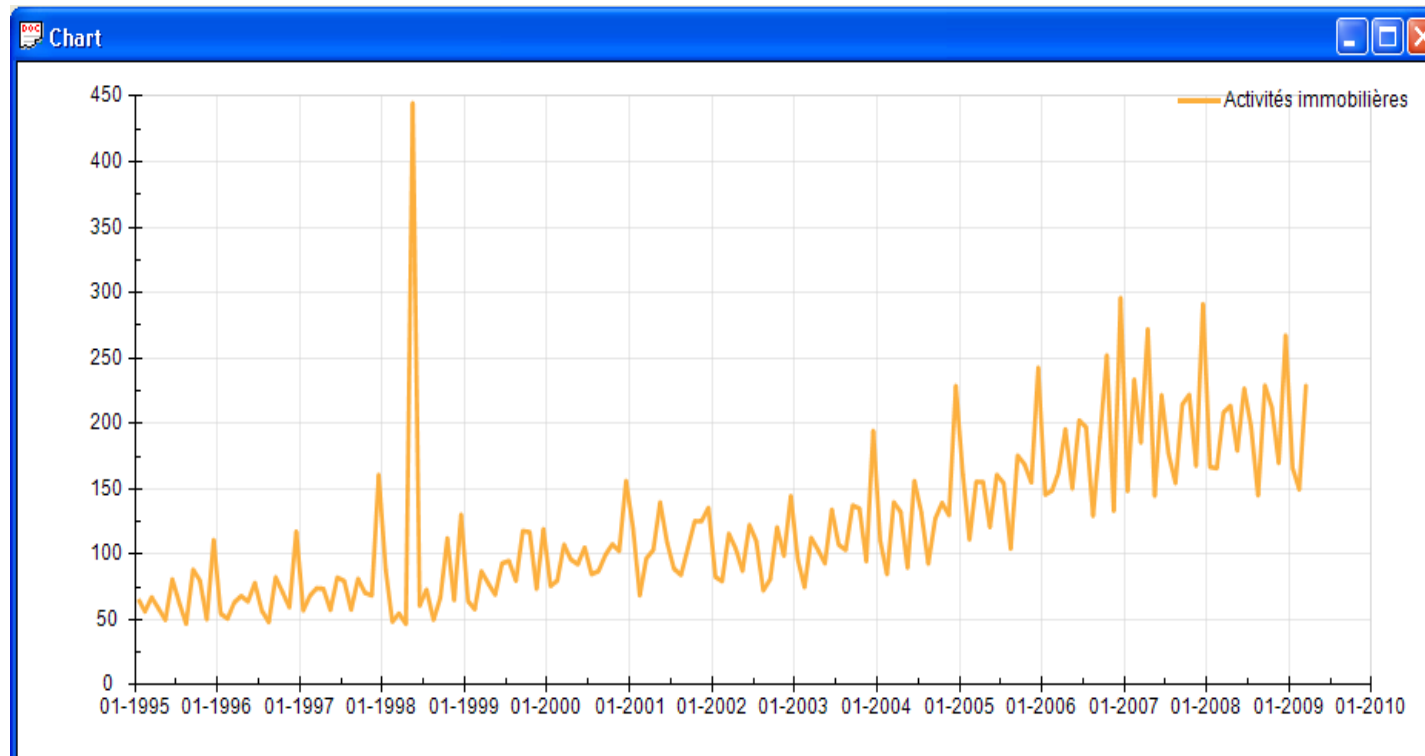
Level Shift (LS): starting from a given time period, the level of the Time Series undergoes a permanent change. Causes could include: change in concepts and definitions of the survey population, in the collection method, in the economic behavior, in the legislation or in the social traditions

Outliers



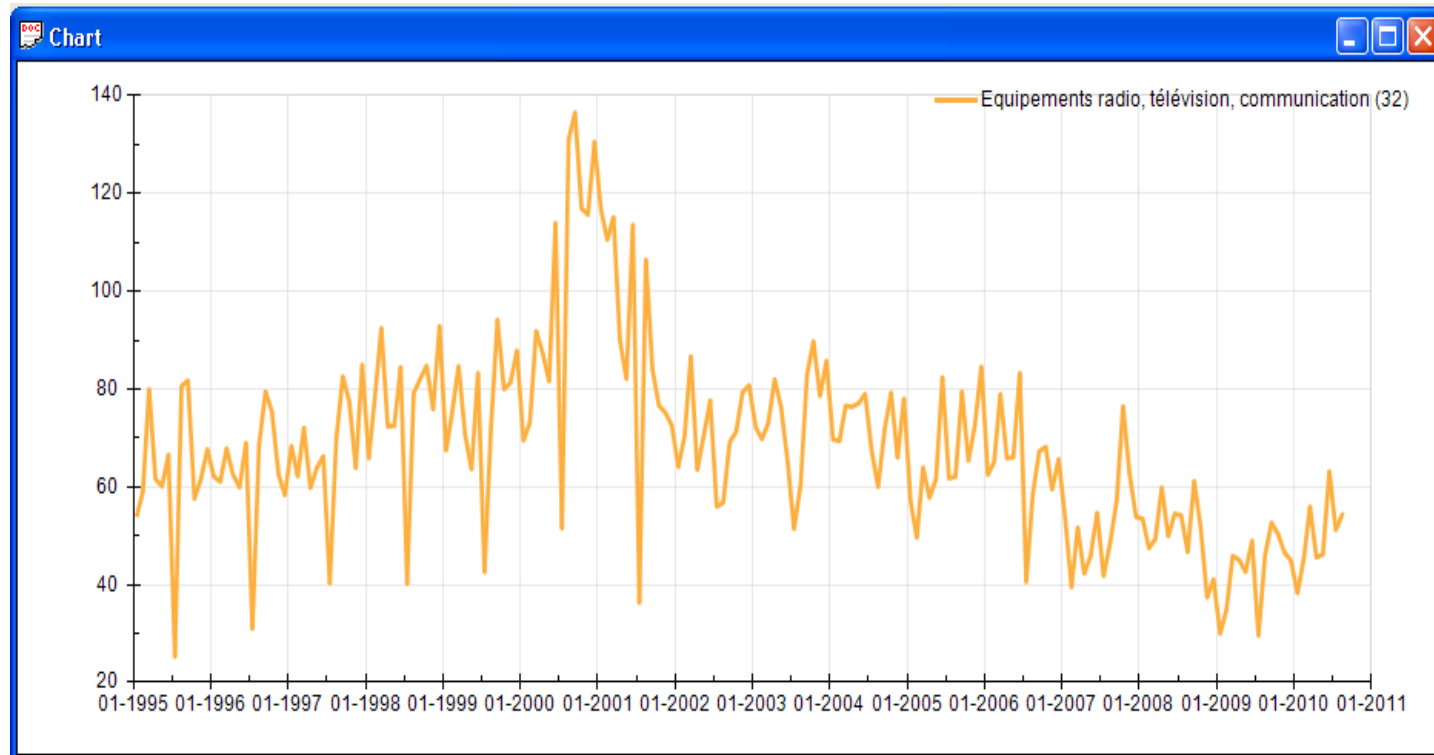
Outliers

AO - Estate agency activity



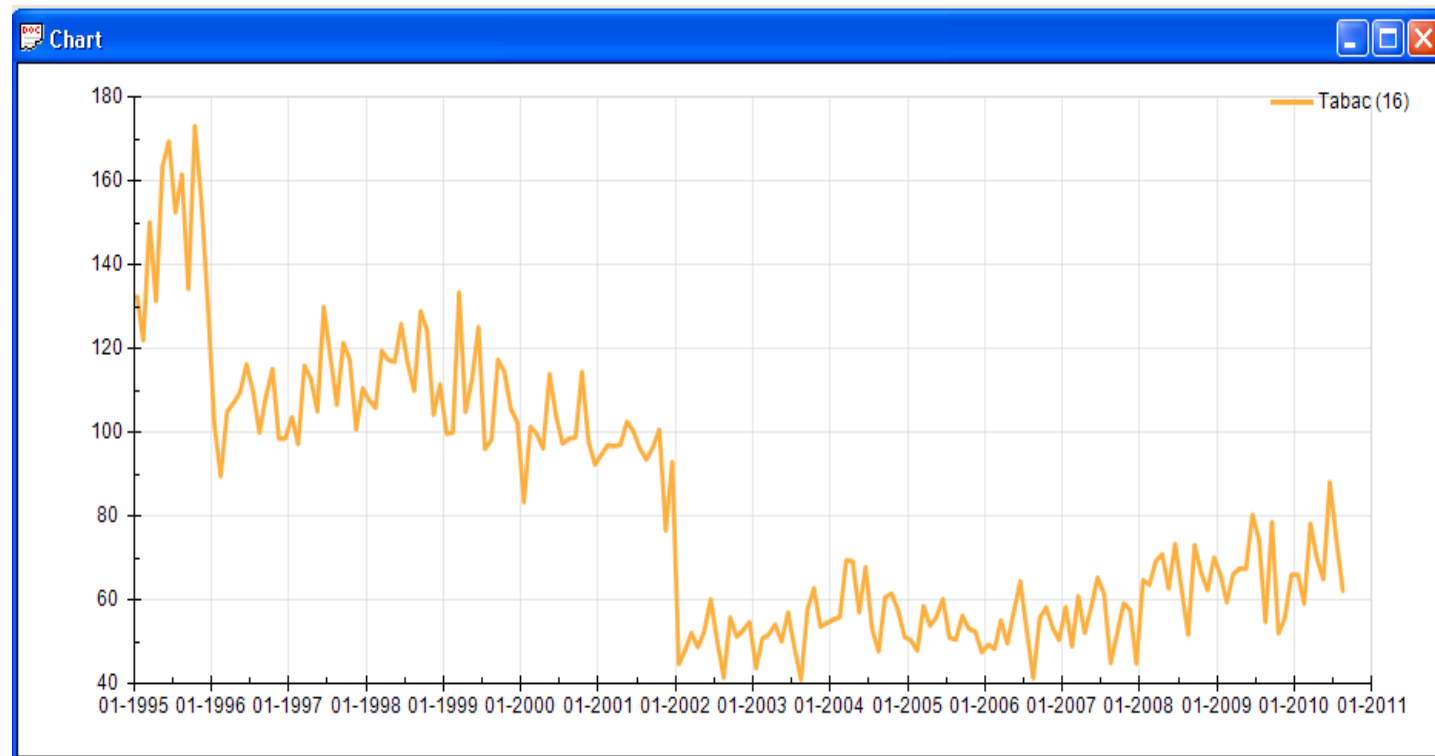
Outliers

TC – Business machine

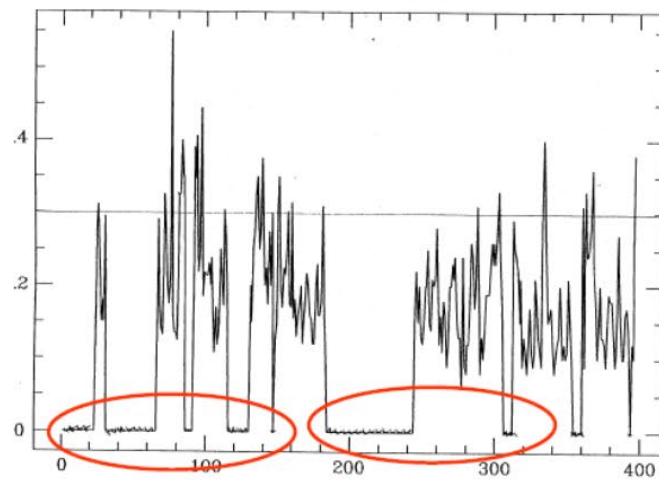


Outliers

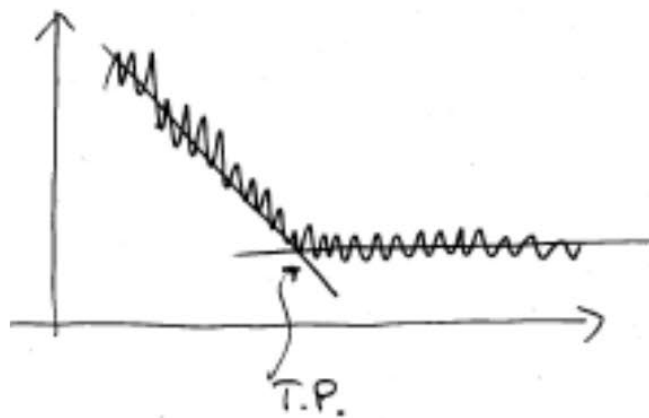
LS – Tobacco



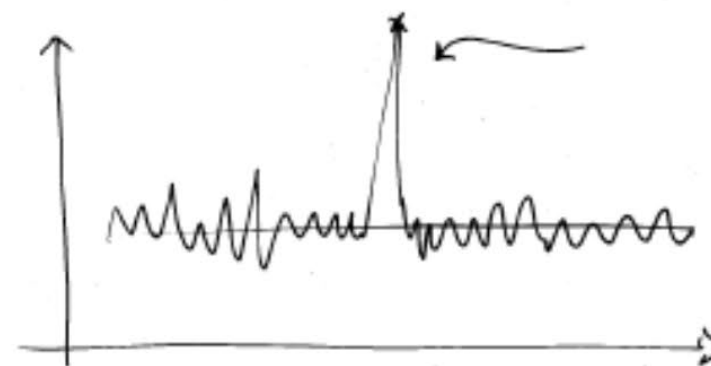
Series with gaps



Series with a turning point



Series with an outlier



Outliers

Additive Outliers:

Unusual high or low singular values in the data series

Transitory Changes:

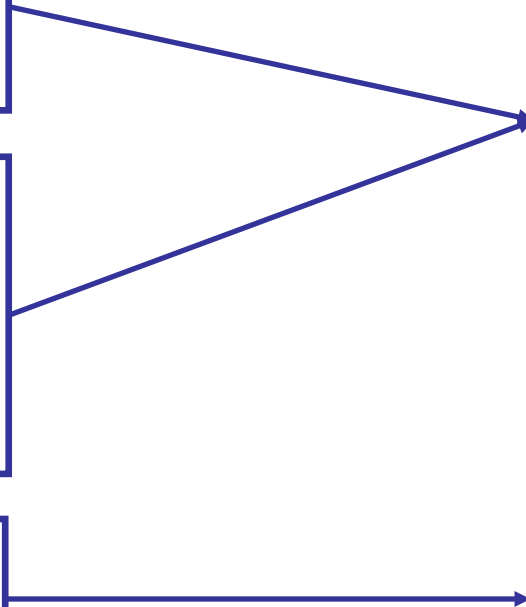
Transitory changes in the trend, followed by slow comebacks to the initial tendency

Level Shift:

Clear changes of the trend

**Assimilated to
the Irregular
Component**

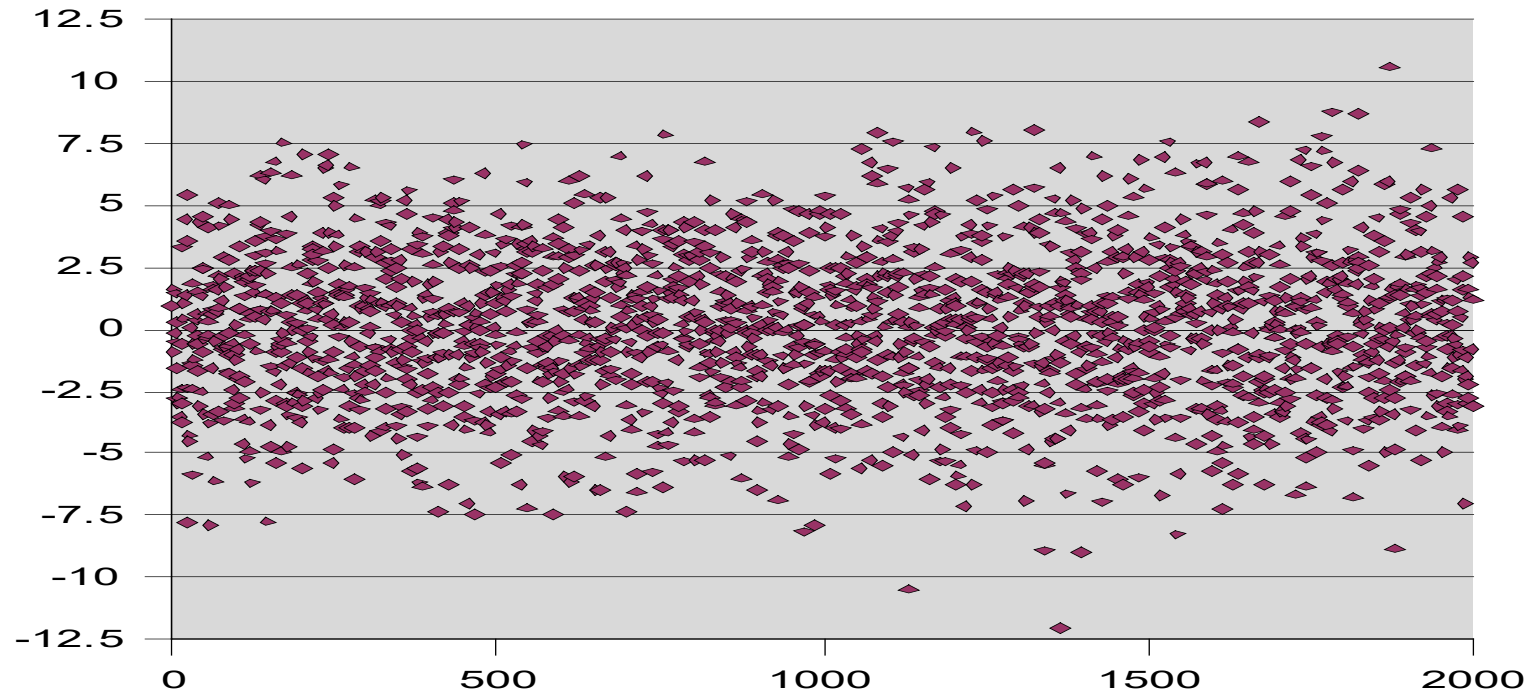
**Assimilated to
the Trend
Component**



Common Models

- White Noise
- AR
- MA
- ARMA
- ARIMA
- SARIMA
- ARMAX
- Kalman Filter
- Exponential Smoothing, trend, seasons

White Noise



mean=zero, Finite std.dev., often unknown
Often, but not always, Gaussian

Autoregressive Models (AR)

An **autoregressive model** is when a value from a time series is regressed on previous values from that same time series. For example, y_t on y_{t-1} :

$$y_t = \beta_0 + \beta_1 y_{t-1} + \epsilon_t.$$

In this regression model, the response variable in the previous time period has become the predictor and the errors have our usual assumptions about errors in a simple linear regression model. The **order** of an autoregression is the number of immediately preceding values in the series that are used to predict the value at the present time. So, the preceding model is a first-order autoregression, written as AR(1).

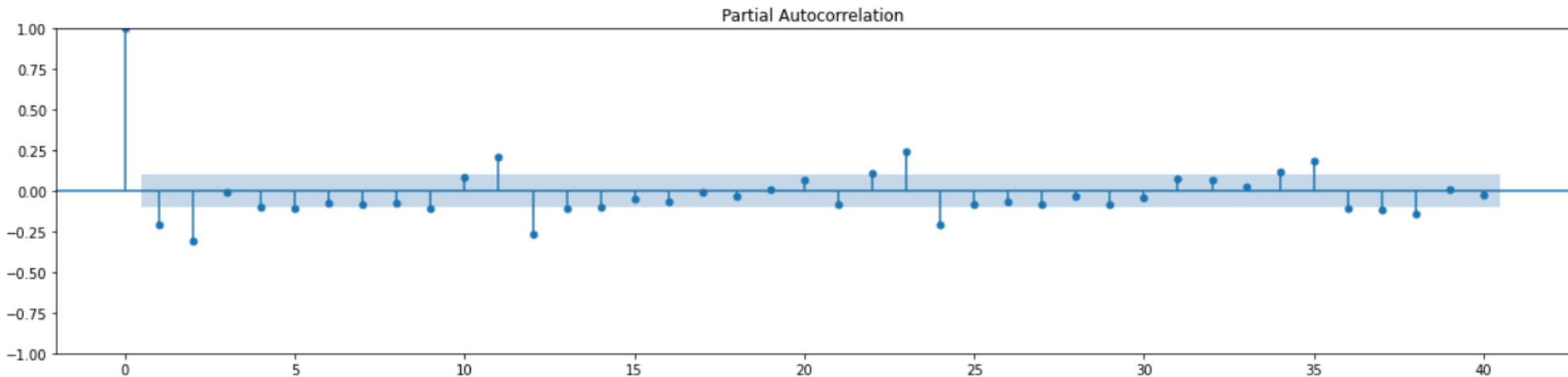
If we want to predict y using measurements of the previous two observations (y_{t-1}, y_{t-2}), then the autoregressive model for doing so would be:

$$y_t = \beta_0 + \beta_1 y_{t-1} + \beta_2 y_{t-2} + \epsilon_t.$$

This model is a second-order autoregression, written as AR(2), since the value at time t is predicted from the values at times $t-1$ and $t-2$.

How to Determine AR Order?

PACF (Partial Autocorrelation Function) is typically used for determining p values. For a given observation in a time series X_t , it may be correlated with a lagged observation X_{t-3} which is also impacted by its lagged values (e.g. X_{t-2} , X_{t-1}). PACF visualizes the **direct** contribution of the past observation to the current observations. For example, the PACF below when lag = 3 the PACF is roughly -0.60, which reflects the impact of lag 3 on the original data point, while the compound factor of lag 1 and lag 2 on lag 3 are not explained in the PACF value. The p values for the AR(p) model is then determined by when the PACF drops to below significant threshold (blue area) for the first time, i.e. $p = 4$ in this example below.



Moving Average (MA)

Moving average model (MA) model generates the current values based on the **ERRORS** from the past forecasts instead of using the past values like AR. Past errors are analyzed to produce the current value. Perfecting a baking recipe will be like a moving average model. You will do adjustments for needed sugars or butter for today's baking depending on the previous days' amount to perfect the recipe.

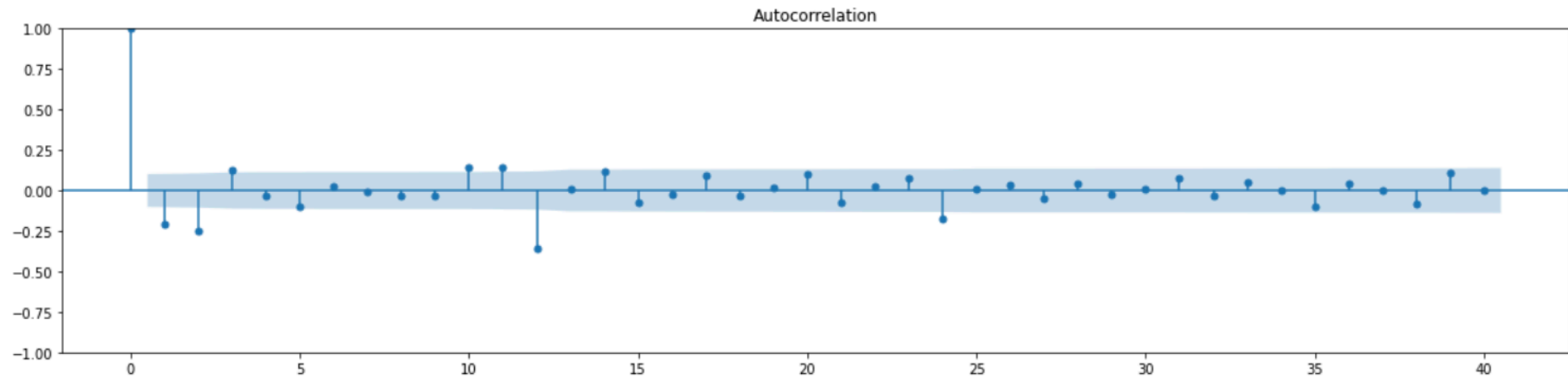
$$MR(X_t, q) = e_t + \theta_1 e_{t-1} + \theta_2 e_{t-2} + \dots + \theta_q e_{t-q}$$

MR(q)

Moving average model, MR(q) adjusts the model based on the average predictions errors from previous q observations, which can be stated as below, where e represents the error terms and ϑ represents the weights. q value determines the number of error terms to include in the moving average window.

How to determine the q value in MA?

ACF can be used for determining the q value. It is typically selected as the first lagged value of which the ACF drops to nearly 0 for the first time. For example, we would choose $q=4$ based on the ACF plot below.



Time Series Analysis

Time Series is a unique type of machine learning where time plays a critical role in model predictions.

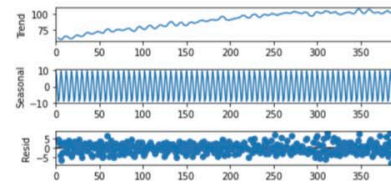
Time Series Components:

1. Trend
2. Seasonality
3. Residual

How to Describe Time Series?

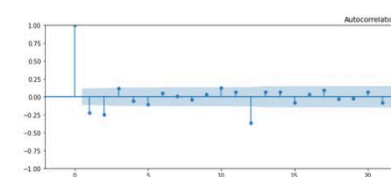
Decomposition

`seasonal_decompose(df)`



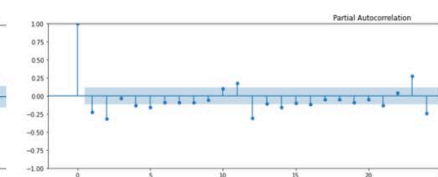
ACF

`sm.graphics.tsa.plot_acf()`



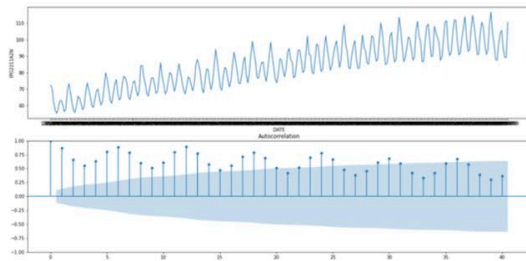
PACF

`sm.graphics.tsa.plot_pacf()`

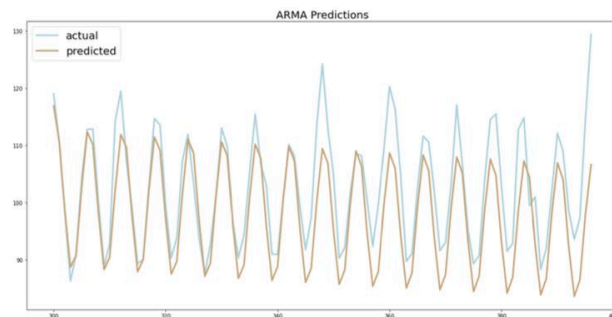


ARMA

Autoregressive
+ Moving Average

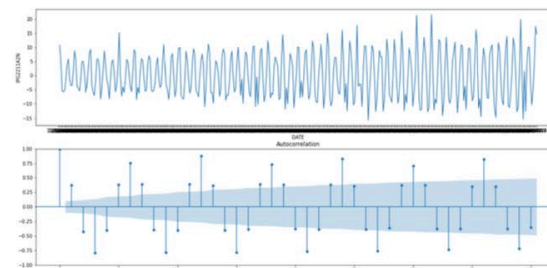


$ARIMA(p, 0, q)$

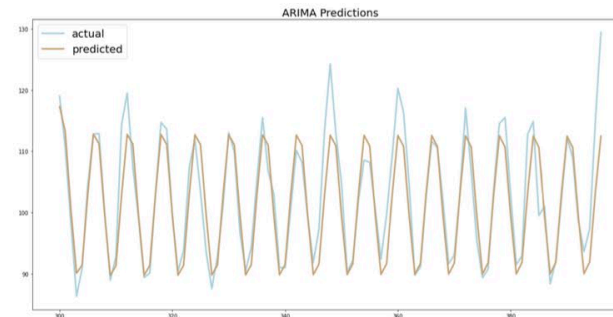


ARIMA

Autoregressive + Moving Average +
Trend Differencing

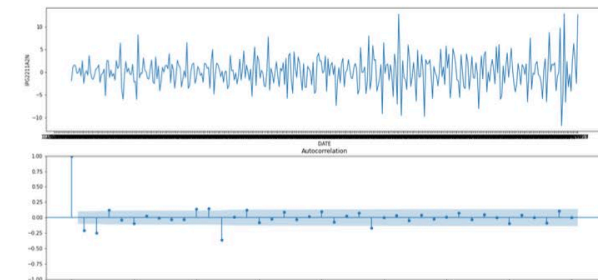


$ARIMA(p, d, q)$

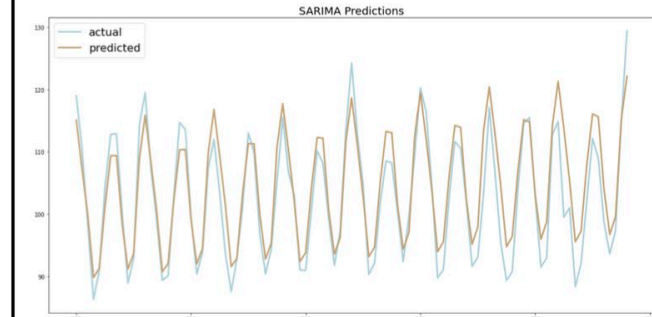


SARIMA

Autoregressive + Moving Average +
Trend Differencing + Seasonality Differencing



$SARIMAX(p, d, q) \times (P, D, Q, s)$



Autoregressive Integrated Moving Average

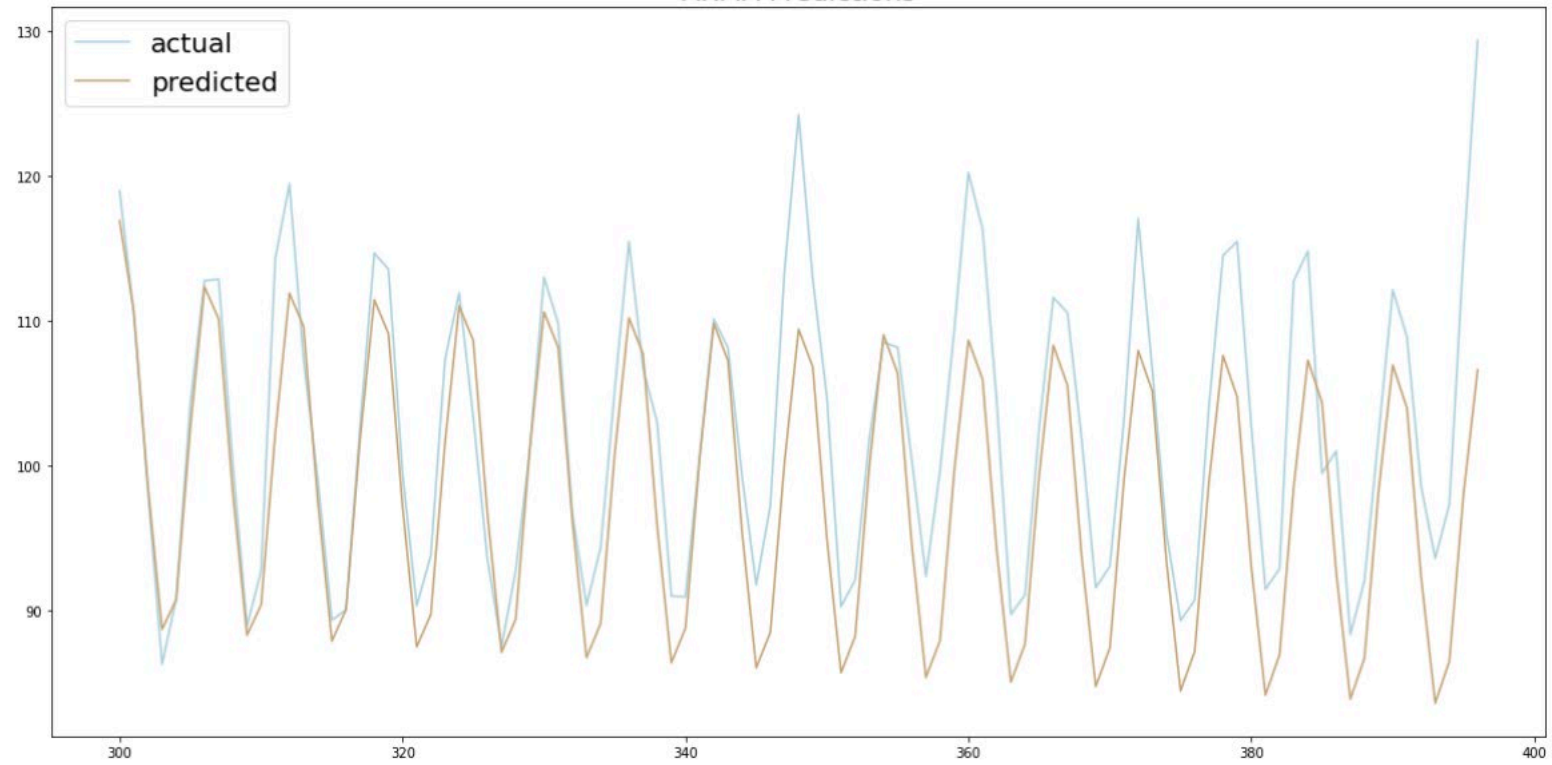
- ARIMA stands for **Autoregressive Integrated Moving Average**, which extends from ARMA model and incorporates the integrated component (inverse of differencing).
- ARIMA builds upon autoregressive model (AR) and moving average model (MA) by introducing degree of differencing components (specified as the parameter d) - [ARIMA \(p, d, q\)](#).

Seasonal ARIMA

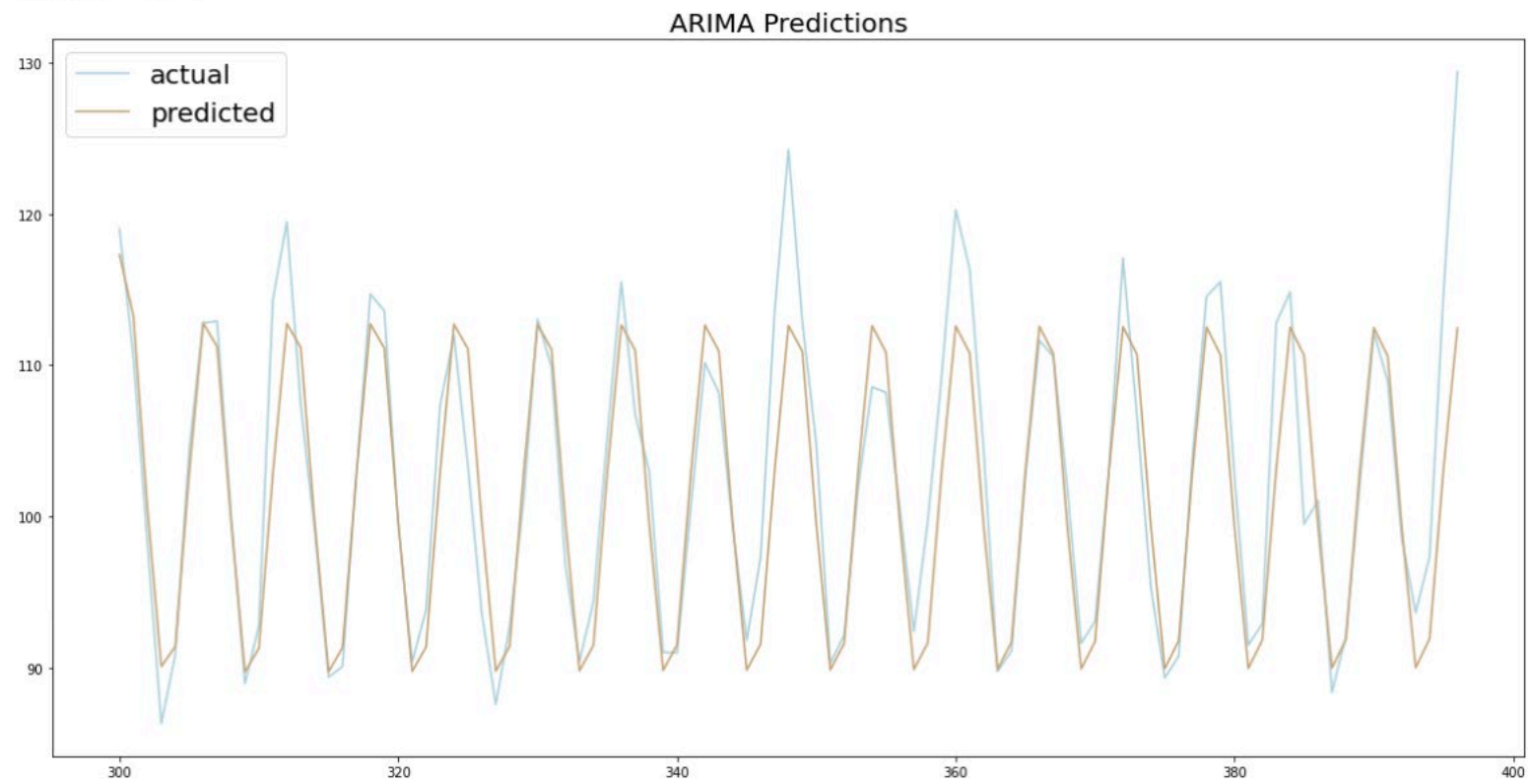
SARIMA stands for Seasonal ARIMA which addresses the periodic pattern observed in the time series. Previously we have introduced how to use seasonal differencing to remove seasonal effects. SARIMA incorporates this functionality to predict seasonally changing time series and we can implement it using $\text{SARIMAX}(p, d, q) \times (P, D, Q, s)$.

The first term (p, d, q) represents the order of the ARIMA model and (P, D, Q, s) represents the seasonal components. P, D, Q are the autoregressive, differencing and moving average terms of the seasonal order respectively. s is the number of observations in each period.

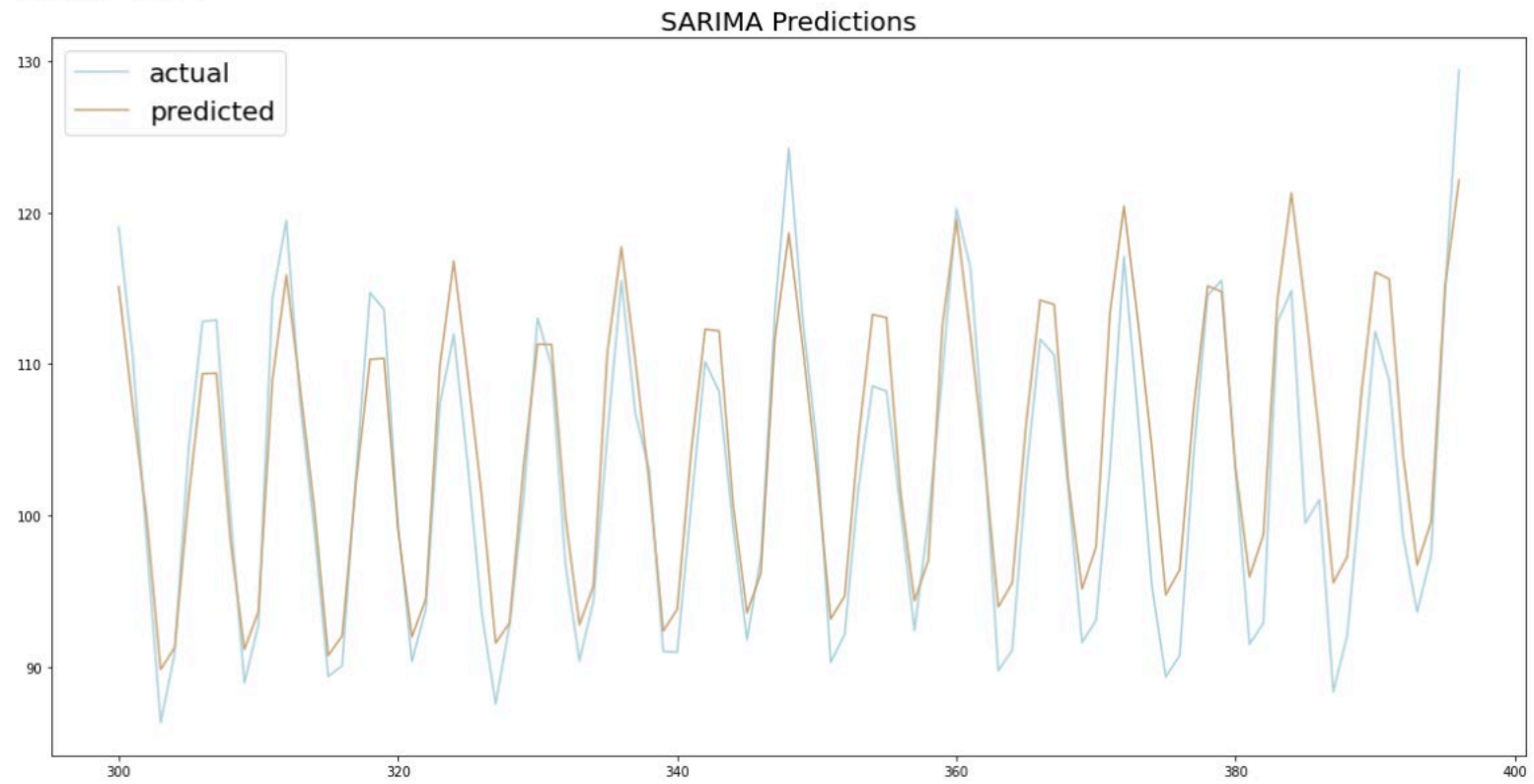
ARMA Predictions



RMSE: 4.35



RMSE: 4.04



Reading

1. *Durbin, J., & Koopman, S. J. (2012). Time series analysis by state space methods. Oxford university press.*
 2. *Giovanni Petris & Sonia Petrone (2011), State Space Models in R, Journal of Statistical Software*
 3. *G Petris, S Petrone, and P Campagnoli (2009). Dynamic Linear Models with R. Springer*
 4. *Hyndman, R. J., & Athanasopoulos, G. (2018). Forecasting: principles and practice. OTexts.*
-