# DATA MINING 2
## Time Series - Introduction & Preprocessing

Riccardo Guidotti

a.a. 2024/2025
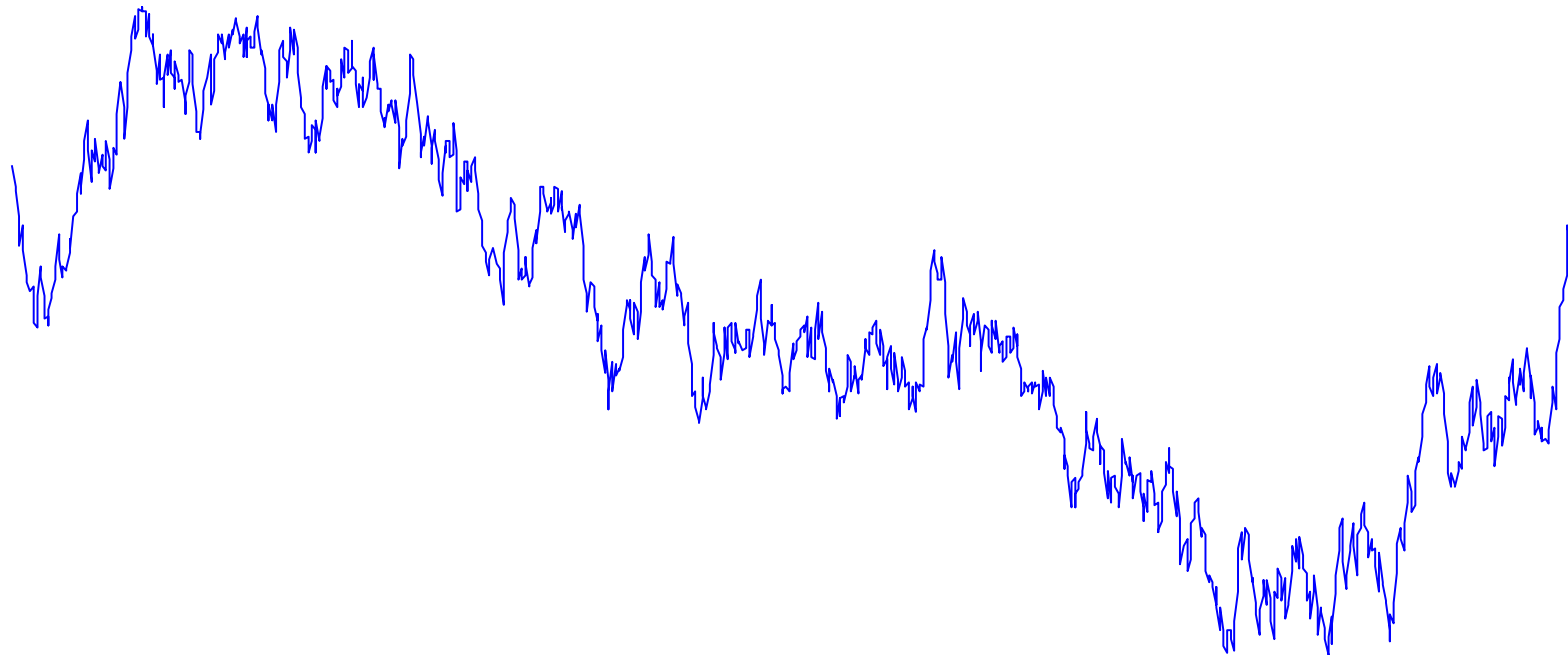
UNIVERSITÀ DI PISA

# What is a Time Series?

- A time series is a collection of observations made sequentially in time, generally at constant time intervals.



| |
|---|
| 25.1750 |
| 25.2250 |
| 25.2500 |
| 25.2500 |
| 25.2750 |
| 25.3250 |
| 25.3500 |
| 25.3500 |
| 25.4000 |
| 25.4000 |
| 25.3250 |
| 25.2250 |
| 25.2000 |
| 25.1750 |
| ... |
| 24.6250 |
| 24.6750 |
| 24.6750 |
| 24.6250 |
| 24.6250 |
| 24.6250 |
| 24.6750 |
| 24.7500 |

# What is Time Series Analysis?

- Time Series Analysis (TSA):
  - TSA is a way of analyzing data points over an interval of time
  - Data points recorded at consistent intervals over a set period
  - Not just intermittent or random data collection

- Unique aspects of TSA:
  - Shows how variables change over time
  - Time is a crucial variable
  - Provides additional information and dependencies between data points

- Requirements for TSA:
  - Large number of data points for consistency and reliability
  - Ensures representative sample size
  - Helps cut through noisy data
  - Identifies trends, patterns, and accounts for seasonal variance

- Applications:
  - Forecasting: predicting future data based on historical data
  - Classification/Regression: predicting exogenous variable based on historical data

# Why Companies Benefits From TSA?

- Benefits of time series analysis:
  - Understands underlying causes of trends or systemic patterns over time
  - Uses data visualizations to see seasonal trends and explore reasons behind them
  - Modern analytics platforms offer advanced visualizations beyond line graphs

- Predictive capabilities:
  - Analyzes data at consistent intervals for time series forecasting
  - Predicts likelihood of future events
  - Identifies seasonality and cyclic behavior for better understanding and forecasting

- Advantages of modern technology:
  - Ability to collect massive amounts of data daily
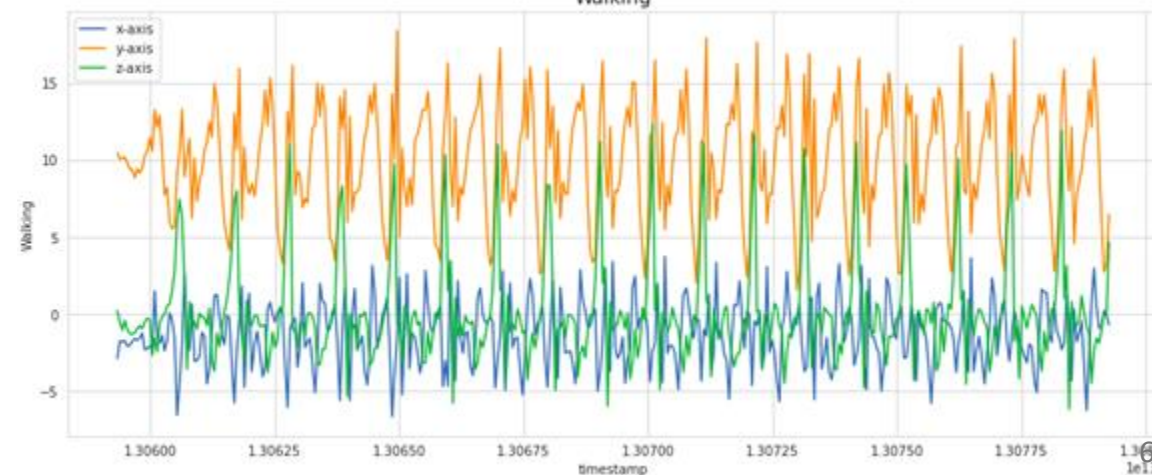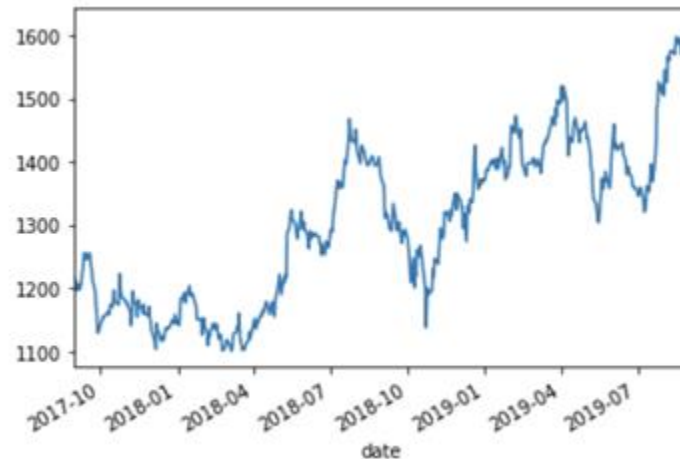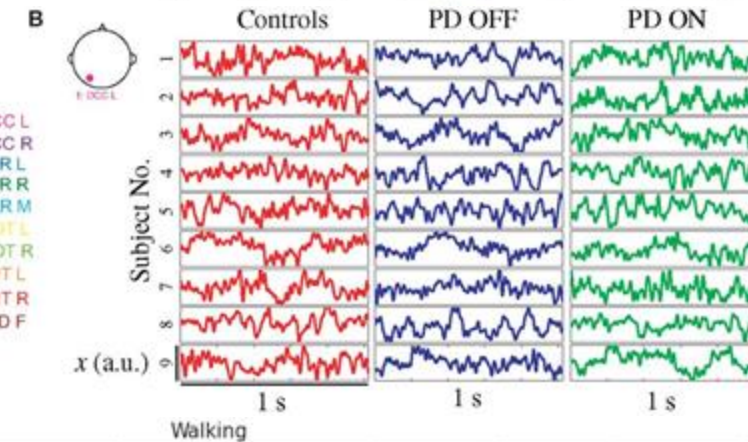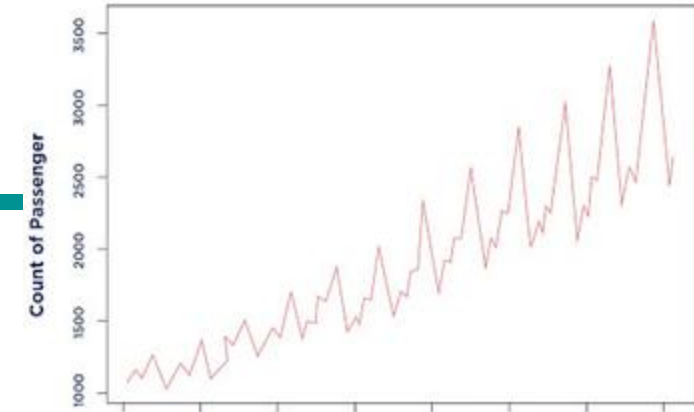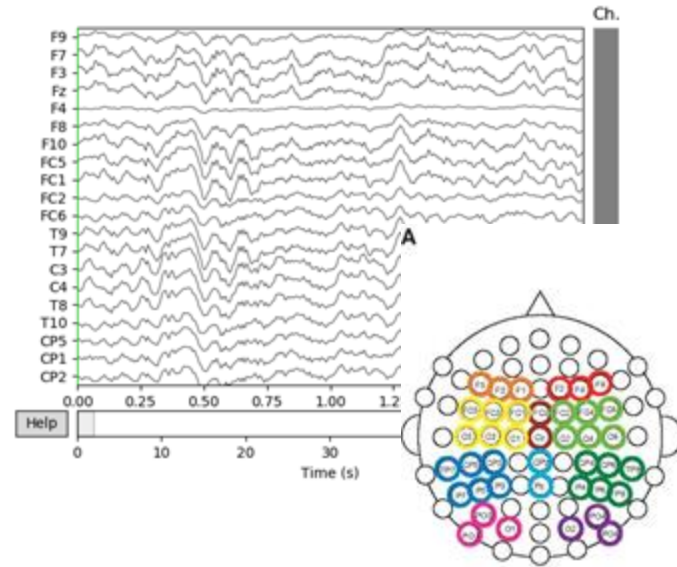  - Easier to gather enough consistent data for comprehensive analysis

# Intervals

- Time: Milliseconds, Seconds, Minutes, Hours, Days, Months, Years, ...

- Spatial: Locations, Positions, Machines, ...

- Relative: one cm left, two cm left, ...

- The important point is to have an ordered variable that provides a direction for its values.

- Then from a given observation $x_i$ it is easy to identify the past, i.e., what came before $x_i$, and the future, i.e., what comes after $x_i$

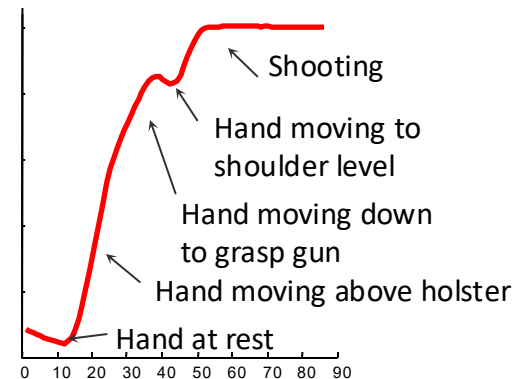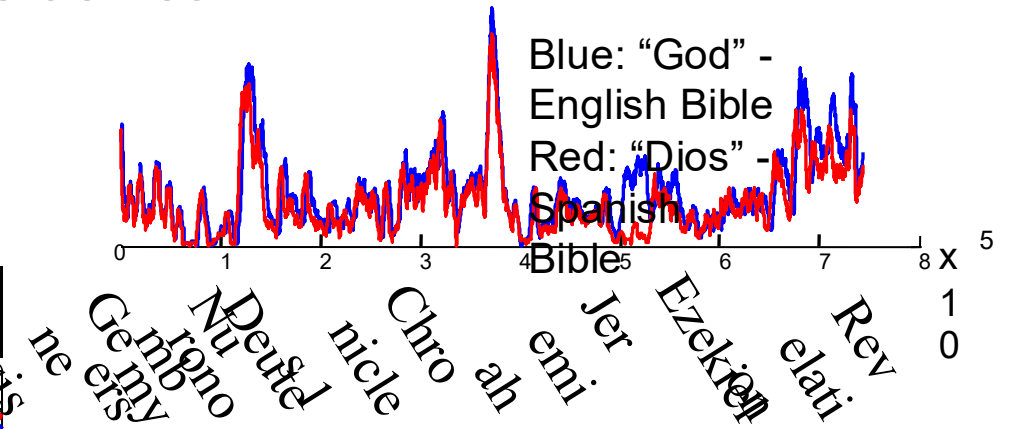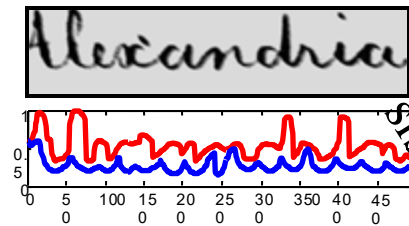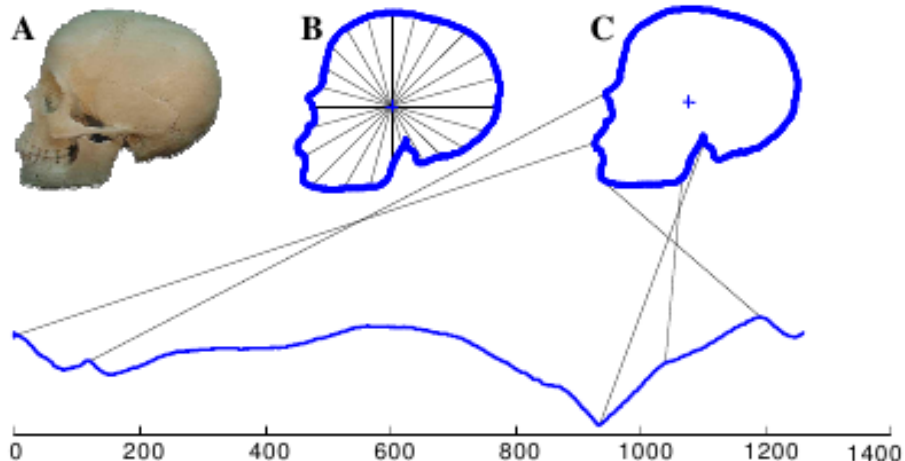| |
|---|
| 25.1750 |
| 25.2250 |
| 25.2500 |
| 25.2500 |
| 25.2750 |
| 25.3250 |
| 25.3500 |
| 25.3500 |
| 25.4000 |
| 25.4000 |
| 25.3250 |
| 25.2250 |
| 25.2000 |
| 25.1750 |
| ... |
| 24.6250 |
| 24.6750 |
| 24.6750 |
| 24.6250 |
| 24.6250 |
| 24.6250 |
| 24.6750 |
| 24.7500 |

# Time Series are Ubiquitous

- Blood pressure
- Politics popularity rating
- The annual rainfall in Pisa
- Passengers of a company
- Accelerations on different axes
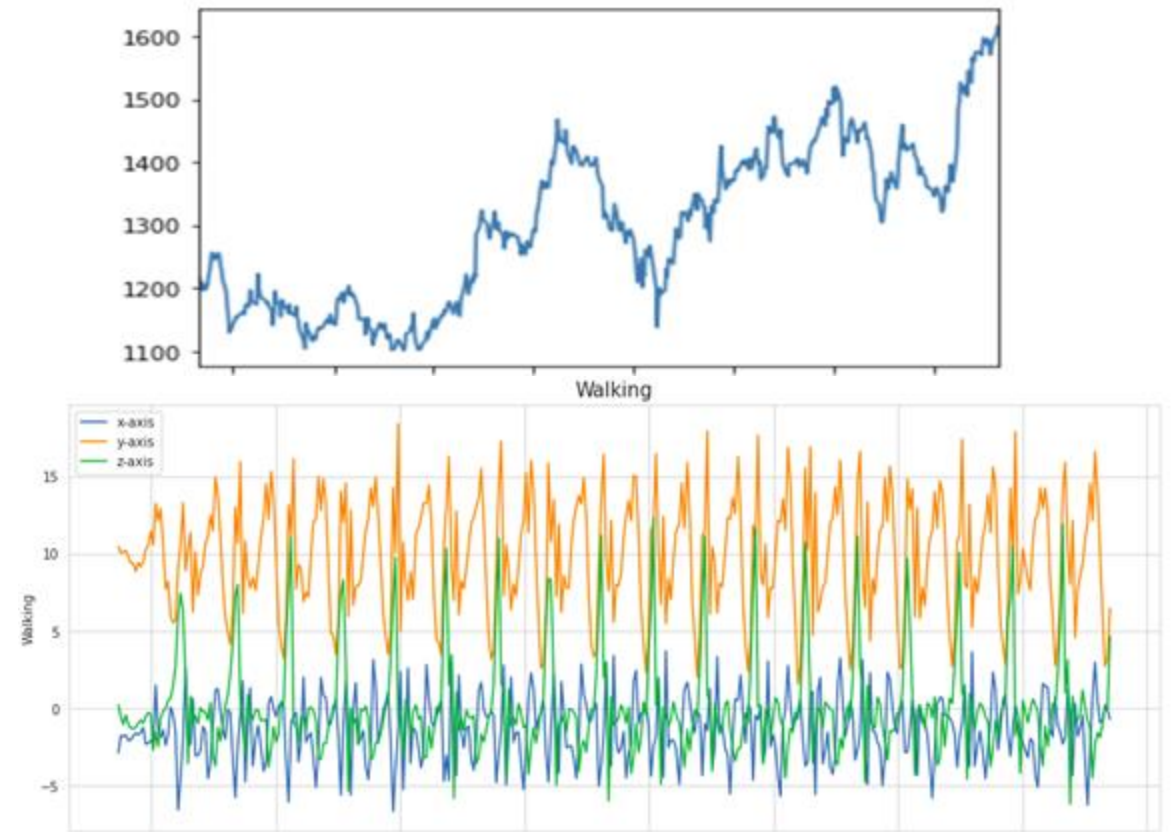- The value of your stocks
- EEG and ECG

# Time Series are Ubiquitous

- Other data types can be modeled as time series
  - Text data: words count
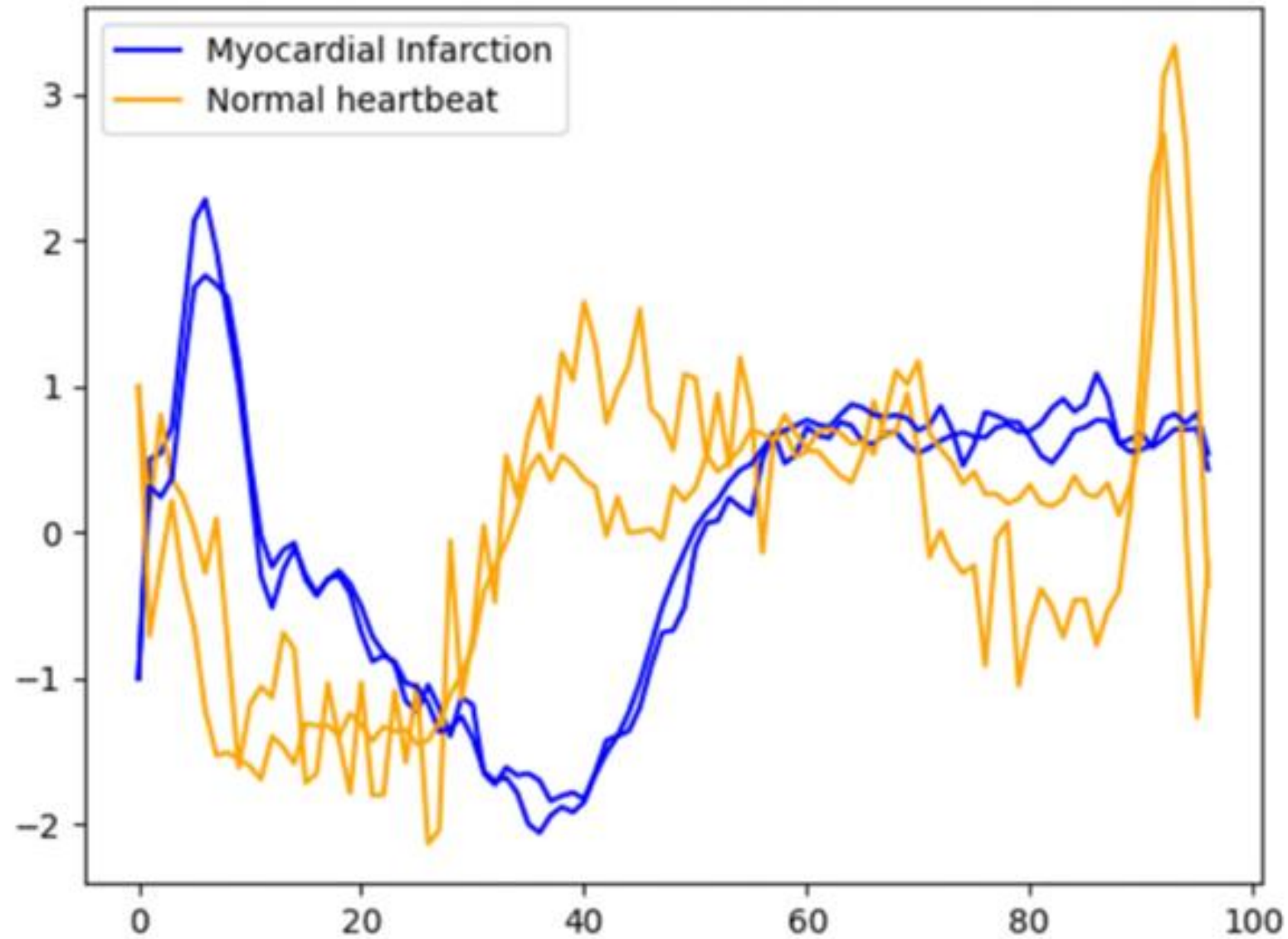  - Images: edges displacement
  - Videos: object positioning



Blue: "God" - English Bible
Red: "Dios" - Spanish Bible

Shooting
Hand moving to shoulder level
Hand moving down to grasp gun
Hand moving above holster
Hand at rest

# Time Series Data

- Time Series $T = \{T_1, ..., T_c\}$ is a collection of $c$ signals (or channels) each one with $m$ observations $x_i$, i.e., $T_j = \{x_1, ..., x_m\}$

- Univariate Time Series *(c = 1)*

- Multivariate Time Serie *(c > 1)*

# Time Series Data

- A Time Series Dataset is a collection of time series $X = \{T_1, \ldots T_n\}$.
- Associated to each time series $T_i$ we can have exogenous variables
  - Categorical Features
    - Accelerometers: Type of Movement, Crash vs NoCrash, Car Model, etc.
    - Machine Sensors: Failure vs NoFailure, Product in Production, Type of Machine, etc.
    - EEG: Syntoms, Seizure vs NoSeizure
    - Students Marks: Student Name, Course, Degree, University, Background
  - Continuous Features
    - Accellerometers: Age, Weight, Engine Temperature, etc.
    - Machine Sensors: Number of Products Realized, System Temperature, etc.
    - EEG: Age, Weight, Height, etc.
    - Students Marks: Age, Family Income, Weight, etc.

Dataset ECG200

# Time Series in Datasets

- Stored ''horizontally'', typically univariate TS

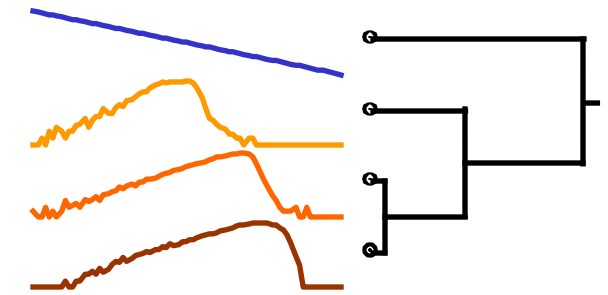|        | $t_0$ | $t_1$ | $t_2$ | $t_3$ | $t_4$ | $t_5$ | $t_6$ | $t_7$ | $t_8$ | $t_9$ |
|--------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| w4r5   | 111   | 110   | 111   | 112   | 120   | 123   | 116   | 118   | 119   | 123   |
| e2t6   | 89    | 76    | 79    | 85    | 67    | 56    | 78    | 97    | 45    | 78    |
| q1e4   | 91    | 95    | 89    | 87    | 110   | 120   | 125   | 130   | 111   | 90    |
| w23t   | 110   | 111   | 112   | 122   | 123   | 112   | 118   | 119   | 123   | 119   |

# Time Series in Datasets

- Stored ''vertically'', typically multivariate TS. Wide Dataset

|  |  | temp | speed | rot | product | failure | sys temp |
|---|---|---|---|---|---|---|---|
| w4r5 | $t_0$ | 20.1 | 111 | 3 | A | 0 | 32.5 |
| w4r5 | $t_1$ | 18.6 | 110 | 4 | A | 0 | 32.5 |
| w4r5 | $t_2$ | 19.4 | 111 | 3 | A | 1 | 32.5 |
| w4r5 | $t_3$ | 20.4 | 112 | 5 | A | 0 | 32.5 |
| w4r5 | $t_4$ | 21.5 | 120 | 6 | A | 0 | 32.5 |
| e2t6 | $t_0$ | 12.7 | 89 | 29 | B | 0 | 34.6 |
| e2t6 | $t_1$ | 19.8 | 76 | 45 | B | 0 | 34.6 |
| e2t6 | $t_2$ | 17.4 | 69 | 34 | B | 0 | 34.6 |
| e2t6 | $t_3$ | 8.4 | 85 | 22 | B | 1 | 34.6 |
| e2t6 | $t_4$ | 7.9 | 65 | 19 | B | 1 | 34.6 |

# Time Series in Datasets

- Stored ''vertically'', typically multivariate TS.
- Long Dataset

| | | feat | value |
|---|---|---|---|
| w4r5 | $t_0$ | temp | 20.1 |
| w4r5 | $t_1$ | temp | 18.6 |
| w4r5 | $t_2$ | temp | 19.4 |
| w4r5 | $t_3$ | temp | 20.4 |
| w4r5 | $t_4$ | temp | 21.5 |
| w4r5 | $t_0$ | speed | 111 |
| w4r5 | $t_1$ | speed | 110 |
| w4r5 | $t_2$ | speed | 111 |
| w4r5 | $t_3$ | speed | 112 |
| w4r5 | $t_4$ | speed | 120 |
| … | … | … | … |

# Time Series Analytics Tasks

- Classification
- Regression
- Forecasting
- Clustering
- Anomaly Detection
- Pattern Mining

# Time Series Classification - TSC

- Given a dataset $X = \{T_1, \ldots T_n\}$, TSC is the task of training a model $f$ to predict an exogenous <u>categorical</u> output $y$ for each time series $T$, i.e., $f(T) = y$.

Labelled training series

Classify unlabelled series

?

# Time Series Classification - Examples

- **Predictive variables**: multivariate time series as accelerometers on the x, y, z axes and speed coming from cars black boxes. **Target variable**: means of transport, crash vs noCrash, car model, etc.

- **Predictive variables**: multivariate time series as accelerometers on the x, y, z axes, heart rate and number of steps per sec coming from personal smart devices. **Target variable**: type of movement, type of device, next location.

- **Predictive variables**: multivariate time series as machine sensors such has temperature, humidity, number of rotations, accelerations, etc. **Target variable**: failure vs nofailure, product in production, type of machine, etc.

- **Predictive variables**: multivariate time series as EEG or ECG. **Target variable**: symptoms, disease, treatment, etc.

- **Predictive variables**: univariate or multivariate time series as students' marks. **Target variable**: student background, university, etc.

# Time Series Extrinsic Regression - TSER

- Given a dataset $X = \{T_1, \dots T_n\}$, TSER is the task of training a model $f$ to predict an exogenous <u>continuous</u> output $y$ for each time series $T$, i.e., $f(T) = y$.
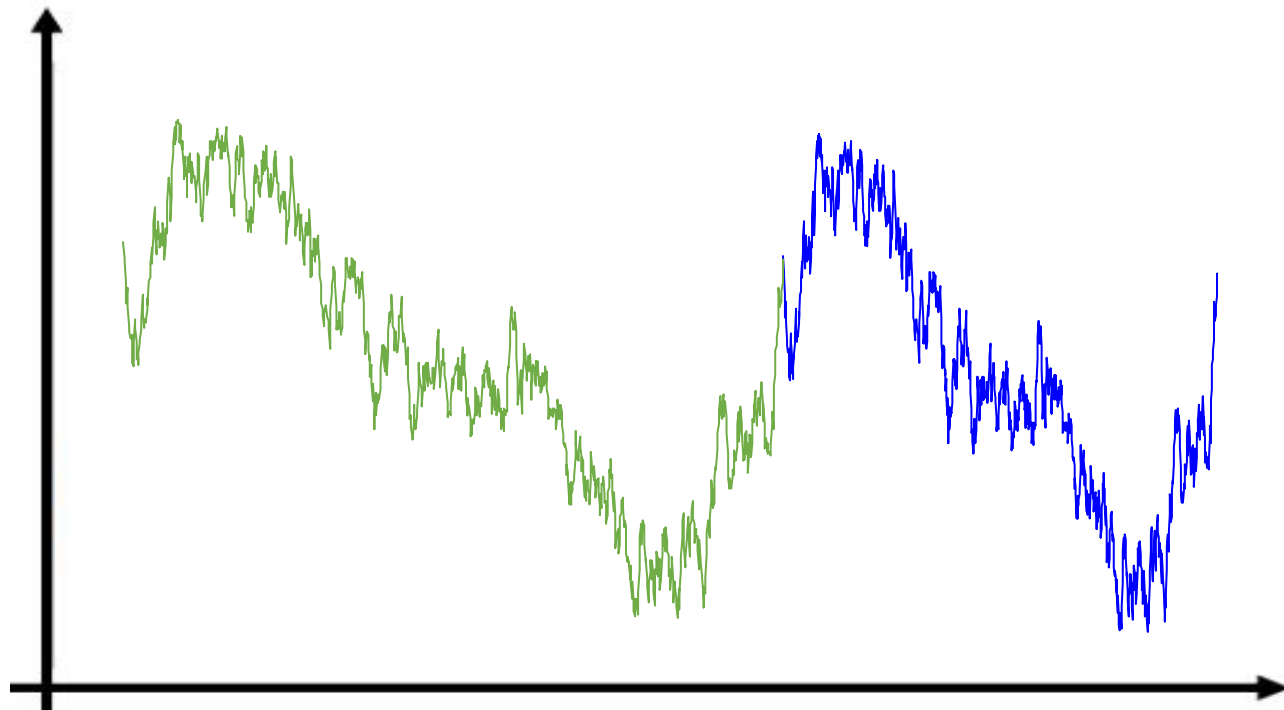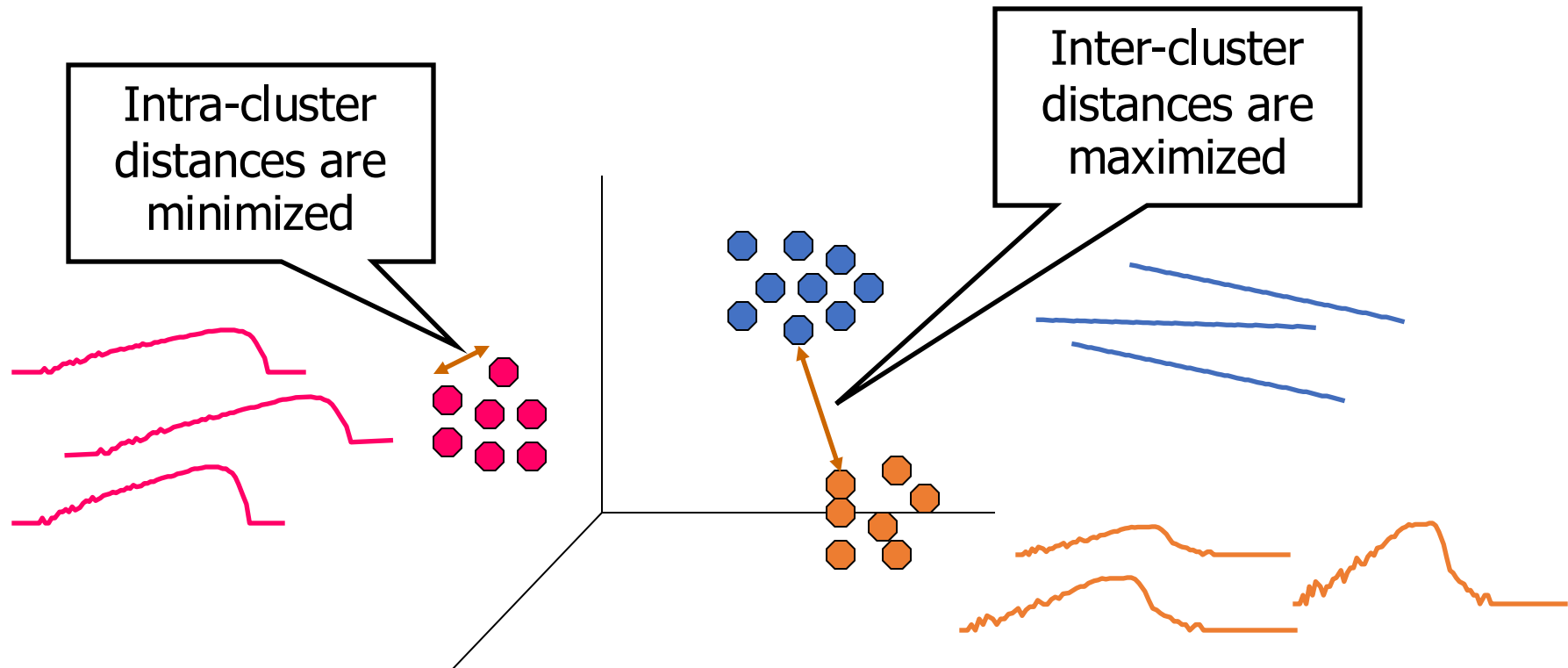
Labelled training series

3.2
2.9
5.6
6.1

Classify
unlabelled series

?

# Time Series Regression - Examples

- ***Predictive variables***: multivariate time series as accelerometers on the x, y, z axes and speed coming from cars black boxes. **Target variable**: engine temperature, number of car repairs, etc.

- ***Predictive variables*** multivariate time series as accelerometers on the x, y, z axes, heart rate and number of steps per sec coming from personal smart devices. ***Target variable***: age, number of times gym is made.

- ***Predictive variables***: multivariate time series as machine sensors such has temperature, humidity, number of rotations, accelerations, etc. ***Target variable***: number of products realized, number of failures.

- ***Predictive variables***: multivariate time series as EEG or ECG. ***Target variable***: patient age, patient weight, etc.

- ***Predictive variables***: univariate or multivariate time series as students' marks. ***Target variable***: student age, family income.

# Time Series Forecasting - TSF

- Given a dataset $X = \{T_1, \dots T_n\}$, TSF is the task of training a model $f$ to predict an endogenous <u>continuous</u> output $y$ for each time series $T$, i.e., $f(T) = y$.
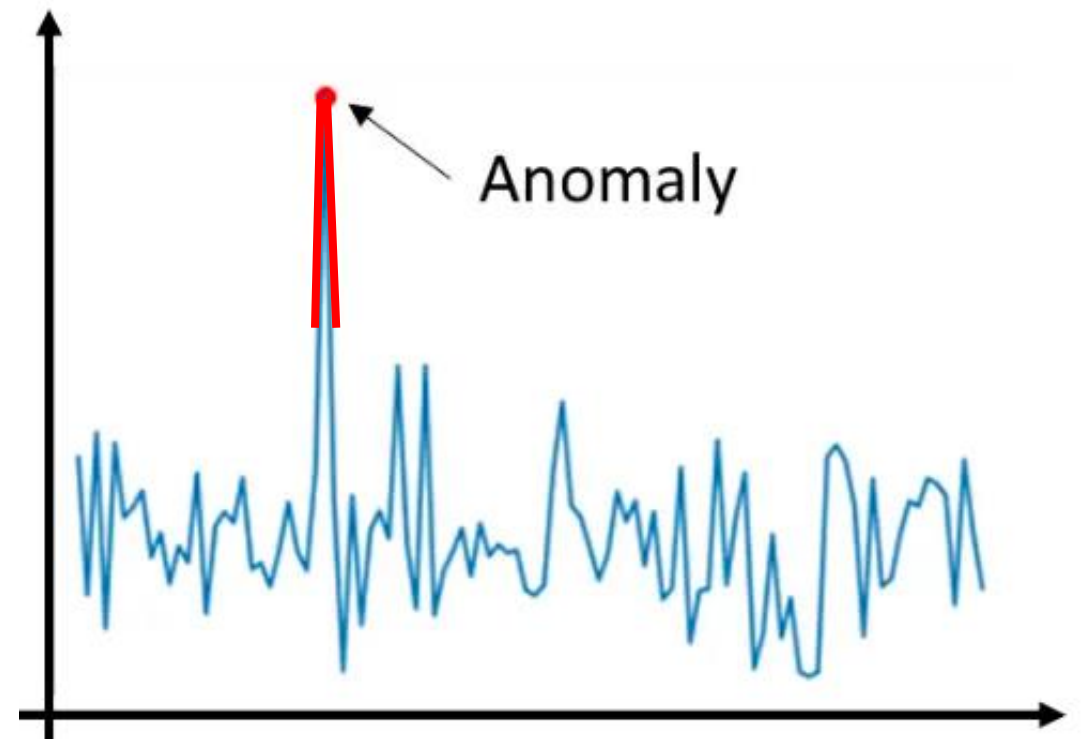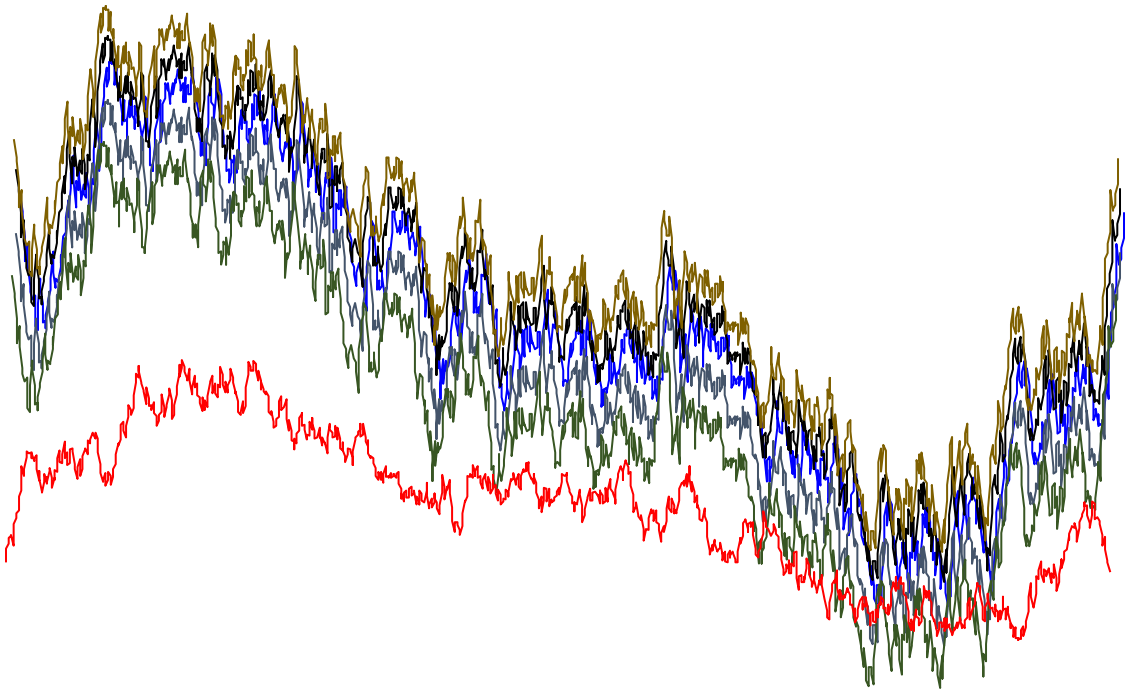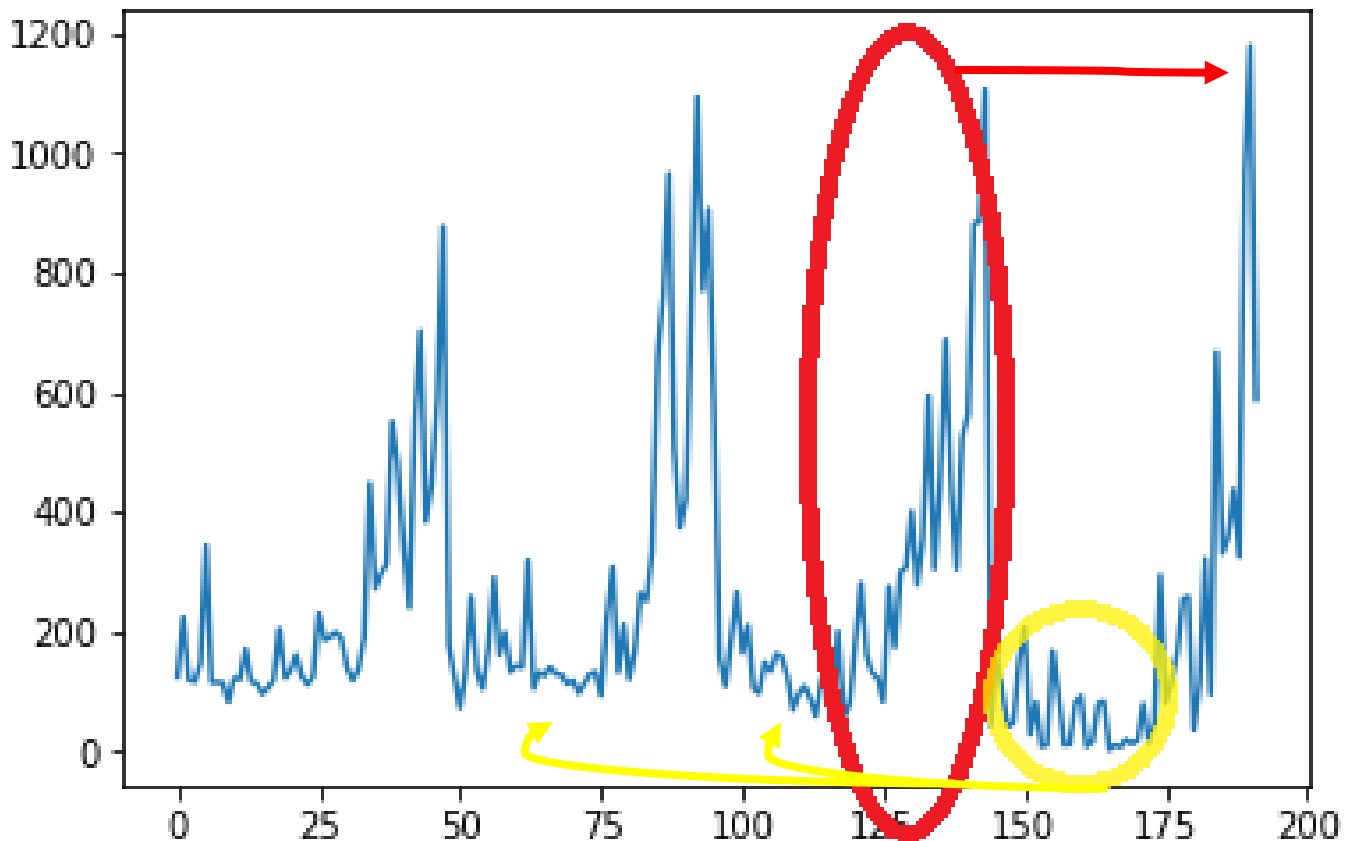
# Time Series Forecasting - Examples

- *Predictive variables*: multivariate time series as accelerometers on the x, y, z axes and speed coming from cars black boxes. *Target variable*: time series as accelerometers on the x, y-, z axes and speed.

- *Predictive variables*: multivariate time series as accelerometers on the x, y, z axes, heart rate and number of steps per sec coming from personal smart devices. *Target variable*: time series as accelerometers on the x, y, z axes, heart rate and number of steps per sec .

- *Predictive variables*: multivariate time series as machine sensors such has temperature, humidity, number of rotations, accelerations, etc. *Target variable*: temperature, humidity, number of rotations, accelerations.

- *Predictive variables*: multivariate time series as EEG or ECG. *Target variable*: future values of EEG or ECG, etc.

- *Predictive variables*: univariate or multivariate time series as students' marks. *Target variable*: future students' marks.

# Time Series Clustering

- Given a dataset $X = \{T_1, \dots T_n\}$, Time Series Clustering is the task of grouping similar time series such that those in a group are similar to one another and different from the time series in other groups.

Intra-cluster distances are minimized

Inter-cluster distances are maximized

# Time Series Clustering - Examples

- *Variables*: multivariate time series as accelerometers on the x, y, z axes and speed coming from cars black boxes. *Result*: groups of similar cars.

- *Variables*: multivariate time series as accelerometers on the x, y, z axes, heart rate and number of steps per sec coming from personal smart devices. *Result*: groups of similar users.

- *Variables*: multivariate time series as machine sensors such has temperature, humidity, number of rotations, accelerations, etc. *Results*: groups of similar products produced.

- *Variables*: multivariate time series as EEG or ECG. *Results*: groups of similar patients.

- *Variables*: univariate or multivariate time series as students' marks. *Results*: groups of similar students.

# Anomaly Detection

- Given a dataset $X = \{T_1, \dots T_n\}$, Anomaly Detection is the task of:
  a) identifying anomalous time series within the set $X$
  b) identifying anomalous time stamps for each time series $T$



Anomaly

# Time Series Anomaly Detection - Examples

- *Variables*: multivariate time series as accelerometers on the x, y, z axes and speed coming from cars black boxes. *Result*: anomalous cars or anomalous time stamps.

- *Variables*: multivariate time series as accelerometers on the x, y, z axes, heart rate and number of steps per sec coming from personal smart devices. *Result*: anomalous users or anomalous time stamps.

- *Variables*: multivariate time series as machine sensors such has temperature, humidity, number of rotations, accelerations, etc. *Results*: anomalous products/machines or anomalous time stamps.

- *Variables*: multivariate time series as EEG or ECG. *Results*: anomalous patients or anomalous time stamps.

- *Variables*: univariate or multivariate time series as students' marks. *Results*: anomalous students or anomalous time stamps.

# Pattern Mining

- Given a dataset $X = \{T_1, \dots T_n\}$, Pattern Mining is the task of identifying repeated subsequences in each time series $T$.

# Problems in Working with Time Series

- Large amount of data

- Similarity is not easy to estimate

- Different data formats

- Different sampling rates

- Noise

- Missing values

- …

# Assumption

We assume that for a given dataset $X = \{T_1, \dots T_n\}$, all the time series have aligned time stamps, i.e., the *i-th* time stamp of a time series $T_a$ corresponds to the *i-th* time stamp of another time series $T_b$ and this is true for all the time series in $X$.

# Assumption

We assume that for a given dataset $X = \{T_1, ... T_n\}$, all the time series have aligned time stamps, i.e., the *i-th* time stamp of a time series $T_a$ corresponds to the *i-th* time stamp of another time series $T_b$ and this is true for all the time series in $X$.

# Time Series Visualization

A plot can reveal

- Trend: upward or downward pattern
- Periodicity: repetition of behavior in a regular pattern
- Seasonality: periodic behavior with a known period (hourly, weekly, monthly...)
- Heteroskedasticity: changing variance
- Dependence: positive (successive observations are similar) or negative (successive observations are dissimilar)
- Outliers: anomalous time stamps, anomalous subsequences
- Missing data: missing values at a certain time stamp or longer subsequences

# Example 1

- Airline passengers from 1949-1961

# Example 1

- Airline passengers from 1949-1961
- Upward trend
- Seasonality on a 12 month interval
- Increasing variability

# Example 2

- U.S.A. population at ten year intervals from 1790-1990

# Example 2

- U.S.A. population at ten year intervals from 1790-1990
- Upward trend
- Slight change in shape/structure
- Nonlinear behavior

# Example 3

- Monthly Beer Production in Australia

# Example 3

- Monthly Beer Production in Australia
- No trend in last 100 months
- No clear seasonality

# Example 4

- Yield from a controlled chemical batch process

# Example 4

- Yield from a controlled chemical batch process
- Negative dependence: successive observations tend to lie on opposite sides of the mean.

TS Histogram

TS BoxPlot

Time in not considered in these plots!!!

# Missing Values

# Missing Values

- Individual values for a single time stamp can be missing
- Contiguous values for sequential time stamps can be missing

# Missing Values Imputation Methods

- Fill with constant value
- Linear interpolation
- Forecasting
- Random

# Missing Values Imputation

# Missing Values Imputation: Padding

- Fill with constant value: **last** value (pad)

# Missing Values Imputation: BackFilling

- Fill with constant value: **next** value (backfill)

# Missing Values Imputation: Mean

- Fill with constant value: **mean/median** value

# Missing Values Imputation: Nearest Value

- Fill with constant value: **nearest** value

# Missing Values Imputation: Interpolation

- Interpolate the last and first not missing values to get the missing ones.

# Missing Values Imputation: Forecasting

• Interpolate using a forecasting or regressive model (see next lectures)

# Anomalies

# Anomalies and Outliers in Time Series

- An outlier (or anomaly) is a value or an observation that is distant from other observations, a data point that differ significantly from other data points.

- Outlier: "An observation which deviates so much from other observations as to arouse suspicions that it was generated by a different mechanism." [Hawkins 1980]

- Sometimes it makes sense to formally distinguish two classes of outliers: Extreme values and mistakes.

- Mistakes can be considered as missing values and removed.

- Extreme values might be considered in the analysis.

# Outlier Detection Methods

- Outlier detection methods **create a model of the normal patterns in the data**, and then compute an outlier score of a given data point based on the deviations from these patterns.

- Different models make different assumptions about the "normal" behavior.

- The outlier score of a data point is then computed by evaluating the quality of the fit between the data points and the model.

- In practice, the choice of the model is often dictated by the analyst's understanding of the kinds of deviations relevant to an application.

# Outlier Types in Time Series

Outlier detection methods may differ depending on the type of outliers:

- **Point Outlier**: A value that behaves unusually in a specific time instant when compared either to the other values in the time series (global outlier) or to its neighboring points (local outlier).

- **Subsequences**: Consecutive points in time whose joint behavior is unusual, although each observation individually is not necessarily a point outlier

- **Instance**: Entire time series can also be outliers, but they can only be detected when the input data is a dataset of time series.



53

# Examples of Anomalies in Time Series

- **Level shifts:** when the signal moves to zero/default value for a short period of time.
- **Unexpected growth/decrease** in a short period of time that looks like a spike.

# Outlier Detection Method

- An outlier is a point that significantly deviates from its expected value.

- Given a univariate time series $x$, a point at time $t$ can be declared an outlier if the distance to its expected value is higher than a predefined threshold.

- **Estimation method**: if $\bar{x}_t$ is obtained using previous and subsequent observations to $x_t$ (past, current, and future data).

- **Prediction method**: if $\bar{x}_t$ is obtained relying only on previous observations to $x_t$ (past data).

# Histogram

# Boxplot

- Data represented with a **box**
- The ends of the box are at the
  - **Q1: 1$^{st}$ quartiles** (25%-quantile or 25$^{th}$ percentile)
  - **Q3: 3$^{rd}$ quartiles** (75%-quantile or 75$^{th}$ percentile)
- **Median**: value in the middle is the **Q2: 2$^{nd}$ quartile** (50%-qua,, 50$^{th}$ perc.)
- The height of the box is **Interquartile range** (**IQR**): Q3 - Q1
- **Whiskers**: two lines outside the box extended from:
  - 1$^{st}$, or 5$^{th}$, or 10$^{th}$ percentile, or Q1 – $k$ IQR (with k = 1.5)
  - 99$^{th}$, or 95$^{th}$, or 90$^{th}$ percentile, or Q3 + $k$ IQR (with k = 1.5)
- **Outliers**: are points beyond whiskers
- In general, p%-quantile (0 < p < 100): Is the value $x$ s.t. p% of the values are smaller and 100-p% are larger.



Outlier

90$^{th}$ percentile

75$^{th}$ percentile

50$^{th}$ percentile

25$^{th}$ percentile

10$^{th}$ percentile

# Inter-Quartile Range Filter

- The IQR criteria means that all observations
- above Q3 + 1.5 * IQR or
- below Q1 − 1.5 * IQR
- where Q3 and Q1 correspond to third ($q_{0.75}$) and first ($q_{0.25}$) quartile respectively, and $IQR$ is the difference between the third and first quartile) are considered as potential outliers.

+ ← Outlier

+
+

← 90th percentile

← 75th percentile

← 50th percentile

← 25th percentile

← 10th percentile

# Hampel Filter

- Consider as outliers the values outside the interval ($I$) formed by the median, plus or minus 3 Median Absolute Deviations (MAD):

- $I=[median_X - 3 * MAD; median_X + 3 * MAD]$

- where $MAD$ is the median absolute deviation and is defined as the median of the absolute deviations from the data's median, i.e.,

- $MAD = median(|x_i - median_X|)$

# Grubbs' Test



- Detect outliers in univariate data

- Assume data comes from normal distribution

- Detects one outlier at a time, remove the outlier, and repeat
  - $H_0$: There is no outlier in data
  - $H_A$: There is at least one outlier

- Grubbs' test statistic:
  - one-sided test with alpha/N
  - two-sided test with alpha/2N

- Reject null hypothesis $H_0$ of no outliers if:

mean

std dev

$$G = \frac{\max\left|X - \overline{X}\right|}{s}$$

alpha significance
t – Student's distribution

$$G > \frac{(N-1)}{\sqrt{N}} \sqrt{\frac{t_{(\alpha/N, N-2)}^2}{N - 2 + t_{(\alpha/N, N-2)}^2}}$$

degrees of freedom

upper critical value of t-distribution

# Handling Outliers in Time Series

Once that anomalies have been identified

a)  they can be removed and treated like missing values, i.e., replaced according to different strategies.

b)  they can be maintained in the analysis.

c)  they can be distinguished between "mistakes" and "extreme values" and treated according the most preferrable strategy.

# Normalizations

# Problem with Time Series Distortions

- A common tool of TSA consists in calculating distances among time series.

- Many ML-tools require that the input data is represented in the same range of values or that values are comparable.

- Distance calculations and ML models are very sensitive to "distortions" in the data.

- These distortions are dangerous and should be removed.

- Most common distortions:
    - Offset Translation
    - Amplitude Scaling
    - Linear Trend
    - Presence of Noise

- These distortions can be removed by using the appropriate normalizations.

# Offset Translation: Mean Removal



$D(Q,C)$

$Q = Q - mean(Q)$

$C = C - mean(C)$
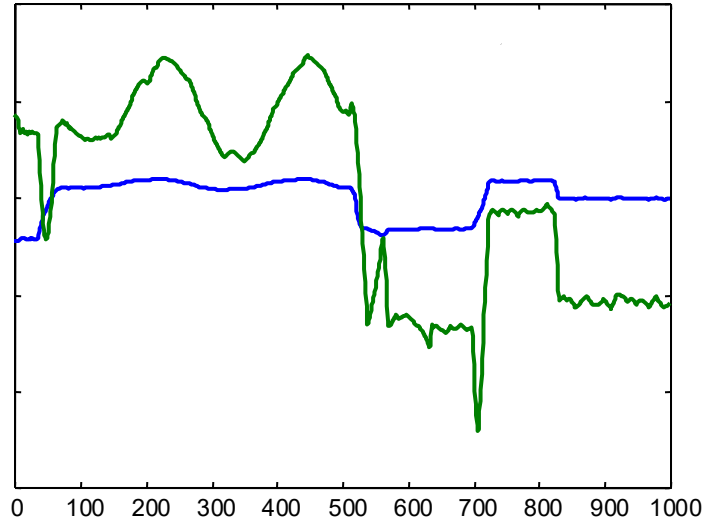
# Offset Translation: Min-Max Normalization



$$D(Q,C)$$

$$Q = (Q - min(Q)) / (max(Q) - min(Q))$$

$$C = (C - min(C)) / (max(C) - min(C))$$

# Amplitude Scaling: Z-Score Normalization



$Q = (Q - mean(Q)) / std(Q)$
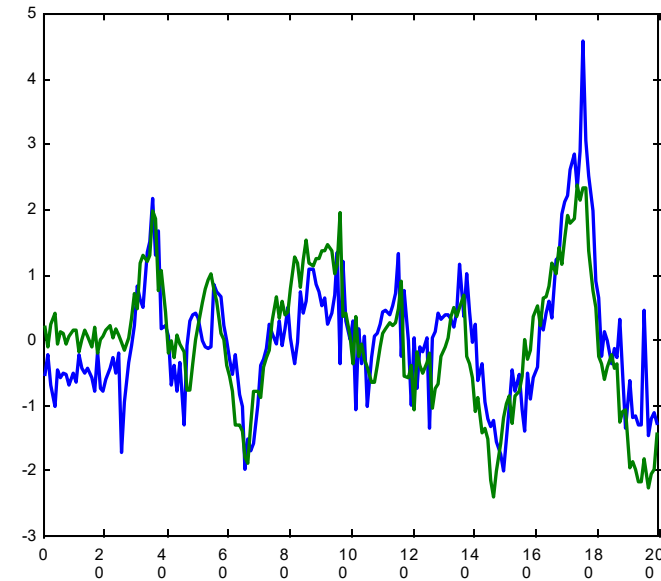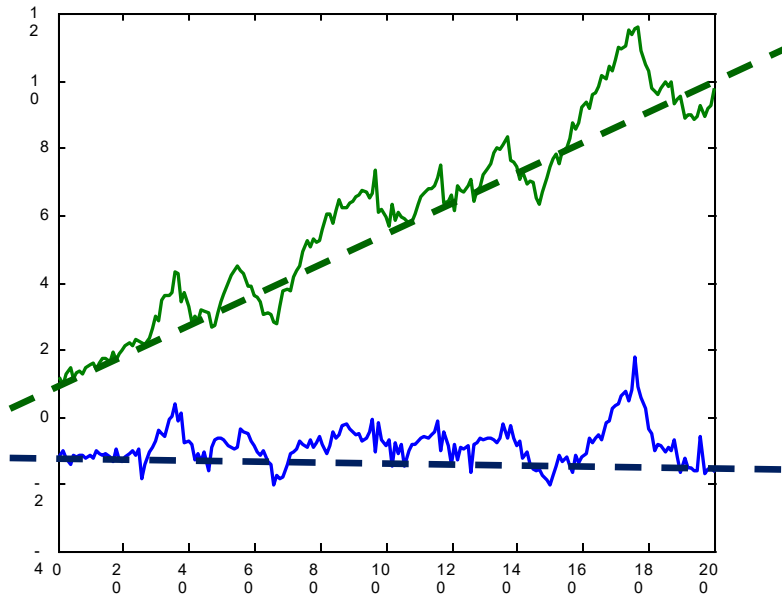
$C = (C - mean(C)) / std(C)$
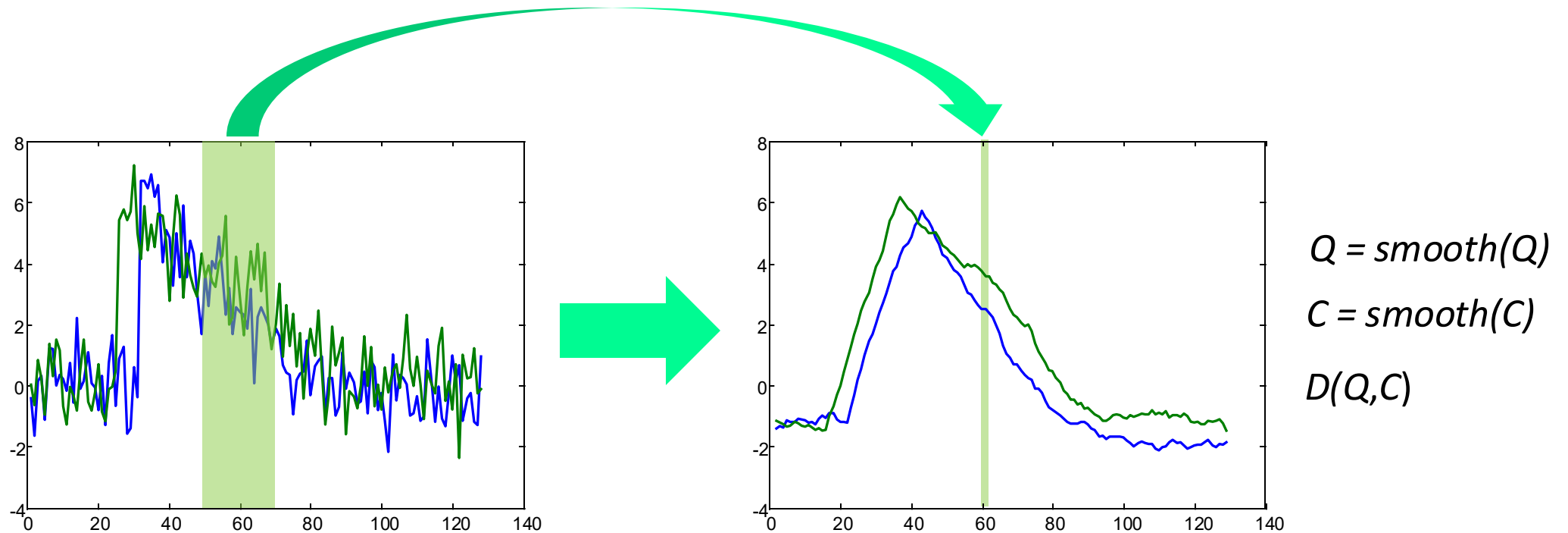
# Linear Trend: Detrending

- Fit the best fitting straight line to the time series, then subtract that line from the time series.



Removed linear trend, offset translation, amplitude scaling

# Noise Removal: Mean Smoothing

- The intuition behind removing noise is to average each datapoints value with its neighbors.



$Q = smooth(Q)$
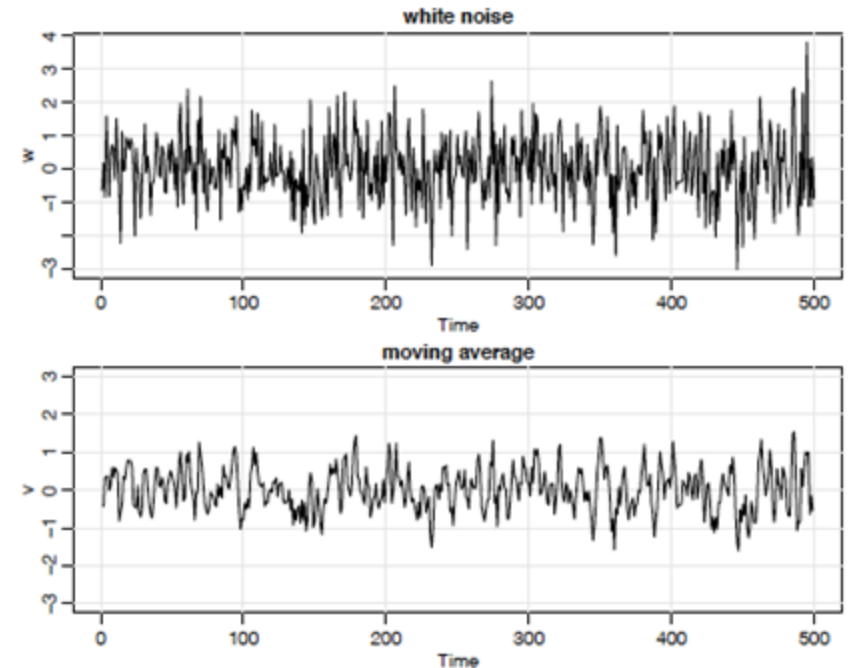
$C = smooth(C)$

$D(Q,C)$

# Moving Average

| time | value | | ma |
|------|-------|---|------|
| t1 | 20 | | - |
| t2 | 24 | | 22.0 |
| t3 | 22 | | 24.0 |
| t4 | 26 | | 24.3 |
| t5 | 25 | | - |

- Noise can be removed by a **moving average** (MA) that smooths the TS.

- Given a window of length $w$ and a TS $t$, the MA is applied as follows

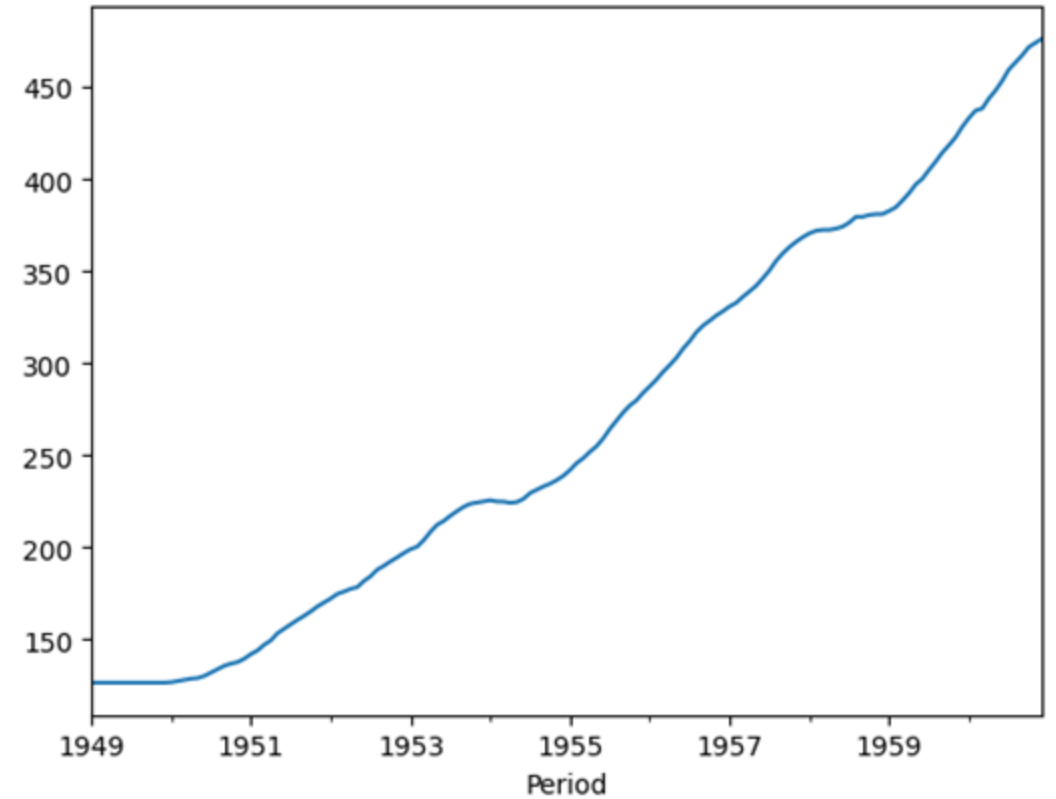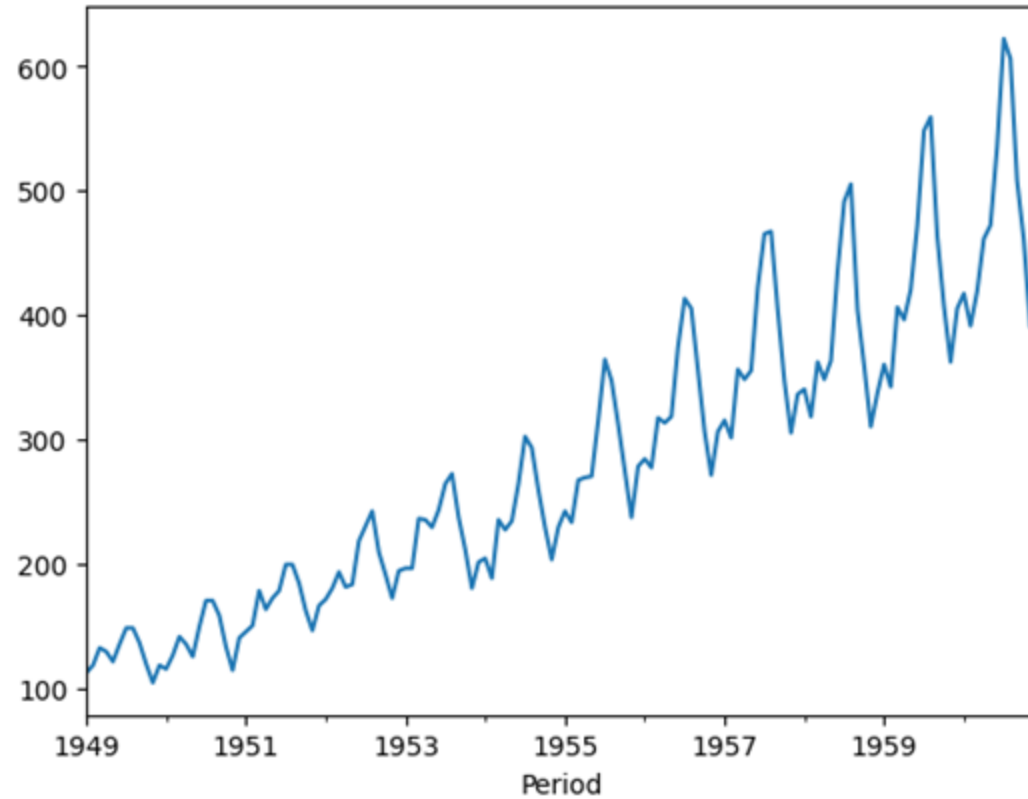- $t_i = \frac{1}{w}\sum_{j=i-w/2}^{w/2} t_j$ for $i = 1, \ldots, n$

- For example, if w=3 we have

- $t_i = \frac{1}{3}(t_{i-1} + t_i + t_{i+1})$



white noise



moving average

69

# Moving Average Example

# Log Transformation
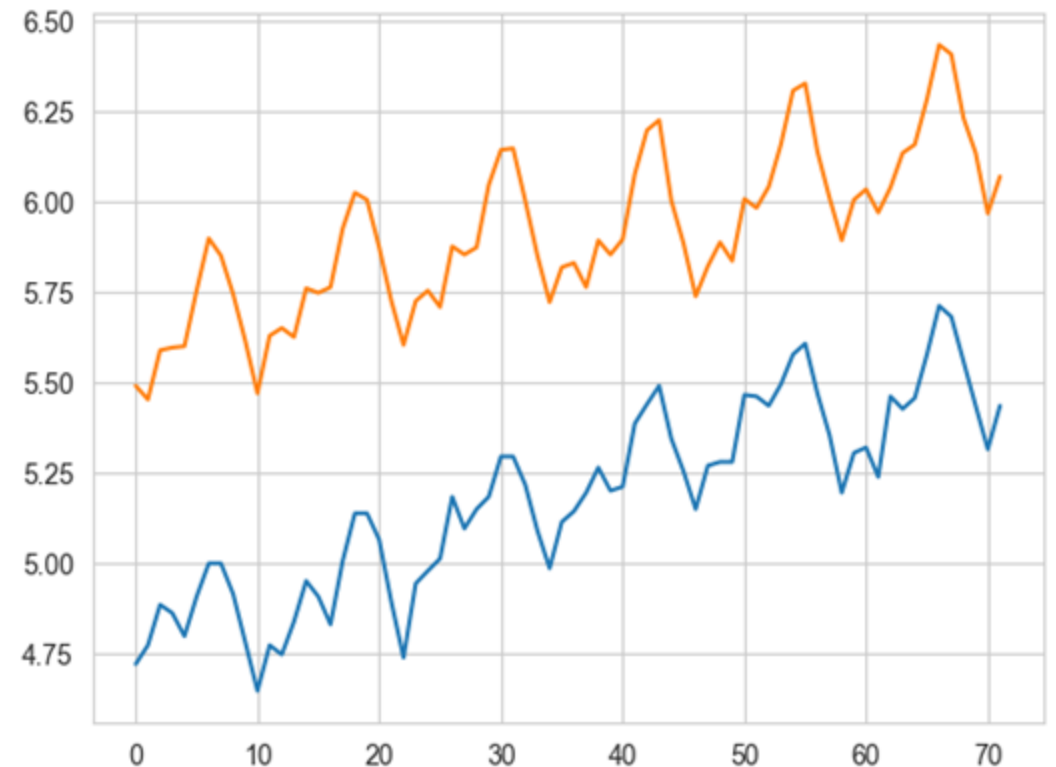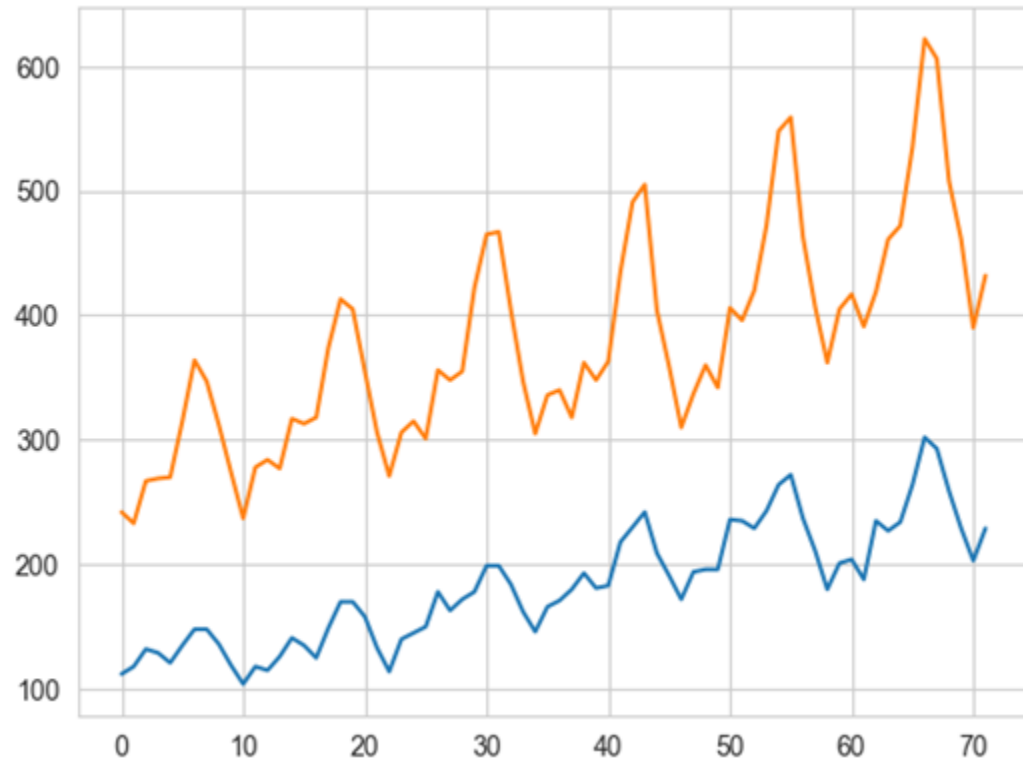
- We apply the logarithm to each value of the TS.

*Log*
$Q = log(Q)$
$C = log(C)$

*Log1p*
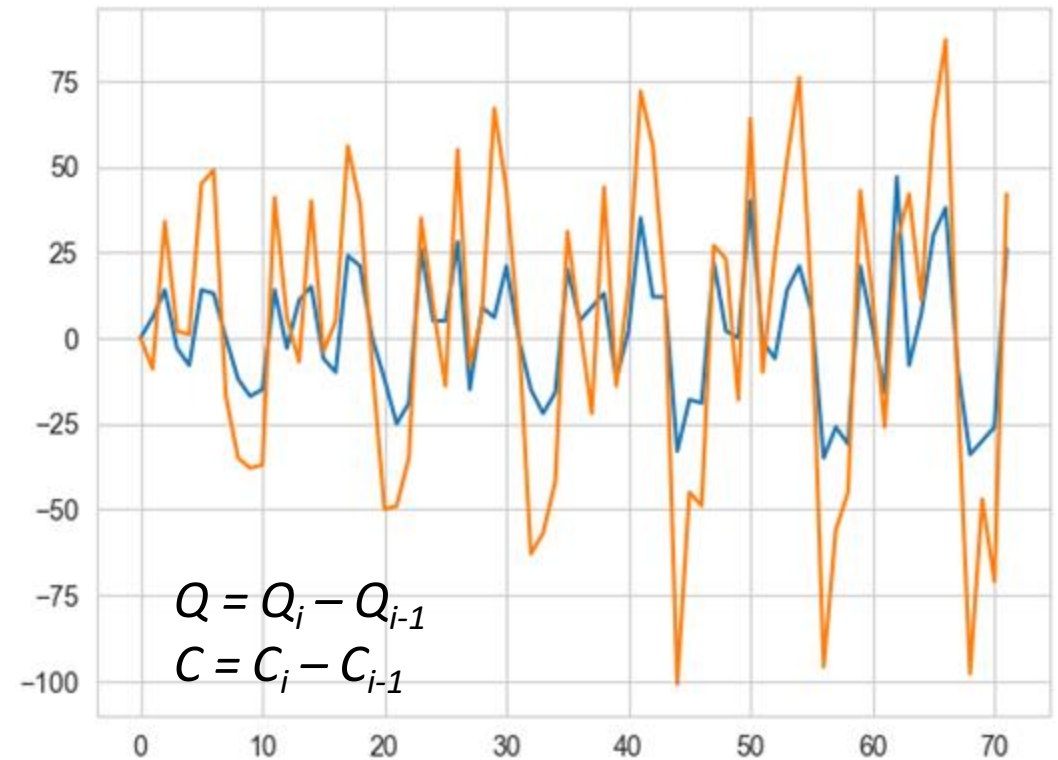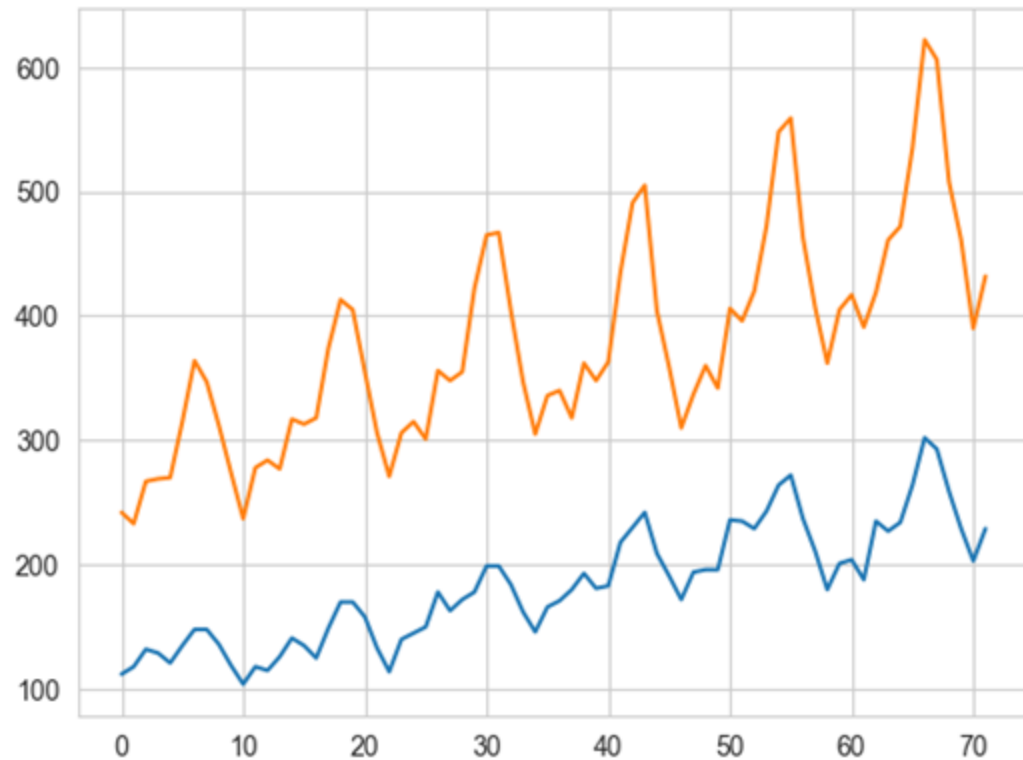$Q = log(Q + 1)$
$C = log(C + 1)$

# Differencing Transformation

- **Differencing**: we take the difference of the observation at a particular instant with that at the previous instant.
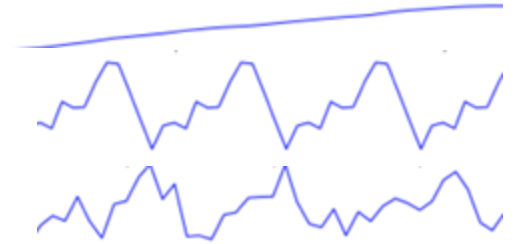


$Q = Q_i - Q_{i-1}$
$C = C_i - C_{i-1}$

# Time Series Components

# Assumption

- Various Data Mining, Machine Learning and TSA models assume that variables are Independent and Identically Distributed (IID)
  - In times series values are (usually) not independent
  - Trend and seasonality might be present
  - Variance may change significantly (heteroskedasticity)
- A first goal in TSA is to reduce the time series to a simpler case:
  - Eliminate trend
  - Eliminate seasonality
  - Eliminate heteroskedasticity
- Then we model the remainder as dependent but identically distributed variables.

# Time Series Components

- A given TS consists of three systematic components including level, trend, seasonality, and one non-systematic component called noise.
    - **Level**: The average value in the series.
    - **Trend**: The increasing or decreasing value in the series.
    - **Seasonality**: The repeating short-term cycle in the series.
    - **Noise**: The random variation in the series.

- A **systematic** component have consistency or recurrence and can be described and modeled.

- A **Non-Systematic** component cannot be directly modeled.
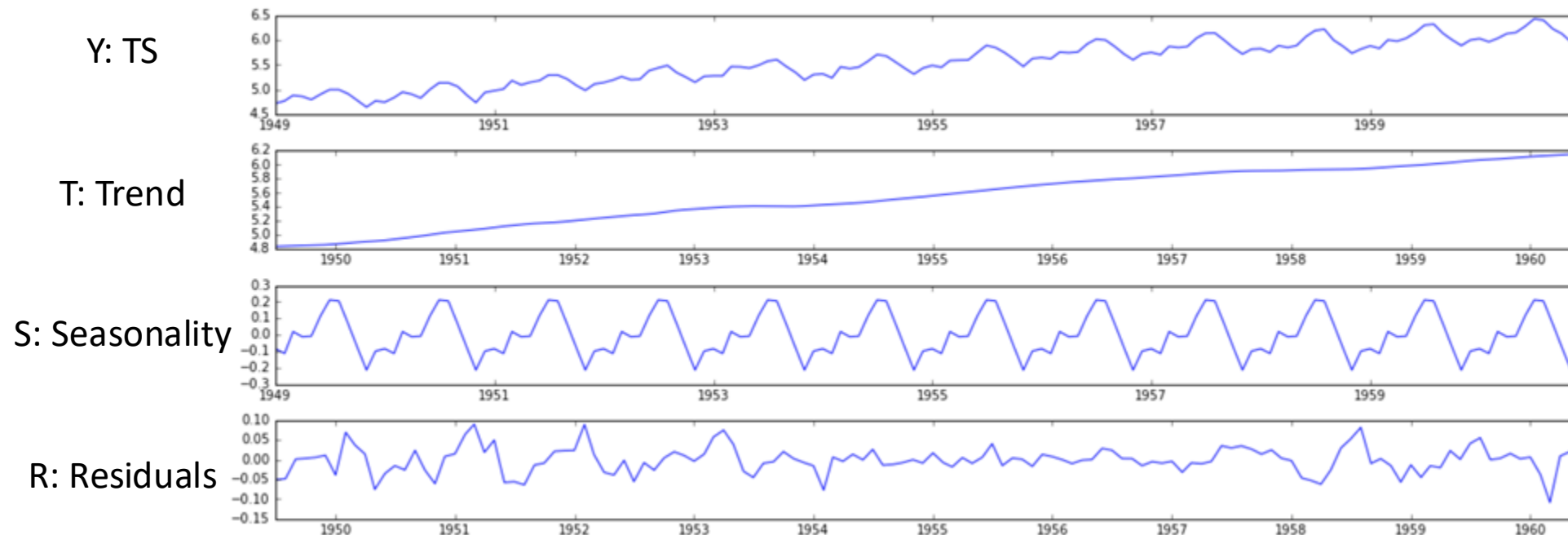
# Time Series Components

- A TS can be modeled as an aggregate or combination of these four components.
- All series have a level and noise. The trend and seasonality components are optional.
- Level can be omitted if Offset Translation with Mean Removal is applied

- **Additive Model**: Y = Level + Trend + Seasonality + Noise/Residuals
  - Changes over time are consistently made by the same amount
  - A linear trend is a straight line.
  - A linear seasonality has the same frequency (width of cycles) and amplitude (height of cycles).

- **Multiplicative Model**: Y = Level * Trend * Seasonality * Noise/Residuals
  - A multiplicative model is nonlinear, such as quadratic or exponential.
  - Changes increase or decrease over time.
  - A nonlinear trend is a curved line.
  - A non-linear seasonality has an increasing/decreasing frequency and/or amplitude over time.

# Time Series Components

Components meaning

- Trend: upward or downward pattern that might be extrapolated into the future.
- Seasonality: periodic behavior with a known period (hourly, weekly, monthly...).
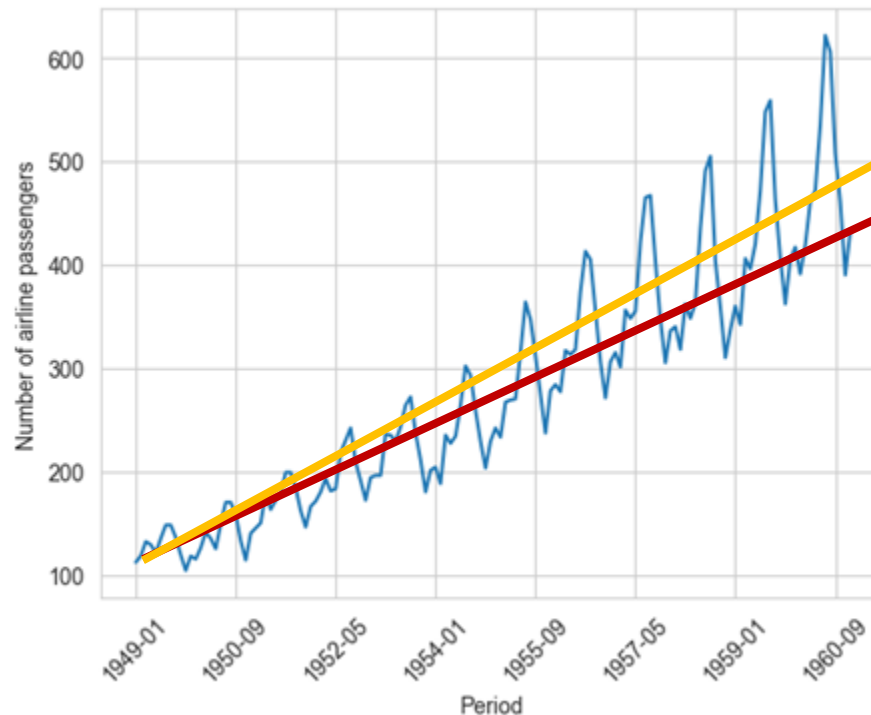- Residuals: containing anything else, i.e., noise.
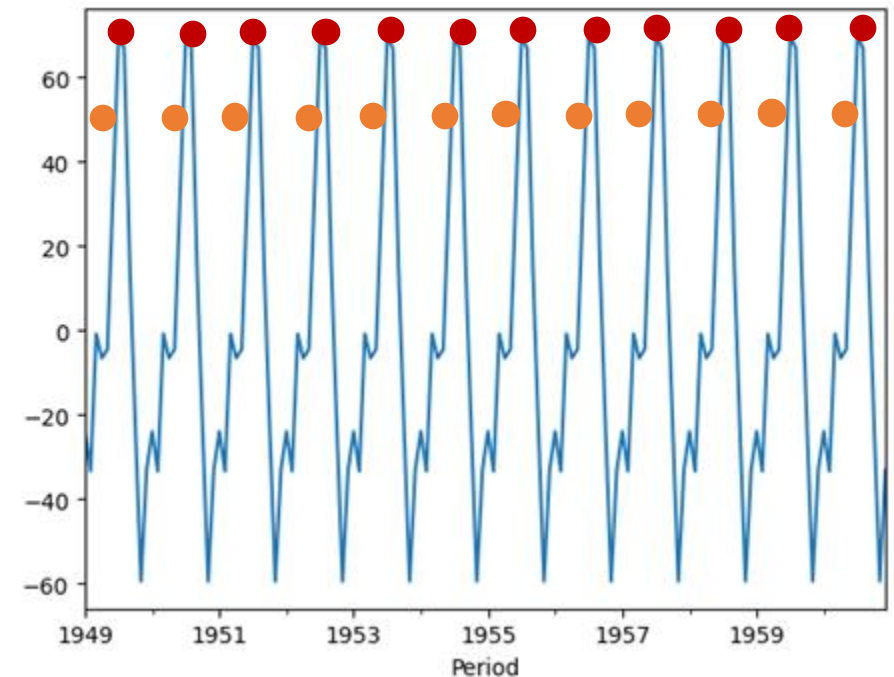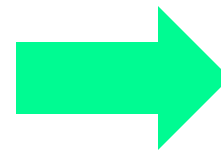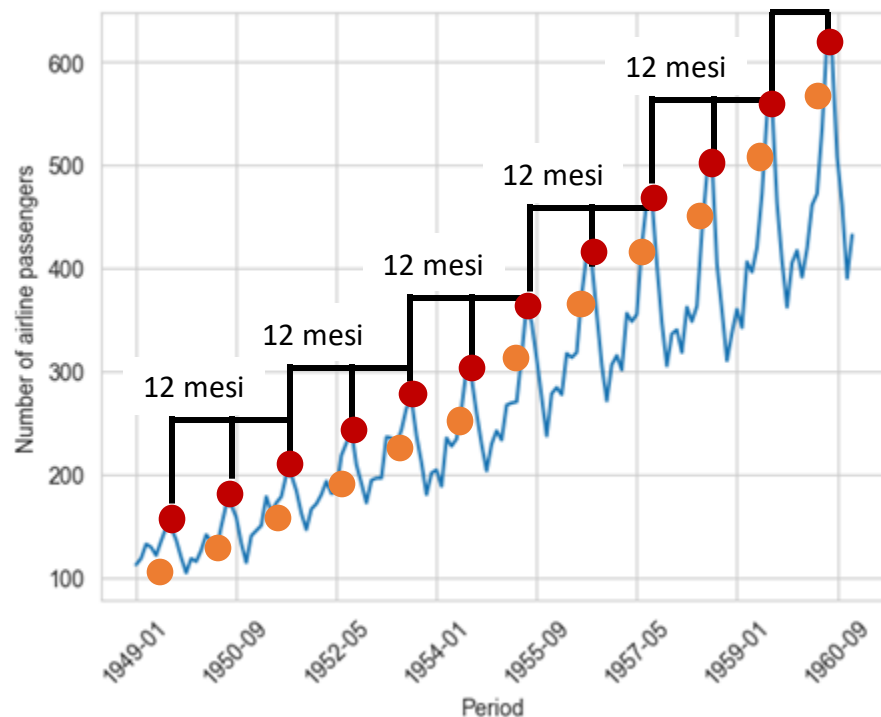


**Additive Model**

$$Y = T + S + R$$

# Time Series Decomposition: Trend

- Trend as a straight line between the first and the last point of the TS.
- Given a window *w*, trend as the moving average along TS with size *w*.
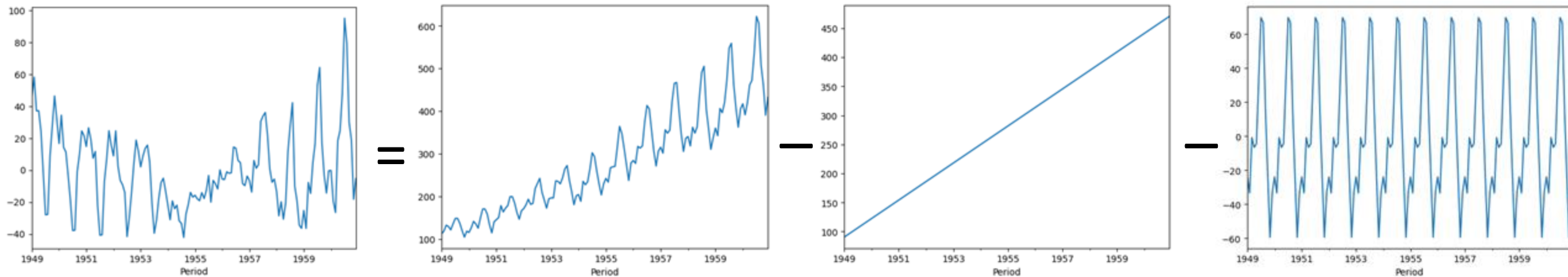- Trend as the best fitting straight line obtained with a linear regression.

# Time Series Decomposition: Seasonality

- Given the estimated number of periods $p$, and a detrended time series D = Y – T considering TS modeled with additive models, seasonality can be calculated as the mean value of for each time stamps among the various periods $p$ repeated $p$ times.
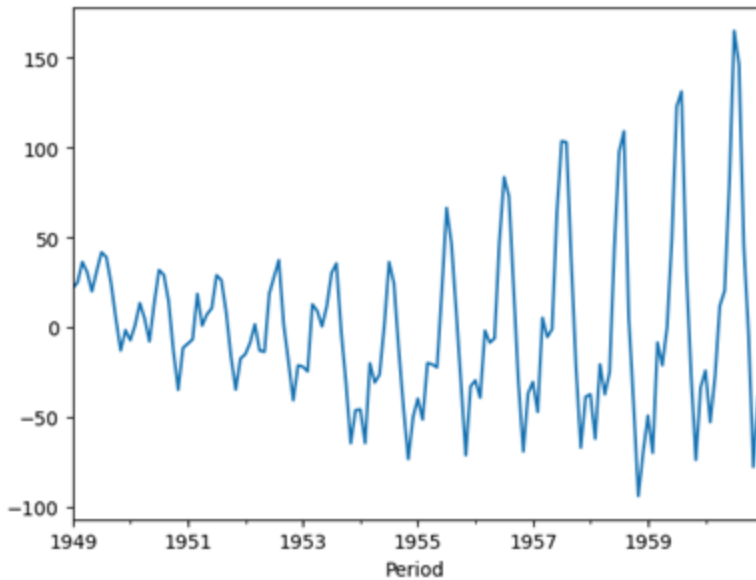
# Time Series Decomposition: Residuals

- Given the time series Y, its trend T, its seasonality S, the residuals R are obtained as R = Y − T − S.
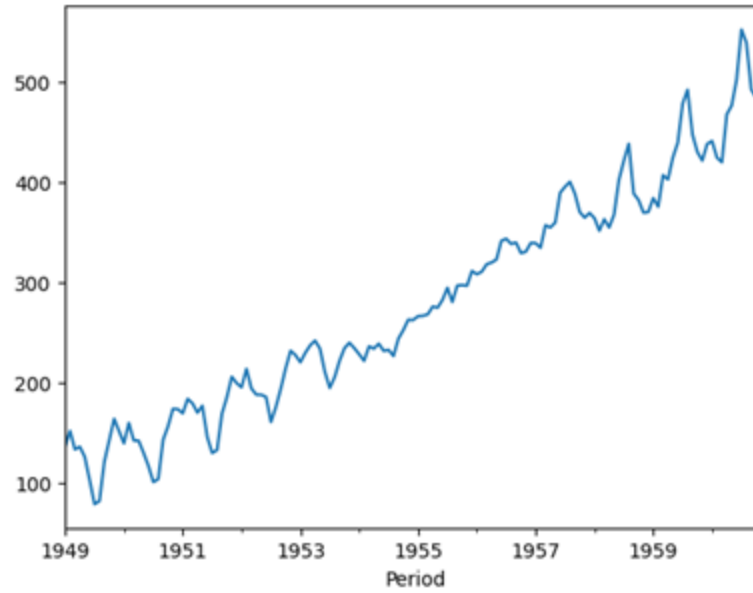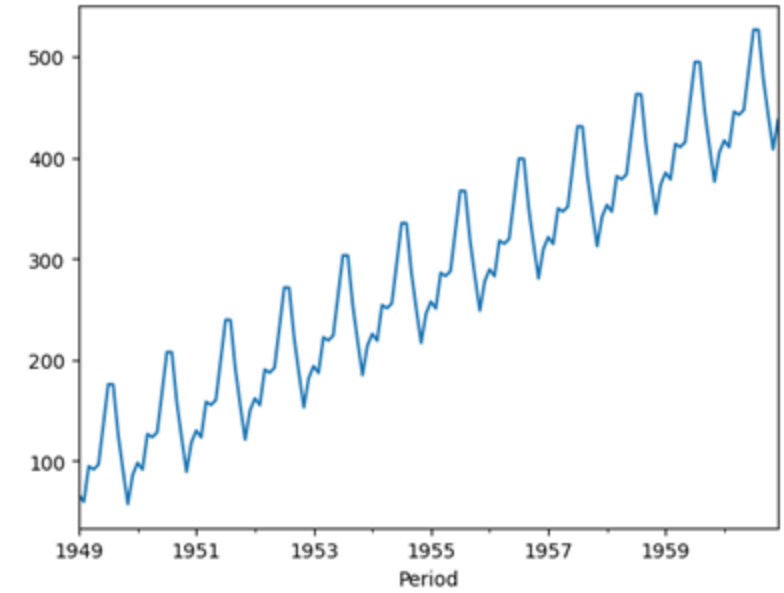
# Component-based Normalizations

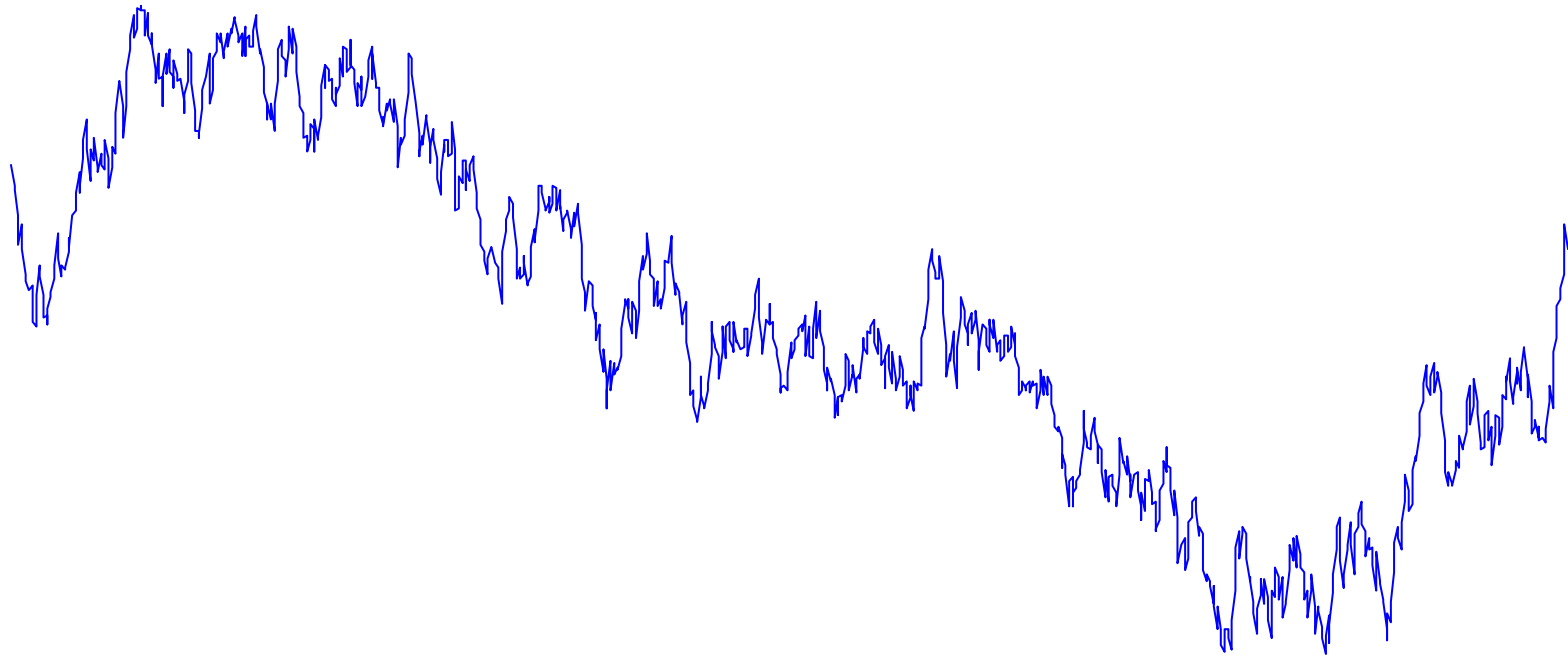Detrend
- Y' = Y − T

Deseasonalize
- Y' = Y − S

Denoise
- Y' = Y - R

# Time Series Preprocessing Remarks

- Depending on the TSA task different normalizations might be required but not necessarily all of them!!!

# References

- Forecasting: Principles and Practic. Rob J Hyndman and George Athanasaopoulus. (https://otexts.com/fpp2/)

- Time Series Analysis and Its Applications. Robert H. Shumway and David S. Stoffer. 4th edition.(http://www.stat.ucla.edu/~frederic/415/S23/tsa4.pdf)

- Time Series Analysis in R (https://s-ai-f.github.io/Time-Series/)

- Mining Time Series Data. Chotirat Ann Ratanamahatana et al. 2010. (https://www.researchgate.net/publication/227001229_Mining_Time_Series_Data)