

Analyse und Visualisierung von Zeitreihen- Daten in Python



Hmetrica



Analyse und Visualisierung von Zeitreihen- Daten in Python

Tag 1

Erwartungen

- Kurze Vorstellung
- Hintergrund, Vorerfahrungen Data Science, Programming
- Wie arbeite ich mit Zeitreihen?
- Was würde ich gern mitnehmen aus dem Kurs?
- Welcher Themenbereich interessiert mich am meisten?



Hmetrica GmbH

www.hmetrica.com |
info@hmetrica.com

Data Science Beratung

- Statistik, Machine Learning und Künstliche Intelligenz.
- Schwerpunkte: Datenvisualisierung, Machine Learning, Zeitreihenanalyse

Inhalte – Was haben wir vor?

Tag 1

1. Einführung in Zeitreihendaten in Python
2. Zeitreihen und ihre Merkmale visualisieren
3. Zeitreihen vorhersagen (Statistik I): Exponentielle Glättung und Holt-Winters
4. Zeitreihen vorhersagen (Statistik II): ARIMA-Modelle

Tag 2

5. Einblick in andere Zeitreihenmodelle
6. Machine Learning für Zeitreihen: Überblick, Vorbereitung und Klassifikation
7. Machine Learning für Zeitreihen: Clustering
8. Deep Learning für Zeitreihen (Einblick)

Fragen?



Einführung in Zeitreihendaten in Python

Session 1 (Montag 09:15 – 10:45)




Was haben wir vor?

1. Einführung in Zeitreihendaten in Python
 - 1.1 Einführung in Python: pandas, matplotlib
 - 1.2 Einführung in Zeitreihendaten: Definitionen, einfache Merkmale

Einführung in Python

Python ist eine allgegenwärtige Skriptsprache, benutzt für:

- Webentwicklung,
- App-Entwicklung,
- wissenschaftliche und numerische Bereiche,
- Geschäftsanwendungen,
- GUI-Design,
- Automatisierung,
- künstliche Intelligenz,
- maschinelles Lernen
- und vieles mehr...



Python Benefits:

- Back-end and front-end development
- cross-platform language
- open-source
- Strong community base
- Plethora of tools
- Fewer and simple lines of codes

Ab ins Jupyter
Notebook



Einführung in Zeitreihendaten



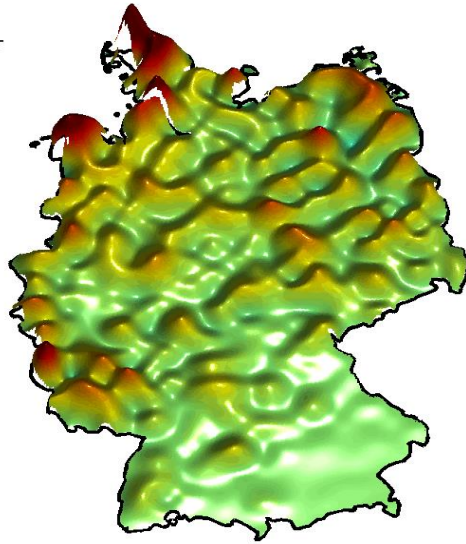
Beispiele für Zeitreihenanalysen



Bsp: Energiewende

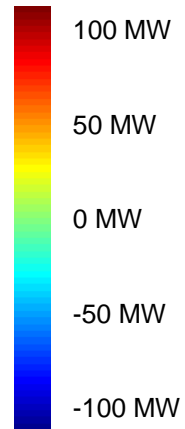
$$(PV) + (Wind) - (Consumption) = (Residual Load) = f(t)$$

Tag 130
00:00Uhr



**Most likely scenario
for 2025**

© IFHT



Was sind Zeitreihen?

- Umgangssprachlich: Daten, die sich auf aufeinanderfolgende Zeitpunkte (oder Zeiträume) beziehen
- Statistik betrachtet: Zeitreihe als “Realisation” eines “stochastischen Prozesses”

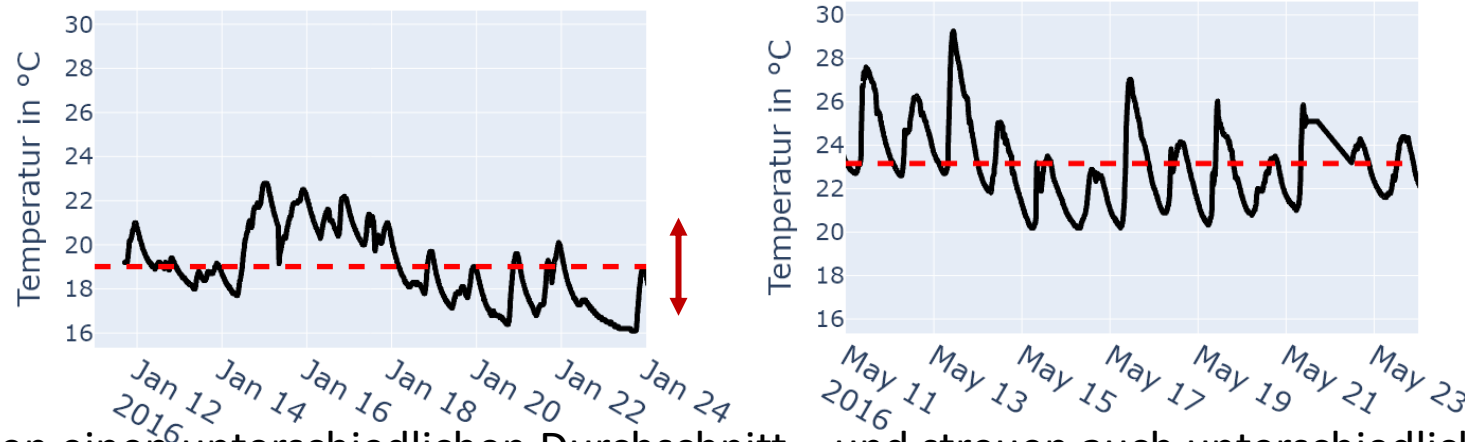
Stochastischer Prozess $\{X_t\}_{t \in T}$ ist eine Folge von Zufallsvariablen (Funktionen $X_t: \Omega \rightarrow \mathbb{R}$), die einen Index t aus einer Indexmenge T (meist Zeit gemessen als \mathbb{N}_0 oder \mathbb{R}_+) haben

Zeitreihen beschreiben

Beispiel

Sie beschreiben die Temperatur in einem Gebäude im Januar und Mai

Wie unterscheiden sich diese beiden Zeitreihen?



Sie haben einen unterschiedlichen Durchschnitt – und streuen auch unterschiedlich darum herum

Stochastischer Prozess X_1, X_2, X_3, \dots

Erwartungswert $\mu_t = \mathbb{E}[X_t]$

Varianz $\text{Var}(X_t) = \mathbb{E}[(X_t - \mu_t)^2] = \gamma_t(0)$

Autokovarianz

$\text{Cov}(X_{t+h}, X_t) = \mathbb{E}[(X_{t+h} - \mu_{t+h})(X_t - \mu_t)] = \gamma_t(h)$

Daten $x_1, x_2, x_3, \dots, x_t$

Mittelwert $\bar{x}_t = \hat{\mu}_t = \frac{1}{t} \sum_{i=1}^t x_i$

Stichprobenvarianz $\hat{\gamma}_t(0) = \frac{1}{t} \sum_{i=1}^t (x_i - \bar{x})^2$

Stichprobenautokovarianz

$\frac{1}{t} \sum_{i=1}^{t-|h|} (x_{i+|h|} - \bar{x})(x_i - \bar{x})$

$\hat{\gamma}_t(h) =$

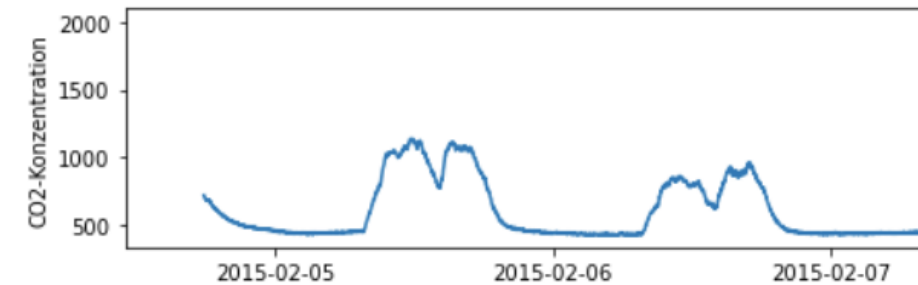
Hängt von
der Zeit ab!

Warum sind Zeitreihen besonders?

- Bei Zeitreihen:
- Beobachtung zu Zeitpunkt t : *könnte etwas zu tun haben* mit Beobachtungen zu anderem Zeitpunkt $t + h$
- Deshalb: **Kovarianz** der Zeitreihe mit sich selbst (zeitverzögert) wichtig!

Beispiel

- (A) Ich messe heute Mittag die CO²-Konzentration in 72 Büros in München (**keine Zeitreihe**)
- (B) Ich messe 72 Stunden lang die CO²-Konzentration in meinem München (**Zeitreihe**)



Stochastischer Prozess $X_1, X_2, X_3, \dots, X_n$

Daten $x_1, x_2, x_3, \dots, x_n$

Kovarianz

$$\text{Cov}(X_{t+h}, X_t) = \mathbb{E}[(X_{t+h} - \mu_{t+h})(X_t - \mu_t)] = \gamma_t(h)$$

Stichprobenkovarianz

$$\hat{\gamma}(h) = \frac{1}{n} \sum_{t=1}^{n-|h|} (x_{t+|h|} - \bar{x})(x_t - \bar{x})$$

Die **Autokorrelationsfunktion ACF** (engl. Auto Correlation Function) ist definiert als

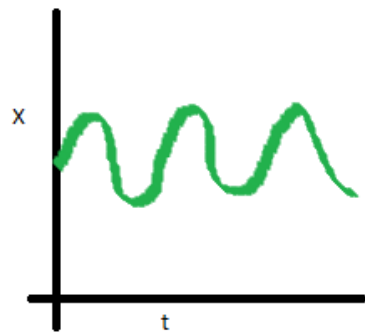
$$\rho_t(h) = \frac{\gamma_t(h)}{\gamma_t(0)}$$

Zeitreihen modellieren

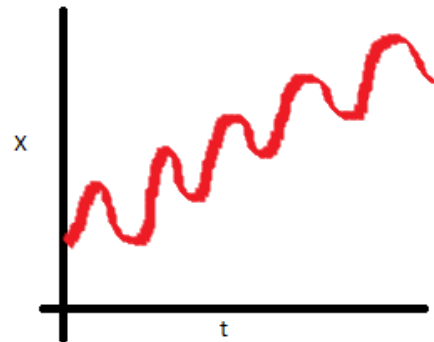
Eine Zeitreihe heißt **stationär**, wenn

$$\mu_t = \mu \quad \text{und} \quad \gamma_t(h) = \gamma(h)$$

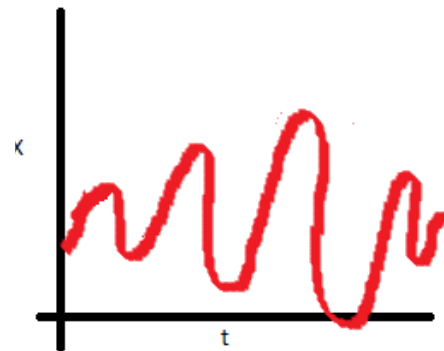
- Mittelwertfunktion und die Autokovarianzfunktion (und somit auch die Varianzfunktion) nicht abhängig von t sind.



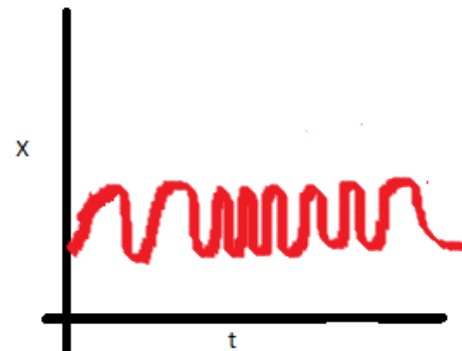
Stationär



Nicht stationär



Nicht stationär



Nicht stationär

Zeitreihen modellieren

Eine Zeitreihe heißt **stationär**, wenn

$$\mu_t = \mu \wedge \gamma_t(h) = \gamma(h)$$

- Mittelwertfunktion und die Autokovarianzfunktion (und somit auch die Varianzfunktion) nicht abhängig von t sind
- Wichtiger **Baustein**, um Zeitreihen zu modellieren:

Eine Zeitreihe $\{X_t\}_{t \in T}$ heißt **weißes Rauschen**, wenn alle X_t unabhängig identisch verteilt und

$$\mu_t = \mu \wedge \gamma_t(h) = \begin{cases} \sigma^2, & h = 0 \\ 0, & h \neq 0 \end{cases}$$

- Wir schreiben $\epsilon_t \sim IID(0, \sigma^2)$
- Sonderfall: Gaußsches weißes Rauschen $\epsilon_t \sim N(0, \sigma^2)$

Fragen?



Ab ins Jupyter
Notebook

