

Datenanalyse

Rohdaten von Sensoren können in den seltensten Fällen ohne Vorverarbeitung für maschinelles Lernen verwendet werden. In diesem Schritt wird daher, in Form einer umfassenden Datenanalyse, diese wichtige Vorarbeit behandelt.

Aufgabe 1: Daten laden

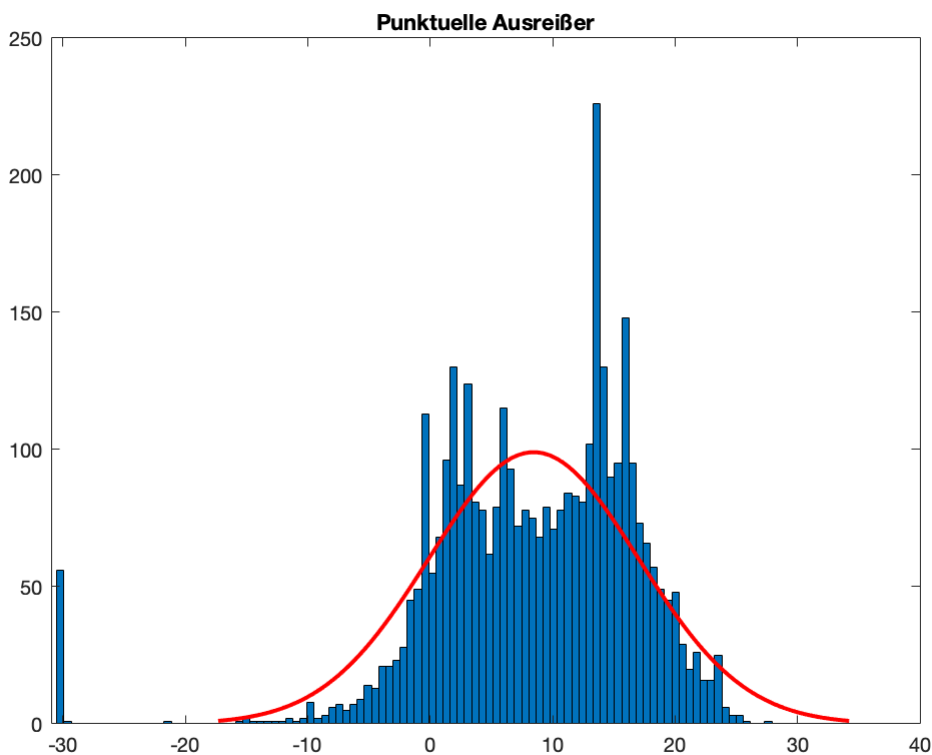
Zuerst wird die Datei `aussentemperatur_ucb.csv` mit den Daten mit Hilfe der Funktion `readtable` geladen.

```
clearvars;  
clc;  
  
% to do: Daten laden
```

Aufgabe 2: Punktueller Ausreißer - Histogramm Plot

Ein Histogramm Plot mit einer Normalverteilungskurve eignet sich gut um punktuelle Ausreißer in Daten zu finden. Hierfür wird die Funktion `histfit` verwendet. Punktueller Ausreißer können dann einfach an den Rändern des Histogramms gefunden werden.

```
% to do: Histogramm Plot mit Normalverteilung
```



Aufgabe 3: Kontextuelle Ausreißer - Zeitreihenzerlegung (Time Series Decomposition)

Für kontextuelle Ausreißer soll eine Zeitreihenzerlegung (Time Series Decomposition) durchgeführt werden. Bei dieser Methode wird jeder Wert eines Datensatzes in drei Komponenten (Trend T_t , Saisonwert S_t , Restfehler I_t) zerlegt, die entweder addiert (additive Methode) oder multipliziert (multiplikative Methode) den ursprünglichen Messwert ergeben. Da die multiplikative Methode besser zu Zeitreihen mit einem exponentiellen Wachstum passt, soll hier die additive Methode verwendet werden. Die Zerlegung ergibt sich also wie folgt: $y_t = T_t + S_t + I_t$

Bei der Bestimmung der Trendkomponente reicht es nicht nur einen Werte zu betrachten. Stattdessen müssen für jeden Wert auch mehrere Werte davor und danach betrachtet werden. Hier werden für jeden Wert immer auch die Werte der 15 Tage davor und danach betrachtet.

Die Trendkomponente ergibt sich wie folgt: $T_t = y \cdot p - 1 \text{ mal } \begin{bmatrix} \frac{1}{2p} \\ \frac{1}{p} \\ \vdots \\ \frac{1}{p} \\ \frac{1}{2p} \end{bmatrix}$

Hinweise:

- mit der Funktion `repmat(a, r, c)` kann eine Matrix der Größe r, c mit dem Element a erzeugt werden
- die Funktion `conv(u, v, 'same')` gibt den mittleren Teil der Faltung u, v mit der Größe u zurück

% to do: Trendkomponente berechnen

Die Saisonkomponente ergibt sich wie folgt: $S_t = \frac{1}{p} \cdot \sum_{i=t-\frac{p}{2}}^{t+\frac{p}{2}} T_i$

Hinweise:

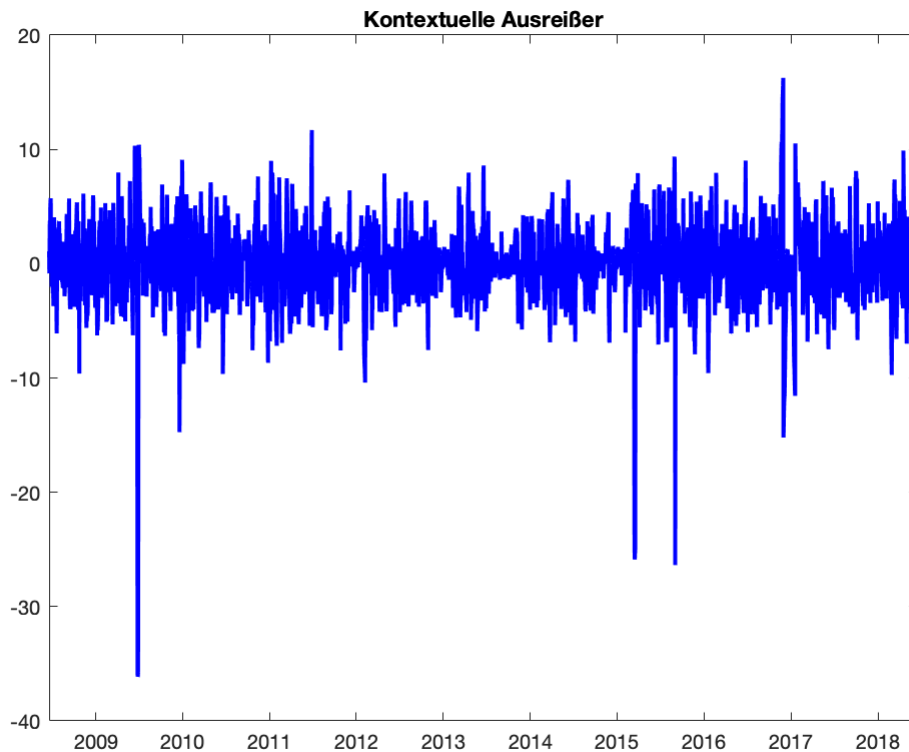
- den Tag eines Jahres liefert die Funktion `day(d, 'dayofyear')`
- mit der Funktion `cellfun(@x, f, c)` kann die Funktion f auf jedes Element von x angewandt werden welches die Bedingung c erfüllt

% to do: Trendkomponente berechnen

Anschließend wird Restfehler durch Subtraktion der Trend- sowie der Saison-Komponente von den Temperaturwerten bestimmt: $I_t = y_t - T_t - S_t$

Die kollektiven Ausreißer können dann am Plot der Restfehlerwerte über die Zeitachse abgelesen werden.

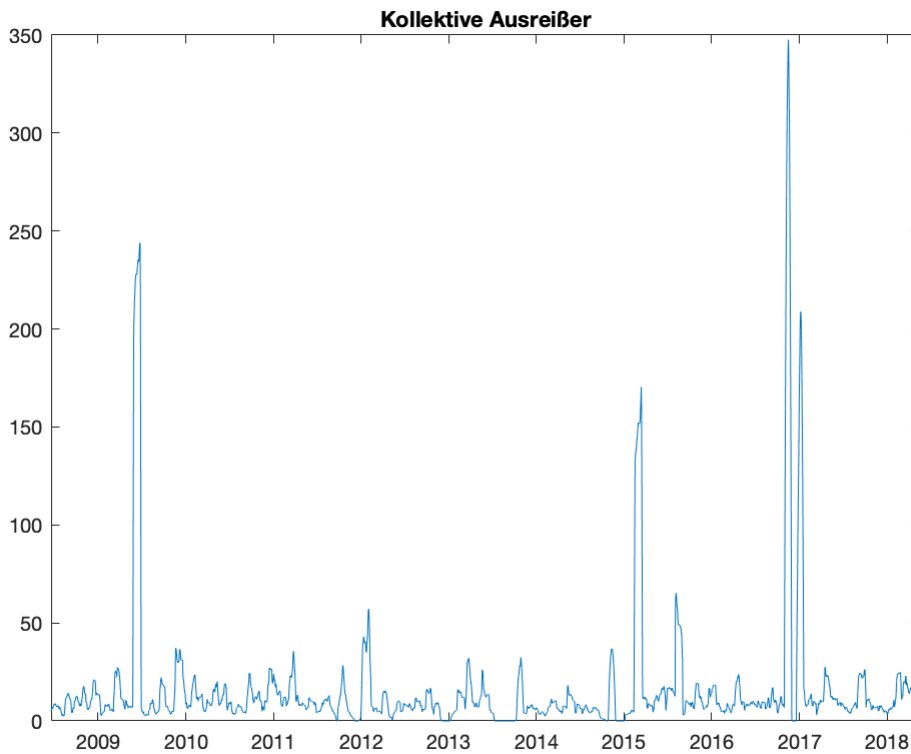
```
% to do: Restfehler plotten
```



Aufgabe 4: Kollektive Ausreißer - Varianzanalyse

Kollektive Ausreißer können durch die Ermittlung der Varianz in einem festen Zeitfenster, welches schrittweise über die Daten bewegt wird, gefunden werden. Als Fenstergröße werden 31 Tage gewählt. Die Varianz kann mittels der Funktion `var` bestimmt werden. Anschließend können kollektive Ausreißer an einem Plot der Varianzwerte über die Zeitachse abgelesen werden.

```
% to do: Varianzwerte berechnen und plotten
```



Aufgabe 5: Ausreißer bereinigen und plotten

Zum Schluss sollen die zuvor gefundenen Ausreißer aus den Messwerten entfernt werden. Dann sollen die bereinigten Messwerte geplottet und mit den zu Beginn geplotteten Messwerten verglichen werden.

Hinweise:

- Zum Löschen sollen die zuvor gespeicherten logischen Indizes verwendet werden.
- Außerdem müssen entsprechend der Fenstergröße der Varianzanalyse Daten am Ende entfernt werden.
- Um Lücken im Plot zu erzeugen kann die Funktion `cellfun(timetable, 'daily')` verwendet werden.

Konnten alle Ausreißer entfernt werden?

```
% to do: Messwerte mit Ausreißern plotten,  
% Ausreißer entfernen,  
% Messwerte ohne Ausreißer plotten
```