

Interdisziplinäre und Machine Learning -Grundlagen

Meilenstein 1

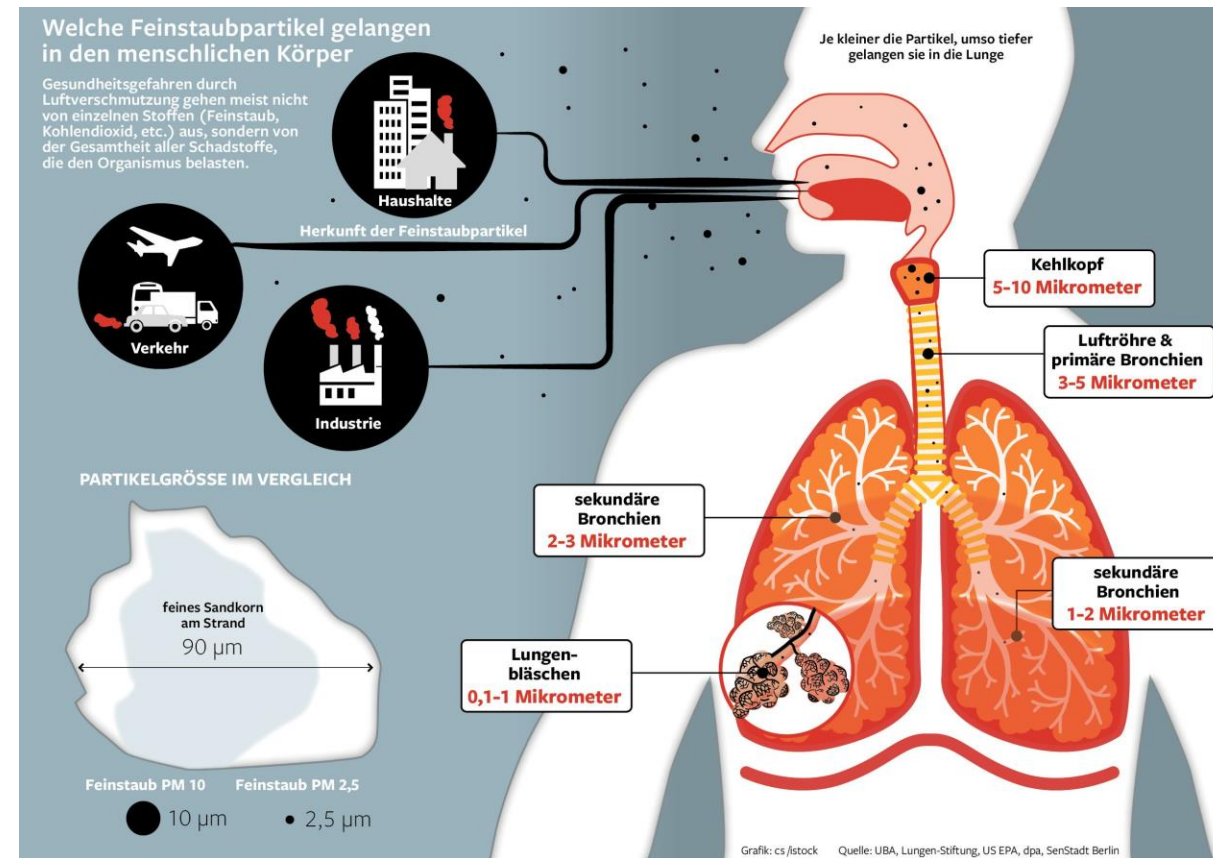
Inhalt

- Motivation
 - Feinstaub
 - Untersuchungsgebiet Polen
 - Messdaten
- Related Work
- Forschungsfrage
- Methodik
 - CNN
 - RNN
 - Modell

Motivation

Feinstaub

- Feinstaub wird nach Größe unterteilt
 - PM_{10} bezeichnet Partikel mit einem Durchmesser $< 10\ \mu m$
 - $PM_{2.5}$ ist eine Teilmenge von PM_{10} mit Partikel von einem Durchmesser $< 2.5\ \mu m$
- Grenzwerte

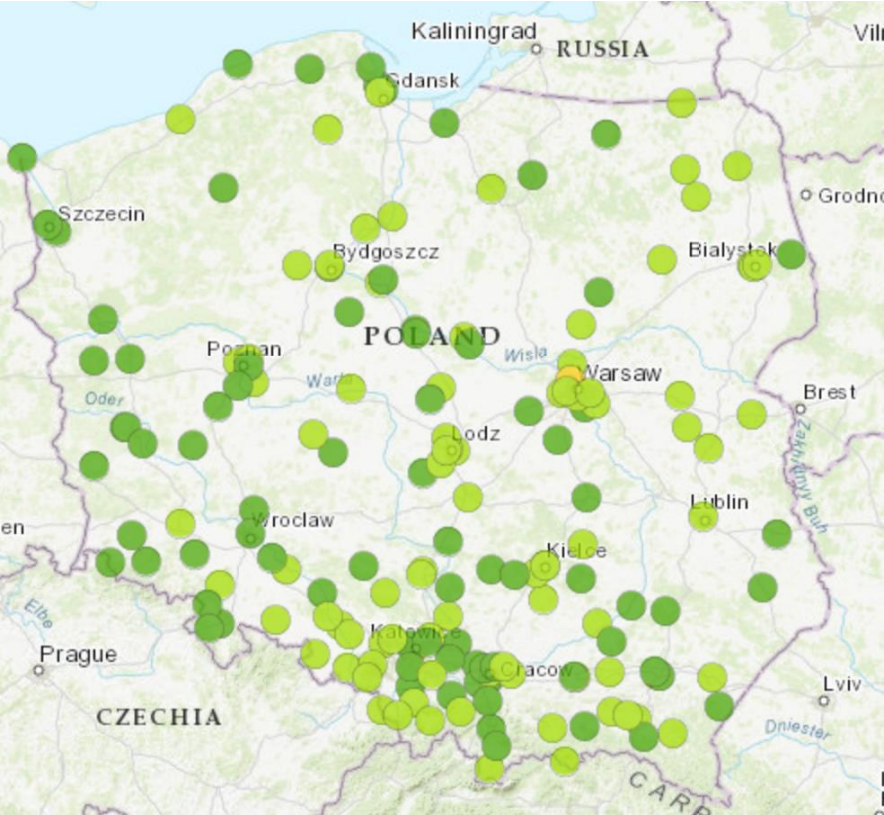


Feinstaub	EU	WHO	Mitteilungszeitraum
PM_{10}	$40\ \mu m^3$ $50\ \mu m^3$ 35 Tage/Jahr	$15\ \mu m^3$ $45\ \mu m^3$ 3-4 Tage/Jahr	1 Jahr 24 Stunden Erlaubte Überschreitung
$PM_{2.5}$	$25\ \mu m^3$ - -	$5\ \mu m^3$ $15\ \mu m^3$ 3-4 Tage/Jahr	1 Jahr 24 Stunden Erlaubte Überschreitung

Motivation

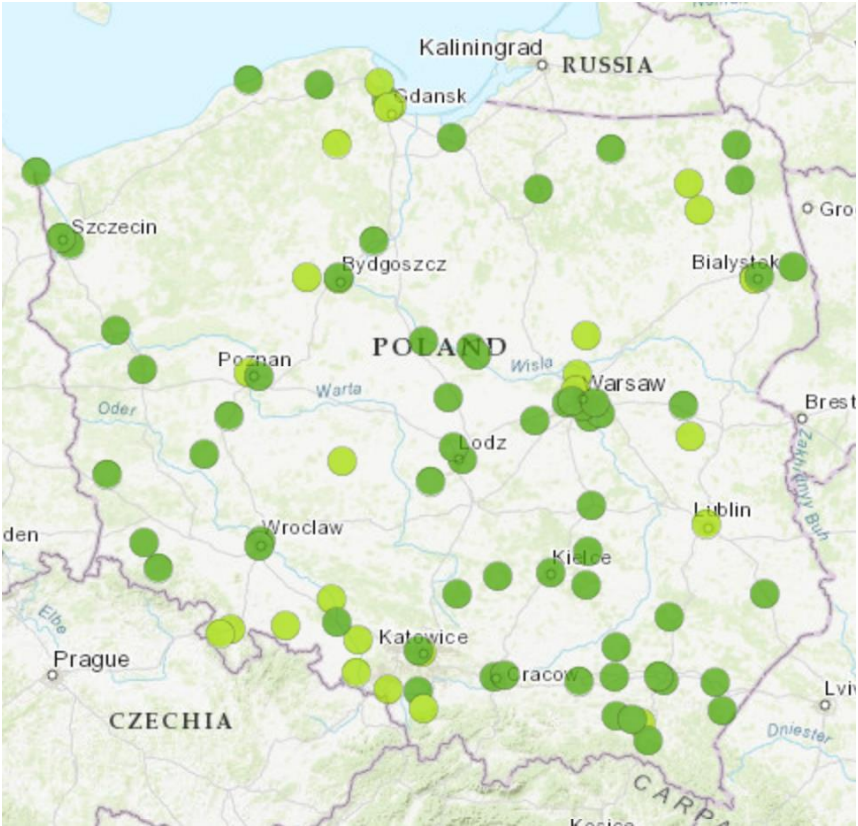
Untersuchungsgebiet Polen

PM₁₀



0 - 20 µg/m³	Very good
20.1 - 50 µg/m³	Good
50.1 - 80 µg/m³	Moderate
80.1 - 110 µg/m³	Sufficient
110.1 - 150 µg/m³	Bad
> 150 µg/m³	Very bad

PM_{2.5}



0 - 13 µg/m³	Very good
13.1 - 35 µg/m³	Good
35.1 - 55 µg/m³	Moderate
55.1 - 75 µg/m³	Sufficient
75.1 - 110 µg/m³	Bad
> 110 µg/m³	Very bad

Motivation

Wetterdaten

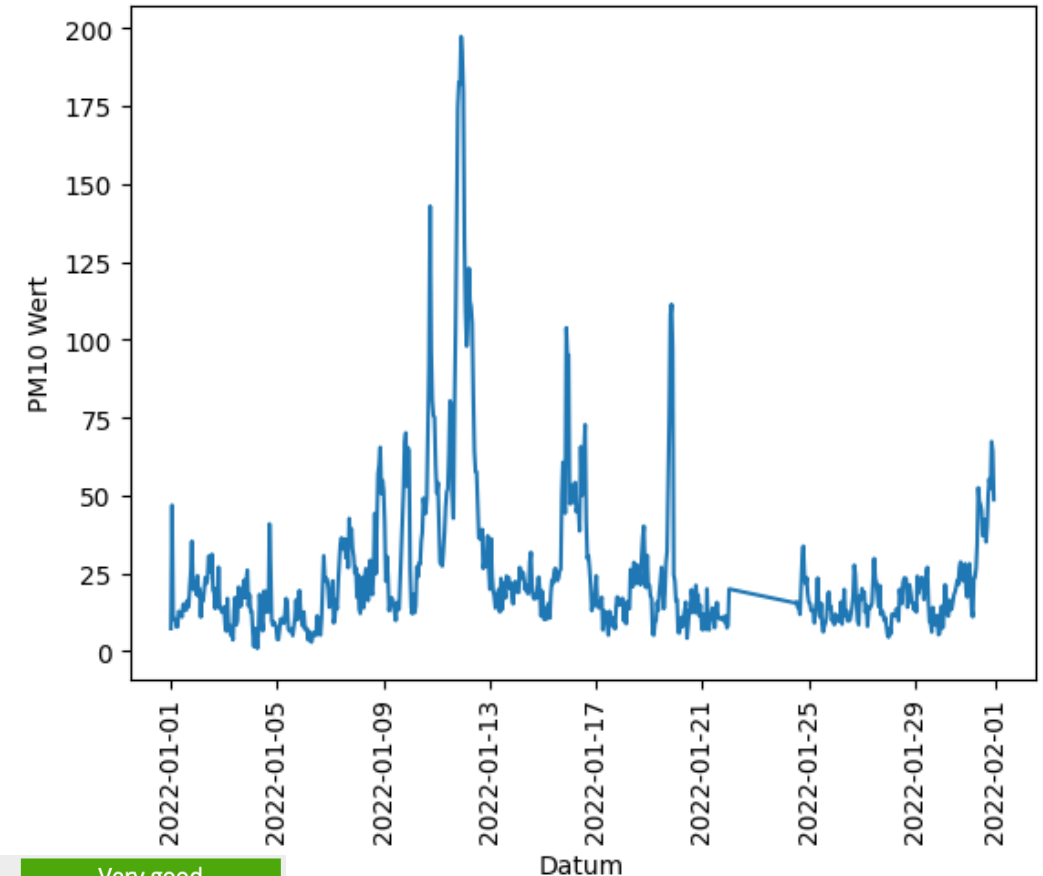
- Temperatur 
- Wind/ Luftströmung 
- Niederschlag 
- Luftfeuchtigkeit 

Inversionswetter

- Die oberen Luftschichten sind wärmer als die unteren
- Die kalte Luftschicht ist schwer und durchmischt sich nicht mit der oberen
- Die obere Luftschicht sperrt den Feinstaub wie eine Käseglocke ein

Beispiel Messdaten PM₁₀

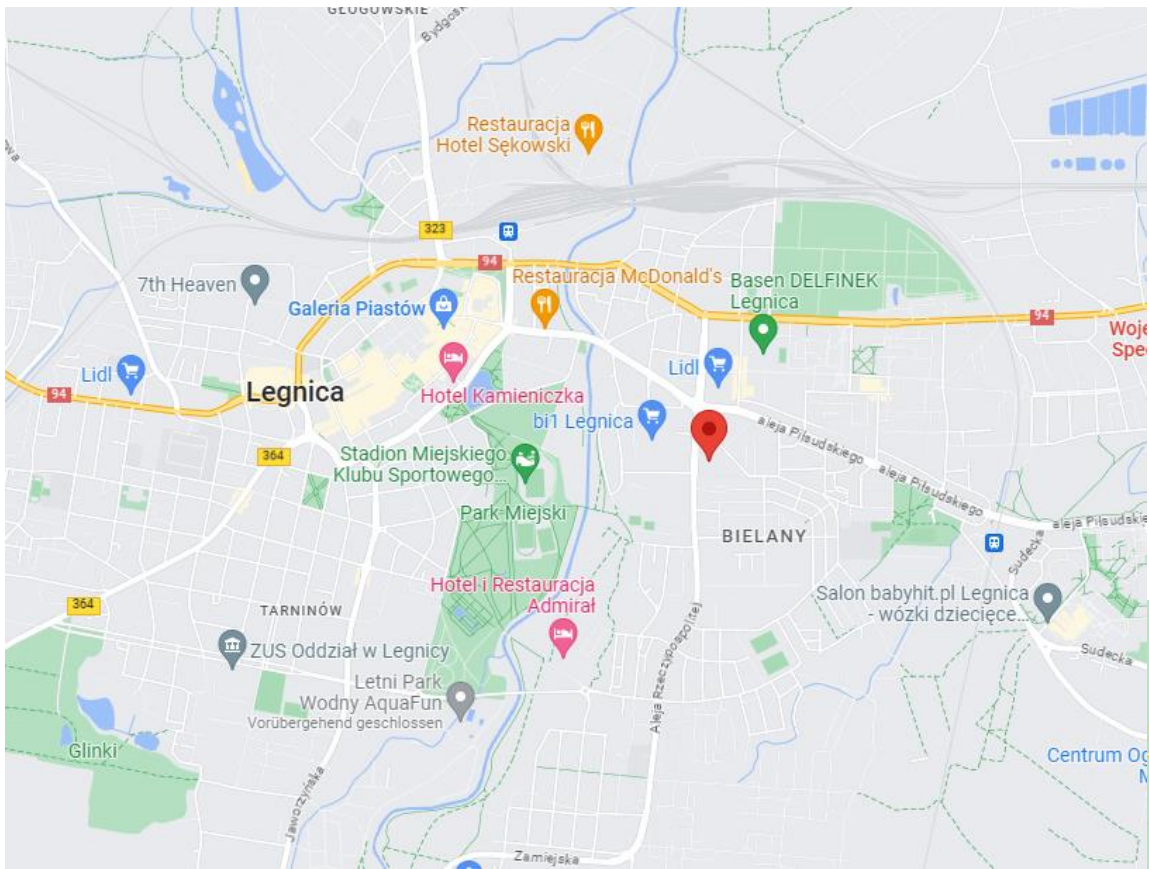
Station 52 SensorID 14397



0 - 20 µg/m ³	Very good
20.1 - 50 µg/m ³	Good
50.1 - 80 µg/m ³	Moderate
80.1 - 110 µg/m ³	Sufficient
110.1 - 150 µg/m ³	Bad
> 150 µg/m ³	Very bad

Beispiel Messdaten PM₁₀

The chart displays the daily count of new COVID-19 cases in Berlin. The y-axis is labeled with values 0, 25, 50, 75, and 100. The x-axis is labeled 'Datum' and shows dates from 2022-01-01 to 2022-02-01. The data is represented by a blue line that fluctuates significantly. Notable peaks occur around January 11 (near 100 cases), January 16 (near 100 cases), and January 20 (near 100 cases). There are also smaller peaks around January 9 and January 29. The overall trend shows high volatility with frequent daily changes in case counts.



Related Work

Daten

Revisiting air quality forecasting: a regression approach (2018)

- Vorhersage PM_{10} für nächsten Tag pro Messstation
- Daten:
 - PM_{10}
 - Temperatur
 - Luftfeuchtigkeit
 - Tag
 - Monat

Air-pollution prediction in smart city, deep learning approach (2021)

- Stündliche Vorhersage $PM_{2.5}$
- Daten:
 - PM_{10}
 - $PM_{2.5}$
 - Temperatur
 - Luftdruck
 - Regen
 - Windrichtung

Modelle

Revisiting air quality forecasting: a regression approach (2018)

- Regression
- Ein Modell pro Messstation
- Vergleich von Modellen
 - Lineare Regression
 - ANN (1 Hidden Layer (10 Neuronen))
 - Random Forest
 - 1.LR, 2. RF, 3. ANN
 - ANN recht klein

Air-pollution prediction in smart city, deep learning approach (2021)

- Regression
- Hybridmodell:
 - CNN (räumliche Merkmale) + LSTM (zeitliche Merkmale)
 - 3 Convolution Layer mit Batch Normalization
 - Maxpooling Layer
 - 2 LSTM Layer (100 und 50 Units)
 - 1 Dense Layer

Ziele

- Prognose von stündlichen Feinstaubdaten innerhalb der nächsten 14 Tage
 - Prädiktion von PM_{10}
 - optional $PM_{2.5}$
 - Aufteilen des Gebietes in Bereiche
 - Prädiktion pro Station
 - Prädiktion für ein Gebiet
 - Wie sehen unsere Prognosen im Vergleich mit denen des polnischen Umweltamts aus? (für einen Tag)
 - Besteht ein Zusammenhang zwischen PM_{10} und $PM_{2.5}$?

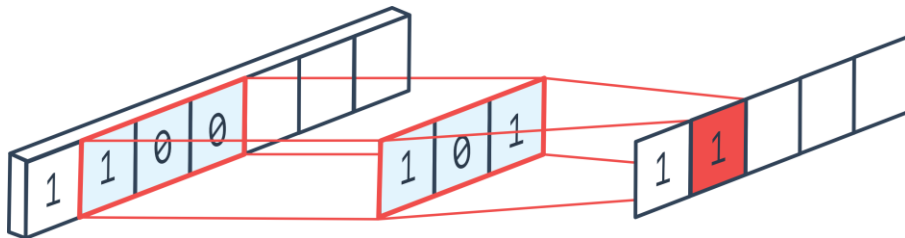
Forschungsfragen

- Lässt sich mit Hilfe eines neuronalen Netzes unter Verwendung einer CNN-LSTM Kombination eine stündliche Prognose von Feinstaubdaten für die nächsten 14 Tage realisieren?
 - ⑩ Ist es möglich den PM_{10} Wert mit einem MAE unter 10 vorherzusagen?
 - ⑩ Ist es möglich den $PM_{2.5}$ Wert mit einem MAE unter 10 vorherzusagen?
 - ⑩ Gibt es einen Zusammenhang zwischen PM_{10} und $PM_{2.5}$, sodass $PM_{2.5}$ mit dem Modell für PM_{10} vorhergesagt werden kann?
- Wie sehen unsere Prognosen im Vergleich mit denen des polnischen Umweltamts aus? (für einen Tag)
- Ist es sinnvoll, Stationen zu Gebieten zusammenzufassen, sodass die Aussagekräftigkeit der Prädiktion im Vergleich zu den einzelnen Stationen gleich bleibt oder verbessert wird?

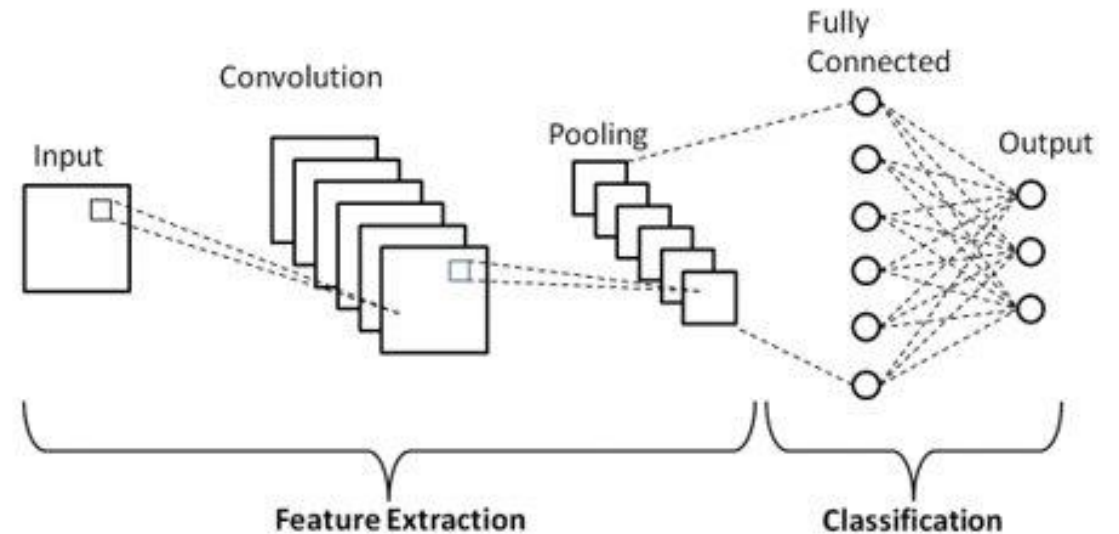
Methodik

Convolutional neuronale Netze (CNN)

- maschinelles Verarbeiten von Daten
- Beruht auf Faltungsoperation (Convolution)
- Convolution
 - Berechnung des neuen Feldes durch Maske
 - Filterung der Daten



Quelle: <https://ai.stackexchange.com/questions/28767/what-does-channel-mean-in-the-case-of-an-1d-convolution>



Quelle: <https://www.upgrad.com/blog/basic-cnn-architecture/>

Rekurrente neuronale Netze (RNN)

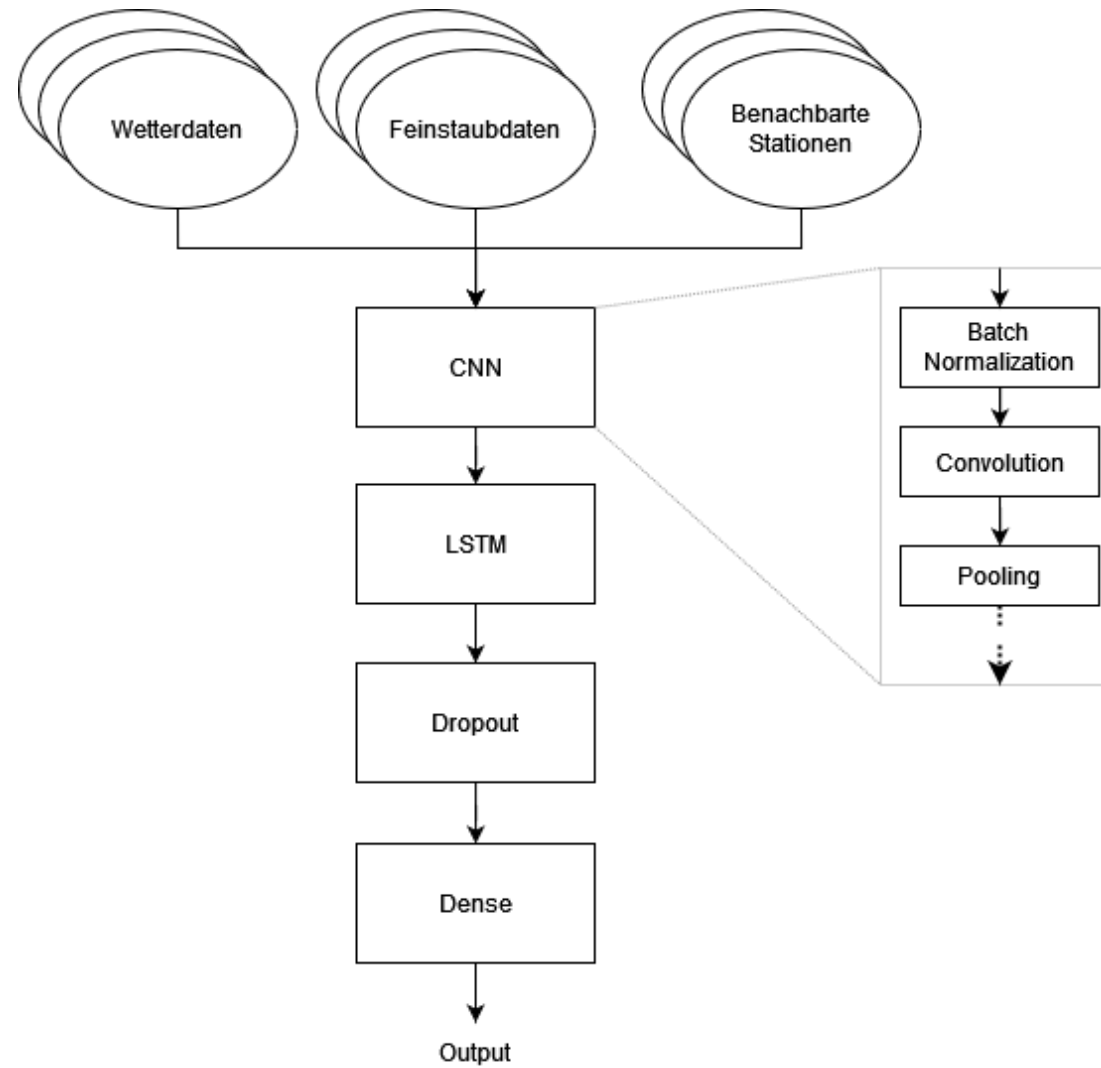
- Neuronale Netze mit Rückkopplungen
 - Erzeugung von Gehirn im Netz
- Besonders geeignet für zeitliche und sequentielle Daten
- Bekannteste Varianten GRU und LSTM
- LSTM (Long short-term memory network)
 - 1997 von Hochreiter & Schmidhuber entwickelt
 - Lernen von Langzeitabhängigkeiten in sequentiellen Daten
 - Verarbeitet Sequenz von Daten
 - Besteht aus Memory Cells
 - Fungieren als Gedächtnis → Speichern Informationen der vorangegangenen Daten
 - Besonders gut geeignet für Feinstaubprädiktion

Modell

- Supervised ML Problem
- CNN-LSTM Kombination zur Lösung des Regressionsproblems
 - Multi-Step Modell
- Multi-Step Forecasting (PM_{10})
 - Input: Merkmalsvektor
 - Output: PM_{10} Wert
- Modell mit MSE als Fehlerfunktion trainieren
- Bewertung des Modells mittels RMSE, MSE und MAE
- Einbeziehen von 3 benachbarten Stationen
 - Gewichtung nach Entfernung
- Trainieren und erstellen des Netzes mit Tensorflow
- Trainingsdaten und Testdaten jeweils 1 Jahr

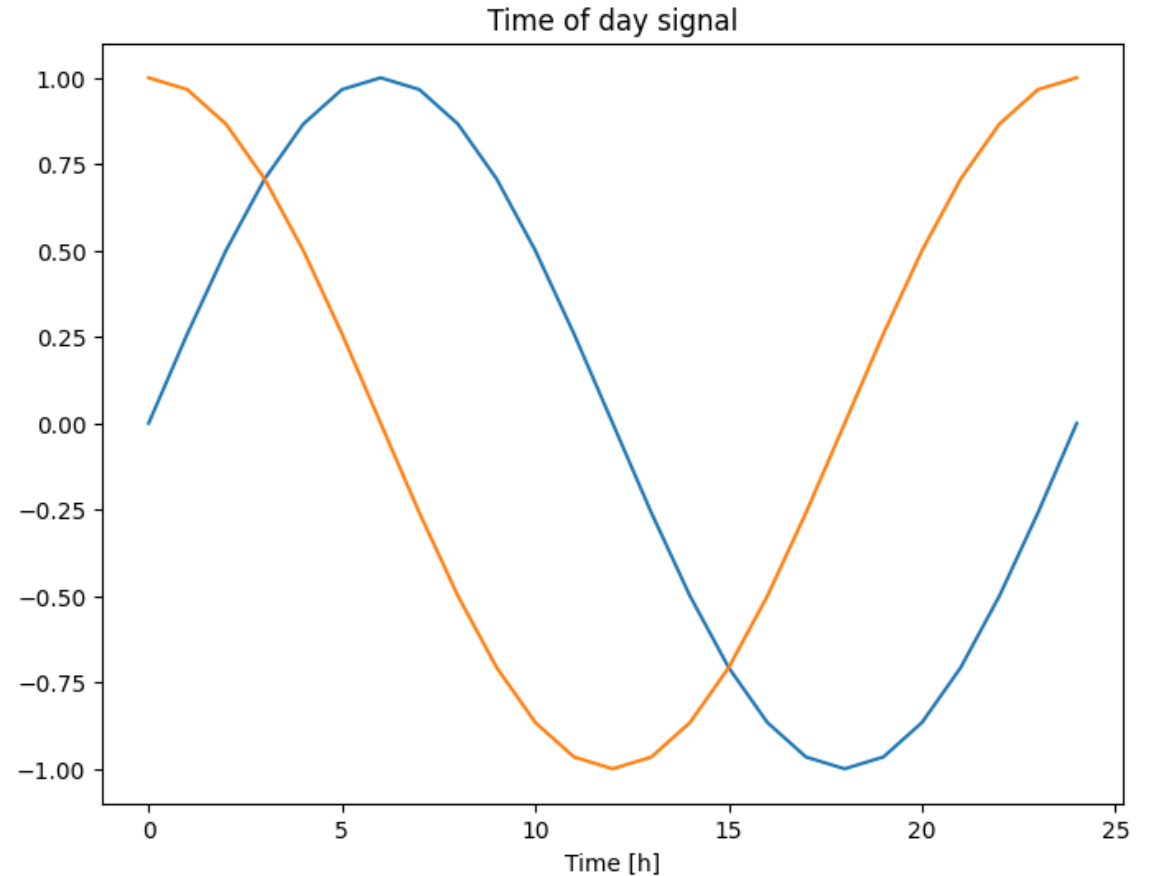
Methodik

Modell



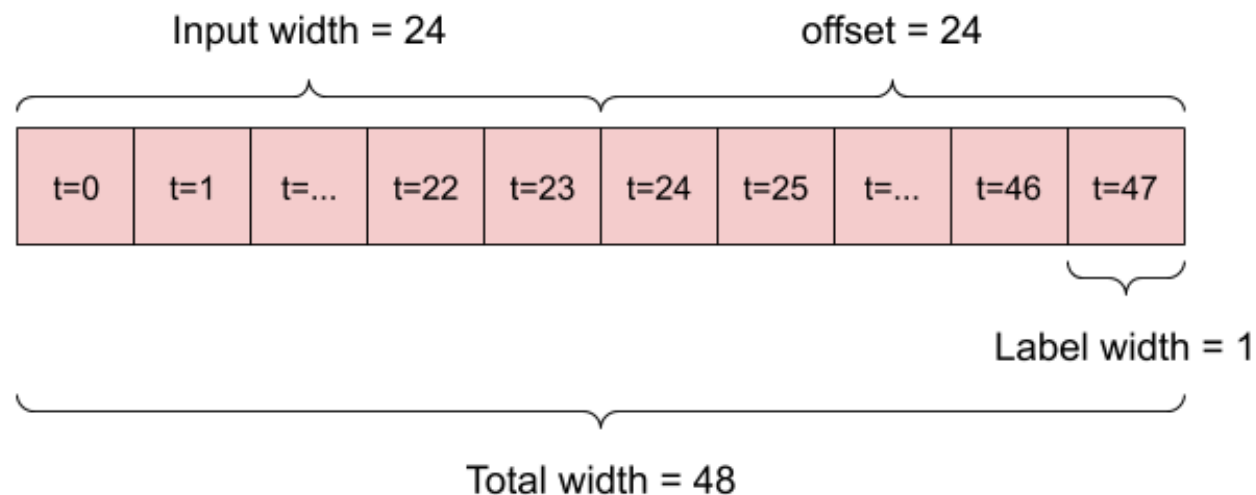
Feature Engineering - Zeit

- Zeit periodisch angeben statt absolut
- Perioden für Tag und Jahr
- Für Modell besser verwendbar
- Periode durch Sin und Cos
 - Eindeutigkeit der Werte durch Kombination
- Gleiches auch für Windrichtung relevant
 - Grad in Radiant umwandeln
 - Abstand zwischen 360° und 0° nicht vorhanden in Radiant



Data Windowing

- Definition des Offset der Vorhersage
 - Definition der Inputdaten-Zeitspanne
 - Beispiel:
 - Ergebnis, das 24 Stunden in der Zukunft liegt
 - Historie an 24 Stunden Daten
- In unserem Fall:
 - 14-tägige Frames



Multi-step model

1. Single shot predictions

- Gesamte Vorhersage durch einmalige prediction

2. Autoregressive predictions

- Single step prediction
- Ergebnis zur Bestimmung der nächsten prediction
- Verwendet bei Korrelationen der vorherigen Werte