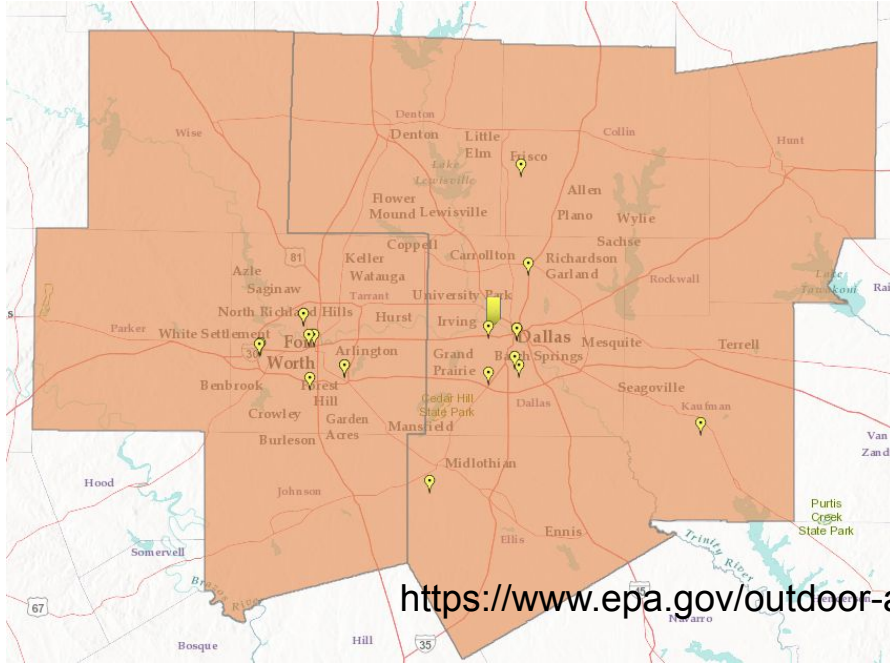


Air-Quality (PM2.5) prediction with Machine-learning models

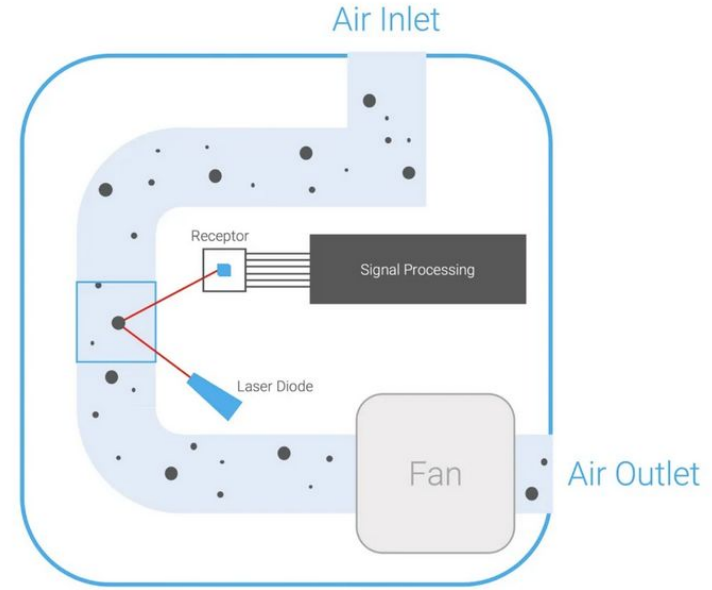
Prahlad Siwakoti
NSS Data Science Cohort 7

Overview

- ❑ Local air-quality variations are difficult to monitor.
- ❑ Inexpensive optical-scattering sensors offer a solution.
 - ❑ Granular insights from a denser array of sensors



<https://www.epa.gov/outdoor-air-quality-data/interactive-map-air-quality-monitors>



Laser sensors schematic

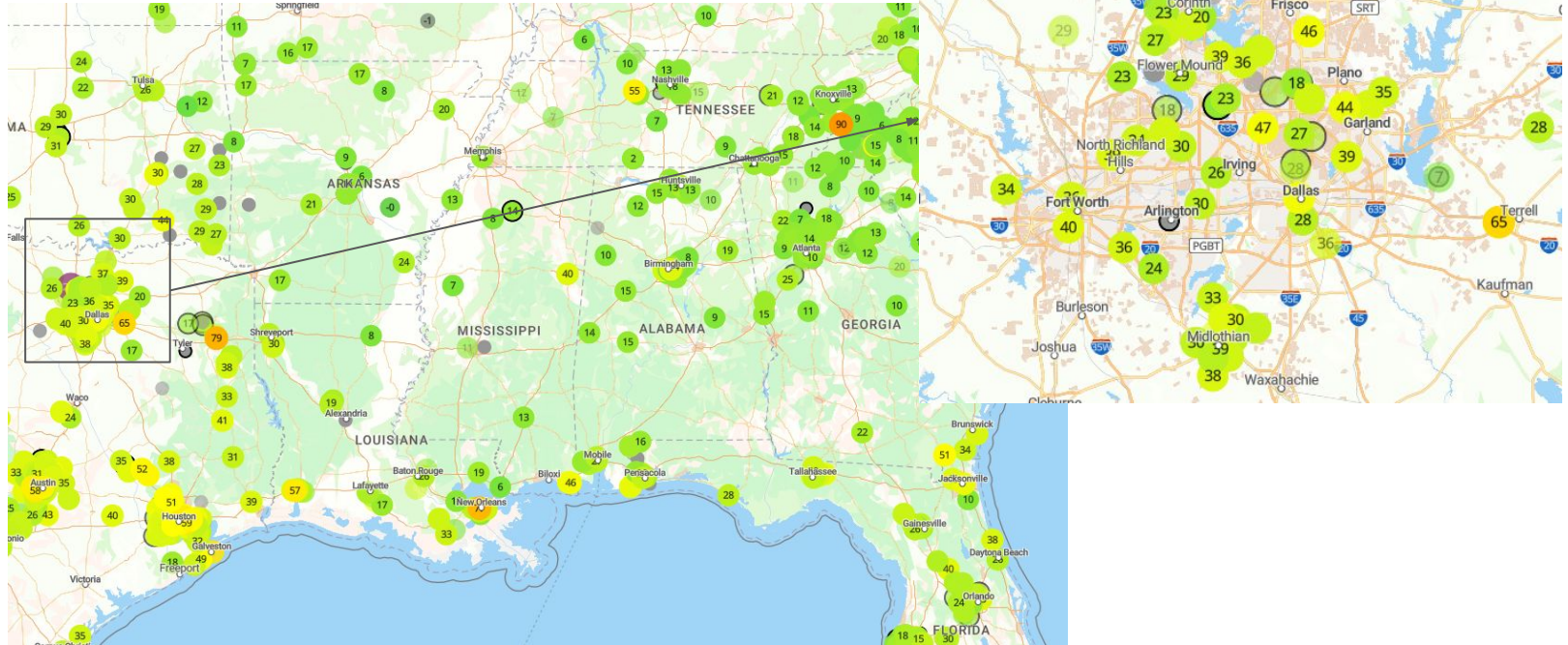
Data Importing and Cleaning

Used PurpleAir api₁ to collect historical data for Dallas Metropolitan region for the date range : 2022/04/01 to 2024/04/01

- ❖ First a list of sensors in the area defined by a geographical bounding box were identified.
- ❖ Historical data for a selected date range requested
- ❖ Filter, clean and load into a sqlite database
- ❖ Weather data for the date range is gathered from NOAA₂

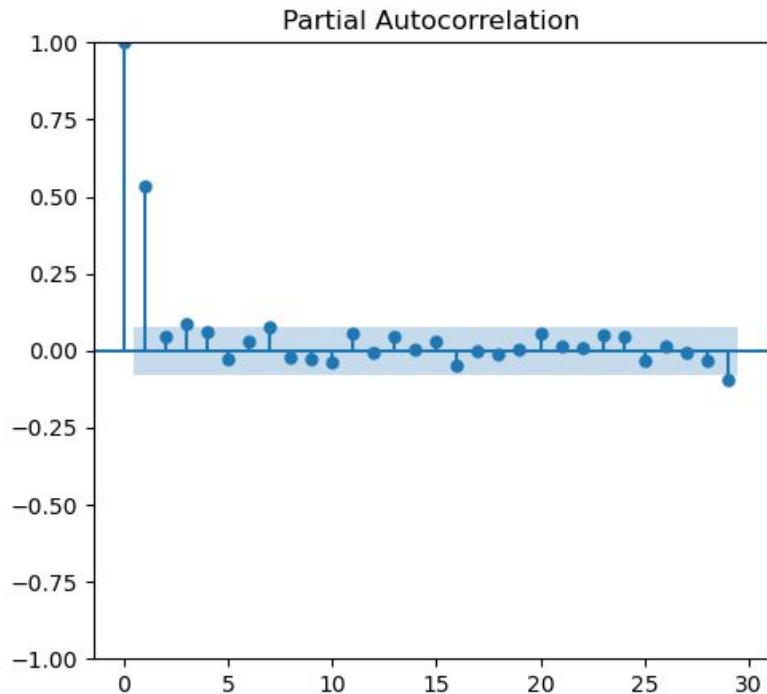
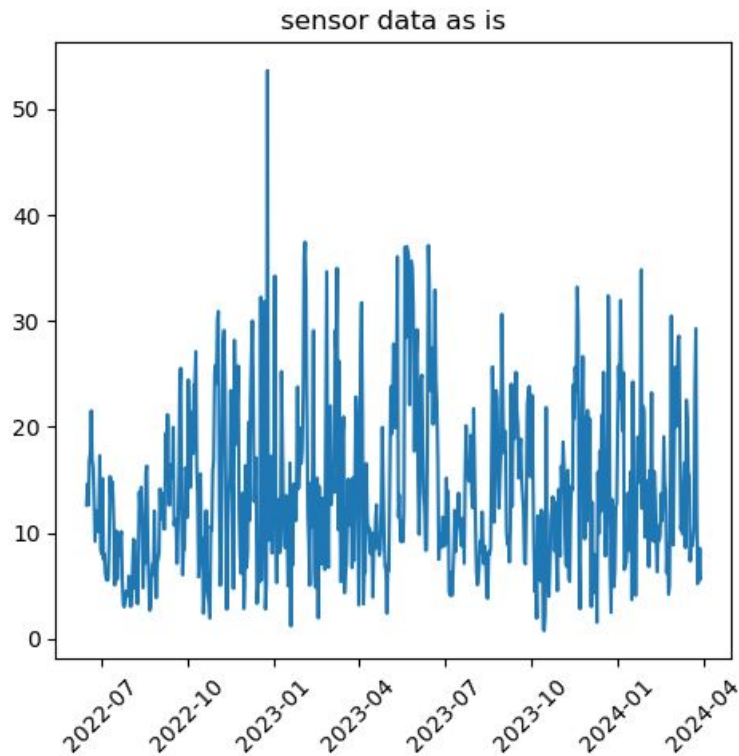
- ◆ 1. <https://api.purpleair.com>
- ◆ 2. <https://www.ncei.noaa.gov/cdo-web/api/v2/{endpoint}>

Real-time map with PM2.5 values



Daily PM2.5 data for a single sensor

- ☐ stationary
- ☐ dependent



Approach

Modelling Approaches

Temporal model: Time series forecasting

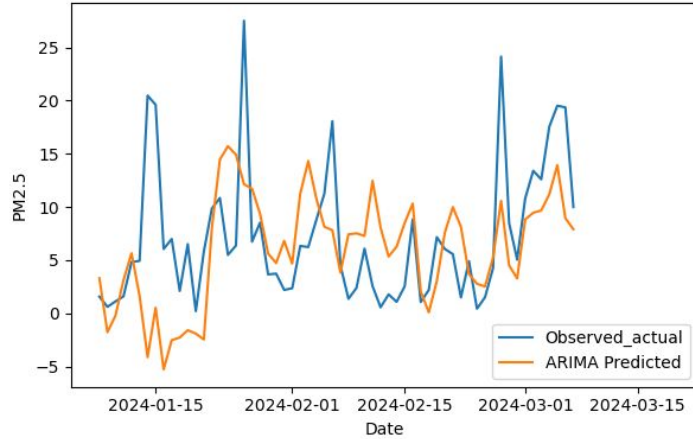
- Weather data included : Avg temp and Weather type (e.g. Fog, Rain etc)
- Multivariate ARIMA (1 0 1)
- XGBoost (with Temporal lag feature)
- No access to the spatial dimension

Spatial model:

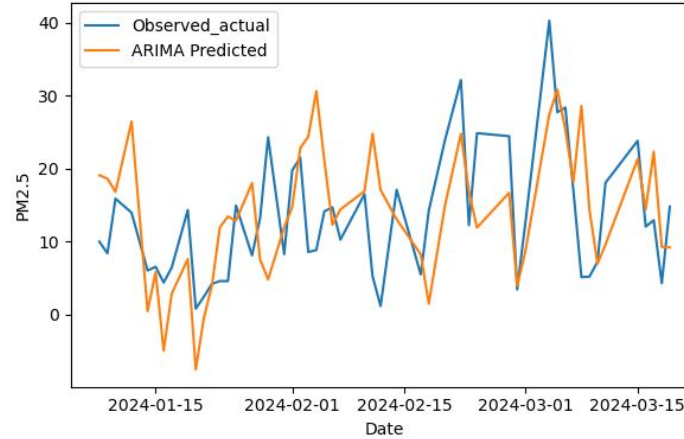
- Feature engineering to create spatial_feature to fit a regression model:
 - Use spatial_weights to create spatial_features corresponding to the PM2.5 values.
 - Train by leaving a sensor out
- Does not have access to the weather data

Time Series Forecasting with ARIMA

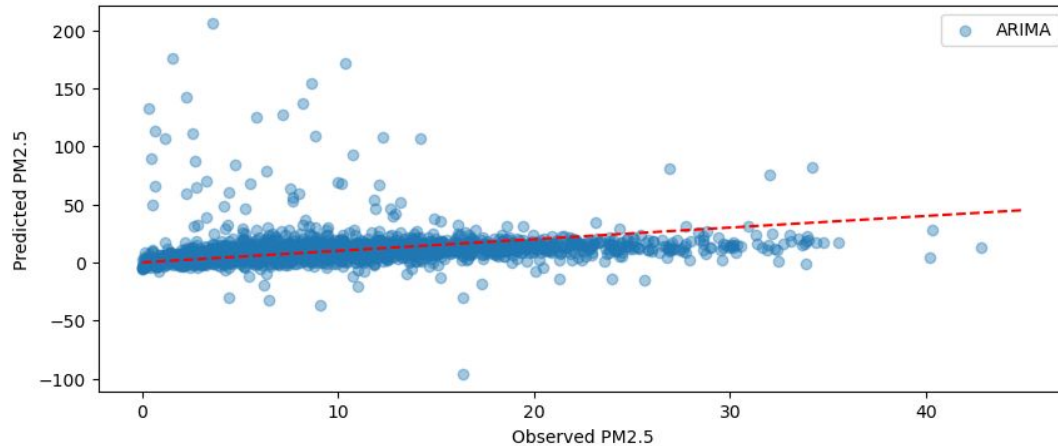
Time Series of PM2.5 for Sensor 127059



Time Series of PM2.5 for Sensor 112984



Observed vs Predicted PM2.5

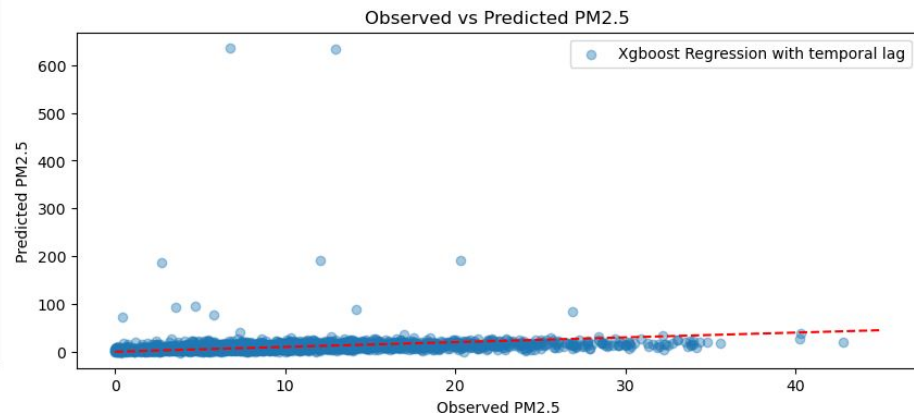
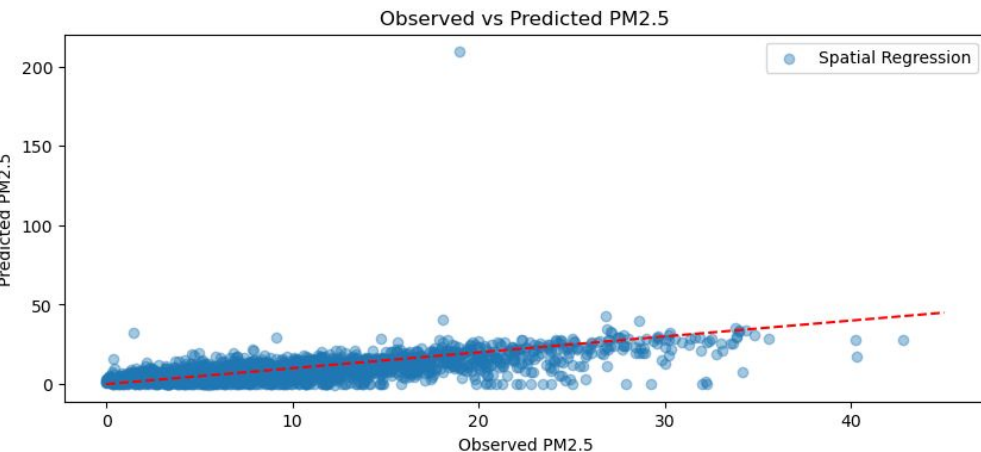
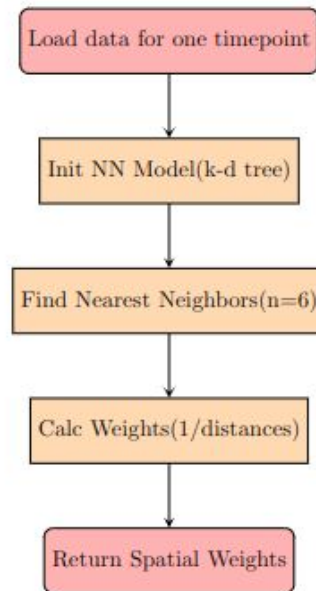


Humidity is the only variable that is significant at the 5% level

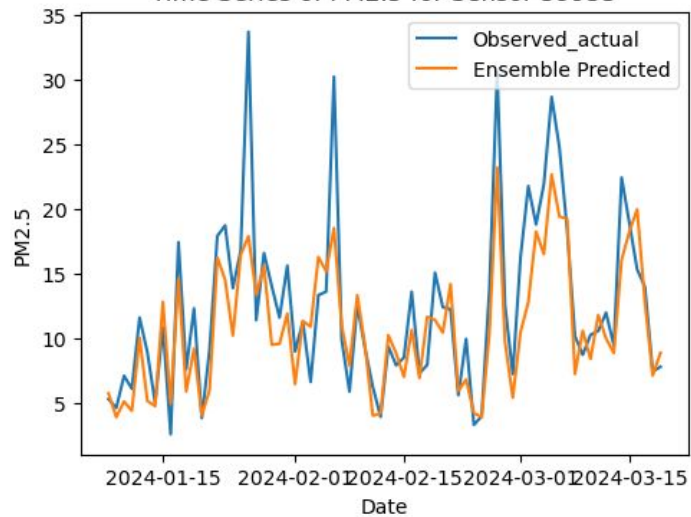
Spatio-Temporal prediction with XGBoost Regression

Feature Engineering :

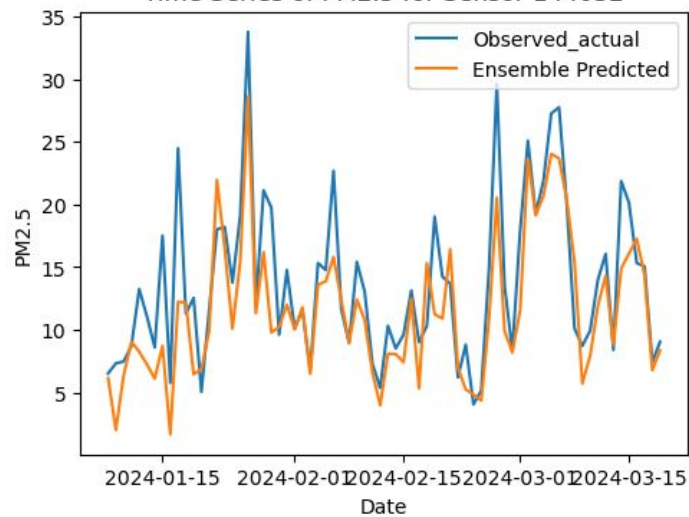
- 1) Spatial lag feature based on proximity to sensors in the geospatial dimension:
 - a) Does not have weather data
- 2) Temporal lag with a shift of 1 :
 - a) Has weather data
- 3) Spatio-Temporal prediction : Average of the two



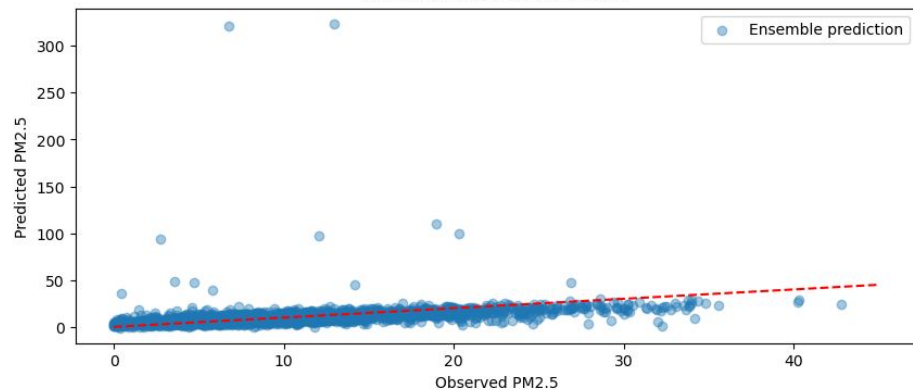
Time Series of PM2.5 for Sensor 59953



Time Series of PM2.5 for Sensor 144032

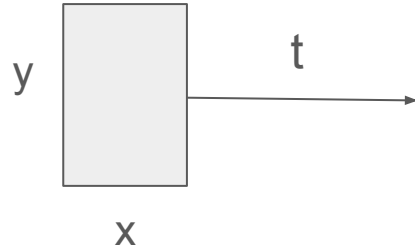


Observed vs Predicted PM2.5

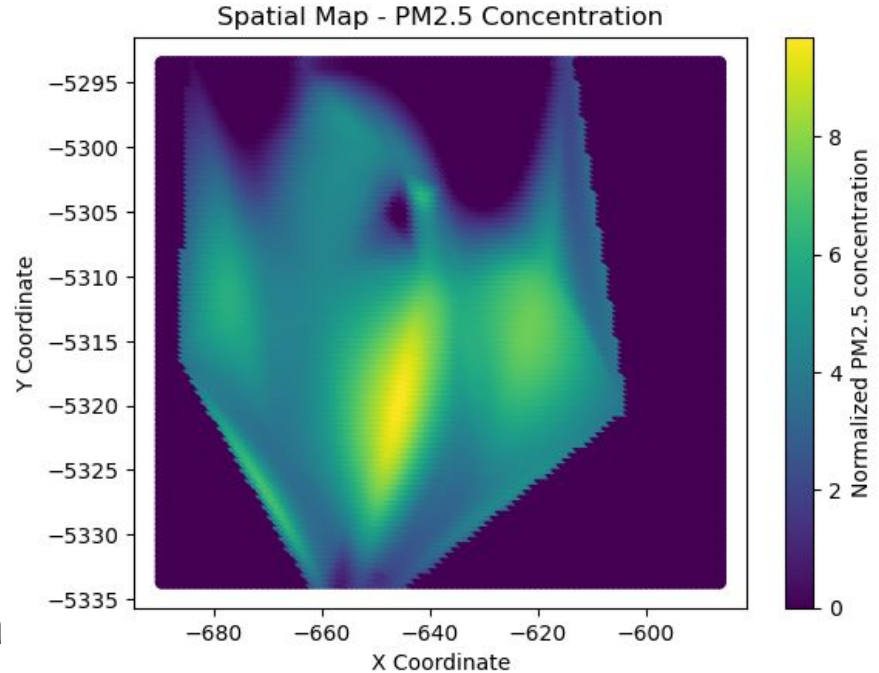


3D CNN for spatio-temporal models

- Temporal sequence of two-dimensional snapshots of values in geographical space



- Stack 2d images over time
- Capture both spatial and temporal tendencies



Single day map

Design and Workflow

Spatial Map Generation:

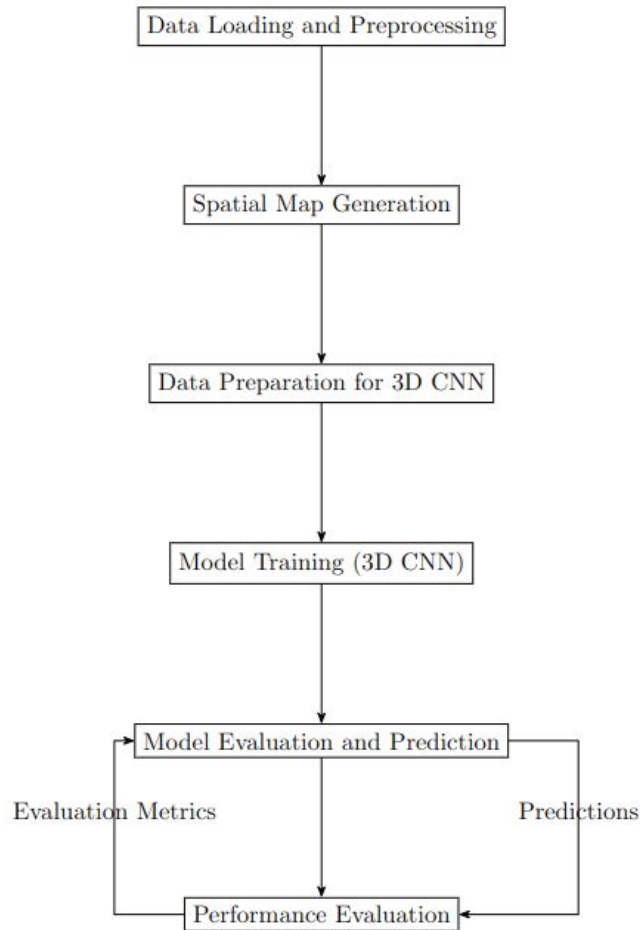
- Generate a series of spatial maps for each timestamp, capturing the spatial distribution of PM2.5 levels.

Data Preparation for 3D CNN:

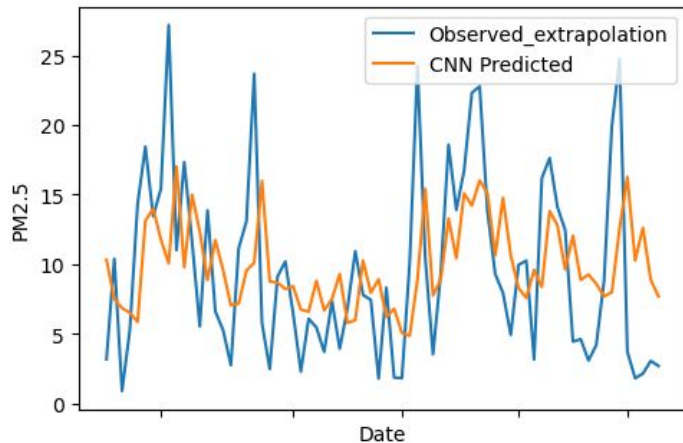
- Create sequences of spatial maps for input into the 3D CNN.
- Split the data into training and testing sets.

Model Training:

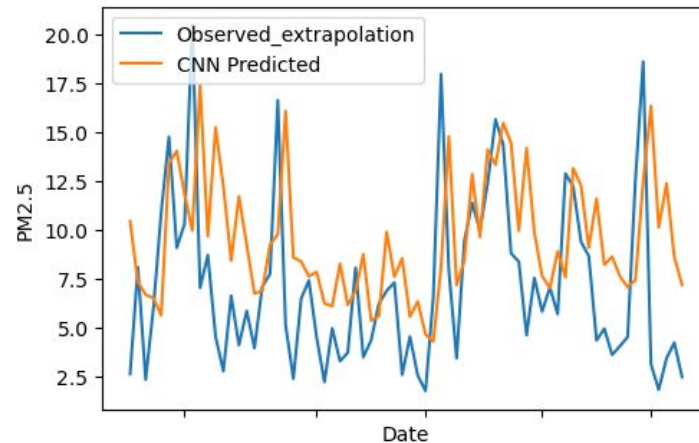
- Define and compile the 3D CNN architecture.
- Train the model on the training set.



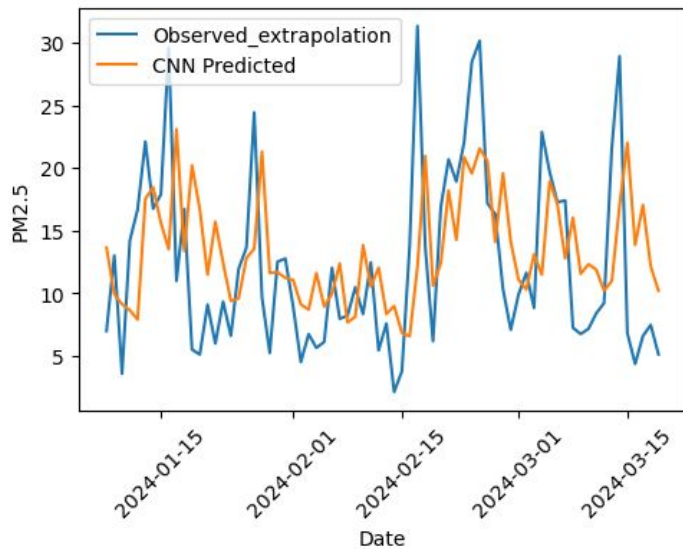
Time Series of PM2.5 for Sensor 123453



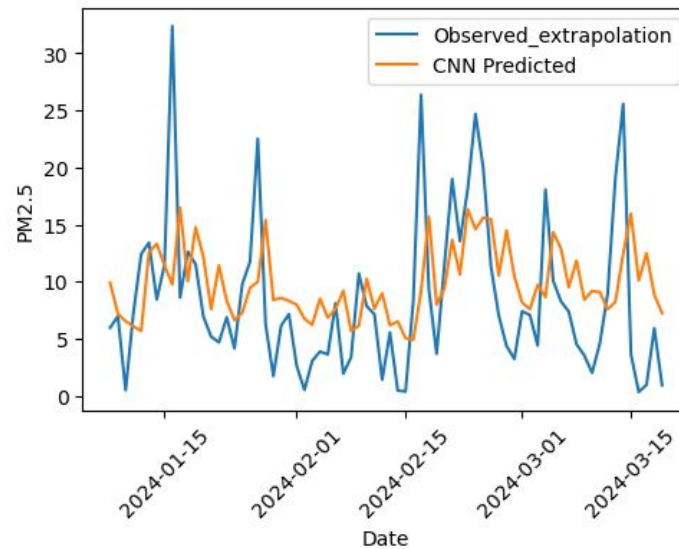
Time Series of PM2.5 for Sensor 99595



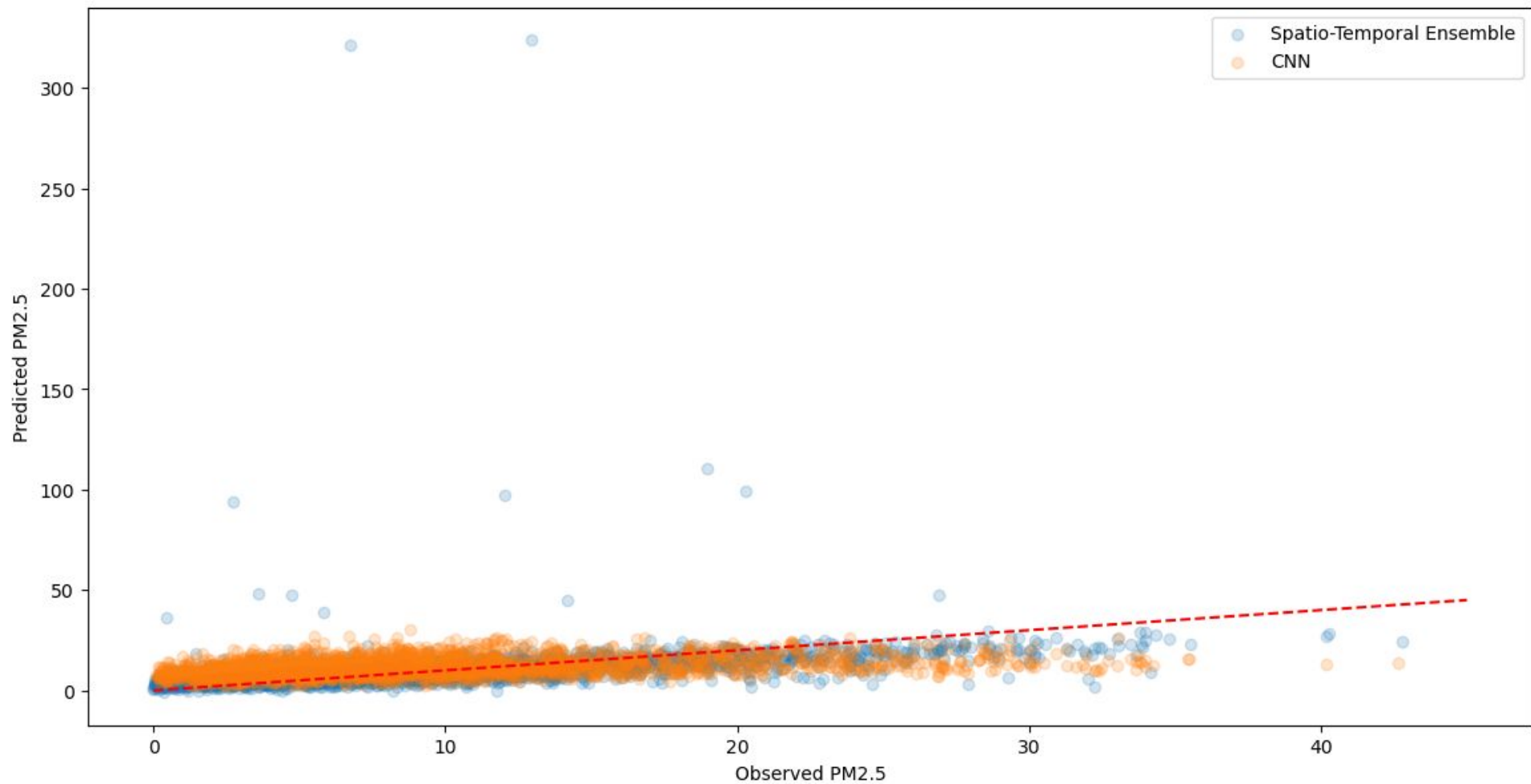
Time Series of PM2.5 for Sensor 87019



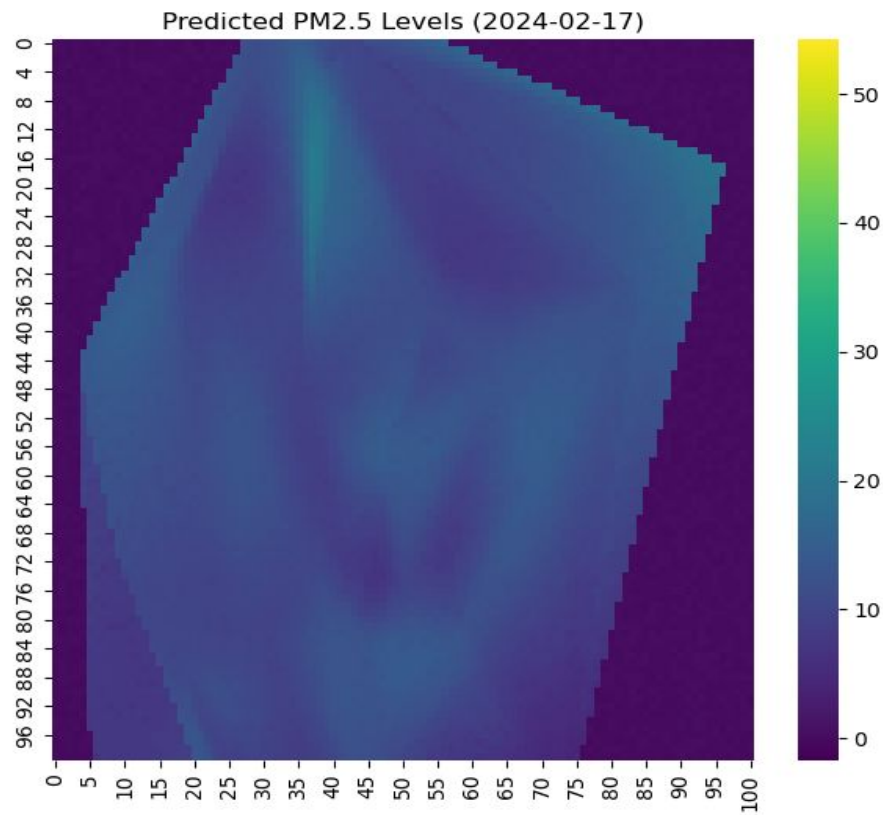
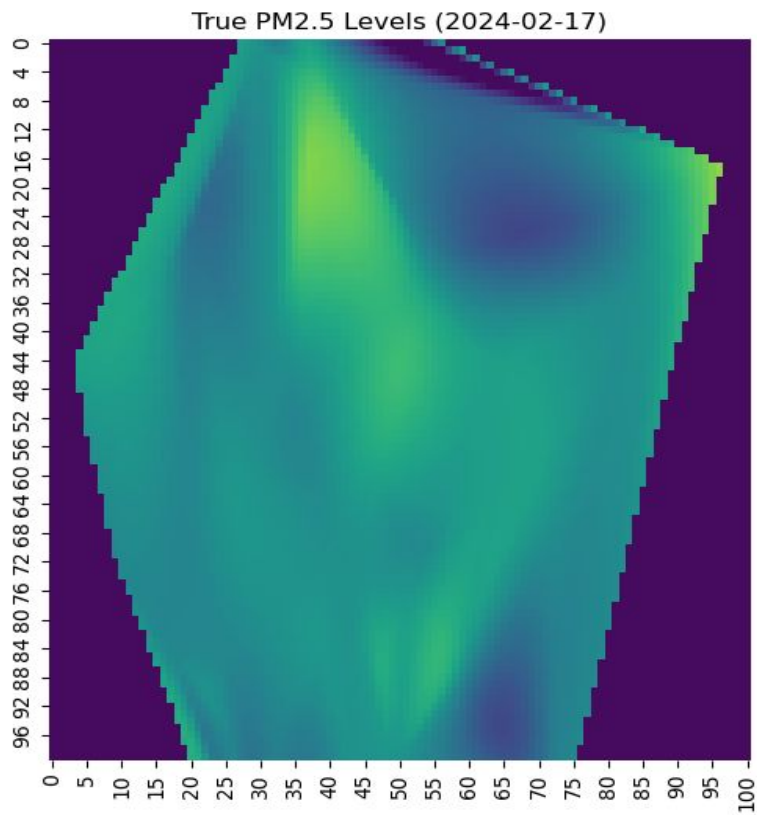
Time Series of PM2.5 for Sensor 120681



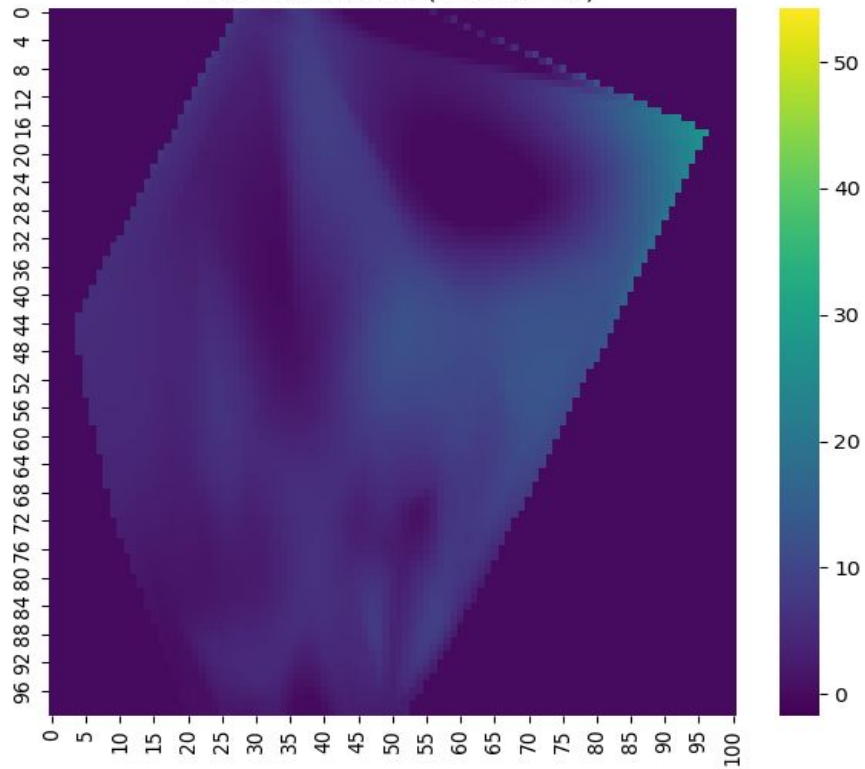
Observed vs Predicted PM2.5



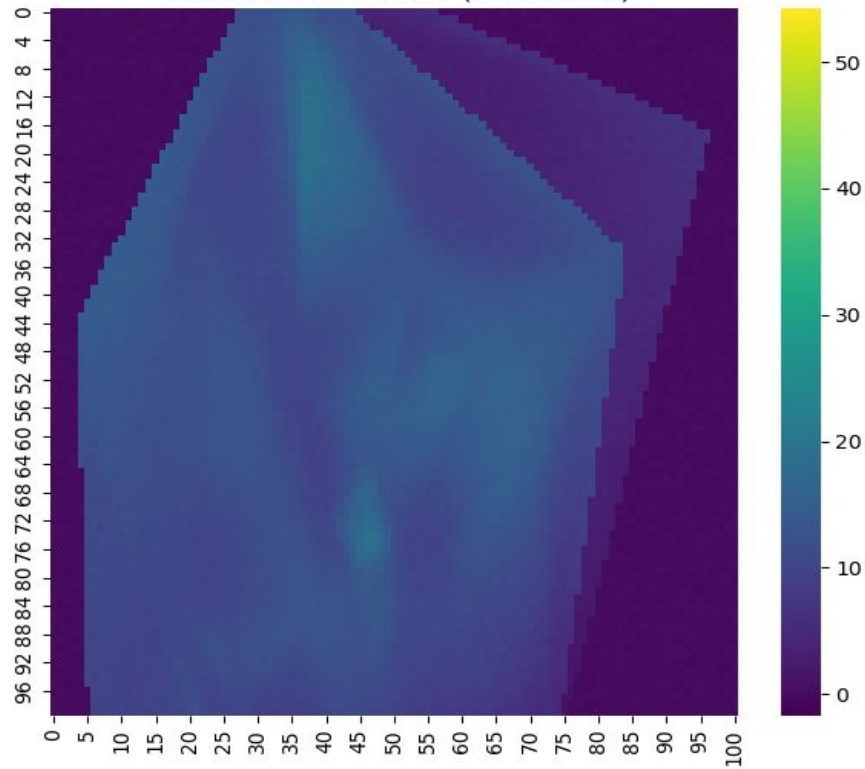
3D CNN prediction and observed (interpolated) PM2.5 values



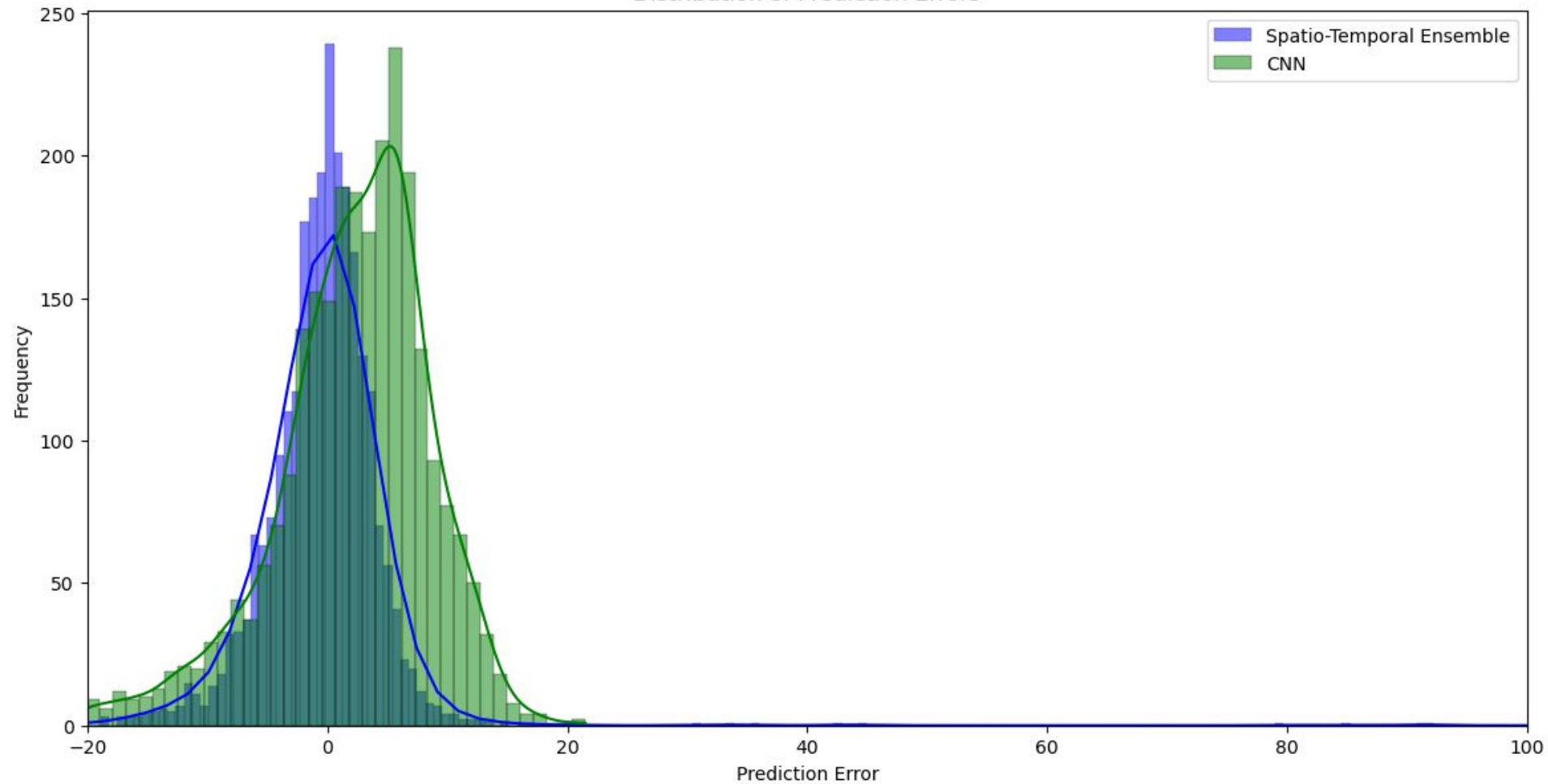
True PM2.5 Levels (2024-01-08)



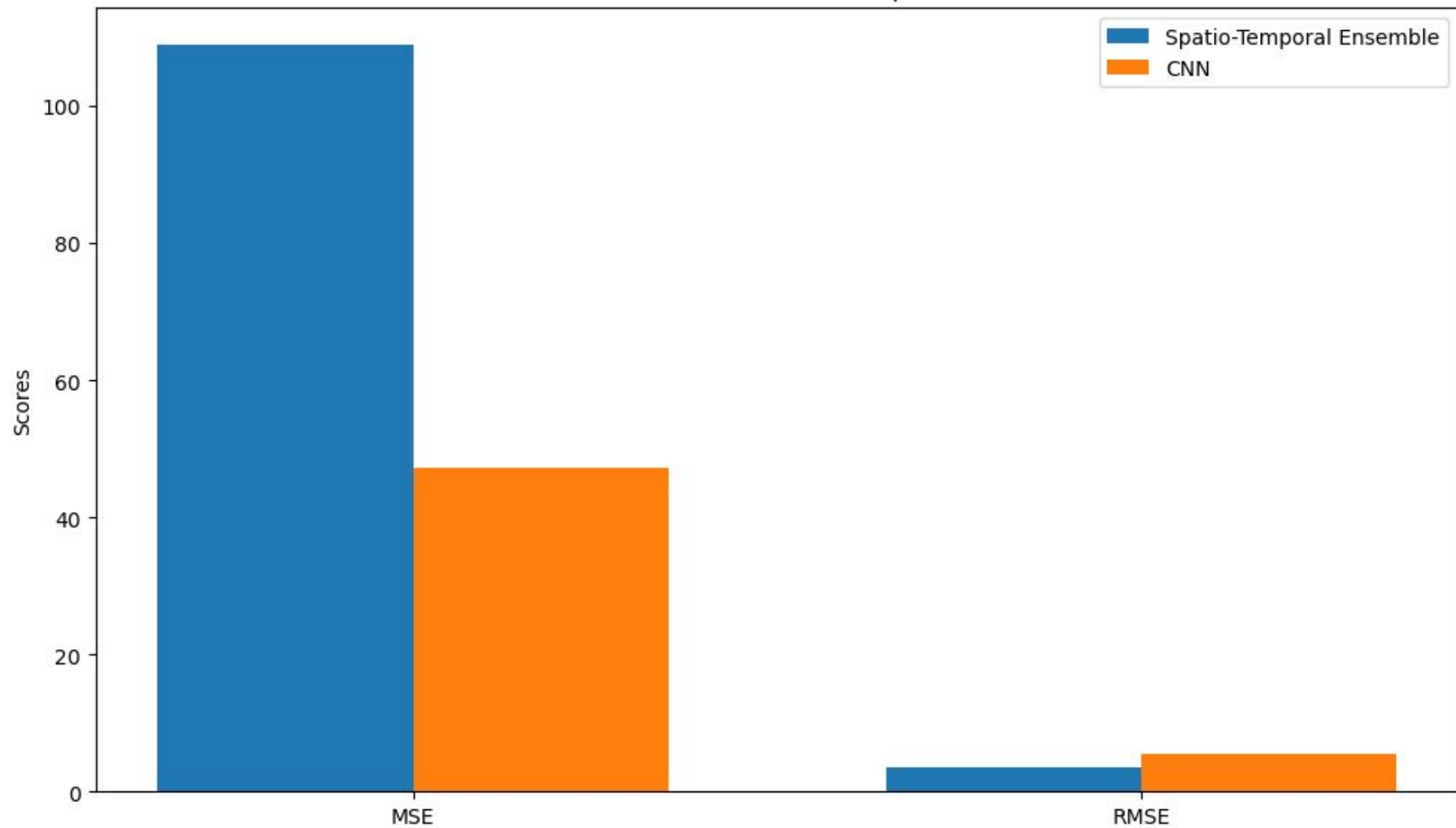
Predicted PM2.5 Levels (2024-01-08)



Distribution of Prediction Errors



Performance Metrics Comparison



Conclusions:

- ❑ XGBoost regression provided the closest prediction to test data
- ❑ 3D CNN underperformed compared to spatio-temporal regression
- ❑ Average wind speed, the generated lag features (both spatial and time), max and min temperatures, and precipitation most important features in regression
- ❑ For ARIMA model, only humidity was found to be significant