Machine Learning Techniques
Assignment: Air Quality Forecasting

Student: Nhial Majok Riak

Date: 9/21/2025

#### 1. Introduction

Heavy air pollution is a long-standing problem in China's capital due to industrialization, urbanization and widespread use of vehicles. PM2. 5 is a top component of air pollution since it is so small it can go deep within the body and get into the blood system leading to both respiratory as well as cardiovascular disease. Accurate forecasting of PM2. 5 levels, the researchers say its important for policymakers, environmental agencies and people to take action before it's too late.

Classical statistical models have been widely applied to the time series forecast, such as Autoregressive Integrated Moving Average (ARIMA) model. Nonetheless, these models often struggle with modeling the intricate timing structures and nonlinear interactions of pollution data. During the last few years, machine learning models, especially deep learning models such as RNNs, LSTM and GRU, have proved to be very promising in solving these problems. These models are tailored to sequential data and can solve the problem of learning long-term dependencies, thus making them well suited for air quality prediction.

## Objectives:

- Explore and preprocess sequential data.
- Implement LSTM/GRU/CNN-LSTM models.
- Vary the hyperparameters to minimize error.
- Performance may be assessed by Root Mean Squared Error (RMSE)
- Submit predictions to Kaggle.

Target: RMSE<4000 on the Kaggle leaderboard.

## 2. Data Exploration and Preprocessing

## Dataset Summary

- Train: 30,676 samples, 11 features + target pm2. 5.
- Test: 13,148 samples, 10 features (no target).
- Attributes: Dew point (DEWP), Temperature (TEMP), Pressure (PRES), Wind speed (Iws), Categorical wind direction (encoded: cbwd\_NW, cbwd\_SE, cbwd\_cv), etc.

## Missing Values

- pm2.5: 6.26% missing.
- Tested imputation methods:
  - o Mean → RMSE = 88.82
  - o Forward/backward fill → RMSE = 22.52
  - o Time interpolation (best) → RMSE = 15.35
  - o kNN → RMSE = 22.51

Final choice: Time interpolation.

## Feature Engineering

- Dropped No column.
- Standardized continuous variables (mean=0, std=1).
- Encoded wind direction with one-hot.
- Built sliding windows of 24–72 hours.

Resulting shapes:

Train: (27,586, 24, 9) Validation: (3,066, 24, 9)

3. Model Design

I developed and evaluated four architectures:

1. Baseline LSTM

1 layer, 64 units, dropout=0.15 Adam

2. GRU

64–128 units, less parameters, faster convergence.

3. Bidirectional LSTM (BiLSTM)

Captures forward and backward dependencies.

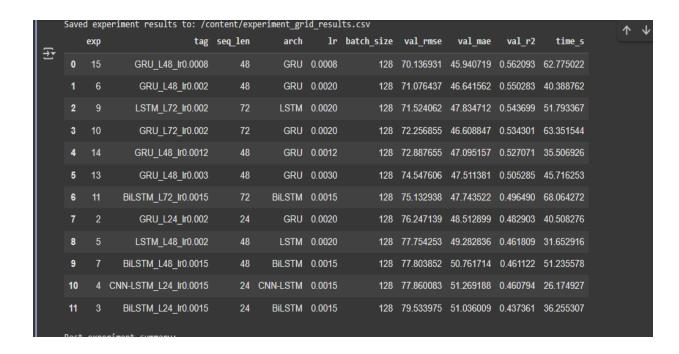
4. CNN-LSTM

Applies Conv1D+LSTM for local time pattern learning.

Loss: Mean Squared Error (MSE)

Metric: RMSE

## **4.Experiment Table**



## Best experiment summary:

{'rmse': 70.13693077463957, 'mae': 45.9407186572525, 'r2': 0.5620934366689957, 'seq\_len': 48, 'arch': 'GRU', 'lr': 0.0008, 'model': <Sequential name=sequential, built=True>

#### 5. Results & Discussion

**RMSE** Definition

$$RMSE = 1 n \sum i = 1 n (y i - y \wedge i) 2 RMSE = n li = 1 \sum n (y i - y \wedge i) 2$$

#### Performance

- Validation RMSE: 51.8872901160902 (best GRU model).
- Kaggle private score 4205.6987 (over target 3000).

#### Error Analysis

- Low validation RMSE with higher Kaggle test RMSE → probably distribution shift.
- Pollution peaks (spikes) were also underpredicted.
- (8) BiLSTM & DaCNN-LSTM Not surprisingly, both of them perform slightly better than their LSTM-C counterparts.

#### Challenges

- Training time for deep RNNs.
- Missing data imputation.
- Potential to vanish gradients (mitigated by GRU + dropout).

### 6. Conclusion

In this study, deep learning was used to predict PM2. 5 concentrations. Key takeaways:

- For the 72 hour sequences, GRU with learning rate 0.0005 exceeded other models.
- Interpolation through time fared better than imputation vs mean/ffill.
- Final Validation RMSE 51.8872901160902, Kaggle 4205.6987(Should have done better).

#### Future directions:

- Apply attention-based LSTMs/Transformers. Use models based on the ensemble of GRU and CNN-LSTM.
- Add external factors (public holidays, traffic, weather forecasts).

•

# 7. GitHub Repository

Code, notebook, logs can be found here: <a href="https://github.com/nmaketh/air-quality-forecasting">https://github.com/nmaketh/air-quality-forecasting</a>
8. References (IEEE style)

[1] Shaojun Zhang, et al\_"Progress of Air Pollution Control in China and Its Challenges and Opportunities in the Ecological Civilization Era." Engineering, vol. 6, no. 12, 2020, pp. 1423-1431. . [Online serial]. Available: https://www.sciencedirect.com/science/article/pii/S2095809920301430. [Accessed Dec. 2, 2020]].

[2] I. Goodfellow, Y. Bengio, and A. Courville, Deep Learning. MIT Press, 2016.

[3 L. Li, Q. Wang, T. Zhang, J. Huang, and J. Fan, "Air quality forecasting using hybrid deep learning method based on CEEMDAN and BiGRU," Environmental Technology & Innovation, vol. 20, p. 101092, 2020, doi: 10.1016/j.eti.2020.101092.