

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/350719854>

Evaluation of Time Series Forecasting Models for Estimation of PM_{2.5} Levels in Air

Preprint · April 2021

CITATIONS

0

READS

120

2 authors:



Satvik Garg

Jaypee University of Information Technology

11 PUBLICATIONS 9 CITATIONS

SEE PROFILE



Himanshu Jindal

Jaypee University of Information Technology

19 PUBLICATIONS 116 CITATIONS

SEE PROFILE

Evaluation of Time Series Forecasting Models for Estimation of PM2.5 Levels in Air

*Note: This paper is accepted and presented in the 6th I2CT 2021 conference. The final version of this paper will appear in the conference proceedings.

Satvik Garg¹, Himanshu Jindal²

Department of Computer Science, Jaypee University of Information Technology
Solan

Email: ¹satvikgarg27@gmail.com, ²himanshu.jindal@juitsolan.in

Abstract—Air contamination in urban areas has risen consistently over the past few years. Due to expanding industrialization and increasing concentration of toxic gases in the climate, the air is getting more poisonous step by step at an alarming rate. Since the arrival of the Coronavirus pandemic, it is getting more critical to lessen air contamination to reduce its impact. The specialists and environmentalists are making a valiant effort to gauge air contamination levels. However, it's genuinely unpredictable to mimic sub-atomic communication in the air, which brings about off-base outcomes. There has been an ascent in using machine learning and deep learning models to foresee the results on time series data. This study adopts ARIMA, FBProphet, and deep learning models such as LSTM, 1D-CNN, to estimate the concentration of PM2.5 in the environment. Our predicted results convey that all adopted methods give comparative outcomes in terms of average root mean squared error. However, the LSTM outperforms all other models with reference to mean absolute percentage error.

Index Terms—Air pollution, PM2.5, Forecasting, Time series, Machine learning, LSTM, CNN, ARIMA, FBProphet

I. INTRODUCTION

From the smog looming over metropolitan regions to pollution inside the home, air defilement represents a huge threat to well-being and the atmosphere. Air contamination accounts for an expected 4.2 million deaths per year because of stroke, coronary illness, cellular breakdown in the lungs, intense and constant respiratory sickness [1].

Air toxins can be present in the air anyplace, whether it is inside or outside. It incorporates gaseous toxins, for example, Carbon Monoxide (CO), Sulfur Dioxide (SO₂), Nitrogen Dioxide (NO₂), Ozone(O₃), and particle matters like PM2.5, PM10. The health impacts from particulate matter are controlled by its size, composition, source, solubility, and capacity to create reactive oxygen. PM2.5 represents air particulate matter with a diameter under 2.5 micrometers, represents 3% the diameter of human hair [2].

An examination was published by the Journal of American Medical Association to evaluate the connection between long-term exposure to fine particulate air contamination with the cellular breakdown in the lungs and cardiopulmonary mortality. It recommends that for every 10-microgram/m³ increase in particulate air contamination, the risk raised by an average of 4-8% of heart stroke and lung cancer mortality [3].

An early assessment of air contamination levels encourages the policymakers to choose the time specific strategies that shall need to be executed for supporting the residents, for example, sponsoring the utilization of public transport, offering free defensive facial covers, and financing clinical tests for asthma patients so that they can plan monetarily and successfully. Notwithstanding, these advantages and strategies are dependent upon knowing pollution levels in advance.

The forecasting of time-related data is a difficult problem because of the unknown changes in air contamination level patterns and conditions. In this research, we explored for finding the solutions that offered the best outcomes concerning lower prediction errors. In this regard, we used the stochastic model ARIMA [4], additive model FBProphet [5], and deep learning models LSTM [6], 1D-CNN [7].

This research work is categorized as follows: Section II provides the literature survey related to time series analysis. Section III describes the framework adopted in this work for forecasting. The analysis and evaluation of predictions using various metrics is presented in Section IV. Section V concludes the paper.

II. LITERATURE SURVEY

Forecasting of time series data is well known and excellent choice for analyzing the economy, stocks, human activity data, traffic, climate, sales, social media mining, and much more. Time series forecasting analyzes lags known as path observation to gain useful features from data to forecast future values using past data.

In early times, time-series data was forecasted using standard regressive models like AR, MA, ARMA, ARIMA [8]. These models are quite usual in forecasting economic and financial data. Still, they had some limitations as the models were not meant to analyze the non-linear behavior between variables and also when the data exhibit conditional covariance implies change in variance over time. However, one could solve this problem by integrating it with the Generalized Auto-regressive Conditional Heteroskedasticity (GARCH), but it's quite difficult to optimize its parameters [9].

Weitao Wang et al. [11] proposed a simple and effective hybrid model for forecasting traffic flow. This study finds out that traffic flow analysis is subjective to both linear and

non-linear relationships of data. The proposed hybrid model consists of the ARIMA model for linear fitting and Radial basis function artificial neural network (RBF-ANN) for the non-linear fitting of data.

Examining the requirement for fast and compelling techniques for performing web forecast of network traffic, Hao Yin et al. [12], proposed an adaptive autoregression (AAR) model. The idea was to integrate an adaptive order-selection and memory shortening technique to uphold the online forecast of dynamic network traffic data. The model achieved good exactness and low computation cost dissimilar to conventional Box–Jenkins time arrangement models like AR, MA, ARMA, ARIMA, etc.

The methods in deep learning was used and developed to deliver the difficulties identifying with the forecasting models. Considering the significance of CNN in various fields, Zhaoyi Xu et al. [13], adopted CNN for forecasting stock indexes, and the impacts of chronicled factors on the model were dissected. Finally, a couple of stock indexes were anticipated to confirm legitimacy and viability of the proposed model. The author also generated a hybrid model combined with CNN, which further improved the CNN network model.

The study [14] examines the use of deep neural networks; gradient boosted trees, random forest, and a simple ensemble of the models to forecast the S&P 500. The authors reported that random forest achieved better results than gradient boosted trees and deep neural networks.

Recently, the Facebook data science team deployed an open-source additive time series model called FBProphet. Assessing the open-source algorithm, faster results, and accuracy, Alabi et al. [15] used FBProphet to estimate COVID-19 passings and confirmed cases. The accuracy achieved by prophet was 79.6% for the data from World Health Organization.

From the survey, we aim to show the time series analysis from basic models like ARIMA to recently developed models like FBProphet. We also discussed deep learning approaches like CNN, and machine learning techniques like the random forest, gradient boosting. The problems were also reviewed using the ARIMA model and how one can solve them using GARCH. Generally, all studies focused on creating a slight change in the existing model and modify them to provide outcomes. However, a proper comparison was somehow limited, covering various approaches for evaluation provided in this research.

III. METHODOLOGIES

The dataset used in this research is Beijing Multi-Site Air-Quality, taken from the UCI Machine Learning Repository [16]. It incorporates hourly air contamination information from 12 broadly controlled air-quality observing locales (stations) between the timeframe March 1st, 2013 to February 28th, 2017. It contains 18 features, which are given in Fig. 1. The target variable is pollution measured in PM2.5. For each station, the dimensions are 18 by 35064, so for a total of 12 stations, it is equivalent to 18 by 420768.

Attribute	Description	Attribute	Description	Attribute	Description
No	Row id	TEMP	Temperature in degree Celsius	PM2.5	PM2.5 concentration in $\mu\text{g}/\text{m}^3$
year	Year of observation	PRES	Atmospheric pressure in hectopascals	PM10	PM10 concentration in $\mu\text{g}/\text{m}^3$
month	Month of observation	DEWP	Dew point temperature in degree Celsius	SO2	SO2 concentration in $\mu\text{g}/\text{m}^3$
day	Day of observation	RAIN	Precipitation in millimeter	CO	CO concentration in $\mu\text{g}/\text{m}^3$
hour	Hour of observation	wd	Wind direction	O3	O3 concentration in ($\mu\text{g}/\text{m}^3$)
station	Air-quality monitoring site name	WSPM	Speed of wind in meter per second	NO2	NO2 concentration in ($\mu\text{g}/\text{m}^3$)

Fig. 1: Attribute Information.

Fig. 2 explains and describes the framework adopted in this research for measuring the PM2.5 levels. We divided the framework into three phases, namely, Data Preprocessing, Modeling, and Evaluation phase.

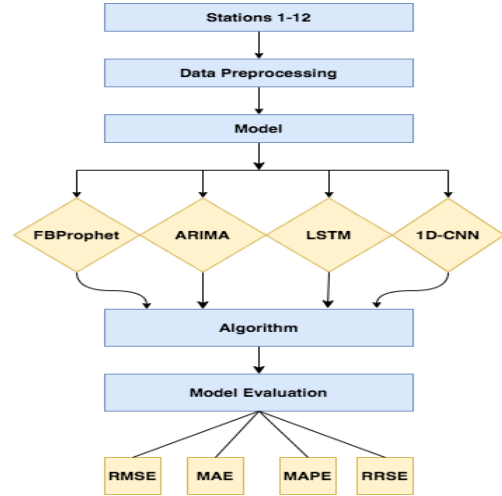


Fig. 2: Flowchart of the framework adopted.

A. Data Preprocessing

Prior training the models it is extremely important to clean the data. We used some basic data preprocessing techniques in this work. Instead of dropping rows containing null values, we applied the 'fillna' method using pandas library [17] to forward fill the empty values in the dataset. The 'wd' feature was converted from categorical to numerical values using the label encoder function from the sci-kit learn toolkit [18]. We took an hour, day, month, year attributed from the data and formed an ordered day by day level 'DateTime' column, which makes the total row count equals 1421 for each station. The data was then divided into a training and a test set. The training data consists from March 1st, 2013 to February 28th, 2016, which accounts for 75% of data and test data from March 1st, 2016 to February 28th, 2017 (25%).

For deep learning methods, LSTM, 1D-CNN, considering the necessity for validation data [19], 20 percent of training data was used for validation purposes.

In fundamental terms, there is a stationary (uniform) and non-stationary (non-uniform) time series. A uniform time series is one whose quantifiable properties such as the mean, variance, and autocorrelation, are congruent with time. Henceforth, a non-uniform time series represents a change in properties over time which shows the presence of a trend. Non-uniform time series should be first changed over into detrended arrangement prior to applying the models. If the time arrangement is dependably extending or diminishing over time, the example mean and the variance will form with the size of the example, and this will reliably deprecate the mean and variance in future periods.

To examine in-case the time series is uniform or not, Augmented Dickey-Fuller (ADF) test [20] is used. Prior going to the ADF test, we should initially comprehend, 'what is the Dickey-Fuller test.' It is a unit root test that checks for the null hypothesis ($\alpha = 1$) in the equation below:

$$y_t = c + \beta t + \alpha y_{t-1} + \phi \Delta Y_{t-1} + e_t \quad (1)$$

where,

y_{t-1} = first Lag of time series

ΔY_{t-1} = first difference of the series at time (t-1)

In simple words, if the value of α equals 1, it indicates the presence of unit root. Thereby, the series has taken to be non-stationary. The ADF test allows higher-order regressive processes in the model by including ΔY_{t-p} .

$$y_t = c + \beta t + \alpha y_{t-1} + \phi_1 \Delta Y_{t-1} + \phi_2 \Delta Y_{t-2} + \dots + \phi_p \Delta Y_{t-p} + e_t \quad (2)$$

The null hypothesis is similar to the dickey fuller test. However, the p-value acquired ought to be less than the significance level of 5% to dismiss the null hypothesis that is the presence of unit root. We applied this test on our target variable PM2.5 and discovered that our series is stationary. Figures 3 and 4 show the subplots of six main air pollutants and six external relevant meteorological variables on Aotizhongxin station, respectively. One can clearly see that there is a seasonality present in various features like Temperature, Dew point Temperature, Precipitation, CO, NO2, O3 and SO2.

B. Modeling

Four different types of models, namely, FBProphet, ARIMA, LSTM, 1D-CNN has been adopted in this work to evaluate time series forecasting PM2.5 levels.

1) *FBProphet*: Facebook prophet [5] is an open-source tool developed by Facebook used for time series analysis derive from a decomposable additive model. It takes into account holidays, and it usually fits nonlinear data with yearly, weekly and daily seasonality. Prophet uses time as a regressor and fits various linear and nonlinear functions of time as components, which is combined given in this equation:

$$y(t) = g(t) + s(t) + h(t) + e(t) \quad (3)$$

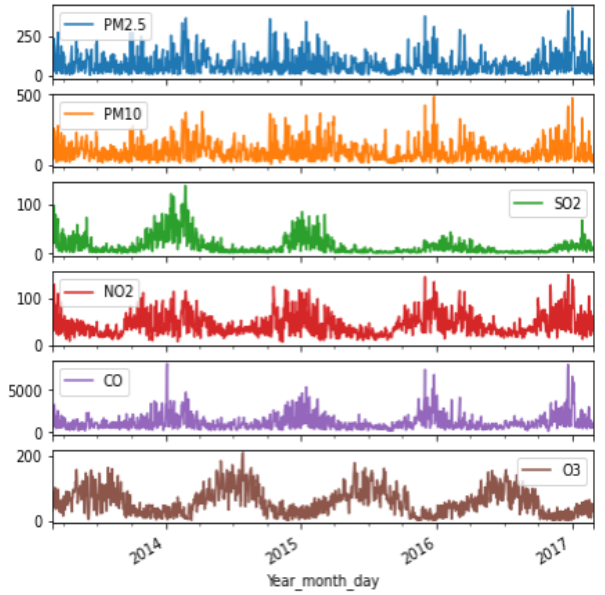


Fig. 3: Plot of 6 main atmospheric pollutants concentration with respect to days on station Aotizhongxin.

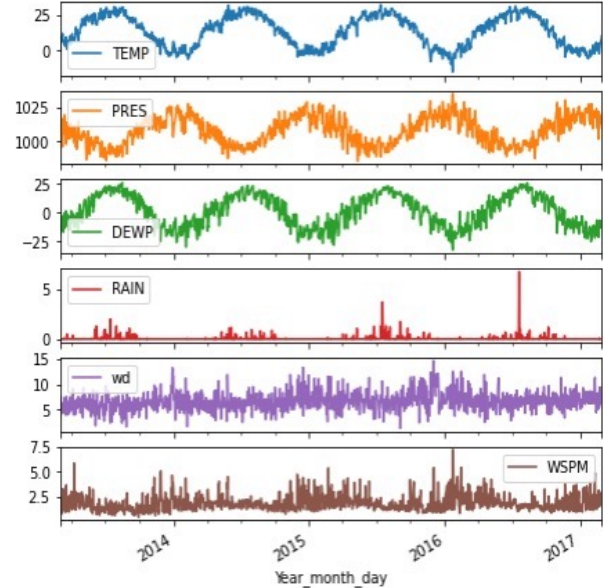


Fig. 4: Plot of 6 meteorological factors with respect to time on station Aotizhongxin

where,

$y(t)$: predictions (forecast).

$g(t)$: trend alludes to changes throughout a significant stretch of time

$s(t)$: refers to seasonality for example, weekly, daily, yearly.

$h(t)$: holidays

$e(t)$: error term represents any surprising changes not obliged by the model

The Facebook prophet has multiple trend, seasonality, change-points, and holiday parameters which needs to be tuned for

better results. Parameters and its values that optimized in this research is shown in Fig. 5. Change point refers to those points at which time series experience an abrupt change. These progressions can be because of anything, for instance, unexpected catastrophe, new government guidelines, and so forth. The parameter optimization resulted in 144 possible model counts, assessed by evaluation metrics.

Parameter	Value	Parameter	Value	Parameter	Value
Changeout_prior_scale	0.1, 0.2, 0.3	N changepoints	[100, 150]	Seasonality_mode	additive
Holidays_prior_scale	0.1, 0.2, 0.3	Weekly_seasonality	True, False	-	-
Daily_seasonality	True, False	Yearly_seasonality	True, False	-	-

Fig. 5: Hyperparameters FBProphet

2) *ARIMA*: ARIMA(p,d,q) [4], composition of autoregression, integrated and moving average is a regressive model used to forecast time series data where p, d, q referred as autoregression order. Auto Regression, AR(p), is a part of the ARIMA model based on the idea that it uses its own lags (past values) as predictors, where p is a boundary of the number of lags taken in

$$y_t = \alpha + \sum_{i=0}^p \beta_i Y_{t-i} + \epsilon_t \quad (4)$$

where, α represents intercept, ϵ_t adures white noise, $\sum_{i=0}^p \beta_i$

are the coefficients of the past values (lags) given by $\sum_{i=0}^p Y_{t-i}$

which are calculated by the model.

Moving Average, MA(q), utilizes residual error of past time points to foresee current and future predictions. Moving normal (MA) eliminates arbitrary developments from a time series. The parameter q is the number of lags forecast errors that utilized to compute current values.

$$y_t = \alpha + \sum_{i=0}^q \phi_i \epsilon_{t-i} + \epsilon_t \quad (5)$$

$\sum_{i=0}^q \epsilon_{t-i}$ represent error terms of the respective lags.

Integrated, I(d), property helps make time-series data stationary to eliminate time dependency and trend. Parameter 'd' represents the degree of difference, which means the number of times the data was differenced. If a time series is stationary, then its degree of difference is zero.

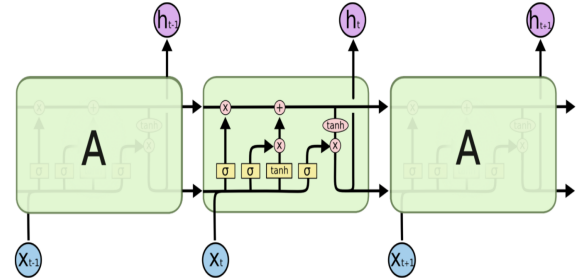
The ARIMA(p,d,q) model is described as:

$$y_t = \alpha + \sum_{i=0}^p \beta_i Y_{t-i} + \sum_{i=0}^q \phi_i \epsilon_{t-i} + \epsilon_t \quad (6)$$

The standard ARIMA models manually take input values (p,d,q) using autocorrelation, and various other statistical tests. In this study, we use Auto Arima from the pmdarima package

[10], which automatically evaluates p, d, q values on its own. Considering our time series is stationary, we kept d=0 and attempted with the p, q values ranging from 0 to 5 to optimize the ideal value for the model.

3) *LSTM*: LSTM [6] represents long short term memory networks. It is a model or design that expands the short-term memory of recurrent neural networks that help determine the time series problems effectively. RNN deals with the current input by considering only the previous yield (feedback) and putting it in its memory for a brief timeframe (short-term memory). Thus, neglecting to store data for a long time implies the inefficiency of RNNs for dealing with long-term dependencies. Different issues with RNNs are vanishing and exploding gradients problems during the training phase by backtracking. This would halt network from learning since the updated weights become smaller and smaller or bigger and bigger. The network should be rebuilt in such a way so that it scales down the scaling factor to one. It should be possible by utilizing different gate units in memory blocks, associated through layers, as appeared in Fig. 6 and called LSTM. This research utilizes just a single LSTM layer with only 128 neurons in it, followed by a fully connected layer for prediction.



Source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Fig. 6: LSTM repeating network layer

4) *CNN*: Convolutional neural network models [7] were mainly produced for picture classification problems. The convolutional layers extract the features from a two-dimensional input, referred to as feature engineering. This equivalent process could outfit on one-dimensional arrangements of data, for example, sequence problems. The advantage to implement CNNs for sequence data is that they can learn from the raw time series data straightforwardly as it provides an architecture to perform smoothing parameters. The underlying two layers of CNN are typically a convolutional layer and a pooling layer. Both perform smoothing, followed by a Dense layer to use the smoothed data, and perform well on a forecasting task. A flatten layer is also used between the max-pooling layer and the fully connected layer to reduce generated multidimensional features to a one-dimensional vector.

C. Algorithms

As shown in Algorithm 1 and 2, input to the proposed algorithm is PM2.5 series (Data) and Features series set of

12 stations. Features represent external variables like NO₂, PM₁₀, SO₂, CO, wd, DEWP, and many more. These features help to build the model efficiently. After building the model, both algorithms output the results of models using four error metrics, RMSE, MAE, MAPE, RRSE. The first step is to divide the dataset for training and testing purposes. Algorithm 2 followed a similar train test split like Algorithm 1 (75:25). However, the 20% of training data and features taken as validation data appeared in lines 7 and 10. For deep learning methods, it is good practice to normalize the data to the range of 0 to 1 since the networks are sensitive to the scale of input data. The MinMaxScaler preprocessing function is used as appeared in lines 1-4.

During the model fitting, for FBProphet as shown in Algorithm 1, applied various hyperparameters that were shown in Fig. 5 using a simple loop. The model was then fitted using a fit() method in line 10. Each fitted model predicted the test data using Test Features as given in line 11.

ARIMA followed the similar approach like FBProphet. However, we didn't use any loop as the values of hyperparameters are generated automatically. The auto_arima method in line 14 will try various hyperparameters on its own and returns the best model having the lowest Akaike information criterion (AIC) score [22].

Algorithm 1 Algorithm for FBProphet and ARIMA

Input: Data, Features
Output: Evaluation metrics of the predicted data
Data split : 75% train and 25% test data
Train
1: $count \leftarrow \text{length}(\text{Data}) * 0.75$
2: $X \leftarrow \text{Data}(0 : count)$
3: $Z \leftarrow \text{Features}(0 : count)$
Test
4: $x \leftarrow \text{Data}(count :)$
5: $z \leftarrow \text{Features}(count :)$
Model fitting FBProphet
6: $parameters \leftarrow \text{Hyperparameters}$
7: **for each** p in $parameters$ **do**
8: $model \leftarrow \text{Prophet}(p, \text{interval_width} = 0.95)$
9: $model.add_regressor(Z)$
10: $model.fit(X)$
11: $forecast \leftarrow model.predict(z)$
12: **return** $rmse, mae, mape, rrse$
13: **end for**
Model fitting ARIMA
14: $model \leftarrow \text{auto_arima}(X, \text{exogenous}=Z)$
15: $model.fit()$
16: $forecast \leftarrow model.predict(n_periods = \text{len}(x), z)$
17: **return** $rmse, mae, mape, rrse$

As shown in Algorithm 2, the deep learning methods, LSTM, 1-D CNN, trained independently for 200, 400, 600, 800, 1000 epochs to check for underfitting and overfitting. However, for LSTM additionally applied the 'tanh' and 'relu' activation function to check for the nonlinearity of 'relu' to help in improving the model. These parameters resulted in the formation of two loops which made a sum of ten distinct models per station for LSTM. 1-D CNN alike to LSTM except there is only one loop used for epochs hyperparameter resulting in five models per station. We use 'adam' optimizer for model compilation as appear in lines 20 and 34.

Algorithm 2 Algorithm for LSTM and 1D-CNN

Input: Data, Features
Output: Evaluation metrics of the predicted data
Data split : 75% training + validation and 25% test data
Normalization
1: $object1 \leftarrow \text{MinMaxScaler}()$
2: $object2 \leftarrow \text{MinMaxScaler}()$
3: $Data2 \leftarrow object1.fit_transform(Data)$
4: $Features2 \leftarrow object2.fit_transform(Features)$
Train and Validation
5: $count \leftarrow \text{length}(Data2) * 0.75$
6: $X \leftarrow Data2(0 : count)$
7: $V \leftarrow X * 0.20$
8: $X1 \leftarrow X - V$
9: $Z \leftarrow Features2(0 : count)$
10: $v \leftarrow Z * 0.20$
11: $Z1 \leftarrow Z - v$
Test
12: $x \leftarrow Data(count :)$
13: $z \leftarrow Features(count :)$
Hyperparameters
14: $epochs \leftarrow [200, 400, 600, 800, 1000]$
15: $activation \leftarrow ['tanh', 'relu']$
Model fitting LSTM
16: **for each** i in $activation$ **do**
17: $model \leftarrow \text{Sequential}()$
18: $model.add(\text{LSTM}(128), \text{activation} = i)$
19: $model.add(\text{Dense}(1))$
20: $model.compile(\text{loss} = 'mae', \text{optimizer} = 'adam')$
21: **for each** j in $epochs$ **do**
22: $model.fit(Z1, X1, \text{validation} = (v, V), \text{epochs} = j)$
23: $forecast \leftarrow model.predict(z)$
24: $forecast \leftarrow object1.inverse_transform(forecast)$
25: **return** $rmse, mae, mape, rrse$
26: **end for**
27: **end for**
Model fitting 1D-CNN
28: $model \leftarrow \text{Sequential}()$
29: $model.add(\text{Conv1D}(\text{filters} = 128, \text{kernel_size} = 2, \text{activation} = 'relu'))$
30: $model.add(\text{MaxPooling1D}(\text{pool_size} = 2))$
31: $model.add(\text{Flatten}())$
32: $model.add(\text{Dense}(64, \text{activation} = 'relu'))$
33: $model.add(\text{Dense}(1))$
34: $model.compile(\text{loss} = 'mae', \text{optimizer} = 'adam')$
35: **for each** i in $epochs$ **do**
36: $model.fit(Z1, X1, \text{validation_data} = (v, V), \text{epochs} = i)$
37: $forecast \leftarrow model.predict(z)$
38: $forecast \leftarrow object1.inverse_transform(forecast)$
39: **return** $rmse, mae, mape, rrse$
40: **end for**

D. Metrics

Let the letter be, $RMSE$: root mean squared error, MAE : mean absolute error, $MAPE$: mean absolute percentage error, $RRSE$: root relative squared error, x_i : actual data, y_i : predicted data, n : length of test data. The following metrics are used in this research to analyze the forecasted results:

$$MAPE = \left(\frac{100}{n}\right) \sum_{i=0}^n \frac{|x_i - y_i|}{|x_i|} \quad (7)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - y_i)^2} \quad (8)$$

$$MAE = \left(\frac{1}{n}\right) \sum_{i=1}^n |(x_i - y_i)| \quad (9)$$

$$RRSE = \sqrt{\frac{\sum_{i=1}^n (y_i - x_i)^2}{\sum_{i=1}^n (x - x_i)^2}} \quad \text{where } x = \left(\frac{1}{n}\right) \sum_{i=1}^n (x_i) \quad (10)$$

IV. RESULTS

In this evaluation, four distinct models were assessed, to be specific LSTM, CNN, ARIMA, FBProphet for the time Series examination on gauging PM2.5 levels on 12 Stations. We utilized test information for forecasts that represent 25 percent of the complete dataset. For deep learning methods, LSTM, CNN, 20% of the training data adopted for validation purposes. RMSE, RRSE, MAPE, MAE scores were used to assess model execution.

MAE is simply the mean of absolute error. MAE calculates errors on a similar scale which infers that it treats big and small errors equally. It isn't sufficient to investigate the forecasts appropriately. MAPE is sort of standardized absolute error, which permits the errors to be contrasted across data with various scales. MAPE is processed over each information point and averaged, and accordingly catches more mistakes and exceptions. It is also helpful to punish negative errors as the estimation of the actual data would be more smaller than forecasted data. On the other hand, RMSE is helpful for penalizing larger errors as the errors are squared which comparatively gives higher weight to larger values. RRSE simply measures the relativeness of the predictions with the average of actual values.

Table I and II shows the forecasting metrics results for examining PM2.5 levels on all 12 stations. The order for ARIMA is referenced alongside with the results to bring about in the tables. From the outcomes, one can without much of a stretch see that LSTM has better execution compared to different models for all evaluation metrics. LSTM has performed well in all stations yet for some stations like, Dingling, Dongsi, Gucheng, and Wanliu, the RMSE values are somewhat higher than ARIMA and FBProphet. For station Dingling, the RMSE and RRSE values are on the higher side compared to different stations. The best MAPE score is 17.9, which is for Dongsi station by LSTM. Best RMSE and MAE scores are 16.6 and 11.5, respectively, for station Tiantan by LSTM. The best RRSE score is 0.23 for station Tiantan and Wanshouxigong by LSTM and CNN.

It has been observed that aside from station Dingling, the distinction in the RMSE and MAE estimations for all stations is near 6 units for each of the four models. It shows that there is some variation present in the value of the errors and extremely huge errors are not likely to have happened. Notwithstanding, for station Dingling the difference is 16 units which recommends the presence of enormous mistakes which is being punished by RMSE. The more prominent the distinction between the RMSE and MAE esteems, the bigger the change in the individual mistakes in the sample.

The average errors of all stations are given in Table III to assess overall model performance. The deep learning based model, LSTM, CNN, outperforms ARIMA, and FBProphet in terms of mean absolute percentage error. At the same time, LSTM has the lowest error in all assessment metrics compared with other models. The MAPE value is highest for FBProphet, which is 34.7.

TABLE I: Predicted Results on 12 stations

Station	Model	RMSE	MAE	MAPE	RRSE
Aotizhongxin	FBProphet	20.1	13.2	29.8	0.27
	ARIMA(1,0,3)	21.0	14.2	33.4	0.29
	LSTM	19.2	12.4	21.2	0.26
	CNN	25.0	15.9	25.5	0.34
Changping	FBProphet	18.9	13.2	37.0	0.30
	ARIMA(2,0,0)	19.2	13.3	34.8	0.30
	LSTM	18.8	12.7	28.0	0.30
	CNN	20.2	14.0	37.1	0.32
Dingling	FBProphet	35.6	16.5	45.3	0.53
	ARIMA(2,0,0)	35.1	16.2	34.8	0.52
	LSTM	35.6	15.2	22.8	0.53
	CNN	36.7	16.4	29.5	0.54
Dongsi	FBProphet	20.5	15.1	33.2	0.26
	ARIMA(1,0,0)	19.8	14.1	27.5	0.25
	LSTM	21.3	13.2	17.9	0.27
	CNN	22.0	16.1	26.9	0.28
Guanyuan	FBProphet	20.2	15.0	34.5	0.27
	ARIMA(2,0,0)	20.2	14.7	31.5	0.27
	LSTM	18.8	12.5	20.3	0.25
	CNN	21.3	14.2	25.6	0.28
Gucheng	FBProphet	20.7	15.0	36.2	0.27
	ARIMA(1,0,0)	20.6	14.6	31.2	0.27
	LSTM	22.2	15.2	23.9	0.29
	CNN	23.2	14.3	22.5	0.30

TABLE II: Predicted Results on 12 stations

Station	Model	RMSE	MAE	MAPE	RRSE
Huairou	FBProphet	18.5	12.2	35.3	0.30
	ARIMA(1,0,1)	19.9	12.3	27.6	0.33
	LSTM	17.2	11.9	32.5	0.29
	CNN	18.9	14.0	38.7	0.31
Nongzhanguan	FBProphet	20.7	15.3	35.4	0.27
	ARIMA(1,0,1)	20.3	14.4	29.4	0.27
	LSTM	18.9	13.2	20.3	0.25
	CNN	19.8	14.0	21.7	0.26
Shunyi	FBProphet	19.8	13.9	32.6	0.28
	ARIMA(4,0,0)	19.9	13.9	29.0	0.28
	LSTM	19.8	12.6	23.9	0.28
	CNN	22.2	14.0	25.6	0.31
Tiantan	FBProphet	20.6	14.7	32.6	0.28
	ARIMA(1,0,0)	19.6	13.4	26.4	0.27
	LSTM	16.6	11.5	19.6	0.23
	CNN	18.4	12.5	19.4	0.25
Wanliu	FBProphet	19.6	14.7	38.3	0.27
	ARIMA(1,0,0)	19.3	14.2	34.0	0.27
	LSTM	23.9	15.9	23.3	0.33
	CNN	23.4	13.9	22.6	0.33
Wanshouxigong	FBProphet	19.3	13.9	26.4	0.24
	ARIMA(1,0,0)	19.5	14.0	26.4	0.25
	LSTM	17.8	12.5	18.7	0.23
	CNN	18.3	12.8	21.1	0.23

TABLE III: Averaged predicted Results

Model	RMSE	MAE	MAPE	RRSE
FBProphet	21.2	14.4	34.7	0.295
ARIMA	21.2	14.1	30.5	0.308
LSTM	20.8	13.2	22.7	0.292
CNN	22.4	14.3	26.3	0.312

We can see from Fig. 7 that there isn't much difference in RMSE and MAE values. However, this difference widened for MAPE values, which show that ARIMA and FBProphet neglected to anticipate the PM2.5 levels consist of smaller peaks.

Since LSTM is giving the best results, it is important to know which activation function performed better in forecasting. The average results for the 'tanh' and 'relu' activation function used in LSTM to forecast PM2.5 levels in all 12 stations are given in Table IV. The results clearly state that LSTM performance improved with the nonlinearity of relu [21].

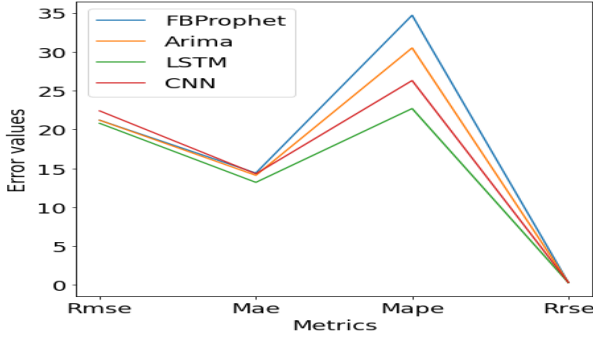


Fig. 7: Line Chart of Averaged predicted results

TABLE IV: Averaged forecast results using tanh and relu activation on LSTM

Activation	RMSE	MAE	MAPE	RRSE
tanh	22.9	14.4	25.9	0.320
relu	21.2	13.3	22.9	0.297

The fitting of PM2.5 actual and forecasted levels by all four models is featured for visualization. For this analysis, we only choose the Aotizhongxin station. The results by all four models, FBProphet, ARIMA, LSTM, and CNN, are shown in Fig. 8. For the FBProphet model, the actual data pictured as black points. The blue line represents the forecasted data with lower and upper confidence intervals in a light blue region, while for all other models, the training and test data are represented as a red line and green line, respectively. The forecast is an orange line overlapping the green line of test data. From the plots, we can easily analyze that, except for CNN, all models covered the peaks well, proving the RMSE value of CNN is slightly lower than all other models.

Subsequently, we can analyze that the LSTM and CNN covered the small peaks well compared to Arima and FBProphet, which suggests the low MAPE error in LSTM, CNN, and high MAPE error in ARIMA and FBProphet.

For clear visualization, the test data plots from March 1, 2016, to February 28, 2017, are also pictured, as shown in Fig. 9.

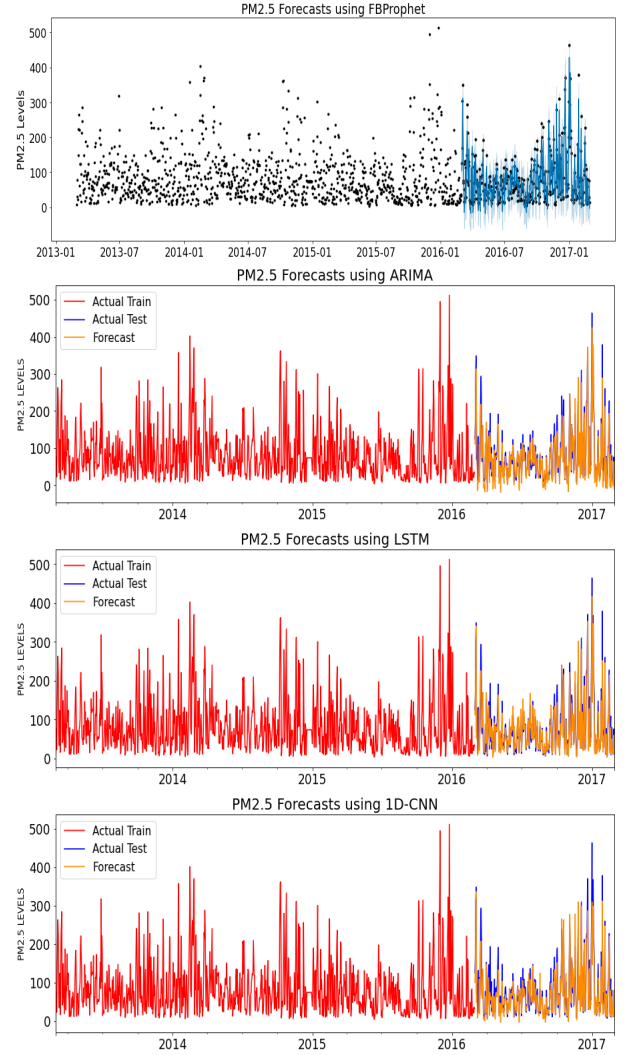


Fig. 8: Actual and forecasted values subplots of FBProphet, ARIMA, LSTM, 1D-CNN on Aotizhongxin station using full dataset.

V. CONCLUSION AND FUTURE WORK

The Heart and respiratory-related diseases like COPD, stroke, lung cancer are correlated with air pollution. Not just humans, it also affects other living organisms and damages the natural environment. This paper provided an analysis and prediction study of the PM2.5 levels on 12 station sites using four models; ARIMA, FBProphet, LSTM, and CNN. For most of the stations, LSTM performed better than all other models across RRSE, RMSE, MAPE, and MAE evaluation metrics. We had also shown that LSTM with relu activation function performed better than tanh, which opens up the possibility of spending more time analyzing hyper optimization techniques. The trend analysis shows that ARIMA and FBProphet did not predict smaller values correctly, resulting in high MAPE error. The adopted methods achieved good results. However, it restricts our examination to the model's adequacy, which can be additionally improved by considering an ensemble

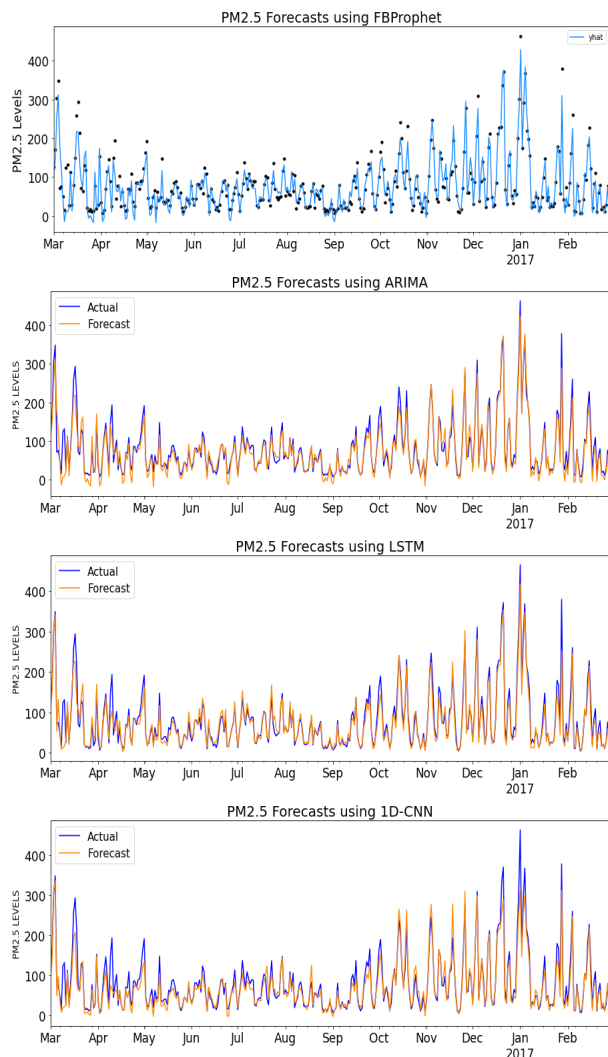


Fig. 9: Actual and forecasted values subplots of FBProphet, ARIMA, LSTM, 1D-CNN on Aotizhongxin station using test data only.

of various forecast models. The obtained forecasting results can be improved by taking feature engineering and better hyperparameter optimization into account, which will be a part of the future work.

REFERENCES

- [1] Air pollution, <https://www.who.int/news-room/air-pollution>.
- [2] Particulate Matter, <https://www.epa.gov/pm-pollution/particulate-matter-pm-basics>.
- [3] Dominici, Francesca et al. "Fine particulate air pollution and hospital admission for cardiovascular and respiratory diseases." JAMA vol. 295,10 (2006): 1127-34. doi:10.1001/jama.295.10.1127
- [4] Wulff, Shaun. (2017). Time Series Analysis: Forecasting and Control, 5th edition. Journal of Quality Technology. 49. 418-419. 10.1080/00224065.2017.11918006.
- [5] Taylor, S. and Benjamin Letham. "Forecasting at Scale." PeerJ Prepr. 5 (2017): e3190.
- [6] Alex Sherstinsky. Fundamentals of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) Network, 2018, Elsevier "Physica D: Nonlinear Phenomena" journal, Volume 404, March 2020: Special Issue on Machine Learning and Dynamical Systems; ar Xiv:1808.03314. DOI: 10.1016/j.physd.2019.132306.
- [7] Dongyang Kuang. A 1d convolutional network for leaf and time series classification, 2019; arXiv:1907.00069.
- [8] Samira, Muhammad & Salh, & Ahmed, Salah. (2014). Box –Jenkins Models For Forecasting The Daily Degrees Of Temperature In Sulaimani City. International Journal of Engineering Research and Applications (IJERA). 4. 2248-9622.
- [9] Salman Mohamadi, Farhang Yeganegi and Nasser M Nasrabadi. Detection and Statistical Modeling of Birth-Death Anomaly, 2019; arXiv:1906.11788.
- [10] Auto Arima, https://alkaline-ml.com/pmdarima/modules/generated/pmdarima.arima.auto_arima.html
- [11] Wang, Weitao & Bai, Yuebin & Yu, Chao & Gu, Yuhao & Feng, Peng & Wang, Xiaojing & Wang, Rui. (2018). A network traffic flow prediction with deep learning approach for large-scale metropolitan area network. 1-9. 10.1109/NOMS.2018.8406252
- [12] Yin, Hao & Lin, Chuang & Berton, Sebastien & Li, Bo & Min, Geyong. (2005). Network traffic prediction based on a new time series model. Int. J. Communication Systems. 18. 711-729. 10.1002/dac.721.
- [13] Xu, Z., Zhang, J., Wang, J. et al. Prediction research of financial time series based on deep learning. Soft Comput 24, 8295–8312 (2020). <https://doi.org/10.1007/s00500-020-04788-w>
- [14] Krauss, Christopher & Do, Xuan Anh & Huck, Nicolas, 2016. "Deep neural networks, gradient-boosted trees, random forests: Statistical arbitrage on the S&P 500," FAU Discussion Papers in Economics 03/2016, Friedrich-Alexander University Erlangen-Nuremberg, Institute for Economics.
- [15] Alabi, Rasheed & Siemuri, Akpojoto & Elmusrati, Mohammed. (2020). COVID-19: Easing the coronavirus lockdowns with caution. 10.1101/2020.05.10.20097295.
- [16] Zhang, Shuyi, et al. "Cautionary tales on air-quality improvement in Beijing." Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences 473.2205 (2017): 20170457.
- [17] Fillna, <https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.fillna.html>
- [18] Label Encoder, <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.LabelEncoder.html>
- [19] Validation data, <https://www.kdnuggets.com/2017/11/create-good-validation-set.html>
- [20] David N. DeJong, John C. Nankervis, N.E. Savin, Charles H. Whiteman, The power problems of unit root test in time series with autoregressive errors, Journal of Econometrics, Volume 53, Issues 1–3, 1992, Pages 323-343, ISSN 0304-4076, [https://doi.org/10.1016/0304-4076\(92\)90090-E](https://doi.org/10.1016/0304-4076(92)90090-E).
- [21] Sachin S. Talathi and Aniket Vartak. Improving performance of recurrent neural network with relu nonlinearity, 2015; arXiv:1511.03771.
- [22] Burnham, K. P., & Anderson, D. R. (2004). Multimodel Inference: Understanding AIC and BIC in Model Selection. Sociological Methods & Research, 33(2), 261–304. <https://doi.org/10.1177/0049124104268644>