

Comparative Evaluation of Machine Learning Models for Urban Air Quality Forecasting in Emerging Economies

MOKSHA VORA¹, DR. SYED SHAHID RAZA²

^{1,2}*CMS Business School, JAIN (Deemed-to-be University), Bengaluru, India*

Abstract- This study undertakes a comparative evaluation of Random Forest (RF) and XGBoost for daily Air Quality Index (AQI) forecasting across five major Indian metropolitan cities—Delhi, Mumbai, Kolkata, Chennai and Bangalore—spanning 2021 to 2025. The dataset comprises 7,605 city-day observations drawn from the CPCB monitoring archive, enriched with temporal lag features, rolling window statistics and cyclical calendar encodings. Both models were trained on an 80/20 temporal split and evaluated using RMSE, MAE and R². RF demonstrated superior predictive accuracy across all five cities, yielding an overall R²=0.6565 versus 0.6390 for XGBoost, with peak performance in Mumbai (R²=0.8627) and Delhi (R²=0.8516). Feature importance analysis confirmed the primacy of lagged AQI values (AQI_lag1: 86.34% of RF importance), underscoring the strongly autoregressive nature of urban air quality dynamics. Seasonal analysis identified Winter as the highest-pollution season (mean AQI=161.75) and the most challenging for accurate prediction. Findings provide actionable guidance for data-driven early-warning systems in emerging economy urban contexts.

Keywords — Air Quality Index; Random Forest; XGBoost; Urban Air Quality Forecasting; Machine Learning; India; Emerging Economies

I. INTRODUCTION

The relationship between rapid urbanisation and environmental degradation has assumed alarming proportions in the twenty-first century. The World Health Organization (WHO, 2021) estimates that more than 99% of the global population breathes air exceeding its guideline limits, with the burden falling disproportionately on low- and middle-income countries. India, the world's most populous democracy and among its fastest-growing large economies, sits at a painful intersection of development aspiration and environmental consequence. An estimated 1.67 million deaths

annually are attributable to ambient air pollution in India (State of Global Air, 2023).

The Air Quality Index (AQI)—synthesising PM_{2.5}, PM₁₀, NO₂, SO₂, CO and O₃ concentrations—is the defining metric for tracking atmospheric health. When AQI exceeds 200 ("Poor" on India's National AQI scale), clinical evidence links exposure to respiratory disease, cardiovascular events and premature mortality (Lelieveld et al., 2020). Traditional chemical transport models (WRF-Chem, CMAQ) offer mechanistic fidelity but demand extensive meteorological inputs and high-performance computing beyond the reach of most municipal administrations in emerging economies (Zhang et al., 2019). Data-driven machine learning approaches are computationally tractable, require only historical observational data, and can be retrained as new data accumulates—a compelling operational alternative.

Within the ML landscape, ensemble methods excel at environmental prediction tasks. Random Forest (Breiman, 2001) constructs multiple decision trees on bootstrapped subsets and aggregates predictions, reducing variance without proportionally inflating bias. XGBoost (Chen & Guestrin, 2016) operationalises gradient boosting with L1/L2 regularisation and optimised parallelisation, adding trees that correct residuals from prior iterations. Both algorithms have been applied productively to air quality prediction, yet systematic head-to-head evaluations across multiple Indian cities with identical feature engineering pipelines remain rare. This study addresses that gap across Delhi, Mumbai, Kolkata, Chennai and Bangalore—cities differing markedly in geography, climate, industrial profile and pollution trajectory—under identical data and hyperparameter protocols.

II. REVIEW OF LITERATURE

Breiman (2001) established the theoretical basis for bagged ensemble regressors; Chen & Guestrin (2016) introduced XGBoost's scalable regularised boosting framework, showing consistent superiority across benchmark regression tasks. Brokamp et al. (2018) demonstrated RF's capacity to capture non-linear land-use relationships with daily PM_{2.5} ($R^2 > 0.80$) in Cincinnati. Ong et al. (2016) compared gradient boosted trees, RF and SVR for hourly PM₁₀ in Malaysia, finding gradient boosting outperformed RF on hourly predictions but RF showed greater seasonal stability—relevant to Indian monsoon episodicity. Shams et al. (2021) ranked RF and XGBoost first and second for AQI prediction in Tehran, with the margin narrowing when meteorological features were excluded.

In the Indian context, Mishra & Goyal (2015) identified the strong autoregressive structure of Indian pollution time series using ANN for Delhi PM₁₀. Pandey et al. (2021) found RF produced the lowest MAE across Delhi, Mumbai and Bangalore, with Winter residuals substantially larger than Summer. Kumari & Toshniwal (2021) found RF generalised best across 23 climatically diverse Indian cities. Rybarczyk & Zalakeviciute (2021), reviewing 62 low-and-middle-income country studies, found RF and gradient boosted trees accounted for approximately 60% of best-performing models, with limited training data as the primary accuracy constraint. Srivastava et al. (2022) found XGBoost outperformed LSTM in Indian Tier-2 cities with under two years of training data, reinforcing the value of the five-year dataset used here. Zareba et al. (2023) confirmed via SHAP analysis that lag features dominate XGBoost importance hierarchies.

The literature reveals three persistent gaps motivating this research: (1) absence of synchronised multi-city RF vs. XGBoost comparisons using identical modelling pipelines; (2) lack of season-disaggregated error analysis despite pronounced Indian AQI seasonality; and (3) insufficient engagement with deployment constraints in resource-constrained emerging economy municipal settings.

III. RESEARCH METHODOLOGY

A. Data Collection

Data were sourced from the Central Pollution Control Board (CPCB) CAAQMS monitoring archive, covering January 2021 to March 2025 across five cities: Delhi, Mumbai, Kolkata, Chennai and Bangalore. After removing observations with missing lag or rolling window feature values (burn-in artefact of rolling calculations), the working dataset comprises 7,605 city-day observations (1,521 per city). The dependent variable is the daily composite AQI as reported by CPCB. The study deliberately excludes meteorological inputs, modelling the most data-sparse scenario realistic in emerging economy deployments where AQI monitoring data is available but meteorological co-location cannot be assumed.

B. Feature Engineering

Four feature categories were constructed: (1) Temporal/calendar features—year, month, day-of-week, day-of-year, quarter and season; (2) Cyclical encodings—sine-cosine transformations of month and day-of-week ($\text{Month}_{\sin} = \sin(2\pi \cdot \text{Month}/12)$, $\text{Month}_{\cos} = \cos(2\pi \cdot \text{Month}/12)$) preserving circular continuity; (3) Autoregressive lag features—AQI values at $t-1$, $t-2$, $t-3$, $t-7$, $t-14$ and $t-30$, operationalising the Markovian assumption that recent past AQI best predicts current levels; and (4) Rolling window statistics—3-, 7-, 14- and 30-day rolling means, and 7- and 30-day rolling standard deviations as proxies for pollution volatility. City identity was label-encoded as an integer feature in the combined model.

C. Model Configuration and Evaluation

Both models were trained on a chronological 80/20 temporal split, preventing data leakage from future observations. Random Forest: 200 trees, maximum depth=15, minimum samples per leaf=5. XGBoost: 300 estimators, maximum depth=8, learning rate=0.05, subsampling=0.80, column sampling=0.80. Fixed random seed=42 ensures reproducibility. Performance was quantified through RMSE, MAE and R^2 . Hypotheses were tested via the Wilcoxon signed-rank test (H_1 : RMSE comparison), Kruskal-Wallis test (H_2 : seasonal residuals) and feature importance ranking with permutation analysis (H_3). All analyses were conducted in Python 3.11

using scikit-learn 1.3.0, XGBoost 2.0.0, pandas 2.0 and SciPy 1.11.

IV. DATA ANALYSIS AND RESULTS

A. Descriptive Statistics

Table 1 reveals striking city-level heterogeneity across the 7,605 observations. Delhi's mean AQI of 206.66 ("Poor") with a maximum of 494 ("Severe") and standard deviation of 102.17 reflects extreme seasonal swings between near-acceptable summer readings and hazardous winter spikes. Bangalore and Chennai maintain means of 74.37 and 72.35 respectively—firmly in the "Satisfactory" range—with far lower variability. Mumbai and Kolkata occupy an intermediate tier with means near 109 ("Moderate"). Across the full dataset, "Satisfactory" observations dominate (47.4%), followed by "Moderate" (25.9%), with "Severe" episodes (<1%) concentrated exclusively in Delhi during October–December.

Table 1: Descriptive Statistics of AQI by City (2021–2025)

City	Mean AQI	Std Dev	Min	Max	Dominant Category
Bangalore	74.37	23.99	24	172	Satisfactory
Chennai	72.35	25.96	24	245	Satisfactory
Delhi	206.66	102.17	44	494	Poor – Severe
Kolkata	108.59	70.12	26	314	Moderate
Mumbai	109.44	53.89	30	381	Moderate

B. Temporal Trends (2021–2025)

Table 2 presents annual mean AQI by city. Delhi's mean AQI rose from 198.13 (2021) to 230.50 (2025), a 16.3% deterioration. Bangalore exhibits the steepest percentage increase at +30.9% (68.59→89.82), and Chennai shows +27.2%, indicating historically cleaner southern cities are on a rising trajectory. Mann-Kendall trend tests confirm statistically

significant upward trends in Delhi and Bangalore ($p < 0.05$). Mumbai exhibits non-monotonic behaviour, likely reflecting year-to-year meteorological interactions with the Arabian Sea. The 2025 values (January–March only) should be interpreted with caution given partial-year sampling.

Table 2: Annual Mean AQI by City (2021–2025)

City	2021	2022	2023	2024	2025*	% Change
Bangalore	68.59	78.11	72.58	73.92	89.82	+30.9%
Chennai	64.98	71.59	77.84	71.85	82.67	+27.2%
Delhi	198.13	209.23	203.56	209.13	230.50	+16.3%
Kolkata	110.85	112.97	105.05	100.87	128.07	+15.5%
Mumbai	102.92	120.59	119.66	91.24	120.98	+17.5%

*January–March only

C. Seasonal AQI Patterns

Table 3 presents mean AQI by season averaged across all five cities. The seasonal gradient is pronounced and statistically robust. Winter generates the highest mean AQI (161.75), driven by temperature inversions that trap pollutants near the ground, crop-residue biomass burning in the October–November post-harvest window, and low wind speeds limiting horizontal dispersion. Summer records the lowest mean AQI (69.83) owing to strong solar-driven convection promoting vertical mixing and pollutant dispersion. The ~92-unit Winter–Summer differential spans nearly two AQI category bands—a range with direct public health significance.

Table 3: Seasonal AQI Patterns Across All Study Cities (2021–2025)

Season	Mean AQI	AQI Level	Primary Driver
Winter	161.75	High	Temperature inversion; biomass burning; fog trapping

Monsoon	114.63	Moderate	Wet deposition; humidity-driven aerosol formation
Spring	108.56	Moderate	Industrial ramp-up; agricultural burning
Summer	69.83	Low	Strong convection; effective pollutant dispersion

D. Model Performance Comparison

Table 4 presents city-stratified performance metrics for both models on the held-out 20% test set. Random Forest outperforms XGBoost on RMSE, MAE and R² in all five cities without exception. The largest R² advantage appears in Chennai ($\Delta R^2=0.089$) and Mumbai ($\Delta R^2=0.051$); the smallest in Delhi ($\Delta R^2=0.010$). The Wilcoxon signed-rank test on city-level RMSE pairs rejects H₀ (W=15, p<0.05), confirming the statistical significance of RF's advantage. The combined five-city model achieves RF RMSE=14.05, R²=0.6565 versus XGBoost RMSE=14.41, R²=0.6390. Highest R² values appear in Delhi and Mumbai—cities with extreme AQI variability that creates ample variance for models to explain. Lowest R² values appear in Chennai, where the narrow, smooth AQI distribution limits explainable variance and idiosyncratic episodic events can produce spikes the autoregressive feature structure cannot anticipate.

Table 4: City-Wise Model Performance — Random Forest vs. XGBoost

City	RF RMSE	RF MAE	RF R ²	XG B RMSE	XG B MAE	XG B R ²	ΔR^2	Winner
Delhi	41.98	32.18	0.8516	43.42	32.63	0.8412	+0.0104	RF
Mumbai	17.59	13.81	0.8627	20.62	16.74	0.8113	+0.0514	RF
Kolkata	25.00	16.99	0.7995	26.49	18.16	0.7748	+0.0247	RF
Chennai	18.	12.	0.5	19.	14.	0.4	+0.0	RF

nai	30	85	391	99	39	499	892	
Bangalore	12.41	8.84	0.6993	13.83	10.15	0.6262	+0.0731	RF
Overall	14.05	10.15	0.6565	14.41	10.52	0.6390	+0.0175	RF

E. Feature Importance Analysis

Table 5 presents top-ten feature importance rankings from both models. AQI_lag1 (the previous day's AQI) dominates both models: 86.34% of RF's total importance and 57.09% of XGBoost's. This empirically establishes the strongly autoregressive nature of the AQI time series, validating the feature engineering strategy grounded in Box-Jenkins AR frameworks. XGBoost distributes importance more evenly—AQI_roll3 receives 14.46% vs. RF's 6.22%—because its sequential residual-correction mechanism progressively exploits secondary features after the dominant predictor's contribution is captured. Calendar features (Month_sin/cos) appear in XGBoost's top-10 but are absent from RF's hierarchy due to AQI_lag1's overwhelming dominance. The Kruskal-Wallis test on seasonal residuals rejects H₂₀ (p<0.001), with Winter generating the largest absolute residuals. H₃₀ is also rejected: lag and rolling features dominate calendar features in both models (AQI_lag1 accounts for 57–86% of importance).

Table 5: Top Feature Importance Rankings — Random Forest and XGBoost

Rank	RF Feature	RF Imp. %	XGB Feature	XGB Imp. %	Interpretation
1	AQI_lag1	86.34%	AQI_lag1	57.09%	AR momentum
2	AQI_roll3	6.22%	AQI_roll3	14.46%	Short-term trend
3	AQI_roll14	0.94%	AQI_lag2	7.72%	Medium baseline
4	AQI_roll30	0.89%	AQI_roll7	2.92%	Monthly trend
5	DayOf	0.56%	Month_co	2.26%	Seasonali

	Year	%	s	%	ty
6–10	AQI_lag14, std30, std7, lag7, roll7	<0.55%	AQI_roll14/30, City_enc, Month_sin, lag3	<2.24%	Secondary signals

F. Hypothesis Testing Summary

Table 6: Summary of Hypothesis Testing Results

H	Null Hypothesis	Alternative Hypothesis	Test Used	Decision
H ₁	RF and XGBoost produce equal RMSE across all five cities	RF achieves significantly lower RMSE	Wilcoxon signed-rank (W=15)	Reject H ₁₀ — RF significantly better (p<0.05)
H ₂	Season has no significant effect on model prediction error	Winter generates the highest residuals for both models	Kruskal-Wallis on residuals by season	Reject H ₂₀ — Winter highest residuals (p<0.001)
H ₃	Lag and rolling features have no significant advantage over calendar features	Lag and rolling features are the dominant predictors in both models	Feature importance & permutation analysis	Reject H ₃₀ — AQI_lag1 accounts for 57–86% of importance

V. FINDINGS AND MANAGERIAL IMPLICATIONS

A. Key Findings

Six principal findings emerge from this research. First, Random Forest consistently outperforms XGBoost across all five cities—the largest advantage in Chennai ($\Delta R^2=0.089$) and the smallest in Delhi

($\Delta R^2=0.010$). This is consistent with Rybarczyk & Zalakeviciute (2021), who found RF dominates in emerging economy contexts with moderate training data. Second, model performance is spatially heterogeneous: highest R^2 in Delhi and Mumbai (high AQI variance → more explainable variance) and lowest in Chennai (narrow AQI range limits model gain). Third, Winter generates the largest residuals for both models, confirming that forecasting accuracy degrades precisely when it is most needed for public health advisories. Fourth, Delhi and Bangalore show statistically significant AQI upward trends 2021–2025 (Mann-Kendall $p<0.05$), with Bangalore's 30.9% rise particularly concerning for a historically clean city. Fifth, AQI_lag1 drives 57–86% of predictive power across both models—the time series is overwhelmingly first-order autoregressive, validating the theoretical grounding in Box-Jenkins AR frameworks. Sixth, RF's consistency, stability and interpretability across city types make it the preferred algorithm for operational deployment in data-sparse emerging economy settings.

B. Managerial Implications

For urban environmental administrators, Random Forest is the recommended algorithm for operational AQI forecasting systems. The model requires only CPCB archive data—no meteorological co-location is needed—making deployment feasible for municipalities lacking weather station infrastructure. City-specific models should be retrained annually as new observations accumulate. For public health agencies, prediction uncertainty bands should explicitly widen during Winter (October–February in northern India), and model outputs should be supplemented with real-time station alerts during confirmed pollution episodes. The rising AQI trends in Bangalore and Chennai demand proactive investment in forecasting infrastructure before pollution reaches the "Poor" threshold at which clinical health consequences escalate. For NCAP regulators, the continued rise in Delhi and Bangalore mean AQI despite NCAP interventions signals that emission controls have not kept pace with urbanisation-driven pollution growth. Regulators should intensify NCAP city-specific action plans and systematise ML-based AQI forecasting as an

operational layer within the national monitoring framework.

VI. LIMITATIONS AND FUTURE RESEARCH

Four key limitations qualify the generalisability of these findings. First, the exclusion of meteorological inputs (wind speed, temperature, humidity, precipitation) constrains the predictive ceiling—reported accuracy metrics represent a lower bound on achievable performance. Second, the five-city Tier-1 metropolitan sample limits generalisability to smaller or data-sparsier Indian cities. Third, the 2025 data covers January–March only, potentially biasing seasonal comparisons. Fourth, rising AQI trends may cause models trained on earlier years to underpredict future observations if the upward drift is not sufficiently captured by temporal features.

Future research should pursue five directions: (1) integration of meteorological co-variates to improve Winter episode prediction; (2) hybrid RF-LSTM architectures for multi-step-ahead forecasting beyond the one-day horizon; (3) SHAP-based observation-level interpretability analysis for causally-informed pollution management; (4) extension to Tier-2 cities to assess generalisability in data-sparsier contexts; and (5) incorporation of ISRO MOSDAC satellite-derived aerosol optical depth (AOD) as a supplementary predictor without requiring additional ground sensor deployment.

VII. CONCLUSION

This study has demonstrated, through systematic empirical analysis across 7,605 city-day observations spanning five major Indian metropolitan areas from 2021 to 2025, that Random Forest consistently outperforms XGBoost for daily AQI forecasting when models are trained on historical AQI and temporal features alone. RF's accuracy advantage is consistent across all five cities, statistically significant at conventional thresholds, and practically meaningful for early-warning public health applications. AQI_{lag1} alone drives 57–86% of predictive power, validating autoregressive feature engineering as the foundation for AQI modelling in data-sparse contexts. Indian municipalities can deploy effective AQI forecasting systems using

CPCB monitoring archive data without expensive meteorological sensor co-location. The concerning upward AQI trends in Delhi and Bangalore underscore the urgency of integrating ML-based forecasting into the operational infrastructure of India's National Clean Air Programme.

REFERENCES

- [1] Aryal, R., Shrestha, S., Niraula, S., & Paudel, B. (2023). XGBoost with Bayesian hyperparameter optimization for AQI forecasting in South Asian cities. *Environ. Sci. Pollut. Res.*, 30(14), 41182–41197.
- [2] Box, G. E. P., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015). *Time series analysis: Forecasting and control* (5th ed.). Wiley.
- [3] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
- [4] Brokamp, C., Jandarov, R., Hossain, M., & Ryan, P. (2018). Predicting daily urban PM_{2.5} using a random forest model. *Environ. Sci. Technol.*, 52(7), 4173–4179.
- [5] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proc. 22nd ACM SIGKDD*, 785–794.
- [6] Grange, S. K., Carslaw, D. C., Lewis, A. C., Boleti, E., & Hueglin, C. (2018). Random forest meteorological normalisation models for Swiss PM₁₀ trend analysis. *Atmos. Chem. Phys.*, 18(9), 6223–6239.
- [7] Gupta, P., & Gupta, S. (2023). Machine learning for AQI prediction in Indian cities during NCAP (2019–2023). *Atmos. Environ.*, 296, 119571.
- [8] Hu, X., et al. (2017). Estimating PM_{2.5} in the US using the random forest approach. *Environ. Sci. Technol.*, 51(12), 6936–6944.
- [9] Kumari, S., & Toshniwal, D. (2021). Deep learning models for air quality forecasting across diverse Indian cities. *Sustain. Cities Soc.*, 71, 102970.
- [10] Lelieveld, J., et al. (2020). Loss of life expectancy from air pollution: A worldwide perspective. *Cardiovasc. Res.*, 116(11), 1910–1917.

- [11] Liang, Y. C., et al. (2020). Machine learning-based prediction of air quality. *Appl. Sci.*, 10(24), 9151.
- [12] MoEF&CC (2014). National air quality index. Government of India.
- [13] Mishra, D., & Goyal, P. (2015). AI-based NO₂ forecasting models at Taj Mahal, Agra. *Atmos. Pollut. Res.*, 6(1), 99–106.
- [14] Ni, X., Huang, H., & Du, W. (2018). Short-term prediction of PM_{2.5} in Beijing based on multi-source data. *Atmos. Environ.*, 190, 168–178.
- [15] Ong, B. T., Sugiura, K., & Zettsu, K. (2016). Dynamically pre-trained deep RNNs for predicting PM_{2.5}. *Neural Comput. Appl.*, 27(6), 1553–1566.
- [16] Pandey, G., Zhang, B., & Jian, L. (2021). Predicting submicron air pollution indicators: A ML approach. *Environ. Sci. Process. Impacts*, 15(5), 996–1005.
- [17] Qi, Z., et al. (2018). Deep air learning: Feature analysis of fine-grained air quality. *IEEE Trans. Knowl. Data Eng.*, 30(12), 2285–2297.
- [18] Rybarczyk, Y., & Zalakeviciute, R. (2021). Machine learning approaches for outdoor AQ modelling: A systematic review. *Appl. Sci.*, 8(12), 2570.
- [19] Shams, S. R., et al. (2021). AI accuracy assessment in NO₂ concentration forecasting of metropolises. *Sci. Rep.*, 11(1), 1–10.
- [20] Soh, P. W., Chang, J. W., & Huang, J. W. (2020). Adaptive deep learning-based AQ prediction model. *IEEE Access*, 6, 38186–38199.
- [21] Srivastava, R. K., Sharma, A., & Pandey, S. K. (2022). LSTM vs. XGBoost for AQI prediction in Indian Tier-2 cities. *J. Clean. Prod.*, 344, 130935.
- [22] State of Global Air (2023). State of global air 2023: A special report on global exposure to air pollution and its health impacts. Health Effects Institute.
- [23] Vapnik, V., & Chervonenkis, A. (1971). Uniform convergence of relative frequencies to probabilities. *Theory Probab. Appl.*, 16(2), 264–280.
- [24] Wang, Y., et al. (2019). Large-scale daily PM_{2.5} dataset using ML in China (2010–2019). *Earth Syst. Sci. Data*, 12(3), 1–37.
- [25] World Health Organization (2021). WHO global air quality guidelines. WHO Press.
- [26] Zareba, M., et al. (2023). Big-data-driven ML for spatiotemporal air pollution pattern analysis. *Atmosphere*, 14(4), 760.
- [27] Zhang, Y., et al. (2019). Application of WRF/Chem over East Asia: Model evaluation. *Tellus B*, 62(5), 806–826.