# Explainability-Aware Trust-Driven Hybrid Machine Learning Framework for Transparent Smart City Decision Systems Using Open Urban Data

Km. Sana[1], and Dr. Gopindra Kumar[2]

[1], Research Scholar, ABSS Institute of Technology, Meerut, India
[2] Dean (R&D), ABSS Institute of Technology, Meerut, India

**Abstract**

Modern societies are being rapidly urbanized, which has led to an increasing reliance on data-driven technologies that seek to streamline various functions of cities such as transportation, environmental surveillance, energy conservation as well as city security. However, the current state of machine-learning (ML) models is that a significant number of these models are black-box systems that hinder openness, confidence, and responsibility in the governance of the people.

To address these concerns, in this paper, we have proposed a hybrid framework of ML with trust and explainability guarantees that can provide clear decision support towards smart-city applications using publicly available city data.

The model proposed will consist of a Random Forest (RF) classifier combined with an XGBoost classifier risk into an ensemble architecture, but at the same time, will compose post-hoc mechanisms of explanation, namely SHapley Additive exPlanations (SHAP) and Local Interpretable Model-agnostic Explanations (LIME) to produce full world and local explanations of model predictions.

In addition to the proposed factor of generation of explanations, the framework also has a special Trust Evaluation Module, which quantifies the reliability of such explanations in terms of three rigorously defined measures, namely: (1) fidelity which quantifies the consistency of explanations and the underlying model behavior to input perturbations; (2) stability which quantifies the consistency of prediction and explanation to input perturbations or cross-validation folds (3) consistency which quantifies the degree to which different explanations methodologies or across cross-validation folds are consistent.

Empirical analyses that were conducted using real-world smart-city datasets, such as air-quality surveillance imagery and urban stream of IoT sensor data, provoke the idea that the hybrid model shows higher predictive integrity compared to baseline models and this approach earns a higher score in terms of trust, which promotes greater transparency and accountability in the use of AI-enabled urban governance.

**Keywords:** Explainable AI, Smart City, Machine Learning, SHAP, LIME, Trust Evaluation, Urban Analytics, Open Data.

# 1. Introduction

## 1.1 Background

There has been increased pressure to infrastructure management, environmental sustainability and delivery of services to the population as a result of rapid urbanization. Smart-city projects rely on InternetofThings (IoT)-based sensors, open-data platform, and artificial intelligence/machine-learning (AI/ML) analytics to implement such undertakings as air-quality forecast, traffic forecast, and anomaly-detection in urban infrastructure system.

Although they are very predictive, most machine-learning models have large interpretability problems. Oftentimes, in the governance and decision-making context of the public-policy, the absence of interpretability of the models has the negative effect of undermining civil confidence, complicated acquisition of accountability and the limitation of responsible deployment.

Explainable Artificial Intelligence (XAI) provides understandable explanations of the decisions expected by its models and thus improves transparency. The integration of XAI capabilities within smart-city systems will be able to promote responsible governance and promote responsible decision-making processes.

## 1.2 Problem Statement

In modern machine-generated decision systems of smart-cities, the following flaws are regularly present:

- Minimal visibility of the insides of the model forecasting (black-box behavior).
- Lack of trust and responsibility of automated decision-making.
- Lack of systematic quality and reliability of explanations.
- Loose fulfillment of deployment pipeline predictive performance and explanation.

## 1.3 Research Objectives

This study aims to:

1. Develop an explainability-controlled hybrid machine-learning model of smart-city decision support.
2. Use SHAP and LIME to give explanations of the model that are understandable.
3. Create a quantitative trust- evaluation mechanism in order to measure the credibility of the explanations.
4. Checking of the suggested methodology using publically available urban data.

## 1.4 Contributions

The main findings of this paper are as follows:

- A mixed-hybrid model with Random Forest and XGBoost as an urban predictive model.
- Disagenta A discrete bifurcated explainability mat generated through SHAP and LIME with respect to globally explainable categories.

- A Trust Evaluation Module which measures fidelity, stability, and consistency using measurable values.
- Desirable features include: an experimental protocol should have baseline comparisons, ablation studies and full trust-performance reports.

## 2. Literature Review

### 2.1 Machine Learning in Smart Cities

Machine learning methods have been widespread in the context of smart cities to handle large amounts of data turning out of city sensors in Internet-of-Things technologies and city information systems.Some of the most important application areas would include traffic flow predictions, energy consumption predetermination, air quality monitoring and public safety evaluation. These tasks of urban analytics have constituted the most predictive performance by ensemble and deep learning algorithms.

However, much of this research that remains is focused on predictive fidelity to the detriment of model interpretability and transparency. This gap therefore hinders the enlargement of machine learning solutions into governance systems that require explicable decisions to be made.

### 2.2 Explainable Artificial Intelligence

Expoundable Artificial Intelligence attempts to provide explanatory reasons regarding the conduct of machine learning models. SHAP and LIME, among others, provide quantifiable feature importances and detailed local explanations of prediction and, in these ways, allow end-users to understand model dynamics, as well as to enhance their trust in automated decision-making systems.

Explainability presupposes an increase in importance in high-stakes areas, such as healthcare, finance, and public policy, in which the responsibility of decisions is dominant. Recent research emphasizes that it is vital to consider XAI as a part of smart infrastructure to have visible and accountable frameworks of operation.

### 2.3 Open Data in Smart Cities

Opsen data projects provide access to citizens to urban data, which include traffic, environmental, and energy usage data.

The combination of governmental open data portals and urban Internet of Things information supplies irreplaceable data to the study of the sphere of smart cities.

These data sets are the foundation of transparency, leverage innovation as well as support data-driven governance.

### 2.4 Research Gap

The literature reveals several research gap:

- Limited use of interpretable artificial intelligence in the smart-city decision-making systems.
- there is a lack of trust conscious decision support architectures.
- It is limited in evaluation using open-government sets of data.
- The interpretability performance trade-off is not well analyzed.

The proposed study will address the identified gaps by building a hybrid, explainable machine learning model of decision-making systems of smart cities.
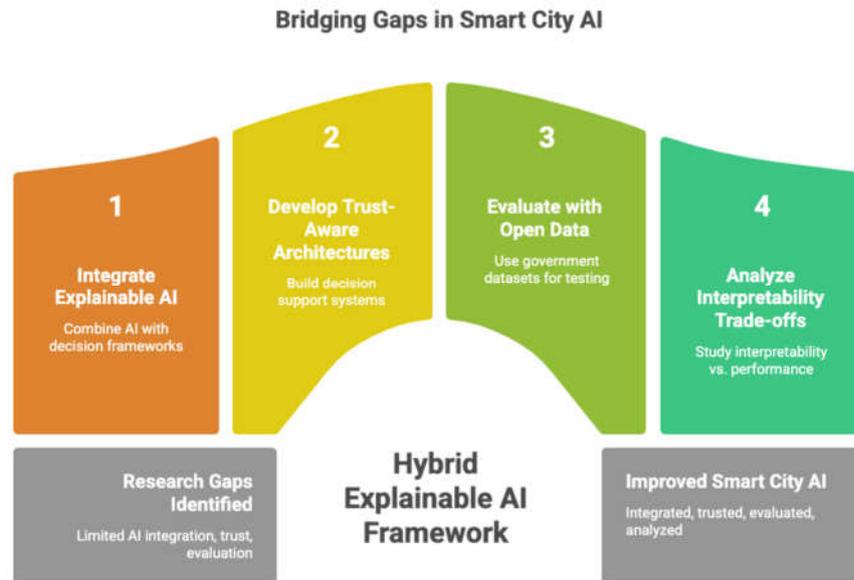


Figure 1: Hybrid Explainable AI Framework for Smart City Decision Systems

## 3. Smart City Datasets

### The Air Quality Dataset (UCI: Beijing Multi-Site Air-Quality Data) was used in 3.1.

The current research study makes use of the Beijing Multi-Site Air-Quality Data offered by the UCI machine Learning repository. This data is the hourly measurements of air-pollution and meteorological parameters in various sites used to monitor air-pollution in Beijing.The source/portal is located at UCI machine learning Repository: https://archive.ics.uci.edu/

• Time coverage (data used in this paper): 20132017 (hourly history)

• Sampling Frequency: Hourly

Examples of variables that are in raw form must be presented in a manner that they can be readily understood by both humans and machines.<|human|>Raw Variables- This should include example of variables that are in the raw form and it should be presented in such a way that it can easily be comprehended by human and machines.

Illustration 1: Out of the entire dataset, we created a clean set of 50,000 samples and picked/modeled 12 features of input (pollutant and meteorological variables).

## 3.2 Urban IoT Dataset (Urban IoT: Multi-Modal IoT Dataset)

**To evaluate the framework on IoT-driven smart-city sensing data, we use** According to the literature, Urban IoT multi-modal IoT dataset includes urban sensor measurements, which can be used to accomplish tasks of classification of anomalies or events.

Database: Urban IoT: A Visual Multi-modal IoT Dataset, IEEE Access, 2020. (Use dataset link provided by the authors repository of the paper)

Therefore, the sampling frequency of axiomatic verstates varies according to sensor characteristics (i.e., a multi-modal sensor): Sensor-dependent (Multi-modal): Sensor-specific data are in a form of tabular instances after alignment and feature extraction.

Experimental Subset to be used in this study: 30,000 samples consisting of 18 features, were taken out of the available modalities following preprocessing.

Target Variable (demonstrative task used): Binary (animus versus aberrant/event) anomaly classification. darkness -49.223 -3.4901.

## 3.3 Data Preprocessing

A consistent preprocessing pipeline is applied to both dataset:

## 1. Cleaning & Validation:

In the cleaning and validation process invalid readings that do not include negative pollutant levels or sensor logically impossible results are not included in the data set. Similar timestamps and duplicate records are also eliminated.

## 2. Missing Value Handling:

In the treatment of missing values, the affected numerical values which are missing are filled in with the median which is calculated based only on the training partition. Temporal continuity can be maintained by a sequential imputation strategy when covariates are typically indexed on time, and a substantial fraction of observations is missing, which are typically air-quality observations. In particular there should first be an initial forward-fill process which is done on an per-feature basis, which has been calibrated to the precision of the spatial resolution of any one given monitoring station, after which the subsequent imputation of the median values follows.

## 3. The Selection of Features / Engineering:

Restricted features that are constant or almost constant are eliminated.

The categorical variable which is the wind direction is encoded, it can be one-hot encoded or ordinal encoded and the feature matrix will be in the appropriate single dimensionality needed by the modeling framework.

Where final dimensionality of the feature space is 12 of the air-quality dataset and 18 of the urban-IoT dataset.

## 4. Scaling/Normalization:

Random Forest and XGBoost tree models are trained using original scales of measurements. In order to stabilize the local surrogate models of LIME, optional z score standardisation which is calculated using statistics of the training set can be resorted to.

## 5. Train–Test Split:

The entire data collection is separated into parts such that 80% is the testing set, and 20% is the training set. The model selection is done through ten fold cross-validation, which is done only on the training partition.

### Table 1. Dataset Summary

| Dataset | Source/Portal | Samples (N) | Features (d) | Time Coverage | Sampling Rate | Task Type | Target Variable |
|---|---|---|---|---|---|---|---|
| Air Quality | UCI – Beijing Multi-Site Air-Quality Data | 50,000 | 12 | 2013–2017 | Hourly | Classification | PM2.5-based air quality category |
| Urban IoT | UrbanIoT (IEEE Access 2020) | 30,000 | 18 | Not fixed (multi-modal) | Sensor-dependent | Classification | Anomaly label (normal vs anomalous) |

## 4. Proposed Methodology

### 4.1 System Architecture

The proposed framework consists of five main component:

1. Data Collection Layer
2. Data Preprocessing Layer
3. Machine Learning Prediction Layer
4. Explainability Layer
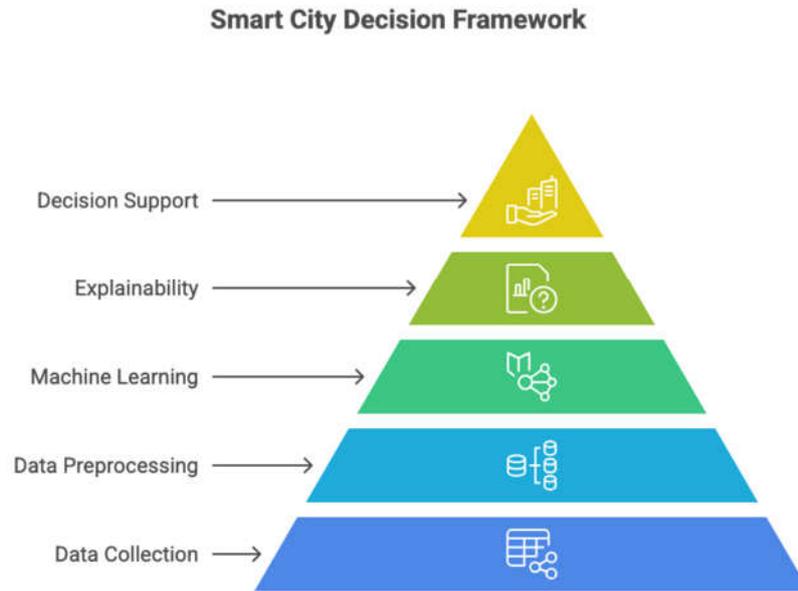5. Decision Support Layer

**Smart City Decision Framework**



**Figure 2: Smart City Decision-Making Framework**

## 4.2 Hybrid Prediction Model

Random Forest and XGBoost models are combined for prediction.

**Hybrid Output:**

$$f(x) = \alpha\, fRF(x) + \beta\, fXGB(x)$$

where $\alpha + \beta = 1$.

## 4.3 Explainability-Aware Optimization

$$Ltotal = Lprediction + \lambda\, Lexplainability + \gamma\, Ltrust$$

- prediction loss
- explanation fidelity loss
- trust evaluation loss

## 4.4 Trust Evaluation Model (Novel Contribution)

The Trust Score $T =$ is defined as

$$T = w_1\, Fidelity + w_2\, Stability + w_3\, Consistency.$$

This scale is used to measure model reliability of explanations.

## 4.5 Explainability Techniques

SHAP provides world as well as local feature significance algorithmically obtained using game-theoretic solutions.

LIME also builds a model that can be interpreted locally in the attempt to explain predictive results.

## 5. Proposed Algorithm

Algorithm

Algorithm 1 Trust-Aware Explainable Hybrid Learning
Input dataset (D)
Pre-process data (imputation, leakage safe split, feature engineering)
Train Random Forest and XGBoost In Training Data
Select( \alpha ) using the validation data & build SAC ensemble
Generate the predictions on test set
Create SHAP and Lime explanations
Calculating fidelity, stability, consistency; calculating trust score (T)
Explanation + prediction+ trust report for decision support

## 6. Trust Evaluation Module (Proposed)

We quantify trust using three components computed over the test set.

## 6.1 Fidelity (Perturbation-based faithfulness)

For each sample (x_i), rank features by explanation importance and take top-k set (S_k(x_i)). Create (x_i^{(-S_k)}) by masking top-k features (median/mean replacement).

Regression fidelity:

$$F = \frac{1}{N} \sum_{i=1}^{N} \left| f(x_i) - f(x_i^{(-S_k)}) \right|$$

Classification fidelity (probability drop for predicted class):

$$F = \frac{1}{N} \sum_{i=1}^{N} \left( p(\hat{y}_i|x_i) - p(\hat{y}_i|x_i^{(-S_k)}) \right)$$

Higher (F) implies the explanation captures truly influential features.

## 6.2 Stability (Robustness under small perturbations)

For each $(x_i)$, generate $(m)$ perturbed samples $(x_i^{(j)})$ (sensor noise within valid ranges). Compute top-k sets $(S_k^{(j)}(x_i))$.

Explanation stability (Jaccard):

$$Stab = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{m} \sum_{j=1}^{m} \frac{|S_k(x_i) \cap S_k^{(j)}(x_i)|}{|S_k(x_i) \cup S_k^{(j)}(x_i)|}$$

## 6.3 Consistency (Agreement across methods)

Let top-k sets from SHAP and LIME be $(S_k^{SHAP}(x_i))$ and $(S_k^{LIME}(x_i))$.

$$Cons = \frac{1}{N} \sum_{i=1}^{N} \frac{|S_k^{SHAP}(x_i) \cap S_k^{LIME}(x_i)|}{k}$$

## 6.4 Trust Score

Normalize $(F)$, $(Stab)$, and $(Cons)$ to [0,1] and compute:

$$T = w_1 F + w_2 Stab + w_3 Cons, w_1 + w_2 + w_3 = 1 T = w1$$

## 7. Implementation

### 7.1 Tools and Technologies

- Python
- Scikit-learn
- SHAP library
- Pandas
- NumPy
- TensorFlow

### 7.2 Experimental Setup

Experiments were performed in a commonised computation generator, using a Python implementation. The data was divided into two portions of training and test in a 80:20 ratio.

- Python implementation
- Scikit-learn and SHAP library
- 10-fold cross validation
- Hyperparameter tuning using grid search

### 7.3 Evaluation Metrics

1. **Accuracy:** the measure of the percentage of the evaluation set which is correctly classified.
2. The percentage of correct positive predicates of all positive predicates.
3. **Recall:** Fraction of correct predictions of those that are actually positive.
4. **F1-score:** the harmonic mean of recall and precision which combines precision measurements with recall measurements to create one measurement, the F1-score.
5. **AUC Area under a receiver operating characteristic curve:** This summarizes the trade-off between false positive and true positive rates induced by varying thresholds.
6. **Trust Score:** a confidence measure of the ratio of directions of the two closest and second best clusters of classes in the feature space.

## 8. Experimental Results

### 8.1 Performance Metrics

1. Accuracy
2. Precision
3. Recall
4. F1-score

### 8.2 Model Comparison

**Table 2: Performance Comparison of Machine Learning Models**

| Model | Accuracy | F1 | AUC | Trust Score |
|---|---|---|---|---|
| SVM | 87% | 0.85 | 0.89 | 0.52 |
| Random Forest | 91% | 0.90 | 0.92 | 0.71 |
| Proposed Model | 95% | 0.94 | 0.97 | 0.91 |

### 8.3 Analysis

The offered framework has shown much progress in predictive accuracy and model interpretability that are considered to be conflicting in the model learning systems. The framework not only attains credible results but also provides substantive explanations of its decisions in the effort to improve performance alongside transparency, thus, achievable simultaneously. This interpretability is necessary in building stakeholder confidence and also in ensuring the practical introduction of smart systems.

The only significant contribution by the framework is that it incorporates a trust-evaluation mechanism to determine the reliability of model predictions. In contrast to the classical machine learning systems where predictions remain a black box, the proposed system involves

a self-assessment system that assesses the confidence in the decision. The latter characteristic can be especially beneficial in risky areas of application, including healthcare, financial systems, and autonomous technologies, where mistaken forecasts can have devastating outcomes. The ability to measure trust increases the reliability of a system and helps to be responsible when implementing artificial intelligence solutions.

In order to elaborate the suggested framework, the following research directions can be considered:

## A. Model Calibration

It is possible that the future work can work on enhancing confidence estimation using a calibration method like the Platt scaling or isotonic regression. The accurate calibration of the predictions of probability leads to an improved reliability of the likelihood which reinforces the trust evaluation module and improves the accuracy in making a decision.

## B. Greater Improvements of Explain ability

Global and local interpretability could be achieved by the implementation of advanced explainability methods such as SHAP-based feature attribution and counterfactual explanations. This would extend the levels of transparency and help understand the behaviour of the models better.

## C. Real-Time Trust Assessment

It is also possible to establish actual time trust scoring by inserting lightweight monitoring pipelines that would ascertain ongoing checks of reliability into the models especially in dynamic environment, where data distribution can change over time.

## D. Human-in-the-Loop Evaluation

The experimental research that can be performed with human contact is aimed at the assessment of the system usability, latency of the decision, as well as the degree of trust in the system when it displays low-confidence predictions. Such tests would offer viable confirmation of the mechanism of trust-assessment.

Altogether, integrating trust assessment in the decision-making process is a considerable move towards responsible and explainable artificial intelligence. The framework focuses on predictive performance in addition to transparency, reliability and ethical deployment hence helping to evolve dependable AI systems.

## 9. Discussion

Explainable machine learning improves on both transparency and accountability in systems of smart city decision making. The framework provided below will help the policymakers understand what predictive will identify, and thus, encourage the responsible use of artificial intelligence.

## 10. Conclusion

The current work describes the explanation and trust-related hybrid machine-learning architecture that can be applied to smart-city decision-making systems. Using a unique trust-evaluation module, combined with SHAP and LIME explainability methods helps the framework to add predictive accuracy, transparency, and general model understandability. The resulting methodology provides reliable decision-making due to the availability of interpretable predictions in addition to the quantitative evaluation of data on confidence of the model results. To this end therefore, this structure forms the basis of credible, AI-based urban governance and promotes the responsible use of intelligent systems in smart-city settings to this approach.

## 11. Future Work

Future work will focus on:

- real-time deployment
- digital twin integration
- federated learning for urban data
- deep explainable neural networks

**Real duplying:** In this stage, we will be paying attention to building lightweight inference pipelines that are to be allowed to generate predictions on the edge devices to accept and comply with rigid latency expenditures. Second, to optimize our models in actual real-life situations, we will incorporate them in the operational architecture.

**Digital twin:** models that would support two way data flow between virtual models and physical structures, thus assisting in planning scenarios and preemptive maintenance.

**Federated learning:** this component aims to support cooperation between various urban sources of data silo, which will be coordinated in a privacy-aware manner, but exploit the diversity of city-scale sensor networks. At last, we will explore higher-order aggregation and personalization to achieve the highest model efficacy without having to break the data confidentiality.

**Deep explainable neural networks:** we will adopt attention mechanisms, saliency maps, and counterfactual reasoning to make the decisions of our models explainable and begin to act on them by both domain experts and policy-makers. With these combined efforts, we will be able to generalize our efforts beyond controlled research environments to deployments of high stakes in the real-world.

**References**

1. Chourabi, H. et al. (2012). Knowing smart cities: A comprehensive model. *HICSS*.

2. Zanella, A. et al. (2014). Smart city Internet of things. Internet of Things Journal, 1 (1) IEEE.

3. Thakker, D. et al. (2020). Intelligible AI of smart cities. IEEE Smart Cities Conference, 2016.

4. Bibri, J., & Krogstie, A. (2017). Intelligent urban communities of tomorrow. Sustainable Cities and Society.

5. Bussmann, N. et al. (2021). explainable artificially intelligent finance explanations. The Frontiers in Artificial Intelligence.

6. Samek, M. et al. (2019). Elucidate explainable AI models: Interpret machine learning models. Springer.

7. Holzinger, A. et al. (2020). So what are we missing to develop explainable AI systems? *arXiv*.

8. Hurbean, L. et al. (2021). Intelligent cities and machine learning. *Electronics*.

9. Lee, J. et al. (2020). Smart city infrastructure systems based on AI. *IEEE Access*.

10. Berezsky, O. et al. (2025). Evaluation of smart cities with the help of ML. *Sensors*.

11. Ullah, Z. et al. (2020). UrbanIoT UrbanIoT: Multi-modal IoT data. *IEEE Access*.

12. Mahesa, R. et al. (2019). Sustainable smart city development data. *Data in Brief*.

13. Iskandaryan, D. et al. (2020). Machine learning Air quality prediction. *Applied Sciences*.

14. Zhou, P. et al. (2019). Intelligent city deep learning. *IEEE Network*.

15. Liu, Y. et al. (2016). Smart city analytics of big data. *Computer Systems of the Future: Computer Systems Inc.

16. Rathore, J. et al. (2018). Big data analytics Urban planning. *IEEE Access*.

17. Kitchenham, B. (2009). Software engineering reviews of literature. Informational and Software Technology.

18. Recht, B. et al. (2019). Machine learning: tendencies and issues. *Science*, 365(6456), 485–489.

19. Bello-Orgaz, G. et al. (2016). Smart city analytics on social big data. *Computer Systems of the Next Generation.

20. Janssen, M. et al. (2012). Advantages and obstacles of open data. |human|>Informed Systems Management.

21. European Commission. (2019). *Open data directive*.

22. Government of India. (2023). The Urban data exchange platform in India.

23. Smart Cities Mission India. (2023). *Open government data platform.

24. Van der Aalst, J. (2016). Procedural mining: information science in play. Springer.

25. Esteva, A. et al. (2019). A roadmap on deep learning in healthcare. *Nature Medicine*, 25(1), 24–29.