

Suvarga: Promoting A Healthy Society

Priya R. L¹, Gayatri Patil², Gaurav Tirodkar³, Yash Mate⁴, Nikhil Nagdev⁵

^{#1}Computer Department, Vivekanand Education Society's Education of Society
Chembur, Mumbai-74

¹priya.rl@ves.ac.in, ²2017.gayatri.patil@ves.ac.in, ³2017.gaurav.tirodkar@ves.ac.in, ⁴2017.yash.mate@ves.ac.in,

⁵2017.nikhil.nagdev@ves.ac.in

Abstract In a country like India, poverty is one of the major issues. Over 22% of the population is below the poverty line. This poverty pushes people on streets which in the future transforms into slums. These slums, as are not planned, lack certain necessities like electricity, sanitary services, and basic hygiene resources leading to a hub for the spread of diseases. Another important cause is the lack of awareness about cleanliness among slum residents. In essence, the primary aim of this paper is to identify the leading causes of diseases in slum areas of Mumbai using data collected from IoT modules, health checkup drives, and various government authorities. With this information, the concerned civic authorities and slum residents will be alerted regarding the danger so that necessary action can be taken. This, in turn, promotes the healthier society in various slum regions of India.

Keywords Internet of Things (IoT), Slum Management, Sanitation, Decision Tree, LSTM, Air Quality Index, Water Quality Index.

1 Introduction

According to the United Nations (UN, 2009) estimates, only 4% of the terrestrial surface is occupied by cities [4]. Though the percentage is so low more than half the world population stays in these cities and which eventually generates a huge imbalance in the world resources as this section consumes three-quarters of the world's natural resources.

For upgrading these slums, commonly, the first action taken is to demolish the slums and reallocate the residents but, since 1970 there have been multiple recommendations by authors such as Turner (1972) which suggest otherwise. This gives birth to the concept of upgrading slums and their residents to a better standard of habitation [3]. The paper uses the approach of data analysis and deep learning to have a better understanding of this approach and provide solutions for implementing the same.

2 Literature Survey

Slum management has always been a major issue in the city like Mumbai [5]. Current research by SRA maps the information of residents on a website with the help of drones. The government has gathered data regarding the current groundwater, the nearby reservoirs, precipitation, water usage domestically, industrially, and agriculture [6]. Although this method covers many data fields, major data fields like the health factors, pollution measures, etc. are ignored.

Improved infrastructure can prove to be a major catalyst for achieving major sustainable goals. Taking this aspect into consideration, the UN-Habitat Opinion Survey method which was based on the nature of social reality and the perspective of the researchers, was used in the slum residents of Africa. The data was collected through an expert opinion survey and the result after analysis displayed that the infrastructure in Africa can be primarily developed with the help of proper water supply, road networks, and telecommunication. [6]

To understand the positive and negative implications of upgrading slums, a case study was conducted in Moravia's Neighborhood, Medellin. The principles of urban design strategies and urban rehabilitation programs were identified through technical documents, qualitative and quantitative data which was collected through surveys at the community level. [7] This research didn't include the up gradation strategies, which are important as it would be difficult to provide complete rehabilitation.

A non-integrated framework was adopted to evaluate the suitability of the interior design of a low-income multipurpose apartment to provide enhanced IEQ. The research plan has 5 sequential steps: Data Collection, random sampling, Simulations, and calculations for scenarios, multi-objective optimization, and participation of stakeholders. Here expert opinion survey was taken, AHP TOPSIS was performed and the final optimized solution was generated. [8] This research included only the interior design and not the other parameters like pollution and health and other geographical aspects.

3 Proposed System

3.1 Overview

To aid the health situation in the country, the model proposed a novel approach by building an Internet of Things (IoT) based Intelligence system. The model will provide regular updates about the chances of epidemics, few symptoms to spot those diseases, and also emergency contacts of concerned doctors along with a few home remedies. It also alerts the government authorities regarding it aiming to create the necessary awareness among the government authorities and the slum residents.

3.2 System Architecture

The system is composed of various modules such as Data Collection from various sources (IoT, BMC Health Data, Water data, sanitation, survey data), Preprocessing, Feature Extraction, Training, and Testing Models, displaying the final output on the web and mobile application as listed below and shown in Fig. 1.

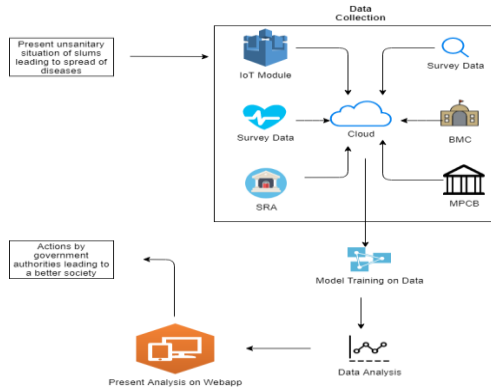


Fig. 1 System Architecture of Suvarga.

Data Collection

To ensure that data used in the model is authentic, health, sanitation, and water data were collected from government authorities like MPCB, SRA, BMC, etc. Health data was provided by the Medical Officer of Health (MoH) of the Chembur ward.

1. Drive for Demographic Data

A health checkup drive was carried out in the slums as shown in fig.2 to get accurate health parameters like Hemoglobin, Weight, Height, BMI Index, etc. of the residents.



Fig. 2: Slum Drive

2. Air Quality Monitoring Device

Air Quality monitoring device is a system deployed using IoT devices to test the air quality parameters with the high tech IoT sensors. It gives a graphical representation of how each gas emission affects the overall air quality in the surrounding.

3. Preprocessing

Data collected from the IoT modules as well as from government authorities was preprocessed to fit into the necessary format. All the missing values were filled using average, mean, or null values according to the data collected. After filling the missing values, the data were combined in a single CSV file by sorting each file month-wise and combining the required parameters from each file.

4. Prediction Model

To achieve a high accuracy the model was trained using LSTM and Decision Tree algorithms. Dataset procured from government authorities was 2407 tuples from 2010-20. This data was split into training and testing datasets in the ratio of 80-20% and trained using Machine Learning Algorithms

5. Test Model

The prediction model was tested using 20% of the unknown dataset. The model was evaluated using various performance measures such as R2 score, Mean Squared Error (MSE) and Mean Absolute Error (MAE)

6. Website or Mobile Application

Awareness and progress can happen only if power is given in the hands of the masses. With the website and app, not only the government authorities but also the commoners can monitor the cleanliness parameters. With the app or website, people will be alerted about the chances of epidemics, its respective cures, and emergency contacts of doctors.

4 Detailed Architecture of Suvarga

The detailed workflow of Suvarga as shown in Fig. 3 describes the data collection from heterogeneous sources with data preprocessing to build a better prediction model.

4.1 Air Quality Monitoring System

To calculate the air quality index of specific regions we have built an IoT based module.

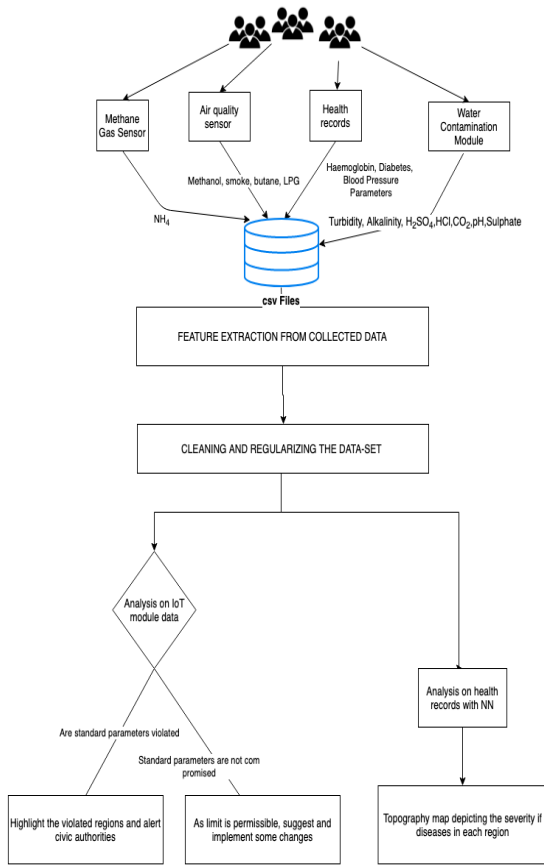


Fig. 3 Detailed architecture of Suvarga

Each module is composed of MQ series sensors (MQ135, MQ2, MQ3) to measure the air quality index. These sensors are placed on an ESP8266.

Name Of The Component	Features and Description
MQ 135	Gases, including NH ₃ , NO _x , alcohol, benzene, smoke, and CO ₂ are detected by this Air Quality Sensor
MQ2	Combustible gas and smoke at concentrations from 300 to 10,000ppm are detected by the semiconductor gas sensor
MQ3	This sensor is used to detect leakage of flammable gases (LPG), methane.

Table 1 Component in Air Quality Monitoring Device.

4.2 Water Quality Monitoring System

Poor water quality is a big issue, especially in slum regions. To test the water quality, the model uses the BMC water quality monitoring module. The

data gives out the pH, Dissolved Oxygen, B.O.D., C.O.D., etc. Over ten years of data are collected.

4.3 Sanitation

Sanitation data from BMC gives the distribution of toilets for men and women in the Chembur region. It provides information on the number of toilets with respect to the number of people.

4.4 Algorithms Used

The model was trained using algorithms like LSTM and Decision Trees. Later it was compared to choose the best algorithm.

LSTM

LSTM known as the Long Short Term Memory algorithm is much more complicated than others as along with a simple tanh it has a lot of new layers as shown in fig.4 below.

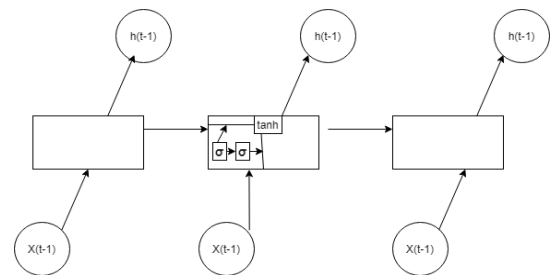


Fig. 4 LSTM Algorithm

Decision Tree

Decision Tree algorithm is a predictive modeling tool and uses a tree-like model that is generated via an algorithmic method. It searches for ways to split a data set based on different conditions.

It is one of the supervised learning approaches which are used for classification as well as regression problems. The rules for decisions are generally in the form of if-then-else statements.

4.5 Intelligent System

The proposed system aims to build an intelligent system for promoting healthy living in various slum regions of India. It consists of mainly two components such as the Prediction or Analysis model and Data Visualization Model. Such components are designed to display the analysis obtained from the analysis model and visualized in graphical formats via the web application.

Prediction Model

Using the data collected from the IoT Module for reading air quality parameters and others collected from various other sources via government authorities such as BMC, MPCB, etc. an analysis algorithm is run to predict the values of Air/ Water, Correlation among various features of the dataset.

LSTM algorithm is applied to data and the correlation among features is found using Pearson's Correlation formula available in the pandas' library.

Data Visualization

The final outcome of all the analysis needs to be presented to the layman in terms understandable, hence a user-friendly Web App is built for the government authorities as well as the slum residents. Each can access multiple features of the web app like predictive analysis of air and water quality in the future, basic care, and home available remedies to prevent oneself and loved ones from epidemics.

5 Implementation

5.1 Need for Real-Time Monitoring

From an overall high-level perspective, the data aggregation throughout the day provided by the government authorities is very useful in data analysis of the resources and the environment over the long term, but it does not provide a mechanism to handle such mishaps. Hence, to address this problem, Suvarga has developed a network of IoT devices that would be installed in the slums. These IoT devices would act as a network, continuously monitoring the various air quality parameters.

5.2 Experimental Setup of Air Quality Monitoring Device

The IoT device consists of a microcontroller called 'NodeMCU'. It is capable of being interfaced with gas sensors, and also transmits data over a Wi-Fi network. The sensors that are interfaced with this microcontroller are MQ135, MQ3 and MQ2 sensors.

Three different IoT devices are fitted at the three corners of the slum area. All of them are connected to a common network provided by the hotspot of the mobile device. The devices act in unison, forming a mesh, transferring data to a common device acting as the source which forms a server. The server would keep a track of the gas levels, continuously monitoring if the levels have crossed a threshold permissible value.

The data received over the IoT module is sent to this centralized data visualization Cayenne server.

The data could be visualized in real-time and plotted on a live- graph which constantly oscillates between a range of values. The set trigger is activated when a sensor gives a value that crosses a threshold, indicating that an accident has taken place. An instant notification in the form of SMS and email alert to the concerned government authority is sent simultaneously. With this, the government can send instant relief or could take necessary actions to pacify the toxic environment. As shown in fig.5, it describes

the experimental setup of the air Quality real-time monitoring device.



Fig. 5 Experimental Setup of the Air Quality real-time Monitoring device

5.3 Slum Health Drive Data Analysis

A Slum health drive was conducted for the residents of the slum adjoining VESIT on the 25th of January 2020. When enquired about the issues faced by the residents on a daily basis, it was brought into the notice of the team that there was an outbreak of malaria and snake bites in the region.

The following are the parameters as depicted in fig. 6 of the data that the team obtained from the health drive.

Name	object
Gender	object
Age	int64
Weight	int64
Height	object
Poverty Status	object
Toilet	object
Drainage linked to the house	object
Waste Collection System	object
Compost pit	object
Source of Water	object
Washing Clothes and Utensils	object
Alcohol/ smoking	object
Cooking	object
Diabetes	object
Hypertension (Blood Pressure)	object
Cholestrol	object
Level of education	object
Level of education.1	object
Adhaar Card	object
Authorized Electricity Connection to Household	object
Bank Account	object
Computer Literate	object
Source of Income	object
Signature	object
dtype:	object

Fig. 6: Parameters obtained during Slum Drive

Slum Drive analysis

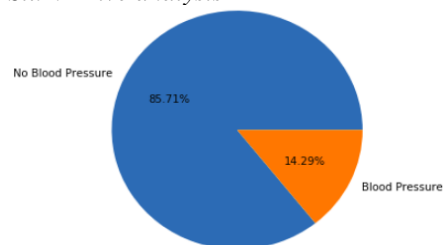


Fig. 7: Blood Pressure Ratio

Fig. 7 shows the ratio of blood pressure people.

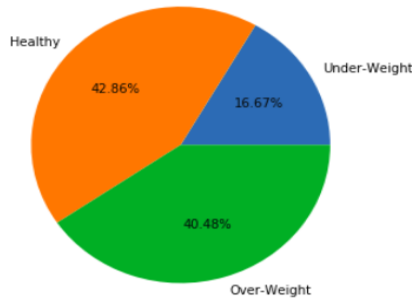


Fig. 8 Weight Ratio.

Fig. 8 represents the percentage of people that were healthy by weight, or overweight, or underweight.

5.4 Sanitation Data Analysis

Open defecation has been an onerous issue for a while, causing a variety of health-related issues. The team of researchers decided to collect sanitation data from the government authorities of the Chembur region through BMC offices, which keep a record of all parameters of that region. The dataset obtained was in the form of a CSV file encompassing parameters displayed in fig.12. The dataset comprises 172 rows (records) and 7 columns (attributes).

While it is essential to have a sufficient number of toilets overall in the entire region it is crucial that every ward has enough sanitation facilities in order to avoid any strain on resources. Hence, ward by ward analysis is done. The group by method in pandas (a library in Python that deals with Data Frames) is used to group the data by each ward. Fig. 9 shows the grouping by ward number 150.

Sr.0	Ward 0.	Toilet Address	Run By	Male	Female	Total
3	7	150 MATOSHRI CHARITABLE TRUST, TEMBREE BRIDGE	SANSTHA	6	2	8
4	8	150 MALEKAR WADI, BHIM SHAKTI MAHILA MANDAL	SANSTHA	7	6	13
5	9	150 EKTA NAGAR, P.L. LOKHANDE MARGE	MCGM	0	14	14
6	10	150 JYOTI NAGAR, P.L. LOKHANDE MARGE	MCGM	12	14	26
7	12	150 SAMBHAJI NAGAR, P.L. LOKHANDE MARGE	MCGM	3	3	6
8	13	150 JAGTAP CHAWL, P.L. LOKHANDE MARGE	MCGM	12	12	24
9	14	150 JADHAV CHAWL, P.L. LOKHANDE MARGE	MCGM	12	12	24
10	15	150 PATHAN CHAWL, P.L. LOKHANDE MARGE	MCGM	3	3	6
11	16	150 MUNJAL NAGAR, BHIMSHAKTI MAHILA MANDAL	SANSTHA	10	10	20
12	17	150 SATAM CHAWL	MCGM	6	6	12
13	18	150 AMAR MITRA MANDAL	MCGM	6	6	12
14	19	150 BHIM SHAKTI JEST NAGRIK	Unknown	0	0	0
15	20	150 SEWA SENSTHA	SANSTHA	6	6	12
16	21	150 AMIR BAUG(PUNJAB CHAWL)	MCGM	12	12	24
17	22	150 AMIR BAUG (KACHARA KUNDI)	MCGM	12	12	24
18	23	150 KRUSHAN MEHAN (P.Y. THORAT)	SANSTHA	16	8	24

Fig. 9: Group No. 150

To find out, if all wards have toilets commensurate with each other, a pie chart has been plotted that shows the distribution of the toilets in the region. A discrepancy has been observed in the distribution. While the number of toilets in ward number 154 soar as high as 32, the number of toilets in ward number 149 is a meager 3.

The total toilets in the region are 2669. According to research, the number of people per toilet is 100.

Taking into consideration the population of the region by the population data obtained, the number of persons per toilet is close to 457(Fig. 10.).

```
print(total_population)
1219613

people_per_toilet=total_population/total_toilets

#Printing the number of people using a particular toilet
print(people_per_toilet)
456.955039340577
```

Fig. 10: People per Toilet

5.5 Health Data Analysis

As the team of researchers meandered in the slums and interviewed quite a few residents it was brought to the team's notice that the slum was an abode to diseases like Malaria & Snakebites and residents were suffering from respiratory infections. With the air and water quality data at the fingertips, the team decided to obtain the health data from the hospitals in the vicinity of the area where such cases happened sporadically. The dataset obtained from BMC offices was month-wise historical data of years 2017 and 2018 and had the name of the dispensary, the month under consideration, and the average levels of health-related parameters like the total number of people suffering from asthma, malaria, URTIs, heart diseases to name a few. The data in the form of a CSV file was analyzed with pandas library. The visualizations offered through seaborn - a data visualization library, make the inference of the results less tedious. Refer Fig. 11

Name of Dispensary	Date -Year	Month	OPD_Total	OPD_New	OPD_Old	Fever	Gastro	URTI	Hepatitis	Asthma	Total	HeartNew	
0	Anik Nagar Disp.	03-01-2017	Jan-17	3325.0	2324.0	1001.0	322.0	47.0	120.0	0.0	...	15.0	0.0
1	Ayodhya Nagar Disp.	04-01-2017	Jan-17	3165.0	2148.0	1017.0	328.0	1.0	399.0	0.0	...	18.0	0.0
2	Gawarpada Disp.	05-01-2017	Jan-17	1280.0	750.0	530.0	123.0	21.0	131.0	0.0	...	20.0	0.0
3	Deonar Disp.	06-01-2017	Jan-17	1977.0	1377.0	600.0	216.0	3.0	395.0	0.0	...	62.0	0.0
4	Maharashtra Nagar Disp. / Cheeta	07-01-2017	Jan-17	2526.0	1910.0	616.0	264.0	62.0	421.0	0.0	...	31.0	0.0

Fig. 11: Health Analysis

The air and water quality data obtained encompassed historic data for the past four years. Whereas the health data procured from the BMC authorities attributed records from 2017 to 2018. To avoid any aberrations, air, and water quality data of 2018 and 2017 are taken into account along with health data of the same two years.

The table in Fig. 12 displays the combined air and water quality data.

```
In [301]: merged_data=pd.merge(df_water, df_health, on='New_Dates')
```

```
In [302]: merged_data.head()
```

```
Out[302]:
```

	Month- Year	pH	Dissolved Oxygen	B.O.D.	C.O.D.	Nitrate	Fecal Coliform	WQI	SO2	NOx	...	AsthmaTotal	HeartNew	HeartOld
0	January 2017	7.4	3.1	18.0	232.0	1.04	920.0	42.56	17.63	35.17	...	15.0	0.0	0.0
1	February 2017	7.7	3.8	19.0	252.0	2.10	540.0	45.98	18.10	53.74	...	12.0	0.0	0.0
2	March 2017	7.2	3.7	17.0	224.0	1.70	540.0	48.41	18.53	52.47	...	15.0	0.0	0.0
3	April 2017	7.4	4.2	16.0	180.0	3.20	540.0	50.89	13.10	39.77	...	15.0	0.0	0.0
4	May 2017	7.4	3.8	16.0	220.0	1.50	1600.0	46.69	11.00	35.97	...	13.0	0.0	0.0

Fig. 12: Preprocessed Air & Water Quality Data

Correlation between all the parameters is obtained, the results of which are stored in a correlation matrix.

All the correlations are stored in the matrix and then sorted in ascending order using quick sort. The variables having the most correlation are shown in Fig 13.

B.O.D.	WQI	0.747056
HypertensionNew	NonCommunicableNew	0.751363
NonCommunicableNew	HypertensionNew	0.751363
C.O.D.	Dissolved Oxygen	0.759451
Dissolved Oxygen	C.O.D.	0.759451
PsychiatricNew	PsychiatricTotal	0.768541
PsychiatricTotal	PsychiatricNew	0.768541
NonCommunicableNew	DiabetesNew	0.795260
DiabetesNew	NonCommunicableNew	0.795260
Fever	OPD_Total	0.796926
OPD_Total	Fever	0.796926
OPD_Old	OPD_New	0.807801
OPD_New	OPD_Old	0.807801
Fever	OPD_New	0.821314
OPD_New	Fever	0.821314
WQI	C.O.D.	0.821821
C.O.D.	WQI	0.821821
Dissolved Oxygen	WQI	0.840527
WQI	Dissolved Oxygen	0.840527
Dissolved Oxygen	B.O.D.	0.871425
B.O.D.	Dissolved Oxygen	0.871425
OPD_Total	OPD_Old	0.909766
OPD_Old	OPD_Total	0.909766
AsthmaTotal	AsthmaOld	0.915464
AsthmaOld	AsthmaTotal	0.915464

Fig. 13: Correlation matrix

However, the correlations above show the relation between the attributes of the same table. A correlation between air quality parameters and URTIs is established and in the same way, Malaria is correlated with the Biological Oxygen Demand in the water.

```
Correlation between URTI and RSPM is 0.5478534282519916
Correlation between Malaria and B.O.D. is -0.5941003965138847
```

Fig. 14: Correlation between Different Tables

According to research, there has been an association between Upper Urinary Tract Infection and Respirable Suspended Particulate Matter

(RSPMs) [9]. An increase in particulates is detrimental to health as it causes a variety of conditions related to respiratory tracts. The analysis conducted (Fig. 15) also states the direct positive association between the two factors.

```
In [324]: sns.regplot(x='URTI', y='RSPM', data=merged_data)
```

```
Out[324]: <matplotlib.axes._subplots.AxesSubplot at 0x7fcd1cd09cf8>
```

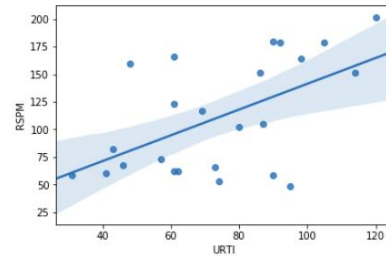


Fig. 15: URTI vs. RSPM

6 Results and Analysis

The comparative study established between the regression algorithms suggests that Decision Tree Regression achieves the lowest error-rate when evaluated with error measurement metrics for regression comprising of Mean Squared Error (MSE), Mean Absolute Error (MAE) and R2 Score as compared to the other regression algorithms being tested on these metrics like Lasso Regression, Lasso Lars Regression, Bayesian Regression, and Random Forest Regression.

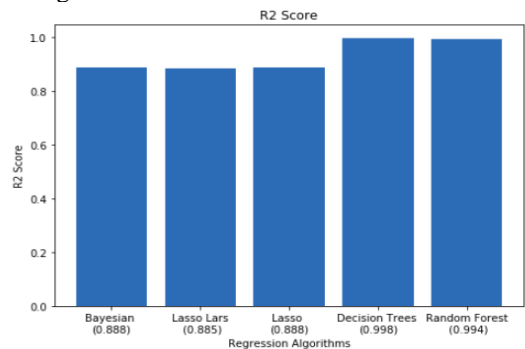


Fig. 16: R2 Score

Mean Squared Error

Mean Squared Error is a metric that tells how close the predicted points are to the actual points on the regression line.

$$L = \frac{1}{N} \left[\sum (\hat{Y} - Y)^2 \right] \quad (1)$$

Where: L - Loss, \hat{Y} - Output Y - Actual Value, N - Samples

The results of MSE indicate that Lasso Lars Regression gives the maximum error of 381.74, while Decision Tree Regression gives the least MSE of 5.65. The MSEs of the other algorithms fall between

these two. Bayesian Regression gives an MSE of 381.74, Lasso Regression produces an MSE of 381.77. Random forests perform better than most of the algorithms except Decision Trees giving an MSE of 17.19. The results are shown in Fig. 17.

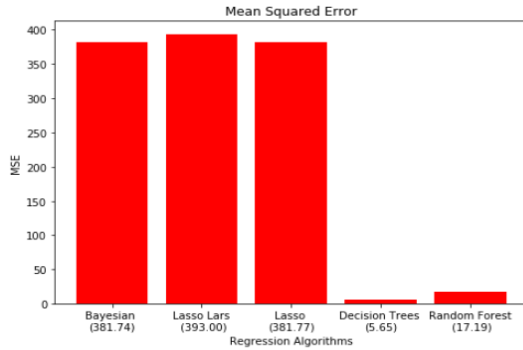


Fig. 17: Mean Squared Error

Mean Absolute Error

Mean Absolute Error is the measure of the difference in the actual value and the predicted value.

$$MAE = \frac{1}{n} \sum_{i=1}^n |x_i - \hat{x}| \quad (2)$$

n = the number of errors, Σ = symbol for summation, $|x_i - \hat{x}|$ = the absolute errors.

From the chart in Fig 19 below, we can infer that the MAE for Decision Trees is the least with 0.64, while Bayesian and Lasso Regression give the most R2 score among the five indicating poor performance. The values of MAE of other algorithms lie between these two, represented by the bar chart

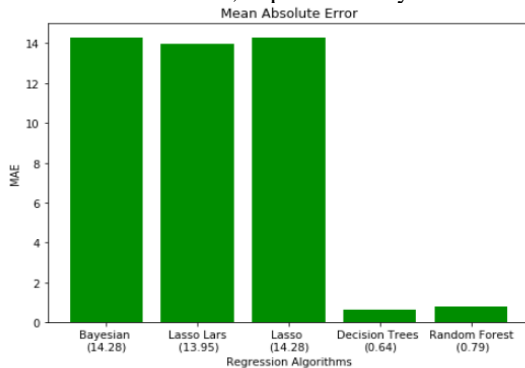


Fig 18: Mean Absolute Error

7 Conclusion and Inferences:

The research undertaken particularly focuses on improving the health and sanitation facilities on the slums in Chembur, Mumbai region. Harnessing the potential of Artificial Intelligence, Data Analysis, and Internet of Things (IoT), the proposed system predicts the patterns in Air Quality, Water Quality, sanitary facilities and builds a strong interdependence of these environmental and sanitary factors on the

health of the individuals residing there. The governmental authorities can formulate laws and policies and take actions if these predicted values cross a certain threshold and mitigate the ill-effects of environmental changes on the health of the residents.

Through sanitation data it was found out that the number of people per toilet was 456.955 which were a lot higher than the ideal ratio which is 100 people per toilet. A correlation is also established between the numbers of Malaria patients in the hospitals in the vicinity of the slums to the Water Quality Index. It was concluded that URTI and RSPM have correlation of 0.547853, thus having a significant correlation. Data from the NGO health drive conducted indicated that a 57.15% of residents were either overweight or underweight. Moreover, the R2 score obtained from decision trees has a value of 0.99, indicating almost perfect prediction. The government could use the findings of the research to take appropriate actions to assuage the detrimental effects of poor well-being, unclean surroundings, and polluted environment on the dwellers of slums.

References

- [1] Fundamentals of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) Network
- [2] Balaban O, Puppim de Oliveira JÁ, Sustainable Buildings for Healthier Cities: Assessing the Co-benefits of Green Buildings in Japan, Journal of Cleaner Production (2016), DOI:10.1016/j.jclepro.2016.01.086.
- [3] Building Resilience of urban slums in Dhaka, Bangladesh, Iftexhar Ahmed/ Procedia-Social and Behavioral Sciences 218(2016)
- [4] United Nations (2009). World Population Prospects: 2009 revision, Population and Development Division, Department of Economics, and Social affairs.
- [5] Sanjay Dikhle, Rakesh Lakhena, GIS-Based Slum Information Management System, 17th Esri India User Conference 2017
- [6] Ben Arimah, Infrastructure as a catalyst for the prosperity of African cities, Urban Transitions Conference, Shanghai, September 2016.
- [7] Katila Vilar, Ivan Cartes, Urban Design, and Social Capital in Slums. Case Study_ Moravia's Neighborhood, Medellin, 2004-2014
- [8] Ahana Sarkar, Ronita Bardhan, Improved indoor environment through an optimised ventilator and furniture positioning_ A case of slum rehabilitation housing, Mumbai, India, accepted 1 December 2019.
- [9] Y.R.Li, C.C.Xiao, J.Li, J.Tang, X.Y.Geng, L.J.Cui, J.X.Zhai, Association between air pollution and upper respiratory tract infection in hospital outpatients aged 0-14 years in Hefei, China: a time series study, Public Health Volume 156, March 2018.