

## A Novel Integrated Machine & Business Intelligence Framework for Sensor Data Analysis

Kalyani. S <sup>1</sup>, Dr. Venkat Rao <sup>2</sup>

1 Vignan's Institute of Engineering for women, Vishakhapatnam,

2 Andhra University Visakhapatnam India

contact author: kalyaniakiri3009@gmail.com

**Abstract:** Increased smart devices in various industries is creating numerous sensors in each of the equipment prompting the need for methods and models for sensor data. Current research proposes a systematic approach to analyze the data generated from sensors attached to industrial equipment. The methodology involves data cleaning, preprocessing, basics statistics, outlier, and anomaly detection. Present study presents the prediction of RUL by using various Machine Learning models like Regression, Polynomial Regression, Random Forest, Decision Tree, XG Boost. Hyper Parameter Optimization is performed to find the optimal parameters for each variable. In each of the model for RUL prediction RMSE, MAE are compared. Outcome of the RUL prediction should be useful for decision maker to drive the business decision; hence Binary classification is performed, and business case analysis is performed. Business case analysis includes the cost of maintenance and cost of non-maintaining a particular asset. Current research is aimed at integrating the machine intelligence and business intelligence so that the industrial operations optimized both in resource and profit.

**Keywords:** RUL Prediction, Sensors, IOT, Aircraft Engine, Business Intelligence

### 1. Introduction

Advances in Internet of things has made possible to connect numerous machines and devices to internet, so that information about machine operation and health is available for analysis. With increased network access it has become easy to capture the data using sensor, there is a need for models and methods for analyzing the data captured. Products in industries like Aerospace, Defense, Automotive when embedded suitable sensors that are miniature and self-powered can make the products both are being developed These sensors are self-powered and self-intelligent to capture information, also process the same as much as possible so that the a reaction mechanism can be activated using the local processing. Fig. 1 depicts the growth of the number of devices that are being connected to internet from 2012 to 2020, and it is predicted that more machines will be connected to internet in near future. Hence there is need for Algorithms needs to be developed which can be deployed both for cloud and edge processing of the data. These algorithms differ from the usual data analysis models as sensor data reflects the challenges of real operating environment of the system. The sensor is self-powered and exposed to harsh working environments hence the data being captured can be faulty and erroneous. One of the challenges of the current IoT system is to generating business insights from IoT platforms.

Nomenclature	
Cycle	Operating cycle
Id	Machine Index
S	Sensor Index

tff	Time To failure
Setting	Operational Setting
Train	Training Data Set
Test	Testing Data Set
Bcl	Binary Classification Label
Mcl	Multi Class Classification Label
TP	True Positives
FN	False Negatives

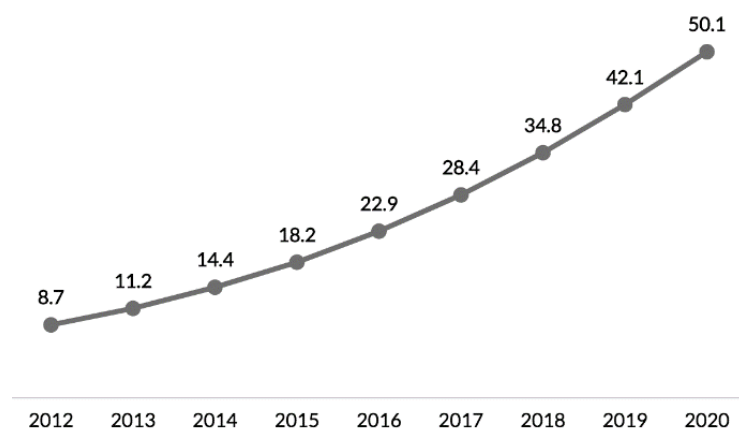


Figure 1. Increasing Trend of Smart Connected Devices

Referring to Fig. 1 each of the device that is connected to internet adds more data Huge volume of data generated from these sensors needs to be integrated with artificial intelligence system for timely action by optimizing the business resources. There is need for systems that are intelligent and self-improving. With the advent of edge and cloud computing it is becoming important to increase the speed and accuracy of predictions.

## 2. Literature Review

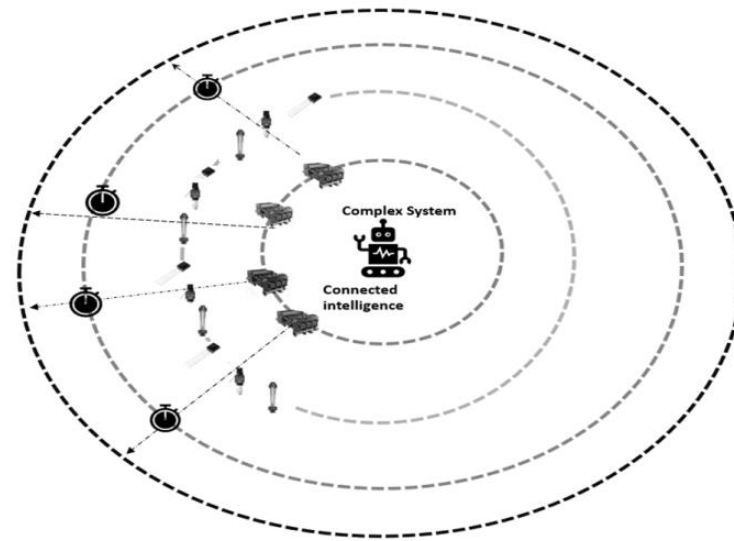
Data Driven Modelling has gained momentum with numerous publications in the field. NASA has been enduring people to use the data available through CAMPSS [1] et al proposed a propagation model. A goal et al used surrogate modelling technique to find the samples of varying size to help model the machine life. Pragma et al [3] used the neural network approach for predicting the condition of industrial equipment like refrigerator. Pranav Kumar et al [6] used genetic algorithm-

based approach for optimizing the industrial design parameters using the genetic algorithm, the research is aimed at using machine learning algorithm for optimization of various parameters. The research so far focused on initial approaches for the modelling the data related to design and operation of industrial components. There is need for better data driven models to be useful for industrial IOT applications. There are studies focusing on anomaly detection using time series models there is need to work on real time anomaly detection for edge computing and safety critical equipment. Alexandria [3] et al used machine learning models for sensor analysis in using a sensor node approach. Temperature ( $^{\circ}\text{C}$ ), Humidity (%), Light (Lux), Pressure (hPa) are sensed using the node and the same are processed by internet gateway. The experiment is conducted in a lab environment with human presence to include additional data of number of people. The approach is limited in terms of variables considered. Taha et al [4] explored machine learning based methods and models for sensor data that is generated from an aircraft engine. The engine consists of multiple subsystems with parameters like pressure, temperature, flow rate. The proposed methodology uses only a set of total number of sensors for predicting the remaining useful life. RUL prediction can be improved if the number of parameters considered can be maximized. The model is more developed by excluding the extreme values of each of the sensors. Okoh et al [5] presented a review of existing remaining useful life estimation computational techniques. Ellefsen et al [6] proposed a deep learning based

architecture, his paper investigates the effect of unsupervised pre-training in RUL predictions utilizing a semi-supervised setup. Additionally, a Genetic Algorithm (GA) approach is applied to tune the diverse number of hyper-parameters in the training procedure. Zhou et al [7] proposed a residual lifetime estimation on function data. Hanachi et al [8], modelled RUL using performance-based sensor data of engineering systems. S.K. Singh [9] propose a novel soft computing method for RUL prediction. Al-Dahidi et al [10] proposed methodology for RUL prediction for a heterogenous fleet operating under normal and severe operating conditions. Saxena [11] proposed RUL prediction model by including run to failure data. The model proposed RUL prediction and Binary and Multi class classification. Zhang [12] et al proposed a RUL estimation for a rotating element like bearing, this method is significant as it is useful for many rotating equipment. Kim [13] et al proposed a RUL prediction for multi sensor data set. Wang et al [14] A generic probabilistic framework for structural health prognostics and uncertainty management. Song et al [15] proposed Statistical degradation modeling and prognostics of multiple sensor signals via data fusion: A composite health index approach. The method uses most data analytics techniques. Liu [16] Integration of data fusion methodology and degradation modeling process to improve prognostics. There are many approaches proposed in the literature to find the RUL prediction, these all methods are focused on machine intelligence. Due to this industrial system have limitation of direct application into industrial needs. Hence there is need for integrating the machine intelligence and business intelligence

### 3. Problem Statement

The data generated from finance, insurance, and other allied industries are popularly analyzed by methods and models developed over a period of past few decades. Usually these models consider cannot be directly used for sensor data. Whereas in sensor data analysis the way data being captured makes it different from other data. These sensors are located at various healthy and unhealthy condition order to apply such models for industrial equipment a new framework is needed as the data is time dependent in nature and also there is a need for integrated approach by integrating all forms of information available like asset life, asset make, operation data.. Initially the frequency of t data collection is fixed at device level and further it will be stored in cloud data base for further processing. Data that is available at central cloud can be grossly viewed as. Context awareness understanding the operator behavior. As shown in the fig 2 each of the machine is connected to other machine using a unique id for each of the machine, and also each of the machine consists of numerous sensors which are self-powered along with unique identity for each of them. All the sensors collect the data at various time stamps starting from seconds. During the operation, the data will be captured from a healthy state to degradation state. Hence current it is important to connect the data both in sensors level and time stamp level.



**figure. 2.** Sensor Data Analysis Framework

- Time series Collection:

In this approach the data is collected in seconds, microseconds. These values are then averaged for minute hour data wise and week wide. A machine with multiple of sensors product huge amount of data which need to average for suitable time window for further analysis. Timeseries data needs to be checked for frequency of captured and healthy data capture.

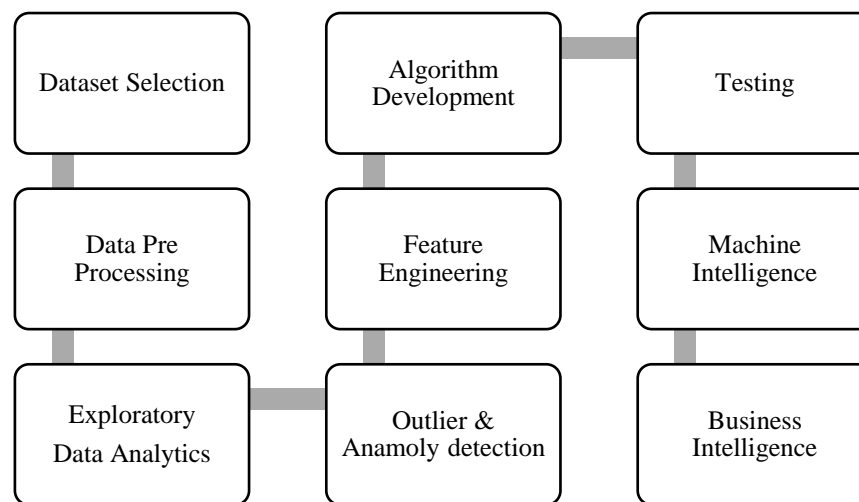
- Cyclic series Collection:

Usually many machines are operated with regular frequency. Typically engineering systems like aircraft consider each flight as a cycle of operation with specific time of flight for each of the cycle with specific load. For example, the health condition of an engine may vary in tens of hours, whereas its dynamic response is on the order of seconds. Identification of the current state of the engine health at slow-time epochs is very important for maintenance engineers because necessary repairs must be carried out much before the engine becomes seriously damaged or permanently non-operable. Thus, it is essential to monitor slow-time-scale anomalies for gas turbine engines from the time series data of the engine observables on a fast time scale. To this end, a generic gas turbine engine simulation test bed has been used to validate the anomaly detection method.<sup>2</sup> The features of the engine simulation model are like those of the engine model

#### 4. Exploratory Data Analysis

Data driven models need data generated from sensors that can capture physics related parameters like Pressure, Temperature, Flow etc. Current research uses the public dataset of NASA PCOE. The contains operational history of 100 aircraft engines. The dataset considered for current study is provided by NASA PCOE, the dataset contains observations of the 21 sensors like pressure temperature flow rate etc. The data is averaged values of these parameters each cycle of operation. Typically, the data is captured till the end of the life of the engine, i.e. run to

failure data of the engine is available. The data set contains the operational history of the aircraft engines which is recorded over a period of operations from 100 engines. The data is available in csv and text format which needed to smoothen to make it useful for model consumption. The data is processed using python analysis package pandas. The data is converted from standard CVS format to panda's framework understanding the data nomenclature, missing values, null entries etc. The data contains around 100 cycles on the average for each of the engine. Fig 3 Depicts the overall data analysis procedure adapted in current research.



**Figure.3.** Analytics framework for Sensor Data Analysis

Usually the data is available in standard enterprise database collection system that collect sensor data and store the data in frameworks like SAP, etc. These data files further saved as excel, csv, txt for facilitating the ease of data analysis on specific modules with smaller data samples. Following table depicts a portion of training data that is captured for engine data. Initially all data set are checked for presence of null values, NA values, empty cells so that the model do not have interfere with model processing.

**Table.1.** Dataset of Aircraft Sensor Data

id	cycle	setting1	setting2	setting3	s1	s2	s3	s4	s5	...	s12	s13	s14	s15	s16	s17	s18	s19	s20	s21	
0	1	1	-0.0007	-0.0004	100.0	518.67	641.82	1589.70	1400.60	14.62	...	521.66	2388.02	8138.62	8.4195	0.03	392	2388	100.0	39.06	23.4190
1	1	2	0.0019	-0.0003	100.0	518.67	642.15	1591.82	1403.14	14.62	...	522.28	2388.07	8131.49	8.4318	0.03	392	2388	100.0	39.00	23.4236
2	1	3	-0.0043	0.0003	100.0	518.67	642.35	1587.99	1404.20	14.62	...	522.42	2388.03	8133.23	8.4178	0.03	390	2388	100.0	38.95	23.3442
3	1	4	0.0007	0.0000	100.0	518.67	642.35	1582.79	1401.87	14.62	...	522.86	2388.08	8133.83	8.3682	0.03	392	2388	100.0	38.88	23.3739
4	1	5	-0.0019	-0.0002	100.0	518.67	642.37	1582.85	1406.22	14.62	...	522.19	2388.04	8133.80	8.4294	0.03	393	2388	100.0	38.90	23.4044

Id: Identity number of the engine

Cycle: count of cycle of operation

Set 1: Operation mode 1

Set 2: Operation Mode 2

Set 3: Operation mode 3

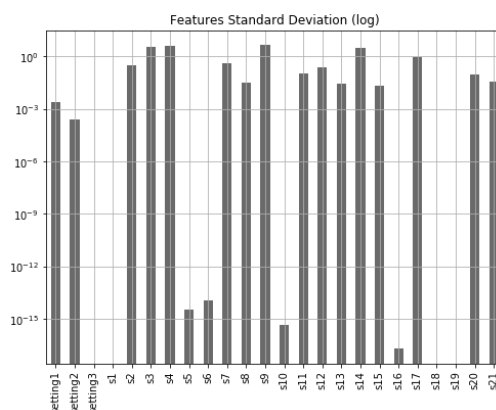
S 1-21: Values of captured parameters for 21 sensors of various physical parameters

In Data wrangling step dataset is checked for presence unwanted data element like Null or NAN, INF. Exploratory Data Analysis is conducted to understand the basic statistical distribution of sensor values. In this analysis maximum, minimum, range, mean and standard deviation of all the variable is found. The same is presented in below table.

Table.2. Key Statistics of the variables

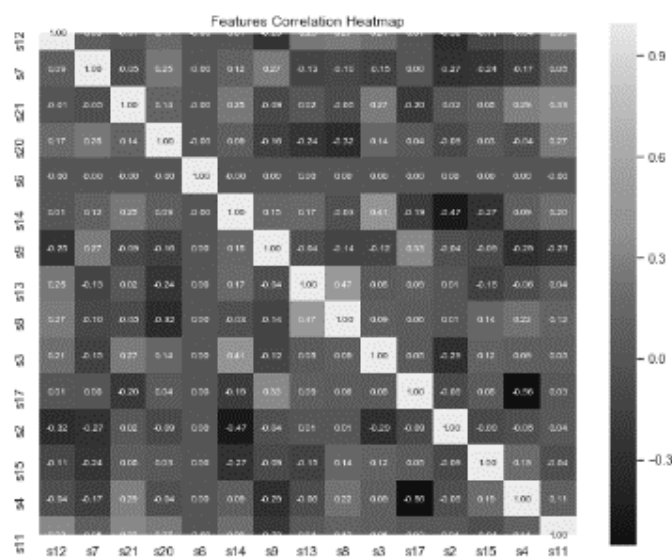
	Id	cycle	setting1	setting3	s1	s2	s3	s4
Count	20631	20631	20631	20631	20631	20631	20631	20631
Mean	51.50657	108.8079	-8.9E-06	100	518.67	642.6809	1590.523	1408.934
Std	29.22763	68.88099	0.002187	0	6.54E-11	0.500053	6.13115	9.000605
Min	1	1	-0.0087	100	518.67	641.21	1571.04	1382.25
25%	26	52	-0.0015	100	518.67	642.325	1586.26	1402.36
50%	52	104	0	100	518.67	642.64	1590.1	1408.04
75%	77	156	0.0015	100	518.67	643	1594.38	1414.555
Max	100	362	0.0087	100	518.67	644.53	1616.91	1441.49

One of the key parameters of the above statistics is standards deviation, the standard deviation value of the variable indicates the probability of that sensor having significant role in predicting the target variable. Referring to Fig 3 standard deviation of the of sensor 9 is maximum in comparison to all other sensors. Standard deviation can also help anomaly and outlier detection.



**Figure. 4.** Standard Deviation of all variables

As the data involves the multi variables and the relation between these variables is key phenomenon to develop predictive model. A correlation among the sensors is found using pandas correlogram. Correlation matrix represents the interdependency of the each of the parameter on the other parameter. A higher correlation 1 represent they have stronger interdependence. Following figure indicate that the relation between sensor 9 and sensor 11 similarly there are correlation values for each of the sensors



**Figure. 5.** Correlation Matrix among all variables

## 5. Solution Approach

Predicting the future value of a given variable is based on numerous statistical and machine learning models. These are models that can predict for single variable to multiple variables. Usually the models are built from a simple linearly to complex neural network models. These models and methods start with fundamental mathematical relation between single or multiple independent variables to dependent variable. Linear Regression is a machine learning algorithm based on supervised learning. It performs a regression task. Regression models a target prediction value based on independent variables. It is mostly used for finding out the relationship between variables and forecasting. Different regression models differ based on – the kind of relationship between dependent and independent variables, they are considering, and the number of independent variables being used. Learn regression models are more suitable for parameters which are linearly dependent on each other. Current research aims at finding the time to failure as accurately as possible hence initially time to failure for all the sensor values for all engines. Hence in current problem time to failure is treated as dependent variable.

For any predictive analysis linear regression model is a basic approach to solve for getting initial estimate in current study. Segment training and test data into features data frame and labels series.

$$AX + B = y \quad (1)$$

Where

$$X = \begin{bmatrix} x1 \\ x2 \\ x3 \end{bmatrix} \quad (2)$$

Vector x represents the variables that represent sensor values. There will be n number of variables corresponding to number of sensors

$$A = \begin{bmatrix} a11 & a12 & a13 \\ a21 & a22 & a23 \\ a31 & a32 & a33 \end{bmatrix} \quad (3)$$

Matrix A represents the co-efficient corresponding to each of the sensor. These coefficients indicate the effect of variable on the dependent variable.

$$B = \begin{bmatrix} b1 \\ b1 \\ b1 \\ b2 \\ b3 \end{bmatrix} \quad (4)$$

Matrix B represents the intercept or constant value of linear model. This will be a constant value that will contribute for the dependent variable

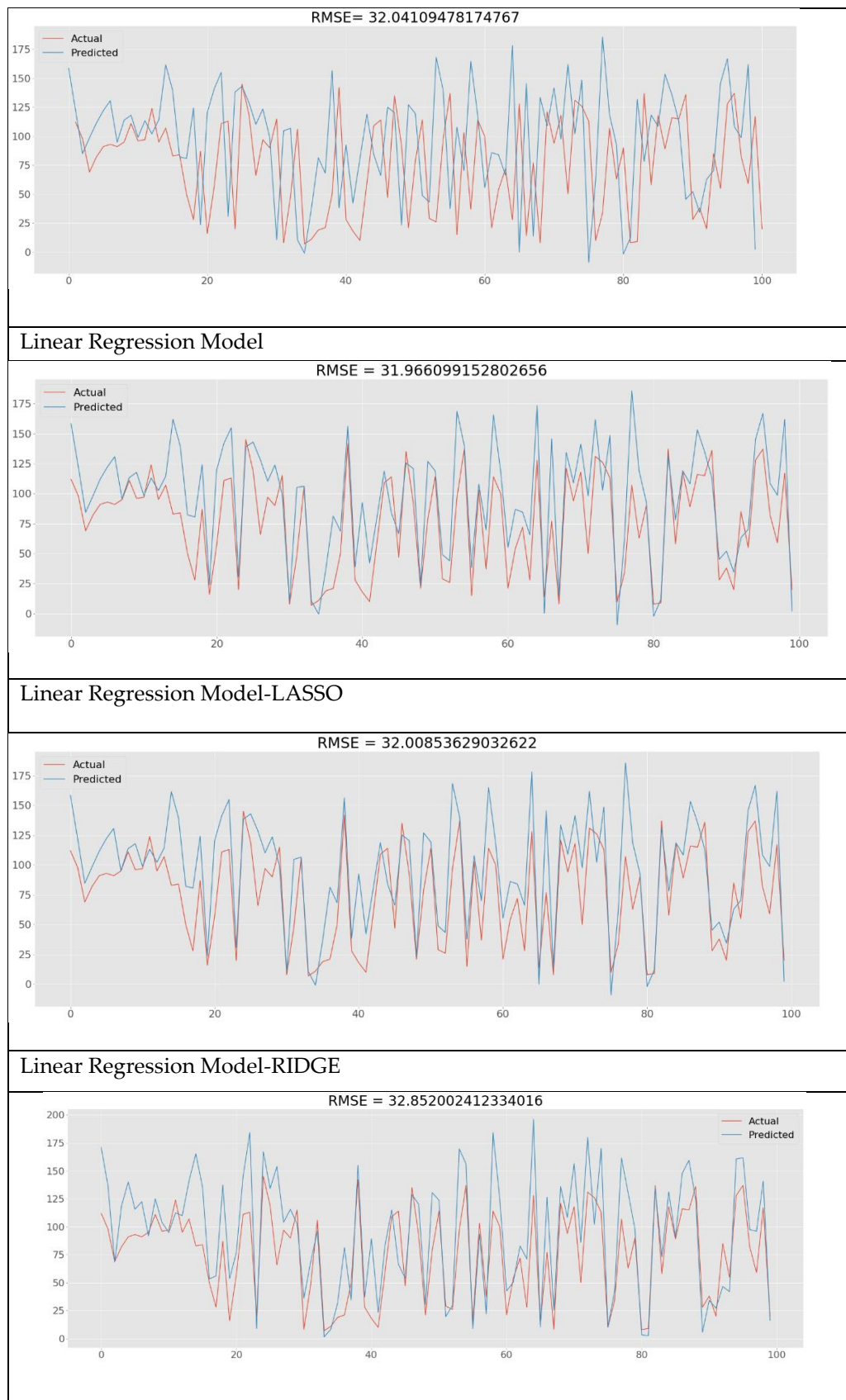
## 6. Prediction of RUL (Remaining Useful Life)

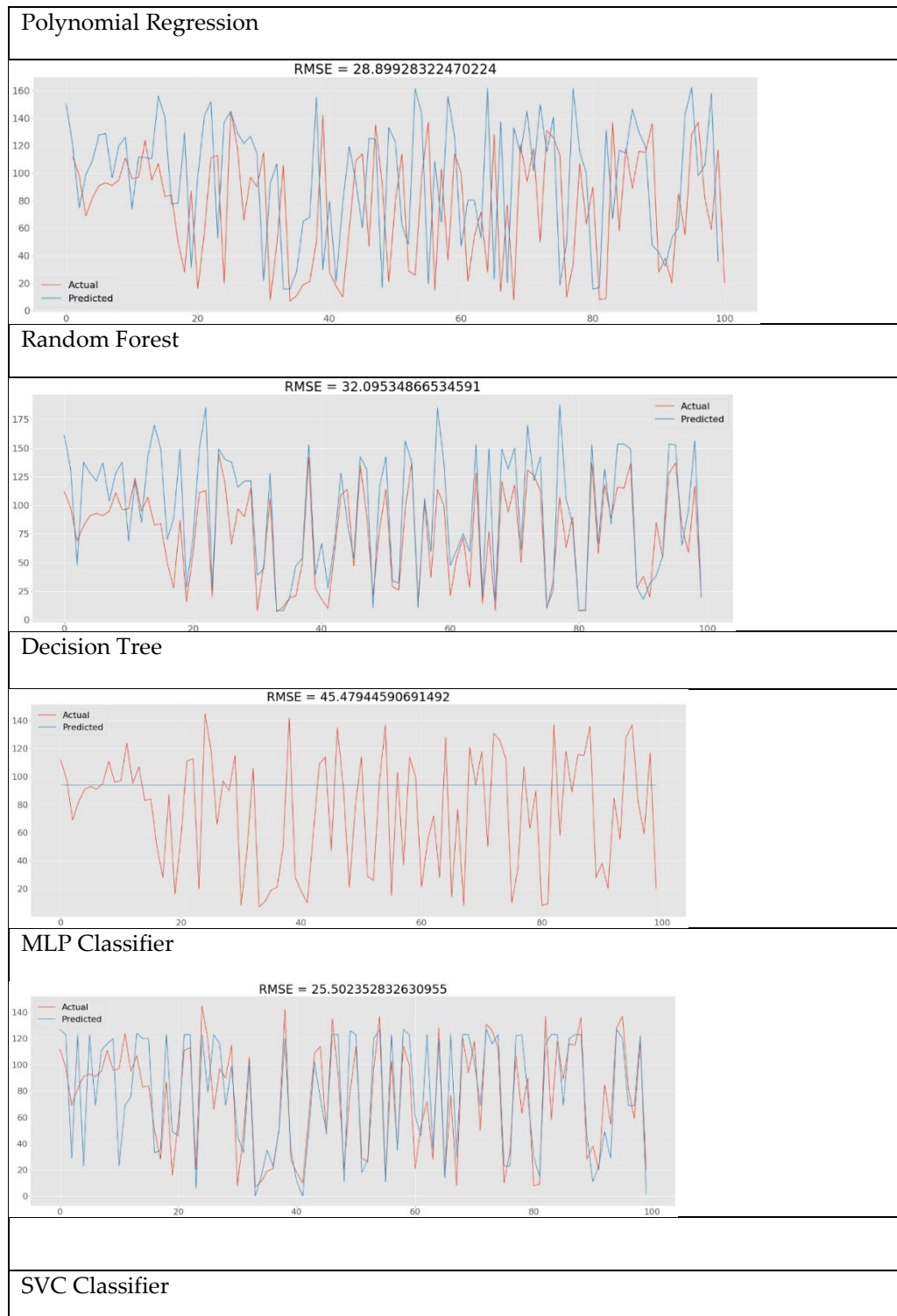
Remaining Useful life of the component is predicted by using following models which includes many regression and neural network models. There are combination of regression and combined models like Linear, LASSO, RIDGE, Polynomial Regression, Random Forest, Decision Forest, Decision Tree, SVC, MLP Classifier.

	RMSE	R2	MAE
Linear Regression Model	32.041	0.4054	25.591780
Linear Regression Model-LASSO	31.966099	0.408275	25.551808
Linear Regression Model-RIDGE	32.008536	25.570256	0.406703
Polynomial Regression	32.852002	0.375023	25.106250
Random Forest	28.899283	0.516368	23.360043
Decision Tree	40.917478	0.030477	26.340000
MLPClassifier	45.479446	-0.197763	36.480000
SVC	25.502353	0.623382	19.690000



Predicting Time to Failure is one of the key requirements for predictive maintenance algorithm. Various Models for time to failure are analyzed among all models SVC could predict the time to failure of 26 cycles, and MLP Model could predict RMSE of 45. These predictions are an initial estimate of remaining useful life. Improved Machine Learning methods are further developed in this research as further presented in later chapters.



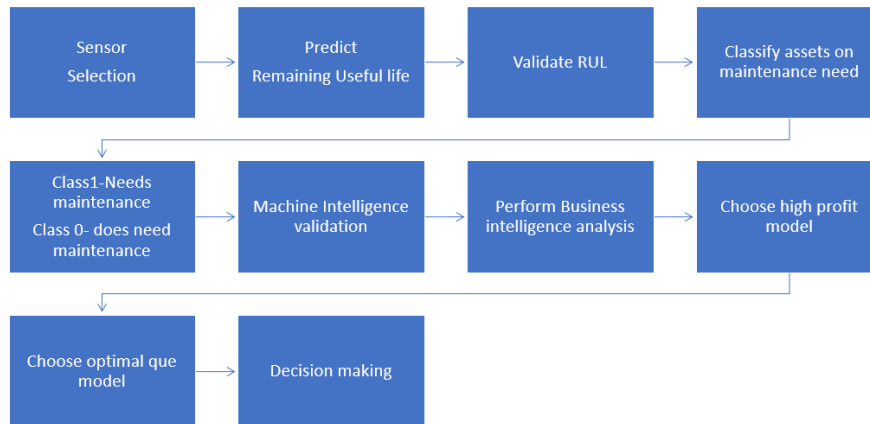


**Figure.6.** RMSE Comparison of Actual & Predicted values Various Models

## 7. Integrated Machine & Business algorithm

One of the challenges of widespread application IIoT system is lack of business intelligence models along with data intelligence models. Any deployment of large number of sensors and cloud-based data storage systems needs to be considered keeping in view of the cost benefit analysis. Proposed is unique and novel framework that includes both the machine intelligence and business intelligence to generate immediate alerts and autonomously plan the operations and maintenance. For

transportation industry like air transportation current framework can automatically plan the operations and maintenance at each engine so that there is no cost associated with manual error in choosing the machine to be sent for maintenance. This framework also generate alerts based on the delayed maintenance.



**Figure 7** Integrated Business & Machine Intelligence framework

## 8. Binary Classification

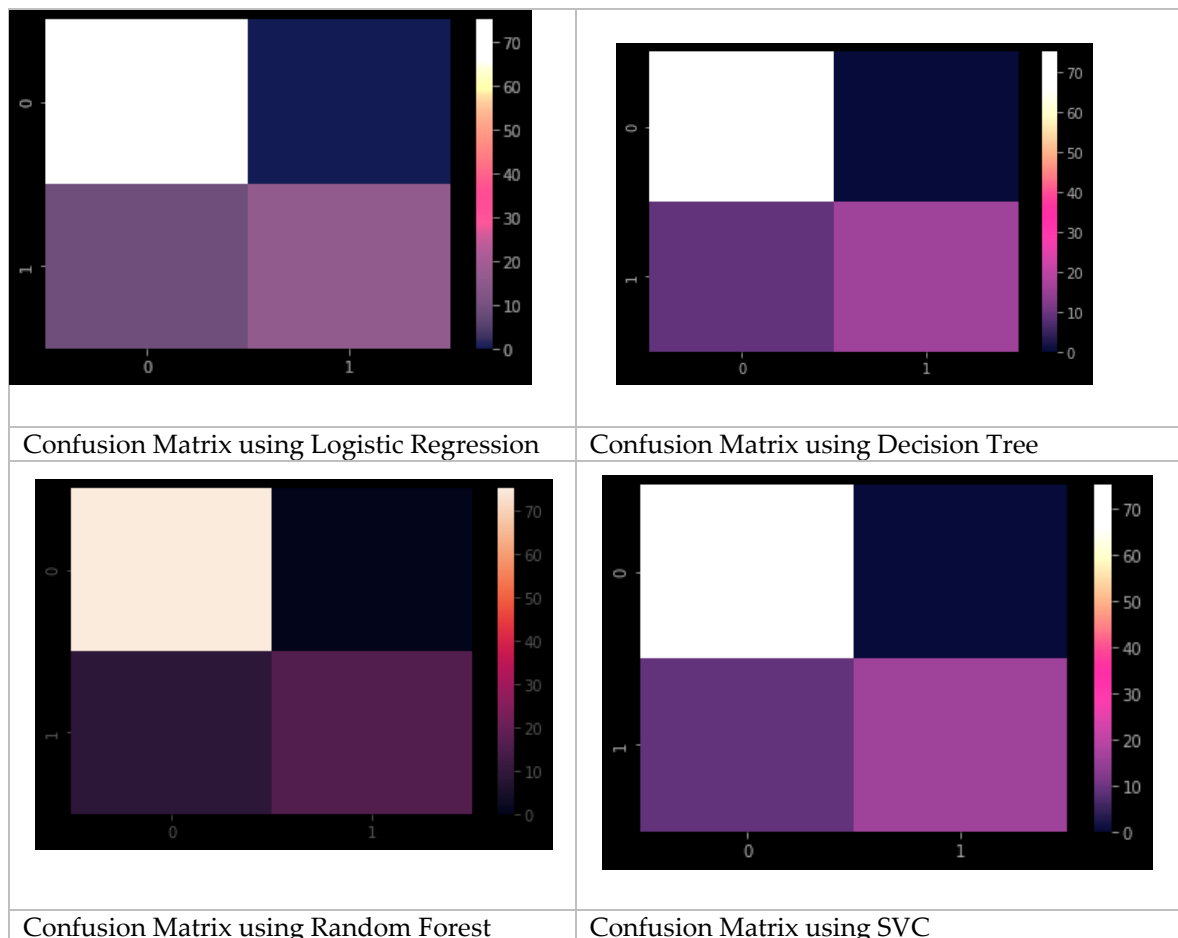
Binary classification is key aspect of the remaining life estimation in this approach a binary classifier is modelled to find the in this approach the modelling is done to classify the engines that will fail within a stipulated period let us say 30 days/cycle. Binary classification will help the decision maker to decide the which of the engines will fail in each time. In current analysis a period of 30 cycles is taken. Binary classification when coupled with business decision making algorithm will provide an integrated platform. Following models are considered for binary classification and key metrics are reported in the below table.

	Logistic Regression	Decision Tree	Random Forest	SVC	SVC Linear	KNN	Gaussian
Accuracy	0.910000	0.880000	0.910000	0.7500	0.770000	0.910000	0.940000
Precision	1.000000	0.933333	0.944444	0.0000	1.000000	0.944444	0.827586
Recall	0.640000	0.560000	0.680000	0.0000	0.080000	0.680000	0.960000
F1 Score	0.780488	0.700000	0.79069	0.0000	0.148148	0.790698	0.888889

ROC-AUC	0.979733	0.945067	0.98026	0.8096	0.971733	0.935200	0.987733

Referring to above table model like Logistic Regression, Decision Tree, Random Forest, SVC, SVC Linear applied to classify the engines as whether they need maintenance and do not need maintenance. Out of all the variables SVC Linear is optimal model in terms of accuracy and ROC-AUC. Hence the same is chosen as Binary classifier. As described in the figure 7 and below table Confusion matrix for all models is presented graphically, as we can see as the models. Linear Regression and SVC Linear Gave maximum number of True positives and maximum True Negatives. These models are optimal for current maintenance algorithm because they will reduce the cost of maintenance of healthy engines which would otherwise operate a greater number of cycles.

	TN	FP	TP	FN
Linear Regression	75	0	9	16
Decision Tree	74	1	11	14
Random Forest	74	1	8	17
SVC	75	0	25	0
SVC Linear	75	0	23	2



**Figure 8** Comparison of Confusion Matrix

## 9. Business Intelligence Analysis

Decision making is critical aspect of any industrial IoT system. It is important that the algorithm should give the outcome in terms of profits the business can reap. In current analysis KNN model gave maximum business benefit. This is based an assumed cost and profit for maintenance and

operation. The same is presented in below table. In many of MRO operation the capacity to few engines is redecided hence if the que is taken as criteria then the logistic regression model would be most suitable.

	Profit	Model	Que	Threshold	TP	FP	TN	FN
0	19	KNN B	0.31	0.090529	25	0	69	6
1	17	Random Forest B	0.33	0.097536	25	0	67	8
2	15.971	Logistic Regression B	0.24	0.21081	22	3	73	2
3	13.054	SVC Linear B	0.27	-1.13252	22	3	70	5
4	10.7027	KNN B	0.26	0.307692	21	4	70	5
5	7.816	Decision Tree B	0.29	0.082857	21	4	67	8
6	-3.92736	SVC B	0.31	-0.99891	18	7	62	13

If the capacity is fixed, then the at 31% then the KNN model is ideal for the maintenance. These models can be integrated with IIOT system to help take the business quick decisions based on the

	Profit	Model	Que	Threshold	TP	FP	TN	FN
0	15.97368421	Logistic Regression B	0.24	0.210810489	22	3	73	2
1	10.7027027	KNN B	0.26	0.307692308	21	4	70	5
2	13.05479452	SVC Linear B	0.27	-1.1325216	22	3	70	5
3	7.816901408	Decision Tree B	0.29	0.082857143	21	4	67	8
4	19	KNN B	0.31	0.090529294	25	0	69	6
5	-3.927536232	SVC B	0.31	-0.99890808	18	7	62	13
6	17	Random Forest B	0.33	0.097536479	25	0	67	8

## 10. Conclusion & Future Work

Industrial IoT is an emerging technology which can significantly help in improving the safe and intelligent operation of many machines. A successful application of such technology needs to be integrated with machine learning models and business benefit. Current research is aimed at developing an integrated framework for industrial IoT. In this study machine learning models like KNN, Decision Tree, Logistic Regression, SVC are explored to develop a framework which can combine both machine intelligence and business intelligence. Many of the current problems in industries lack unified and integrated decision support system for applying the industrial internet. Current research can further improve to integrate the multiple needs of the industrial plans.

## References

- [1] B. D. Minor, J. R. Doppa and D. J. Cook, "Learning Activity Predictors from Sensor Data: Algorithms, Evaluation, and Applications," *IEEE transactions on knowledge and data engineering*, vol. 29, no. 12, pp.2744–2757, 2017.
- [2] K. Ravi, B. Prajna, "Engineering Optimization using Artificial Neural Network," *International Journal of Innovations in Engineering & Technology (IJJET)*, , vol. 4, no. 3, pp. 63–72, 2014.
- [3] M. Alexandra, P. Marko, P. Maria, F. Carolina and M Dunja, "Using Machine Learning on Sensor Data", *Journal of Computing and Information Technology* vol. 18, no. 4, pp. 341–347, 2010.
- [4] H.A. Taha, A.H. Sakr, Y. Soumaya, "Aircraft Engine Remaining Useful Life Prediction Framework for Industry 4.0", in *Proc. ICIEOMT*, Canada, pp. 23-25, 2019.
- [5] C. Okoh, R. Roy, J. Mehnen, L. Redding "Overview of Remaining Useful Life Prediction Techniques in Through-Life Engineering Services" in *Proc. CIRP 16* pp.158 – 163, 2014.
- [6] A. L. Ellefsen, B. Emil, A. Vilmar, S. Ushakov, H. Zhanga "Remaining useful life predictions for turbofan engine degradation using semi-supervised deep architecture" *Reliability Engineering and System safety* 183, pp.240-251, 2019.
- [7] R R. Zhou, N. Serban N. Gebraeel, "Degradation modeling applied to residual lifetime prediction using functional data analysis", *The Annals of Applied Statistics*, Vol. 5, no. 2B, pp. 1586–1610, 2011
- [8] H. Hanachi, C. Mechefske, J. Liu, A. Banerjee, Y. Chen, "Performance-based gas turbine health monitoring, diagnostics, and prognostics: A survey" *IEEE Trans. Reliab.* Vol. 67, no. 3, pp.1340–1363 . 2018.
- [9] S.K. Singh, S.Kumar, J.P Dwivedi, "A novel soft computing method for engine RUL prediction", *Multimed. Tools Appl* Vol. 78, pp. 4065–4087, 2019.
- [10] S. Al-Dahidi, F. Di Maio, P. Baraldi, E. Zio, "Remaining useful life estimation in heterogeneous fleets working under variable operating conditions," *Reliab. Eng. Syst. Saf.*, Vol 156, pp. 109–124. 2016.
- [11] A. Saxena, K. Goebel, D. Simon, N. Eklund, "Damage Propagation for Aircraft Engine Run-To-Failure Simulation," in *Proc. ICPHM*, Denver, CO, USA, PP. 6–9 2008.
- [12] B. Zhang, L.J. Zhang, J.W. Xu, "Degradation feature selection for remaining useful life prediction of rolling element bearings," *Qual. Reliab. Eng. Int.*, no.32, pp. 547–554. 2016.
- [13] M. Kim, C. Song, K. Liu, "A generic health index approach for multisensor degradation modeling and sensor selection," *IEEE Trans. Autom. Sci. Eng.*, no.16, pp.1426–1437, 2019.
- [14] P.F. Wang, B.D. Youn, C. Hu, " A generic probabilistic framework for structural health prognostics and uncertainty management," *Mech. Syst. Signal Process.*, no. 28, pp. 622–637. 2012.
- [15] C.Y. Song, K.B Liu, " " *IIEE Trans.* no. 50, pp. 853–867 ,2018.
- [16] Liu, K.B.; Huang, S. *IEEE Trans. Autom. Sci. Eng.*, no. 13, pp. 344–354, 2016.