

**Data Mining Approaches in Healthcare**  
**Operations to Improve Decisions and Activity**  
**Workflow**



*Submitted By*

Avinash Yadav (ITM2016004)

**INDIAN INSTITUTE OF INFORMATION  
TECHNOLOGY, ALLAHABAD**

(A UNIVERSITY ESTABLISHED UNDER SEC.3 OF UGC ACT, 1956 VIDE NOTIFICATION NO.  
F.9-4/99-U.3 DATED 04.08.2000 OF THE GOVT. OF INDIA)

A CENTRE OF EXCELLENCE IN INFORMATION TECHNOLOGY ESTABLISHED BY  
GOVT.OF INDIA

**May 2021**

## **1. Abstract:**

Healthcare organizations have the potential to improve their existing system in many different ways with available resources. They need modern techniques to deliver cost-effective and efficient services to their customers. Due to advancements in information technology, every record is stored in an electronic format. This advancement has increased the complexity in many domains of healthcare organizations. However, many advanced software systems still exist in various fields of the healthcare system. Still, the exploration in the advancement of the Operations domain is needed. Data mining algorithms are versatile and fit to complex examples, might be especially appropriate to taking care of these issues. Two significant benefits of machine learning, the force of building solid models from countless feebly prescient highlights, and the capacity to distinguish key elements in complex capabilities have an especially immediate association with the foremost operational difficulties. The principal objective of this work was to consider this relationship utilizing two significant sorts of operational issues: predicting operational activities delay, and recognizing key work process drivers. Utilizing pragmatic models, we show how machine learning can improve human capacity to comprehend and oversee medical care tasks, prompting more effective medical services.

## **2. Introduction:**

Healthcare Organizations play an important role in our life. The healthcare system in the US is unique and way different than other countries. The expenditure on the US healthcare system is much more than the entire GDP of some countries. With the introduction of many acts and amendments, the US healthcare system is one of the advanced systems in the world which runs on a hybrid system of multi-payers and multi-providers. 89.6 percent of the population pays around 30 percent of their income for their health insurance [1]. Healthcare organizations have the potential to improve their existing system in many different ways with available resources. They need modern techniques to deliver cost-effective and efficient services to their customers. Due to advancements in information technology, every record is stored in an electronic format. This advancement has increased the complexity in many domains of healthcare

organizations. However, many advanced software systems still exist in various fields of the healthcare system. Still, the exploration in the advancement of the Operations domain is needed. The healthcare system in the United States has grown its conventional methods of management. Therefore, it has expanded the dimensions of complexity and uncertainty in healthcare operations. To solve the real-time issues, there is a need for adaptive models with the help of medical records and machine learning algorithms.

The potential of artificial intelligence and machine learning models for medical care tasks and operations has been neglected to a great extent. Until this point, inside healthcare services, machine learning algorithms and models were used for clinical applications instead of healthcare activities, for example, Speech conversion of patient's record and prediction of diseases from the available symptoms from the previous record of the patient [2]. However, machine learning techniques could be revolutionary for healthcare operations management. Currently, the operations of the healthcare system are very complex and events from the registration window to post-service are cascading. Delays or operational errors will impact all further events, thus customer's experience. These intrinsic difficulties have been compounded by expanding work and capital expenses, a deficiency of qualified clinical staff, and restricted offices, prompting an assortment of serious operational issues, for example, restricted admittance and overcrowding [3]. Additionally, operational disappointments in healthcare have both clinical and monetary expenses, straightforwardly affecting patient wellbeing and prosperity.

These difficulties have been supplemented by the raising multifaceted nature and broadness of information, developing electronic medical records with over 25,000 petabytes of data [4]. Electronic Medical records (EMR) have been furnishing clinical focuses with long stretches of operational information, for example, handling timestamps, booking and managing records, assessment types, and different asset qualities, recorded regularly by emergency Hospital Information Systems (HIS). These systems are defined in the standard format provided by US regulatory agencies. The

Standard healthcare records are stored in the format of Health Level 7 (HL7), Fast Healthcare Interoperability Resources (FHIR), or Digital Imaging and Communications in Medicine (DICOM) [5]. The data stored in these standard formats are further divided into Hospital Information Systems, medical imaging, and other sources. These data points describe the healthcare operations that are very complex to analyze for managers due to many attributes. Thus, there is much scope of research in the healthcare operations domain.

### **3. Literature review:**

Due to the complexity of healthcare operations, some experiments had been conducted using simulation and modeling techniques to understand the potential of problems and drilling down deep to find some results. Several simulation experiments were conducted in the exploration of healthcare potential areas in the Netherlands [6], discrete event simulation design for an emergency department to optimize the staff level operations [7]. Expert Systems were also designed with constraints for the analysis of patient records [8]. The expert systems are data-driven and rule-based which are developed on many assumptions with the help of domain experts. Thus, they require time to design rules and collect relevant information before developing an expert system. After all these efforts, there are so many deviations in results compared to real-time operations due to non-adaptability to the changing environment. Therefore, these expert systems are not used in real-time predictions. Apart from it, some experiments were more data-oriented. Analysis of variance in medical records has been done to maximize the operational efficiency using the ambulance surgery model [9], and for the development of highly efficient rooms for operational work [10], Medical team training [11] was utilized to receive patient waiting time information to upgrade working room effectiveness. The discussed cases follow the strategies that depend on inferential measurements, are somewhat restricted in degree because of presumptions about the fundamental information structure, and are restricted to help personalized understanding. Therefore, more strong models are needed to understand the operational problems in healthcare.

To be helpful, the healthcare operations facility ought to run progressively which requires minimum presumptions, and consequently recognize key highlights liable for all events and should be adjusted according to the environment due to variability in the healthcare operations activities. Due to the tremendous amount of information and variability in nature, machine learning algorithms can be used for exploration in this domain. Similar kinds of situations have been tackled previously in different industries. Predicting the resources consumed by some APIs using machine learning [12], problems detection in traffic control along with the dynamic time required in routes using machine learning algorithms [13], predicting the time required in the manufacturing process using machine learning [14], environment adaptability and scheduling with genetic algorithms [15] has been done in various industries, which infers us to explore machine learning algorithm in healthcare operation management to identify critical factors and design a workflow to improve the decision-making process.

#### **4. Objective:**

The objective of this study is to solve significant healthcare problems in operations management with the help of machine learning algorithms that can be implemented in real-time hospital information systems to improve the workflow and decision-making process. We'll be focusing on two significant features by predicting the delay in operational activities (waiting time for patients in walk-in/ scheduled facility) and identifying significant factors in operational events in the workflow.

#### **5. Methodology:**

##### **4.1 Data:**

The secondary data is used for this study. This dataset contains the healthcare operational attributes for 4 different healthcare providers collected through real HISs. Hence, to take care of the confidentiality of information, augmentation of data has been done by putting the future value of timestamps in arrival and scheduled time such that their relative order can be maintained to explore the operational events. Three of them have the facility of scheduling appointments for the customers, while one of them has a

walk-in facility. This data contains patient records who visited these facilities in 3 to 4 years. Each row represents the scheduled arrival or walks in the patient record (one line per each patient visit event). This data is available for the operational experiment by a medical group of Massachusetts general hospital to explore more.

### **Prediction of delay in operational activities (waiting time for patients in walk-in/ scheduled facility) :**

In the workflow of the cascading operational activities, the activities are dependent on several factors, these factors are not sufficiently correlated so that we can neglect the remaining factors. For example, the workflow delay in healthcare organizations may rely upon its staffing, scheduling, and flow of arrival of customers, time, season, frequency of assessments performed, the operational environment, climate, and occasional patterns, cafeteria lunch hours, and numerous different elements. Every one of these highlights, whenever considered independently, would generally have little impact on our prediction function. This functionality refers us to use machine learning lazy approaches [16] due to less correlation and abundance of factors in operational activities. The accuracy in the predicted delay will alert management to adjust operational activity accordingly thus improve the customer satisfaction level [17].

### **Identification of significant factors in operational events:**

The insightful capacity of machine learning gets from ideal component subset determination, a typical part of ML streamlining [18], ideal component subset determination can utilize different component choice calculations like pruning, elimination process, drilling down into data, and elimination to reject the most insignificant factors until the model finds the significant and most prescient arrangement of key highlights. This technique answers the most important questions in healthcare operational workflow which were unanswered or misinterpreted.

## 6. Discussion:

Prediction of delay in operational activity is determined as the contrast between two timestamps promptly accessible in HISs: for a stroll in arrangements, it is the distinction between understanding appearance and the start of the assessment; for planned arrangements, it is the contrast between the booked beginning and the real beginning (better communicated as 'delay'). While walk-in, stand-by times show exceptionally intricate and profoundly factor fleeting patterns, making them difficult to foresee with a straightforward normal or by taking a gander at the past quiet stand-by times. Also, ceaseless ecological changes, for example, non-weekend days versus ends of the week, distinctive staffing, and changes in office assets and patient volume call for versatile models. Machine learning lazy approaches have been tried to achieve the results. To accomplish the results there were two phases:

### Phase 1:

Identification of significant factors in operational activity, which is done by ideal component subset determination by correlation test of each facility independently. These factors were considered as the base for determining the prediction of delay in operational activity. The correlation test tells us which factors should be focused more on operational management by adjusting the time spent on the least significant factors to offer better services to the customers.

### Phase 2:

After determining ideal components in operational activity, some machine learning models have been used using python and sklearn libraries to accomplish the critical task of this study is the prediction of delay in operational activity. In this phase, models have been tested with elimination as well as consideration of highly correlated factors. These models are as follows; The Linear regression model considering ( most recent wait, moving average most recent waits, moving average based on queue length, best/all features with/without intercept), Linear Robust regression (TheilSenRegressor) using all attributes with/without intercept, ElasticNet regression using all attributes with/without intercept, Decision Tree, Boosted random-forest with 10/20/30 splits,

Neural Network single/multiple Layer. Here, Fig 1 shows the appointment data points vs waiting for time (mins) for the scheduled facility. The most accurate model with the best fit from ( boosted random forest with 30 splits/ ThielSen Regressor/ some Linear regression models ) is used to demonstrate the results. The blue lines are the actual waiting time while visiting the scheduled facility, while the best fit model can predict the waiting time close to the original one before past delays and the visit by the patient. This will surely help in improving the operational activity of the hospital and enhance the customer experience.

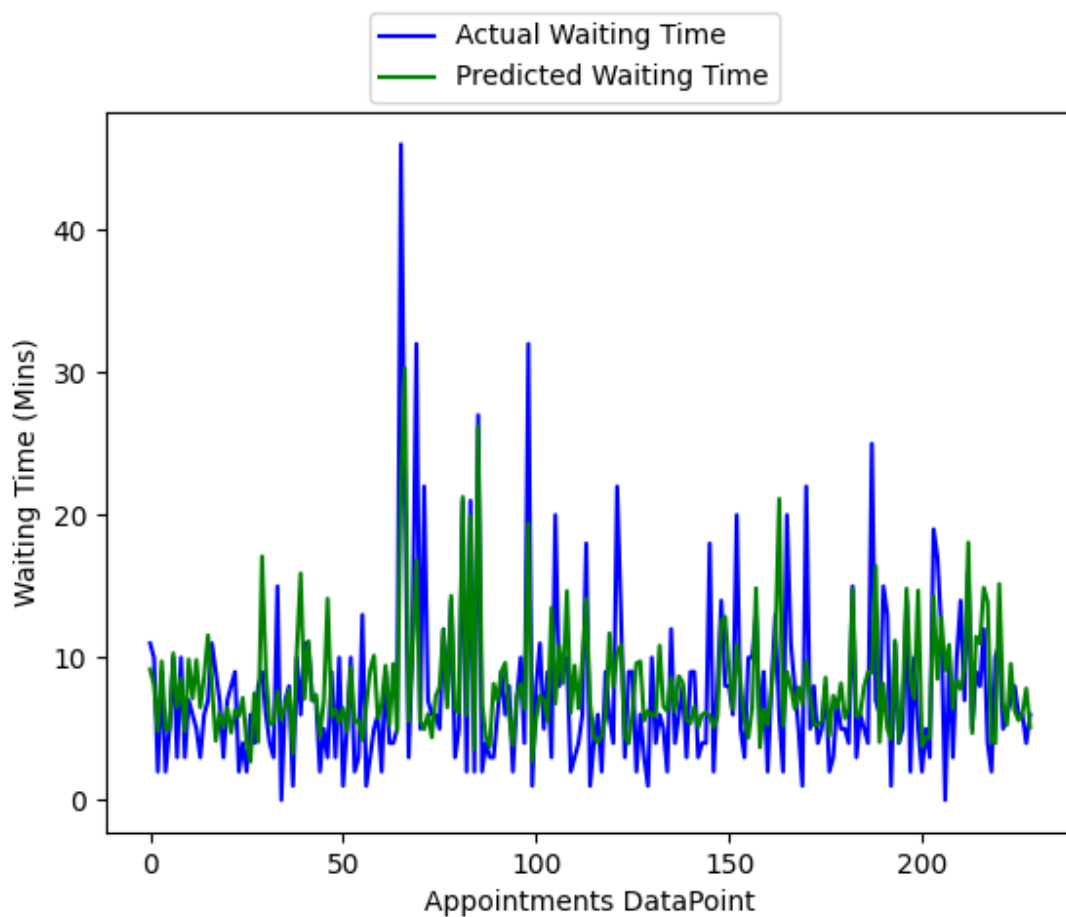


Fig. 1: Appointment data points vs waiting time (mins) for the scheduled facility

Other than the prediction of delay in operational activity, the models have been tested based on the inclusive dependent activities dedicated to several independent activities in the hospital facilities like the abdominal test, Thoracic test, cardiac test, and



so on. Operational activities are critical and dependent on many factors. The testing facilities in any hospital will impact the core operational activities. Hence, our model should predict accurately on time series data of these existing facilities to contribute to core operational activities. Here, Fig. 2, 3 are showing the fitness of the model discussed in the phase 2 section. The plot is drawn between the probability of getting a queue in the abdominal and Thoracic testing facility respectively vs time series data on an hourly basis. The black lines are showing the probability of getting a long queue whose length is greater than 6 at any timestamp in both plots. Green spikes are probability values of prediction of getting a queue whose length is greater than 6 at any timestamp in the abdominal testing facility, while blue spikes are probability values of prediction of getting a queue whose length is greater than 6 at any timestamp in the abdominal testing facility.

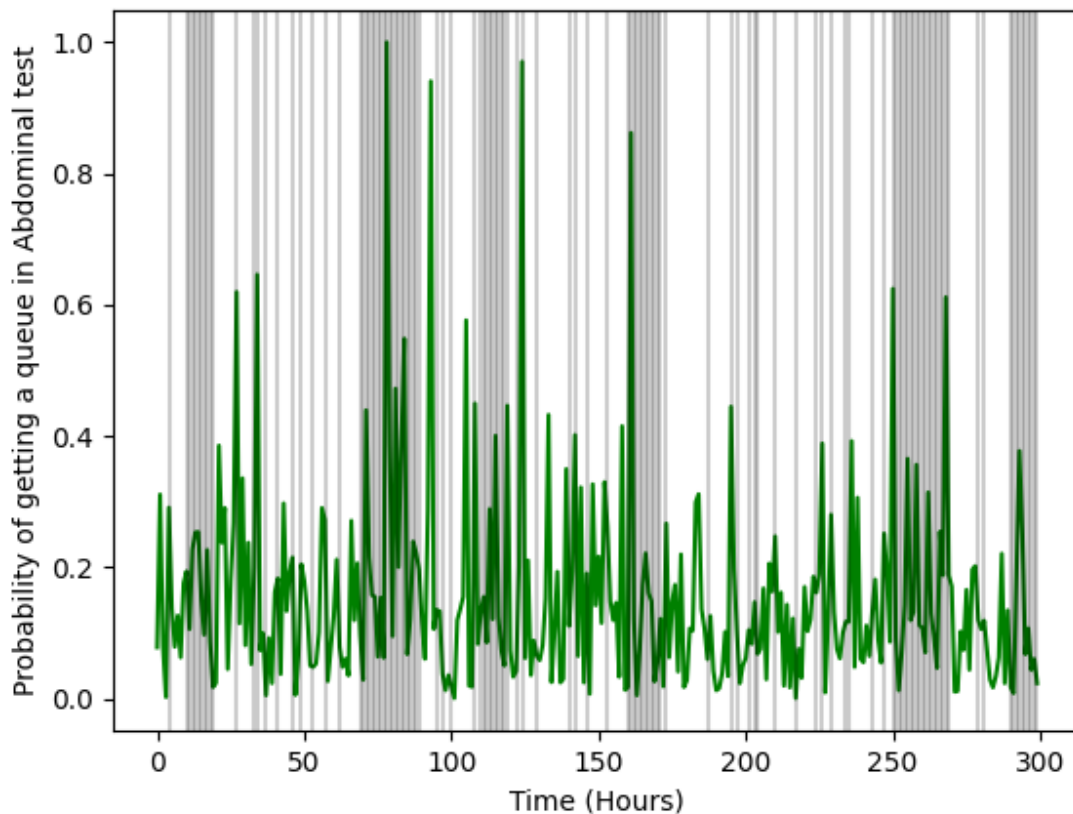


Fig. 2: F1 Facility Abdominal Test plot for the probability of getting a queue length greater than 6.

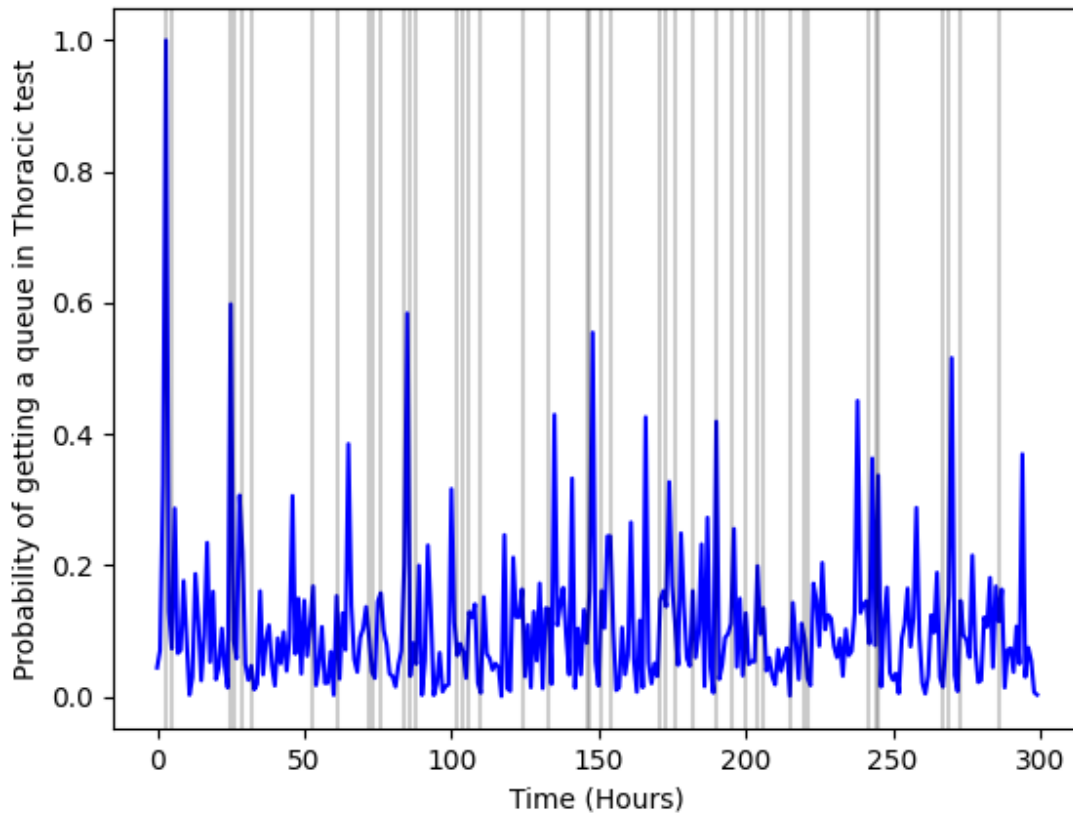


Fig. 3: F3 Facility Thoracic Test plot for the probability of getting a queue length greater than 6.

The upsides of Machine learning can be represented with Fig. 4, 5, looking at the precision of ML-based delay forecast of operational activity to the most far and wide pre-ML draws near (anticipating from past delays and their moving averages). To evaluate model quality, we utilized both  $R^2$  and mean absolute error (MAE), where a decrease in MAE was figured comparative with the MAE of the first stand by time. Discussed models were utilized to exhibit the advantages of lazy learning of machine learning. All model errors were figured on the test datasets, to prohibit overfitting. As Fig. 4, 5 illustrates, the utilization of ML brought about considerably higher model exactness. Expanding both the model intricacy and the model list of capabilities improves machine learning results. All things considered, anticipating hang tight occasions for that four facilities without machine learning would be incomprehensible.

There are cutoff points to how precise operational models could be? Indeed, based on our experience, we would describe them as main operational change abilities, recorded underneath.

1. Transient inconstancy: This sort of changeability is related to endeavors to foresee for longer timeframes, which inevitably makes expectations less precise (arbitrary occasions accumulating during a more extended interaction). It very well may be found in facility F2, where patient assessments take longer than other facilities.

2. Work process changeability: This kind of fluctuation is related to profoundly flighty and problematic occasions, where the most unsurprising F4 is the lone office taking arbitrary walk-in patients.

These two restrictions are established into their particular work processes, and not into the insufficiencies of the machine learning models. Profoundly random or abstract work processes would be difficult to demonstrate even with the most progressive machine learning. Subsequently, any operational cycles with helpless consistency ought to be first checked for operational consistency; this should prompt either refining the handling rationale or finding some new ML highlights to make it more unsurprising.



Fig. 4:  $R^2$  value for ML models in F1, F2, F3, F4 facilities

$$R^2 = 1 - \frac{RSS}{TSS}$$

$R^2$  = coefficient of determination

$RSS$  = sum of squares of residuals

$TSS$  = total sum of squares

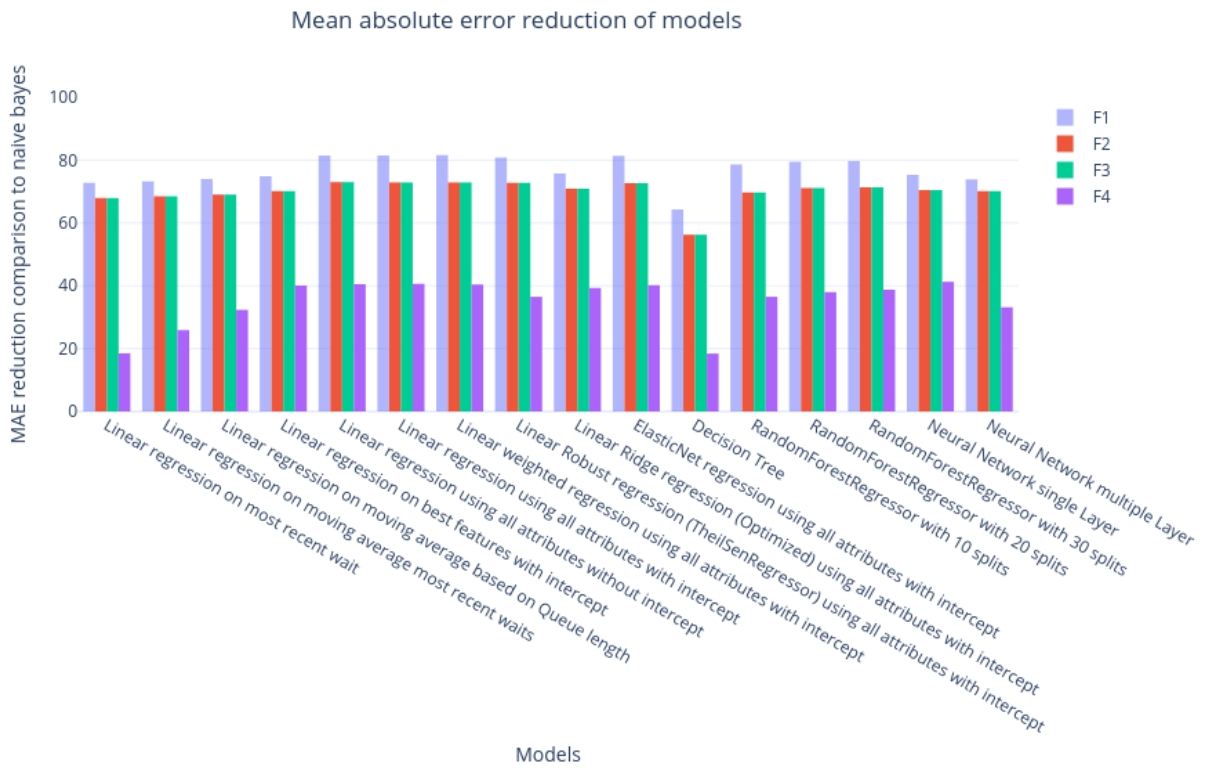


Fig. 5: MAE reduction compared to naive Bayes for ML models in F1, F2, F3, F4 facilities

$$MAE = \frac{\sum_{i=1}^n |y_i - x_i|}{n}$$

$MAE$  = mean absolute error

$y_i$  = prediction

$x_i$  = true value

$n$  = total number of data points

Here, Fig. 4 demonstrates the coefficient of determination value for tested models in each facility linear regression and boosted random forest models found to be the best among models. Similarly, Fig. 5 demonstrates the improvement in the MAE value in comparison to the naive Bayes model. The collective results along with the standard deviation in  $R^2$  and MAE test and training value have been analyzed with the proximity of prediction in delay, the finalized model has been presented in this section.

## **7. Managerial Implications:**

Worldwide normalization of medical services information. Since the 1990s, all medical services information components (operational included) are recorded utilizing a couple of major advanced principles (HL7 or FHIR). Accordingly, all operational information is gathered in a uniform and autonomous way. Consequently, porting an operational ML model to another clinical office doesn't need reevaluating the model or its list of capabilities. All things considered, the model necessities just to be retrained into the nearby work process designs. Adaptivity of machine learning is dissimilar to numerous pre-machine learning models dependent on static guidelines or coefficients, machine learning calculations are normally versatile. They needn't bother with human help to learn from their information. Keeping operational ML on an ordinary, robotized retraining plan guarantees that the models stay exact as their surroundings advance.

Utilizing delay in operational activities prediction, managers can explore the novel patient identifiers related to high anticipated stand-by times to examine the basic reasons for these deferrals. This can drive future operational changes (for instance, if forecasts are high despite apparently low asset usage). With work process models, one gets ready to anticipate queues in quiet lines some hours ahead of time. This was a management-driven edge, this time was the measure of time expected to bring in help staff to help in subduing that queues before they could bottleneck the framework. ML furnishes medical services directors with a novel, new chance to make educated, precise, ongoing choices. The effect of this methodology can be straightforwardly seen in the information, which affirms the precision of the model, yet in addition the general machine learning effect of the model-driven operational enhancements.

## **8. Conclusion:**

This study aims to exhibit the benefits of machine learning for health care activities the operations. To do as such, we bring machine learning into two significant classes of medical services operational issues as predicting delay in operational activities and recognizing key operational features. This study has demonstrated how two explicit properties of machine learning; lazy learning and ideal subset determination make machine learning especially reasonable for tackling operational issues. At last, the benefits and impediments of these arrangements and represent them with a few applications our gathering has effectively evolved and carried out in enormous consideration places.

Despite the fact that predicting delay in operational activities and recognizing key operational features are two headspaces of operational issues, other challenges like booking slots enhancement, execution and efficiency analysis, process exploration, quality considerations, and that's just the beginning can profit by machine learning methods. Also, numerous inflexible and not-taking operational calculations from the past can be upgraded with machine learning configuration, making them significantly more relevant in complex operational conditions. Subsequently, mixing area ideal calculations with machine learning information learning opens an incredibly encouraging course for tackling genuine operational issues, where both ML and traditional operational methodologies may improve and supplement one another, beating their current constraints. From multiple points of view, this 'learning' way to deal with tasks the management should change the whole facility's worldview. Operational cycles regularly carry on as characteristic wonders, where the standards should be found instead of constrained. As our models illustrate, machine learning lazy learning and ideal subset choice make these disclosures conceivable, straightforwardly adding to more proficient medical services.

## 9. References:

- [1] Department for Professional Employees, 2016. The U.S. HealthCare System: An International Perspective — Department for Professional Employees, AFL-CIO. [online] Department for Professional Employees, AFL-CIO. Available at: <<https://www.dpeaflcio.org/factsheets/the-us-health-care-system-an-international-perspective>>
- [2] Choy, G., Khalilzadeh, O., Michalski, M., Do, S., Samir, A. E., Pinykh, O. S., ... & Dreyer, K. J. (2018). Current applications and future impact of machine learning in radiology. *Radiology*, 288(2), 318-328.
- [3] Zhao, Y., Peng, Q., Strome, T., Weldon, E., Zhang, M., & Chochinov, A. (2015). Bottleneck detection for improvement of emergency department efficiency. *Business Process Management Journal*.
- [4] Kruse, C. S., Goswamy, R., Raval, Y. J., & Marawi, S. (2016). Challenges and opportunities of big data in health care: a systematic review. *JMIR medical informatics*, 4(4), e38.
- [5] Bender, D., & Sartipi, K. (2013, June). HL7 FHIR: An Agile and RESTful approach to healthcare information exchange. In *Proceedings of the 26th IEEE international symposium on computer-based medical systems* (pp. 326-331). IEEE.
- [6] Barjis, J. (2011). Healthcare simulation and its potential areas and future trends. *SCS M&S Magazine*, 2(5), 1-6.
- [7] Ghanes, K., Jouini, O., Jemai, Z., Wargon, M., Hellmann, R., Thomas, V., & Koole, G. (2014, December). A comprehensive simulation modeling of an emergency department: A case study for simulation optimization of staffing levels. In *Proceedings of the Winter Simulation Conference 2014* (pp. 1421-1432). IEEE.

[8] Brockway, M., Carlson, G. M., Kadhiresan, V., & Kovtun, V. (2008). U.S. Patent No. 7,433,853. Washington, DC: U.S. Patent and Trademark Office.

[9] Tayne, S., Merrill, C. A., Saxena, R. C., King, C., Devarajan, K., Ianchulev, S., & Chilingirian, J. (2018). Maximizing operational efficiency using an in-house ambulatory surgery model at an academic medical center. *Journal of Healthcare Management*, 63(2), 118-129.

[10] Attarian, D. E., Wahl, J. E., Wellman, S. S., & Bolognesi, M. P. (2013). Developing a high-efficiency operating room for total joint arthroplasty in an academic setting. *Clinical Orthopaedics and Related Research®*, 471(6), 1832-1836.

[11] Wolf, F. A., Way, L. W., & Stewart, L. (2010). The efficacy of medical team training: improved team performance and decreased operating room delays: a detailed analysis of 4863 cases. *Annals of Surgery*, 252(3), 477-485.

[12] Matsunaga, A., & Fortes, J. A. (2010, May). On the use of machine learning to predict the time and resources consumed by applications. In 2010 10th IEEE/ACM International Conference on Cluster, Cloud and Grid Computing (pp. 495-504). IEEE.

[13] Elhenawy, M. M. Z. (2015). Applying Machine and Statistical Learning Techniques to Intelligent Transport Systems: Bottleneck Identification and Prediction, Dynamic Travel Time Prediction, Driver Run-Stop Behavior Modeling, and Autonomous Vehicle Control at Intersections (Doctoral dissertation, Virginia Tech).

[14] Priore, P., Gomez, A., Pino, R., & Rosillo, R. (2014). Dynamic scheduling of manufacturing systems using machine learning: An updated review. *Ai Edam*, 28(1), 83-97.



[15] Jakobović, D., & Budin, L. (2006, April). Dynamic scheduling with genetic programming. In European Conference on Genetic Programming (pp. 73-84). Springer, Berlin, Heidelberg.

[16] De Mantaras, R. L., & Armengol, E. (1998). Machine learning from examples: Inductive and Lazy methods. *Data & Knowledge Engineering*, 25(1-2), 99-123.

[17] Holbrook, A., Glenn Jr, H., Mahmood, R., Cai, Q., Kang, J., & Duszak Jr, R. (2016). Shorter perceived outpatient MRI wait times associated with higher patient satisfaction. *Journal of the American college of radiology*, 13(5), 505-509

[18] Bertsimas, D., King, A., & Mazumder, R. (2016). Best subset selection via a modern optimization lens. *Annals of Statistics*, 44(2), 813-852.

## **10. Appendix:**

Raw data: [download sheet here](#)

Results: [download sheet here](#)

Code: [Available here](#)