

# Real-time prediction of inpatient length of stay for discharge prioritization

RECEIVED 2 February 2015  
 REVISED 18 May 2015  
 ACCEPTED 31 May 2015  
 PUBLISHED ONLINE FIRST 7 August 2015

Sean Barnes<sup>1</sup>, Eric Hamrock<sup>2</sup>, Matthew Toerper<sup>3</sup>, Sauleh Siddiqui<sup>4</sup>, Scott Levin<sup>5</sup>



## ABSTRACT

**Objective** Hospitals are challenged to provide timely patient care while maintaining high resource utilization. This has prompted hospital initiatives to increase patient flow and minimize nonvalue added care time. Real-time demand capacity management (RTDC) is one such initiative whereby clinicians convene each morning to predict patients able to leave the same day and prioritize their remaining tasks for early discharge. Our objective is to automate and improve these discharge predictions by applying supervised machine learning methods to readily available health information.

**Materials and Methods** The authors use supervised machine learning methods to predict patients' likelihood of discharge by 2 p.m. and by midnight each day for an inpatient medical unit. Using data collected over 8000 patient stays and 20 000 patient days, the predictive performance of the model is compared to clinicians using sensitivity, specificity, Youden's Index (i.e., sensitivity + specificity – 1), and aggregate accuracy measures.

**Results** The model compared to clinician predictions demonstrated significantly higher sensitivity ( $P < .01$ ), lower specificity ( $P < .01$ ), and a comparable Youden Index ( $P > .10$ ). Early discharges were less predictable than midnight discharges. The model was more accurate than clinicians in predicting the total number of daily discharges and capable of ranking patients closest to future discharge.

**Conclusions** There is potential to use readily available health information to predict daily patient discharges with accuracies comparable to clinician predictions. This approach may be used to automate and support daily RTDC predictions aimed at improving patient flow.

**Keywords:** length of stay, patient flow, machine learning, operational forecasting

## BACKGROUND AND SIGNIFICANCE

Hospitals and clinician objectives require balance in treating each patient's condition effectively while efficiently distributing healthcare resources to patient populations over time.<sup>1</sup> A key determinant of hospital capacity and resource management is linked to patient flow, a common indicator of patient safety, satisfaction, and access.<sup>2–5</sup> Optimal patient flow facilitates beneficial treatment, minimal waiting, minimal exposure to risks associated with hospitalization, and efficient use of resources (e.g., of beds, clinical staff, and medical equipment). Patient flow is also a determinant of access to specialized inpatient services. These patients request admission from external sources (e.g., other health services) and internal sources such as the emergency department (ED), procedural areas, or peri-anesthesia care units (PACUs).<sup>6,7</sup>

Evidence supporting the benefits of improved flow has been mounting.<sup>8</sup> For patients, prolonged hospital stays increase the risk of adverse events, such as hospital-acquired infections, adverse drug events, poor nutritional levels, and other complications.<sup>9–12</sup> Extended stays have also been associated with poor patient satisfaction.<sup>11,13–17</sup> For hospitals, economic pressures to deliver efficient and accessible care are at unprecedented highs. Healthcare costs as a percentage of gross domestic product (GDP) (17.9% in 2012) have been rising faster than anticipated,<sup>18</sup> and approximately 30–40% of these expenditures have been attributed to “overuse, underuse, misuse, duplication, system failures, unnecessary repetition, poor communication, and inefficiency.”<sup>19</sup> These factors impede patient flow, prolong patient stays, and increase the cost of care per patient.<sup>14,20–22</sup>

Patient flow, and by association, bed and capacity management, is a common focus area for operations management methods applied to healthcare. Discrete-event simulation, optimization, and Lean Six Sigma approaches have been applied successfully in various settings to improve patient flow by either redesigning care delivery processes or more efficiently matching staff and other resources (e.g., operating rooms, medical equipment) to demand.<sup>23–35</sup> These patient flow evaluations are data-driven and inform long-term operational decision-making. Outcomes of these studies include improved patient and staff scheduling strategies, new bed management policies, and reduction in care process variability.

A more recent advancement in patient flow management alternatively focuses on short-term operational decisions. Real-time demand capacity management (RTDC) is a new method developed by the Institute for Healthcare Improvement that has shown promising but variable results when pilot tested in hospitals. The RTDC process involves 4 steps: 1) predicting capacity, 2) predicting demand, 3) developing a plan, and 4) evaluating the plan.<sup>36</sup> The RTDC process centers around a morning clinician huddle to predict which and how many patients will be discharged that same day. Given daily predictions for demand, the group then attempts to match their supply of beds to the demand from new patients (i.e., admissions) by prioritizing current patients able to be discharged. RTDC implementation signifies a culture change as the hospital staff dedicates time each morning to coordinate and focus on patient flow. The developers of RTDC have demonstrated that this approach, after an initial learning period, may reduce

Correspondence to Department of Decision, Operations & Information Technologies, Robert H. Smith School of Business, 4352 Van Munching Hall, University of Maryland, College Park, MD 20742, USA; sbarnes@rhsmith.umd.edu; Tel: (301) 405-9679

© The Author 2015. Published by Oxford University Press on behalf of the American Medical Informatics Association.

All rights reserved. For Permissions, please email: journals.permissions@oup.com

For numbered affiliations see end of article.

key indicators of patient flow across hospitals. These indicators include ED boarding time, patients who left without being seen, and overnight holds in the PACU.

RTDC does have limitations that we strive to address in this work. First, RTDC requires a daily clinician huddle that demands dedicated time that may be expedited or even eliminated through automated prediction. Second, the prediction process is subjective and thus susceptible to high variability. In addition, the most successful implementations of RTDC have occurred in surgical units, where patient conditions and clinical pathways are commonly well prescribed.<sup>37,38</sup> Thus, previous results on RTDC predictive performance for surgical patients may not be generalized to patients in medical units where clinical pathways are less defined and length of stay (LOS) is more variable.

Our objective for this study is to support automation of the RTDC process by producing daily predictions of patient discharge times for a single inpatient medical unit. We hypothesize that a predictive model using readily available health information may perform comparable or better than predictions made by clinicians during their daily huddle. We apply supervised machine learning methods to predict the probability of patient discharge by 2 p.m. and by the end of each day (i.e., midnight). By discharging a patient earlier in the day (e.g., prior to 2 p.m.), there is an increased likelihood of admitting a new patient to the same bed, thereby facilitating increased bed utilization and patient flow.<sup>36</sup>

### Operational Forecasting Literature

Forecasting models have been leveraged to predict hospital occupancy, patient arrivals and discharges, and other unit-specific operational metrics. These models have been derived from variations of autoregressive moving average approaches,<sup>39–44</sup> exponential smoothing,<sup>39,40,42</sup> Poisson regression,<sup>45</sup> neural networks,<sup>42</sup> and discrete-event simulation methods.<sup>46–48</sup> More recent studies have employed logistic regression<sup>49</sup> and survival analysis<sup>50</sup> methods to predict patient LOS. Through accurate prediction of patients' LOS, clinical staff may efficiently schedule future appointments or admissions to avoid backlog or denials.<sup>44,51–53</sup> Our approach builds on this research by focusing on short-term predictions (i.e., daily discharges) for patients.

### Supervised Machine Learning Literature

Tree-based supervised machine learning algorithms have been applied in healthcare to 1) predict a continuous-valued outcome or 2) classify patients into one or more clinical subgroups. Example applications for the former include using linear regression or regression trees to predict the range of motion for orthopedic patients,<sup>54</sup> costs,<sup>55,56</sup> and utilization.<sup>57</sup> An example of the latter typically involves logistic regression or classification trees, and has been used to differentiate between benign and malignant tumors,<sup>58</sup> identify patients most likely to benefit from screening procedures,<sup>59,60</sup> identify high risk patients,<sup>61</sup> and predict specific clinical outcomes.<sup>62–64</sup> Recently, more advanced machine learning techniques such as bagging, boosting, random forests, and support vector machines have been applied to healthcare problems such as classifying heart failure patients<sup>65</sup> and predicting healthcare costs.<sup>66</sup> We build upon this research by adding a new application area for these methods and leveraging the benefits of ensemble learning techniques such as bagging and boosting.

## MATERIALS AND METHODS

### Setting and Data Sources

The study was conducted within a single, 36-bed medical unit in a large, mid-Atlantic academic medical center serving an urban population. The unit is staffed with hospitalist physicians with no teaching responsibilities. Patient flow data (i.e., admission and discharge times),

demographics, and basic admission diagnoses data was collected for 9636 patient visits over a 34-month study period from January 1, 2011 to November 1, 2013. After excluding incomplete and erroneous records, we retained data for 8852 patient visits and converted these visits to  $N=20\,243$  individual patient days. These data are summarized in Table 1.

These data are standardized and readily accessible in most hospital information systems, thus reproducible in other hospitals. The demographic and clinical predictors are static model inputs that are known at the time of admission and do not change during a patient's stay. Other predictors such as patient census, day of the week, and elapsed length of stay are dynamic and are continuously updated during a patient's stay. Elapsed length of stay, age, and patient census are numerical variables, whereas the remaining predictors are modeled using binary indicator variables (i.e., 0 or 1 indicating the absence or presence of a specific category).

The reason for visit is determined according to the International Classification of Diseases-9 diagnoses structure.<sup>67</sup> A physician identifies this condition (via an electronic pick list) at the time of admission to the unit, thereby documenting the primary reason for hospitalization. Observation status is assigned to patients who are cared for within the unit and must be monitored and evaluated before they are eligible for safe discharge.<sup>68</sup> Administratively, hospitals may use this designation to bill Medicare for the patient under the outpatient service category. However, a large study by the Department of Health and Human Services found that observation patients have the same health conditions as those who are fully admitted.<sup>69</sup>

In addition to the patient data listed in Table 1, a novel aspect of this study is our collection of clinician predictions of patients to be discharged at 2 p.m. and midnight. We recorded predictions from daily morning huddles for 8 overlapping months between March 18, 2013 and November 1, 2013. The prediction team was comprised of a charge nurse, case manager, and physicians. Members of this team were directly involved with the administrative or clinical aspects of a subset of patients in the unit each day. These team members had access to substantially more information for their predictions than what was accessible to our prediction models. This unique data facilitated a comparative study between the machine-learning and clinician predictions. In order to compare performance directly to the RTDC process, the machine-learning models were designed to produce predictions based on data available at 7 a.m. each day, which is analogous to the clinician huddle times.

### Analytic Methods

We applied and evaluated several supervised machine learning algorithms to predict patient discharge and compare with clinician predictions. These algorithms are considered 'supervised' because they are fit to labeled training data (i.e., the outcomes are known in retrospect) and then independently evaluated on separate labeled test data to estimate their performance in practice. For each patient day, these algorithms produced two predictions representing the probability of a patient being discharged by either 2 p.m. or midnight that same day. These predictions were generated for every patient in the unit for each day of their stay, which simulated the clinician-based RTDC prediction process. This model is designed to support real-time predictions of expected bed capacity and could be adapted to predict patient discharges for any time interval (e.g., hourly instead of specifically 2 p.m. and the end of the day).

Systematic experiments applying common supervised machine-learning methods to predict individual patient discharges were performed. These methods included logistic regression (i.e., reference method), classification and regression trees, and tree-based ensemble

Table 1: Patient flow and prediction data.

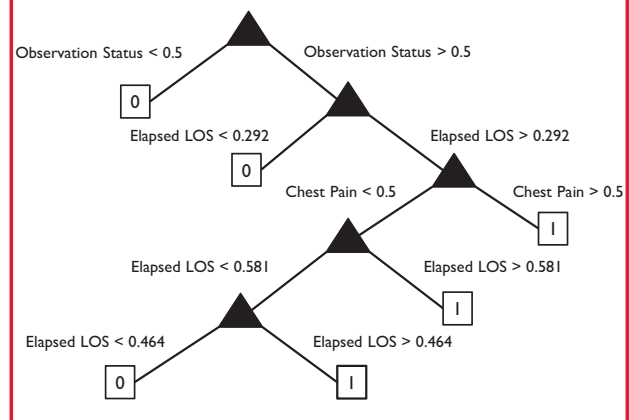
Patient Flow (Outcomes)	
Length of stay: mean, median (IQR)	Mean: 52 h, median: 37 h (42.40 h)
Discharge timing	27.4% of discharges by 2 p.m., 88.8% by 7 p.m.
Demographic predictors	
Gender	57.1% female
Ethnicity	68.9% Caucasian, 26.7% African American, 4.4% Other
Age: mean, median (IQR)	58.4 years, 57.0 years (24 years)
Insurance	45.3% Private, 39.5% Medicare, 9.4% Uninsured, 5.9% Medicaid
Clinical predictors	
Reason for visit	Chest pain 39.4%, syncope 11.4%, abdominal pain 3.8%, chronic obstructive pulmonary disorder (COPD) 3.6%, congestive heart failure 1.8%
Observation status	52.5% of patients
Pending discharge location	Home 85%, other healthcare facilities: 15%
Unit workload predictors	
Patient census: mean, median (IQR)	20.87, 21 (7)
Timing predictors	
Day of week	Monday 14.7%, Tuesday 15.4%, Wednesday 16.0%, Thursday 16.2%, Friday 16.8%, Saturday 11.4%, Sunday 9.5%
Elapsed length of stay	Changes dynamically

learning methods. Model parameters and thresholds were tuned for optimal performance with respect to the predictive measures described in the following section.<sup>70</sup> We compared the results among these predictive methods, and then selected the best method to compare to the clinician predictions.

Logistic regression is a commonly used classification method for clinical applications, and it is an effective approach for providing a baseline for how effectively the available data can be leveraged to predict the primary outcome. However, logistic regression is sometimes limited for predictive applications, especially with large, highly dimensional data where interactions may exist.<sup>70</sup> Despite this, logistic regression is a well-established method that will serve as a baseline for comparison with more robust supervised machine learning algorithms.

Next, we applied tree-based methods to predict patient discharge outcomes, which iteratively partition the data into groups with similar characteristics and outcomes.<sup>71–73</sup> These methods are adept at recognizing predictor variable interactions and identifying sub-cohorts of patients that are more likely to have a positive outcome.<sup>71</sup> Tree-based models have the distinct advantage of providing more practical insight than regression-based methods. The highest levels of these trees can be translated into effective decision rules that may be interpreted

Figure 1: Classification tree for end-of-day discharge predictions.



by clinicians, and embedded within existing clinical information technology infrastructure to facilitate implementation and uptake in practice.<sup>68</sup>

Figure 1 shows a visualization of the classification tree for end-of-day discharge predictions, which provides an example of how our results could be shared with clinicians. In this figure, we observe that observation status is the most critical predictor relative to predicting the outcome. Patients who are not on observation status (i.e., an intermittent care status to evaluate whether they need to be admitted) tend to stay, regardless of any of their other conditions. Patients who are on observation status are more likely to be discharged only when their elapsed LOS exceeds approximately 12 h or if they reported chest pain as their chief complaint.

Despite these benefits, simple tree-based learning methods have the potential to over-fit training data – causing the in-sample performance to exceed the out-of-sample predictive performance, which occurs when the model is overly complex and mistakes noise for key underlying relationships.<sup>72,74</sup> Tree-based ensemble learning methods, such as bagging and boosting, address over-fitting by training large numbers of “weak learners” and leveraging the diversity across learners (i.e., individual trees) to produce stable out-of-sample predictions. Diverse trees are aggregated in some form (e.g., by selecting majority votes or averaging) to classify or estimate risk for an outcome.<sup>72,74–76</sup> Bagging utilizes a bootstrapping process (i.e., sampling with replacement) to train numerous trees to produce model-averaged predictions.<sup>74</sup> The random forest is a common bagging method that increases diversity across each individual tree.<sup>72,75,76</sup> The random forest approach also facilitates evaluation of weak learners by selecting a random subset of predictors at each candidate split in the learning process for predictor variable importance in classifying outcomes. Boosting implements optimization algorithms to re-weight misclassified observations in an attempt to improve overall prediction accuracy.<sup>72</sup>

### Training and Evaluation

Our ultimate goal was to implement the most parsimonious model with high accuracy that utilizes data that can be automatically extracted from electronic medical record systems. Predictor variables (see Table 1) for each patient day (collected at 7 a.m.) were linked to binary outcomes indicating whether the patient was discharged by 2 p.m. or the end of day, respectively. We trained each model on data collected over 26.6 months from the start of the study (January 1, 2011) until the date when the clinician predictions began (18 March, 2013). Model

predictions were generated for the following 9.4 months (19 March 2013 to 31 December, 2013) out-of-sample; purposefully overlapping with the same period clinician predictions were collected. This facilitated cross-validation (78% training set and 22% testing set) of the model and direct comparison to the clinician predictions. Predictive performance was then estimated for model-based and clinician predictions using binary classification statistics. Measures from the confusion matrix (i.e., true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN)) were used to calculate:

- Sensitivity =  $TP / (TP + FN)$
- Specificity =  $TN / (TN + FP)$
- Youden's Index  $J = \text{Sensitivity} + \text{Specificity} - 1$

These measures, along with the number of positive predictions (i.e., discharges), captured different aspects of model performance. Youden's Index<sup>77</sup> is a global accuracy measure whereas sensitivity and specificity capture how well the predictions perform with respect to a specific outcome (i.e., discharge or stay). Youden's Index is a common metric used to evaluate the performance of diagnostic tests,<sup>78–80</sup> and similar to these studies, we used it to optimize the cut-off thresholds (i.e., to determine a discharge or stay prediction) for each algorithm. We also calculated these metrics for near future outcomes (i.e., outcomes for the next time period). For example, how many patients predicted to be discharged by 2 p.m. were discharged by the end of the day? Similarly, how many patients predicted to be discharged by the end of the day were discharged by the end of the next day? These measures include correct predictions but also incorrect predictions that were correct within one day, which provides some insight as to the magnitude of predictive errors.

We conducted hypothesis testing ( $\alpha = 0.05$ ) to detect statistically significant differences between the predictive capabilities for the model and the clinicians. Each of these tests can be applied readily to paired samples, which is required for comparing predictions for the same patient day(s). We used McNemar's test with Yates correction for continuity<sup>81–82</sup> to test the null hypothesis  $H_0: p_{\text{model}} = p_{\text{clinicians}}$ , where  $p$  represents the proportion of correct predictions relative to the total number of positive (sensitivity) or negative (specificity) outcomes. For Youden's Index  $J$ , we apply the method devised in Chen et al.<sup>83</sup> for paired samples to test the null hypothesis  $H_0: J_{\text{model}} = J_{\text{clinicians}}$ . For both tests, we used two-sided alternative hypotheses.

In addition, we evaluated the results in two additional measures that are relevant to model implementation. First, we aggregated (i.e., summed) the number of expected daily discharges from the model and clinician predictions, and then compared these results to the actual number of discharged patients for each day using paired hypothesis tests at the  $\alpha = 0.05$  level. This aggregate measure is useful for proactive patient flow management by anticipating available (i.e., unoccupied) beds in advance and facilitating accept/reject admission decisions related to capacity constraints. Separately, we used model predictions to rank the patients in order of their likelihood for discharge each day. Spearman rank correlation coefficients were computed to evaluate the accuracy of these rankings with observed patients remaining length of stay.<sup>84</sup> This latter approach provided some indication as to whether this model could be used to accurately prioritize patients for discharge.

## RESULTS

We applied tree-based supervised machine learning methods to predict discharge by 2 p.m. and the end of the day for each patient day.

The regression random forest (RRF) method proved most accurate from systematic experiments of several tree-based algorithms with parameters tuned for optimal predictive performance. Figure 2, summarizing variable importance from the training of the RRF models, provides evidence that elapsed LOS and observation status are substantially more important than all of the other predictors, followed by Sunday, chest pain, disposition, age, and syncope. Predictors such as gender, ethnicity, weekdays (Monday through Thursday), and less common reasons for visit (e.g., abdominal pain, COPD, congestive heart failure) had little to no predictive power for this patient population.

### Individual Predictions

We compared the RRF model predictions to the clinician predictions in order to evaluate the potential of this approach in practice. We have also included comparisons to logistic regression as a reference to better understand the efficacy of the ensemble learning approach. Results for both 2 p.m. and end-of-day outcomes are summarized in Table 2 for the specific days when clinician predictions were available (19 March, 2013 to 31 December, 2013), which included 4833 patient days.

Overall, the logistic regression and RRF model were more aggressive in predicting discharge than the clinicians for both the 2 p.m. and end-of-day outcomes. This resulted in the automated models predicting discharges with higher sensitivity ( $P < .01$ ) and lower specificity ( $P < .01$ ) compared to the clinicians. Thus the models predicted a higher proportion of discharges, but at the cost of producing more false positives. The logistic regression baseline model consistently demonstrated more sensitive behavior than the RRF model.

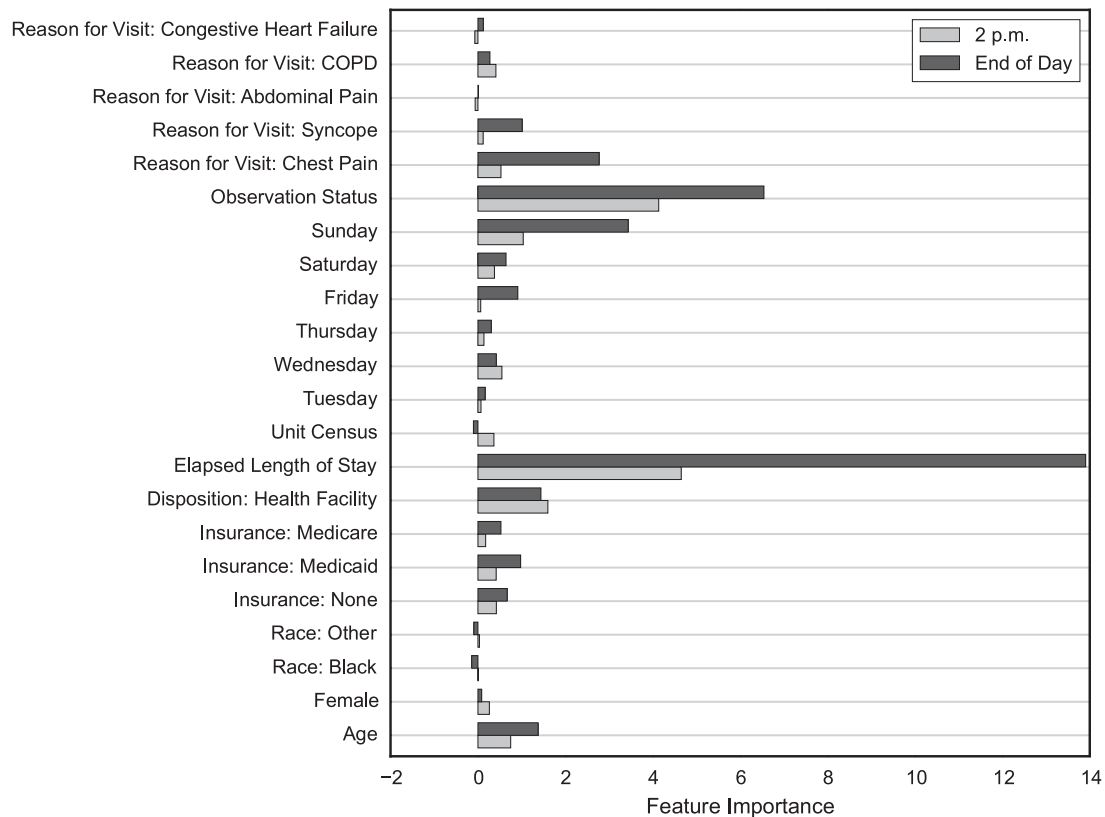
Despite differences in sensitivity and specificity, the RRF model and clinicians predictive performance were comparable for Youden's Index, our global accuracy measure. There were no significant differences for same-day predictions (2 p.m. RRF model: 26.0%, clinicians 26.5%,  $P = .81$ ; end-of-day RRF model: 34.0%, clinicians 34.0%,  $P = .84$ ) or the near-future end-of-day outcome (RRF model: 26.6%, clinicians 25.2%,  $P = .62$ ). However, the RRF model did perform significantly better for the near-future 2 p.m. outcome (RRF model: 31.3%, clinicians 19.0%,  $P < .01$ ). Predictions across models and clinicians were more accurate for the end of the day than for 2 p.m.

### Daily Aggregate Predictions

We also compared the aggregated RRF model and clinician predictions for the total number of discharges each day. For the model predictions, we summed the predicted probabilities across all patients to calculate the expected value of discharged patients. The results for 2 p.m. and the end of day for the overlapping date range are summarized in Figures 3 and 4, which show that the model outperforms clinician estimates of the average number of patients to be discharged early and by the end of the day.

In Figure 3, we plot the distribution of the actual number of discharges per day against the distributions for the RRF model and clinician predictions for each outcome. The actual number of discharges per day averaged 2.36 by 2 p.m. and 8.29 patients by the end of the day. Discharges predicted each day by the RRF model were 2.45 ( $P = .37$ ) and 8.51 ( $P = .13$ ), respectively, demonstrating no statistically significant difference from the actual number of discharges. In comparison, clinician predictions were accurate for 2 p.m. discharges (2.16,  $P = .19$ ), but significantly deviated from the actual number of discharges for the end of the day (6.54,  $P < .01$ ). The model predicted the total number of discharges within 2 patients 82% of days for 2 p.m. and 105 63% of days for the end of the day, whereas the

**Figure 2:** Variable importance summary from training regression random forest models for the 2 p.m. and end-of-day discharge predictions.



clinicians predicted the same for 52% and 32% of days, respectively (see Figure 4). In addition, large residuals (i.e., errors of 5 discharges or more) were much more frequent for the clinician predictions.

### Daily Rank Predictions

In practice, our prediction model could be used to rank patients daily—based on their likelihood of being discharged—in order to prioritize the remaining tasks for the most likely patients. We computed Spearman's rank correlation between the remaining LOS for each patient and their respective RRF model predicted probabilities for each day (trained on the full data set). Figure 5 shows a histogram of these correlations (in intervals of 0.05) for the 2 p.m. and end-of-day models. The mean rank correlation for 2 p.m. was 0.4816 and for the end of the day was 0.4489. Both plots show that the correlations are almost exclusively positive and moderately large, which suggests that the rank of the model predictions was moderately correlated with the actual discharge order. For some days, the model nearly predicted the exact order in which the patients were discharged.

## DISCUSSION

Improving patient flow continues to be a top priority in the acute-care setting, where patients with longer lengths of stay are less satisfied and exposed to the risk of adverse events (e.g., hospital-acquired infections, complications). Hospitals have aligned incentives to improve patient flow because of the rising demand for services and the economic pressures to reduce costs and improve resource utilization. Our

approach is designed to empower clinicians and hospital administrators with analytical tools to increase their collective efficiency.

Our approach could be operationalized in three distinct ways. First, it could be used to identify individual patients who are most likely to be discharged on a given day. Hospital staff could prioritize these patients in order to discharge them as early as possible—without negatively affecting their care—so that other patients can be admitted in their place. Supervised machine learning methods may be used to rank patients concurrently in a hospital (or specific unit) according to their discharge probabilities. We have shown that our models perform well for prediction or ranking. Alternatively, patients who are most likely to be discharged may not be impacted significantly by an effort to prioritize their remaining tasks. Therefore, another approach could be to identify a second tier of patients who are only mildly to moderately likely to be discharged. Prioritizing the remaining tasks for these patients may have a more significant global impact on the number of patients discharged over a given time period. Setting a range of predicted probabilities immediately below the classification threshold would identify these patients. The last approach is to aggregate the predicted discharge probabilities into daily discharge predictions. This approach does not facilitate prioritization for patients likely to be discharged, but it supports bed capacity planning for the unit. We have shown that these types of predictions are very accurate and would provide the staff with a good idea of how many in-use beds are likely to be available by a specific time of day.

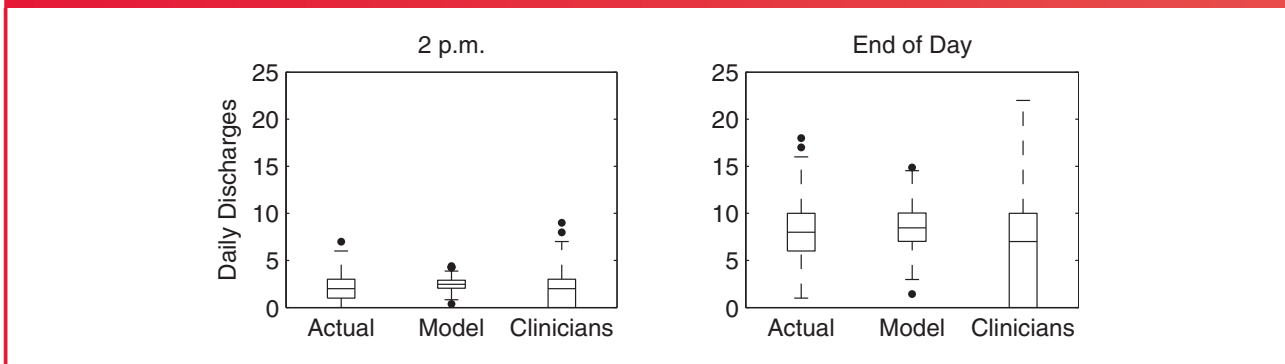
There are three important limitations to highlight for this study. First, the performance of the models may improve with a

Table 2: Performance comparison summary between logistic regression, regression random forest, and clinician predictions performance comparison.

Performance Measure	2 p.m.			End of Day		
	Logistic Regression	Regression Random Forest	Clinicians	Logistic Regression	Regression Random Forest	Clinicians
Positive Predictions (Discharge)	1749	1781	485	1870	2175	1471
Sensitivity (%)	65.9	60.0	33.6	71.5	66.1	51.3
	$P < .01$	$P < .01$		$P = .02$	$P < .01$	
Specificity (%)	52.8	66.0	92.9	54.9	68.3	82.7
	$P < .01$	$P < .01$		$P < .01$	$P < .01$	
Youden's Index (%)	18.7	26.0	26.5	26.4	34.0	34.0
	$P < .01$	$P = .81$		$P < .01$	$P = .84$	
Near Future Sensitivity (%)	73.4	56.1	21.7	81.9	54.1	39.0
	$P < .01$	$P < .01$		$P < .01$	$P < .01$	
Near Future Specificity (%)	51.5	75.2	97.3	49.3	72.5	86.1
	$P < .01$	$P < .01$		$P < .01$	$P < .01$	
Near Future Youden Index (%)	24.9	31.3	19.0	31.2	26.6	25.2
	$P = .06$	$P < .01$		$P = .02$	$P = .62$	

The  $P$ -value listed beneath each performance measure represents the results of McNemar's test for sensitivity and specificity or the method described in<sup>83</sup> for Youden's index for the estimated difference between each model and the clinicians. In each case, the null hypothesis of equality is tested against the two-sided alternative.

Figure 3: Comparison of actual, model prediction, and clinician prediction distributions of the average number of patients discharged from the unit each day.



larger data set, in terms of the number of patient days used to train the models and also with respect to being collected from multiple sites. Increasing the size of the training data set would improve our ability to detect more complex patterns in patient length of stay. Similarly, the patterns detected in our training data set may not be generalizable to other hospitals or units. Our intention is for each hospital and/or unit to replicate our prediction model and train it on its own data. Finally, we compared the performance of continuous (i.e., probability-based) predictions from our models to binary predictions (i.e., exit or stay) from the clinicians. Ideally, the fairest comparison would be between continuous predictions from both methods, however these predictions would be difficult to generate and collect in practice.

### CONCLUSIONS

We applied supervised machine learning algorithms to readily available health information to predict daily discharge outcomes as part of the RTDC process. We directly compared model predictions to clinician predictions using several performance metrics. The model predicted discharges with higher sensitivity and lower specificity compared to the clinicians, and the two methods were comparable (i.e., not statistically significantly different) for our global accuracy measure (Youden's Index). However, the model did outperform the clinicians for some near-future and aggregate prediction metrics. Thus there is high potential for these models to automate and expedite the RTDC prediction process, thereby eliminating the need for daily clinician huddles or supporting more accurate clinician predictions. Furthermore, these

Figure 4: Histogram summary of differences between model, clinician, and actual daily discharges.

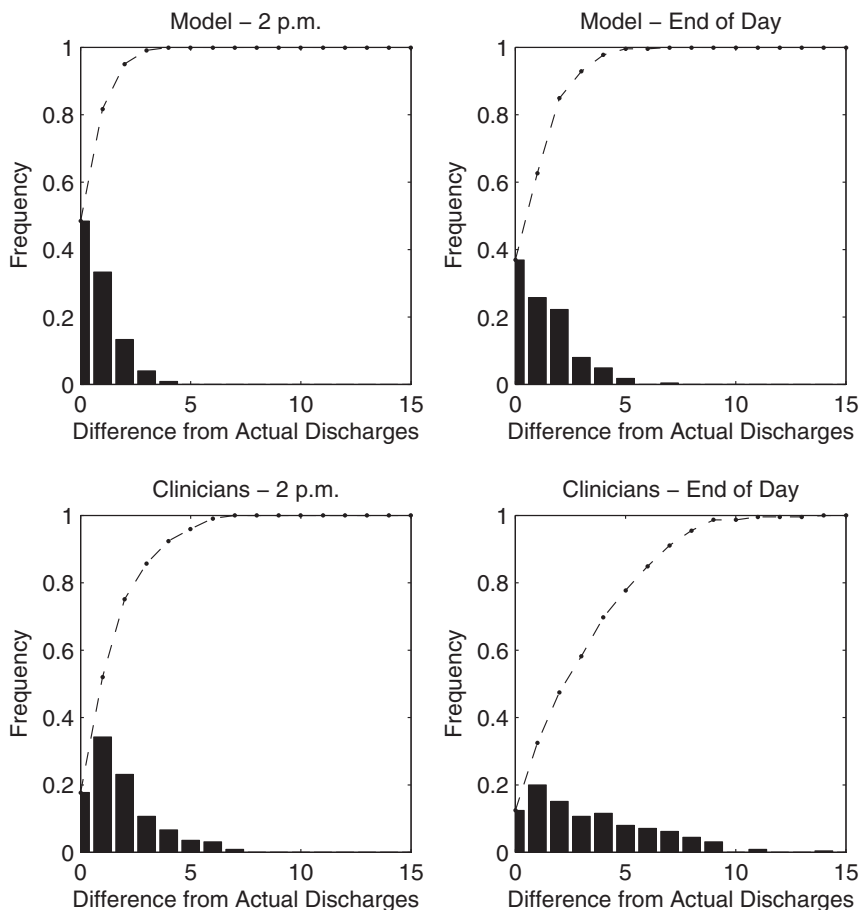
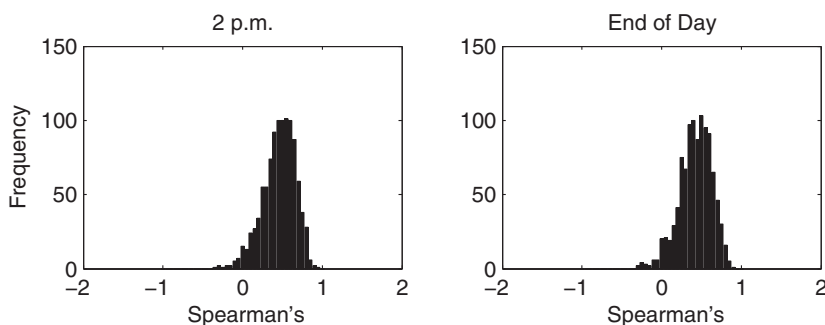


Figure 5: Histogram summary of daily Spearman's rank correlation between remaining LOS and model prediction scores for 2 p.m. and the end of the day.



models were applied to simple and readily accessible information that can be quite powerful, and easily replicated across acute-care environments using electronic information systems.

**FUNDING**

This work was supported in part by the National Science Foundation grant number 0927207.

**CONTRIBUTORS**

S.B., S.L., and E.H. made significant contributions to the conception and design of the work. M.T., E.H., and S.L. made significant contributions to the acquisition and processing of the data used for the analysis. S.B., S.L., and S.S. made significant contributions to the analysis and interpretation of the data for the work. All authors have contributed to either drafting or revising the article, have approved the version





46. Hoot NR, LeBlanc LJ, Jones I, et al. Forecasting emergency department crowding: a discrete event simulation. *Ann Emerg Med.* 2008;52(2):116–125.
47. Hoot NR, LeBlanc LJ, Jones I, et al. Forecasting emergency department crowding: a prospective, real-time evaluation. *JAMIA.* 2009;16(3):338–345.
48. Hoot NR, Zhou C, Jones I, et al. Measuring and forecasting emergency department crowding in real time. *Ann Emerg Med.* 2007;49(6):747–755.
49. Levin SR, Harley ET, Fackler JC, et al. Real-time forecasting of pediatric intensive care unit length of stay using computerized provider orders. *Crit Care Med.* 2012;40(11):3058–3064.
50. Clark DE, Ryan LM. Concurrent prediction of hospital mortality and length of stay from risk factors on admission. *Health Services Res.* 2002;37(3):631–645.
51. Peck JS, Benneyan JC, Nightingale DJ, et al. Predicting emergency department inpatient admissions to improve same-day patient flow. *Acad Emerg Med.* 2012;19(9):E1045–1054.
52. Sun Y, Heng BH, Tay SY, et al. Predicting hospital admissions at emergency department triage using routine administrative data. *Acad Emerg Med.* 2011;18(8):844–850.
53. Peck JS, Gaehde SA, Nightingale DJ, et al. Generalizability of a simple approach for predicting hospital admission from an emergency department. *Acad Emerg Med.* 2013;20(11):1156–1163.
54. Ritter MA, Harty LD, Davis KE, et al. Predicting range of motion after total knee arthroplasty clustering, log-linear regression, and regression tree analysis. *J Bone Jt Surg.* 2003;85(7):1278–1285.
55. Raebel MA, Malone DC, Conner DA, et al. Health services use and health care costs of obese and nonobese individuals. *Arch Int Med.* 2004;164(19):2135–2140.
56. Gregori D, Petrinco M, Bo S, et al. Regression models for analyzing costs and their determinants in health care: an introductory review. *Intern J Qual Health Care.* 2011;23(3):331–341.
57. de Boer AG, Wijker W, de Haes HC. Predictors of health care utilization in the chronically ill: a review of the literature. *Health Policy.* 1997;42(2):101–115.
58. Sauerbrei W, Madjar H, Prompeler H. Differentiation of benign and malignant breast tumors by logistic regression and a classification tree using Doppler flow signals. *Methods Inform Med.* 1998;37(3):226–234.
59. Barriga KJ, Hamman RF, Hoag S, et al. Population screening for glucose intolerant subjects using decision tree analyses. *Diabet Res Clin Pract.* 1996;34 (Suppl.):S17–S29.
60. McGrath JS, Ponich TP, Gregor JC. Screening for colorectal cancer: the cost to find an advanced adenoma. *Am J Gastroenterol.* 2002;97:2902–2907.
61. Goldman L, Cook EF, Johnson PA, et al. Prediction of the need for intensive care in patients who come to the emergency departments with acute chest pain. *New Engl J Med.* 1996;334:1498–1504.
62. Falconer JA, Naughton BJ, Dunlop DD, et al. Predicting stroke inpatient rehabilitation outcome using a classification tree approach. *Arch Phys Med Rehabil.* 1994;75(6):619–625.
63. Germanson T, Lanzino G, Kassell NF. CART for prediction of function after head trauma. *J Neurosurg.* 1995;83:941–942.
64. Temkin NR, Holubkov R, Machamer JE, Winn HR, Dikmen SS, et al. Classification and regression trees (CART) for prediction of function at 1 year following head trauma. *J Neurosurg.* 1995;82(5):764–771.
65. Austin P, Tu J, Levy D. Using methods from the data-mining and machine-learning literature for disease classification and prediction: a case study examining classification of heart failure subtypes. *J Clin Epidemiol.* 2013;66:398–407.
66. Robinson JW. Regression tree boosting to adjust health care cost predictions for diagnostic mix. *Health Services Res.* 2008;43(2):755–772.
67. Centers for Medicare and Medicaid Services. *Medicare Claims Processing Manual: Chapter 23 – Fee Schedule Administration and Coding Requirements.* 2015. <http://www.cms.gov/Regulations-and-Guidance/Guidance/Manuals/downloads/clm104c23.pdf> Accessed April 8, 2015.
68. Centers for Medicare and Medicaid Services. *Are You a Hospital Inpatient or Outpatient? CMS Product No. 11435.* 2014. <https://www.medicare.gov/Pubs/pdf/11435.pdf> Accessed April 8, 2015.
69. Wright, S. *Memorandum Report: Hospitals' Use of Observation Status and Short Inpatient Stays for Medicare Beneficiaries, OEI-02-12-00040.* 2013. <https://kaiserhealthnews.files.wordpress.com/2013/07/oei-02-12-00040.pdf> Accessed April 8, 2015.
70. Perlich C, Provost F, Simonoff JS. Tree induction vs. logistic regression: a learning-curve analysis. *J Mach Learn Res.* 2003;4:211–255.
71. Breiman L, Friedman J, Olshen R, et al. *Classification and Regression Trees, 1st edn.* Monterey, CA: Chapman and Hall; 1984.
72. Hastie T, Tibshirani R, Friedman J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction.* New York, NY: Springer; 2001.
73. Barner M. *Principles of Data Mining.* London: Springer; 2007.
74. Breiman L. Bagging predictors. *Mach Learn.* 1996;24(2):123–140.
75. Breiman L. Random forests. *Mach Learn.* 2001;45(1):5–32.
76. Schapire R. The strength of weak learnability. *Mach Learn.* 1990;5(2):197–227.
77. Youden WJ. Index for rating diagnostic tests. *Cancer.* 1950;3:32–35.
78. Greiner M, Sohr D, Göbel P. A modified ROC analysis for the selection of cut-off values and the definition of intermediate results of serodiagnostic tests. *J Immunol Methods.* 1995;185(1):123–132.
79. Fluss R, Faraggi D, Reiser B. Estimation of the Youden Index and its associated cutoff point. *Biometrical J.* 2005;47(4):458–472.
80. Schisterman EF, Perkins NJ, Liu A, Bondell H. Optimal cut-point and its corresponding Youden Index to discriminate individuals using pooled blood samples. *Epidemiology.* 2005;16(1):73–81.
81. McNemar Q. Note on the sampling error of the difference between correlated proportions or percentages. *Psychometrika.* 1947;12(2):153–157.
82. Yates F. Contingency tables. *J Royal Stat Soc.* 1934;1:217–235.
83. Chen F, Xue Y, Tan MT, Chen P. Efficient statistical tests to compare Youden index: accounting for contingency correlation. *Stat Med.* 2015;34:1560–1576.
84. Spearman C. The proof and measurement of association between two things. *Am J Psychol.* 1904;15:72–101.

## AUTHOR AFFILIATIONS

<sup>1</sup>Department of Decision, Operations & Information Technologies, Robert H. Smith School of Business, 4352 Van Munching Hall, University of Maryland, College Park, MD 20742, USA

<sup>2</sup>Department of Operations Integration, Johns Hopkins Health System, Baltimore, MD, USA

<sup>3</sup>Department of Emergency Medicine, Johns Hopkins Hospital, Baltimore, MD, USA

<sup>4</sup>Departments of Civil Engineering and Applied Mathematics & Statistics, Johns Hopkins Systems Institute, Johns Hopkins University, Baltimore, MD, USA

<sup>5</sup>Department of Emergency Medicine and Civil Engineering, Johns Hopkins Systems Institute, Johns Hopkins University, Baltimore, MD, USA