

Patient-Specific Readmission Prediction and Intervention for Health Care

Yan Zhang

*Microsoft, New England Research Center,
1 Memorial Dr., Cambridge, MA, 02142, USA
zhangya@microsoft.com*

ABSTRACT

Hospital readmission is often associated with unfavorable patient outcomes and a large cost of resources. Therefore, preventing avoidable re-hospitalizations is imperative. To target this problem, one important metric that researchers and practitioners strive to reduce is the 30-day hospital readmission rate. In this paper, we introduce a general decision support system that utilizes machine learning (ML) based patient-specific prediction to guide the suggestion of patient intervention program assignment, with the objective of minimizing the readmission cost for hospitals. This work has three major contributions. First, the proposed solution is highly scalable by using PySpark. Second, we outline solution architecture components including (1) data injection (both real-time sensor reading and data at rest), processing, and analysis; and (2) ML model building, evaluation, deployment and scoring. Third, we discuss how the ML prediction results can be taken into account in a decision support system by presenting a rich visualization.

1. INTRODUCTION

The U.S. health care cost is unarguably top ranked comparing to other developed countries, one major cause of which is inefficiency in the health-care system. Inefficiencies come from poorly management of care experience, insights from research, and available evidence (Rumsfeld, Joynt, & Maddox, 2016). More specifically, the opportunities are missed and the resources are wasted during the process of data collection (not capturing previous care experience), data analysis (not capable of analyzing the collected data), and taking actions (not effective of making usage of the analysis result).

Among many areas where efficiency can be improved, hospital readmission reduction is one of them. In health care domain, a hospital readmission is defined as an admission to

a hospital within a certain time frame following the discharge from the original hospital stay. Hospital readmission is often associated with unfavorable patient outcomes together with a large cost of resources. Research around hospital readmission is mainly in following three focuses. The first focus is to establish evidence base supporting effective and safe care. For instance, one study specifically analyzes how much the impact of HbA1c measurement on hospital readmission rates on inpatient diabetes encounters (Strack et al., 2014). In another study, the objective is to identify the risk factors and to enhance the medical procedures for elective primary total joint arthroplasty (Ramkumar et al., 2015). The second focus is to compare ML modeling methodology of hospital readmission rate prediction. In the study conducted by (Mortazavi et al., 2016), the authors compare the effectiveness of several classification methods (random forests, boosting, etc.) to predict 30-day and 180-day all-cause readmissions due to heart failure. Another similar research proposes various ML approaches to identify the readmission risk group for individual heart failure patient (Zheng et al., 2015). The goal of these studies is often to recognize the best ML model that outperforms others in terms of accuracy or sensitivity. The third focus, also the least researched area, is to transform the analysis result into actions. It is crucial because we must identify the owner who is capable and motivated to take this action in order for any ML/AI research having an impact in real applications. Taking following study as an example. A classifier is first built to categorize individual patient's readmission risk. With these prediction results, a decision analysis is then proposed to guide the allocation of post-discharge support to congestive heart failure patients (Bayati et al., 2014). In this case, hospitals have the strong motivation to carry the decision analysis forward, since hospitals face a fine on 30-day readmissions of certain disease causes according to the federal readmission penalty program (McIlvennan, Eapen, & Allen, 2015). In (Shickel, Tighe, Bihorac, & Rashidi, 2018), the authors give a survey about applying machine learning in predictive care in the EHR (electronic health record) analysis.

In this paper, we propose a real-world decision support sys-

Yan Zhang et al. This is an open-access article distributed under the terms of the Creative Commons Attribution 3.0 United States License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

tem that realizes the decision analysis outlined in (Bayati et al., 2014). With the objective of reducing the 30-day hospital readmission rate, a common approach is to assign post-discharge intervention programs for supporting discharged patients. Example programs include post-discharge care coordination, patient education and self-management support, scheduled outpatient visits, tele-medicine, etc. However, it is not economically feasible to provide such programs to all patients. Therefore, a decision need to be made for assigning intervention programs to high risk patients.

The rest of this paper is organized as follows. In Section 2.1, we present the sample data and ML modeling approach in the proposed decision system. In Section 2.2, we outline the architecture components for data injection, processing, and the ML model deployment pipeline. In Section 2.3, we illustrate the visualization dashboard showing the decision analysis for patient-specific intervention program recommendation.

2. METHOD

In this section, we present a general decision support system that utilizes machine learning (ML) based patient-specific prediction to assist doctors to make the decision on whether to enroll a patient into an intervention program or not. The solution implementation including the source code and visual dashboard can be obtained from (Zhang, Bleik, & Wahl, 2017).

2.1. Data and Modeling

We use two sources of data for building the ML model: data at rest and data in motion. The former dataset is obtained from UCI Data Repository - Diabetes 130-US hospitals for years 1999-2008 Data Set (Strack et al., 2014). In (Strack et al., 2014), the authors show that the HbA1c measurement is a useful predictor of readmission rates for patients with diabetes mellitus. This dataset represents 10 years of clinical care at 130 US hospitals and integrated delivery networks. It includes over 50 features representing patient and hospital outcomes. Some features include patient number, race, gender, age, admission type, time in hospital, medical specialty of admitting physician, number of lab test performed, HbA1c test result, diagnosis, number of medication, diabetic medications, number of outpatient, inpatient, and emergency visits in the year before the hospitalization, etc. The latter data is simulated glucose sensor readings, which mimics the scenario where real-time information is needed to take into account for model prediction.

We construct a classifier based on patients’ demographic information (age, gender, zip code), medical record history (e.g. present illness, physical exam results, and medication), and glucose readings to stratify patients into different readmission risk groups. This classifier is trained offline using pySpark with Random forest algorithm with evaluation accuracy

```

Train the model
trained_model = RandomForestClassifier(featuresCol='features', labelCol='label').fit(train)

Store the model:

model_filename = 'wasb://model@{}.blob.core.windows.net/model'.format(account_name)
trained_model.save(model_filename)

Evaluate the model
mc_evaluator = MulticlassClassificationEvaluator(labelCol='label')

predictions = trained_model.transform(test)
precision = mc_evaluator.evaluate(predictions, {mc_evaluator.metricName: "weightedPrecision"})
recall = mc_evaluator.evaluate(predictions, {mc_evaluator.metricName: "weightedRecall"})
accuracy = mc_evaluator.evaluate(predictions, {mc_evaluator.metricName: "accuracy"})
print('Accuracy: {:.3f}\nWeighted precision: {:.3f}\nWeighted recall: {:.3f}\n'.format(accuracy, precision, recall))

Accuracy: 0.883
Weighted precision: 0.836
Weighted recall: 0.883
    
```

Figure 1. Model and Evaluation

88.3%. Please see detailed results in Figure 1 (Zhang et al., 2017).

It should be noticed that, the datasets we have used is for demonstration purpose only. It does not show solid insights (e.g. feature importance, etc) applicable in real clinical studies.

2.2. Solution Architecture

We show the solution architecture diagram to illustrate the scoring pipeline in Figure 2. The simulated patient glucose data is generated via the Data Simulator Web Job and ingested into the Event Hub. It is then being aggregated in Stream Analytics and sinks in Azure Storage Blob. These data then serves the raw scoring data, which goes through sequential steps such as data preprocessing, feature engineering, and scoring with the pre-trained classifier (get predictions of the probability of 30-day readmission). All these procedures are implemented using PySpark. Azure Data Factory is used to schedule scoring at daily basis for patients at discharge every day, as well as copying the results from Azure Storage Blob to Azure SQL database.

2.3. Solution Dashboard

In Power BI visualizations, we present a decision system that uses the predictions from the classifier to guide decisions about post-discharge interventions. With a set of tunable parameters such as per patient readmission cost, cost of the intervention program, and efficacy of the intervention programs (Bayati et al., 2014), we report the performance of the methodology and show the overall expected value of employing a real-time decision system.

There are two report pages included in the dashboard. The CMIO report aims to serve the management team such as the chief medical information officer for high level performance of this patient-specific prediction and interventions assignment system. By comparing different intervention programs with varying cost and efficacy, the manager can obtain the expected cost savings or lost. The doctor report, as shown

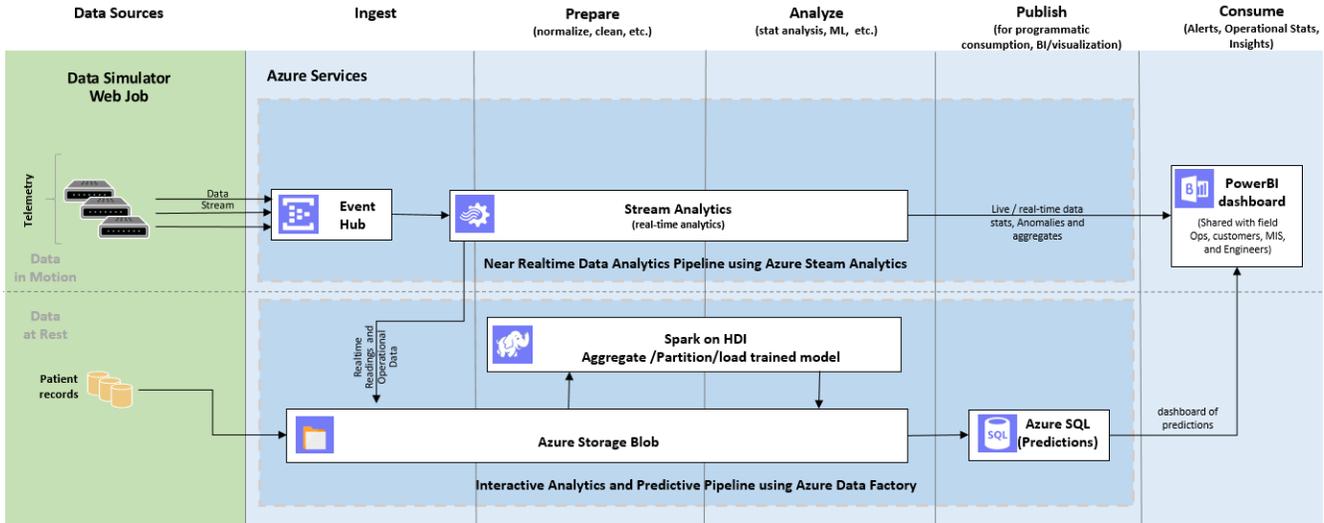


Figure 2. Solution Architecture Diagram

in Figure 3, shows the patient level information to doctors who need to make decision on patient-specific intervention program suggestions.

Between the two extreme cases, assign-to-all-patients and assign-to-no-patients, the decision system aims to find a balance point – driven by total cost to the hospital. The potential costs include: (1) Cost of Readmission (per patient) e.g. \$10,000; (2) Cost of Intervention program e.g. \$1800. We also take into account of the Efficacy of Intervention program (e.g. 35%). Based on all these information, the decision system calculates the threshold of prediction score $p = 0.376$ in above example. In another word, if a patient has a probability of 37.6% or above to get re-admitted then assign the intervention program. The doctor can take this suggestion as an input information to make the final decision.

In the real scenario, ideally we make usage of data available within the EHR system for building the classifier. When the dashboard is being used for real-time clinical use, the doctor can obtain findings on specific patients without spending additional time. The automatic enrollment recommendations will be provided based on the computed risk score. This solution will greatly improve the work efficiency comparing to the typical method of filling out a paper form and manually evaluating each patient.

3. DISCUSSION

Machine Learning (ML) and Artificial Intelligence (AI) are making tremendous progress not only in academia but also in industry. Data driven ML solution, which now is a topic that every enterprise is taken into serious consideration. An ML based solution is exciting, but not always successful. It leads to a failure largely because following questions are not clarified in advance.

- Where does the ML solution fit into the existing business pipeline?
- Who are the users of the solution?
- How much benefit for the potential users will gain?
- How much cost for establishing the ML solution. Does it pay off?

When we study the hospital readmission problem, the basic ML question is to build a model that predicts the probability of a patient at discharge to be readmitted within 30 days. Why should a data scientist care about the business problem? The reason is simple. Because the data we model on is tied to the business goal. It matters whether the research is on all-cause readmission or cause-specific readmission, since we must choose relevant data based on the answer to this question. It matters whether the business goal is to save the hospital's readmission cost, or to improve the clinical procedures. If the answer to above questions is the latter, data scientists may want to choose to use interpretable modeling approaches such as decision tree, rather than the more accurate but black-box deep neural networks.

U.S. health-care system is a very complex system, which consists of many players. It is unrealistic to improve its efficiency in one shot. When talking about saving the medical cost, is the goal to save cost for hospitals, for patients (hospital cost, doctor's cost, lab procedure cost, drug cost, etc.), or for insurance companies? When talking about improving patients' care, is the goal to reduce unnecessary clinical procedures, or to improve doctor/nurse's efficiency/efficacy? For each of these goals, a different set of data is needed for building an ML model.

Doctor Daily Report



This report shows the suggested decision on whether to recommend intervention programs to a patient at discharge time.

Threshold Probability

37.6%

Patient-specific Decision Suggestion

Patient ID	Probability of Readmission	Expected cost of readmission	Expected cost of applying intervention program	Expected savings from intervention program	Suggest Intervention ?
101487563	40.9 %	\$5,595	\$5,437	\$158	Yes
102249958	37.1 %	\$5,075	\$5,099	(\$24)	No
105608490	32.1 %	\$4,391	\$4,654	(\$263)	No
113768967	39.1 %	\$5,348	\$5,277	\$72	Yes
117401721	41.7 %	\$5,704	\$5,508	\$196	Yes
118734959	33.5 %	\$4,582	\$4,779	(\$196)	No
119811218	34.7 %	\$4,747	\$4,885	(\$139)	No
121019127	35.3 %	\$4,829	\$4,939	(\$110)	No
134554254	31.5 %	\$4,309	\$4,601	(\$292)	No
135426674	29.8 %	\$4,076	\$4,450	(\$373)	No
137787363	37.0 %	\$5,061	\$5,090	(\$29)	No
140009504	36.9 %	\$5,048	\$5,081	(\$33)	No
144882752	37.1 %	\$5,075	\$5,099	(\$24)	No
146265098	38.9 %	\$5,321	\$5,259	\$62	Yes
148322486	33.6 %	\$4,596	\$4,787	(\$191)	No
151229840	27.6 %	\$3,775	\$4,254	(\$479)	No

Programs: program1, program2, program3, program4, program5, program6

Intervention Cost: \$1,800 | Efficacy: 35 %

Readmission Cost: \$10,000, \$13,679, \$15,000

Patient ID: 101487563, 102249958, 105608490, 113768967, 117401721, 118734959

Expected savings from intervention program by Patient ID

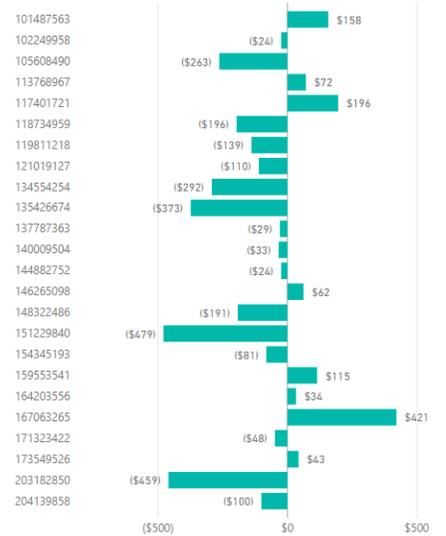


Figure 3. Daily Report Dashboard for Doctors

4. CONCLUSION

In this paper, we proposed a solution to reduce hospital readmission cost for hospitals. We introduced a general ML based patient-specific decision support system using a cloud based platform, Microsoft Azure. We choose to use PySpark as the ML modeling language in order for the potential scalability needs. We further use a visualization dashboard to show the decision making process for assigning a patient to post-discharge intervention program or not.

This work has three major contributions. First, we choose to use pySpark as the ML modeling language in order for the potential scalability needs. Second, we outline solution architecture components including all the data science processes for a real ML system. Third, we discuss the topics around how the ML modeling results can be taken a step forward into real-world health care application.

REFERENCES

Bayati, M., Braverman, M., Gillam, M., Mack, K. M., Ruiz, G., Smith, M. S., & Horvitz, E. (2014). Data-driven decisions for reducing readmissions for heart failure: General methodology and case study. *PLoS ONE*, 9(10), 1–9.

McIlvennan, C. K., Eapen, Z. J., & Allen, L. A. (2015).

Hospital readmissions reduction program. *Circulation*, 131(20), 1796–1803.

Mortazavi, B. J., Downing, N. S., Bucholz, E. M., Dharmarajan, K., Manhapra, A., Li, S.-X., ... Krumholz, H. M. (2016). Analysis of machine learning techniques for heart failure readmissions. *Circulation: Cardiovascular Quality and Outcomes*, CIRCOUTCOMES–116.

Ramkumar, P. N., Chu, C. T., Harris, J. D., Athiviraham, A., Harrington, M. A., White, D. L., ... Li, L. T. (2015). Causes and rates of unplanned readmissions after elective primary total joint arthroplasty: a systematic review and meta-analysis. *American journal of orthopedics*, 44(9), 397–405.

Rumsfeld, J. S., Joynt, K. E., & Maddox, T. M. (2016). Big data analytics to improve cardiovascular care: promise and challenges. *Nature Reviews Cardiology*, 13(6), 350.

Shickel, B., Tighe, P. J., Bihorac, A., & Rashidi, P. (2018). Deep ehr: A survey of recent advances in deep learning techniques for electronic health record (ehr) analysis. *IEEE journal of biomedical and health informatics*, 22(5), 1589–1604.

Strack, B., DeShazo, J. P., Gennings, C., Olmo, J. L., Ventura,

S., Cios, K. J., & Clore, J. N. (2014). Impact of hba1c measurement on hospital readmission rates: analysis of 70,000 clinical database patient records. *BioMed research international*, 2014.

Zhang, Y., Bleik, S., & Wahl, M. (2017). *Patient-specific readmission prediction and intervention for health care*. <https://github.com/YanZhangADS/cortana-intelligence-population-health-management/tree/master/Spark>.

Zheng, B., Zhang, J., Yoon, S. W., Lam, S. S., Khasawneh, M., & Poranki, S. (2015). Predictive modeling of hospital readmissions using metaheuristics and data mining. *Expert Systems with Applications*, 42(20), 7110–7120.

BIOGRAPHIES



Yan Zhang is a Senior Data Scientist in New England Research Center, Microsoft, Cambridge MA, USA. She builds predictive analytics models and generalize machine learning solutions on Microsoft’s analytical platforms. She worked data science project in various verticals including customer market segmentation, predictive maintenance in manufacture, and hospital readmission prediction in health-care. She received her Ph.D. in data mining, computer science. She holds a Ph.D in data mining and has 10+ years of technical experience on machine learning and business intelligence. Before joining Microsoft, she was a research faculty at Syracuse University.