Contents lists available at ScienceDirect



Data Science and Management



journal homepage: www.keaipublishing.com/en/journals/data-science-and-management

Research article

Machine learning-based prediction models for patients no-show in online outpatient appointments



Guorui Fan^a, Zhaohua Deng^{b,*}, Qing Ye^{a,c}, Bin Wang^d

^a School of Medicine and Health Management, Tongji Medical College, Huazhong University of Science and Technology, 430074, Wuhan, China

^b School of Management, Huazhong University of Science and Technology, 430074, Wuhan, China

^c Department of Information Management, Tongji Hospital, Tongji Medical College, Huazhong University of Science and Technology, 430074, Wuhan, China

^d Robert C. Vackar College of Business and Entrepreneurship, University of Texas Rio Grande Valley, Edinburg, TX78539, USA

ARTICLE INFO

Keywords: Online health Online outpatient appointment Patient no-show Prediction model Machine learning

ABSTRACT

With the development of information and communication technologies, all public tertiary hospitals in China began to use online outpatient appointment systems. However, the phenomenon of patient no-shows in online outpatient appointments is becoming more serious. The objective of this study is to design a prediction model for patient no-shows, thereby assisting hospitals in making relevant decisions, and reducing the probability of patient no-show behavior. We used 382,004 original online outpatient appointment records, and divided the data set into a training set (N₁ = 286,503), and a validation set (N₂ = 95,501). We used machine learning algorithms such as logistic regression, k-nearest neighbor (KNN), boosting, decision tree (DT), random forest (RF) and bagging to design prediction models for patient no-show in online outpatient appointments. The patient no-show rate of online outpatient appointment was 11.1% (N = 42,224). From the validation set, bagging had the highest area under the ROC curve and AUC value, which was 0.990, followed by random forest and boosting models, the area under ROC and AUC values of the logistic regression, decision tree, and k-nearest neighbors were lower at 0.597, 0.499 and 0.843, respectively. This study demonstrates the possibility of using data from multiple sources to predict patient no-shows. The prediction model results can provide decision basis for hospitals to reduce medical resource waste, develop effective outpatient appointment policies, and optimize operations.

1. Introduction

In recent years, outpatient clinics have taken center stage in healthcare systems due to an emphasis on preventive medical practices, shorter hospital stays, and service provision on an outpatient basis (Hooshangi-Tabrizi et al., 2020; Pan et al., 2021; Srinivas and Ravindran, 2020). Outpatient appointment scheduling systems have become an important component for efficient care delivery and demand management in outpatient clinics (Javid et al., 2017; Kuiper et al., 2021). And "too difficult to see a doctor" is a persistent problem in China (Liu, 2009; Yip and Hsiao, 2008; Zhang et al., 2014). To regulate hospital service capacity and satisfy patient demand, outpatient appointment scheduling systems are widely used by healthcare providers (Liu, 2016). An effective outpatient appointment system can improve the efficiency of hospital operation and the delivery of medical services, as well as enhancing patient satisfaction and improving the economic and social benefits of hospitals (Lee et al., 2018).

In order to streamline healthcare facilities' operations, provide better medical service to patients in remote areas, and create a good outpatient appointment diagnosis and treatment services system, online appointment systems as new Internet-based appointment systems have received extensive attention from hospitals in China (Cao et al., 2011). With the development of information and communication technology and support from the Ministry of Health of China, all public tertiary hospitals in China began to use online appointment systems in 2009 (Zhang et al., 2014). However, with the growing popularity of online appointment scheduling in outpatient clinics, patient no-show behavior in outpatient online appointment systems has become more serious.

Patient no-show refers to the case where a patient does not come to the clinic or hospital at the appointment time or cancels the appointment

* Corresponding author.

https://doi.org/10.1016/j.dsm.2021.06.002

Received 15 April 2021; Received in revised form 21 June 2021; Accepted 22 June 2021 Available online 24 June 2021

2666-7649/© 2021 Xi'an Jiaotong University. Publishing services by Elsevier B.V. on behalf of KeAi Communications Co. Ltd. This is an open access article under the CC BY-NC-ND license (http://creativecommons.org/licenses/by-nc-nd/4.0/).

Peer review under responsibility of Xi'an Jiaotong University.

E-mail address: zh-deng@hust.edu.cn (Z. Deng).

shortly before the appointment time (usually within one day), resulting in the appointment time slot not being filled (Ding et al., 2018; Huang and Hanauer, 2014). At present, the no-show rate of outpatient appointments in most hospitals in China is 10%-20%, with those in a few hospitals exceeding 30% (Cronin et al., 2013; Ding et al., 2018; Distelhorst et al., 2018; Kheirkhah et al., 2016; Lenzi et al., 2019). The high rate of appointment no-shows for outpatients leads to a waste of medical resources, reduces the opportunity for other patients in obtaining the healthcare they need (Fiorillo et al., 2018), reduces patients' health conditions (Kheirkhah et al., 2016; Kurasawa et al., 2016; Lenzi et al., 2019), and leads to the suboptimal allocation of medical resources (Alaeddini et al., 2015; Nuti et al., 2012). Patient no-show behavior will disturb the normal order of outpatient service, lead to the extension of work hours of outpatient doctors, increase the burden of doctors' work, reduce hospital operating efficiency, and is not conducive to the construction of normal medical order (Lee et al., 2018). At the same time, patient no-show behavior will not only exert a negative impact on hospital management but also have a significant impact on the financial performance and resource utilization of the healthcare system (Liu, 2016)

The purpose of this study is to predict patient no-show behavior in online outpatient appointment scheduling systems to reduce the adverse consequences of such behavior, improve the operating efficiency of hospitals' outpatient appointment scheduling systems, and achieve better usage of medical resources. Previous research on patient no-shows have relied on the use of offline predictors (Ding et al., 2018; Lekham et al., 2020; Lenzi et al., 2019). In online outpatient appoint scheduling systems, not only offline factors but also online factors such as the doctor's online reputation may affect patient no-show behavior. Our research is the first to combine both online and offline factors in the prediction of patient no-shows. In addition, machine learning is widely used to develop risk prediction models (Goldstein et al., 2017) such as predicting patient readmission (Bardhan et al., 2014), identifying patient adverse events (Rochefort et al., 2015), and developing inpatient mortality predictive models (Tabak et al., 2014). Thus, this study applies machine learning algorithms to construct prediction models for patient no-show behavior. Effective prediction results could assist hospitals in making relevant decisions, reducing the probability of patient no-shows for appointments made through online outpatient appointment systems, improving the operational efficiency of hospital outpatient clinics, and improving the economic and social benefits of hospitals.

2. Literature review

2.1. Online outpatient appointment scheduling

Providing high-quality and high-efficiency medical services to meet patients' emerging health needs is one of the main challenges facing the Chinese medical industry (Yip and Hsiao, 2014). The rapid growth of information and communication technologies and the increasing use of the Internet and mobile devices have contributed to a flourishing online medical services industry in China. Currently, China's tertiary public hospitals provide various medical services through online platforms, including the WeChat platform, self-developed applications, and various services provided in cooperation with third-party platforms. The services provided mainly include online consultation (Kamsu and Foguem, 2014; Wu and Lu, 2017), outpatient appointment scheduling (Colaci et al., 2016; Hoque, 2016), electronic medical prescription, and online payment. Among these services, the outpatient appointment scheduling is the most heavily utilized, which is the focus of our study.

Online outpatient appointment scheduling systems provide medical services through the Internet. Patients can make outpatient appointments through websites or mobile phone applications (Mold and Lusignan, 2015). Online outpatient appointments break down fixed time and distance barriers. Through online outpatient appointments, hospitals can provide effective access to medical services for patients in rural and

remote areas, expanding their coverage of medical services (Ahmed et al., 2014; Xie et al., 2017). The increase in demand for outpatient services, coupled with shortages in the supply of physicians, will lead to a supply-demand mismatch, such as resource burnout, medical errors, decrease in productivity, longer patient waiting times, and longer appointment delays. Thus, both patient satisfaction and resource utilization will be negatively impacted (Srinivas and Ravindran, 2018). Meanwhile, online outpatient appointments have been suggested to improve the workflow, thereby reducing the wait time and improving the patient experience (Mey and Sankaranarayanan, 2013).

2.2. Patient no-show behavior

Previous studies have shown that the phenomenon of patient no-show is not accidental, and the occurrence of patient no-show behavior will be affected by a variety of factors (Cronin et al., 2013; Daggy et al., 2010; Dantas et al., 2018). Existing studies have shown that patients' demographic characteristics can influence their no-show behavior, and younger patients are more likely to be no-show in outpatient appointments (Cronin et al., 2013; Fiorillo et al., 2018; Lehmann et al., 2007), but some studies suggest that younger patients have a lower rate of no-show in outpatient appointments (Zhou et al., 2018). At the same time, some studies suggest that age has no significant influence on patient no-show in an appointment (Dantas et al., 2019). In terms of the influence of gender on patient no-show behavior, previous studies suggest that compared with men, women have a lower risk of no-show in outpatient appointments (Liu, 2016), with other studies revealing opposite results (Kheirkhah et al., 2016). In addition, patients' income, marital status, race, and other characteristics can affect the patient no-show behavior (Daggy et al., 2010; Lee et al., 2005; Lehmann et al., 2007). In terms of the appointment time, patient with a long lead time of appointment is at a higher risk of no-show in outpatient appointments (Chang et al., 2015; Peng et al., 2016; Rosenbaum et al., 2018). Since the outpatient system is usually not open on weekends, patients have a higher rate of no-show on Mondays, and afternoon appointments have a higher rate of no-shows (Cronin et al., 2013; Peng et al., 2016). Some studies have found that the weather condition on the day when patients visit outpatient clinics is also an important predictor of patient no-show behavior (DeFife et al., 2010; Peng et al., 2016). Moreover, patient distance from the hospital (Daggy et al., 2010; Dantas et al., 2019; Whiting et al., 2015), previous patient visit experience (Fiorillo et al., 2018; Jain and Chou, 2000), medical insurance (Peng et al., 2016; Whiting et al., 2015), appointment channels (Zhou et al., 2018) and appointment specialties (Jain and Chou, 2000; Rosenbaum et al., 2018) also affect the patient no-show behavior in outpatient appointments.

In order to effectively predict patient no-show behavior and reduce the negative impact of patient no-shows on hospitals, this study uses machine learning algorithms to predict patient no-show behavior using online outpatient appointment data.

2.3. Predicting patient no-shows

In the past few years, patient no-shows have been a subject of extensive study, with growing applications of machine learning algorithms in these studies. Based on regularized logistic regression and lasso regression, Ding et al. (2018) constructed three levels of prediction models and showed that specific clinical models are better than general models. Lenzi et al. (2019) identified the previous no-show rate and whether it was an appointment for the day as the most important predictors of patient no-shows using mixed-effects logistic. Simsek et al. (2021) built a Tree Augmented Naive Bayes-based model to predict no-shows of minority patients. Srinivas and Salah (2021) used random forests, stochastic gradient boosting, and deep neural networks to predict patient no-show. And found the stochastic gradient boosted classification tree has the best performance. Another study uses new wrapper methods based on opposition-based, self-adaptive, and cohort intelligence

algorithms. The results showed better performance in terms of dimensionality reduction and convergence speed with similar area under the curve (AUC), high sensitivity and specificity scores (Aladeemy et al., 2020). In outpatient primary care in rural areas, Lekham et al. (2020) found that appointment lead time is an important predictor for patient no-shows by using logistic regression, decision tree, and tree-based ensemble classifiers. By using logistic regression, random forest, gradient booster, and artificial neural network, Srinivas (2020) found that a single model cannot perform best in predictive performance, training time, and interpretability. The important predictors of patient no-show are previous no-show history, age, and afternoon appointments.

2.4. Research gaps and contribution to the literature

Our study identifies and addresses the research gap in the literature pertaining to the prediction for patient no-shows. In online outpatient appointments, patients are more susceptible to the online reputation of doctors. However, existing studies of patient no-shows only focus on characteristics of the outpatient appointment system and ignore the influence of online information on patient no-show behavior, such as the online reputation of doctors. Our study addresses this gap in the patient no-show prediction literature by including online doctor rating as one of the features. To the best of our knowledge, this work is one of the first for predicting patient no-show behavior using online doctor rating, in addition to patient information and appointment information. Based on the data on outpatient appointments, our study innovatively integrates data generated online as the predictor variables to better understand the patient no-show behavior in the context of online outpatient appointments. This study demonstrates the possibility of using data from patients and outpatient appointment systems to predict patient no-shows.

3. Materials and methods

3.1. Research context and available data

This study draws data from a large general hospital in central China. We extracted data from outpatient appointment records in EMR systems from May to August 2019, with a total of 454,217 outpatient visit records. Then, we selected patients who used the online outpatient appointment system as the sample of our study. The ratings of the hospital's outpatient doctors were drawn from the Good Doctor Online (www.haodf.com), given that it is the earliest and largest online doctor review and online healthcare community website in China (Hao, 2015). These two data sets were merged using the doctors' names. Any incomplete or erroneous records were removed to ensure the reliability of the data. After removing invalid records, we obtained a sample size of 382,004.

3.2. Dependent variable definition

Previous studies defined patients' no-show behavior as the case where "the patient did not show up according to a scheduled appointment or canceled an appointment when it is close to the time of the scheduled appointment" (such as canceling the appointment on or before the day of the scheduled appointment), and thus, the outpatient appointment cannot be reassigned to another patient (Ding et al., 2018; Huang and Hanauer, 2014; Lenzi et al., 2019). Based on the outpatient appointment scheduling mechanism of the hospital, we define a patient's no-show behavior as the patient's failure to attend a scheduled appointment or cancellation of an outpatient appointment after 6 a.m. on the day of the scheduled appointment.

3.3. Predictor variables

In addition to age and gender, the main predictors of this study are appointment lead time (Leadtime), appointment time, weekday appointment, online doctor rating (DOC RATING), appointment doctor

type (DOC TYPE), the patient's distance from the hospital (Distance), previous outpatient visit experience (EXP) and other predictive variables for a total of 15 predictors. Appointment lead time is the number of days between the patient's outpatient appointment creation time and appointment time. Patient visit time is the time of the patient visit (morning or afternoon). Weekday appointment is the patient's scheduled appointment day in a week (Monday, Tuesday, Wednesday, Thursday, Friday, or weekend). The online doctor rating is the outpatient doctor's online rating and ranges from 0 to 5. Reputation is a vital quality factor in health care delivery (Huang and Zuniga, 2014) and is considered the most valuable attribute of a physician (Romano and Baum, 2014). An increasing number of Chinese consumers have used online doctor reviews to rate their doctors or to look for a particular doctor to tend to their health care concerns (Huang and Hanauer, 2016). Thus, the reputation of outpatient doctors will affect patients' no-show behavior. We used online doctor ratings to measure the outpatient doctor's reputation. Appointment doctor type is measured by the type of doctor appointment (expert or regular). The patient's distance from the hospital is measured by the distance from the patient's location to the hospital (less than 300 km or more than 300 km) (Ye et al., 2019). Previous outpatient visit experience is measured by whether the patient was visiting the outpatient clinic for the first time. The hospital district captured three districts with different geographical locations, sizes, and doctor compositions. Patient registration time was the date the patient registered in the appointment system. Appointment creation time was the date the patient created an outpatient appointment in the appointment system. Appointment doctor title was the medical title of the doctor (e.g., professor, associate professor, or others). The patient's province reflects the province the patient resided in, and appointment time was the date of the patient's appointment time with the outpatient clinic.

3.4. Features selected and development of predictive models

To eliminate redundant collinear features, we performed feature selection by recognizing the most predictive variables using the least absolute shrinkage and selection operator (LASSO) (Fu et al., 2018; Gao et al., 2020; Liang et al., 2020). LASSO added the L1 norm of the feature coefficients as a penalty term to the loss function, which forced the coefficients corresponding to those weak features to be zero (Gao et al., 2020). Therefore, we considered features with zero coefficients as redundant features and deleted them. Among the 17 features obtained from the outpatient appointment system, we obtained 15 features, including 12 categorical features and 3 continuous features that underwent feature selection by LASSO as shown in Fig. 1.

Next, we randomly divided the data set into a training data set ($N_1 = 286,503$) and a validation data set ($N_2 = 95,501$). The study uses the training set to train the patient no-show prediction model and uses the validation set to evaluate model performance and ensure the robustness of the prediction model results. We used ten-fold cross validation in the training set to ensure the accuracy of the results. In this process, the training set was split into ten subsets with nine of the subsets being used as the training set and the remaining subset being used as the testing data set. This process was repeated ten times where each of the ten subsets was used as the testing set once.

In this study, machine learning algorithms including logistic regression (LR), k-nearest neighbor (KNN), decision tree (DT), random forest (RF), bagging, and boosting were used to design the patient no-show prediction model. The LR algorithm is a widely used classifier in medical settings (Srinivas, 2020; Srinivas and Ravindran, 2018). As a lazy learning algorithm, KNN is one of the most fundamental and simple classification methods. Recently, decision tree has also played a significant role in medical decision support such as predicting patient no-shows and clinical diagnosis (Alves et al., 2021; Lekham et al., 2020). Ensemble methods, which combine the predictions of multiple algorithms, are known to reduce the variance and yield a superior predictive performance (Polikar, 2012). Thus, our study also considered using random



Fig. 1. Feature selection by LASSO.

forest, bagging, and boosting methods in building our predictive models.

We used the R language to design patient no-show prediction models. After the manual screening, this study specified the nearest known classification sample point K = 3 in the KNN model. For random forest, the number of variables used in the binary tree in the designated node was set to 4. For the bagging method, the number of variables used in the binary tree in the designated node was set to 15. In boosting, we set the compression parameter $\lambda = 0.1$. The rest of the model parameters were the default values. In our study, the categorical variables were one-hot encoded using the R Mltools package. Due to the uneven distribution of the dependent variables in this study, the sample is unbalanced. Therefore, we used the R DMwR package Synthetic Minority Oversampling Technique (SMOTE) on the training data set to obtain a more balanced sample (Torgo and Torgo, 2013).

3.5. Predictive models evaluation metrics

Model performance is typically evaluated via precision, recall, accuracy, F1, and area under the receiver operating characteristic (ROC) curve. The metrics are calculated using the outputs of the confusion matrix for each class: TP for true positives, TN for true negatives, FN for false negatives, and FP for false positives. Our evaluated metrics are defined as follows:

$$Precision = \frac{TP}{TP + FP} \tag{1}$$

$$Recall = \frac{TP}{TP + FN}$$
(2)

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$
(3)

$$F1 = \frac{2 \times TP}{2 \times TP + FP + FN}$$
(4)

4. Results

4.1. Descriptive statistics

Table 1 shows a summary of the variables' descriptive statistics. We report the frequency distribution for each categorical variable and the mean and standard deviation for each continuous variable. In outpatient online appointments, patients' no-show rate was 11.1%. The mean age of outpatients was 36.6 years (SD = 19.6). Moreover, 40.9% of the patients were male, whereas 59.1% were female. In online outpatient appointments, the average appointment lead time was 4.9 days (SD = 19.6), which indicates that the online outpatient appointment has facilitated

Tal	ble	1	
	-		-

Sample descriptive statistics

Categorical variables	Frequency	Percentage (%)
Patient no-shows		
No-show	42,224	11.1
Show	339,780	88.9
Gender		
Male	156,307	40.9
Female	225,697	59.1
Patient visit time		
Morning	251,414	65.8
Afternoon	130,590	34.2
Weekday appointment		
Monday	85,837	22.5
Tuesday	76,420	20.0
Wednesday	70,257	18.4
Thursday	68,786	18.0
Friday	46,046	12.0
Weekend	34,658	9.1
Appointment doctor type		
Expert	317,788	83.2
Ordinary	64,216	16.8
Patient's distance from the hospital		
Less than 300 km	327,191	85.7
More than 300 km	54,813	14.3
Appointment doctor title		
Professor	196,727	51.5
Associate professor	121,061	31.7
Others	64,216	16.8
Previous outpatient visit experience		
First visit	256,506	67.2
Not first visit	125,498	32.9
Hospital district		
District 1	295,710	77.4
District 2	54,618	14.3
District 3	31,676	8.3
Patient registration time ^a		
Appointment creation time ^a		
Appointment time ^a		
Patient's source province ^a		
Continuous variables	Mean	Standard deviation
Age (in years)	36.6	19.6
Appointment lead time (in days)	4.9	4.7
Online doctor rating	3.9	0.2

Notes.

^a See Appendix A.

patients to some extent. As the online outpatient appointment system breaks the distance barrier, the average distance of patients with online appointments to the state hospital was 184.6 km (SD = 313.5).

4.2. Prediction model performance

The performance of the prediction model are shown in Table 2. For the validation data set, the prediction accuracy of the logistic regression model is 68.79%. However, the precision rate of the logistic regression model is only 15.34%, its recall rate is 40.33%, and the F1 values of logistic regression is 0.222. The performance of logistic regression as a prediction model is poor. The prediction accuracy of the KNN is 90.14%, slightly higher than the prediction accuracy of the logistic regression model. The precision and recall of the prediction model constructed according to the KNN are 53.77% and 76.75%, respectively, with an F1

Table 2

Prediction model results on the validation data se	set.
----------------------------------------------------	------

	Precision (%)	Recall (%)	Accuracy (%)	F1	AUC
Logistic	15.34	40.33	68.79	0.222	0.597
KNN	53.77	76.75	90.14	0.632	0.843
Boosting	98.30	96.57	99.44	0.974	0.976
Decision tree	12.62	13.58	80.06	0.131	0.499
Random forest	79.72	96.70	96.92	0.874	0.987
Bagging	90.52	96.54	98.50	0.934	0.990

value of 0.632. The prediction accuracy of the boosting method is 99.44%, higher than the prediction accuracy of the logistic regression model and the KNN method. The precision and recall of boosting are 98.30% and 96.57%, respectively, and the F1 is 0.974. These results show that the boosting model has a better performance. The accuracy, precision and recall rate of the DT prediction model are 80.06%, 12.62% and 13.58%, respectively, and its F1 value is 0.131. Although DT has a high accuracy rate, its precision and recall rates are low, so the performance of the model is poor. As for the RF prediction model, the accuracy is 96.92%, and the precision of RF is slightly lower, which is 79.72%. Due to the RF prediction model has a better recall and F1 value, which are 96.70% and 0.874, respectively, the RF has a good predictive performance. As for bagging, the precision, recall and accuracy are 90.52%, 96.54% and 98.50%, respectively, and the F1 value is 0.934. Therefore, bagging has a good performance as a predictive model of patient no-show behavior. This result shows that the boosting, RF and bagging method have higher precision, recall, accuracy rates and F1 values, and these prediction models perform better than the logistic regression, KNN and DT

Fig. 2 shows the ROC curve and the AUC value of each prediction model. The bagging algorithm has the highest area under the ROC curve and AUC value, which is 0.990. RF and booting also have good areas under the ROC curves and AUC values, with AUC values of 0.987 and 0.976, respectively, followed by the KNN method with an AUC value of 0.843. However, compared with the bagging, RF, boosting and KNN prediction models, the areas under ROC and AUC values of the logistics regression and DT are lower at 0.597 and 0.499, respectively. The DT has the lowest area under ROC and AUC values, thus it is not suitable for the design of predictive models of patient no-show behavior.

Thus, boosting, RF and bagging prediction models have the best precision, recall, accuracy, F1 and AUC values, and they are suitable to design prediction models of patient no-show. Although the accuracies of the logistic regression, DT and KNN prediction models are relatively high, these algorithms are not suitable to design prediction models for patient no-show with online outpatient appointments due to the low precision and recall rates.

4.3. Important predictor variables

Finally, we assessed important predictor variables of the prediction models. The predictor variable importance of each model was calculated using the varImp function in the caret R package. Fig. 3 shows the importance of all predictors. Among the three best prediction models, patient registration time, patient visit time, and appointment creation time are important predictors. In addition, appointment doctor type and the patient's distance from the hospital are important predictors of patient no-shows. Due to the different characteristics of the outpatient online appointment systems, our study innovatively found that online doctor rating is an important predictor of patient no-shows.

5. Discussion and contributions

5.1. Discussion

Patient no-show behavior reduces revenues and impair the delivery of quality healthcare. Many studies have been devoted to predicting patient no-shows in outpatient appointment systems. Different from traditional outpatient appointment systems, online outpatient appointment systems overcome the limitation of fixed time and space, thus improving the accessibility and fairness of medical services. At the same time, online outpatient appointment systems also have their own features such as a wide range of patients from different geographical areas and a long appointment lead time. Therefore, it is important to effectively predict patient no-show behavior in online outpatient appointments to improve the operating efficiency of hospital outpatient services. The main purpose of this study is to use machine learning algorithms such as boosting, random forest, and bagging algorithms to design reasonable prediction models for patient no-show behavior in online outpatient appointments. The robustness of the prediction models is tested on the validation set.

Consistent with previous research, we found patient appointment time, prior visit experience, and patient's age to be important features predicting patient no-shows (Cronin et al., 2013; Dantas et al., 2018; Peng et al., 2016). Earlier studies had also established an association between a patient's distance from the hospital and the likelihood of



Fig. 2. ROC curves of the prediction models on the validation data set. Figure legend: The long-solid line in each figure is the ROC curve of the coefficient of each prediction model. The blue area indicates the area under the ROC curve (AUC).



Fig. 3. Important predictors for each prediction model. Figure legend: The sizes of the circles represent the relative importance. The different colors of circles represent feature importance in different models.

patient no-shows (Daggy et al., 2010; DeFife et al., 2010; Miller et al., 2015; Zhou et al., 2018), and our study results further substantiated this as patient's distance from the hospital was found to be a critical predictor of patient no-show behavior. The reason patient's distance from the hospital is an important predictor may be because that the online outpatient appointment system overcomes the limitation of space and improves the coverage of outpatient services. In addition, our research also divulged new critical predictors, such as appointment doctor type and online doctor rating, as potentially related to patient no-show behavior in the online outpatient appointment system in an online environment, appointment doctor type and online doctor rating, as important sources of doctor personal reputation and patients' perceived quality (Huang and Hanauer, 2016), are also important predictors of patient no-shows.

The integration of both online and offline features in the prediction models will enable a better digital healthcare system, an important transition necessary for transforming processes and achieving costeffectiveness (Duggal et al., 2018; Khan and Sakamura, 2015). A machine learning algorithm with good predictive performance will provide tremendous value to clinics and healthcare practitioners in delivering targeted interventions and achieving efficient planning. As online outpatient appointment systems are used by more and more patients, the construction of predictive models not only needs to consider the data in outpatient appointment systems but also the impact of online health information on patient no-show behavior. Our study innovatively uses online doctor rating as a predictor variable and achieves the integration of online data and outpatient appointment data. The results also confirmed the potential of using data from multiple sources to predict patient no-shows.

5.2. Contributions

This study had the following practical contributions. First, we successfully used machine learning algorithms to design prediction models for patient no-show behavior in online outpatient appointments based on online outpatient appointments data and online doctor reputation data. Our results demonstrate the potential of developing effective predictive models using large amounts of data in outpatient online appointment systems. Since the boosting, random forest and bagging models accurately predicted patient no-show behavior, the prediction results can provide a reference for hospitals to formulate a reasonable outpatient appointment and treatment system. For example, hospitals can adjust parameters such as the number of patients absent from appointments given by the patient no-show prediction models, the number of appointment sources on the day and in the future to optimize appointment scheduling and formulate a reasonable overbooking policy to reduce the impact of patient no-shows on hospital efficiency. In addition, to reduce patients' no-show behavior, which may be caused by temporal

and distance factors, hospitals can adopt effective notification measures to reduce the probability of patients' no-show behavior. They also can provide different numbers of appointment sources on different weekdays to reduce the adverse consequences of patients' no-show behavior in online outpatient appointment systems.

This study also has the following theoretical contributions. Different from the traditional outpatient appointment environment, the online outpatient appointment has its unique background. In this study, we designed prediction models for patient no-show behavior based on online outpatient appointment systems, and considered the impact of online health information on patient no-show behavior. Our study also confirmed the potential of using online health information to predict patient no-show behavior. This study not only helps us to deepen the understanding of patient no-show behavior but also enriches and supplements the existing knowledge on prediction models of patient noshow behavior in online outpatient appointments.

6. Limitations and conclusion

This study has the following limitations. First, the data used in this study are provided by a specific hospital, thus the generalizability of the research results may be limited. Future research can extend the data source to other hospitals to cross-validate our results. Second, the predictor variables we used include the patient's gender, which may raise ethical concerns. However, only anonymized data without specific identity-revealing personal information were used in our data analysis. Third, online doctor rating is the result of patient participation, and may have potential self-selection bias. Fourth, only 15 variables were included as predictors in the construction of the prediction models of this study. Hence, the impacts of predictors not included in our analysis were not controlled. We will consider adding more predictors in future studies to improve the performance of the prediction models.

Since the phenomenon of patient no-show behavior in online outpatient appointments is becoming more serious, it is necessary for hospitals to make effective predictions on patient no-show behavior. Using machine learning algorithms, this study confirms the possibility of using a large amount of online outpatient appointment data to build predictive models for patient no-show behavior. The prediction models can assist hospitals in optimizing their outpatient appointment systems by predicting the potential patient no-show behavior, making flexible and reasonable adjustments to outpatient appointment systems to improve their overall efficiency.

Declaration of competing interest

The authors declare that there are no conflicts of interest.

Acknowledgements

This work is supported by the National Natural Science Foundation Program of China [No. 71971092], [No. 71671073] and [71810107003]. The authors thank cooperative medical institutions for providing the data used for this study.

Appendix A. Supplementary data

Supplementary data to this article can be found online at https://doi.org/10.1016/j.dsm.2021.06.002.

References

Ahmed, T., Lucas, H., Khan, A.S., Islam, R., Bhuiya, A., Iqbal, M., 2014. eHealth and mHealth initiatives in Bangladesh: a scoping study. BMC Health Serv. Res. 14 (1), 260.

- Aladeemy, M., Adwan, L., Booth, A., Khasawneh, M.T., Poranki, S., 2020. New feature selection methods based on opposition-based learning and self-adaptive cohort intelligence for predicting patient no-shows. Appl. Soft Comput. 86 (Jan.), 105866.
- Alaeddini, A., Yang, K., Reeves, P., Reddy, C.K., 2015. A hybrid prediction model for noshows and cancellations of outpatient appointments. IIE Trans. Healthc. Syst. Eng. 5 (1), 14–32.
- Alves, M.A., Castro, G.Z., Oliveira, B.A.S., Ferreira, L.A., Ramírez, J.A., Silva, R., Guimarães, F.G., 2021. Explaining machine learning based diagnosis of COVID-19 from routine blood tests with decision trees and criteria graphs. Comput. Biol. Med. 132 (May.), 104335.
- Bardhan, I., Oh, J.-H., Zheng, Z., Kirksey, K., 2014. Predictive analytics for readmission of patients with congestive heart failure. Inf. Syst. Res. 26 (1), 19–39.
- Cao, W., Wan, Y., Tu, H., Shang, F., Liu, D., Tan, Z., Sun, C., Ye, Q., Xu, Y., 2011. A webbased appointment system to reduce waiting for outpatients: a retrospective study. BMC Health Serv. Res. 11, 318.
- Chang, J.T., Sewell, J.L., Day, L.W., 2015. Prevalence and predictors of patient no-shows to outpatient endoscopic procedures scheduled with anesthesia. BMC Gastroenterol. 15 (1), 123.
- Colaci, D., Chaudhri, S., Vasan, A., 2016. mHealth interventions in low-income countries to address maternal health: a systematic review. Annals of Global Health, Current Topics in Global Health 82 (5), 922–935.
- Cronin, P.R., DeCoste, L., Kimball, A.B., 2013. A multivariate analysis of dermatology missed appointment predictors. Jama Dermatology 149 (12), 1435–1437.
- Daggy, J., Lawley, M., Willis, D., Thayer, D., Suelzer, C., DeLaurentis, P.-C., Turkcan, A., Chakraborty, S., Sands, L., 2010. Using no-show modeling to improve clinic performance. Health Inf. J. 16 (4), 246–259.
- Dantas, L.F., Fleck, J.L., Cyrino Oliveira, F.L., Hamacher, S., 2018. No-shows in appointment scheduling – a systematic literature review. Health Pol. 122 (4), 412–421.
- Dantas, L.F., Hamacher, S., Cyrino Oliveira, F.L., Barbosa, S.D.J., Viegas, F., 2019. Predicting patient no-show behavior: a study in a bariatric clinic. Obes. Surg. 29 (2), 40–47.
- DeFife, J.A., Conklin, C.Z., Smith, J.M., Poole, J., 2010. Psychotherapy appointment noshows: rates and reasons. Psychother. Theor. Res. Pract. Train. 47 (3), 413–417.
- Ding, X., Gellad, Z.F., Mather, C., Barth, P., Poon, E.G., Newman, M., Goldstein, B.A., 2018. Designing risk prediction models for ambulatory no-shows across different specialties and clinics. J. Am. Med. Inf. Assoc. 25 (8), 924–930.
- Distelhorst, K., Claussen, R., Dion, K., Bena, J.F., Morrison, S.L., Walker, D., Tai, H.-L., Albert, N.M., 2018. Factors associated with adherence to 14-day office appointments after heart failure discharge. J. Card. Fail. 24 (6), 407–411.
- Duggal, R., Brindle, I., Bagenal, J., 2018. Digital healthcare: regulating the revolution. BMJ 360, k6.
- Fiorillo, C.E., Hughes, A.L., I-Chen, C., Westgate, P.M., Gal, T.J., Bush, M.L., Comer, B.T., 2018. Factors associated with patient no-show rates in an academic otolaryngology practice. Laryngoscope 128 (3), 626–631.
- Fu, H., Zhu, Y., Wang, Y., Liu, Z., Zhang, J., Xie, H., Fu, Q., Dai, B., Ye, D., Xu, J., 2018. Identification and validation of stromal immunotype predict survival and benefit from adjuvant chemotherapy in patients with muscle-invasive bladder cancer. Clin. Canc. Res. 24 (13), 3069–3078.
- Gao, Y., Cai, G.-Y., Fang, W., Li, H.-Y., Wang, S.-Y., Chen, L., Yu, Y., Liu, D., Xu, S., Cui, P.-F., Zeng, S.-Q., Feng, X.-X., Yu, R.-D., Wang, Y., Yuan, Y., Jiao, X.-F., Chi, J.-H., Liu, J.-H., Li, R.-Y., Zheng, X., Song, C.-Y., Jin, N., Gong, W.-J., Liu, X.-Y., Huang, L., Tian, X., Li, L., Xing, H., Ma, D., Li, C.-R., Ye, F., Gao, Q.-L., 2020. Machine learning based early warning system enables accurate mortality risk prediction for COVID-19. Nat. Commun. 11 (1), 5033.
- Goldstein, B.A., Navar, A.M., Pencina, M.J., Ioannidis, J.P.A., 2017. Opportunities and challenges in developing risk prediction models with electronic health records data: a systematic review. J. Am. Med. Inf. Assoc. 24 (1), 198–208.
- Hao, H., 2015. The development of online doctor reviews in China: an analysis of the largest online doctor review website in China. J. Med. Internet Res. 17 (6), e134.
- Hooshangi-Tabrizi, P., Contreras, I., Bhuiyan, N., Batist, G., 2020. Improving patient-care services at an oncology clinic using a flexible and adaptive scheduling procedure. Expert Syst. Appl. 150 (Jul.), 113267.
- Hoque, M.R., 2016. An empirical study of mHealth adoption in a developing country: the moderating effect of gender concern. BMC Med. Inf. Decis. Making 16 (1), 51.
- Huang, Y., Hanauer, D.A., 2014. Patient No-show predictive model development using multiple data sources for an effective overbooking approach. Appl. Clin. Inf. 5 (3), 836–860.
- Huang, Y., Zuniga, P., 2014. Effective cancellation policy to reduce the negative impact of patient no-show. J. Oper. Res. Soc. 65 (5), 605–615.
- Huang, Y.-L., Hanauer, D.A., 2016. Time dependent patient no-show predictive modelling development. Int. J. Health Care Qual. Assur. 29 (4), 475–488.
- Jain, S., Chou, C.L., 2000. Use of an orientation clinic to reduce failed new patient appointments in primary care. J. Gen. Intern. Med. 15 (12), 878–880.
- Kamsu-Foguem, B., Foguem, C., 2014. Telemedicine and mobile health with integrative medicine in developing countries. Health Policy and Technology 3 (4), 264–271.
- Khan, M.F.F., Sakamura, K., 2015. Fine-grained access control to medical records in digital healthcare enterprises. In: 2015 International Symposium on Networks, Computers and Communications (ISNCC). Presented at the 2015 International Symposium on Networks. Computers and Communications (ISNCC), pp. 1–6.
- Kheirkhah, P., Feng, Q., Travis, L.M., Tavakoli-Tabasi, S., Sharafkhaneh, A., 2016. Prevalence, predictors and economic consequences of no-shows. BMC Health Serv. Res. 16 (1), 13.
- Kuiper, A., de Mast, J., Mandjes, M., 2021. The problem of appointment scheduling in outpatient clinics: a multiple case study of clinical practice. Omega 98 (Jan.), 102122.

Ahmadi-Javid, A., Jalali, Z., Klassen, K.J., 2017. Outpatient appointment systems in healthcare: a review of optimization studies. Eur. J. Oper. Res. 258 (1), 3–34.

G. Fan et al.

Kurasawa, H., Hayashi, K., Fujino, A., Takasugi, K., Haga, T., Waki, K., Noguchi, T., Ohe, K., 2016. Machine-learning-based prediction of a missed scheduled clinical appointment by patients with diabetes. J. Diabetes Sci. Technol. 10 (3), 730–736.

- Lee, S.J., Heim, G.R., Sriskandarajah, C., Zhu, Y., 2018. Outpatient appointment block scheduling under patient heterogeneity and patient no-shows. Prod. Oper. Manag. 27 (1), 28–48.
- Lee, V.J., Earnest, A., Chen, M.I., Krishnan, B., 2005. Predictors of failed attendances in a multi-specialty outpatient centre using electronic databases. BMC Health Serv. Res. 5 (1), 51.
- Lehmann, T.N.O., Aebi, A., Lehmann, D., Balandraux Olivet, M., Stalder, H., 2007. Missed appointments at a Swiss university outpatient clinic. Publ. Health 121 (10), 790–799.
- Lekham, L.A., Wang, Y., Hey, E., Lam, S.S., Khasawneh, M.T., 2020. A multi-stage predictive model for missed appointments at outpatient primary care settings serving rural areas. IISE Transactions on Healthcare Systems Engineering 11 (2), 79–94.
- Lenzi, H., Ben, A.J., Stein, A.T., 2019. Development and validation of a patient no-show predictive model at a primary care setting in Southern Brazil. PloS One 14 (4), e0214869.
- Liang, W., Liang, H., Ou, L., Chen, B., Chen, A., Li, C., Li, Y., Guan, W., Sang, L., Lu, J., Xu, Y., Chen, G., Guo, H., Guo, J., Chen, Z., Zhao, Y., Li, S., Zhang, N., Zhong, N., He, J., 2020. Development and validation of a clinical risk score to predict the occurrence of critical illness in hospitalized patients with COVID-19. Jama Internal Medicine 180 (8), 1081–1089.
- Liu, N., 2016. Optimal choice for appointment scheduling window under patient no-show behavior. Prod. Oper. Manag. 25 (1), 128–142.
- Liu, Y., 2009. Reforming China's health care: for the people, by the people? Lancet 373 (9660), 281–283.
- Mey, Y.S., Sankaranarayanan, S., 2013. Near field communication based patient Appointment. 2013 International Conference on Cloud & Ubiquitous Computing & Emerging Technologies. IEEE, Pune, India, pp. 98–103.
- Miller, A.J., Chae, E., Peterson, E., Ko, A.B., 2015. Predictors of repeated "no-showing" to clinic appointments. Am. J. Otolaryngol. 36 (3), 411–414.
- Mold, F., De Lusignan, S., 2015. Patients' online access to their primary care electronic health records and linked online services: implications for research and practice. J. Personalized Med. 5 (4), 452–469.
- Nuti A., L., Lawley, M., Turkcan, A., Tian, Z., Zhang, L., Chang, K., Willis R., D., Sands P., L., 2012. No-shows to primary care appointments: subsequent acute care utilization among diabetic patients. BMC Health Serv. Res. 12 (1), 304.
- Pan, X., Geng, N., Xie, X., 2021. Appointment scheduling and real-time sequencing strategies for patient unpunctuality. Eur. J. Oper. Res. 295 (1), 246–260.
- Peng, Y., Erdem, E., Shi, J., Masek, C., Woodbridge, P., 2016. Large-scale assessment of missed opportunity risks in a complex hospital setting. Inf. Health Soc. Care 41 (2), 112–127.
- Polikar, R., 2012. Ensemble learning. In: Zhang, C., Ma, Y. (Eds.), Ensemble Machine Learning. Springer, Boston, pp. 1–34.
- Rochefort, C.M., Verma, A.D., Eguale, T., Lee, T.C., Buckeridge, D.L., 2015. A novel method of adverse event detection can accurately identify venous thromboembolisms

(VTEs) from narrative electronic health record data. J. Am. Med. Inf. Assoc. 22 (1), 155–165.

- Romano, R., Baum, N., 2014. Reputation management. J. Med. Pract. Manag. 29 (6), 369–372.
- Rosenbaum, J.I., Mieloszyk, R.J., Hall, C.S., Hippe, D.S., Gunn, M.L., Bhargava, P., 2018. Understanding why patients no-show: observations of 2.9 million outpatient imaging visits over 16 years. J. Am. Coll. Radiol. 15 (7), 944–950.
- Simsek, S., Dag, A., Tiahrt, T., Oztekin, A., 2021. A Bayesian Belief Network-based probabilistic mechanism to determine patient no-show risk categories. Omega 100 (Apr.), 102296.
- Srinivas, S., 2020. A machine learning-based approach for predicting patient punctuality in ambulatory care centers. Int. J. Environ. Res. Publ. Health 17 (10), 3703.

Srinivas, S., Ravindran, A.R., 2018. Optimizing outpatient appointment system using machine learning algorithms and scheduling rules: a prescriptive analytics framework. Expert Syst. Appl. 102 (Jul.), 245–261.

- Srinivas, S., Ravindran, A.R., 2020. Designing schedule configuration of a hybrid appointment system for a two-stage outpatient clinic with multiple servers. Health Care Manag. Sci. 23 (Sep.), 360–386.
- Srinivas, S., Salah, H., 2021. Consultation length and no-show prediction for improving appointment scheduling efficiency at a cardiology clinic: a data analytics approach. Int. J. Med. Inf. 145 (Jan.), 104290.
- Tabak, Y.P., Sun, X., Nunez, C.M., Johannes, R.S., 2014. Using electronic health record data to develop inpatient mortality predictive model: acute Laboratory Risk of Mortality Score (ALaRMS). J. Am. Med. Inf. Assoc. 21 (3), 455–463.

Torgo, L., Torgo, M.L., 2013. Package DMwR. In: Comprehensive R Archive Network.

- Whiting, P.S., Greenberg, S.E., Thakore, R.V., Alamanda, V.K., Ehrenfeld, J.M., Obremskey, W.T., Jahangir, A., Sethi, M.K., 2015. What factors influence follow-up in orthopedic trauma surgery? Arch. Orthop. Trauma Surg. 135 (3), 321–327.
- Wu, H., Lu, N., 2017. Online written consultation, telephone consultation and offline appointment: an examination of the channel effect in online health communities. Int. J. Med. Inf. 107 (Nov.), 107–119.
- Xie, X., Zhou, W., Lin, L., Fan, S., Lin, F., Wang, L., Guo, T., Ma, C., Zhang, J., He, Y., Chen, Y., 2017. Internet hospitals in China: cross-sectional survey. J. Med. Internet Res. 19 (7), e239.
- Ye, Q., Deng, Z., Chen, Y., Liao, J., Li, G., Lu, Y., 2019. How resource scarcity and accessibility affect patients' usage of mobile health in China: resource competition perspective. JMIR Mhealth Uhealth 7 (8), e13491.
- Yip, W., Hsiao, W., 2014. Harnessing the privatisation of China's fragmented health-care delivery. Lancet 384 (9945), 805–818.
- Yip, W., Hsiao, W.C., 2008. The Chinese health system at a crossroads. Health Aff. 27 (2), 460–468.
- Zhang, M., Zhang, C., Sun, Q., Cai, Q., Yang, H., Zhang, Y., 2014. Questionnaire survey about use of an online appointment booking system in one large tertiary public hospital outpatient service center in China. BMC Med. Inf. Decis. Making 14 (1), 49.
- Zhou, Y., Dong, D., Jiang, W., 2018. Influence factors of patient no show in a outpatient department. IOP Conf. Ser. Mater. Sci. Eng. 439 (3), 032047.