

Research and Applications

A machine learning–based clinical decision support system to identify prescriptions with a high risk of medication error

Jennifer Corny ¹, Asok Rajkumar,¹ Olivier Martin,² Xavier Dode,^{3,4}
Jean-Patrick Lajonchère,⁵ Olivier Billuart,⁶ Yvonnick Bézie,¹ and Anne Buronfosse⁶

¹Pharmacy Department, Groupe Hospitalier Paris Saint Joseph, Paris, France, ²Lumio Medical, Paris, France, ³Centre National Hospitalier d'Information sur le Médicament, Paris, France, ⁴Pharmacy Department, Hospices Civils de Lyon University Hospital, Lyon, France, ⁵Groupe Hospitalier Paris Saint Joseph, Paris, France, ⁶Medical Information Department, Groupe Hospitalier Paris Saint Joseph, Paris, France

Corresponding author: Jennifer Corny, PharmD, Pharmacy Department, Groupe Hospitalier Paris Saint Joseph, 185 rue Raymond Losserand, 75014 Paris, France; jcorny@hpsj.fr

Received 3 April 2020; Revised 10 June 2020; Editorial Decision 19 June 2020; Accepted 30 June 2020

ABSTRACT

Objective: To improve patient safety and clinical outcomes by reducing the risk of prescribing errors, we tested the accuracy of a hybrid clinical decision support system in prioritizing prescription checks.

Materials and Methods: Data from electronic health records were collated over a period of 18 months. Inferred scores at a patient level (probability of a patient's set of active orders to require a pharmacist review) were calculated using a hybrid approach (machine learning and a rule-based expert system). A clinical pharmacist analyzed randomly selected prescription orders over a 2-week period to corroborate our findings. Predicted scores were compared with the pharmacist's review using the area under the receiving-operating characteristic curve and area under the precision-recall curve. These metrics were compared with existing tools: computerized alerts generated by a clinical decision support (CDS) system and a literature-based multicriteria query prioritization technique. Data from 10 716 individual patients (133 179 prescription orders) were used to train the algorithm on the basis of 25 features in a development dataset.

Results: While the pharmacist analyzed 412 individual patients (3364 prescription orders) in an independent validation dataset, the areas under the receiving-operating characteristic and precision-recall curves of our digital system were 0.81 and 0.75, respectively, thus demonstrating greater accuracy than the CDS system (0.65 and 0.56, respectively) and multicriteria query techniques (0.68 and 0.56, respectively).

Discussion: Our innovative digital tool was notably more accurate than existing techniques (CDS system and multicriteria query) at intercepting potential prescription errors.

Conclusions: By primarily targeting high-risk patients, this novel hybrid decision support system improved the accuracy and reliability of prescription checks in a hospital setting.

Key words: supervised machine learning, electronic prescribing, clinical pharmacy information systems, medication errors, decision support systems, clinical

INTRODUCTION

Medical errors are a major public health problem and a leading cause of mortality. With some 250 000 deaths per year in the United States, medical errors now rank after heart disease and cancer as the third leading cause of death.¹ Even back in 1999, the Institute of Medicine highlighted the need for technologies to prevent the estimated 44 000 to 98 000 annual deaths resulting from medical errors.² The problem is global, and the findings from the United States are readily supported by data from other countries that have also reported substantial rates of health care–related adverse events.^{3–7} During hospitalization, the majority of adverse events are attributed to invasive procedures, hospital-acquired infections, and health products (drugs and medical devices).⁷ While at least 30% of adverse events are likely easily preventable, one study showed that adverse events associated with negligence (including medical errors) were more likely to be associated with injury to patients than other adverse events.⁸

Improving patient safety by reducing medication errors has therefore become a top priority. Although mistakes can occur at any stage of the medical treatment process, prescribing and administration of drugs are the most frequent sources of medication errors.^{9,10} Numerous tools have been developed in an attempt to improve point-of-care prescription processes, including electronic prescribing—termed computerized prescriber order entry (CPOE)—and computerized clinical decision support (CDS) alert systems. However, CPOE is known to generate other types of prescribing errors such as the wrong drug being prescribed,¹¹ and CDS systems are known to generate numerous unnecessary alerts, leading to alert fatigue (desensitization) and subsequent inefficiency.^{12,13}

Medication review by clinical pharmacists is currently the gold-standard method of verification and is increasingly recognized as a critical step in the prevention of potential adverse drug events.^{14,15} Review can trigger pharmaceutical interventions, which are actions implemented in response to—or to prevent—a drug-related problem in an individual patient. However, the process is time-consuming and, as with any human process, not always reproducible. In addition, interventions of this type need to be targeted to ensure that patients who have the greatest risk of errors in their prescription orders are given priority. We recently showed that polypharmacy, patient age, and impaired renal function were associated with more frequent drug-related problems, and consequently, more frequent pharmaceutical interventions.¹⁶ However, neither individual predictor strategies nor multivariate model-based strategies have reliably demonstrated an ability to detect high-risk patients.^{17,18} While these models have marginally improved the efficacy of medication review by targeting patients at the greatest risk, other parameters—including clinical presentation or progression of the patient, laboratory findings, drug regimens, and drug interactions—also need to be taken into consideration. The increasingly widespread availability of electronic health records and the development of big data analytics are currently paving the way for the use of artificial intelligence, which relies on sophisticated algorithms with the capacity to analyze vast quantities of data to make medication review more effective and thus help pharmacists predict or intercept drug-related problems, and therefore make potential medication errors more accurately.¹⁹

In this context, we tested the accuracy of Lumio Medication (developed by Lumio Medical, Paris, France), a hybrid artificial intelligence decision support system—combining machine learning and a rule-based expert system—in a typical hospital setting. This system

is making prediction at the patient level, rather than through predictions about individual prescription orders.

MATERIALS AND METHODS

Setting

The study was conducted in a large, private, nonprofit hospital (592 beds) in Paris that provides both surgical and medical activities. With the exception of neonatology and intensive care units that rely on a dedicated software program, all patient files are typically digitized and recorded using DxCare medical software (Dedalus, Le Plessis-Robinson, France). Prescription orders and medication reviews by clinical pharmacists are managed with the same software, which includes all medical and nursing notes, laboratory results, and vital signs. All drugs can be prescribed (formulary or nonformulary). The software also includes a CDS system and is interfaced with a national drug database. However, given the associated workload, systematic comprehensive medication reviews are typically restricted to certain medical wards—namely cardiology, endocrinology, internal medicine, and rheumatology—where patients are considered to have a high risk of exposure to prescription errors.

Datasets

Data collection

Data collected over an 18-month period from January 2017 to August 2018 were extracted directly from the electronic health records by the hospital Medical Information Department: medication order data, laboratory reports, demographics, medical history, and vital signs. Together these data comprise the development dataset used to train the algorithm.

Throughout the study, all data were stored on a secure local server. Patient identities were pseudonymized: each patient was allocated an identification number that was saved on an identification table on the server.

Prediction target-pharmacist interventions

All the prescription orders included in the development dataset and used to train the algorithm (labeled data) had already undergone at least 1 medication order review. Details of pharmacist interventions were routinely affixed to each medication order as a comment. Each intervention was categorized by a clinical pharmacist according to the ACT-IP classification of the French Society of Clinical Pharmacy.²⁰

Score design

Input features

A large number of datasets were compiled to provide the most comprehensive overview of the context of each prescription order and the patient's medical condition, and thus mitigate any possible bias in the selection of features. Each feature is commonly used for medication review in routine clinical pharmacy practice, and corresponded to a set of individual patient data: laboratory reports (eg, renal function, serum potassium, international normalized ratio), demographics (sex, age), medical history (information on allergies extracted from free text comment fields), and physiological data (vital signs, such as weight, heart rate, and arterial pressure).

For each prescription order, a specific set of rule-based alerts relating to the medication (eg, dosage, frequency, route) were used. The number of times each type of alert was raised represented a

categorical feature (discrete number). Alerts were either extracted from the French online drug database Th  riaque (edited by Centre National Hospitalier d'Information sur le M  dicament),²¹ or built specifically by the hospital pharmacy team on the basis of published literature, and avoiding any overlap with the drug database. For example, a prescription order of a concentration >4 g/L of potassium chloride was considered as inappropriate, and a rule was built to address this rare but potentially harmful prescribing error.

Classifier architecture

We sought to generate a score representing the risk for a given patient's prescription order to contain at least 1 drug-related problem (considered as a medication error), thus answering the question: "Does this patient's set of active orders require a pharmacist review?" To this end, we trained a binary classifier to identify patients who were likely to have at least 1 drug-related error in their prescription order. The classifier selected was derived from LightGBM, a gradient-boosting framework based on decision tree algorithms and developed by Microsoft (Redmond, WA). The 2 types of data were combined as inputs: patient-related data and prescription-related alerts.

For each prescription order in the development dataset, the labeling was therefore binary: 1 = a pharmaceutical intervention was carried out; 0 = no pharmaceutical intervention.

Preprocessing

The binary classifier used 25 features engineered from heterogeneous inputs: numerical quantities, date/time objects, categorical values, and natural language open text fields. During preprocessing, all numerical features were calibrated for outliers, standardized and imputed, while categorical features were encoded, using scikit-learn and scikit-learn machine learning-compatible libraries.

Testing protocol

We tested the performance of the tool on a separate validation dataset (independent of the one used for model development). The accuracy of the algorithm was assessed and compared with the prescription orders reviewed by pharmacists and also with classic techniques: a CDS alert system and a multicriteria query approach.

Methodology

Over a 2-week period, a fully trained clinical pharmacist routinely analyzed randomly selected patient prescription orders on all wards and made a note of the interventions that followed. The selection of prescription orders came from an automatic daily extraction from the medical software documenting all patients who had at least 1 drug prescription. The pharmacist reviewed as many patients as possible over the 2-week period.

The data scientists were blinded to the actual pharmaceutical interventions. Data relating to these prescription orders were then used as inputs for the algorithm, which was thus tested on all wards. All predicted scores (a continuous variable: probability for a prescription order to contain errors) were then compared with the binary score: 1 = a pharmaceutical intervention was carried out during medication review; 0 = no pharmaceutical intervention.

For drug-related problems that were not intercepted by the tool (false negatives), a group of 2 physicians and 2 pharmacists ranked the level of severity (from 1 [minor] to 4 [life-threatening]). To determine this risk, the patient's file was reviewed and the potential immediate or midterm harm was assessed.

Metrics

Prior to the start of the study, several widely acknowledged metrics were selected to assess the algorithm's capacity to predict the risk of an individual patient prescription to contain at least 1 drug-related problem (binary classification). The area under the receiver-operating characteristic curve (AUROC) does not rely on prevalence and is therefore the most widely used metric for model comparisons in statistics. The receiver-operating characteristic (ROC) curve is created by plotting the true positive rate (also known as sensitivity or recall) against the false positive rate at various threshold settings. However, the area under the precision-recall curve (AUCPR) provides a more accurate representation of the impact the algorithm can have on the pharmacist's work,²² as the prevalence of pharmaceutical interventions on prescription orders was 3.6%; it is created by plotting the positive predictive value (also known as precision) against the true positive rate at various threshold settings.

We also calculated the F1 score (ie, the harmonic mean of precision and recall). The 95% confidence interval was calculated for AUCPR and AUROC scores using a bootstrap method after sampling with replacement from 10 000 random samples.²³

Permutation tests (also called exact tests) were computed with 10 000 permutations to assess the statistical significance of the results compared.

Comparators

Two other prioritization processes were used as comparators to evaluate the performance of our algorithm: one based on patient-related data (the multicriteria query), and the other based on medication orders (CDS alert system). The latter makes use of a certified drug database to provide alerts after analyzing patient prescription orders. The CDS alert system thus raised alerts relating to drug interactions, dosage errors, and contraindications for renal insufficiency.

We also used a multicriteria query strategy based on 4 easily available and widely recognized criteria to target high-risk patients (ie, age, renal function [glomerular filtration rate], serum potassium, and international normalized ratio).^{17,18} The score was determined using the following thresholds with the calculation of the number of alerts raised:

- Age
 - >75 years: 1
 - Other: 0
- Estimated glomerular filtration rate (Modification of Diet in Renal Disease) formula:
 - <30 mL/min: 1
 - Other: 0
- Serum potassium level
 - <3 mmol/L: 1
 - >5 mmol/L: 1
 - Other: 0
- International normalized ratio
 - <5: 0
 - Other: 1

For example:

- a 40-year-old patient with no other criteria: score, 0
- an 80-year-old patient with an international normalized ratio of 4: score, 1
- a 76-year-old patient with a glomerular filtration rate of 27 mL/min: score, 2

Ethics approval

Institutional review board approval was obtained from the local ethics committee. Considering the type of study, international review board approval was not required.

RESULTS

Data collection and development

Over an 18-month period, data were collected on 94 720 hospitalizations and a total of 61 611 patients (mean length of stay, 4.1 days; mean age, 69 years; female/male, 49.8%/50.2%), with a mean of 9.4 prescription orders per hospitalization.

During this period, pharmacists reviewed a total of 10 716 individual patient files (133 179 prescription orders), along with each patient's individual data (laboratory findings, demographics, medical history, and vital signs) as part of their daily practice, and these data were used to train the algorithm. This comprised the dataset used for model development.

A pharmaceutical intervention was recommended for 2163 individual patients (20%), meaning 20% of the patients whose all prescription orders were reviewed were at risk of a potential medication error. Based on the ACT-IP classification, the main drug-related problems requiring pharmaceutical intervention were overdosing (33%; dose adaptation suggested), underdosing (16%; dose adaptation suggested), and noncompliance with the drug formulary (16%; replacement by a therapeutic equivalent suggested). Only 3.6% of prescription orders required a pharmaceutical intervention.

In all, 72.7% of the input to the development dataset was from the endocrinology, cardiology, rheumatology, internal medicine, and vascular medicine wards.

Comparative performance of the different techniques

Of the 412 individual patients (3364 prescription orders) that were randomly selected for analysis in the validation dataset, at least 1 pharmaceutical intervention was recommended in 174 (42%) patients. In all, there were 211 pharmaceutical interventions (ie, 6.3% of all prescription orders). In the validation dataset, 64.7% of the input was from the accident and emergency department, and vascular medicine, maternity, orthopedics and endocrinology wards.

Performance of the algorithm vs classic prescription order analysis tools

The accuracy of the hybrid decision support algorithm was compared with the CDS alert system and the multicriteria query (Table 1).

For continuous scoring, such as the output of the hybrid algorithm (a probability), recall and precision were calculated by selecting the classification threshold that maximizes the F1 score.

The decision support algorithm outperformed classic systems in its capacity to both detect patients with a medication error (recall, also known as sensitivity), and to limit the number of false alerts (precision, also called the positive predictive value).

Figures 1 and 2 show the results with regard to AUCPR and AUROC.

Accuracy of medication review using the algorithm

In the independent validation dataset, with the classification threshold that maximizes the F1 score, the Lumio Medication algorithm intercepted 74% of all prescription orders that required pharmacist intervention while also demonstrating 74% precision. Of the remaining 26% prescription orders that required pharmacist intervention (false negatives) that were not intercepted by the algorithm, none were life-threatening. The ensuing F1 score thus showed 15.6% greater accuracy than multicriteria query techniques and 21% greater accuracy than the CDS alert system.

The AUCPR and AUROC scores showed the greater accuracy (+33% and +19%, respectively) of the algorithm vs multicriteria query techniques, and vs the CDS alert system (+33% and +24%, respectively).

Statistical significance

The permutation tests showed the statistically significant superiority of the AUCPR and AUROC scores of the hybrid machine learning system over those of random variables and vs the comparator techniques: $P < .00001$ vs the CDS alert system, and $P = .001$ (AUCPR) and $P = .0152$ (AUROC) vs the multicriteria query technique.

DISCUSSION

This study presents a hybrid machine learning–based decision support system for reviewing the accuracy of medical prescription orders. Our findings confirm that the algorithm outperformed classic systems in its capacity to limit the number of false alerts without overlooking patients with prescription order errors.

CPOE has proven its efficacy in reducing medication errors. A meta-analysis published in 2014 found that compared with paper order entry, CPOE was associated with 50% less medication errors, although new types of medication errors were associated with CPOE.¹¹

CDS systems included in CPOE are a known source of alert fatigue. Ancker et al²⁴ found that the likelihood of alerts being

Table 1. Comparative performance of classic prescription order analysis tools versus the Lumio Medication algorithm in terms of recall, precision, and F1 scores, as well as AUCPR and AUROC

Metric	CDS alert system	Multicriteria query	Lumio Medicationalgorithm
Recall	0.69	0.66	0.74 ^a
Precision	0.54	0.62	0.74 ^a
F1	0.61	0.64	0.74 ^a
AUCPR	0.56 (95% CI, 0.50-0.62; $P < .00001$)	0.56 (95% CI, 0.51-0.61; $P = .001$)	0.75 (95% CI, 0.70-0.80) ^a
AUROC	0.65 (95% CI, 0.61-0.69; $P < .00001$)	0.68 (95% CI, 0.64-0.72; $P = .0152$)	0.81 (95% CI, 0.78-0.84)

AUCPR: area under the precision-recall curve; AUROC: area under receiver-operating characteristic curve; CDS: clinical decision support; CI, confidence interval.

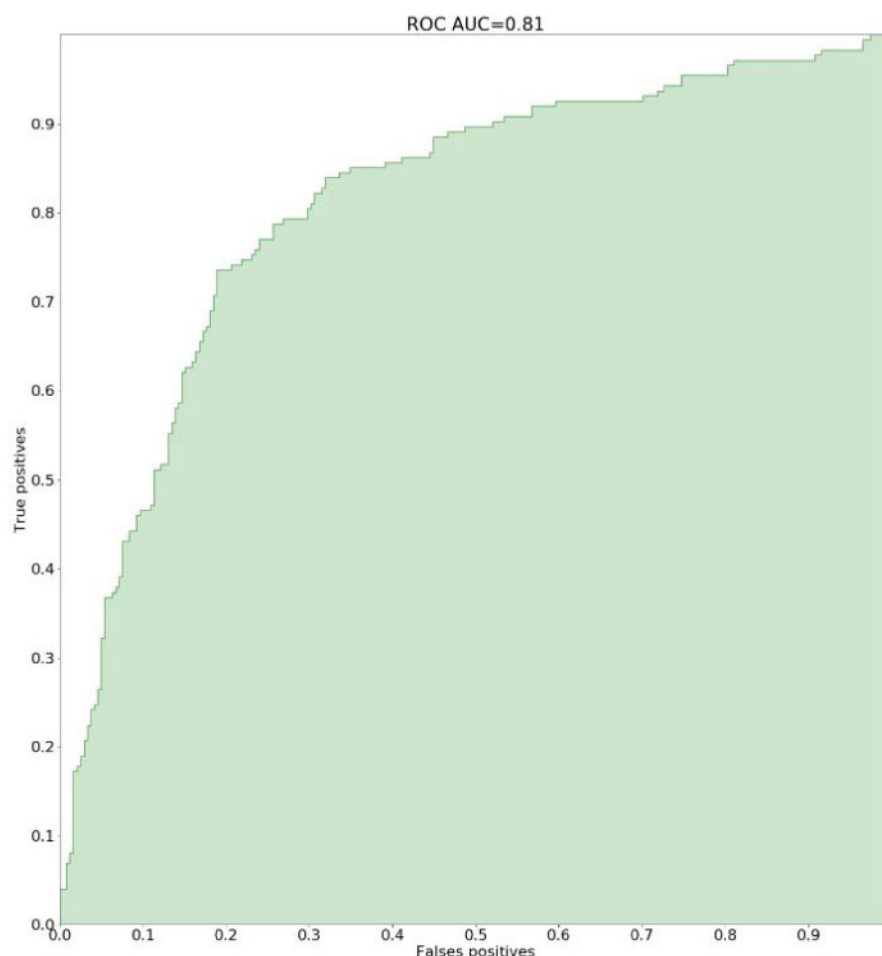


Figure 1. Accuracy of the decision support system: ROC AUC: area under receiver operating characteristic curve.

accepted decreased by 30% for each additional reminder received, and by 10% for each 5% increase in the number of repeated alerts.

Given that current day processes of medication review do not have the capacity to cover all medical prescription orders, the process is in urgent need of improvement. Pharmaceutical interventions are still relatively scarce and therefore prioritization of medication reviews, based on the likelihood of drug-related problems and therefore medication errors, is essential. This study provided compelling evidence of the accuracy of artificial intelligence in identifying those patients with the greatest risk of errors in their prescription orders (ie, prioritizing patients in whom medication review is justified). Earlier studies identified several risk factors for medication errors, including patient age, renal dysfunction, and the number of drugs prescribed, to help pharmacists target interventions more effectively. However, we previously found that these risk factors accounted for only 34% of the variations in the number of pharmaceutical interventions.¹⁶ A multicriteria model-based strategy was also developed to identify patients whose prescription orders presented a high risk of containing errors. This model was based on 11 predictors, of which patient age and the number of drugs on the prescription were the most significant, with a C-statistic of 0.72.¹⁷ Nevertheless, of a total of 303 individual patients, 6 still needed to be reviewed for a drug-related problem to be detected, demonstrating the need for innovative approaches to make this activity more effective. Other studies reported a C-score model to detect only previously identified

adverse events, with interesting results.^{25,26} However, these only focused on selected adverse events and did not consider all potential medication errors.

Our hybrid decision support system combining machine learning with a rule-based expert system was notably more accurate at detecting medication errors compared with other tools described in the literature. Two of 3 individual patient prescription orders reviewed by our tool triggered a pharmaceutical intervention, a figure that compares very favorably with the 20% in our development dataset or approximately 17% in the study by Nguyen et al.¹⁷ The sensitivity of our tool was also significantly higher than that of the CDS alert system or multicriteria query techniques. The hybrid model we have developed uses both knowledge-driven (expert system) and data-driven (machine learning) approaches. It can therefore be expected to overcome the main shortcomings of both these techniques by (1) not overfitting and consequently reproducing the same error patterns that occurred in the development dataset; (2) addressing the issue of certain infrequent though critical medical errors such as the so-called never-events, that is, serious incidents that are wholly preventable, as highlighted by the French Agency for the Safety of Health Products; and (3) reducing the number of false positive alerts typically seen with tools such as CDS alert systems.

This tool can also be easily adjusted by the addition of specific rules to account for noisy or conflicting categories that the algorithm has not yet learnt to deal with.

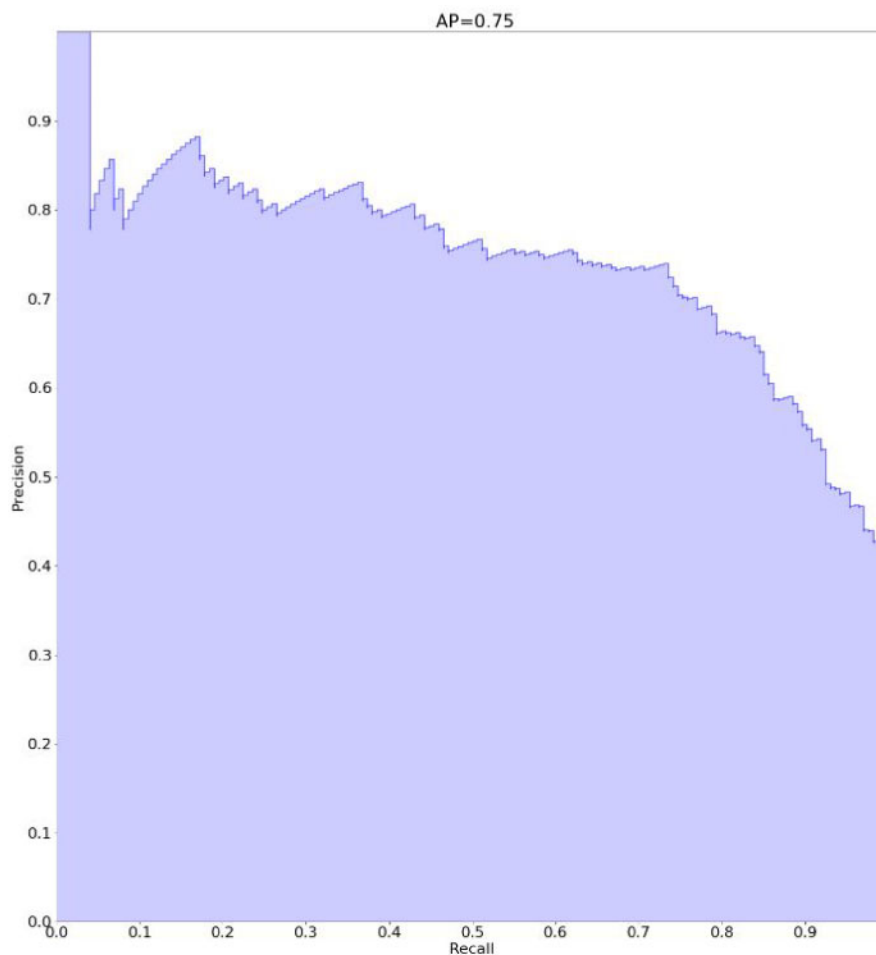


Figure 2. Accuracy of the decision support system: AP: area under the precision–recall curve.

One recently published study presented the results of an outlier detection machine learning–based tool in a real-life setting that exhibited 89% accuracy in terms of the alerts raised.²⁷ However, only 0.4% of the prescription orders generated an alert, whereas in our study, 6.3% of prescription orders were associated with a pharmaceutical intervention. In addition, the authors did not report any data on the sensitivity of this tool; it is therefore legitimate to speculate that because the system focused solely on outlier detection, other common medication errors likely went undetected.

The capacity of our hybrid system to constantly be adjusted as more prescription orders are reviewed by clinical pharmacists and more potential error patterns identified gives it a significant advantage over existing CDS systems.

Our findings are currently limited in scope in as far as the study was conducted in a single hospital setting. In addition, neonatology and intensive care unit patients were not included because they are managed by a different medical software. Consequently, evidence of the accuracy of the algorithm to identify prescription order errors in these units has yet to be demonstrated and our results can therefore not be applied to these patients. Importantly, more pharmaceutical interventions were recommended during the test phase than in the development dataset. There are 2 possible explanations: (1) in the validation dataset, prescription orders were reviewed by an experienced clinical pharmacist, whereas several pharmacists (junior and senior) with different levels of experience were involved in the medi-

cation review over the 18-month data collection and development period; and (2) input with regard to pharmaceutical interventions for the validation dataset and the development dataset came from different wards.

The next step is to deploy our system throughout other hospitals, thus extending the patient population covered. This will also enable us to benefit from the experience of a greater number of clinical pharmacists to confirm our findings. Adjustments are currently being made to the algorithm to integrate unstructured data to further improve the performance of this tool. As an example, while the algorithm used in this study does not yet identify wrong-patient errors, adjustments presently underway will enable it to address potential errors on medical notes (free text) associated with CPOE. Finally, a real-life evaluation is currently being conducted to assess the performance of this tool in daily medication reviews.

CONCLUSION

A hybrid machine learning–based decision support system has been developed to intercept prescription orders with a high risk of containing at least 1 medication error. Given that it is based on machine learning– and rule-based alerts, this decision support system has the advantage of not overfitting errors, of decreasing alert fatigue, and also of addressing infrequent but nevertheless potentially critical

