






# Using Machine Learning to Identify Patients at High Risk of Inappropriate Drug Dosing in Periods with Renal Dysfunction

Benjamin Skov Kaas-Hansen <sup>1-3</sup>, Cristina Leal Rodríguez <sup>2</sup>, Davide Placido <sup>2</sup>, Hans-Christian Thorsen-Meyer <sup>2,4</sup>, Anna Pors Nielsen <sup>2</sup>, Nicolas Dérian<sup>5</sup>, Søren Brunak<sup>2</sup>, Stig Ejdrup Andersen<sup>1</sup>

<sup>1</sup>Clinical Pharmacology Unit, Zealand University Hospital, Roskilde, Denmark; <sup>2</sup>NNF Center for Protein Research, University of Copenhagen, Copenhagen, Denmark; <sup>3</sup>Section for Biostatistics, Department of Public Health, University of Copenhagen, Copenhagen, Denmark; <sup>4</sup>Department of Intensive Care Medicine, Copenhagen University Hospital (Rigshospitalet), Copenhagen, Denmark; <sup>5</sup>Data and Development Support, Region Zealand, Sorø, Denmark

Correspondence: Benjamin Skov Kaas-Hansen, Clinical Pharmacology Unit, Zealand University Hospital, Munkevej 18, Roskilde, 4000, Denmark, Tel +45 60 19 68 02, Email epiben@hey.com

**Purpose:** Dosing of renally cleared drugs in patients with kidney failure often deviates from clinical guidelines, so we sought to elicit predictors of receiving inappropriate doses of renal risk drugs.

**Patients and methods:** We combined data from the Danish National Patient Register and in-hospital data on drug administrations and estimated glomerular filtration rates for admissions between 1 October 2009 and 1 June 2016, from a pool of about 2.6 million persons. We trained artificial neural network and linear logistic ridge regression models to predict the risk of five outcomes ( $>0$ ,  $\geq 1$ ,  $\geq 2$ ,  $\geq 3$  and  $\geq 5$  inappropriate doses daily) with index set 24 hours after admission. We used time-series validation for evaluating discrimination, calibration, clinical utility and explanations.

**Results:** Of 52,451 admissions included, 42,250 (81%) were used for model development. The median age was 77 years; 50% of admissions were of women.  $\geq 5$  drugs were used between admission start and index in 23,124 admissions (44%); the most common drug classes were analgesics, systemic antibacterials, diuretics, antithrombotics, and antacids. The neural network models had better discriminative power (all AUROCs between 0.77 and 0.81) and were better calibrated than their linear counterparts. The main prediction drivers were use of anti-inflammatory, antidiabetic and anti-Parkinson's drugs as well as having a diagnosis of chronic kidney failure. Sex and age affected predictions but slightly.

**Conclusion:** Our models can flag patients at high risk of receiving at least one inappropriate dose daily in a controlled in-silico setting. A prospective clinical study may confirm that this holds in real-life settings and translates into benefits in hard endpoints.

**Keywords:** predictive modelling, kidney failure, machine learning, risk markers, inappropriate drug dosing, renal risk drugs

## Introduction

Renal diseases affect patients' susceptibility to, and modify the effects of many drugs, and they reduce renal clearance exposing patients to higher steady-state concentrations when given standard doses. Kidneys excrete active forms and/or metabolites of many drugs, so renal dysfunction necessitates dose-adjustment of renally cleared drugs with narrow therapeutic indices to prevent adverse events and accidental over-dosing.

Inadequate dose-adjustment of such drugs has been linked to polypharmacy<sup>1,2</sup> and can cause noxious events<sup>3</sup> or accidental over-dosing.<sup>4</sup> Although not a new issue,<sup>5,6</sup> deviating from guidelines is widespread with prevalence estimates up to 70%.<sup>1,2,7-9</sup> Despite large inter-individual variability in clearance and response, dose adjustment for many drugs is crude and based on the estimated glomerular filtration rate (eGFR), for example, halving the dose when eGFR  $<60$  mL/min/1.73 m<sup>2</sup>.

Appropriate alerts in order-entry systems may facilitate rational clinical decision-making,<sup>10,11</sup> and convincing examples have showcased how computerized systems can underpin rational pharmacotherapy.<sup>4,12</sup> However, downsides of extensive computerization of healthcare emerge,<sup>13</sup> alert fatigue<sup>14</sup> is particularly problematic, and strategies and interventions have been proposed to mitigate its negative effects.<sup>15</sup>

At Danish hospitals, prescriptions are mostly dispensed and administered by nurses who record detailed meta-data.<sup>16</sup> Prescriptions are usually made and revised by physicians regularly during clinical rounds, typically in the morning or early afternoon. Electronic decision support is generally immature and neither prescribing physicians nor dispensing nurses are warned if dose-adjustment be advised or even required.

We suspect that the need for dose-adjustment in patients with renal dysfunction often goes unrecognized. Thus, with this paper, we study its predictability to inform clinicians and health-care personnel upfront about which patients with renal dysfunction are at elevated risk of inappropriate drug dosing. To this end, we used and compared predictive modelling methods from classical statistical modelling and machine learning as the richer and more complex models in the latter group may capture more complex relationships than (and thus outperform<sup>17</sup>) those in the former.<sup>18</sup>

## Methods

### Study Design, Patients and Data

We conducted a register-based prediction study with prospective data<sup>19</sup> for patients admitted to 12 public hospitals in two Danish regions comprising about 2.6 million persons (more than half the Danish population). We collected diagnosis data from the Danish National Patient Register, demographic data from the Danish Civil Registration System,<sup>20</sup> as well as medication and biochemical data from electronic patient records. Diagnoses were encoded using the 10th revision of the International Classification of Diseases (ICD-10), drugs with the Anatomical and Therapeutic Chemical classification (ATC).

The units of analysis were inpatient admissions, defined as chains of successive in-hospital visits at most 24 hours apart. We included admissions starting between 1 October 2009 and 1 June 2016, with at least one eGFR measurement  $\leq 30$  during the first 24 hours of admission. We excluded minors (age  $< 18$  years). Admission time uses hour resolution (an admission starting at 9:54 is recorded as starting at 9:00) so to ensure at least 24 hours of observation time before inclusion, index was set at hour of admission + 25 hours. Prior sample-size estimation was foregone.

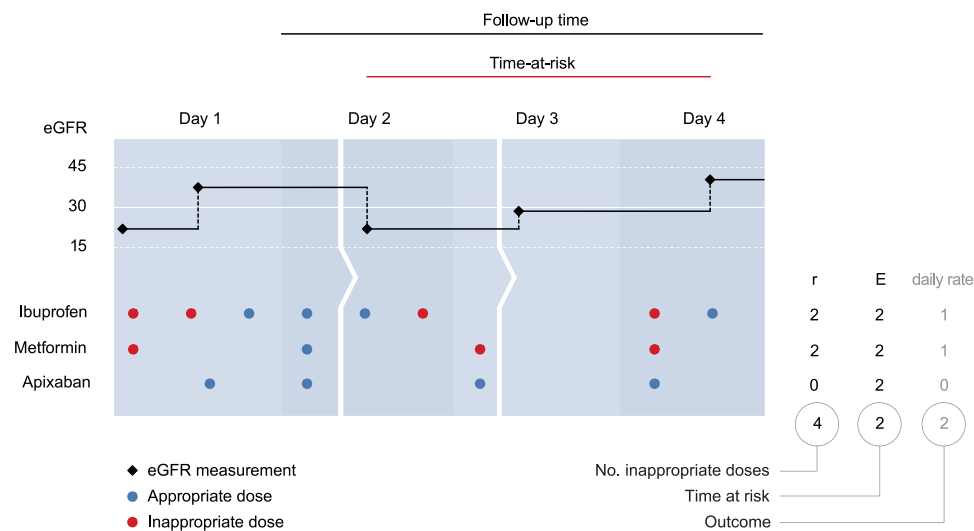
### Outcomes

The outcome variables were based on the daily rate  $= r/E$  of inappropriate doses during follow-up, capped at 30 days.  $r$  is the number of given inappropriate doses of select drugs cleared mainly renally and with narrow therapeutic indices;  $E$  the time-at-risk (Figure 1). To obtain well-defined times-at-risk, we set the eGFR threshold to  $\leq 30$  mL/min/1.73m<sup>2</sup> (unit omitted from here onward) and used the rules in [Supplementary Table S1](#) for counting the number of inappropriate doses, based on the official reference guidelines for Danish physicians (pro.medicin.dk) as of January 2021.

We used two rules, one definitive (maximum daily dose = 0 mg) and one of dose-adjustment (reduced daily dose). Operationalization of the definitive rule is straightforward: if the last eGFR  $\leq 30$ , there should be no administrations until an eGFR  $> 30$  is measured. The dose-adjustment rule is slightly more involved as inappropriate dosing comes in two forms: (a) on a given day, there are more than one eGFR measurements, of which at least one is  $\leq 30$ , and the cumulative daily dose surpasses the threshold in the period(s) between above-threshold measurements, or (b) all eGFR measurements of a given day are  $\leq 30$  and the cumulative daily dose surpasses the threshold.

### Variables and Features

Variables are original data (eg sex and age at admission) and features the results of rendering the variables appropriate as model inputs (eg one-hot-encoded day of admission). Based on clinical and pharmacological experience we hand-picked pertinent variables likely to be informative to the prediction problem and realistically available in the clinical setting. These fall into three categories. Demographic: age at admission (numeric), sex (binary). Clinical: number of distinct drugs (ATC level 5) administered between admission and index (numeric); therapeutic drug classes (ATC level 2) used



**Figure 1** Deriving the outcome variables. This exemplary admission is composed of three successive in-patient visits (ie the patient has been transferred twice represented by the arrows). The admission is eligible because it spans more than 24 hours and an eGFR  $\leq 30$  was measured before index. Here, apixaban was given while the patient's eGFR was  $\leq 30$ , but dose reduction rendered these administrations appropriate.

between admission and index (one-hot-encoded); the Elixhauser score at admission (numeric, AQHR adaptation);<sup>21</sup> ICD-10 chapters of diagnoses recorded in the past five years before admission (one-hot-encoded); record of chronic kidney failure in the past five years before admission (ICD-10 N18\* diagnoses, one-hot-encoded). Contextual: hour of admission (numeric, transformed as  $f(t) = \text{abs}(12 - t)$ ; see [Supplementary Figure S1](#)); weekday of admission (one-hot-encoded); number of admissions in the past 5 years before admission (numeric). In all, there were 98 features.

Missing values, only present for hour of admission and discharge, were imputed by sampling from the empirical distributions of valid values.

## Models and Training

We tried two model architectures (linear logistic ridge regression and artificial neural network) with several binary outcomes defined by increasing thresholds of the daily rate of inappropriate doses ( $>0$ ,  $\geq 1$ ,  $\geq 2$ ,  $\geq 3$  and  $\geq 5$ ). The neural network models were multi-layer perceptrons (MLPs) enabling speedy training and evaluation. We compared these full models with sparser versions (reference models with 9 features known to precipitate renal dysfunction, see [Supplementary Figures S79-S82](#)) for the outcomes  $>0$  and  $\geq 1$  inappropriate daily doses, to assess the added benefit with respect to predictive performance of the full (richer) feature set.

All admissions starting before 1 July 2015 were assigned to the development set (42,250 admissions [81%] of 27,253 patients) and the rest to the independent hold-out test set (10,201 admissions [19%] of 8412 patients). Because admissions constitute the unit of analysis, some patients likely appear in both the development and test sets. Information may leak between the sets,<sup>22</sup> so as a sensitivity analysis, we evaluated the performance also in the subset of test-set patients not in the development set.

We used the multivariate *TPESampler* from *Optuna*<sup>23</sup> to find the best-performing hyperparameters by sampling 100 configurations, each using 5-fold stratified-and-grouped cross-validation, from the following proposal distributions (discrete values in round brackets, bounds of log-uniform distributions in squared): optimizer (Adam, RMSprop), learning rate [ $10^{-6}$ ,  $10^{-1}$ ], activation function (tanh, sigmoid), L2 penalty [ $10^{-6}$ ,  $10^{-2}$ ], number of hidden layers (1, 2, 3, 4), number of nodes per hidden layer (16, 32, 65, 128), batch size (32, 64, 128, 256, 512), class handling (see below). We used the binary cross-entropy loss function throughout.

Only relevant hyperparameters were sampled and we ran *Optuna* on linear and MLP models separately because they have disparate hyperparameter sets. MLP models with more hidden layers and more nodes therein can learn more complex relationships but become prone to overfitting which we countered with early stopping<sup>24</sup> and L2 regularization

(handles collinearity better than L1 regularization).<sup>18,25</sup> The batch size is the number of observations from which the model learns at a time; small batches can give outliers undue influence while full-batch training (batch size = number of units) can become computationally impractical.<sup>22</sup> Class imbalances in binary outcomes can misguide training, so we tested the following remedies: synthetic minority oversampling technique (SMOTE), random over-sampling of minority class, NearMiss, random under-sampling of majority class, class weighting, and none. SMOTE creates a dataset similar to the minority class but of the same size as the majority class;<sup>26</sup> NearMiss downsizes the majority class in a systematic way to retain as much information as possible in fewer data points.<sup>27</sup> Class weighting retains the original data but gives more weight to minority-class observations.

Hyperparameter optimization models trained for maximum 500 epochs with 50-epoch patience on improvement in the validation loss. The final models were trained on the full development set until the loss reached that obtained in the best cross-validation fold for the best configuration.<sup>24</sup>

## Evaluation and Explanation

Discrimination was assessed with receiver operating characteristic (ROC) curves and areas under the ROC curves (AUROC), calibration-in-the-small by plotting decile-binned predicted probabilities against corresponding bin-wise observed event proportions<sup>28</sup> with 95% Jeffrey intervals;<sup>29</sup> results from a perfectly calibrated model fall on the diagonal. We used the decision-curve analytic framework to gauge the models' potential clinical utility.<sup>30,31</sup> In the interest of transparency, all performance figures of the reference models are included in the supplement.

For explanation and scrutiny of prediction drivers, we used the SHAP DeepExplainer yielding one shap value per feature per unit.<sup>32</sup> The shap value for a risk prediction model is the absolute change in risk of a given unit's value for each feature: the cohort-wide mean risk plus the sum of one unit's shap values equals that unit's risk.

## Analysis and Ethics

The full analytical pipeline was built with Snakemake<sup>33</sup> (schematic overview in [Supplementary Figure S2](#)) to facilitate transparency and reproducibility; blinding was impractical and so foregone, but all analytic code is available online (DOI: 10.5281/zenodo.4560078). Univariate distributions were summarized by median (inter-quartile range) and count (proportion), as appropriate. This report adheres to pertinent items in the MINIMAR guideline<sup>34</sup> and TRIPOD statement;<sup>35</sup> the latter is available in the supplement.

All data have been marshalled on Computerome, a secure high-performance Danish computing infrastructure, after obtaining approval from the Danish Patient Safety Authority (3–3013-1723; then competent authority for ethical approval), the Danish Data Protection Agency (DT SUND 2016–48, 2016–50, 2017–57) and the Danish Health Data Authority (FSEID 00003724).

## Results

**Table 1** shows univariate summary statistics of the 52,451 admissions (42,250 + 10,201) of 35,665 patients (27,253 + 8412) included in the study (see [Supplementary Table S2](#) for extended version with all features). Patients in the test sets were similar to those in the development set with some notable exceptions. Fewer had received inappropriate doses, especially in the test-set patients not part of the development set who also had fewer previous admissions.

In the development set, the median age was 77 years (IQR: 67–85) and 20,743 admissions (49%) were of 13,759 women (50%). The median time at risk was 3.5 days (inter-quartile range: 1.7–7.7) and at least one inappropriate dose was given in 3786 admissions (9.0%);  $\geq 1$  inappropriate dose daily was given in 5.3% of admissions and  $\geq 5$  inappropriate doses daily were given in 0.9%. The events-per-feature ratios in the development set varied between 39 (3.786/98, for at least one inappropriate dose) and 3.7 (366/98,  $\geq 5$  inappropriate doses daily). The target drugs most commonly given in inappropriate doses were ibuprofen (M01AE01, 4.1%) and metformin (A10BA02, 3.4%); inappropriate doses of the other target drugs were given in <1% of admissions.

Patients in 4988 admissions (12%) had no admissions in the 5 years before inclusion; 13,960 (33%) had  $\geq 7$  previous admissions. The most common drug classes used between admission and index were analgesics (N02, 37%), systemic antibacterials (J01, 35%), diuretics (C03, 33%) antithrombotics (B01, 28%), and antacids (A02, 25%). Previous

**Table 1** Univariate Summary Statistics of Select Features. Values are Median (Inter-Quartile Range) and Count (Proportion) as Appropriate. *Distinct Patients* and *Distinct Women* Show Counts of Actual Patients (as a Patient Can Contribute More Than One Unit)

| Variate                                 | Development Set<br>(N = 42,250) | Test Set<br>(N = 10,201) | Test Set (Not in Devel. Set)<br>(N = 5980) |
|---|---------------------------------|--------------------------|--|
| Women                                   | 20,743 (49%)                    | 4854 (48%)               | 2940 (49%)                                 |
| Distinct patients                       | 27,253                          | 8412                     | 5341                                       |
| Distinct women                          | 13,759 (50%)                    | 4049 (48%)               | 2629 (49%)                                 |
| Time at risk, days                      | 3.5 (1.7–7.7)                   | 3.5 (1.7–7.2)            | 2.9 (1.5–6.4)                              |
| Inappropriate doses (outcomes)          |                                 |                          |  |
| >0 (at least one)                       | 3786 (9.0%)                     | 1080 (11%)               | 740 (12%)                                  |
| ≥1 daily                                | 2241 (5.3%)                     | 588 (5.8%)               | 333 (5.6%)                                 |
| ≥2 daily                                | 1236 (2.9%)                     | 288 (2.8%)               | 108 (1.8%)                                 |
| ≥3 daily                                | 783 (1.9%)                      | 171 (1.7%)               | 56 (0.9%)                                  |
| ≥5 daily                                | 366 (0.9%)                      | 64 (0.6%)                | 9 (0.2%)                                   |
| Admissions 5 years before admission     |                                 |                          |  |
| None                                    | 4988 (12%)                      | 1082 (11%)               | 1074 (18%)                                 |
| 1–2                                     | 10,100 (24%)                    | 2367 (23%)               | 1873 (31%)                                 |
| 3–4                                     | 7712 (18%)                      | 1919 (19%)               | 1232 (21%)                                 |
| 5–6                                     | 5490 (13%)                      | 1303 (13%)               | 685 (12%)                                  |
| ≥7                                      | 13,960 (33%)                    | 3530 (35%)               | 1116 (19%)                                 |
| Drugs used between admission and index  |                                 |                          |  |
| None                                    | 6165 (15%)                      | 1228 (12%)               | 762 (13%)                                  |
| 1–2                                     | 9111 (22%)                      | 1984 (19%)               | 1254 (21%)                                 |
| 3–4                                     | 8761 (21%)                      | 2078 (20%)               | 1355 (23%)                                 |
| 5–6                                     | 7197 (17%)                      | 1852 (18%)               | 1095 (18%)                                 |
| ≥7                                      | 11,016 (26%)                    | 3059 (30%)               | 1514 (25%)                                 |
| Any diagnosis of chronic kidney failure | 13,470 (32%)                    | 3391 (33%)               | 732 (12%)                                  |
| Top-5 ICD-10 chapters <sup>†</sup>      |                                 |                          |  |
| Cardiovascular (IX)                     | 25,757 (61%)                    | 6392 (63%)               | 3283 (55%)                                 |
| Genitourinary (XIV)                     | 23,025 (55%)                    | 5819 (57%)               | 2306 (39%)                                 |
| Lesions, external causes, etc. (XIX)    | 20,275 (48%)                    | 4749 (47%)               | 2481 (42%)                                 |
| Metabolic-endocrine (IV)                | 19,716 (47%)                    | 5096 (50%)               | 2415 (40%)                                 |
| Symptoms/abnormal findings (XVIII)      | 18,663 (44%)                    | 5711 (56%)               | 2882 (48%)                                 |
| Top-5 drug classes <sup>‡</sup>         |                                 |                          |  |
| Analgesics (N02)                        | 15,740 (37%)                    | 4367 (43%)               | 2506 (42%)                                 |
| Systemic antibacterials (J01)           | 14,719 (35%)                    | 3257 (32%)               | 1938 (32%)                                 |
| Diuretics (C03)                         | 13,966 (33%)                    | 3672 (36%)               | 1951 (33%)                                 |
| Antithrombotics (B01)                   | 11,842 (28%)                    | 3181 (31%)               | 1795 (30%)                                 |
| Antacids (A02)                          | 10,635 (25%)                    | 2776 (27%)               | 1407 (24%)                                 |

**Notes:** <sup>†</sup>ICD-10 chapters (Roman numbering) of diagnoses recorded in the last 5 years before admission. <sup>‡</sup>Drug classes (ATC level 2) administered between admission and index.

**Abbreviations:** ICD-10, 10th version of the international classification of disease; ATC, anatomical therapeutic chemical classification of medicines.

diagnoses were most commonly cardiovascular (chapter IX, 61%), genitourinary (XIV, 55%), related to i.a. lesions and external causes (XIX, 48%), endocrine-metabolic (IV, 47%), and symptoms/abnormal findings (XVIII, 44%).

Table 2 shows the hyperparameters of the best configurations with performance metrics of the final and reference models (see also [Supplementary Figures S3–S16](#)). Generally, multi-layer perceptron (MLP) models performed slightly better than their linear counterparts, all obtaining AUROC's between 0.77 and 0.81 in the test set (ROC curves in [Supplementary Figures S17–S30](#)). The MLP models more consistently showed good calibration in the development set. For daily rates >0, ≥1 and ≥2 both MLP and linear models were very well calibrated in the test set ([Supplementary Figures S31–S44](#)). The full models had better discrimination and were better calibrated than the corresponding reference

**Table 2** Performance Metrics of Final Models and Results of Optuna Hyperparameter Optimization

| Parameter               | Daily Rate >0         |                       | Daily Rate ≥1         |                       | Daily Rate ≥2         |                       | Daily Rate ≥3         |                       | Daily Rate ≥5         |                       |
|-------------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
|                         | Linear                | MLP                   | Linear                | MLP                   | Linear                | MLP                   | Linear                | MLP                   | Linear                | MLP                   |
| AUROC <sup>†</sup>      |                       |                       |                       |                       |                       |                       |                       |                       |                       |                       |
| Development set         | 0.80 (0.71)           | 0.81 (0.73)           | 0.81 (0.72)           | 0.83 (0.75)           | 0.81                  | 0.84                  | 0.82                  | 0.83                  | 0.82                  | 0.83                  |
| Test set                | 0.77 (0.70)           | 0.79 (0.70)           | 0.78 (0.69)           | 0.79 (0.70)           | 0.79                  | 0.79                  | 0.81                  | 0.81                  | 0.78                  | 0.80                  |
| Test set (new patients) | 0.78 (0.68)           | 0.79 (0.68)           | 0.82 (0.70)           | 0.83 (0.71)           | 0.86                  | 0.86                  | 0.89                  | 0.90                  | 0.82                  | 0.79                  |
| Hyperparameters         |                       |                       |                       |                       |                       |                       |                       |                       |                       |                       |
| Batch size              | 512                   | 128                   | 512                   | 32                    | 32                    | 64                    | 256                   | 256                   | 64                    | 64                    |
| Class handling          | Undersample           | SMOTE                 | NearMiss              | NearMiss              | Oversample            | SMOTE                 | Oversample            | NearMiss              | Oversample            | None                  |
| L2 penalty              | $1.28 \times 10^{-6}$ | $1.66 \times 10^{-6}$ | $3.02 \times 10^{-6}$ | $1.43 \times 10^{-6}$ | $4.38 \times 10^{-6}$ | $1.39 \times 10^{-6}$ | $1.43 \times 10^{-6}$ | $1.30 \times 10^{-6}$ | $1.09 \times 10^{-5}$ | $3.94 \times 10^{-6}$ |
| Learning rate           | $1.79 \times 10^{-2}$ | $1.20 \times 10^{-4}$ | $1.92 \times 10^{-2}$ | $3.45 \times 10^{-4}$ | $6.73 \times 10^{-3}$ | $2.71 \times 10^{-4}$ | $3.76 \times 10^{-2}$ | $3.08 \times 10^{-4}$ | $2.11 \times 10^{-2}$ | $4.86 \times 10^{-4}$ |
| Optimizer               | Adam                  | Adam                  | Adam                  | Adam                  | Adam                  | Adam                  | Adam                  | Adam                  | Adam                  | Adam                  |
| Activation function     | —                     | tanh                  | —                     | sigmoid               | —                     | tanh                  | —                     | sigmoid               | —                     | sigmoid               |
| No. hidden layers       | —                     | 3                     | —                     | 1                     | —                     | 1                     | —                     | 1                     | —                     | 2                     |
| Nodes per hidden layer  | —                     | 8                     | —                     | 8                     | —                     | 32                    | —                     | 32                    | —                     | 8                     |

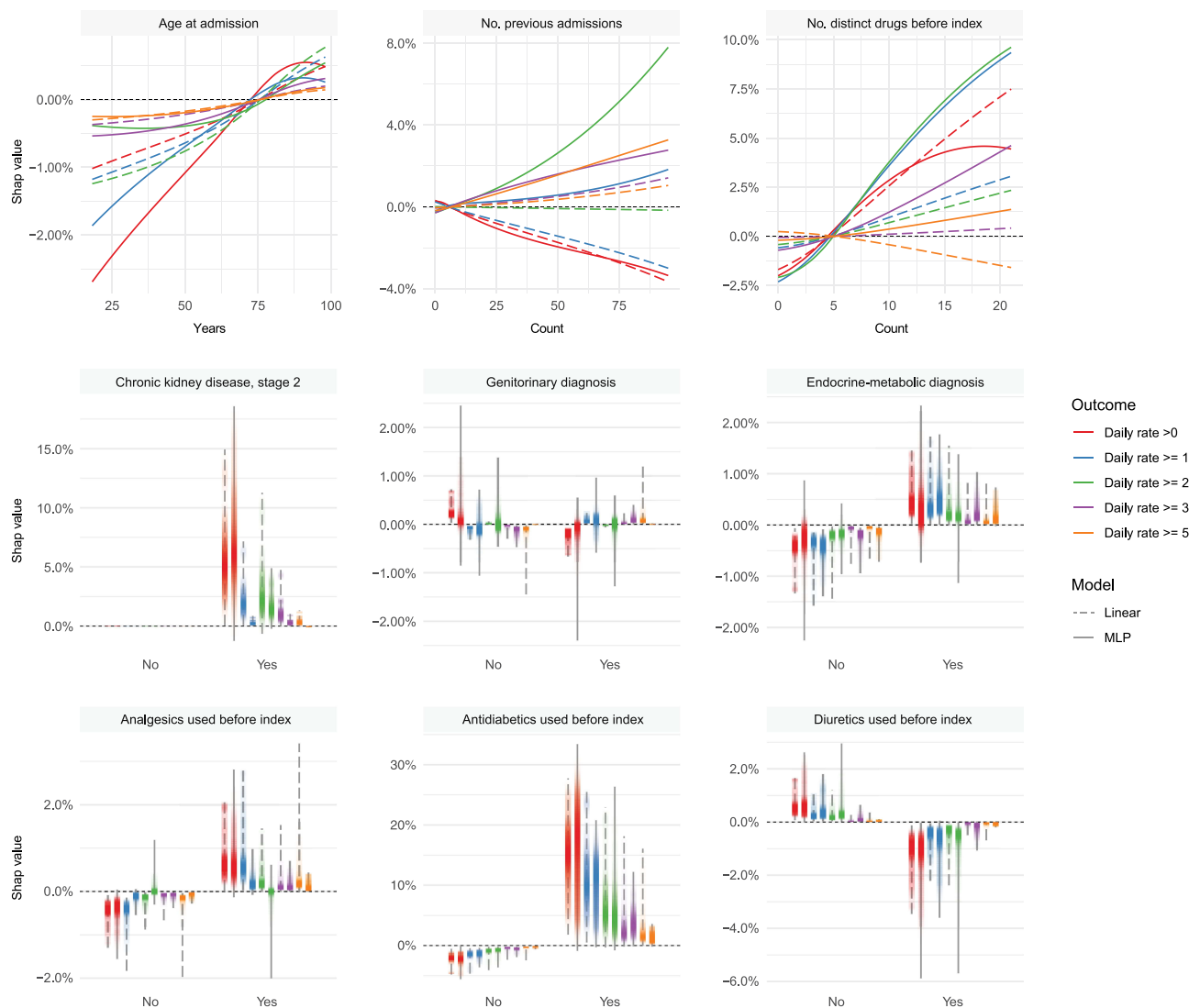
**Notes:** <sup>†</sup>Values in parentheses pertain to the reference models (see text for specification). Undersample: random sample of the size of the minority class, from the majority class. Oversample: randomly sample (with replacement) from the minority class until reaching a sample size equal to the size of the majority class.<sup>26</sup> NearMiss: a method for non-random, systematic downsampling of the majority class while retaining as much information as possible.<sup>27</sup>

**Abbreviations:** AUROC, area under the receiver operating characteristic curve; MLP, multi-layer perceptron; SMOTE, synthetic minority oversampling technique.

models. The decision curves did not suggest the clinical utility of the MLP models be superior to that of the linear (Supplementary Figures S45–S54).

The model-specific shap values offer some insights (Supplementary Figures S59–S68). First, many features contribute substantively to the predictions of daily rate  $>0$  and  $\geq 1$  outcomes, while few features almost entirely drive the predictions for the other outcomes. Second, few features are the dominant prediction drivers across outcomes and models: use of anti-inflammatory, antirheumatic and antidiabetic drugs as well as diagnoses of chronic kidney failure. Third, sex and age contribute little to predictions. Fourth, using more distinct drugs (reflecting various levels of polypharmacy) pushes the risk up and using fewer drugs pulls the risk down. Fifth, the linear models tend to give most weight to relatively few features whereas the MLP models spread out the contributions across more features. Finally, the number of previous admissions (a proxy for frailty) became an increasingly important driver with increasing rarity of the outcome, in the MLP models.

Figure 2 shows the relationships between values of select features and their shap values and illustrates how MLP models capture highly non-linear effects and near-linear effects as appropriate (eg the effects of age at admission and number of previous admissions for daily rate  $>0$ .)



**Figure 2** Bivariate relationships between values of select features (x axis) and their corresponding shap values (y axis). The continuous features are summarized by locally estimated scatterplot smoothing (LOESS), binary features by vertical density bands.



## Discussion

This study reveals that 9.0% of patients with reduced kidney function are exposed to inappropriate doses of selected renal risk drugs in the follow-up period. Our models performed quite well with AUROC's between 0.77 and 0.81 with good calibration-in-the-small for daily rates  $>0$  and  $\geq 1$ , in the test set. For rarer outcomes (daily rates  $\geq 2$ ,  $\geq 3$  and  $\geq 5$ ) calibration suffered and clinical utility is unlikely to be substantive.

Apt intervention necessitates comprehension of the nature and extent of the problem. Use of renal risk drugs and associated problems, including inappropriate dosing, in patients with renal dysfunction is well described.<sup>36–40</sup> A cross-sectional study of 83,000 American outpatient Veterans found that 32% of patients with creatinine clearance between 15 and 29 were given drugs at excessive doses considering their kidney function.<sup>41</sup> Medication burden had the strongest cooccurrence with inappropriate dosing and metformin was a prominent drug among those with inappropriate doses. This agrees with our findings although our study design has clearer temporality.

Some have called for a prediction tool to identify elderly at elevated risk of adverse drug reactions,<sup>42</sup> a notion similar to ours in spirit but different in scope. Studies of factors associated with inadequate dose adjustment are few and often of retrospective nature eliciting relationships with characteristics after inappropriate doses have already been given. One study seeking to elicit factors associated with dosing appropriateness, using a logistic regression, reported the statistically strongest association to be with severity of chronic kidney failure (p-value = 7%).<sup>43</sup> A similar study found dosing errors in 33% of the patients; *age* (odds ratio, OR: 1.05), *number of drug prescriptions* (OR: 1.1) and *number of drugs requiring dose adjustment* (OR: 2.0) were associated with dosing errors.<sup>44</sup> A third study found that, in patients with chronic kidney failure, *late-stage chronic kidney disease*, *number of prescribed drugs* and *presence of comorbidity* were associated with dosing errors. Ill-defined indices and times-at-risk render such enquiries of little use for a priori prediction and risk stratification: the ability to intervene presupposes a reliable estimate of risk in advance, before the event happens.

Carey et al found only few factors to be genuinely predictive of potentially inappropriate prescribing in elderly outside the hospital setting.<sup>45</sup> Our models had AUROC's (0.77–0.81), slightly higher than that of their model (0.76). In a prospective study from Norway<sup>37</sup> of internal-medicine patients with a mean age of 71 years, 35% received suboptimal doses; a composite variable (*number of clinical/pharmacological risk factors*) was quite strongly associated with non-optimal dosing (RR: 1.33), less so *number of drugs at admission* (RR: 1.09), whereas *sex* and *age* were not predictive of non-optimal dosing. Our results agree quite well with that finding, probably because the information captured by age and sex (essentially, proxies of comorbidity) is expressed explicitly in our feature set.

As such, our models fare quite well with performance metrics superior to those of other published models even though ours came from an independent and temporally distinct test set. Many studies employing machine learning models for predicting medical outcomes use normal split-sample validation, putting aside a random sample of the observations for testing. This has several logical and practical implications, perhaps most notably that a model developed with data collected between, say, 2005 and 2015 will likely perform better in a test case from 2013 than in one from 2017. The subset of our test set with patients not part of the development set is a conceptually appealing way to gauge how the model might perform in a new population. It does, however, distort the data and somewhat delink it from the clinical reality: some patients have previous admissions and those admitted for the first time are probably different from the rest.

## Strengths

Here, we highlight six principal strengths of this study. First, this is by far the largest study of its kind to date. Second, time-series validation yielded realistic performance evaluation in distinct (future) data<sup>46</sup> vis-a-vis many articles on predictive modelling, perhaps most clearly seen in the surge of COVID-19 papers.<sup>47</sup> Third, our data were richer than in any other study in this area thanks to the combined diversity and reliability of longitudinal diagnostic data from the National Patient Register and deep phenotypic in-hospital data. Fourth, our summary statistics are well aligned with descriptive studies of deviations from dosing recommendations, and the nature of the general patient population to which a model as ours would be applied.<sup>48</sup> Fifth, we found that the full models performed substantively better than the sparser reference models with few features known to precipitate renal dysfunction. Finally, the shap-value analysis suggests that the models picked up clinically relevant information without undue influence of individual predictors.



## Limitations

Like any study, this has potential limitations. First, albeit simple and elegant, using *only* eGFR as a proxy for kidney function is not always advisable.<sup>49</sup> It is, however, considered a reasonable metric for medicinal dosing<sup>50</sup> and used in Danish guidelines. Second, eGFR can be estimated in several ways<sup>51</sup> and both the 4-variable MRDR Study and CKD-EPI equations were used in our data. However, clinicians use the reported eGFR estimate as-is and both equations perform well for low eGFR values.<sup>52</sup> Third, hard thresholds on eGFR are arbitrary: the difference in kidney function between eGFRs of 29 and 31 is minuscule, but the cutoff must be set somewhere. Again, we stayed loyal to the guidelines as these are, nevertheless, what should support clinicians' prescribing decisions. Fourth, many drugs have narrow and intermediate therapeutic indices. We focused on seven drugs cleared primarily by the kidneys and with narrow therapeutic indices that are fairly common in a Danish setting and span several important drug classes. The drugs included also allowed for reasonably harmonized rules of inappropriate dosing. Fifth, our binary outcomes are soft endpoints and do constitute a simplification. Seemingly inappropriate doses could be conscious choices and the outcome variables do not capture information about actual toxicity experienced by the patient. However, the narrow therapeutic indices of the included drugs increase the likelihood of noxious effects without appropriate dose adjustment. Finally, although low eGFR in the first 24 hours of admission was required for inclusion, using data on the kidney function between admission and index as input to the model could have been beneficial; however, perhaps mostly so if modelling the eGFR trajectory with eg a Long Short-Term Memory neural network or a Hidden Markov Model.

## Conclusion

Despite physicians' awareness of the need for dose adjustment in patients with kidney dysfunction, a well-performing clinical decision support tool may help prevent such patients from "flying under the radar" in a busy clinical setting. Indeed, our models can flag patients at high risk of receiving  $>0$  or  $\geq 1$  inappropriate dose daily.

A prospective evaluation is necessary to assess if these results transport to the clinic and if the models can offer genuine clinical utility for the patients. Receiving inappropriate doses is a soft endpoint, so clinical evaluation should consider also hard endpoints, either generic (eg length-of-stay, need for post-discharge rehabilitation and mortality) or specific ones related to the target drugs (eg transfusion and occurrence of known side-effects of these drugs.)

## Data Sharing Statement

Due to the sensitive nature of the data, we can neither offer access to nor share our data with third parties. Data can be obtained from the original sources upon request.

## Funding

The authors would like to thank Innovation Fund Denmark (5153-00002B) and the Novo Nordisk Foundation (NNF14CC0001, NNF17OC0027594) for their financial contribution to BigTempHealth without which this study had not been possible. The funders played no role in designing, conducting, interpreting, or reporting this study.

## Disclosure

SB reports ownerships in Intomics A/S, Hoba Therapeutics Aps, Novo Nordisk A/S, Lundbeck A/S, and managing board memberships in Proscion A/S and Intomics A/S outside the submitted work. Dr Anna Pors Nielsen reports grants from Novo Nordisk Foundation, during the conduct of the study. All other authors report no other conflicts of interest in this work.

## References

1. Saleem A, Masood I. Pattern and predictors of medication dosing errors in chronic kidney disease patients in Pakistan: a single center retrospective analysis. *PLoS One*. 2016;11(7):e0158677. doi:10.1371/journal.pone.0158677
2. Hoffmann F, Boeschen D, Dorks M, Herget-Rosenthal S, Petersen J, Schmiemann G. Renal Insufficiency and Medication in Nursing Home Residents. A cross-sectional study (IMREN). *Dtsch Arztebl Int*. 2016;113(6):92–98. doi:10.3238/arztebl.2016.0092
3. Munar MY, Singh H. Drug dosing adjustments in patients with chronic kidney disease. *Am Fam Phys*. 2007;75(10):1487–1496.
4. Niedrig D, Krattinger R, Jodicke A, Gott C, Bucklar G, Russmann S. Development, implementation and outcome analysis of semi-automated alerts for metformin dose adjustment in hospitalized patients with renal impairment. *Pharmacoepidemiol Drug Saf*. 2016;25(10):1204–1209. doi:10.1002/pds.4062

5. Bernstein JM, Erk SD. Choice of antibiotics, pharmacokinetics, and dose adjustments in acute and chronic renal failure. *Med Clin North Am*. 1990;74(4):1059–1076. doi:10.1016/S0025-7125(16)30536-3
6. Khare AK. Antibiotic dose adjustment in renal insufficiency. *Lancet*. 1992;340(8833):1480.
7. Dorks M, Allers K, Schmiemann G, Herget-Rosenthal S, Hoffmann F. Inappropriate medication in non-hospitalized patients with renal insufficiency: a systematic review. *J Am Geriatr Soc*. 2017;65(4):853–862. doi:10.1111/jgs.14809
8. Getachew H, Tadesse Y, Shibeshi W. Drug dosage adjustment in hospitalized patients with renal impairment at Tikur Anbessa specialized hospital, Addis Ababa, Ethiopia. *BMC Nephrol*. 2015;16:158. doi:10.1186/s12882-015-0155-9
9. Altunbas G, Yazc M, Solak Y, et al. Renal drug dosage adjustment according to estimated creatinine clearance in hospitalized patients with heart failure. *Am J Ther*. 2016;23(4):e1004–8. doi:10.1097/01.mjt.0000434042.62372.49
10. Hillestad R, Bigelow J, Bower A, et al. Can electronic medical record systems transform health care? Potential health benefits, savings, and costs. *Health Aff (Millwood)*. 2005;24(5):1103–1117. doi:10.1377/hlthaff.24.5.1103
11. Stewart WF, Shah NR, Selna MJ, Paulus RA, Walker JM. Bridging the inferential gap: the electronic health record and clinical evidence. *Health Affairs (Millwood)*. 2007;26(2):w181–91. doi:10.1377/hlthaff.26.2.w181
12. Boussadi A, Caruba T, Karras A, et al. Validity of a clinical decision rule-based alert system for drug dose adjustment in patients with renal failure intended to improve pharmacists' analysis of medication orders in hospitals. *Int J Med Inform*. 2013;82(10):964–972. doi:10.1016/j.ijmedinf.2013.06.006
13. Gawande A Why doctors hate their computers. *The New Yorker*. November 12, 2018.
14. Baysari MT, Tariq A, Day RO, Westbrook JI. Alert override as a habitual behavior - a new perspective on a persistent problem. *J Am Med Inform Assoc*. 2017;24(2):409–412. doi:10.1093/jamia/ocw072
15. Kane-Gill SL, O'Connor MF, Rothschild JM, et al. Technologic Distractions (Part 1): summary of approaches to manage alert quantity with intent to reduce alert fatigue and suggestions for alert fatigue metrics. *Crit Care Med*. 2017;45(9):1481–1488. doi:10.1097/CCM.0000000000002580
16. Jensen TB, Jimenez-Solem E, Cortes R, et al. Content and validation of the Electronic Patient Medication module (EPM)—the administrative in-hospital drug use database in the Capital Region of Denmark. *Scand J Public Health*. 2018;48(1):43–48. doi:10.1177/1403494818760050
17. Zhang Z, Ho KM, Hong Y. Machine learning for the prediction of volume responsiveness in patients with oliguric acute kidney injury in critical care. *Crit Care*. 2019;23(1):112. doi:10.1186/s13054-019-2411-z
18. Efron B, Hastie T. *Computer Age Statistical Inference: Algorithms, Evidence, and Data Science*. London, United Kingdom: Cambridge University Press; 2016.
19. Rothman KJ, Lash TL, Greenland S. *Modern Epidemiology*. 3rd ed. Lippincott Williams & Wilkins; 2012.
20. Schmidt M, Schmidt SAJ, Sandegaard JL, Ehrenstein V, Pedersen L, Sørensen HT. The Danish National Patient Registry: a review of content, data quality, and research potential. *Clin Epidemiol*. 2015;7:449–490. doi:10.2147/CLEP.S91125
21. Moore BJ, White S, Washington R, Coenen N, Elixhauser A. Identifying increased risk of readmission and in-hospital mortality using hospital administrative data: the AHRQ Elixhauser Comorbidity Index. *Med Care*. 2017;55(7):698–705. doi:10.1097/MLR.0000000000000735
22. Chollet F. *Deep Learning with Python*. New York, USA: Manning Publications Co; 2018.
23. Akiba T, Sano S, Yanase T, Ohta T, Koyama M. Optuna: a next-generation hyperparameter optimization framework. *arXiv (Unpublished)*. 2019.
24. Goodfellow I, Bengio Y, Courville A. *Deep Learning*. Cambridge (MA), USA: MIT Press; 2016.
25. Hastie T, Tibshirani R, Friedman J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. 2nd ed. New York: Springer; 2009.
26. Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP. SMOTE: synthetic Minority Over-sampling Technique. *J Artif Intell Res*. 2002;16:321–357. doi:10.1613/jair.953
27. Zhang J, Mani I kNN approach to unbalanced data distributions: a case study involving information extraction. In: *Proceedings of the ICML'2003 Workshop on Learning from Imbalanced Datasets*; 2003.
28. Steyerberg EW. *Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating*. New York: Springer; 2009.
29. Brown LD, Cai TT, DasGupta A. Interval estimation for a binomial proportion. *Statistical Sci*. 2001;16(2):101–133. doi:10.1214/ss/1009213286
30. Vickers AJ, Elkin EB. Decision curve analysis: a novel method for evaluating prediction models. *Med Decision Making*. 2006;26(6):565–574. doi:10.1177/0272989X06295361
31. Kerr KF, Brown MD, Zhu K, Janes H. Assessing the clinical impact of risk prediction models with decision curves: guidance for correct interpretation and appropriate use. *J Clin Oncol*. 2016;34(21):2534–2540. doi:10.1200/JCO.2015.65.5654
32. Lundberg SM, Lee S. A unified approach to interpreting model predictions. In: Guyon I, Luxburg UV, Bengio S, et al., editors. *Advances in Neural Information Processing Systems*. Vol. 30. Curran Associated, Inc.; 2017:4765–4774.
33. Köster J, Rahmann S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics*. 2012;28(19):2520–2522. doi:10.1093/bioinformatics/bts480
34. Hernandez-Boussard T, Bozkurt S, Ioannidis JPA, Shah NH. MINIMAR (MINimum Information for Medical AI Reporting): developing reporting standards for artificial intelligence in health care. *J Am Med Inf Assoc*. 2020;6:2011–2015.
35. Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD): the TRIPOD Statement. *Ann Intern Med*. 2015;162(1):55–63. doi:10.7326/M14-0697
36. Saad R, Hallit S, Chahine B. Evaluation of renal drug dosing adjustment in chronic kidney disease patients at two university hospitals in Lebanon. *Pharm Pract (Granada)*. 2019;17(01):1304. doi:10.18549/PharmPract.2019.1.1304
37. Blix HS, Viktil KK, Reikvam A, et al. The majority of hospitalised patients have drug-related problems: results from a prospective study in general hospitals. *Eur J Clin Pharmacol*. 2004;60(9):651–658. doi:10.1007/s00228-004-0830-4
38. Andreu Cayuelas JM, Caro Martínez C, Flores Blanco PJ, et al. Kidney function monitoring and nonvitamin K oral anticoagulant dosage in atrial fibrillation. *Eur J Clin Invest*. 2018;48(6):e12907. doi:10.1111/eci.12907
39. Seiberth S, Bauer D, Schönermarck U, et al. Correct use of non-indexed eGFR for drug dosing and renal drug-related problems at hospital admission. *Eur J Clin Pharmacol*. 2020;76(12):1683–1693. doi:10.1007/s00228-020-02953-6
40. Breton G, Froissart M, Janus N, et al. Inappropriate drug use and mortality in community-dwelling elderly with impaired kidney function—the Three-City population-based study. *Nephrol Dial Transplant*. 2011;26(9):2852–2859. doi:10.1093/ndt/gfq827

41. Chang F, O'Hare AM, Miao Y, Steinman MA. Use of renally inappropriate medications in older veterans: a national study. *J Am Geriatr Soc.* 2015;63(11):2290–2297. doi:10.1111/jgs.13790
42. Parameswaran Nair N, Chalmers L, Peterson GM, Bereznicki BJ, Castelino RL, Bereznicki LR. Hospitalization in older patients due to adverse drug reactions - the need for a prediction tool. *Clin Interv Aging.* 2016;11:497–505. doi:10.2147/CIA.S99097
43. Kalender-Rich JL, Mahnken JD, Wetmore JB, Rigler SK. Transient impact of automated glomerular filtration rate reporting on drug dosing for hospitalized older adults with concealed renal insufficiency. *Am J Geriatr Pharmacother.* 2011;9(5):320–327. doi:10.1016/j.amjopharm.2011.08.003
44. Won H, Chung G, Lee KJ, et al. Evaluation of medication dosing errors in elderly patients with renal impairment. *Int J Clin Pharmacol Ther.* 2018;56(8):358–365. doi:10.5414/CP203258
45. Carey IM, De Wilde S, Harris T, et al. What factors predict potentially inappropriate primary care prescribing in older people? *Drugs Aging.* 2008;25(8):693–706. doi:10.2165/00002512-200825080-00006
46. Steyerberg EW, Harrell FEJ. Prediction models need appropriate internal, internal-external, and external validation. *J Clin Epidemiol.* 2016;69:245–247. doi:10.1016/j.jclinepi.2015.04.005
47. Wynants L, Van Calster B, Collins GS, et al. Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal. *BMJ.* 2020;7:369.
48. Yusuf M, Atal I, Li J, et al. Reporting quality of studies using machine learning models for medical diagnosis: a systematic review. *BMJ Open.* 2020;10:3. doi:10.1136/bmjopen-2019-034568
49. Eppenga WL, Kramers C, Derijks HJ, Wensing M, Wetzels JFM, De Smet PAGM. Drug therapy management in patients with renal impairment: how to use creatinine-based formulas in clinical practice. *Eur J Clin Pharmacol.* 2016;72(12):1433–1439. doi:10.1007/s00228-016-2113-2
50. Rule AD, Glasscock RJ. GFR estimating equations: getting closer to the truth?. *Clin J Am Soc Nephrol.* 2013;8(8):1414–1420.
51. Corsonello A, Onder G, Bustacchini S, et al. Estimating renal function to reduce the risk of adverse drug reactions. *Drug Safety.* 2012;35(Suppl 1):47–54. doi:10.1007/BF03319102
52. Levey AS, Stevens LA, Schmid CH, et al. A new equation to estimate glomerular filtration rate. *Ann Intern Med.* 2009;150(9):604–612. doi:10.7326/0003-4819-150-9-200905050-00006

## Clinical Epidemiology

Dovepress

### Publish your work in this journal

Clinical Epidemiology is an international, peer-reviewed, open access, online journal focusing on disease and drug epidemiology, identification of risk factors and screening procedures to develop optimal preventative initiatives and programs. Specific topics include: diagnosis, prognosis, treatment, screening, prevention, risk factor modification, systematic reviews, risk & safety of medical interventions, epidemiology & biostatistical methods, and evaluation of guidelines, translational medicine, health policies & economic evaluations. The manuscript management system is completely online and includes a very quick and fair peer-review system, which is all easy to use.

Submit your manuscript here: <https://www.dovepress.com/clinical-epidemiology-journal>