

1 Identifying patients at high risk of inappropriate drug
2 dosing in periods with renal dysfunction

3
4 Benjamin Skov Kaas-Hansen MD, MSc^{*1,2} – Cristina Leal Rodríguez MSc²

5 Davide Placido MSc² – Hans-Christian Thorsen-Meyer MD^{2,3}

6 Anna Pors Nielsen MD² – Nicolas Dérian MSc, PhD⁴

7 Søren Brunak MSc, PhD² – Stig Ejdrup Andersen MD, PhD¹

8

9 Affiliations

10 * Corresponding author. Address: Munkesoevej 18, 4000 Roskilde, Denmark.

11 ¹ Clinical Pharmacology Unit, Zealand University Hospital, Roskilde, Denmark

12 ² NNF Center for Protein Research, University of Copenhagen, Denmark

13 ³ Department of Intensive Care Medicine, Copenhagen University Hospital
14 (Rigshospitalet), Copenhagen, Denmark

15 ⁴ Data and Development Support, Region Zealand, Denmark

16 ORCID ID's

17 BSKH: 0000-0003-1023-0371. SEA: 0000-0002-1914-4720. CLR: 0000-0002-3133-0630.

18 ND: 0000-0002-4477-023X. APN: 0000-0001-7903-5051. SB: 0000-0003-0316-5866.

19 **Keywords**

- 20 • Predictive modelling
- 21 • Kidney failure/Renal dysfunction
- 22 • Machine learning
- 23 • Risk markers
- 24 • Inappropriate drug dosing

25 **Abstract**

26 **Introduction**

27 Dosing of renally cleared drugs in patients with kidney failure often deviates from
28 clinical guidelines but little is known about what is predictive of receiving
29 inappropriate doses.

30 **Methods and materials**

31 We combined data from the Danish National Patient Register and in-hospital data
32 on drug administrations and estimated glomerular filtration rates for admissions
33 between 1 October 2009 and 1 June 2016, from a pool of about 2.9 million persons.
34 We trained artificial neural network and linear logistic ridge regression models to
35 predict the risk of five outcomes (>0 , ≥ 1 , ≥ 2 , ≥ 3 and ≥ 5 inappropriate doses daily)
36 with index set 24 hours after admission. We used time-series validation for
37 evaluating discrimination, calibration, clinical utility and explanations.

38 **Results**

39 Of 52,451 admissions included, 42,250 (81%) were used for model development.
40 The median age was 77 years; 50% of admissions were of women. ≥ 5 drugs were
41 used between admission start and index in 23,124 admissions (44%); the most
42 common drug classes were analgesics, systemic antibacterials, diuretics,
43 antithrombotics, and antacids. The neural network models had better
44 discriminative power (all AUROCs between 0.77 and 0.81) and were better
45 calibrated than their linear counterparts. The main prediction drivers were use of

46 anti-inflammatory, antidiabetic and anti-Parkinson's drugs as well as having a
47 diagnosis of chronic kidney failure. Sex and age affected predictions but slightly.

48 Conclusion

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

53 Introduction

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

60 Inadequate dose-adjustment of such drugs has been linked to polypharmacy [1,2]
61 and can cause noxious events [3] or accidental over-dosing [4]. Although not a new
62 issue, [5,6] deviating from guidelines is widespread with prevalence estimates up
63 to 70% [1,2,7-9]. Despite large inter-individual variability in clearance and
64 response, dose adjustment for many drugs is crude and based on the estimated
65 glomerular filtration rate (eGFR), for example, halving the dose when eGFR < 60
66 ml/min/1.73 m².

67 Appropriate alerts in order-entry systems may facilitate rational clinical decision-
68 making, [10,11] and convincing examples have showcased how computerised
69 systems can underpin rational pharmacotherapy [4,12]. However, downsides of
70 extensive computerisation of healthcare emerge [13]; alert fatigue [14] is
71 particularly problematic, and strategies and interventions have been proposed to
72 mitigate its negative effects [15].

73 At Danish hospitals, prescriptions are mostly dispensed and administered by
74 nurses who record detailed meta-data [16]. Prescriptions are usually made and

75 revised by physicians regularly during clinical rounds, typically in the morning or
76 early afternoon. Electronic decision support is generally immature and neither
77 prescribing physicians nor dispensing nurses are warned if dose-adjustment be
78 advised or even required.

79 We suspect that the need for dose-adjustment in patients with renal dysfunction
80 often goes unrecognised. Thus, with this paper we study its predictability to
81 inform clinicians and healthcare personnel upfront about which patients with renal
82 dysfunction are at elevated risk of inappropriate drug dosing. To this end we used
83 and compared predictive modelling methods from classical statistical modelling
84 and machine learning.

85 **Methods**

86 **Study design, patients and data**

87 We conducted a register-based prediction study with prospective data for patients
88 admitted to 12 public hospitals in two Danish regions comprising about 2.9 million
89 persons (more than half the Danish population). We collected diagnosis data from
90 the Danish National Patient Register, demographic data from the Danish Civil
91 Registration System [17], as well as medication and biochemical data from
92 electronic patient records. Diagnoses were encoded using the 10th revision of the
93 International Classification of Diseases (ICD-10), drugs with the Anatomical and
94 Therapeutic Chemical classification (ATC).

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

102 Outcomes

103 The outcome variables were based on the daily rate $= r/E$ of inappropriate doses
104 during follow-up, capped at 30 days. r is the number of given inappropriate doses
105 of select drugs cleared mainly renally and with narrow therapeutic indices; E the
106 time-at-risk (figure 1). To obtain well-defined times-at-risk, we set the eGFR
107 threshold to ≤ 30 ml/min/1.73m² (unit omitted from here onward) and used the
108 rules in supplementary table S1 for counting the number of inappropriate doses,
109 based on the official reference guidelines for Danish physicians (pro.medicin.dk) as
110 of January 2021.

111 We used two rules, one definitive (maximum daily dose = 0 mg) and one of dose-
112 adjustment (reduced daily dose). Operationalisation of the definitive rule is
113 straightforward: if the last eGFR ≤ 30 , there should be no administrations until an
114 eGFR > 30 is measured. The dose-adjustment rule is slightly more involved as
115 inappropriate dosing comes in two forms: (a) on a given day there are more than
116 one eGFR measurements, of which at least one is ≤ 30 , and the cumulative daily

117 dose surpasses the threshold in the period(s) between above-threshold
118 measurements, or (b) all eGFR measurements of a given day are ≤ 30 and the
119 cumulative daily dose surpasses the threshold.

120 Variables and features

121 Variables are original data (e.g. sex and age at admission) and features the results
122 of rendering the variables appropriate as model inputs (e.g. one-hot-encoded day
123 of admission). Based on clinical and pharmacological experience we hand-picked
124 pertinent variables likely to be informative to the prediction problem and
125 realistically available in the clinical setting. These fall into three categories.
126 Demographic: age at admission (numeric), sex (binary). Clinical: number of
127 distinct drugs (ATC level 5) administered between admission and index (numeric);
128 therapeutic drug classes (ATC level 2) used between admission and index (one-
129 hot-encoded); the Elixhauser score at admission (numeric, AQHR adaptation) [18];
130 ICD-10 chapters of diagnoses recorded in the past five years before admission
131 (one-hot-encoded); record of chronic kidney failure in the past five years before
132 admission (ICD-10 N18* diagnoses, one-hot-encoded). Contextual: hour of
133 admission (numeric, transformed as $f(t) = \text{abs}(12 - t)$; see supplementary figure S1);
134 weekday of admission (one-hot-encoded); number of admissions in the past 5
135 years before admission (numeric).
136 Missing values, only present for hour of admission and discharge, were imputed
137 by sampling from the empirical distributions of valid values.

138 Models and training

139 We tried two model architectures (linear logistic ridge regression and artificial
140 neural network) with several binary outcomes defined by increasing thresholds of
141 the daily rate of inappropriate doses (>0 , ≥ 1 , ≥ 2 , ≥ 3 and ≥ 5). The neural network
142 models were multilayer perceptrons (MLPs) enabling speedy training and
143 evaluation.

144 All admissions starting before 1 July 2015 were assigned to the development set
145 (42,250 admissions [81%] of 27,253 patients) and the rest to the independent hold-
146 out test set (10,201 admissions [19%] of 8,412 patients). Because admissions
147 constitute the unit of analysis, some patients likely appear in both the development
148 and test sets. Information may leak between the sets [19] so as a sensitivity
149 analysis, we evaluated the performance also in the subset of test-set patients not in
150 the development set.

151 We used the multivariate *TPEsampler* from *Optuna* [20] to find the best-performing
152 hyperparameters by sampling 100 configurations, each using 5-fold stratified-and-
153 grouped cross-validation, from the following proposal distributions (discrete
154 values in round brackets, bounds of log-uniform distributions in squared):
155 optimiser (Adam, RMSprop), learning rate [10^{-6} , 10^{-1}], activation function (tanh,
156 sigmoid), L2 penalty [10^{-6} , 10^{-2}], number of hidden layers (1, 2, 3, 4), number of
157 nodes per hidden layer [16, 32, 65, 128], batch size (32, 64, 128, 256, 512), class
158 handling (see below).

159 Only relevant hyperparameters were sampled and we ran Optuna on linear and
160 MLP models separately because they have disparate hyperparameter sets. MLP
161 models with more hidden layers and more nodes therein can learn more complex
162 relationships but become prone to overfitting which we countered with early
163 stopping [21] and L2 regularisation (handles collinearity better than L1
164 regularisation) [22,23]. The batch size is the number of observations from which
165 the model learns at a time; small batches can give outliers undue influence while
166 full-batch training (batch size = number of units) can become computationally
167 impractical [19]. Class imbalances in binary outcomes can misguide training, so we
168 tested the following remedies: synthetic minority oversampling technique
169 (SMOTE), random over-sampling of minority class, NearMiss, random under-
170 sampling of majority class, class weighting, and none. SMOTE creates a dataset
171 similar to the minority class but of the same size as the majority class [24];
172 NearMiss downsizes the majority class in a systematic way to retain as much
173 information as possible in fewer data points [25]. Class weighting retains the
174 original data but gives more weight to minority-class observations.
175 Hyperparameter optimisation models trained for maximum 500 epochs with 50-
176 epoch patience on improvement in the validation loss. The final models were
177 trained on the full development set until the loss reached that obtained in the best
178 cross-validation fold for the best configuration [21].

179 **Evaluation and explanation**

180 Discrimination was assessed with receiver operating characteristic (ROC) curves
181 and areas under the ROC curves (AUROC), calibration-in-the-small by plotting
182 decile-binned predicted probabilities against corresponding bin-wise observed
183 event proportions [26] with 95% Jeffrey intervals [27]; results from a perfectly
184 calibrated model fall on the diagonal. We used the decision-curve analytic
185 framework to gauge the models' potential clinical utility [28,29].
186 For explanation and scrutiny of prediction drivers, we used the SHAP
187 DeepExplainer yielding one shap value per feature per unit [30]. The shap value
188 for a risk prediction model is the absolute change in risk of a given unit's value for
189 each feature: the cohort-wide mean risk plus the sum of one unit's shap values
190 equals that unit's risk.

191 **Analysis and ethics**

192 The full analytical pipeline was built with Snakemake [31] (schematic overview in
193 supplementary figure S2) to facilitate transparency and reproducibility; blinding
194 was impractical and so foregone, but all analytic code is available online (DOI:
195 [10.5281/zenodo.4560078](https://doi.org/10.5281/zenodo.4560078)). Univariate distributions were summarised by median
196 (inter-quartile range) and count (proportion), as appropriate. This report adheres
197 to pertinent items in the MINIMAR guideline [32] and TRIPOD statement [33].
198 All data have been marshalled on Computerome, a secure high-performance
199 Danish computing infrastructure, after obtaining approval from the Danish Patient

- 200 Safety Authority (3-3013-1723; then competent authority for ethical approval), the
- 201 Danish Data Protection Agency (DT SUND 2016-48, 2016-50, 2017-57) and the
- 202 Danish Health Data Authority (FSEID 00003724).

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 (N = 42,250)	Test set (N = 10,201)	Test set (not in devel. set) (N = 5,980)
Women	20,743 (49%)	4,854 (48%)	2,940 (49%)
Distinct patients	27,253	8,412	5,341
Distinct women	13,759 (50%)	4,049 (48%)	2,629 (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)	3,786 (9.0%)	1,080 (11%)	740 (12%)
≥ 1 daily	2,241 (5.3%)	588 (5.8%)	333 (5.6%)
≥ 2 daily	1,236 (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	4,988 (12%)	1,082 (11%)	1,074 (18%)
1–2	10,100 (24%)	2,367 (23%)	1,873 (31%)
3–4	7,712 (18%)	1,919 (19%)	1,232 (21%)
5–6	5,490 (13%)	1,303 (13%)	685 (12%)
≥ 7	13,960 (33%)	3,530 (35%)	1,116 (19%)
Drugs used between admission and index			
None	6,165 (15%)	1,228 (12%)	762 (13%)
1–2	9,111 (22%)	1,984 (19%)	1,254 (21%)
3–4	8,761 (21%)	2,078 (20%)	1,355 (23%)
5–6	7,197 (17%)	1,852 (18%)	1,095 (18%)
≥ 7	11,016 (26%)	3,059 (30%)	1,514 (25%)
Any diagnosis of chronic kidney failure	13,470 (32%)	3,391 (33%)	732 (12%)
Top-5 ICD-10 chapters [†]			
Cardiovascular (IX)	25,757 (61%)	6,392 (63%)	3,283 (55%)
Genitourinary (XIV)	23,025 (55%)	5,819 (57%)	2,306 (39%)
Lesions, external causes, etc. (XIX)	20,275 (48%)	4,749 (47%)	2,481 (42%)
Metabolic-endocrine (IV)	19,716 (47%)	5,096 (50%)	2,415 (40%)
Symptoms/abnormal findings (XVIII)	18,663 (44%)	5,711 (56%)	2,882 (48%)
Top-5 drug classes [‡]			
Analgesics (N02)	15,740 (37%)	4,367 (43%)	2,506 (42%)
Systemic antibacterials (J01)	14,719 (35%)	3,257 (32%)	1,938 (32%)

Diuretics (C03)	13,966 (33%)	3,672 (36%)	1,951 (33%)
Antithrombotics (B01)	11,842 (28%)	3,181 (31%)	1,795 (30%)
Antacids (A02)	10,635 (25%)	2,776 (27%)	1,407 (24%)

† ICD-10 chapters (Roman numbering) of diagnoses recorded in the last 5 years before admission.

‡ Drug classes (ATC level 2) administered between admission and index.

204 Results

205 Table 1 shows univariate summary statistics of the 52,451 admissions (42,250 +
206 10,201) of 35,665 patients (27,253 + 8,412) included in the study (see supplementary
207 table S2 for extended version with all features). Patients in the test sets were similar
208 to those in the development set with some notable exceptions. Fewer had received
209 inappropriate doses, especially in the test-set patients not part of the development
210 set who also had fewer previous admissions.

211 In the development set, the median age was 77 years (IQR: 67-85) and 20,743
212 admissions (49%) were of 13,759 women (50%). The median time at risk was 3.5
213 days (inter-quartile range: 1.7–7.7) and at least one inappropriate dose was given in
214 3,786 admissions (9.0%); ≥ 1 inappropriate dose daily was given in 5.3% of
215 admissions and ≥ 5 inappropriate doses daily were given in 0.9%. The target drugs
216 most commonly given in inappropriate doses were ibuprofen (M01AE01, 4.1%)
217 and metformin (A10BA02, 3.4%); inappropriate doses of the other target drugs
218 were given in <1% of admissions.

219 Patients in 4,988 admissions (12%) had no admissions in the 5 years before
220 inclusion; 13,960 (33%) had ≥ 7 previous admissions. The most common drug
221 classes used between admission and index were analgesics (N02, 37%), systemic
222 antibacterials (J01, 35%), diuretics (C03, 33%) antithrombotics (B01, 28%), and
223 antacids (A02, 25%). Previous diagnoses were most commonly cardiovascular
224 (chapter IX, 61%), genitourinary (XIV, 55%), related to i.a. lesions and external

- 225 causes (XIX, 48%), endocrine-metabolic (IV, 47%), and symptoms/abnormal
- 226 findings (XVIII, 44%).

Table 2: Performance metrics of final models and results of Optuna hyperparameter optimisation. AUROC: area under the receiver operating characteristic curve. MLP: multi-layer perceptron. 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. SMOTE: synthetic minority oversampling technique [24]. NearMiss: a method for non-random, systematic downsampling of the majority class while retaining as much information as possible [25].

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										
Development set	0.80	0.81	0.81	0.83	0.81	0.84	0.82	0.83	0.82	0.83
Test set	0.77	0.79	0.78	0.79	0.79	0.79	0.81	0.81	0.78	0.80
Test set (new patients)	0.78	0.79	0.82	0.83	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	1.28×10^{-6}	1.66×10^{-6}	3.02×10^{-6}	1.43×10^{-6}	4.38×10^{-6}	1.39×10^{-6}	1.43×10^{-6}	1.30×10^{-6}	1.09×10^{-5}	3.94×10^{-6}
Learning rate	1.79×10^{-2}	1.20×10^{-4}	1.92×10^{-2}	3.45×10^{-4}	6.73×10^{-3}	2.71×10^{-4}	3.76×10^{-2}	3.08×10^{-4}	2.11×10^{-2}	4.86×10^{-4}
Optimiser	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

227 Table 2 shows the hyperparameters of the best configurations with performance
228 metrics of the final models (see also supplementary figures S3–S12). Generally,
229 multi-layer perceptron (MLP) models performed slightly better than their linear
230 counterparts, all obtaining AUROC's between 0.77 and 0.81 in the test set (ROC
231 curves in supplementary figures S13–S22). The MLP models more consistently
232 showed good calibration in the development set. For daily rates >0 , ≥ 1 and ≥ 2 both
233 MLP and linear models were very well-calibrated in the test set (supplementary
234 figures S23–S32). The decision curves did not suggest the clinical utility of the MLP
235 models be superior to that of the linear (supplementary figures S33–S42).

236 The model-specific shap values offer some insights (supplementary figures S43–
237 S53). First, many features contribute substantively to the predictions of daily rate
238 >0 and ≥ 1 outcomes, while few features almost entirely drive the predictions for
239 the other outcomes. Second, few features are the dominant prediction drivers
240 across outcomes and models: use of anti-inflammatory, antirheumatic and
241 antidiabetic drugs as well as diagnoses of chronic kidney failure. Third, sex and
242 age contribute little to predictions. Fourth, using more distinct drugs (reflecting
243 various levels of polypharmacy) pushes the risk up and using fewer drugs pulls
244 the risk down. Fifth, the linear models tend to give most weight to relatively few
245 features whereas the MLP models spread out the contributions across more
246 features. Finally, the number of previous admissions (a proxy for frailty) became
247 an increasingly important driver with increasing rarity of the outcome, in the MLP
248 models.

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

253 Discussion

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

260 Apt intervention necessitates comprehension of the nature and extent of the
261 problem. Use of renal risk drugs and associated problems, including inappropriate
262 dosing, in patients with renal dysfunction is well-described [34-38]. A cross-
263 sectional study of 83,000 American outpatient Veterans found that 32% of patients
264 with creatinine clearance between 15 and 29 were given drugs at excessive doses
265 considering their kidney function [39]. Medication burden had the strongest
266 cooccurrence with inappropriate dosing and metformin was a prominent drug
267 among those with inappropriate doses. This agrees with our findings although our
268 study design has clearer temporality.

269 Some have called for a prediction tool to identify elderly at elevated risk of adverse
270 drug reactions [40], a notion similar to ours in spirit but different in scope. Studies
271 of factors associated with inadequate dose adjustment are few and often of
272 retrospective nature eliciting relationships with characteristics after inappropriate
273 doses have already been given. One study seeking to elicit factors associated with
274 dosing appropriateness, using a logistic regression, reported the statistically
275 strongest association to be with severity of chronic kidney failure (p-value = 7%)
276 [41]. A similar study found dosing errors in 33% of the patients; *age* (odds ratio,
277 OR: 1.05), *number of drug prescriptions* (OR: 1.1) and *number of drugs requiring dose*
278 *adjustment* (OR: 2.0) were associated with dosing errors [42]. A third study found
279 that, in patients with chronic kidney failure, *late-stage chronic kidney disease*, *number*
280 *of prescribed drugs* and *presence of comorbidity* were associated with dosing errors. Ill-
281 defined indices and times-at-risk render such enquiries of little use for a priori
282 prediction and risk stratification: the ability to intervene presupposes a reliable
283 estimate of risk in advance, before the event happens.

284 Carey et al. found only few factors to be genuinely predictive of potentially
285 inappropriate prescribing in elderly outside the hospital setting [43]. Our models
286 had AUROC's (0.77–0.81) slightly higher than that of their model (0.76). In a
287 prospective study from Norway [35] of internal-medicine patients with a mean age
288 of 71 years, 35% received suboptimal doses; a composite variable (*number of*
289 *clinical/pharmacological risk factors*) was quite strongly associated with non-optimal
290 dosing (RR: 1.33), less so *number of drugs at admission* (RR: 1.09), whereas *sex* and
291 *age* were not predictive of non-optimal dosing. Our results agree quite well with

292 that finding, probably because the information captured by age and sex
293 (essentially, proxies of comorbidity) is expressed explicitly in our feature set.

294 As such, our models fare quite well with performance metrics superior to those of
295 other published models even though ours came from an independent and
296 temporally distinct test set. Many studies employing machine learning models for
297 predicting medical outcomes use normal split-sample validation, putting aside a
298 random sample of the observations for testing. This has several logical and
299 practical implications, perhaps most notably that a model developed with data
300 collected between, say, 2005 and 2015 will likely perform better in a test case from
301 2013 than in one from 2017. The subset of our test set with patients not part of the
302 development set is a conceptually appealing way to gauge how the model might
303 perform in a new population. It does, however, distort the data and somewhat
304 delink it from the clinical reality: some patients have previous admissions and
305 those admitted for the first time are probably different from the rest.

306 **Strengths**

307 Here we highlight five principal strengths of this study. First, this is by far the
308 largest study of its kind to date. Second, time-series validation yielded realistic
309 performance evaluation in distinct (future) data [44] vis-a-vis many articles on
310 predictive modelling, perhaps most clearly seen in the surge of COVID-19 papers
311 [45]. Third, our data were richer than in any other study in this area thanks to the
312 combined diversity and reliability of longitudinal diagnostic data from the
313 National Patient Register and deep phenotypic in-hospital data. Fourth, our

314 summary statistics are well-aligned with descriptive studies of deviations from
315 dosing recommendations, and the nature of the general patient population to
316 which a model as ours would be applied [46]. Finally, the shap-value analysis
317 suggests that the models picked up clinically relevant information without undue
318 influence of individual predictors.

319 Limitations

320 Like any study, this has potential limitations. First, albeit simple and elegant, using
321 *only* eGFR as a proxy for kidney function is not always advisable [47]. It is,
322 however, considered a reasonable metric for medicinal dosing [48] and used in
323 Danish guidelines. Second, eGFR can be estimated in several ways [49] and both
324 the 4-variable MRDR Study and CKD-EPI equations were used in our data.
325 However, clinicians use the reported eGFR estimate as-is and both equations
326 perform well for low eGFR values [50]. Third, hard thresholds on eGFR are
327 arbitrary: the difference in kidney function between eGFRs of 29 and 31 is
328 minuscule, but the cutoff must be set somewhere. Again, we stayed loyal to the
329 guidelines as these are, nevertheless, what should support clinicians' prescribing
330 decisions. Fourth, many drugs have narrow and intermediate therapeutic indices.
331 We focused on seven drugs cleared primarily by the kidneys and with narrow
332 therapeutic indices that are fairly common in a Danish setting and span several
333 important drug classes. The drugs included also allowed for reasonably
334 harmonised rules of inappropriate dosing. Finally, our binary outcomes are soft
335 endpoints and do constitute a simplification. Seemingly inappropriate doses could

336 be conscious choices and the outcome variables do not capture information about
337 actual toxicity experienced by the patient. However, the narrow therapeutic
338 indices of the included drugs increase the likelihood of noxious effects without
339 appropriate dose adjustment.

340 Conclusion

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

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

352 Data availability

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

356 Acknowledgements

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

362 Contributions

363 Conceptualisation: BSKH, SEA. Data curation: BSKH, CLR. Formal analysis: BSKH.
364 Methodology: APN, BSKH, DP, HCTM, ND. Software: BSKH. Code review: CLR,
365 DP, HCTM. Drafting: BSKH. Funding acquisition: SB, SEA. Resources: SB, SEA.
366 Supervision: SEA. Review: All.

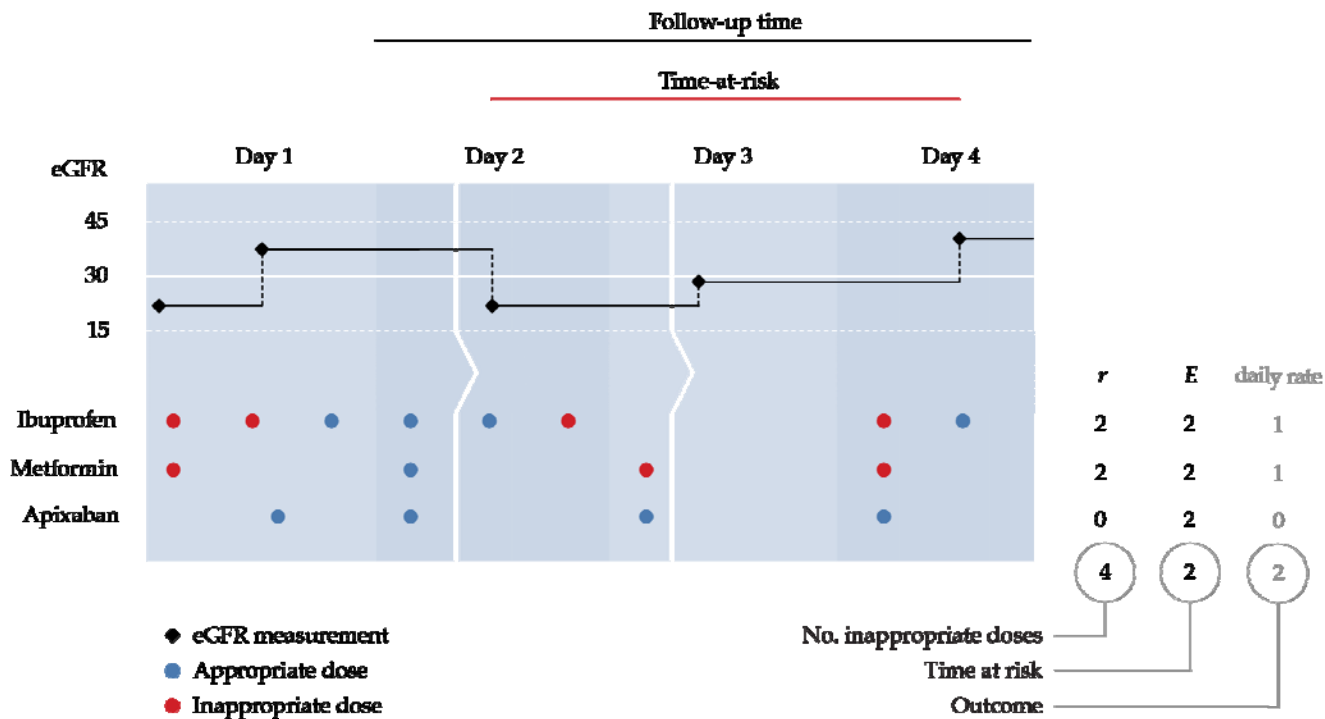
367 Conflicts of interest

368 The authors declare the following competing interests:

- 369 • BSKH: None
- 370 • CRL: None
- 371 • DP: None
- 372 • HCTM: None
- 373 • ND: None
- 374 • APN: None

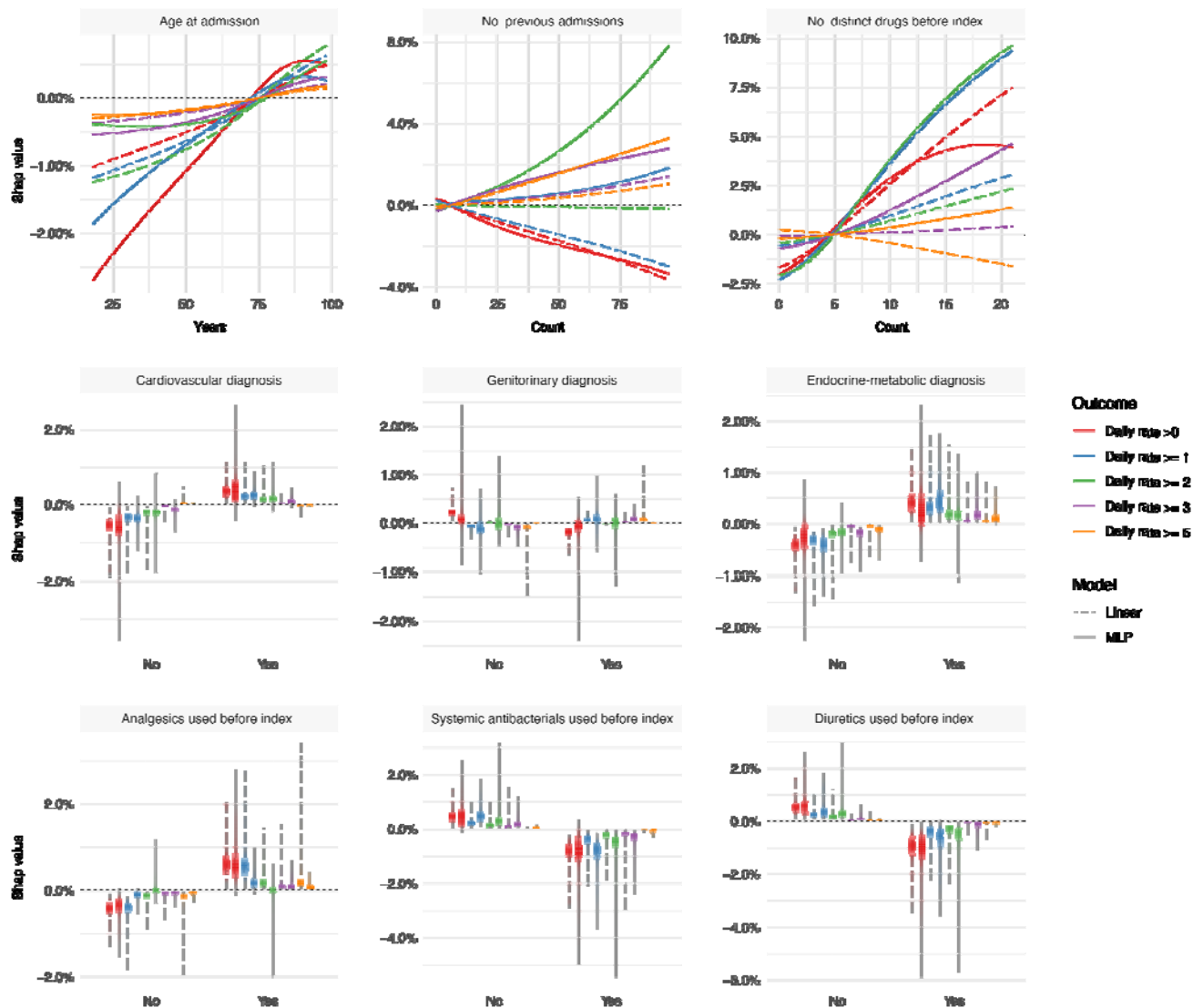
- 375 • SB reports ownerships in Intomics A/S, Hoba Therapeutics Aps, Novo
- 376 Nordisk A/S, Lundbeck A/S, and managing board memberships in
- 377 Proscion A/S and Intomics A/S outside the submitted work
- 378 • SEA: None

79 **Figures**



80

81 **Figure 1:** Deriving the outcome variables. This exemplary admission is composed of three successive in-patient visits
 82 (i.e. the patient has been transferred twice represented by the arrows). The admission is eligible because it spans more
 83 than 24 hours and an eGFR ≤ 30 was measured before index. Here, apixaban was given while the patient's eGFR was
 84 ≤ 30 , but dose reduction rendered these administrations appropriate.



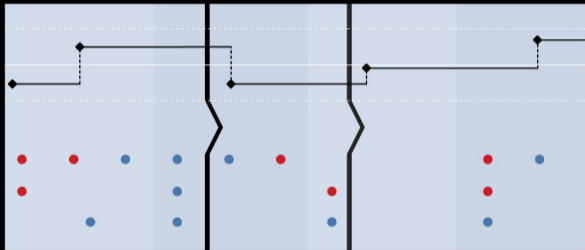
85

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

389 References

- 390 1. Saleem, A. & Masood, I. Pattern and Predictors of Medication Dosing Errors in Chronic
391 Kidney Disease Patients in Pakistan: A Single Center Retrospective Analysis. *PLoS One* **11**,
392 e0158677 (2016).
- 393 2. Hoffmann, F. et al. Renal Insufficiency and Medication in Nursing Home Residents. A Cross-
394 Sectional Study (IMREN). *Dtsch. Arztebl. Int.* **113**, 92–98 (2016).
- 395 3. Munar, M. Y. & Singh, H. Drug dosing adjustments in patients with chronic kidney disease.
396 *Am. Fam. Physician* **75**, 1487–1496 (2007).
- 397 4. Niedrig, D. et al. Development, implementation and outcome analysis of semi-automated
398 alerts for metformin dose adjustment in hospitalized patients with renal impairment.
399 *Pharmacoepidemiol. Drug. Saf.* **25**, 1204–1209 (2016).
- 400 5. Bernstein, J. M. & Erk, S. D. Choice of antibiotics, pharmacokinetics, and dose adjustments in
401 acute and chronic renal failure. *Med. Clin. North. Am.* **74**, 1059–1076 (1990).
- 402 6. Khare, A. K. Antibiotic dose adjustment in renal insufficiency. *Lancet* **340**, 1480 (1992).
- 403 7. Dorks, M., Allers, K., Schmiemann, G., Herget-Rosenthal, S. & Hoffmann, F. Inappropriate
404 Medication in Non-Hospitalized Patients With Renal Insufficiency: A Systematic Review. *J.*
405 *Am. Geriatr. Soc.* **65**, 853–862 (2017).
- 406 8. Getachew, H., Tadesse, Y. & Shibeshi, W. Drug dosage adjustment in hospitalized patients
407 with renal impairment at Tikur Anbessa specialized hospital, Addis Ababa, Ethiopia. *BMC*
408 *Nephrol.* **16**, 158 (2015).
- 409 9. Altunbas, G. et al. Renal Drug Dosage Adjustment According to Estimated Creatinine
410 Clearance in Hospitalized Patients With Heart Failure. *Am. J. Ther.* **23**, e1004-8 (2016).
- 411 10. Hillestad, R. et al. Can electronic medical record systems transform health care? Potential
412 health benefits, savings, and costs. *Health Aff. (Millwood)* **24**, 1103–1117 (2005).
- 413 11. Stewart, W. F., Shah, N. R., Selna, M. J., Paulus, R. A. & Walker, J. M. Bridging the
414 inferential gap: the electronic health record and clinical evidence. *Health Aff. (Millwood)* **26**,
415 w181-91 (2007).
- 416 12. Boussadi, A. et al. Validity of a clinical decision rule-based alert system for drug dose
417 adjustment in patients with renal failure intended to improve pharmacists' analysis of
418 medication orders in hospitals. *Int. J. Med. Inform.* **82**, 964–972 (2013).
- 419 13. Gawande, A. Why doctors hate their computers. *The New Yorker* (2018).
- 420 14. Baysari, M. T., Tariq, A., Day, R. O. & Westbrook, J. I. Alert override as a habitual behavior -
421 a new perspective on a persistent problem. *J. Am. Med. Inform. Assoc.* **24**, 409–412 (2017).
- 422 15. Kane-Gill, S. L. et al. Technologic Distractions (Part 1): Summary of Approaches to Manage
423 Alert Quantity With Intent to Reduce Alert Fatigue and Suggestions for Alert Fatigue Metrics.
424 *Crit. Care Med.* **45**, 1481–1488 (2017).
- 425 16. Jensen, T. B. et al. Content and validation of the Electronic Patient Medication module
426 (EPM)—the administrative in-hospital drug use database in the Capital Region of Denmark.
427 *Scand. J. Public Health* **0**, 1403494818760050 (2018).
- 428 17. Schmidt, M. et al. The Danish National Patient Registry: a review of content, data quality, and
429 research potential. *Clin. Epidemiol.* **7**, 449–490 (2015).
- 430 18. Moore, B. J., White, S., Washington, R., Coenen, N. & Elixhauser, A. Identifying Increased
431 Risk of Readmission and In-hospital Mortality Using Hospital Administrative Data: The
432 AHRQ Elixhauser Comorbidity Index. *Med. Care* **55**, 698–705 (2017).
- 433 19. Chollet, F. *Deep Learning with Python* (Manning Publications Co., New York, USA, 2018).
- 434 20. Akiba, T., Sano, S., Yanase, T., Ohta, T. & Koyama, M. Optuna: A Next-generation
435 Hyperparameter Optimization Framework. Preprint at <http://arxiv.org/abs/1907.10902> (2019).
- 436 21. Goodfellow, I., Bengio, Y. & Courville, A. *Deep Learning* (MIT Press, Cambridge (MA),
437 USA, 2016).
- 438 22. Efron, B. & Hastie, T. *Computer Age Statistical Inference: Algorithms, Evidence, and Data*
439 *Science* (Cambridge University Press, London, United Kingdom, 2016).
- 440 23. Hastie, T., Tibshirani, R. & Friedman, J. *The Elements of Statistical Learning: Data Mining,*
441 *Inference, and Prediction* (2nd ed., Springer, New York, 2009).
- 442 24. Chawla, N. V., Bowyer, K. W., Hall, L. O. & Kegelmeyer, W. P. SMOTE: Synthetic Minority
443 Over-sampling Technique. *J. Artif. Intell. Res.* **16**, 321–357 (2002).
- 444 25. Zhang, J. & Mani, I. kNN approach to unbalanced data distributions: a case study involving

- 445 information extraction. In: *Proceedings of the ICML'2003 Workshop on Learning from*
446 *Imbalanced Datasets* (2003).
- 447 26. Steyerberg, E. W. *Clinical prediction models: a practical approach to development,*
448 *validation, and updating* (Springer, New York, 2009).
- 449 27. Brown, L. D., Cai, T. T. & DasGupta, A. Interval Estimation for a Binomial Proportion.
450 *Statist. Sci.* **16**, 101–133 (2001).
- 451 28. Vickers, A. J. & Elkin, E. B. Decision curve analysis: a novel method for evaluating prediction
452 models. *Med. Decis. Making.* **26**, 565–574 (2006).
- 453 29. Kerr, K. F., Brown, M. D., Zhu, K. & Janes, H. Assessing the Clinical Impact of Risk
454 Prediction Models With Decision Curves: Guidance for Correct Interpretation and Appropriate
455 Use. *J. Clin. Oncol.* **34**, 2534–2540 (2016).
- 456 30. Lundberg, S. M. & Lee, S.-I. A Unified Approach to Interpreting Model Predictions. In:
457 *Advances in Neural Information Processing Systems 30* (2017).
- 458 31. Köster, J. & Rahmann, S. Snakemake—a scalable bioinformatics workflow engine.
459 *Bioinformatics* **28**, 2520–2522 (2012).
- 460 32. Hernandez-Boussard, T., Bozkurt, S., Ioannidis, J. P. A. & Shah, N. H. MINIMAR (MINimum
461 Information for Medical AI Reporting): Developing reporting standards for artificial
462 intelligence in health care. *J. Am. Med. Inform. Assoc.* (2020).
- 463 33. Collins, G. S., Reitsma, J. B., Altman, D. G. & Moons, K. G. M. Transparent Reporting of a
464 multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD): The
465 TRIPOD Statement. *Ann. Intern. Med.* **162**, 55–63 (2015).
- 466 34. Saad, R., Hallit, S. & Chahine, B. Evaluation of renal drug dosing adjustment in chronic
467 kidney disease patients at two university hospitals in Lebanon. *Pharm. Pract. (Granada)* **17**,
468 (2019).
- 469 35. Blix, H. S. et al. The majority of hospitalised patients have drug-related problems: results from
470 a prospective study in general hospitals. *Eur. J. Clin. Pharmacol.* **60**, 651–658 (2004).
- 471 36. Andreu Cayuelas, J. M. et al. Kidney function monitoring and nonvitamin K oral anticoagulant
472 dosage in atrial fibrillation. *Eur. J. Clin. Invest.* **48**, e12907 (2018).
- 473 37. Seiberth, S. et al. Correct use of non-indexed eGFR for drug dosing and renal drug-related
474 problems at hospital admission. *Eur. J. Clin. Pharmacol.* (2020).
- 475 38. Breton, G. et al. Inappropriate drug use and mortality in community-dwelling elderly with
476 impaired kidney function—the Three-City population-based study. *Nephrol. Dial. Transplant.*
477 **26**, 2852–2859 (2011).
- 478 39. Chang, F., O'Hare, A. M., Miao, Y. & Steinman, M. A. Use of Renally Inappropriate
479 Medications in Older Veterans: A National Study. *J. Am. Geriatr. Soc.* **63**, 2290–2297 (2015).
- 480 40. Parameswaran Nair, N. et al. Hospitalization in older patients due to adverse drug reactions -
481 the need for a prediction tool. *Clin. Interv. Aging* **11**, 497–505 (2016).
- 482 41. Kalender-Rich, J. L., Mahnken, J. D., Wetmore, J. B. & Rigler, S. K. Transient impact of
483 automated glomerular filtration rate reporting on drug dosing for hospitalized older adults with
484 concealed renal insufficiency. *Am. J. Geriatr. Pharmacother.* **9**, 320–327 (2011).
- 485 42. Won, H.-J. et al. Evaluation of medication dosing errors in elderly patients with renal
486 impairment. *Int. J. Clin. Pharmacol. Ther.* **56**, 358–365 (2018).
- 487 43. Carey, I. M. et al. What Factors Predict Potentially Inappropriate Primary Care Prescribing in
488 Older People? *Drug Aging* **25**, 693–706 (2008).
- 489 44. Steyerberg, E. W. & Harrell, F. E. J. Prediction models need appropriate internal, internal-
490 external, and external validation. *J Clin Epidemiol* **69**, 245–247 (2016).
- 491 45. Wynants, L. et al. Prediction models for diagnosis and prognosis of covid-19: systematic
492 review and critical appraisal. *BMJ* **369**, (2020).
- 493 46. Yusuf, M. et al. Reporting quality of studies using machine learning models for medical
494 diagnosis: a systematic review. *BMJ Open* **10**, (2020).
- 495 47. Eppenga, W. L. et al. Drug therapy management in patients with renal impairment: how to use
496 creatinine-based formulas in clinical practice. *Eur. J. Clin. Pharmacol.* **72**, 1433–1439 (2016).
- 497 48. Rule, A. D. & Glasscock, R. J. GFR estimating equations: getting closer to the truth? *Clin. J.*
498 *Am. Soc. Nephrol.* **8**, 1414–1420 (2013).
- 499 49. Corsonello, A. et al. Estimating renal function to reduce the risk of adverse drug reactions.
500 *Drug Saf.* **35 Suppl 1**, 47–54 (2012).
- 501 50. Levey, A. S. et al. A new equation to estimate glomerular filtration rate. *Ann. Intern. Med.*
502 **150**, 604–612 (2009).



No. inappropriate doses

Time at risk

Outcome

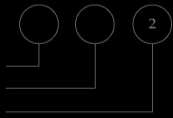
daily rate

1

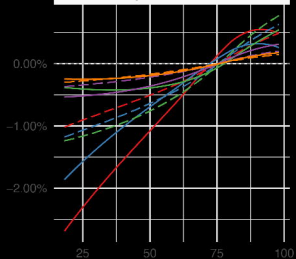
1

0

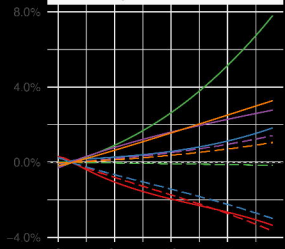
2



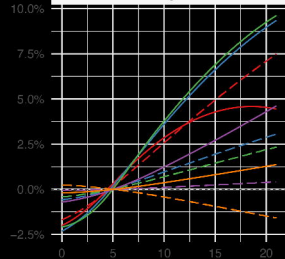
Age at admission



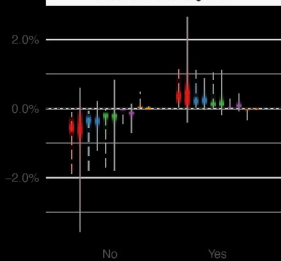
No. previous admissions



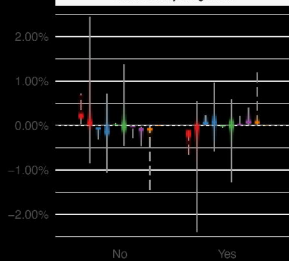
No. distinct drugs before index



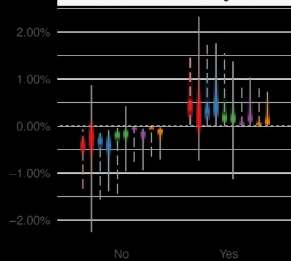
Cardiovascular diagnosis



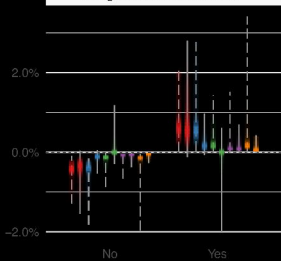
Genitourinary diagnosis



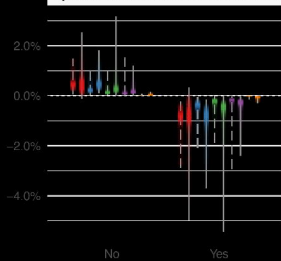
Endocrine-metabolic diagnosis



Analgesics used before index



Systemic antibacterials used before index



Diuretics used before index

