

Contribution of temporal data to predictive performance in 30-day readmission of morbidly obese patients

Petra Povalej Brzan^{1,2}, Zoran Obradovic³, Gregor Stiglic^{Corresp. 1,2}

¹ Faculty of Health Sciences, University of Maribor, Maribor, Slovenia

² Faculty of Electrical Engineering and Computer Science, University of Maribor, Maribor, Slovenia

³ Center for Data Analytics and Biomedical Informatics, Temple University, Philadelphia, PA, United States

Corresponding Author: Gregor Stiglic

Email address: gregor.stiglic@um.si

Background. Reduction of readmissions after discharge represents an important challenge for many hospitals and has attracted the interest of many researchers in the past few years. Most of the studies in this field focus on building cross-sectional predictive models that aim to predict the occurrence of readmission within 30-days based on information from the current hospitalization. The aim of this study is demonstration of predictive performance gain obtained by inclusion of information from historical hospitalization records among morbidly obese patients.

Methods. The California Statewide inpatient database was used to build regularized logistic regression models for prediction of readmission in morbidly obese patients ($n=18,881$). Temporal features were extracted from historical patient hospitalization records in a one year timeframe. Five different datasets of patients were prepared based on the number of available hospitalizations per patient. Sample size of the five datasets ranged from 4,787 patients with more than 5 hospitalizations to 20,521 patients with at least two hospitalization records in one year. 10-fold cross validation was repeated 100 times to assess the variability of the results. Additionally, random forest and extreme gradient boosting were used to confirm the results.

Results. Area under the ROC curve increased significantly when including information from up to three historical records on all datasets. The inclusion of more than three historical records was not efficient. Similar results can be observed for Brier score and PPV value. The number of selected predictors corresponded to the complexity of the dataset ranging from an average of 29.50 selected features on the smallest dataset to 184.96 on the largest dataset based on 100 repetitions of 10-fold cross-validation.

Discussion. The results show positive influence of adding information from historical hospitalization records on predictive performance using all predictive modeling techniques used in this study. We can conclude that it is advantageous to build separate readmission prediction models in subgroups of patients with more hospital admissions by aggregating information from up to three previous hospitalizations.

1 **Contribution of temporal data to predictive performance in 30-day**
2 **readmission of morbidly obese patients**

3

4 P. Povalej Brzan^{1,2}, Z. Obradovic³, G. Stiglic^{1,2}

5 ¹ Faculty of Health Sciences, University of Maribor, Maribor, Slovenia

6 ² Faculty of Electrical Engineering and Computer Science, University of Maribor, Maribor,
7 Slovenia

8 ³ Center for Data Analytics and Biomedical Informatics, Temple University, Philadelphia, PA,
9 USA

10

11 Corresponding Author:

12 Petra Povalej Brzan¹,

13 Zitna ulica 15, Maribor, SI-2000, Slovenia

14 Email address: petra.povalej@um.si

15

16

17

18

19

20

21

22 **ABSTRACT**

23 **Background.** Reduction of readmissions after discharge represents an important challenge for
24 many hospitals and has attracted the interest of many researchers in the past few years. Most of
25 the studies in this field focus on building cross-sectional predictive models that aim to predict the
26 occurrence of readmission within 30-days based on information from the current hospitalization.
27 The aim of this study is demonstration of predictive performance gain obtained by inclusion of
28 information from historical hospitalization records among morbidly obese patients.

29 **Methods.** The California Statewide inpatient database was used to build regularized logistic
30 regression models for prediction of readmission in morbidly obese patients (n=18,881).
31 Temporal features were extracted from historical patient hospitalization records in a one year
32 timeframe. Five different datasets of patients were prepared based on the number of available
33 hospitalizations per patient. Sample size of the five datasets ranged from 4,787 patients with
34 more than 5 hospitalizations to 20,521 patients with at least two hospitalization records in one
35 year. 10-fold cross validation was repeated 100 times to assess the variability of the results.

36 Additionally, random forest and extreme gradient boosting were used to confirm the results.

37 **Results.** Area under the ROC curve increased significantly when including information from up
38 to three historical records on all datasets. The inclusion of more than three historical records was
39 not efficient. Similar results can be observed for Brier score and PPV value. The number of
40 selected predictors corresponded to the complexity of the dataset ranging from an average of
41 29.50 selected features on the smallest dataset to 184.96 on the largest dataset based on 100
42 repetitions of 10-fold cross-validation.

43 **Discussion.** The results show positive influence of adding information from historical
44 hospitalization records on predictive performance using all predictive modeling techniques used
45 in this study. We can conclude that it is advantageous to build separate readmission prediction
46 models in subgroups of patients with more hospital admissions by aggregating information from
47 up to three previous hospitalizations.

48 Introduction

49 Hospital readmission prediction models have been widely studied and deployed worldwide (Zhu
50 et al, 2015; Hao et al, 2015; Stiglic et al, 2015). Different types of prediction models were
51 proposed to predict and potentially prevent hospital readmissions. As described in a review study
52 by Kansagara et al. (2011), we can divide the proposed predictive models into two groups – i.e.
53 models relying on retrospective administrative data and models using real-time administrative
54 data. While the second group usually focuses on data that is collected during hospitalization, the
55 first group of models relies on retrospective data. Although many studies include information on
56 prior hospitalizations to improve the predictive performance of readmission prediction models,
57 they usually do not provide evidence on the level of their contribution to predictive performance.

58 He et al. (2014) demonstrated the importance of a logistic regression predictor representing the
59 number of prior hospitalizations in the past 5 years. This simple variable was selected as
60 significant in both a general and a specific chronic pancreatitis subgroup based predictive model.
61 Walsh and Hripesak (2014) compared predictive performance of readmission prediction models
62 for specific subgroups of patients. Their results show a strong gain in predictive performance
63 when laboratory data and visit history data is included in the predictive model development. A
64 study by Shahn et al. (2015) incorporates relative temporal relationships among multiple health
65 events in the space of predictors to build a Random Relational Forest (RRF) based classifier
66 originally developed in the context of speech recognition. RRF generates informative labeled
67 graphs representing temporal relations among health events at the nodes of randomized decision
68 trees. Although the target of a study by Shahn et al. does not include readmission classification,
69 but focuses on predicting strokes in patients with prior diagnoses of Atrial Fibrillation, it
70 demonstrates the importance of temporal information to achieve meaningful improvements in
71 predictive performance. Similarly, the use of temporal information from electronic health
72 records (EHRs) for prediction of Anastomosis Leakage was demonstrated in a study by
73 Soguero-Ruiz et al. (2016). The predictive performance gain was proven when combining the
74 data from heterogenous data sources (extracted free text, blood samples values, and patient vital
75 signs).

76 Morbidly obese patients represent one of the most complex populations in healthcare systems
77 and are often related to higher treatment costs (Kadry et al., 2014). As outlined by Incavo et al.
78 (2014), hospitals need extra personnel and equipment to lift and transport morbidly obese
79 patients. Additionally, such patients tend to have above average number of comorbidities,
80 including chronic diseases. Consequently, hospital staff is affected from the heavy lifting of
81 patients and healthcare providers need to provide more care with the same reimbursement (Choi
82 & Brings, 2016).

83 In this study, we focus on the following research question: Does the inclusion of additional
84 information from historical patient hospitalization records improve the predictive performance of
85 30-day readmission models? Hospitalization claims data from morbidly obese patients were used
86 to build a readmission prediction model. We hypothesized that a significant improvement in
87 predictive performance can be obtained by inclusion of new predictors based on previous
88 hospitalizations. Least Absolute Selection and Shrinkage Operator (Lasso) regularization

89 (Tibshirani, 2011) was used to allow interpretation of the results by observing the frequency of
90 the predictor inclusion (Stiglic et al., 2013) in the Lasso logistic regression model. We further
91 explored the complexity of built models measured as number of selected features in addition to
92 comparing the results using advanced predictive modeling techniques.

93

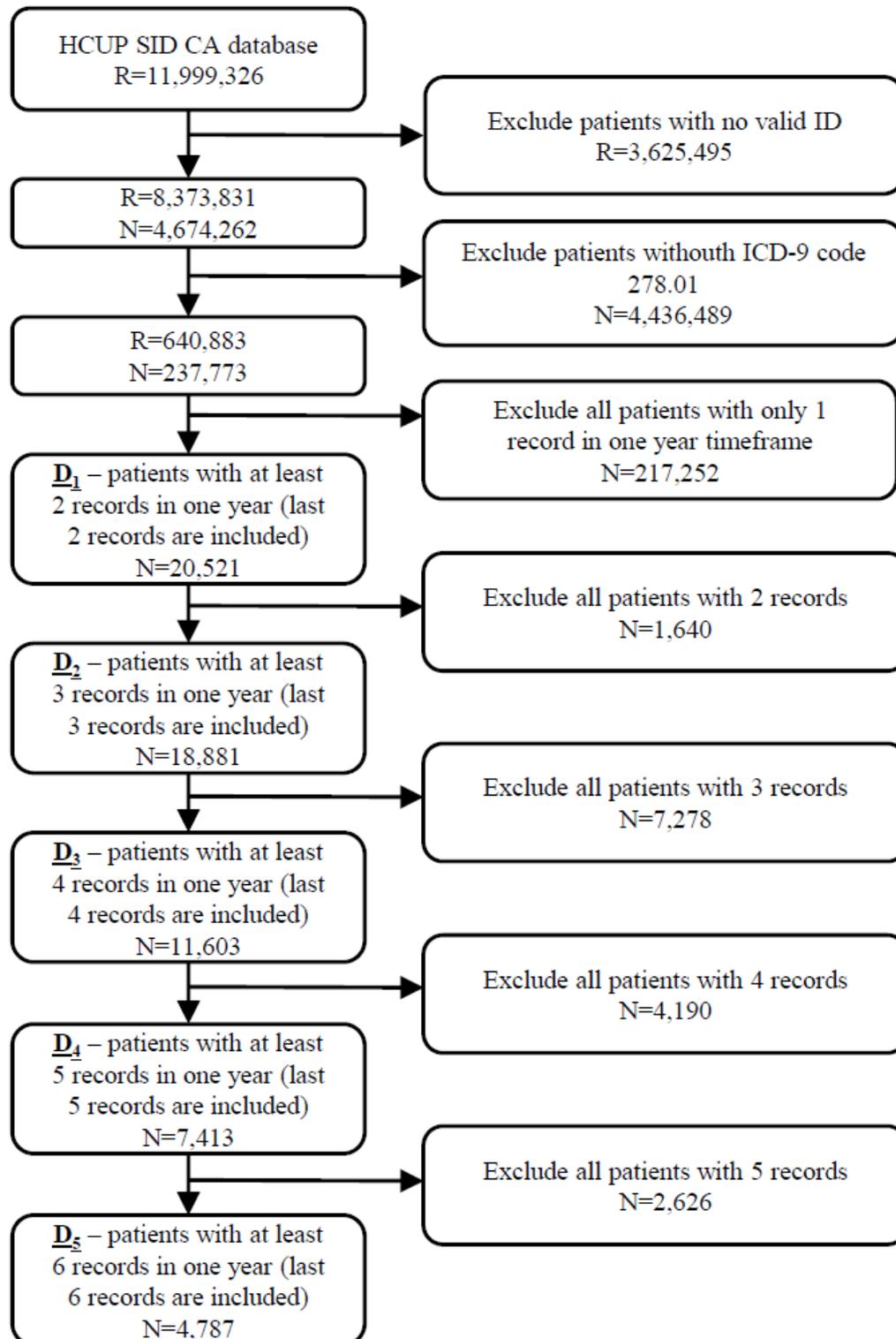
94 **Materials & Methods**

95 **Dataset**

96 The 11,889,326 hospitalization records from the Healthcare Cost and Utilization Project (HCUP)
97 State Inpatient Database for California (SID CA) for the years 2009 to 2011 were used in the
98 study. Each hospitalization record includes demographic information about the patient (age, birth
99 year, sex, race, etc.), information about the hospital stay (length of stay, total charges, type of
100 payment, discharge month, survival information, scheduled visit, etc.), primary diagnosis, up to
101 25 diagnoses and up to 25 procedures performed on a patient during hospitalization. For the
102 purpose of protecting patient privacy the hospitalization records are anonymized and a unique ID
103 value, which enables tracking the patients through several years, is used for each patient. The
104 diagnoses and procedures are described in ICD-9 codes and CSS codes. In our experiment the
105 ICD-9 codes were used.

106 The HCUP SID CA dataset was filtered based on the following inclusion criteria

107 The initial database was first filtered on hospitalization records with valid patient ID on
108 8,373,831 hospitalization records from 4,674,262 patients. In the next step 237,773 patients
109 (640,883 hospitalization records) with ICD-9 code 278.01 (morbid obesity) in at least one
110 hospitalization were extracted from the database. The most recent hospitalization with ICD-9
111 code 278.01 from 2010-2011 was selected for each patient. Additionally, all historical records in
112 a one year timeframe were added for each patient. In further experiments we position the patient
113 in one hospitalization before the last one (index hospitalization) with the purpose of predicting
114 the last hospitalization. All scheduled predicted admissions were excluded. For that purpose,
115 five different datasets of patients were prepared based on the number of available
116 hospitalizations per patient. In further text let D_1 denote the dataset of 20,521 patients with at
117 least two hospitalization records in one year timeframe. Similarly, D_2 consists of 18,881 patients
118 with more than two hospitalizations, D_3 dataset represents 11,603 patients with more than 3
119 hospitalizations and D_4 7,413 patients with more than 4 hospitalization records. The smallest
120 dataset D_5 includes 4,787 patients with more than 5 hospitalization records in a one year
121 timeframe (Figure 1).



122

123 *Figure 1: Extraction of datasets from the original HCUP SID CA database including number of*
 124 *records (R) and patients (N).*

125 The percentage of patients readmitted in 30 days increases from 34.36% in the largest dataset D_1
 126 to 43.85% in the smallest dataset of the most complex patients (D_5). Note that the purpose of this
 127 study is to evaluate the contribution of additional information from historical records on the
 128 prediction of the 30-days readmission of morbid obesity patients. Therefore we used the patients
 129 that were readmitted at least once, because these patients represent a cohort of patients that are
 130 costly and have a strong impact on hospital and global performance indicators such as waiting
 131 lists, mortality, planned care, etc.

132 All available diagnoses (7,196) and procedures (2,488) in the year 2009 were first ordered by
 133 frequency to arbitrarily select a cut-off value for selection of diagnoses and procedures used in
 134 later stages. Due to the long tail distribution of frequencies the cut-off was set to 3%, resulting in
 135 the final set of 217 diagnoses and 75 procedures that were further used as dichotomous variables.

136

137 **Statistical Analysis**

138 The following features from consecutive hospital records were obtained for each patient on the
 139 discharge day: total number of hospital days in all hospitalizations, total number of
 140 hospitalizations, total number of procedures in all hospitalizations, total number of chronic
 141 diseases in all hospitalizations, mean number of chronic diseases from all hospitalizations, mean
 142 number of hospital days from all hospitalizations, mean number of procedures from all
 143 hospitalizations. Additionally, a total number of hospital days, total number of hospitalizations,
 144 mean number of chronic diseases, mean number of hospital days, and mean number of
 145 procedures were calculated for the last 30/60/90/180/270 days prior to index hospitalization.

146 The number of occurrences of each diagnosis and procedures from all historical patient records
 147 were added as new features for each patient hospitalization record. The patient hospitalization
 148 record used in further analysis therefore consisted of patient's demographic information (age,
 149 birth year, sex, race, etc.) and a set of features from historical hospital records obtained as
 150 described above.

151 **Predictive modeling**

152 A generalized linear model via penalized maximum likelihood combining L1-norm (lasso) and
 153 L2-norm (ridge) regularization was used as defined by Friedman, Hastie and Tibshirani (2013).
 154 The process of predictive modeling can be significantly simplified by using regularized logistic
 155 regression methods, since the feature selection step is integrated in the process.

156 A generalized linear model via penalized maximum likelihood can be described as:

$$157 \min_{\beta_0, \beta} \frac{1}{N} \sum_{i=1}^N w_i l(y_i \beta_0 + \beta^T x_i) + \lambda \left[(1 - \alpha) \|\beta\|_2^2 / 2 + \alpha \|\beta\|_1 \right]$$

158 where i represents observations and its negative log-likelihood contribution is noted as $l(y, \eta)$.
 159 The regularization path λ is computed for a grid of values for the regularization parameter which
 160 controls the overall strength of the penalty and is in our case calculated for lasso ($\alpha = 1$), since
 161 our initial experiments did not show any significant gain with elastic-net.
 162

163 Additionally, the results of the lasso predictive model were compared to random forest (using
164 200 decision trees) and XGBoost method using default parameter values and the same
165 experimental settings. Random forest is often applied in healthcare prediction problems since it
166 offers a high level of robustness (Zhou et al., 2016). XGBoost is a powerful implementation of
167 gradient boosting first proposed by Friedman (Friedman, 2001) designed for speed and
168 performance.

169

170 [Experimental Setting](#)

171 Each experimental run included 10-fold cross-validation repeated 100 times with the same
172 random number generator seed values in all experiments. By using 100 repetitions of 10-fold
173 cross validation, we were able to assess the variance of the results under different cross-
174 validations.

175 The following experimental setup was used:

- 176 1. Select the dataset of patients with at least R hospitalization records in a one year
177 timeframe.
- 178 2. Predict 30-days readmission for the last hospitalization record from the temporal features
179 for 1, 2, ..., R historical records, using 10-fold cross-validation.
- 180 3. Repeat step 2 for 100 times.
- 181 4. Repeat steps 1 - 3 for $R = \{1, 2, 3, 4, 5\}$

182 Predictions were obtained for each validation sample using the model derived on the derivation
183 sample. The predictive accuracy of each model was summarized by area under the receiver
184 operating characteristics (ROC) curve (AUC), where 1 represents perfect predictive performance
185 and 0.5 represents random performance. In addition to AUC, Brier's score, which allows more
186 efficient evaluation of probabilistic predictions, was used. In contrast to AUC, lower Brier score
187 represents better predictive performance. Due to imbalanced datasets, sensitivity, specificity,
188 positive predictive value (PPV) and negative predictive value (NPV) were calculated. All
189 evaluation measures were calculated for each 10-fold run on validation samples and then the
190 mean value was calculated. The average of 10-fold mean over 100 repetitions and 95%
191 confidence interval for each experiment is presented in all figures.

192 The independent samples t test was used to test the difference in mean evaluation measures
193 between two samples. The ANOVA test was used for comparing the mean values for different
194 datasets. P value less than 0.05 was considered statistically significant.

195 All experiments were conducted using the R language and environment for statistical computing
196 (R Core Team, 2016).

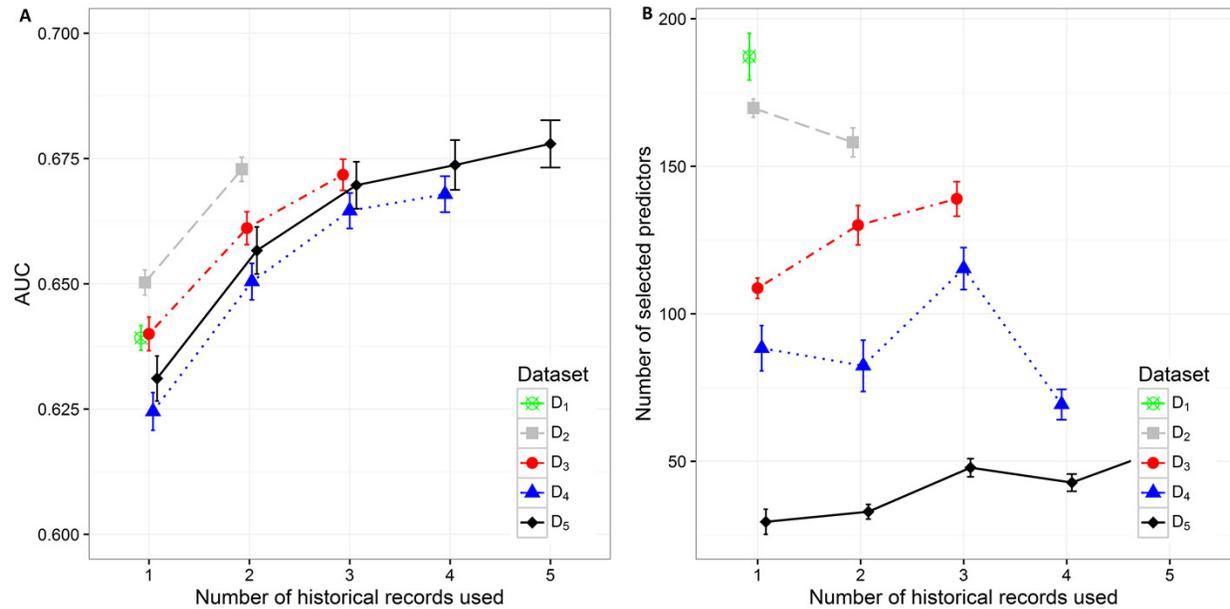
197

198

199 Results

200 The results in the text are presented as mean AUC (95% CI). The detailed results are described in
201 Supplement 1. Number of selected features and AUC on different datasets are presented in
202 Figure 2. The AUC on the smallest dataset (D_5) of 4,787 patients with at least 5 historical
203 hospitalizations in a one-year timeframe increased from 0.631 (0.627, 0.636), when only current
204 hospitalization was considered, to 0.657 (0.652, 0.661) when temporal features from the last 2
205 hospitalizations were included. Information from an additional (third) historical hospitalization
206 again significantly improves AUC, which increases to 0.670 (0.665, 0.674). However, when the
207 fourth and fifth historical hospitalization is added, the AUC increases only for 0.04 in each case
208 (AUC=0.674 (0.669,0.679) for 4 hospitalizations and AUC=0.678 (0.673,0.683) for 5
209 hospitalizations). Similar trends can be seen when we reduce the required minimum number of
210 historical hospitalization records in a one year timeframe to 4, 3 and 2 records (D_4 , D_3 , D_2) and
211 consequently increase the sample size. The most significant difference in mean AUC can be
212 observed when at least one historical hospitalization record is included. The AUC increases from
213 0.650 (0.648, 0.653) to 0.673 (0.670, 0.675) on the largest dataset D_2 . Similar observations can
214 be made on D_3 , D_4 , and D_5 . As seen on the smallest dataset D_5 , the statistical difference in AUC
215 between the classifiers with two hospitalization records and three hospitalization records
216 considered can be observed on bigger datasets (D_4 and D_3) as well. For databases D_4 and D_5 the
217 experiment was repeated by adding a fourth hospitalization record. The increase in AUC was
218 only 0.03 for D_4 (from 0.665 (0.661,0.668) to 0.668 (0.664,0.671)) and 0.04 for D_5 (from 0.670
219 (0.665,0.674) to 0.678 (0.673,0.683)).

220 By increasing the dataset size, the percentage of positive samples (patients that were re-
221 hospitalized in 30 days) is decreasing. This is the most plausible reason for higher initial AUC
222 when only current hospitalization is considered for prediction of re-hospitalization on bigger
223 datasets. The highest mean value of AUC (0.650 (0.648, 0.653)) can be observed on the dataset
224 D_2 with 18,881 patients and the lowest AUC (0.625 (0.621, 0.628)) on the dataset D_4 with 7,413
225 patients. When adding one more historical hospitalization record, the mean AUC value ranges
226 between 0.650 (0.647, 0.654) on D_4 and 0.673 (0.670, 0.675) on D_2 . The difference is significant
227 between all datasets except between D_3 (AUC=0.661 (0.658, 0.664)) and D_5 (AUC=0.657
228 (0.652, 0.661)) and between D_4 (AUC=0.650 (0.647, 0.654)) and D_5 (AUC=0.657 (0.652,
229 0.661)). However, when the third historical record is added, the mean AUC value ranges
230 between 0.665 (0.661, 0.668) on D_4 and 0.672 (0.669, 0.675) on D_3 . The difference in mean
231 AUC is significant between D_3 and D_4 . On the other hand, adding additional historical records
232 (more than 3) does not result in significantly better AUC.

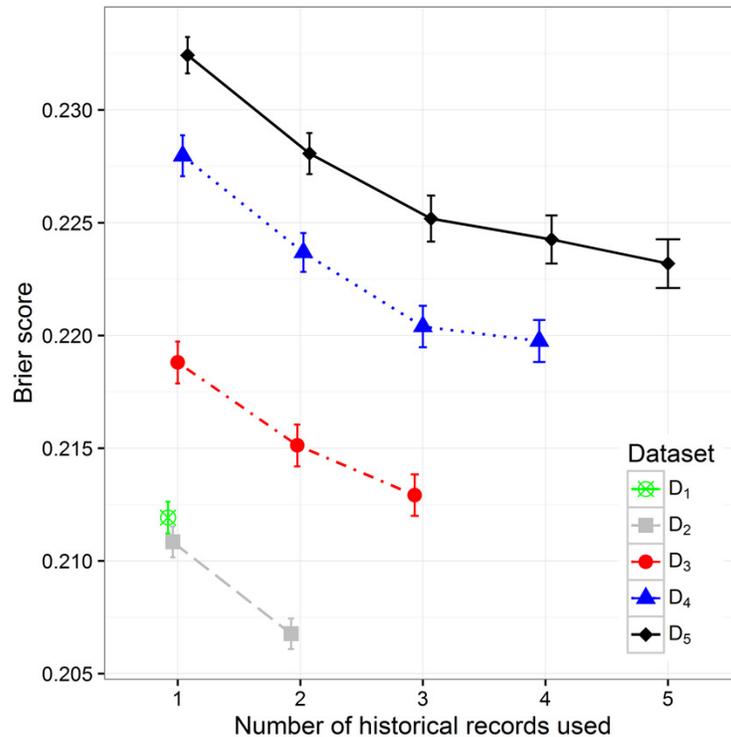


233

234 *Figure 2: Mean value of the (A) AUC and (B) number of selected predictors with corresponding*
 235 *95% CI on different datasets with different number of historical hospitalization records.*

236 By increasing the size of the dataset we can observe the decreasing trend of the Brier score
 237 (mean squared error), which can be expected considering that with the size of the dataset the
 238 percentage of positive samples decreases (Figure 3). When observing separate datasets, we can
 239 also see the decrease in Brier score when including additional historical hospitalization records.
 240 The decrease is statistically significant in models that include from one to three historical
 241 records. The inclusion of more historical records does not further decrease the score in a
 242 significant way.

243

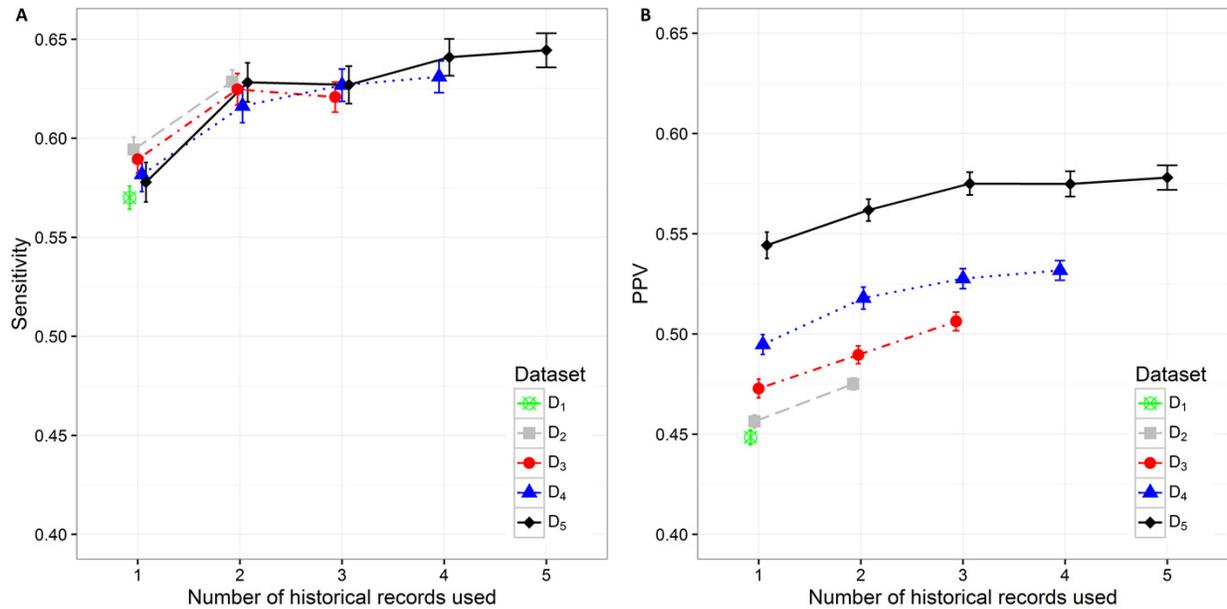


244

245 *Figure 3: The Brier score mean value and 95% CI on different datasets with different number of*
 246 *historical hospitalization records.*

247 The sensitivity value and PPV value trends can be observed in Figure 4 (A) and (B). The lowest
 248 sensitivity (0.570 (0.564, 0.576)) and PPV value (0.448 (0.445, 0.452)) are shown on the biggest
 249 dataset D₁. The PPV value increases significantly when decreasing the size of the dataset. The
 250 highest PPV value (0.544 (0.538, 0.551)) when only current hospitalization data are included is
 251 therefore achieved on the smallest dataset D₅. The sensitivity and PPV increase also when
 252 additional information about historical hospitalizations is included. The highest difference can be
 253 observed when at least one additional hospitalization record is added to the current
 254 hospitalization data. In all datasets this difference is significant. When adding the third
 255 hospitalization record the mean PPV value statistically still increases on all datasets, however the
 256 mean sensitivity value does not change significantly. The highest mean PPV value of 0.578
 257 (0.572, 0.584) was achieved in the smallest dataset (D₅) when the data from all five historical
 258 hospitalization records was included. However, it does not change significantly from the PPV
 259 value for D₅ when only three historical records are considered (0.575 (0.569, 0.581)).

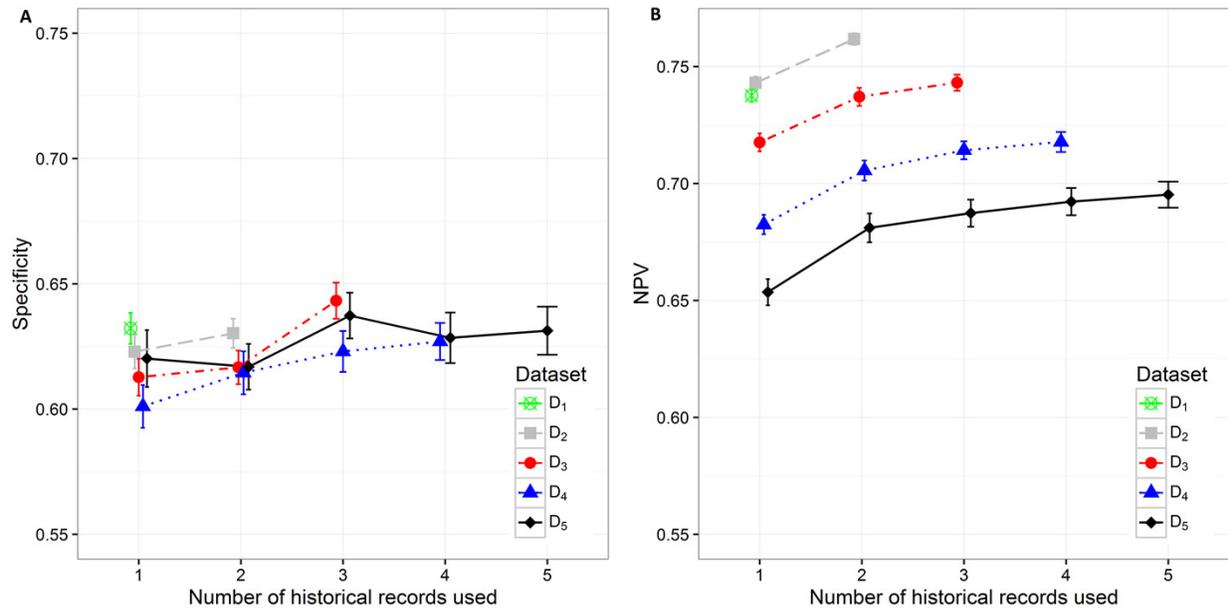
260 The sensitivity changes significantly only when adding one historical record to the current
 261 hospitalization data on all datasets. The differences vary from 0.034 on D₂ to 0.050 on D₅. The
 262 highest sensitivity value of 0.644 (0.636, 0.653) was achieved on the smallest dataset D₅ when
 263 using the data from all five hospitalizations. However, the observation has to be made that it is
 264 not significantly different from the sensitivity value for D₅ when only two historical records are
 265 considered (0.628 (0.618, 0.638)).



266

267 *Figure 4: Mean value of the (A) sensitivity and (B) PPV with corresponding 95% CI on different*
 268 *datasets with different number of historical hospitalization records.*

269 Specificity varies from minimum 0.601 (0.593, 0.610) on D₄ with only current hospitalization
 270 record considered to maximum 0.643 (0.636, 0.650) on D₃ with historical information on all
 271 three records included. The increase in specificity while adding new historical information can
 272 be observed on all datasets except in D₅. As already observed in other measures, the statistically
 273 significant increase is shown between one and three hospitalization records included on all
 274 datasets. A very similar trend can be observed for NPV values, which also increase when adding
 275 historical information from previous hospitalizations (Figure 5).



276

277 *Figure 5: Mean value of the (A) specificity and (B) NPV with corresponding 95% CI on different*
 278 *datasets with different number of historical hospitalization records.*

279 The number of selected predictors was considered as a measure of classifier complexity, where
 280 lower complexity is considered as positive when interpretation of results from the medical point
 281 of view is needed. Figure 2 (B) shows that the lowest complexity in the terms of the mean
 282 number of selected predictors can be observed on the smallest dataset (D₅) where mean value
 283 increases from 29.50 (25.264, 33.776) to 54.98 (51.864, 58.096) with addition of previous
 284 hospitalization information. A similarly increasing trend can be observed on D₃ (from 108.700
 285 (105.242, 112.158) to 138.950 (133.100, 144.800) features). The decrease in the number of
 286 selected predictors when adding historical records is shown on D₂ (from 169.7 (166.629,
 287 172.811) to 158.130 (153.233, 163.027)) and partly on D₄. On D₄ with three historical records
 288 included the mean number of selected predictors increases from 82.390 (73.700, 91.080) (for two
 289 historical records) to 115.330 (108.193, 122.467) and then decreases again to 69.28 (64.120,
 290 74.440) when all 4 historical records are included. The highest mean number of selected features
 291 was obtained from the largest dataset D₁ (184.960 (179.257, 195.083))

292 Additional experiments using random forest and XGBoost were performed. The results
 293 (Supplement 2, 3) show similar trends as can be seen in the figures above (Figure 2-5). Detailed
 294 results including all performance metrics with corresponding confidence intervals for each
 295 predictive modeling technique can be found in Supplement 1.

296

297 Discussion

298 Although claims data are very limited in information regarding a specific patient, claims
299 databases usually contain large volumes of data and are accessible to researchers. In this study,
300 we demonstrated how information on historical hospitalizations influences the predictive
301 performance of a classifier built with a Lasso regularized generalized linear model on morbidly
302 obese patients.

303 The initial idea was to show that more complex patients, which are more frequently hospitalized,
304 are more similar to each other and therefore separate prediction models should be used for them.
305 In order to make a fair comparison, five datasets were constructed from the initial database based
306 on the minimum number of historical hospitalizations per patient (D_1 to D_5). Then the model for
307 30-days readmission risk prediction was built using the data from one or more historical
308 hospitalization records for each patient in the dataset. Each model was evaluated using a
309 validation sample. The robustness of the models was tested using 10-fold CV, which was
310 repeated 100 times.

311 The complexity of models was expressed by the number of selected predictors. Generally, higher
312 mean AUC, sensitivity, specificity, PPV and NPV was achieved when more historical
313 information was added. However, the initial models using only the data from the current
314 hospitalization performed better on less restrictive datasets.

315 It is rather difficult to provide a guideline as to how many historical records should be included
316 in a predictive model. We expect the optimal model to have an agreement between the
317 performance measures (for example high AUC and low Brier score). However, the complexity of
318 the model should also be considered. The simplest model should be selected for practical
319 purposes of model interpretability.

320 The presented analysis of the influence of adding historical data on the predictive performance of
321 a classifier provided very useful insights. When considering all performance measures and also
322 the complexity of the classifiers, one can observe that it is highly important to build the
323 prediction model using at least information from the last two hospitalization records if available.
324 Including data from more than three historical records did not improve the performance of
325 classifiers significantly. Therefore we can conclude that the inclusion of data for more than three
326 previous hospitalization records per patient is not required. Results obtained with random forest
327 and XGBoost predictive models (Supplement 1) confirm the same trends that can be seen in all
328 performance metrics for lasso logistic regression (Figure 2-5).

329 The analysis of model complexity on the basis of the number of selected predictors shows that
330 building separate models for patients with a higher number of hospitalizations is reasonable,
331 since the models built on more restrictive (homogeneous) datasets gain on simplicity and
332 predictive performance. On the other hand, one can expect higher stability of prediction models
333 on larger sample size (Figure 2a, 4 and 5). Although we obtained the best results, measured by
334 AUC, using random forest, it should be noted that these models used significantly more features
335 in comparison to the other two techniques. XGBoost achieved the best Brier score performance.

336 However, as in random forest the interpretability of XGBoost models is very limited compared
337 to lasso logistic regression.

338 Readmissions represent one of the most important indicators of quality of care in the healthcare
339 environment, resulting in great economic impact. Readmission rates within 30 days are reported
340 as high as 19.6 %, including approximately 76 % of preventable readmissions, resulting in a
341 reduction of about \$25 billion annually in the US (Behara, 2013). Therefore a robust and
342 efficient solution to predict readmissions contributes to higher quality of care and reduces costs.
343 Moreover, this work focuses on a specific subgroup of morbidly obese patients where nurses and
344 nursing assistants manually lifting patients experience the highest rates of back and shoulder
345 musculoskeletal injuries (Choi, 2016), such that an effective predictive model multiplies the
346 benefits on both the patient and hospital staff side.

347 The limitations of this study are related to the characteristics of the claims data, where only a
348 limited set of features is available. Additionally, it is important to realize that some diagnoses are
349 not correctly recorded due to the influence of health insurance policies on costs related to
350 different diagnoses and procedures. In the case of more complete data (full EHR records) the
351 general predictive performance is expected to be higher, however the contribution of the
352 historical data would still be significant. On the other hand the availability and volume of claims
353 data were more important in order to achieve the aim of this study.

354

355 **Conclusions**

356 Existing literature on readmission prediction based on claims data shows relatively low
357 predictive performance. Therefore this study does not only focus on improvement of the
358 predictive performance, but also includes the analysis of how historical information about the
359 patient influences the predictive performance and complexity of the predictive model.

360 The presented results show positive influence, which reflects in statistically significant increase
361 in AUC, sensitivity, specificity, PPV and NPV values and decrease in Brier score when adding
362 up to three historical hospitalization records.

363 As expected, the number of selected predictors increases with the size of the dataset. The
364 homogeneity of the patient's dataset increases when tightening the criteria regarding minimum
365 number of hospitalizations per patient in a one year timeframe.

366 This study can be usefull for data-scientists and software engineers developing similar prediction
367 models using data from hospitalization records. Most already developed readmission prediction
368 models are focused on readmission prediction based on the data from only one hospitalization
369 record and do not include historical records even if they are available. However, from this study
370 we can conclude that it is advantageous to generate separate models for predicting readmissions
371 on more complex patients including the data from their historical hospitalizations, since they
372 form a more homogenous dataset and consequently present a more complex classification
373 problem.

374 **References**

375

376 Behara, R., Agarwal, A., Fatteh, F., Furht, B. 2013. Predicting Hospital Readmission Risk for
377 COPD Using EHR Information. *Handbook of Medical and Healthcare Technologies*, pp. 297–
378 308. Springer, New York.

379 Choi SD, Brings K. 2016. Work-related musculoskeletal risks associated with nurses and nursing
380 assistants handling overweight and obese patients: A literature review. *Work*, 53(2), pp.439-448.

381 Delen D, Oztekin A, Tomak L. 2012. An analytic approach to better understanding and
382 management of coronary surgeries. *Decision Support Systems*. 52(3) (2012) 698-705.

383 Friedman J, Hastie T, Tibshirani R. 2013. Glmnet: Lasso and elastic-net regularized generalized
384 linear models, R package, version 1.9-5. 2013.

385 Gligorijevic D, Stojanovic J, Obradovic Z. 2015. Improving Confidence while Predicting Trends
386 in Temporal Disease Networks, 4th Workshop on Data Mining for Medicine and Healthcare,
387 2015 SIAM International Conference on Data Mining, Vancouver, Canada, April 30 - May 02,
388 2015.

389 Krumholz HM, Wang Y, Mattera JA, Han LF, Ingber MJ, Roman S, Normand SL. 2006. An
390 administrative claims model suitable for profiling hospital performance based on 30-day
391 mortality rates among patients with heart failure. *Circulation*. 113(13) (2006) 1693-1701.

392 Friedman, Jerome H. Greedy function approximation: A gradient boosting machine. *Ann.*
393 *Statist.* 29 (2001), no. 5, 1189-1232.

394 Hao S, Wang Y, Jin B, Shin AY, Zhu C, Huang M, Zheng L, Luo J, Hu Z, Fu C, Dai D. 2015.
395 Development, Validation and Deployment of a Real Time 30 Day Hospital Readmission Risk
396 Assessment Tool in the Maine Healthcare Information Exchange. *PLoS ONE*. 10(10) (2015)
397 e0140271.

398 HCUP State Inpatient Databases (SID). Healthcare Cost and Utilization Project (HCUP). 2009-
399 2011. Agency for Healthcare Research and Quality, Rockville, MD. [www.hcup-](http://www.hcup-us.ahrq.gov/sidoverview.jsp)
400 [us.ahrq.gov/sidoverview.jsp](http://www.hcup-us.ahrq.gov/sidoverview.jsp)

401 He D, Mathews SC, Kalloo AN, Hutfless S. 2014. Mining high-dimensional administrative
402 claims data to predict early hospital readmissions. *Journal of the American Medical Informatics*
403 *Association*. 21(2) (2014) 272-279.

404 Incavo SJ, Derasari AM. 2014. The Cost of Obesity, *The Journal of Bone & Joint Surgery*. 96(9)
405 (2014) e79.

406 Kadry B, Press CD, Alesh H, Opper IM, Orsini J, Popov IA, Brodsky JB, Macario A. 2014.
407 Obesity increases operating room times in patients undergoing primary hip arthroplasty: a
408 retrospective cohort analysis. *PeerJ*, 2, p.e530.

- 409 Kansagara D, Englander H, Salanitro A, Kagen D, Theobald C, Freeman M, Kripalani S. 2011.
410 Risk prediction models for hospital readmission: a systematic review. *Jama*. 306(15) (2011)
411 1688-98.
- 412 R Core Team. 2016. R: A language and environment for statistical computing. R Foundation for
413 Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- 414 Shahn Z, Ryan P, Madigan D. 2015. Predicting health outcomes from high-dimensional
415 longitudinal health histories using relational random forests, *Statistical Analysis and Data
416 Mining: The ASA Data Science Journal*. 8(2) (2015) 128-36.
- 417 Soguero-Ruiz C, Hindberg K, Mora-Jiménez I, Rojo-Álvarez JL, Skrøvseth SO, Godtliebsen F,
418 Mortensen K, Revhaug A, Lindsetmo R, Augestad KM, Jenssen R. 2016. Predicting colorectal
419 surgical complications using heterogeneous clinical data and kernel methods, *Journal of
420 Biomedical Informatics*. 61 (2016) 87-96.
- 421 Stiglic G, Brzan PP, Fijacko N, Fei W, Delibasic B, Kalousis A, Obradovic Z. 2015.
422 Comprehensible Predictive Modeling Using Regularized Logistic Regression and Comorbidity
423 Based Features. *PLoS ONE*. 10(12) (2015) e0144439.
- 424 Stiglic G, Davey A, Obradovic Z. 2013. Temporal Evaluation of Risk Factors for Acute
425 Myocardial Infarction Readmissions, Healthcare Informatics (ICHI), 2013 IEEE International
426 Conference. (2013) 557-562.
- 427 Tibshirani R. 2011. Regression shrinkage and selection via the lasso: a retrospective. *Journal of
428 the Royal Statistical Society: Series B (Statistical Methodology)*. 73(3) (2011) 273-82.
- 429 Walsh C, Hripesak G. 2014. The effects of data sources, cohort selection, and outcome definition
430 on a predictive model of risk of thirty-day hospital readmissions. *Journal of biomedical
431 informatics*. 52 (2014) 418-26.
- 432 Yang Y, Logan J. 2006. A data mining and survey study on diseases associated with
433 paraesophageal hernia. Proceedings American Medical Information Association 2006 Annu
434 Symp, American Medical Informatics Association. (2006) 829-833.
- 435 Zhou SM, Fernandez-Gutierrez F, Kennedy J, Cooksey R, Atkinson M, Denaxas S, Siebert S,
436 Dixon WG, O'Neill TW, Choy E, Sudlow C. 2016. "Defining disease phenotypes in primary care
437 electronic health records by a machine learning approach: a case study in identifying rheumatoid
438 arthritis." *PloS one*, e0154515.
- 439 Zhu K, Lou Z, Zhou J, Ballester N, Kong N, Parikh P. 2015. Predicting 30-day Hospital
440 Readmission with Publicly Available Administrative Database. *Methods of information in
441 medicine*. 54(6) (2015) 560-567.