

1
2
3
4
5
6
7
8
9

Computational Barthel Index: An Automated Tool for Assessing and Predicting Activities of Daily Living Among Nursing Home Patients

Janusz Wojtusiak^{1*}, Negin Asadzadehzanjani¹, Cari Levy², Farrokh Alemi¹, Allison E. Williams³

*Correspondence: jwojtusiak@gmu.edu

¹Department of Health Administration and Policy, George Mason University, Fairfax, VA, USA.

Full list of author information is available at the end of the article

Abstract

Background: Assessment of functional ability, including Activities of Daily Living (ADLs), is a manual process completed by skilled health professionals. In the presented research, an automated decision support tool, the Computational Barthel Index Tool (CBIT), was constructed that can automatically assess and predict probabilities of current and future ADLs based on patients' medical history.

Methods: The data used to construct the tool include the demographic information, inpatient and outpatient diagnosis codes, and reported disabilities of 181,213 residents of the Department of Veterans Affairs' (VA) Community Living Centers. Supervised machine learning methods were applied to construct the CBIT. Temporal information about times from the first and the most recent occurrence of diagnoses was encoded. Ten-fold cross-validation was used to tune hyperparameters, and independent test sets were used to evaluate models using AUC, accuracy, recall and precision. Random forest achieved the best model quality. Models were calibrated using isotonic regression.

Results: The unabridged version of CBIT uses 578 patient characteristics and achieved average AUC of 0.94 (0.93-0.95), accuracy of 0.90 (0.89-0.91), precision of 0.91 (0.89-0.92), and recall of 0.90 (0.84-0.95) when re-evaluating patients. CBIT is also capable of predicting ADLs up to one year ahead, with accuracy decreasing over time, giving average AUC of 0.77 (0.73-0.79), accuracy of 0.73 (0.69-0.80), precision of 0.74 (0.66-0.81), and recall of 0.69 (0.34-0.96). A simplified version of CBIT with 50 top patient characteristics reached performance that does not significantly differ from full CBIT.

Conclusion: Discharge planners, disability application reviewers and clinicians evaluating comparative effectiveness of treatments can use CBIT to assess and predict information on functional status of patients.

Keywords: Machine learning, Supervised learning, Gerontology, Activities of Daily Living

Background

Knowledge about functional abilities and their decline is important for decision making regarding care provided to patients. For example, in a study by Fried [1], it was observed that patients who were aware that they were unlikely to return to their baseline functional status were less likely to proceed with hospital treatment. It is shown that the quality of life is more important than living longer [2]. Quality of life depends on many factors, one of which is patients'

33 functional independence. Functional ability of nursing home patients is assessed by direct observation of a skilled nurse
34 practitioner, which is a time consuming and costly process. The assessments are often reported using the Minimum Data
35 Set (MDS), a standardized patient evaluation instrument collected by nurses through observing patients in consultation
36 with other care team members. In the United States, assessment data are collected by all Medicare and Medicaid-
37 certified nursing homes and entered in MDS Section G [3]. MDS data are typically collected every three months, or
38 whenever a patient status changes. In contrast, similar detailed functional assessments are not routinely collected for
39 most elderly patients outside of nursing homes. To remedy this situation, this paper examines whether functional ability
40 can be assessed and predicted through coded data available in Electronic Health Records (EHRs) or medical claims.
41 Specifically, the focus is on the ability to independently perform activities of daily living (ADLs). Nine out of ten functional
42 abilities in the Barthel Index (Score) were used [4,5] as described in the Data section. The ten items that represent the
43 ability and level of independence in performing activities of daily living include: feeding, bathing, grooming, dressing,
44 bowel incontinence, bladder incontinence, toilet use, transfers (bed to chair and back), mobility (walking), and stairs [6].

45 The ability to automatically derive and predict patients' functional status has several important uses in clinical work
46 and research. Firstly, it may provide a more efficient and cost-effective means of assessing functional status in groups
47 for whom functional status is currently manually assessed. In a recent review that examined functional status quality
48 indicators, the authors concluded that using chart reviews or patient-reports is costly and administratively burdensome
49 [7]. Secondly, it may allow for retrospective assessment of patients' functional status for whom evaluations have not
50 been completed. Thirdly, it can be beneficial for patients who are typically not evaluated for the purpose of comparing
51 care across settings. Finally, predicting functional status up to one year in the future provides a basis for an informed
52 discussion between clinicians and patients/caregivers and may help in planning care for patients.

53 Previously, a set of models capable of predicting trajectories of ADL improvement or decline post-hospitalization [8],
54 as well as sequences of functional decline were constructed [9]. The former focused on predicting if patients are likely
55 to follow one of seven pre-defined trajectories of improvement/decline. Predictions were anchored to the time of
56 hospital discharge and diagnoses were extracted only from inpatient records of the corresponding hospitalization. The
57 method and tool discussed in this paper, called the Computational Barthel Index Tool (CBIT), significantly extends the
58 previous work and is designed to allow for assessment of functional status at any arbitrary moment. The tool that allows
59 for prediction of each ADL up to one year ahead, is based on a larger cohort of patients, and uses both inpatient and
60 outpatient diagnoses. The name is inspired by the original Barthel Index (Score), which is a standardized tool used to
61 evaluate activities of daily living [10]. Computational machine learning methods are used to construct the index. The

62 presented research also extends previous work [8] by incorporating temporal information about when events happened
63 in the patient's medical history, which was not applicable to hospitalization-only data. Many diagnoses present in
64 medical records correlate with the patient's functional ability, with some of these correlations being temporary and
65 others being permanent. For example, some surgical patients have urinary incontinence for a short period after the
66 surgery, while amputation affects the ability to walk permanently. Thus, it is assumed that the codes present in data are
67 time-dependent. It was shown that adding temporal information can improve the accuracy of the constructed CBIT
68 models, as discussed later in the paper.

69 Prediction of functional status and disability is challenging. Researchers in many studies have attempted to
70 automatically assess and predict functional status, including ADLs. Overall, there are three main approaches to assess
71 and predict ADLs by (1) using specific clinical data, (2) using sensor data collected by wearable devices or smarthome
72 environments, and (3) using patient records extracted from EHR or claims data in making assessment and predictions.
73 Despite wide selection of published works, the research presented here is unique in the latter category as its attempts
74 to assess and predict ADLs purely based on diagnoses and demographics present in the patient records. It should be
75 noted that there are a number of published papers that discuss ADLs as predictors of other outcomes such as disease
76 progression and mortality [11,12], while the focus of this study is on predicting ADLs.

77 Many studies attempted to predict ADLs in a specific population, i.e., related to a disease or injury [13, 14, 15], while
78 others are more general. In one study, machine learning (ML) methods were linked to biomedical ontologies to predict
79 functional status [16], achieving predictive accuracy of 0.6. In another work, researchers described a logistic regression-
80 based method to predict mortality and disability post-injury for the elderly [17] with reported R^2 of 0.86. Tarekegn et al.
81 developed a set of models to predict disability as a metric for frailty conditions resulting in models with F-1 scores ranging
82 between 0.74 to 0.76 [18]. Similarly, Gobbens and van Assen examined six standard frailty indicators (gait speed, physical
83 activity, hand grip, body mass index, and fatigue and balance) for assessing ADLs, of which only gait speed was predictive
84 of ADL disabilities [19]; however, no actual predictive accuracy was reported. More recently, Jonkman et al., constructed
85 logistic regression- based models from four datasets to predict decline in five ADLs [20], with the average AUC of 0.72.
86 It is clear that the above studies reported model performances below ones reported here. However, it should be
87 mentioned that these works were performed in different settings thus no direct comparison is meaningful. A systematic
88 review of published works related to assessing ADLs identified several commonly used predictors, including age,
89 cognitive functioning, depression, and hospital length of stay [21]. In the data-driven approach presented here, some of
90 the predictors are the same as those previously reported in the literature.

91 Not surprisingly, several research groups focused on assessing ADLs from sensor data. Assessing ADLs selected by
92 wearable sensors is a reasonable approach as it allows for continuous monitoring rather than a snapshot of activities
93 evaluated by a healthcare provider [21, 22, 23, 24, 25, 26]. In some studies, ambient intelligence and smarthome sensors
94 were used to assess the ability to perform ADLs. These works rely on the use of specific sensors installed in smarthome
95 environment that monitor movement [27, 28], as well as use of specific home devices [29, 30, 31]. Further, beyond the
96 direct application to the elderly population, activity recognition is a well-established field with several review papers
97 available to summarize the works [32, 33, 34].

98 The presented CBIT can be linked to an EHR through a standardized interface and used by clinicians to assess
99 functional abilities at the time of a specific patient visit or in a batch/bulk mode to predict current functional abilities as
100 well as ADL changes for a group of patients. The models used in the tool rely on readily available data in EHR systems or
101 claims data and do not require additional data collection. In addition, a simplified version of the tool was developed
102 based on 50 patient characteristics selected from amongst 578 used in the complete model. The simplified version was
103 used to build an online calculator capable of asking limited number of questions about patients' medical history and
104 presenting the results in a graphical form such as exemplified in Figure 1. In the figure, each line corresponds to one ADL
105 plotted over time for a hypothetical patient. The horizontal axis indicates time and the vertical axis shows the probability
106 of functional independence. It should be mentioned that this probabilistic interpretation of the prediction is not
107 intended to indicate the level of disability, but rather the confidence the models have in predictions. In this example,
108 the hypothetical patient is predicted to have functional independence with high probability in terms of bathing, bladder,
109 dressing, toileting, transferring and walking. In terms of eating and grooming, this patient is predicted to temporarily
110 recover approximately 6 months after the initial assessment and decline afterwards (see Discussion section for more
111 details).

112 <INSERT FIGURE 1 HERE>

113 In the presented work, two cases are considered: when previous functional status of a patient is unknown and only
114 diagnoses and demographics can be used as predictors, and when a patient was previously evaluated and results of that
115 evaluation (nine previous ADL attributes) can be added to the list of predictors. Thus, two sets of models were
116 constructed: *Evaluation models*, M_E^d , in which previous functional status assessment is unknown, and *Re-Evaluation*
117 *models*, M_{RE}^d , in which previous functional status is known. Here d is an ADL (bathing, grooming, etc.), and $\tau \in$
118 $\{0,90,180,365\}$ is the prediction horizon (given as the number of days), i.e., how far ahead in time the value is predicted.
119 As names suggest, M_E^d models are used in situations in which a new patient is being evaluated in terms of ADLs, and

120 M_{RE}^d models are used when an evaluation of the previously assessed patient needs to be refreshed as new information
121 becomes available.

122 The presented research has been initiated as part of a larger IRB-approved project in the Department of Veterans
123 Affairs (VA) with the purpose of assessing the cost and effectiveness of the Medical Foster Home program compared to
124 traditional Community Living Centers (nursing homes) [8, 9, 35]. Determination of patients' functional status was used
125 as one of the characteristics to match residents in both settings for comparison purposes. In this context, the main
126 contributions of the presented work are in (1) the development of models for assessment and prediction of ADLs up to
127 one year ahead; (2) construction of attributes that represent time between diagnosis and prediction; (3) detailed testing
128 and analysis of the developed models, and (4) creation of an online decision support tool.

129 **Methods**

130 **Data**

131 Data from the Department of Veterans Affairs Corporate Data Warehouse were extracted and analyzed within the
132 VA Computing Infrastructure. The original data came from two sources: (1) medical records from the VA's Electronic
133 Medical Record System, and (2) MDS evaluations for nationwide VA nursing homes. Both datasets are collected as part
134 of routine patient care and were provided to the research team in a deidentified form. The data were organized around
135 patient evaluations using Minimum Data Set 2.0 [36], which were mapped to the nine Barthel Index categories using a
136 previously developed procedure [8]. The Barthel Index (or Barthel Score), which measures independence in performing
137 ADLs [4,5] includes 10 items with the total value ranging from 0 to 100 (feeding, bathing, grooming, dressing, bowel
138 incontinence, bladder incontinence, toilet use, transfers, mobility, and stairs). In this research, the last item of the
139 Barthel Score (stairs) was eliminated, which was not consistently assessed and thus difficult to standardize among
140 nursing home residents. Thus, the total considered scale is 0-90 based on the first nine items predicted independently.
141 Each of the items in the Barthel Score has different levels of functional abilities, with highest values indicating full
142 independence (see Additional file 1 for more details). For instance, Barthel Score captures three levels for toileting:
143 dependent (0), needs some help (5), and independent (10). Binary output for each of the ADLs was constructed defined
144 as fully functional vs. any level of dependency.

145 The data consisted of 1,901,354 MDS evaluations completed between 2000 and 2011 from which 1,151,222
146 complete evaluations were retrieved for 295,491 patients. The data were linked to medical records from which
147 demographics and history of diagnoses were extracted. The EHR data are limited to services provided by the VA's health
148 system. The data consisted of 18,912,553 inpatient and 180,123,710 outpatient diagnosis codes using the International

149 Classification of Diseases, ninth edition (ICD-9) standard along with corresponding dates. These codes were transformed
150 into clinically relevant categories using Clinical Classification Software (CCS) from the Agency of Health Research and
151 Quality (AHRQ) resulting in 281 distinct CCS codes representing health comorbidities. All diagnosis codes were combined
152 from inpatient and outpatient records. Distinguishing between inpatient and outpatient codes is important for some
153 applications (inpatient codes are typically treated as more severe). In the presented work, it is assumed that only
154 information about the presence of a diagnosis along with appropriate time was important in the context of predicting
155 disabilities, rather than distinguishing between the specific sources. Demographic information including age, race, and
156 gender was also included. Age was recorded as a continuous variable and race was represented using one-hot vectors
157 (0/1 values are used to indicate the presence or absence of the features). Missing data for age were imputed as mean
158 value in the dataset and no special treatment for missing data for other attributes was needed. Patients with only one
159 MDS evaluation were excluded to allow for modeling of change of patient status over time, resulting in a final dataset
160 of 855,731 evaluations for 181,213 patients. The collected data were organized per MDS evaluation, resulting in the
161 average of 4.72 +/- 6.21 MDS evaluations per patient. Table 1 shows descriptive statistics of the final dataset as counted
162 in analyzed MDS records as well as per patient, and is representative of the overall nursing home population in the VA.
163 Most patients were male and white with an average age of over 71 years and mean Barthel Score (sum of assigned
164 Barthel items) of about 48 out of 90, indicating overall high levels of disability in the studied population. In addition, the
165 average score at the first evaluation was about 52. The average time between MDS evaluations was also about 100 days,
166 which is slightly over three months.

167 **Table 1 Characteristics of data.**

	All Data		Patients with at least 2 MDS Evaluations	
	MDS Records	Patients	MDS Records	Patients
N	1151222	295491	855731	181213
Gender	Male	96.8%	96.9%	96.7%
	Female	3.32%	3.1%	3.3%
Race	Asian	1.5%	1.4%	1.6%
	Black	13%	11.9%	13.4%
	White	58.8%	55.4%	59.9%
	Other	26.7%	31.3%	25.1%
Age	71.89 +/- 12.38	--	72.26 +/- 12.31	--
Age at first MDS	--	70.8 +/-12.51	--	71.05 +/- 12.43
CCS ^{max}	1424.65+/-1215.7	--	1504.17+/- 1123.78	--
CCS ^{min}	619.75+/-867.49	--	663.84+/-889.14	--

Barthel Score	49.00 +/- 29.98	--	47.81 +/- 30.17	--
Score at first MDS	--	52.44 +/- 29.14	--	53.6 +/- 28.8
Time between	--	--	101.93 +/- 234.31	143.66 +/- 374.16

168 In addition, the distribution of values for the nine ADLs is presented in Table 2. With the exception of bladder
169 incontinence, bowel incontinence and eating, the majority of evaluations indicate some level of dependency in
170 performing ADLs. Lack of full independence in terms of walking is the most prominent, with 73% of evaluation records
171 and 80% of patients. While these values are not equal to 50%, the data are reasonably balanced thus no additional
172 resampling or balancing was required.

173 **Table 2: Distribution of the nine considered ADLs.**

N	MDS Records	Patients
	855,731	181,213
Any level of dependency		
Bathing	74.2%	77.5%
Bladder	39.7%	43.0%
Bowels	41.4%	45.3%
Dressing	66.2%	71.4%
Eating	47.8%	54.6%
Grooming	63.0%	67.1%
Toileting	60.9%	66.8%
Transferring	52.1%	60.9%
Walking	73.2%	80.1%

174 The numbers are proportion of data with values indicating any level of dependency.

175 In the used data warehouse, as well as in many administrative datasets, patient medical records often span many
176 years, making it possible to examine temporal relationships between diagnoses and the predicted events. In the
177 presented research, a simple approach to incorporate time was used. Values of attributes corresponding to diagnoses
178 represent time between first known occurrence of a diagnosis code and the time of MDS evaluation.

$$179 \quad ccs_i^{max} = \max_{t_i}(t_p - t_i) \quad (1)$$

180 Here, (t_i) is the time of i -th diagnosis code occurring in the data, and (t_p) is the time of prediction. Note that each diagnosis
181 code may be present in the data multiple times. Another set of attributes represent the last recorded occurrence of the
182 diagnosis code relative to the time of MDS evaluation.

$$183 \quad ccs_i^{min} = \min_{t_i}(t_p - t_i) \quad (2)$$

184 In the original data, diagnoses have associated dates thus days are used as unit of time. This allows counting the
185 difference in time as the number of days. In other words ccs_i^{max} is the number of days separating the first occurrence

186 of the diagnosis and the time of prediction, and ccs_i^{min} is the number of days separating the most recent occurrence of
 187 the diagnosis and the time of prediction.

188 This method of constructing attributes provides information about how long a patient suffers from a given condition
 189 as well as if the condition is still present at the time of assessment (when was the most recent diagnosis of a specific
 190 health condition). The rationale behind this approach is that for many chronic conditions that affect patients' ability to
 191 perform ADLs over time, it is important to know how long the condition is present for the patient. Similarly, for many
 192 acute conditions, their effects on ADLs are temporary, thus only recent occurrences are important to consider. It should
 193 be noted that the chronic/acute status of a condition is not assigned ahead of time and each diagnosis is encoded using
 194 both ccs_i^{max} and ccs_i^{min} . It was observed that the models tend to rank higher ccs_i^{max} codes for chronic conditions and
 195 ccs_i^{min} for acute conditions, yet full validation of this fact is out of scope of this paper.

196 An example of data encoded using the above method is presented in Table 3. The table shows data for two different
 197 fictitious patients. Patient 1 has two MDS evaluations in the data 90 days apart. Patient 2 also has two MDS evaluations
 198 100 days apart. Patient 1 was diagnosed with septicemia only once, 210 days prior to the first evaluation
 199 ($ccs_2^{min}=ccs_2^{max}=210$). The patient has not been diagnosed second time between the evaluations because both columns
 200 representing the first and most recent occurrence increased by the same amount. The patient was diagnosed with
 201 hypertension 18 days prior to the first evaluation ($ccs_{99}^{min}=18$), and for the first time 500 days prior to the first evaluation.
 202 The patient was diagnosed with hypertension again 5 days prior to the second evaluation. Similarly, Patient 2 has been
 203 diagnosed with septicemia twice, 15 and 700 days prior to the first evaluation ($ccs_2^{min}=15$ and $ccs_2^{max}=700$). Patient 2
 204 was also diagnosed with tuberculosis 71 days before the second evaluation ($ccs_1^{min}=ccs_1^{max}=71$). One can also notice
 205 that Patient 1's ADLs declined between the evaluations. Diagnoses not present/recorded in patient's records are coded
 206 as -999999 and 999999.

207

208 **Table 3 Four example records of the data for two patients.**

Demographics		ADLs			Diagnoses							
Pat	...	Age	Feed	Transferring	...	ccs_1^{min}	ccs_1^{max}	ccs_2^{min}	ccs_2^{max}	...	ccs_{99}^{min}	ccs_{99}^{max}
1	...	73	10	5		999999	-999999	210	210		18	500
1	...	73	5	0		999999	-999999	300	300		5	590
2	...	60	10	15		999999	-999999	15	700		999999	-999999
2	...	61	10	15		71	71	115	800		999999	-999999

209 Complete data has 578 columns and 888,731 rows.

210 Negative numbers (-999999) are used for coding of not present diagnoses in ccs_i^{max} columns because that time is
211 intended to capture positive correlation between long-term chronic conditions and disabilities. Intuitively, the longer a
212 patient suffers from a chronic condition (large values for time), the worse the prognosis is. When a condition is not
213 present in the patient’s medical history, it needs to be coded as “much better” than if the patient was just diagnosed;
214 thus, using a large negative number is reasonable. Similarly, positive numbers (999999) are used for coding of not
215 present diagnoses in columns, ccs_i^{min} , because of the negative correlation of time between the most recent occurrence
216 of conditions and disabilities. Full evaluation of this coding method in CBIT is discussed in the Results section.

217 **Construction of models**

218 The presented study followed a standard experimental design used in machine learning. Patients were randomly
219 assigned to training (90%) and testing (10%) sets. The testing set with a sufficiently large sample size (approximately
220 18,000 patients) was used only for final validation of the models. Training dataset was used for 10-fold cross-validated
221 hyperparameter tuning, model selection and final model construction. A selection of machine learning methods was
222 investigated to construct models capable of assessing and predicting ADLs.

223 Machine learning methods are rapidly gaining popularity in medical and health applications [37] and are also
224 applicable to the prediction of ADLs. Machine learning (ML) is an experimental field that provides a large toolset of
225 methods that can be used for prediction. More specifically, the presented work utilizes a set of ML methods called
226 supervised learning. These methods are intended to build models that allow for predicting outcomes for individuals
227 based on their characteristics. The supervision comes in the form of training data in which outcomes are known for
228 historical cases. These historical cases/patients are generalized to allow predictions for new previously unseen cases.

229 In the presented work, selected ML methods (regularized logistic regression, Bayesian networks, decision trees, and
230 random forests) were evaluated in terms of their performance and it was shown that random forest stands out in terms
231 of model quality. Random forests [38] are ensembles of decision trees (typically many), that are inferred from randomly
232 selected subsets of data thus guaranteed to be different on sufficiently large data. Random forests are created by
233 applying bagging (a.k.a., bootstrap aggregation) [39] to both sample and attributes (patient characteristics). Standard
234 top-down decision tree learning algorithms are used to create individual trees. The process is repeated to create multiple
235 trees (typically in the order of tens or hundreds). After a forest is assembled, the final classification decision is made by
236 applying all of the trees to new examples (patients). When there is a disagreement in prediction, the trees vote on the
237 predicted outcome. Random forests output classification scores (in the presented work, they were converted to
238 probabilities) which in the case of the described models represent patients being disabled or functionally independent.

239 These scores are calculated as a proportion of trees voting for a given outcome [40]. In the presented work, 10-fold
240 cross-validated hyperparameter tuning was performed. The tuning led to the selection of random forests consisting of
241 about 100 decision trees (each model was optimized separately, and the numbers of trees were slightly different). Other
242 algorithm parameters, including the number of randomly selected patient characteristics (number of attributes in each
243 tree) and Gini Index [38] as an internal quality criterion were tested and set to default as they did not make any
244 improvements.

245 The models were created to assess functional status at the time of prediction (current status), as well as to predict
246 functional status 3, 6, and 12 months beyond the time of prediction as depicted in Figure 2. Data available prior to the
247 time of prediction were used to construct input attributes for the model. In the constructed models, there are 9 ADLs
248 and 4 time points, thus there are 36 output attributes that are being predicted. Since Evaluation and Re-Evaluation
249 models are considered separately, CBIT consists of a total of 72 models.

250 <INSERT FIGURE 2 HERE>

251 The quality of the constructed models was evaluated in terms of standard statistical measures used in ML, namely,
252 accuracy (percentage of correctly predicted cases), area under the curve (AUC; often referred to as C-statistic), recall
253 (rate of correctly identified patients with functional dependencies), precision (rate of patients with disabilities among
254 those indicated as disabled by the model), and F1-score. Because of probabilistic interpretation of prediction results (see
255 discussion), we consider AUC as the most important metric. Evaluation was applied on the test set of patients not being
256 used in model construction, selection or tuning. In order to provide better insight into the created models, calibration
257 plots (described in later section) and learning curves were also created for all developed models. The learning curves
258 are used here to check if the amount of data used to train the models is sufficient. Curves that get flat on the right side
259 indicate that it is unlikely that more data would improve models, while those steeply growing suggest that models could
260 have been improved if more data were available. Learning curves for the CBIT models are available in supplemental
261 material (See Additional files 4 to 7).

262 It should be mentioned that the presented work does not include clinical validation of the models. Also, note that
263 the created models predict the probability of functional dependence of any level, while the graphical representation or
264 prediction (in the web calculator discussed later and presented figures) shows the probability of functional
265 independence. The conversion between the two is a simple operation, which is one minus probability. The reason for
266 this conversion is that prediction of disability as a target event is conceptually cleaner from a machine learning
267 perspective (assuming that being independent is normal, the abnormal state of disability is predicted). On the other

268 hand, clinicians are used to having higher values represent better status (this can also be the case in in the original
 269 Barthel Index). This conversion has no effect on presented results or modeling and is only reflected in the graphical
 270 representation of results.

271 For the data analysis part of the project, the Microsoft SQL Server was used to preprocess data. The data
 272 preprocessing started with MDS evaluations that were later linked to other data components. Final data were analyzed
 273 using Python programming language with Scikit-learn machine learning library [41] and visualizations were done using
 274 Matplotlib Python library [42].

275 Results

276 Computational Barthel Index Tool (CBIT) consists of a set of 72 random forest models, 36 $M_{E\tau}^d$ and 36 $M_{RE\tau}^d$ models.
 277 The CBIT can assess the level of functional dependency in performing ADLs and predicting functional dependency up to
 278 one year ahead by using demographics, diagnoses, and (if available) last known functional status. Table 4 presents a
 279 summary of the performance of the models for each ADL at the time of prediction, as well as 3, 6 and 12 months ahead
 280 for both $M_{E\tau}^d$ and $M_{RE\tau}^d$ models. The results are presented in terms of average AUC, accuracy, precision and recall of the
 281 nine outcome categories. The CBIT showed very high accuracy in assessing ADLs at a given time. The AUC of assessing if
 282 patients have any level of ADL dependency in $M_{RE^d_0}$ models was on average 0.94 (0.93-0.95), accuracy 0.90 (0.89-0.91),
 283 precision 0.91 (0.89-0.92), and recall 0.90 (0.84-0.95). When predicting functional status up to one year ahead, $\tau \in$
 284 $\{90,180,365\}$, the $M_{RE\tau}^d$ models' accuracy drops to AUC 0.77 (0.73-0.79), accuracy 0.73 (0.69-0.80), precision 0.74 (0.66-
 285 0.81), and recall 0.69 (0.34-0.96). When the previous functional status is unknown (i.e., initial evaluation), the
 286 performance of the current assessment models $M_{E^d_0}$ decreased by about 16% ($p < 0.01$) in terms of AUC. On average,
 287 the obtained results for these models are AUC 0.79, accuracy 0.74, precision 0.74, and recall 0.80. A complete set of
 288 results for individual models is available in Additional file 2.

289 **Table 4 Average +/- standard deviation of accuracy, AUC, precision and recall of models in predicting functional status**

Prediction	Re-Evaluation Models ($M_{RE\tau}^d$)				Evaluation Models ($M_{E\tau}^d$)			
	Accuracy	AUC	Precision	Recall	Accuracy	AUC	Precision	Recall
Time τ								
Current	.900 ± .007	.947 ± .006	.910 ± .011	.907 ± .041	.743 ± .029	.795 ± .010	.743 ± .046	.800 ± .128
3 Months	.815 ± .020	.876 ± .011	.849 ± .019	.816 ± .094	.727 ± .037	.761 ± .006	.734 ± .049	.783 ± .161
6 Months	.759 ± .029	.808 ± .014	.784 ± .029	.737 ± .165	.720 ± .038	.746 ± .009	.721 ± .045	.729 ± .238
12 Months	.737 ± .035	.772 ± .022	.742 ± .049	.699 ± .226	.716 ± .039	.725 ± .016	.696 ± .073	.701 ± .264

290 **Top predictors**

291 Further analysis also identified the top predictors used in the assessment and prediction of ADLs. Average Gini Index
 292 [38] produced by random forest was used to measure the quality of predictors. Gini index is a data impurity measure
 293 used in the presented work by random forest as an internal measure of attribute quality when constructing individual
 294 decision trees. It should not be interpreted as a strength or effect of the variable on the predicted output, but rather to
 295 understand the relative importance of attributes. In general, random forests can use many other attribute quality
 296 measures, but model tuning indicated that Gini index performs the best in CBIT. Top predictors along with their reported
 297 importance (average Gini index over all trees in forest and over all models) are presented in Table 5. Note that all ccs_i^{min}
 298 and ccs_i^{max} codes were included in full models. A longer list of diagnosis codes and previous evaluations are available in
 299 Additional file 3. Not surprisingly, the most predictive attributes in $M_{RE\tau}^d$ models were past functional status, being
 300 responsible for AUC of 0.93. Other most predictive attributes were the time since the most recent diagnosis of delirium,
 301 dementia, and amnestic and other cognitive disorders (CCS 653) and patient age. These were followed by encoded time
 302 of diagnoses/administrative codes for: the urinary tract infections (CCS 159); chronic ulcer of skin (CCS 199); other
 303 connective tissue disease (CCS 211); paralysis (CCS 82); administrative/social admission (CCS 255); alcohol-related
 304 disorders (CCS 660); aspiration pneumonitis; food/vomitus (CCS 129); and schizophrenia and other psychotic disorders
 305 (CCS 659). For most of the diagnoses listed above, it is important when (number of days) a patient was diagnosed with
 306 that condition most recently. For ulcers and aspiration pneumonitis; food/vomitus the first diagnosis is important. In
 307 addition, the table has marked potentially reversible conditions (R), as judged by clinicians, which can be influenced in
 308 the care provided to the patients and affect the outcome.

309 **Table 5 Top ranked predictors of functional status**

Rank	Attributes	Min/Max	Description	R	GINI RE-EVAL	GINI EVAL
1	ccs653	Min	Delirium, dementia, and amnestic and other cognitive disorders		0.0216	0.0310
2	Age		Age at the time of prediction		0.0133	0.0335
3	ccs159	Min	Urinary tract infections	X	0.0128	0.0217
4	ccs199	Max	Chronic ulcer of skin		0.0071	0.0121
5	ccs211	Min	Other connective tissue disease		0.0065	0.0091
6	ccs82	Min	Paralysis	X	0.0062	0.0110
7	ccs255	Min	Administrative/social admission	X	0.0061	0.0107
8	ccs660	Min	Alcohol-related disorders	X	0.0058	0.0110
9	ccs129	Max	Aspiration pneumonitis; food/vomitus		0.0055	0.0072

10	cs659	Min	Schizophrenia and other psychotic disorders	0.0055	0.0089
	...				
337	W		Race White	0.0006	0.0012
341	UR		Unknown Race	0.0006	0.0011
365	B		Race Black	0.0004	0.0009
434	Gender		Gender	0.0002	0.0004
445	A		Race Asian	0.0002	0.0003

310 "GINI RE-EVAL" indicates score of a variable in Re-Evaluation models ($M_{RE^d_\tau}$). "GINI EVAL" indicates score of a variable in in Evaluation
311 models ($M_{E^d_\tau}$). R are potentially reversible or red flag that this person is at risk and needs restorative therapy; Race and Gender variables
312 are included at the bottom of the table for comparison but have very low impact on prediction.

313 Simplified models

314 Further, simplified models (called $MS_{RE^d_\tau}$ and $MS_{E^d_\tau}$) that include only selected top-ranking patient characteristics
315 were developed. Average GINI score was used to rank attributes. As depicted in Figure 3, adding more characteristics
316 beyond the most predictive 41 attributes did not significantly improve the accuracy ($p < 0.05$) of the models in assessing
317 the current functional status ($\tau=0$) as compared to full model. The curves were also similar for predicting up to 12 months
318 ahead, $\tau \in \{90, 180, 365\}$. When using 25 top patient characteristics, models that included previous evaluations ($MS_{RE^d_\tau}$)
319 reached an average AUC of 0.94, accuracy 0.90, precision 0.91, and recall 0.90. Furthermore, the performance of the
320 simplified models with 41 patient characteristics and without previous evaluations ($MS_{E^d_\tau}$) raised to average AUC of 0.79,
321 accuracy 0.74, precision 0.74, and recall 0.78. Note that top predictors for each ADL are different. In the $MS_{RE^d_\tau}$ and $MS_{E^d_\tau}$
322 models, top ranking attributes were included across all models to minimize information needed by CBIT for all ADLs,
323 even though this set of attributes may not be optimal for individual models.

324 <INSERT FIGURE 3 HERE>

325 Temporal coding

326 One important advancement of the presented CBIT is the way it captures time in encoding diagnoses as previously
327 shown in equations (1) and (2) and illustrated in Table 3. The proposed method of constructing attributes for diagnoses
328 was investigated to determine how it would be different from binary attributes (1 when a diagnosis is present in a given
329 patient's record and 0 otherwise) when used in CBIT. All constructed $M_{RE^d_\tau}$, $M_{E^d_\tau}$, $MS_{RE^d_\tau}$ and $MS_{E^d_\tau}$ models were
330 compared in terms of AUC at different time points up to one year ahead. In one experiment, random forest was
331 compared with other algorithms including logistic regression, decision tree, and naïve Bayes.

332 As mentioned earlier, when temporal attributes are used, one needs to assign special values to diagnoses that are
333 not present in data. Therefore, +/- 999999 (6_9) was compared with +/-9999 (4_9), and +/-99999 (5_9) coding across all

334 models (here X_9 indicates 10^X-1). Temporal coding (6_9) was also compared with binary coding to determine any
 335 significant difference. Two-tailed t-test was used to assess all comparisons ($p<0.05$).

336 As summarized in Table 6, both random forest and logistic regression show a significant difference in AUC when
 337 temporal information is applied ($p<0.05$). The results indicated that random forest with the temporal coding performs
 338 significantly better than binary coding, while for the logistic regression the relationship is opposite (the binary coding is
 339 better). However, logistic regression with binary coding is still doing worse than random forest. Decision trees and naïve
 340 Bayes results were also included in the table, but the performance was typically inferior. It was observed that random
 341 forest, decision tree and naïve Bayes are not affected by how the special values were assigned, while the performance
 342 of logistic regression is affected by the coding. The rationale for this result is that for symbolic methods it is irrelevant
 343 how not-present values are coded as long as the value is distinct, while parametric models need to find a coefficient for
 344 each diagnosis code, which is affected by the coding.

345 **Table 6 Comparison of temporal and binary diagnosis coding as part of CBIT construction and evaluation.**

AUC		Current Assessment				3 Month Prediction				6 Month Prediction				12 Month Prediction				
		RF	LR	DT	NB	RF	LR	DT	NB	RF	LR	DT	NB	RF	LR	DT	NB	
$M_{RE}^{d_t}$	Temporary	4_9	0.95 [*]	0.85 ⁺⁺	0.92 ⁺	0.87 ⁺⁺	0.88	0.79 ⁺⁺	0.83 ⁺	0.83 ⁺	0.81	0.77 ⁺⁺	0.74 ⁺	0.78 ⁺	0.77	0.74 ⁺⁺	0.70 ⁺	0.74 ⁺
		5_9	0.95	0.78 ⁺⁺	0.92 ⁺	0.89 ⁺	0.88	0.76 ⁺⁺	0.83 ⁺	0.83 ⁺	0.81	0.74 ⁺	0.74 ⁺⁺	0.78 ⁺	0.77	0.71 ⁺⁺	0.70 ⁺	0.74 ⁺
		6_9	0.95	0.78 ⁺	0.92 ⁺	0.90 ⁺	0.88	0.75 ⁺	0.83 ⁺	0.83 ⁺	0.81	0.74 ⁺	0.74 ⁺	0.78 ⁺	0.77	0.72 ⁺	0.70 ⁺	0.74 ⁺
	Binary	0.94 ⁺	0.94 ⁺	0.91 ⁺⁺	0.87 ⁺⁺	0.87 ⁺	0.87 ⁺⁺	0.82 ⁺⁺	0.80 ⁺	0.81	0.81 ⁺⁺	0.74 ⁺	0.77 ⁺⁺	0.77	0.77 ⁺⁺	0.70 ⁺	0.74 ⁺⁺	
$M_{RE}^{d_t}$	Temporary	4_9	0.95	0.94 ⁺⁺	0.92 ⁺	0.89 ⁺	0.88	0.88 ⁺⁺	0.83 ⁺	0.82 ⁺	0.81 ⁺	0.81 ⁺⁺	0.74 ⁺	0.76 ⁺	0.77	0.77 ⁺	0.70 ⁺	0.72 ⁺
		5_9	0.95	0.93 ⁺⁺	0.92 ⁺	0.89 ⁺	0.88	0.84 ⁺⁺	0.82 ⁺	0.82 ⁺	0.81 ⁺	0.79 ⁺⁺	0.74 ⁺	0.76 ⁺	0.77	0.75 ⁺⁺	0.70 ⁺	0.72 ⁺
		6_9	0.95	0.76 ⁺	0.92 ⁺	0.90 ⁺	0.88	0.72 ⁺	0.83 ⁺	0.82 ⁺	0.81	0.71 ⁺	0.74 ⁺	0.76 ⁺	0.77	0.69 ⁺	0.70 ⁺	0.72 ⁺
	Binary	0.94 ⁺	0.94 ⁺⁺	0.90 ⁺⁺	0.90 ⁺	0.88 ⁺	0.87 ⁺⁺	0.81 ⁺⁺	0.83 ⁺	0.81 ⁺	0.81 ⁺⁺	0.74 ⁺⁺	0.78 ⁺⁺	0.77	0.77 ⁺	0.69 ⁺	0.74 ⁺⁺	
$M_{RE}^{d_t}$	Temporary	4_9	0.79	0.79 ⁺⁺	0.72 ⁺⁺	0.73 ⁺	0.76	0.76 ⁺	0.68 ⁺	0.68 ⁺	0.75 ⁺	0.75 ⁺	0.66 ⁺	0.71 ⁺	0.73	0.72 ⁺	0.64 ⁺	0.69 ⁺
		5_9	0.79	0.78 ⁺⁺	0.71 ⁺	0.73 ⁺	0.76	0.75 ⁺	0.68 ⁺	0.68 ⁺	0.75	0.74 ⁺⁺	0.66 ⁺	0.71 ⁺	0.73	0.71 ⁺⁺	0.64 ⁺	0.69 ⁺
		6_9	0.79	0.78 ⁺	0.72 ⁺	0.73 ⁺	0.76	0.75 ⁺	0.68 ⁺	0.68 ⁺	0.75	0.74	0.66 ⁺	0.71 ⁺	0.73	0.72 ⁺	0.64 ⁺	0.69 ⁺
	Binary	0.78 ⁺	0.78 ⁺	0.70 ⁺⁺	0.73 ⁺	0.76	0.76 ⁺	0.67 ⁺⁺	0.70 ⁺⁺	0.75	0.75 ⁺	0.66 ⁺	0.71 ⁺⁺	0.72 ⁺	0.73 ⁺⁺	0.64 ⁺	0.69 ⁺⁺	
$M_{RE}^{d_t}$	Temporary	4_9	0.79	0.77 ⁺⁺	0.71 ⁺	0.64 ⁺	0.76	0.75 ⁺	0.68 ⁺	0.63 ⁺	0.74	0.73 ⁺⁺	0.66 ⁺	0.60 ⁺	0.72	0.72 ⁺	0.63 ⁺	0.58 ⁺
		5_9	0.79	0.76 ⁺⁺	0.71 ⁺	0.64 ⁺	0.76	0.73 ⁺⁺	0.68 ⁺	0.63 ⁺	0.74	0.72 ⁺⁺	0.66 ⁺	0.60 ⁺	0.72	0.71 ⁺⁺	0.63 ⁺	0.58 ⁺
		6_9	0.79	0.75 ⁺	0.71 ⁺	0.64 ⁺	0.76	0.72 ⁺	0.68 ⁺	0.63 ⁺	0.74	0.71 ⁺	0.66 ⁺	0.60 ⁺	0.72	0.69 ⁺	0.63 ⁺	0.58 ⁺
	Binary	0.76 ⁺	0.77 ⁺⁺	0.68 ⁺⁺	0.74 ⁺⁺	0.74 ⁺	0.74 ⁺	0.65 ⁺⁺	0.71 ⁺⁺	0.73 ⁺	0.73 ⁺	0.64 ⁺⁺	0.71 ⁺⁺	0.71 ⁺	0.72 ⁺⁺	0.63 ⁺	0.69 ⁺⁺	

346 The results are presented in terms of AUC for the current assessment and prediction up to 12 months ahead. Full models that include 578
347 attributes and simplified models with 50 attributes are shown. 4_9, 5_9, and 6_9 indicate the encoding of diagnoses not present in patient's
348 history for +/-9999, +/-99999, and +/-999999, respectively. * indicates significance ($p < 0.05$) of coding systems compared to "6_9" and + indicates
349 significance ($p < 0.05$) of different algorithms compared to random forest.

350 Calibration

351 Calibration allows for the probability interpretation of the output scores from the models, further allowing for
352 frequency interpretation of the results. Thus, all models were calibrated using 5-fold cross-validated isotonic regression.
353 This approach fits a secondary model on top of the created random forest models and attempts to adjust returned
354 scores to make them closer to probabilities. The results showed that the models were well-calibrated with mean squared
355 error of about 3%. Figure 4 shows an example of the calibration curve for the model that assesses bathing at the current
356 time point, $MS_{RE}^{bathing_0}$. Similar curves were also developed for all 144 models and are available in supplemental materials
357 (See Additional files 8 to 11).

358 <INSERT FIGURE 4 HERE>

359 Discussion

360 Methods

361 It was shown that it is possible to assess and predict functional status using machine learning methods. Moreover, it
362 was shown that the inclusion of time between diagnosis and time of prediction is important in constructing attributes
363 in the data. While further work is needed to validate the new way of constructing attributes representing diagnoses and
364 study its limitations, counting days from the first and last known occurrence of a diagnosis code works for the problem
365 at hand.

366 Machine learning methods are gaining popularity in medical and health applications, yet there is no consensus on
367 what validation is needed for their use in clinical settings. There is also no agreement about what information is needed
368 to allow for full reproducibility of ML results, or even what reproducibility in this context means [43]. Models created
369 for CBIT were evaluated using standard measures in ML model testing (cross-validation, independent test set, etc.), and
370 investigated in terms of their calibration and learning curves. There is a need for further validation of the models and
371 their impact on patient care. Such validation focuses on detailed model analysis in terms of accuracy, transparency and
372 the ability to provide explanations, and eventually trust and acceptability by the medical community. A randomized trial
373 to assess outcomes of the model use may be required for full acceptance in clinical settings.

374 In addition, there is an ongoing discussion about the overall validity of applying machine learning methods to the
375 prediction of patient outcomes, and potential bias of the models based on gender, race and socioeconomic status. One
376 needs to clearly understand data limitations and definitions of the prediction problem to understand the drawbacks of

377 the method. Supervised machine learning methods, by definition, learn what they are asked to learn, and may (typically
378 do) propagate biases from training data. Biases in machine learning-based models are typically not caused by machine
379 learning, but by underlying process used to create training data. The key is in the definition and construction of the
380 output attributes of the model and their proper interpretation. One needs to answer a question if models predict events
381 in the real world, or data artifacts that somehow approximate that reality. Similarly, CBIT is intended to mimic the tasks
382 of nurses performing evaluations of ADLs as part of the MDS. Therefore, any biases, inaccuracies, or subjectivity in this
383 process may also be repeated by the CBIT models. However in CBIT, as shown in Table 5, race and gender had only a
384 negligible impact on predictions and were completely dropped from simplified models, which suggests diminishing the
385 potential racial and gender bias in the models. A different set of methods, typically used in health services research, is
386 needed to understand the existence of potential bias in constructed models.

387 The probabilistic interpretation of the prediction results used in the presented work seems to be reasonable.
388 Conceptually, the future can never be predicted with the probability of one (even though for some cases the models
389 may be certain of the future and output 1.0). Instead, the values represent how likely an event (here functional
390 independence) will occur according to the models. Such interpretation has several advantages. It allows end users to
391 interpret the chances of an event happening, and in turn, describes when models are uncertain about the predicted
392 outcomes. In applications such as the presented CBIT, providing probability as a form of explanation makes the
393 predictions more transparent. Knowing how likely an outcome is going to happen can help clinicians, patients and their
394 families make informed decisions related to planning care. This is in contrast with systems in which ML-based models
395 trigger certain event such as alerts within EHR systems. Such triggered events are binary in nature (alert or no alert),
396 thus the final assigned class is most important. The probability results also explain why model accuracy is not 100%
397 when executed on the test data (examples with predicted probability not equal to one, are ambiguous by the definition
398 of probability). The latter is the most evident when analyzing calibration curves, such as one presented in Figure 4. The
399 disadvantage of using the probabilities is that they may be misinterpreted as severity of disability. When presenting
400 results, one needs to specify that the number represent how likely a patient is independent, and not the level of
401 dependency.

402 The created models in this manuscript are based on ICD-9 diagnosis codes, which were mapped to CCS codes. One
403 advantage of this approach is that since all new data are coded with ICD-10 codes, they could easily be mapped to CCS
404 codes making the models applicable to data with the newer coding system. Another important issue is that the diagnosis
405 codes in both EHR and claims data are subject to under- and over-coding, thus affecting the potential reliability of the

406 models. However, it is important to note that our modeling efforts were not intended to understand the effects of
407 diagnoses on ADLs, but rather their use in making prediction. In addition, as long as diagnoses are systematically
408 over/under-coded, they should not affect performance of the models. Despite these limitations, results indicated that
409 our data were appropriate for this purpose.

410 The evaluations presented in this paper are only summaries and examples of detailed results. A detailed examination
411 of all 72 M^d_τ models that are part of CBIT and 72 MS^d_τ models that are part of the limited CBIT was performed, and the
412 results are available in the supplemental materials and through the online calculator [44]. All developed models and
413 source codes are available for everyone who wishes to conduct their testing on independent data from other
414 institutions, i.e. to test cross-institution generalizability.

415 **Clinical and administrative use**

416 Very little evidence exists to address whether measuring functional status can change the quality of life, but our
417 research shows that prior knowledge about functional disability is a key indicator of future functional status. Notably,
418 past research has provided evidence that improvements in functional status are possible over time through therapy [45]
419 by improving, slowing decline, and/or maintaining functional status. The presented CBIT tool which predicts
420 improvement or decline could be used by health professionals as means of identifying patient characteristics that are
421 modifiable and plan care accordingly. It can serve as a basis for an informed discussion between clinicians, patients and
422 caregivers. In addition, these measures could potentially serve as a patient-centered measure for examining the value
423 of the services provided.

424 **Graphical presentation of results and web calculator**

425 A graphical representation of the assessment and prediction of functional status can be used by healthcare
426 professionals and caregivers for decision making regarding the patients' care. Our full models can be integrated as
427 decision support tools within EHR systems or linked to claims data, while the simplified models can operate standalone
428 as an interactive online tool. For example, Figure 5 illustrates CBIT-predicted outcomes for three fictitious patients
429 similar to what was shown in Figure 1. Values indicate the probability of functional independence for each ADL up to
430 one year after the prediction time. The higher the value is, the higher the chance that the patient is functionally
431 independent. One can observe significant differences between the functional dependency trajectories for these
432 patients. Patient (a) is currently likely to be independent but expected to decline within 6 months as the probability of
433 independence decreases. Patient (b) is currently likely to be dependent in most ADLs (probability of independence
434 ranging from 0.2 to 0.6) but predicted to recover in the next 3 months and stay at this level afterward. Patient (c) is

435 independent and predicted to remain independent in terms of walking and is almost certainly disabled in terms of
436 bladder, bowels and eating. The patient is likely to have a temporary decline in terms of other ADLs. Construction of
437 each of the plots requires execution of 36 random forest models (9 ADLs, 4 time points).

438 <INSERT FIGURE 5 HERE>

439 An experimental version of the online calculator that takes patient characteristics and outputs plots is available at
440 <https://hi.gmu.edu/cbit> [44]. It is accessible through a web form or an application programming interface (API). The web
441 calculator is implemented in Python 3 and uses Flask as a web application framework, with Pandas and Scikit-learn
442 libraries performing data analysis. To ensure the performance of the web calculator, all of the models are loaded on the
443 startup and reside in RAM. Additional changes have been made to the calculator to improve clinical use. For example,
444 for the ease of use, the numbers of days associated with diagnoses were discretized to allow users to select them from
445 drop-down menus. Numbers closer to zero are discretized with higher precision than larger numbers, which further
446 improves understandability. Users can enter patient information and are provided with results similar to those shown
447 in Figure 5 along with a data table containing the values of predicted probabilities. An explanation module that provides
448 human-oriented interpretation of the results as well as the reasons for predictions is in development.

449 **Conclusion**

450 This study found that functional status can be assessed and predicted with high accuracy when prior functional status
451 in medical history is available, but also without requiring previous in-person functional assessment. It exemplifies an
452 opportunity of applying machine learning to large data to produce meaningful results. It was hypothesized that a
453 parsimonious model could be developed with variables available in EHRs or claims data and assumed that this model
454 would retain predictive accuracy for up to a year ahead. Our experimental results confirmed this hypothesis. The
455 constructed tool is intended to be used in both clinical and administrative settings and has implications for caregivers,
456 clinicians, and policy makers. Assessment and prediction of functional status may also lead to better care planning for
457 nursing home residents as well as the elderly residing in their own homes. Automated large-scale assessment and
458 prediction of functional status can be used to compare care settings and as a benchmark for provider outcomes.

459 The constructed full model requires a large number of predictors, which makes it impossible to manually enter
460 values. Hence, the full version of CBIT would need to be integrated with an EHR or claims management system to be
461 part of the clinical decision support. Such integration can be achieved using HL7s FHIR interface. The simplified version
462 of the CBIT that uses 50 predictors is available within a web calculator. Beyond the use of EHR data, the constructed

463 CBIT could be enhanced by sensor data allowing for continuous patient monitoring and be integrated with the presented
464 approach. Such data can aid assessment, particularly for ADLs that measure patient movement [46, 47, 48].

465 The presented work has a number of limitations. The tool is not applicable in settings in which longitudinal patient
466 records are not available. Only large health systems with long-established electronic medical records have sufficient
467 longitudinal data to apply models that use temporal diagnosis information. Additionally, the models were developed
468 using data from the US Department of Veterans Affairs (VA), which does not reflect the general population of nursing
469 home residents outside of the VA system. The performance of these models on other datasets, including Medicare
470 claims data are being investigated. It is unclear how the models will perform on a very different population and if the
471 existing CBIT models can be adapted. Finally, random forests are known to be “black box” models that work well but
472 are not well understood by end users. Even though their explanation is easier than other types of models such as neural
473 networks, they are significantly more difficult than linear models, decision trees or decision rules. Instead of trying to
474 explain the entire model (as part of the online calculator), there is an ongoing effort in designing an explanation module
475 that provides users with “reasons” for making specific predictions in one individual case (prediction explanation). The
476 reasons consist of a list of patient characteristics that are the strongest predictors (both confirming and disconfirming)
477 for that individual case. Despite these limitations, CBIT can be used to support clinicians and administrators in decision
478 making. Our novel data coding method, applying machine learning to unique health data, comprehensive model testing,
479 and transparency of the work contribute to the state-of-the-art in ML-based decision support.

480 **Abbreviations**

481 ADL: Activity of Daily Living; CBIT: Computational Barthel Index Tool; VA: Department of Veterans Affairs; MDS: Minimum Data Set; EHR: Electronic Health Record; ML:
482 Machine learning; ICD-9: International Classification of the Diseases, ninth edition; CCS: Clinical Classification Software; AHRQ: Agency of Health Research and Quality
483 Control; AUC: Area under the curve; RF: Random forest; LG: Logistic Regression; DT: Decision tree; NB: Naïve Bayes; ICD-10: International Classification of the Diseases,
484 tenth edition; API: Application programming interface; RAM: Random access memory.

485 **Declarations**

486 **Ethics approval and consent to participate**

487 The project has been approved by the Bay Pines VA IRB #2897. Secondary analysis of deidentified data was completed, thus no consent was required.

488 **Consent for publication**

489 Not applicable.

490 **Availability of data and material**

491 The data used to construct presented CBIT models are individual level and cannot be shared. Access to the original data may be requested through the US Department of
492 Veterans Affairs. All constructed models, source code, and detailed testing results are freely available at the project website.

493 **Competing interests**

494 The authors declare that they have no competing interests.

495 **Funding**

496 The project was funded in part by appropriation #3620160 from the VA Office of Geriatrics and Extended Care. The contents of this article do not represent the views of
497 the Department of Veterans Affairs or the United States Government.

498 **Author's contributions**

499 Dr. J.W is the main author of the manuscript. He also designed and partially implemented the method, supervised experimental evaluation, and implemented online
500 calculator. Ms. N.A partially implemented the method, constructed models, performed experiments, and wrote large section of experimental results. Dr. F.A contributed
501 to writing of the paper and helped in designing the overall experiments. Dr. C.L mapped clinical functional status data and contributed to writing clinical sections of the
502 paper. Dr. A.E.W provided clinical interpretation of the results and conclusions and supervised overall project. All authors have read and approved the manuscript.

503 **Acknowledgments**

504 Not applicable.

505 **Authors' information**

506 ¹Health Informatics Program, Department of Health Administration and Policy, George Mason University, Fairfax, VA, USA. ²Department of Veterans Affairs, Denver, CO,
507 USA. ³Department of Veterans Affairs, Bay Pines, FL.

508 **References**

- 509 1. Fried TR, Bradley EH, Towle VR, Allore H. Understanding the treatment preferences of seriously ill patients. *New*
510 *England Journal of Medicine*. 2002 Apr 4;346(14):1061-6.
- 511 2. McCarthy EP, Phillips RS, Zhong Z, Drews RE, Lynn J. Dying with cancer: patients' function, symptoms, and care
512 preferences as death approaches. *Journal of the American Geriatrics Society*. 2000 May;48(S1):S110-21.
- 513 3. MDS 3.0 Technical Information [Internet]. Available from: [https://www.cms.gov/Medicare/Quality-Initiatives-](https://www.cms.gov/Medicare/Quality-Initiatives-Patient-Assessment-Instruments/NursingHomeQualityInits/NHQIMDS30TechnicalInformation)
514 [Patient-Assessment-Instruments/NursingHomeQualityInits/NHQIMDS30TechnicalInformation](https://www.cms.gov/Medicare/Quality-Initiatives-Patient-Assessment-Instruments/NursingHomeQualityInits/NHQIMDS30TechnicalInformation).
- 515 4. Collin C, Wade DT, Davies S, Horne V. The Barthel ADL Index: a reliability study. *International disability studies*.
516 1988 Jan 1;10(2):61-3.
- 517 5. Shah S, Vanclay F, Cooper B. Improving the sensitivity of the Barthel Index for stroke rehabilitation. *Journal of*
518 *clinical epidemiology*. 1989 Jan 1;42(8):703-9.
- 519 6. THE BARTHEL INDEX [Internet]. Strokecenter.org. [cited 2020 Nov 6]. Available from:
520 <http://www.strokecenter.org/wp-content/uploads/2011/08/barthel.pdf>
- 521 7. Dy SM, Pfoh ER, Salive ME, Boyd CM. Health-related quality of life and functional status quality indicators for
522 older persons with multiple chronic conditions. *Journal of the American Geriatrics Society*. 2013
523 Dec;61(12):2120-7.
- 524 8. Wojtusiak J, Levy CR, Williams AE, Alemi F. Predicting functional decline and recovery for residents in veterans
525 affairs nursing homes. *The Gerontologist*. 2016 Feb 1;56(1):42-51.
- 526 9. Levy CR, Zargoush M, Williams AE, Williams AR, Giang P, Wojtusiak J, Kheirbek RE, Alemi F. Sequence of
527 functional loss and recovery in nursing homes. *The Gerontologist*. 2016 Feb 1;56(1):52-61.
- 528 10. Mahoney FI, Barthel DW. Functional evaluation: the Barthel Index: a simple index of independence useful in
529 scoring improvement in the rehabilitation of the chronically ill. *Maryland state medical journal*. 1965.
- 530 11. Hong HG, An HS, Sarzynski E, Oberst K. New Composite Measure for ADL Limitations: Application to Predicting
531 Nursing Home Placement for Michigan MI Choice Clients. *Medical Care Research and Review*. 2019 Nov
532 8:1077558719886735.
- 533 12. Li QX, Zhao XJ, Wang Y, Wang DL, Zhang J, Liu TJ, Peng YB, Fan HY, Zheng FX. Value of the Barthel scale in
534 prognostic prediction for patients with cerebral infarction. *BMC Cardiovascular Disorders*. 2020 Dec;20(1):1-5.
- 535 13. Veerbeek JM, Kwakkel G, van Wegen EE, Ket JC, Heymans MW. Early prediction of outcome of activities of
536 daily living after stroke: a systematic review. *Stroke*. 2011 May;42(5):1482-8.
- 537 14. Frank B, Schlote A, Hasenbein U, Wallesch CW. Prognosis and prognostic factors in ADL dependent stroke
538 patients during their first in-patient rehabilitation—a prospective multicentre study. *Disability and*
539 *rehabilitation*. 2006 Jan 1;28(21):1311-8.
- 540 15. Tanaka R, Umehara T, Fujimura T, Ozawa J. Clinical prediction rule for declines in activities of daily living at 6
541 months after surgery for hip fracture repair. *Archives of physical medicine and rehabilitation*. 2016 Dec
542 1;97(12):2076-84.
- 543 16. Min H, Mobahi H, Irvin K, Avramovic S, Wojtusiak J. Predicting activities of daily living for cancer patients using
544 an ontology-guided machine learning methodology. *Journal of biomedical semantics*. 2017 Dec 1;8(1):39.
- 545 17. Jeffery AD, Dietrich MS, Maxwell CA. Predicting 1-year disability and mortality of injured older adults. *Archives*
546 *of gerontology and geriatrics*. 2018 Mar 1;75:191-6.

- 547 18. Tarekegn A, Ricceri F, Costa G, Ferracin E, Giacobini M. Predictive Modeling for Frailty Conditions in Elderly
548 People: Machine Learning Approaches. *JMIR Medical Informatics*. 2020;8(6):e16678.
- 549 19. Gobbens RJ, van Assen MA. The prediction of ADL and IADL disability using six physical indicators of frailty: a
550 longitudinal study in the Netherlands. *Current gerontology and geriatrics research*. 2014;2014.
- 551 20. Jonkman NH, Colpo M, Klenk J, Todd C, Hoekstra T, Del Panta V, Rapp K, Van Schoor NM, Bandinelli S,
552 Heymans MW, Mauger D. Development of a clinical prediction model for the onset of functional decline in
553 people aged 65–75 years: pooled analysis of four European cohort studies. *BMC geriatrics*. 2019 Dec
554 1;19(1):179.
- 555 21. Hoogerduijn JG, Schuurmans MJ, Duijnste MS, De Rooij SE, Grypdonck MF. A systematic review of predictors
556 and screening instruments to identify older hospitalized patients at risk for functional decline. *Journal of clinical
557 nursing*. 2007 Jan;16(1):46-57.
- 558 22. Hong YJ, Kim IJ, Ahn SC, Kim HG. Activity recognition using wearable sensors for elder care. In 2008 Second
559 International Conference on Future Generation Communication and Networking 2008 Dec 13 (Vol. 2, pp. 302-
560 305). IEEE.
- 561 23. Liu J, Sohn J, Kim S. Classification of daily activities for the elderly using wearable sensors. *Journal of
562 healthcare engineering*. 2017 Nov 26;2017.
- 563 24. Cook DJ, Schmitter-Edgecombe M, Jönsson L, Morant AV. Technology-enabled assessment of functional
564 health. *IEEE reviews in biomedical engineering*. 2018 Jun 28;12:319-32.
- 565 25. Chatterjee P, Armentano R, Palombi L, Kun L. Editorial Preface: Special issue on IoT for eHealth, elderly and
566 aging.
- 567 26. Akbari A, Jafari R. Personalizing Activity Recognition Models with Quantifying Different Types of Uncertainty
568 Using Wearable Sensors. *IEEE Transactions on Biomedical Engineering*. 2020 Jan 3.
- 569 27. Sridharan M, Bigham J, Campbell PM, Phillips C, Bodanese E. Inferring Micro-Activities Using Wearable Sensing
570 for ADL Recognition of Home-Care Patients. *IEEE journal of biomedical and health informatics*. 2019 May
571 24;24(3):747-59.
- 572 28. Robben S, Englebienne G, Kröse B. Delta features from ambient sensor data are good predictors of change in
573 functional health. *IEEE journal of biomedical and health informatics*. 2016 Jul 22;21(4):986-93.
- 574 29. Ghayvat H, Mukhopadhyay S, Shenjie B, Chouhan A, Chen W. Smart home based ambient assisted living:
575 Recognition of anomaly in the activity of daily living for an elderly living alone. In 2018 IEEE International
576 Instrumentation and Measurement Technology Conference (I2MTC) 2018 May 14 (pp. 1-5). IEEE.
- 577 30. Sasaki W, Fujiwara M, Fujimoto M, Suwa H, Arakawa Y, Yasumoto K. Predicting Occurrence Time of Daily
578 Living Activities Through Time Series Analysis of Smart Home Data. In 2019 IEEE International Conference on
579 Pervasive Computing and Communications Workshops (PerCom Workshops) 2019 Mar 11 (pp. 233-238). IEEE.
- 580 31. Sokullu R, Akkaş MA, Demir E. IoT Supported Smart Home for the Elderly. *Internet of Things*. 2020 Jun
581 6:100239.
- 582 32. Dhiman C, Vishwakarma DK. A review of state-of-the-art techniques for abnormal human activity recognition.
583 *Engineering Applications of Artificial Intelligence*. 2019 Jan 1;77:21-45.
- 584 33. Hussain Z, Sheng QZ, Zhang WE. A review and categorization of techniques on device-free human activity
585 recognition. *Journal of Network and Computer Applications*. 2020 Jun 23:102738.
- 586 34. Banu PN, Kavitha R. Single Activity Recognition System: A Review. *Internet of Things (IoT) 2020* (pp. 257-271).
587 Springer, Cham.
- 588 35. Levy CR, Alemi F, Williams AE, Williams AR, Wojtusiak J, Sutton B, Giang P, Pracht E, Argyros L. Shared homes
589 as an alternative to nursing home care: Impact of VA's medical foster home program on hospitalization. *The
590 Gerontologist*. 2016 Feb 1;56(1):62-71.
- 591 36. Hawes C, Morris JN, Phillips CD, Mor V, Fries BE, Nonemaker S. Reliability estimates for the Minimum Data Set
592 for nursing home resident assessment and care screening (MDS). *The Gerontologist*. 1995 Apr 1;35(2):172-8.
- 593 37. Handelman GS, Kok HK, Chandra RV, Razavi AH, Lee MJ, Asadi H. eD octor: machine learning and the future of
594 medicine. *Journal of internal medicine*. 2018 Dec;284(6):603-19.
- 595 38. Breiman L. Random forests. *Machine learning*. 2001 Oct 1;45(1):5-32.
- 596 39. Breiman L. Bagging predictors. *Machine learning*. 1996 Aug 1;24(2):123-40.
- 597 40. Olson MA, Wyner AJ. Making sense of random forest probabilities: a kernel perspective. *arXiv preprint
598 arXiv:1812.05792*. 2018 Dec 14.
- 599 41. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R,
600 Dubourg V, Vanderplas J. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*.
601 2011 Nov 1;12:2825-30.
- 602 42. Matplotlib: Python plotting — Matplotlib 3.2.2 documentation [Internet]. [cited 2020 Jun 25]. Available from:
603 <https://matplotlib.org/>

- 604 43. Wojtusiak J. Machine Learning and Inference Reporting Criteria. Reports of the Machine Learning and Inference
605 Laboratory, MLI 20-1.2020.
- 606 44. Computational Barthel Index (CBIT) for Activities of Daily Living [Internet]. [cited 2020 Jun 25]. Available from:
607 <https://hi.gmu.edu/cbit>.
- 608 45. Stenholm S, Westerlund H, Salo P, Hyde M, Pentti J, Head J, Kivimäki M, Vahtera J. Age-related trajectories of
609 physical functioning in work and retirement: the role of sociodemographic factors, lifestyle and disease. J
610 Epidemiol Community Health. 2014 Jun 1;68(6):503-9.
- 611 46. Nisar MA, Shirahama K, Li F, Huang X, Grzegorzec M. Rank Pooling Approach for Wearable Sensor-Based ADLs
612 Recognition. Sensors. 2020 Jan;20(12):3463.
- 613 47. Poli A, Scalise L, Spinsante S, Strazza A. ADLs Monitoring by Accelerometer-Based Wearable Sensors: Effect of
614 Measurement Device and Data Uncertainty on Classification Accuracy. In2020 IEEE International Symposium
615 on Medical Measurements and Applications (MeMeA) 2020 Jun 1 (pp. 1-6). IEEE.
- 616 48. Vepakomma P, De D, Das SK, Bhansali S. A-Wristocracy: Deep learning on wrist-worn sensing for recognition of
617 user complex activities. In2015 IEEE 12th International conference on wearable and implantable body sensor
618 networks (BSN) 2015 Jun 9 (pp. 1-6). IEEE.
- 619

620 Figure Legends

621 **Figure 1 Predicted probability visualization of functional independence for a hypothetical patient up to one year**
622 **ahead.**

623 **Figure 2 Prediction timeline.** Past EHR data is used to make an assessment of current functional status as well as
624 prediction 3, 6 months, and 1 year afterwards.

625 **Figure 3 Average AUC for the current assessment based on the number of attributes used.** The blue line refers to
626 MS_{RE}^d models and the orange line refers to MS_E^d models.

627 **Figure 4 Example of the calibration curve for model that predicts bathing at current time point.** The shape of the curve
628 indicates that the model is well-calibrated. Similar curves were created for all models used in CBIT.

629 **Figure 5 Predicted probability visualization of functional independence for three patients up to one year ahead.**

630 Additional Files

631 Additional file 1.doc - Barthel Index categories of functional abilities along with assigned scores. Reproduced from:
632 <http://www.strokecenter.org/wp-content/uploads/2011/08/bartel.pdf>

633 Additional file 2.pdf- Detailed performance of the models

634 Additional file 3.doc -Top ranked predictors of functional status

635 The table includes top 50 attributes across Re-Evaluation and Evaluation models. Previous evaluations results associated
636 with the Re-evaluation models (M_{RE}^d) were included at the beginning of the table. Gender and race along with their ranking
637 were also added at the bottom of the table for comparison.

638 Additional file 4.pdf - Learning_Curves_Full_Evaluation_Models

639 Additional file 5.pdf - Learning_Curves_Full_Re-Evaluation_Models

640 Additional file 6.pdf - Learning_Curves_Simplified_Evaluation_Models

- 641 Additional file 7.pdf- Learning_Curves_Simplified_Re-Evaluation_Models
- 642 Additional file 8.pdf - Calibration_plots_Full_Evaluation_Models.
- 643 Additional file 9.pdf - Calibration_plots_Full_Re_Evaluation_Models
- 644 Additional file 10.pdf -Calibration_plots_Simplified_Evaluation_Models
- 645 Additional file 11.pdf - Calibration_plots_Simplified_Re_Evaluation_Models
- 646