

1
2
3
4
5
6
7
8

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. We investigated the possibility of constructing an automated decision support tool, the Computational Barthel Index Tool (CBIT), that automatically assesses and predicts probabilities of current and future ADLs based on patients' medical history.

Methods: The data used to construct the tool include the demographic information, 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 isonomic 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, 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. This paper examines whether functional ability can be assessed through
40 coded data available in Electronic Health Records (EHR) or medical claims. We focus on nine out of ten functional abilities
41 which are part of the Barthel Index [4,5] as described in the Data section.

42 The ability to automatically derive and predict patients' functional status has several important uses in clinical work
43 and research. Firstly, it may provide a more efficient and cost-effective means of assessing functional status in groups
44 for whom functional status is currently manually assessed. In a recent review that examined functional status quality
45 indicators, the authors concluded that using chart reviews or patient-reports is costly and administratively burdensome
46 [6]. Secondly, it may allow for retrospective assessment of patients' functional status for whom evaluations have not
47 been completed. Thirdly, it can be beneficial for patients who are typically not evaluated for the purpose of comparing
48 care across settings. Finally, predicting functional status up to one year in the future provides a basis for an informed
49 discussion between clinicians and patients/caregivers and may help in planning care for patients.

50 Previously, we constructed a set of models capable of predicting trajectories of ADL improvement or decline post-
51 hospitalization [7], as well as sequences of functional decline [8]. The method and tool discussed in this paper, called
52 the Computational Barthel Index Tool (CBIT), significantly extends the previous work and is designed to allow for
53 assessment of functional status at any arbitrary moment. The tool allows for prediction up to one year ahead. It also
54 extends our previous work [7] by incorporating temporal information about when events happened in the patient's
55 medical history. Many diagnoses present in medical records correlate with the patient's functional ability, with some of
56 these correlations being temporary and others being permanent. For example, some surgical patients have urinary
57 incontinence for a short period after the surgery, while amputation affects the ability to walk permanently. Thus, we
58 hypothesize that the codes present in data are time-dependent. We showed that proper handling of temporal
59 relationships can improve the accuracy of the constructed CBIT models, as discussed later in the paper.

60 Prediction of functional status and disability is challenging. Researchers in many studies have attempted to
61 automatically assess and predict functional status, including ADLs. In one study, machine learning (ML) methods were

62 linked to biomedical ontologies to predict functional status [9], achieving predictive accuracy of 0.6. In another work,
63 researchers described a logistic regression-based method to predict mortality and disability post-injury for the elderly
64 [10] with reported R^2 of 0.86. Tarekegn et al. developed a set of models to predict disability as a metric for frailty
65 conditions resulting in models with F-1 scores ranging between 0.74 to 0.76 [11]. Similarly, Gobbens and van Assen
66 examined six standard frailty indicators (gait speed, physical activity, hand grip, body mass index, and fatigue and
67 balance) for assessing ADLs, of which only gait speed was predictive of ADL disabilities [12]; however, no actual
68 predictive accuracy was reported. It is clear that the above studies reported model performances below ones reported
69 here. However, it should be mentioned that these works were performed in different settings thus no direct comparison
70 is meaningful. A systematic review of published works related to assessing ADLs identified several commonly used
71 predictors, including age, cognitive functioning, depression, and hospital length of stay [13]. In the data-driven approach
72 presented here, some of the predictors are the same as those previously reported in the literature. However, in this
73 study, a more comprehensive approach derived entirely from data is used to identify important patient characteristics.
74 Besides using demographic characteristics and the latest functional status (if available), we included patients' diagnoses
75 collected longitudinally.

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

90 <INSERT FIGURE 1 HERE>

91 In the presented work, we constructed two sets of models: *Evaluation models*, M_E^d , in which previous functional
92 status assessment is unknown, and *Re-Evaluation models*, M_{RE}^d , in which previous functional status is known. Here d is
93 an ADL (bathing, grooming, etc.), and $\tau \in \{0,90,180,365\}$ is the prediction horizon (given as the number of days), i.e.,
94 how far ahead in time the value is predicted. As names suggest, M_E^d models are used in situations in which a new patient
95 is being evaluated in terms of ADLs, and M_{RE}^d models are used when an evaluation of the previously assessed patient
96 needs to be refreshed as new information becomes available.

97 The presented research has been initiated as part of a larger project in the Department of Veterans Affairs (VA) with
98 the purpose of assessing the cost and effectiveness of the Medical Foster Home (MFH) program compared to traditional
99 Community Living Centers (CLC; nursing homes) [7,8,14]. Determination of patients' functional status was used as one
100 of the characteristics to match MFH and CLC residents for comparison purposes. The main contributions of the presented
101 work are in the systematic development of models for prediction of ADLs by incorporating temporal information into
102 encoding diagnoses, detailed testing and analysis of the developed models, and creation of an online tool.

103 **Methods**

104 **Data**

105 We extracted data from the Department of Veterans Affairs Corporate Data Warehouse (CDW) and analyzed within
106 the VA Informatics and Computing Infrastructure (VINCI). The data were organized around patient evaluations using
107 Minimum Data Set 2.0 [15], which were mapped to the nine Barthel Index categories using a procedure used in previous
108 research [7]. The Barthel Index (or Barthel Score), which measures independence in performing ADLs [4,5] includes 10
109 items with the total value ranging from 0 to 100 (feeding, bathing, grooming, dressing, bowels, bladder, toilet use,
110 transfers, mobility, and stairs). In this research, we eliminated the last item of the Barthel Score, stairs, which was not
111 consistently assessed and thus difficult to standardize among nursing home residents. Thus, the total considered scale
112 is 0-90 based on the first nine items predicted independently.

113 The data consisted of 1,901,354 MDS evaluations completed between 2000 and 2011 from which 1,151,222
114 complete evaluations were retrieved for 295,491 patients. The data were linked to medical records from which
115 demographics and history of diagnoses were extracted. The data consisted of 18,912,553 inpatient and 180,123,710
116 outpatient diagnosis codes using the International Classification of Diseases, ninth edition (ICD-9) standard along with
117 corresponding dates. These codes were transformed into clinically relevant categories using Clinical Classification
118 Software (CCS) from the Agency of Health Research and Quality (AHRQ) resulting in 281 distinct CCS codes representing
119 health comorbidities. All diagnosis codes were combined from inpatient and outpatient records. Distinguishing between

120 inpatient and outpatient codes is important for some applications (inpatient codes are typically treated as more severe).
121 In the presented work, we assumed that only information about the presence of a diagnosis along with appropriate time
122 was important in the context of predicting disabilities, rather than distinguishing between the specific sources. We also
123 included demographic information including age, race, and gender. Age was recorded as a continuous variable and race
124 was represented as a set of binary indicators (sometimes referred to as dummies or one-hot coding). Missing data for
125 age were imputed as mean value in the dataset and no special treatment for missing data for other attributes was
126 needed. Patients with only one MDS evaluation were excluded to allow for modeling of change of patient status over
127 time, resulting in a final dataset of 855,731 evaluations for 181,213 patients. The collected data were organized per MDS
128 evaluation, resulting in the average of 4.72 +/- 6.21 MDS evaluations per patient. Table 1 shows descriptive statistics of
129 the final dataset as counted in analyzed MDS records as well as per patient, and is representative of the overall CLC
130 population in the VA. The majority of patients were male and white with an average age of over 71 years and mean
131 Barthel Score (sum of assigned Barthel items) of about 48 out of 90, indicating overall high levels of disability in the
132 studied population. In addition, the average score at the first evaluation was about 52. The average time between MDS
133 evaluations was also about 100 days, which is slightly over three months.

134 **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% | 96.9% |
| | Female | 3.32% | 3.1% | 3.1% |
| Race | Asian | 1.5% | 1.4% | 1.41% |
| | Black | 13% | 11.9% | 12.02% |
| | White | 58.8% | 55.4% | 55.03% |
| | Other | 26.7% | 31.3% | 31.53% |
| 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 |

135 In the VA's CDW, as well as in many administrative datasets, patient medical records often span many years, making
 136 it possible to examine temporal relationships between diagnoses and the predicted events. In the presented research,
 137 we encoded temporal information about diagnoses by calculating the number of days from the first known occurrence
 138 of the i -th, diagnosis code (t_i) to the time of MDS evaluation (t_p),

139
$$ccs_i^{max} = \max_{t_i}(t_p - t_i) \quad (1)$$

140 as well as last recorded occurrence of the diagnosis code relative to the time of MDS evaluation,

141
$$ccs_i^{min} = \min_{t_i}(t_p - t_i) \quad (2)$$

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

150 An example of data encoded using the above method is presented in Table 2. The table shows data for two different
 151 fictitious patients. Patient 1 has two MDS evaluations in the data 90 days apart. Patient 2 also has two MDS evaluations
 152 100 days apart. Patient 1 was diagnosed with septicemia only once, 210 days prior to the first evaluation
 153 ($ccs_2^{min}=ccs_2^{max}=210$). We know that the patient has not been diagnosed second time between the evaluations because
 154 both columns representing the first and most recent occurrence increased by the same amount. The patient has been
 155 diagnosed with hypertension 18 days prior to the first evaluation ($ccs_{99}^{min}=18$), and for the first time 500 days prior to
 156 the first evaluation. The patient has been diagnosed with hypertension again 5 days prior to the second evaluation.
 157 Similarly, Patient 2 has been diagnosed with septicemia twice, 15 and 700 days prior to the first evaluation ($ccs_2^{min}=15$
 158 and $ccs_2^{max}=700$). Patient 2 was also diagnosed with tuberculosis 71 days before the second evaluation
 159 ($ccs_1^{min}=ccs_1^{max}=71$). One can also notice that Patient 1's ADLs declined between the evaluations. Diagnoses not
 160 present/recorded in patient's records are coded as -999999 and 999999.

161

162 **Table 2 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 |

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

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

171 Construction of models

172 Patients were randomly assigned to training (90%) and testing (10%) sets. The testing set with a sufficiently large
173 sample size (approximately 18,000 patients) was used to provide final validation of the models. Training dataset was
174 used for cross-validated parameter tuning and final model construction. We applied a selection of machine learning
175 (ML) methods to construct models capable of assessing and predicting ADLs. Machine learning methods are rapidly
176 gaining popularity in medical and health applications [16] and are also applicable to the prediction of ADLs. Machine
177 learning is an experimental field that provides a large toolset of methods that can be used for prediction. In the
178 presented work, we evaluated selected ML methods (regularized logistic regression, Bayesian networks, decision trees,
179 and random forests) in terms of their performance, and showed that random forest stands out in terms of model quality.
180 Random forests [17] are ensembles of decision trees (typically many), DT_i , that are inferred from randomly selected
181 subsets of data (X_i, Y_i) , thus guaranteed to be different on sufficiently large data.

$$182 \quad RF(X, Y) = \left\{ \begin{array}{l} DT_i(X_i, Y_i): D(X_i) \subset D(X) \\ \wedge R(X_i) = R(Y_i) \subset R(X) = R(Y) \end{array} \right\} \quad (3)$$

183 In equation (3), D indicates a set of attributes in the data and R corresponds to a set of examples in the data. Random
184 forests are created by applying bagging (a.k.a., bootstrap aggregation) [18] to both sample and attributes (patient
185 characteristics). Each decision tree is created on a random sample of data (with repetitions) drawn from the dataset and
186 a random set of attributes drawn from the list of all available attributes. Standard top-down decision tree learning
187 algorithms are used. The process is repeated to create multiple trees (typically in the order of tens or hundreds). After

188 a forest is assembled, the final classification decision is made by applying all of the trees to new examples. When there
189 is a disagreement in prediction, the trees vote on the predicted outcome. Random forests output classification scores
190 (in the presented work they were converted to probabilities) which in the case of the described models represent
191 patients being disabled or functionally independent. These scores are calculated as a proportion of trees voting for a
192 given outcome [19]. In the presented work, we performed 10-fold cross-validated hyperparameter tuning for random
193 forests. The tuning led to the selection of random forests consisting of about 100 decision trees (each model was
194 optimized separately, and the numbers of trees were different). Other algorithm parameters, including the number of
195 randomly selected patient characteristics and Gini Index [17] as an internal quality criterion were tested and set to
196 default as they did not make any improvements.

197 The models were created to assess functional status at the time of prediction (current status), as well as to predict
198 functional status 3, 6, and 12 months beyond the time of prediction as depicted in Figure 2. Data available prior to the
199 time of prediction were used to encode input attributes for the model.

200 <INSERT FIGURE 2 HERE>

201 The quality of the constructed models was evaluated in terms of standard statistical measures used in ML, namely,
202 accuracy (percentage of correctly predicted cases), area under the curve (AUC; often referred to as C-statistic), recall
203 (rate of correctly identified patients with functional dependencies), and precision (rate of patients with disabilities
204 among those indicated as disabled by the model). We applied the evaluation on the test set of patients not being used
205 in model construction, selection or tuning. In order to provide better insight into the created models, we also created
206 calibration plots and learning curves for all 144 developed models. It should be mentioned that the presented work
207 does not include clinical validation of the models. Also, note that the created models predict the probability of functional
208 dependence of any level, while the graphical representation or prediction (in the web calculator discussed later and
209 presented figures) shows the probability of functional independence. The conversion between the two is a simple
210 operation, which is one minus probability. The reason for this conversion is that prediction of disability as a target event
211 is conceptually cleaner from a machine learning perspective (assuming that being independent is normal, we predict the
212 abnormal state of disability). On the other hand, clinicians are used to having higher values represent better status (this
213 can also be the case in in the original Barthel Index). This conversion has no effect on presented results or modeling and
214 is only reflected in the graphical representation of results.

215 For the data analysis part of the project, we used the Microsoft SQL Server to preprocess data within the VINCI. The
216 data preprocessing started with MDS evaluations that were later linked to other data components. Final data were

217 analyzed using Python programming language with Scikit-learn machine learning library [20] and visualizations were
 218 done using Matplotlib Python library [21].

219 Results

220 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.
 221 The CBIT can assess the level of functional dependency in performing ADLs and predicting functional dependency up to
 222 one year ahead by using demographics, diagnoses, and (if available) last known functional status. Table 3 presents a
 223 summary of the performance of the models for each ADL at the time of prediction, as well as 3, 6 and 12 months ahead
 224 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
 225 nine outcome categories. The CBIT showed very high accuracy in assessing ADLs at a given time. The AUC of assessing if
 226 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),
 227 precision 0.91 (0.89-0.92), and recall 0.90 (0.84-0.95). When predicting functional status up to one year ahead, $\tau \in$
 228 $\{90,180,365\}$, the $M_{RE^d_\tau}$ models' accuracy drops to AUC 0.77 (0.73-0.79), accuracy 0.73 (0.69-0.80), precision 0.74 (0.66-
 229 0.81), and recall 0.69 (0.34-0.96). When the previous functional status is unknown (i.e., initial evaluation), the
 230 performance of the current assessment models $M_{E^d_0}$ decreased by about 16% ($p < 0.01$) in terms of AUC. On average,
 231 the obtained results for these models are AUC 0.79, accuracy 0.74, precision 0.74, and recall 0.80.

232 **Table 3 Average +/- standard deviation of accuracy, AUC, precision and recall of models in predicting functional status**

| Prediction | Re-Evaluation Models ($M_{RE^d_\tau}$) | | | | Evaluation Models ($M_{E^d_\tau}$) | | | |
|-------------|--|-------------|-------------|-------------|--------------------------------------|-------------|-------------|-------------|
| | 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 |

233 Top predictors

234 Further analysis also identified the top predictors used in the assessment and prediction of ADLs. We used average
 235 Gini Index [17] produced by random forest to measure the quality of predictors. Gini index is used by random forest as
 236 an internal measure of attribute quality when constructing individual decision trees. It should not be interpreted as a
 237 strength or effect of the variable on the predicted output, but rather to understand the relative importance of attributes.
 238 Top predictors along with their reported importance (average Gini index over all trees in forest and over all models) are

239 presented in Table 4. Note that all ccs_i^{min} and ccs_i^{max} codes were included in full models. However, a comprehensive list
 240 of diagnosis codes and previous evaluations are not shown in the table due to space limitation [see Additional file 1 for
 241 more details]. Not surprisingly, the most predictive attributes in M_{RE}^d models were past functional status, being
 242 responsible for AUC of 0.93. Other most predictive attributes were the time since the most recent diagnosis of delirium,
 243 dementia, and amnestic and other cognitive disorders (CCS 653) and patient age. These were followed by encoded time
 244 of diagnoses/administrative codes for: the urinary tract infections (CCS 159); chronic ulcer of skin (CCS 199); other
 245 connective tissue disease (CCS 211); paralysis (CCS 82); administrative/social admission (CCS 255); alcohol-related
 246 disorders (CCS 660); aspiration pneumonitis; food/vomitus (CCS 129); and schizophrenia and other psychotic disorders
 247 (CCS 659). For most of the diagnoses listed above, it is important when (number of days) a patient was diagnosed with
 248 that condition most recently. For ulcers and aspiration pneumonitis; food/vomitus the first diagnosis is important. In
 249 addition, the table has marked potentially reversible conditions (R), as judged by clinicians, which can be influenced in
 250 the care provided to the patients and affect the outcome.

251 **Table 4 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 |

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

255 **Simplified models**

256 Further, we analyzed the possibility of creating simplified models (called $MS_{RE^d_\tau}$ and $MS_{E^d_\tau}$) that include only selected
257 top-ranking patient characteristics. As depicted in Figure 3, adding more characteristics beyond the most predictive 41
258 attributes did not significantly improve the accuracy ($p<0.05$) of the models in assessing the current functional status
259 ($\tau=0$). The curves were also similar for predicting up to 12 months ahead, $\tau \in \{90,180,365\}$. When using 25 top patient
260 characteristics, models that included previous evaluations ($MS_{RE^d_\tau}$) reached an average AUC of 0.94, accuracy 0.90,
261 precision 0.91, and recall 0.90. Furthermore, the performance of the simplified models with 41 patient characteristics
262 and without previous evaluations ($MS_{E^d_\tau}$) raised to average AUC of 0.79, accuracy 0.74, precision 0.74, and recall 0.78.
263 Note that top predictors for each ADL are different. In the $MS_{RE^d_\tau}$ and $MS_{E^d_\tau}$ models, we included top ranking attributes
264 across all models to minimize information needed by CBIT for all ADLs, even though this set of attributes may not be
265 optimal for individual models.

266 <INSERT FIGURE 3 HERE>

267 **Temporal coding**

268 One important advancement of the presented CBIT is the way it captures time in encoding diagnoses as previously
269 shown in equations (1) and (2) and illustrated in Table 2. We investigated how our proposed method in the encoding of
270 diagnoses would be different from binary coding (1 when a diagnosis is present in a given patient's record and 0
271 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 compared in terms of AUC at
272 different time points up to one year ahead. In one experiment, random forest was compared with other algorithms
273 including logistic regression, decision tree, and naïve Bayes. As described earlier, when temporal coding is used, one
274 needs to assign special values to diagnoses that are not present in data. Therefore, we compared +/- 999999 (6_9) with
275 +/-9999 (4_9), and +/-99999 (5_9) coding across all models (here X_9 indicates 10^X-1). Temporal coding (6_9) was also
276 compared with binary coding to determine any significant difference. Two-tailed t-test was used to assess all
277 comparisons ($p<0.05$).

278 As summarized in Table 5, both random forest and logistic regression show a significant difference in AUC when
279 temporal information is applied ($p<0.05$). The results indicated that random forest with the temporal coding performs
280 significantly better than binary coding, while for the logistic regression the relationship is opposite (the binary coding is
281 better). However, logistic regression with binary coding is still doing worse than random forest. We also included
282 decision trees and naïve Bayes results in the table, but the performance was typically inferior. We observed that random
283 forest, decision tree and naïve Bayes are not affected by how the special values were assigned, while the performance

284 of logistic regression is affected by the coding. The rationale for this result is that for symbolic methods it is irrelevant
 285 how not-present values are coded as long as the value is distinct, while parametric models need to find a coefficient for
 286 each diagnosis code, which is affected by the coding.

287 **Table 5 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** | |
| $MS_{RE}^{d_t}$ | 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** | |
| $MS_{RE}^{d_t}$ | 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** | |

288 The results are presented in terms of AUC for the current assessment and prediction up to 12 months ahead. Full models that include 578
 289 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
 290 history for +/-9999, +/-99999, and +/-999999, respectively. * indicates significance (p<0.05) of coding systems compared to "6_9" and + indicates
 291 significance (p<0.05) of different algorithms compared to random forest.

292 **Calibration**

293 Calibration allows for the probability interpretation of the output scores from the models, further allowing for
 294 frequency interpretation of the results. Thus, we calibrated all models using 5-fold cross-validated isonomic regression.
 295 The results showed that the models were well-calibrated with an average error of about 3%. Figure 4 shows an example
 296 of the calibration curve for the model that assesses bathing at the current time point, $MS_{RE}^{bathing_0}$. Similar curves were
 297 also developed for all 144 models and are available at the project website [22].

299 Discussion

300 Methods

301 We demonstrated that it is possible to assess and predict functional status using machine learning methods.
302 Moreover, we showed that the inclusion of time between diagnosis and time of prediction (temporal information) is
303 important in encoding the data. While further work is needed to validate the new encoding schema and study its
304 limitations, the encoding of the first and last known occurrence of a diagnosis code works for the problem at hand. We
305 are currently investigating this method on independent datasets.

306 Machine learning methods are gaining popularity in medical and health applications, yet there is no consensus on
307 what validation is needed for their use in clinical settings. There is also no agreement about what information is needed
308 to allow for full reproducibility of ML results, or even what reproducibility in this context means [23]. Models created
309 for CBIT were evaluated using standard measures in ML model testing (cross-validation, independent test set, etc.), and
310 investigated in terms of their calibration and learning curves. There is a need for further validation of the models and
311 their impact on patient care. Such validation focuses on detailed model analysis in terms of accuracy, transparency and
312 the ability to provide explanations, and eventually trust and acceptability by the medical community. A randomized trial
313 to assess outcomes of the model use may be required for full acceptance in clinical settings.

314 In addition, there is an ongoing discussion about the overall validity of applying machine learning methods to the
315 prediction of patient outcomes, and potential bias of the models based on gender, race and socioeconomic status. One
316 needs to clearly understand data limitations and definitions of the prediction problem to understand the drawbacks of
317 the method. Supervised machine learning methods, by definition, learn what they are asked to learn, and may (typically
318 do) propagate biases from training data. Biases in machine learning-based models are not caused by machine learning,
319 but by underlying process used to create training data. The key is in the definition and construction of the output
320 attributes of the model and their proper interpretation. One needs to answer a question if models predict events in real
321 world, or data artifacts that somehow approximate that reality. Similarly, CBIT is intended to mimic the tasks of nurses
322 performing evaluations of ADLs as part of the MDS. Therefore, any biases, inaccuracies, or subjectivity in this process
323 may also be repeated by the CBIT models. However in CBIT, as shown in Table 4, race and gender had only a negligible
324 impact on predictions and were completely dropped from simplified models, which suggests diminishing the potential
325 racial and gender bias in the models. A different set of methods, typically used in health services research, is needed to
326 understand the existence of potential bias in constructed models.

327 Further, the probabilistic interpretation of the prediction results used in the presented work seems to be reasonable.
328 Conceptually, the future can never be predicted with the probability of one (even though for some cases the models
329 may be certain of the future and output 1.0). Instead, the values represent how likely an event (here disability) will occur
330 according to the models. Such interpretation has several advantages. It allows end users to interpret the chances of an
331 event happening, and in turn, describes when models are uncertain about the outcomes. It also explains why model
332 accuracy is not 100% when executed on the test data (examples with predicted probability not equal to one, are
333 ambiguous by the definition of probability). The latter is the most evident when analyzing calibration curves, such as
334 one presented in Figure 4. The disadvantage of using the probabilities is that they may be misinterpreted as severity of
335 disability. When presenting results, one needs to specify that the number represent how likely a patient is dependent,
336 and not the level of dependency.

337 The created models in this manuscript are based on ICD-9 diagnosis codes, which were mapped to CCS codes. One
338 advantage of this approach is that since all new data are coded with ICD-10 codes, they could easily be mapped to CCS
339 codes making the models applicable to data with the newer coding system. Another important issue is that the diagnosis
340 codes in both EHR and claims data are subject to under- and over-coding, thus affecting the potential reliability of the
341 models. However, it is important to note that our modeling efforts were not intended to understand the effects of
342 diagnoses on ADLs, but rather their use in making prediction. In addition, as long as diagnoses are systematically
343 over/under-coded, they should not affect performance of the models. Despite these limitations, results indicated that
344 our data were appropriate for this purpose.

345 The evaluations presented in this paper are only summaries and examples of detailed results. We performed a
346 detailed examination 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. Due
347 to very large amount of material to be presented, the results are available through the online calculator [22]. All
348 developed models and source codes are available for everyone who wishes to conduct their testing on independent
349 data from other institutions, i.e. to test cross-institution generalizability.

350 **Clinical and administrative use**

351 Very little evidence exists to address whether measuring functional status can change the quality of life, but our
352 research shows that prior knowledge about functional disability is a key indicator of future functional status. Notably,
353 past research has provided evidence that improvements in functional status are possible over time through therapy [24]
354 by improving, slowing decline, and/or maintaining functional status. The presented CBIT tool which predicts
355 improvement or decline could be used by health professionals as means of identifying patient characteristics that are

356 modifiable and plan care accordingly. It can serve as a basis for an informed discussion between clinicians, patients and
357 caregivers. In addition, these measures could potentially serve as a patient-centered measure for examining the value
358 of the services provided.

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

360 A graphical representation of the assessment and prediction of functional status can be used by healthcare
361 professionals and caregivers for decision making regarding the patients' care. Our full models can be integrated as
362 decision support tools within EHR systems or linked to claims data, while the simplified models can operate standalone
363 as an interactive online tool. For example, Figure 5 illustrates CBIT-predicted outcomes for three fictitious patients
364 similar to what was shown in Figure 1. Values indicate the probability of functional independence for each ADL up to
365 one year after the prediction time. The higher the value is, the higher the chance that the patient is functionally
366 independent. One can observe significant differences between the functional dependency trajectories for these
367 patients. Patient (a) is currently likely to be independent but expected to decline within 6 months as the probability of
368 independence decreases. Patient (b) is currently likely to be dependent in most ADLs (probability of independence
369 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
370 independent and predicted to remain independent in terms of walking and is almost certainly disabled in terms of
371 bladder, bowels and eating. The patient is likely to have a temporary decline in terms of other ADLs. Construction of
372 each of the plots requires execution of 36 random forest models (9 ADLs, 4 time points).

373 <INSERT FIGURE 5 HERE>

374 An experimental version of the online calculator that takes patient characteristics and outputs plots is available at
375 <https://hi.gmu.edu/cbit> [22]. It is accessible through a web form or an application programming interface (API). The web
376 calculator is implemented in Python 3 and uses Flask as a web application framework, with Pandas and Scikit-learn
377 libraries performing data analysis. To ensure the performance of the web calculator, all of the models are loaded on the
378 startup and reside in RAM. Additional changes have been made to the calculator to improve clinical use. For example,
379 for the ease of use, the numbers of days associated with diagnoses were discretized to allow users to select them from
380 drop-down menus. Numbers closer to zero are discretized with higher precision than larger numbers, which further
381 improves understandability.

382 Figure 6 shows part of the interface used to enter data. Users can enter patient information and are provided with
383 results similar to those shown in Figure 5 along with a data table containing the values of predicted probabilities. We

384 are currently working on developing an explanation module that provides human-oriented interpretation of the results
385 as well as the reasons for predictions.

386 <INSERT FIGURE 6 HERE>

387 **Conclusion**

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

397 The constructed full model requires a large number of predictors, which makes it impossible to manually enter
398 values. Hence, the full version of CBIT would need to be integrated with an EHR or claims management system to be
399 part of the clinical decision support. Such integration can be achieved using HL7s FHIR interface. The simplified version
400 of the CBIT that uses 50 predictors is available within a web calculator. Beyond the use of EHR data, the constructed
401 CBIT could be enhanced by sensor data allowing for continuous patient monitoring and be integrated with the presented
402 approach. Such data can aid assessment, particularly for ADLs that measure patient movement [25,26,27].

403 The presented work has a number of limitations. The tool is not applicable in settings in which longitudinal patient
404 records are not available. Only large health systems with long-established electronic medical records have sufficient
405 longitudinal data to apply models that use temporal diagnosis information. Additionally, the models were developed
406 using data from the US Department of Veterans Affairs (VA), which does not reflect the general population of nursing
407 home residents outside of the VA system. We are currently investigating the performance of these models on other
408 datasets, including Medicare claims data. Depending on the results of the evaluation, one can adapt the developed
409 methods to non-VA populations using transfer learning methods [28]. Finally, random forests are known to be “black
410 box” models that work well but are not well understood by end users. As part of the online calculator, we are designing
411 an explanation module that provides users with “reasons” for making specific predictions in one individual case. The

412 reasons consist of a list of patient characteristics that are the strongest predictors (both confirming and disconfirming)
413 for that individual case.

414 Despite these limitations, CBIT can be used to support clinicians and administrators in decision making. Our novel
415 data coding method, applying machine learning to unique health data, comprehensive model testing, and transparency
416 of the work contribute to the state-of-the-art in ML-based decision support.

417 **Abbreviations**

418 ADL: Activity of Daily Living; CBIT: Computational Barthel Index Tool; VA: Department of Veterans Affairs; MDS: Minimum Data Set; EHR: Electronic Health Record; ML:
419 Machine learning; MFH: Medical Foster Home; CLC: Community Living Center; CDW: Corporate Data Warehouse; VINCI: VA Informatics and Computing Infrastructure; ICD-
420 9: International Classification of the Diseases, ninth edition; CCS: Clinical Classification Software; AHRQ: Agency of Health Research and Quality Control; AUC: Area under
421 the curve; RF: Random forest; LG: Logistic Regression; DT: Decision tree; NB: Naïve Bayes; ICD-10: International Classification of the Diseases, tenth edition; API: Application
422 programming interface; RAM: Random access memory.

423 **Declarations**

424 **Ethics approval and consent to participate**

425 Not applicable.

426 **Consent for publication**

427 Not applicable.

428 **Availability of data and material**

429 The data used to construct presented CBIT models are individual level and cannot be shared. All constructed models, source code, and detailed testing results are freely
430 available at the project website.

431 **Competing interests**

432 The authors declare that they have no competing interests.

433 **Funding**

434 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
435 the Department of Veterans Affairs or the United States Government.

436 **Author's contributions**

437 Dr. J.W is the main author of the manuscript. He designed and partially implemented the method, supervised experimental evaluation, and implemented online
438 calculator. Ms. N.A partially implemented the method, constructed models, performed experiments, and wrote large section of experimental results. Dr. F.A contributed
439 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
440 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.

441 **Acknowledgments**

442 Not applicable.

443 **Authors' information**

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

446 **References**

447 1. Fried TR, Bradley EH, Towle VR, Allore H. Understanding the treatment preferences of seriously ill patients. *New*
448 *England Journal of Medicine*. 2002 Apr 4;346(14):1061-6.

2. McCarthy EP, Phillips RS, Zhong Z, Drews RE, Lynn J. Dying with cancer: patients' function, symptoms, and care preferences as death approaches. *Journal of the American Geriatrics Society*. 2000 May;48(S1):S110-21.
3. MDS 3.0 Technical Information [Internet]. Available from: <https://www.cms.gov/Medicare/Quality-Initiatives-Patient-Assessment-Instruments/NursingHomeQualityInits/NHOIMDS30TechnicalInformation>.
4. Collin C, Wade DT, Davies S, Horne V. The Barthel ADL Index: a reliability study. *International disability studies*. 1988 Jan 1;10(2):61-3.
5. Shah S, Vanclay F, Cooper B. Improving the sensitivity of the Barthel Index for stroke rehabilitation. *Journal of clinical epidemiology*. 1989 Jan 1;42(8):703-9.
6. Dy SM, Pfoh ER, Salive ME, Boyd CM. Health-related quality of life and functional status quality indicators for older persons with multiple chronic conditions. *Journal of the American Geriatrics Society*. 2013 Dec;61(12):2120-7.
7. Wojtusiak J, Levy CR, Williams AE, Alemi F. Predicting functional decline and recovery for residents in veterans affairs nursing homes. *The Gerontologist*. 2016 Feb 1;56(1):42-51.
8. Levy CR, Zargoush M, Williams AE, Williams AR, Giang P, Wojtusiak J, Kheirbek RE, Alemi F. Sequence of functional loss and recovery in nursing homes. *The Gerontologist*. 2016 Feb 1;56(1):52-61.
9. Min H, Mobahi H, Irvin K, Avramovic S, Wojtusiak J. Predicting activities of daily living for cancer patients using an ontology-guided machine learning methodology. *Journal of biomedical semantics*. 2017 Dec 1;8(1):39.
10. Jeffery AD, Dietrich MS, Maxwell CA. Predicting 1-year disability and mortality of injured older adults. *Archives of gerontology and geriatrics*. 2018 Mar 1;75:191-6.
11. Tarekegn A, Ricceri F, Costa G, Ferracin E, Giacobini M. Predictive Modeling for Frailty Conditions in Elderly People: Machine Learning Approaches. *JMIR Medical Informatics*. 2020;8(6):e16678.
12. Gobbens RJ, van Assen MA. The prediction of ADL and IADL disability using six physical indicators of frailty: a longitudinal study in the Netherlands. *Current gerontology and geriatrics research*. 2014;2014.
13. Hoogerduijn JG, Schuurmans MJ, Duijnste MS, De Rooij SE, Grypdonck MF. A systematic review of predictors and screening instruments to identify older hospitalized patients at risk for functional decline. *Journal of clinical nursing*. 2007 Jan;16(1):46-57.
14. Levy CR, Alemi F, Williams AE, Williams AR, Wojtusiak J, Sutton B, Giang P, Pracht E, Argyros L. Shared homes as an alternative to nursing home care: Impact of VA's medical foster home program on hospitalization. *The Gerontologist*. 2016 Feb 1;56(1):62-71.
15. Hawes C, Morris JN, Phillips CD, Mor V, Fries BE, Nonemaker S. Reliability estimates for the Minimum Data Set for nursing home resident assessment and care screening (MDS). *The Gerontologist*. 1995 Apr 1;35(2):172-8.
16. Handelman GS, Kok HK, Chandra RV, Razavi AH, Lee MJ, Asadi H. eD octor: machine learning and the future of medicine. *Journal of internal medicine*. 2018 Dec;284(6):603-19.
17. Breiman L. Random forests. *Machine learning*. 2001 Oct 1;45(1):5-32.
18. Breiman L. Bagging predictors. *Machine learning*. 1996 Aug 1;24(2):123-40.
19. Olson MA, Wyner AJ. Making sense of random forest probabilities: a kernel perspective. *arXiv preprint arXiv:1812.05792*. 2018 Dec 14.
20. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*. 2011 Nov 1;12:2825-30.
21. Matplotlib: Python plotting — Matplotlib 3.2.2 documentation [Internet]. [cited 2020 Jun 25]. Available from: <https://matplotlib.org/>
22. Computational Barthel Index (CBIT) for Activities of Daily Living [Internet]. [cited 2020 Jun 25]. Available from: <https://hi.gmu.edu/cbit>.
23. Wojtusiak J. Machine Learning and Inference Reporting Criteria. Reports of the Machine Learning and Inference Laboratory, MLI 20-1.2020.
24. Stenholm S, Westerlund H, Salo P, Hyde M, Pentti J, Head J, Kivimäki M, Vahtera J. Age-related trajectories of physical functioning in work and retirement: the role of sociodemographic factors, lifestyle and disease. *J Epidemiol Community Health*. 2014 Jun 1;68(6):503-9.
25. Nisar MA, Shirahama K, Li F, Huang X, Grzegorzec M. Rank Pooling Approach for Wearable Sensor-Based ADLs Recognition. *Sensors*. 2020 Jan;20(12):3463.
26. Poli A, Scalise L, Spinsante S, Strazza A. ADLs Monitoring by Accelerometer-Based Wearable Sensors: Effect of Measurement Device and Data Uncertainty on Classification Accuracy. In 2020 IEEE International Symposium on Medical Measurements and Applications (MeMeA) 2020 Jun 1 (pp. 1-6). IEEE.
27. Vepakomma P, De D, Das SK, Bhansali S. A-Wristocracy: Deep learning on wrist-worn sensing for recognition of user complex activities. In 2015 IEEE 12th International conference on wearable and implantable body sensor networks (BSN) 2015 Jun 9 (pp. 1-6). IEEE.

28. Pan SJ, Yang Q. A survey on transfer learning. IEEE Transactions on knowledge and data engineering. 2009 Oct 16;22(10):1345-59.

Figure Legends

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

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

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

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

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

Figure 6 Top part of CBIT web calculator screen used to enter patient characteristics. The calculator is available at <https://hi.gmu.edu/cbit> [22].

Additional Files

Additional file 1.doc -Top ranked predictors of functional status

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