

1 **Full Title:** Social Determinants of Healthy Aging: An Investigation using the All of Us

2 Cohort

3 **Short title:** Social Determinants of Healthy Aging

4 Wei-Han Chen¹, MS, Yao-An Lee¹, MS, Huilin Tang¹, MSc, Chenyu Li², Ying Lu¹, Yu

5 Huang³, PhD, Rui Yin, PhD³, Melissa J. Armstrong, MD⁴, Yang Yang, MD⁵, Gregor

6 Štiglic, PhD^{6,7}, Jiang Bian, PhD³, Jingchuan Guo, MD, PhD¹

7 ¹Department of Pharmaceutical Outcomes and Policy, University of Florida, Gainesville,

8 Florida

9 ² Department of Biomedical Informatics, University of Pittsburgh, Pittsburgh,

10 Pennsylvania

11 ³Department of Health Outcomes and Bioinformatics, University of Florida, Gainesville,

12 Florida

13 ⁴Fixel Institute for Neurological Diseases, Department of Neurology, University of

14 Florida, Gainesville, Florida

15 ⁵Databricks Inc.

16 ⁶Faculty of Health Sciences, University of Maribor, Maribor, Slovenia

17 ⁷Usher Institute, University of Edinburgh, Edinburgh, UK

18

19

20 **Corresponding author:**

21 Jingchuan Guo, MD, PhD

22 Assistant Professor

23 Department of Pharmaceutical Outcomes & Policy

24 College of Pharmacy

25 1889 Museum Road, DSIT 6004, Gainesville, FL 32606

26 Phone +1-352-273-6533

27 E-mail: guoj1@ufl.edu

28

29 **Conflict of Interest :**

30 None.

31

32 **Funding Support :**

33 National Institute of Health / National Institute on Aging: R01AG089445

34

35

36

37 **Data sharing statement:** This dataset for this study is sourced from the All of Us
38 Research Program Registered tier dataset v7, which is not publicly available. The
39 authors have obtained all necessary permissions to access and utilize the All of Us
40 Research Program registered tier dataset for this study. Researchers who are interested
41 in accessing the All of Us Research Program dataset can apply for access through the
42 All of Us Research Hub. For more information and to initiate the application process,
43 please visit the All of Us Research Hub website at <https://www.researchallofus.org/>.

44

45 **ABSTRACT**

46 **Introduction**

47 The increasing aging population raises significant concerns about the ability of
48 individuals to age healthily, avoiding chronic diseases and maintaining cognitive and
49 physical functions. However, the pathways through which SDOH factors are associated
50 with healthy aging remain unclear.

51 **Methods**

52 This retrospective cohort study used the registered tier data from the *All of Us* Research
53 Program (AoURP) registered tier dataset v7. Eligible study participants are those aged
54 50 and older who have responded to any of the SDOH survey questions with available
55 EHR data. Three different algorithms were trained (logistic regression [LR], multi-layer
56 perceptron [MLP], and extreme gradient boosting [XGBoost]). The outcome is healthy
57 aging, which is measured by a composite score of the status for 1) comorbidities, 2)
58 cognitive conditions, and 3) mobility function. We evaluate the model performance by
59 area under the receiver operating characteristic curve (AUROC) and assess the fairness of
60 best-performed model through predictive parity. Feature importance is analyzed using
61 SHapley Additive exPlanations (SHAP) values.

62 **Results:**

63 Our study included 99,935 participants aged 50 and above, and the mean (SD) age was
64 74 (9.3), with 55,294 (55.3%) females, 67,457 (67.5%) Whites, 11,109 (11.1%) Hispanic
65 ethnicity, and 44,109 (44.1%) are classified as healthy aging. Most of the individuals
66 lived in their own house (64%), were married (51%), obtained college or advanced
67 degrees (74%), and had Medicare (56.2%). The best predictive model was XGBoost
68 with random oversampler, with a performance of AUROC [95% CI]: 0.793 [0.788-0.796],

69 F1 score: 0.697 [0.692-0.701], recall: 0.739 [0.732-0.748], precision: 0.659 [0.655-
70 0.663], and accuracy: 0.716 [0.712-0.720], and the XGBoost model achieved predictive
71 parity by similar positive and negative predictive values across race and sex groups
72 (0.86-1.06). In feature importance analysis, health insurance type is ranked as the most
73 predictive feature, followed by employment status, substance use, and health insurance
74 coverage (yes/no).

75 **Conclusion**

76 In this cohort study, XGBoost model accurately predicted individuals achieving healthy
77 aging, outperforming LR and MLP. Our findings underscore the significant role of health
78 insurance in contributing to healthy aging.

79

80 INTRODUCTION

81 Aging is an inevitable biological process characterized by a gradual decline in
82 physiological functions, leading to increased vulnerability to diseases and death. The
83 concept of healthy aging has emerged as a critical focus in geriatric research and public
84 health.(1–3) Healthy aging refers to the process of developing and maintaining
85 functional abilities that enable well-being in older age.(4–6) It encompasses physical,
86 mental, cognitive, and social well-being, allowing individuals to live independently and
87 enjoy a good quality of life despite the natural aging process. The increasing proportion
88 of older individuals in the global population has intensified the need to understand and
89 promote healthy aging, making it a vital area of study.(7) In 2020, nearly 1 in 6
90 Americans were 65 years or older, and this group is estimated to constitute 23% of the
91 total US population in 2050.(8,9)

92 Social determinants of health (SDOH) — the conditions where people are born,
93 grow, work, live, and age — play a crucial role in individuals' health, influencing the
94 aging process and the ability to age healthily.(10) SDOH includes factors such as
95 socioeconomic status, education, neighborhood and physical environment,
96 employment, social support networks, and access to healthcare.(11–13) Previous
97 studies have demonstrated that these SDOH can significantly affect an individual's
98 health outcomes by influencing behaviors, exposures, and access to resources
99 necessary for maintaining health.(14–18) Individuals with higher socioeconomic status,
100 better education, and stronger social support tend to have better health outcomes and a
101 higher likelihood of healthy aging.(19–21) Addressing disparities in SDOH is therefore

102 essential for promoting health equity and improving the quality of life for older adults,
103 especially those socioeconomically disadvantaged groups.(22–24)

104 Existing studies on the impact of SDOH on healthy aging are limited.(25) For
105 instance, Sowa et al. have identified a set of predictors using health surveys in Europe,
106 however, they focused only on lifestyle and psychosocial factors and did not consider
107 many other SDOH.(26) On the other hand, the application of machine learning (ML)
108 models has shown great promise in predicting health outcomes.(27) Other studies that
109 applied ML techniques have mainly focused on biological or physiological factors in
110 healthy aging,(28) none have studied SDOH.

111 To fill the gap, the objective of this study is to develop a prediction model of
112 healthy aging by leveraging a large cohort of older adults from the AoU and advanced
113 ML techniques. Understanding the relationship between SDOH and healthy aging holds
114 significant clinical and policy implications. Clinically, this knowledge enables healthcare
115 providers to create more personalized care plans that address both medical and social
116 factors influencing a patient's health. On the policy side, identifying key SDOH linked to
117 healthy aging can guide targeted interventions and resource allocation, fostering public
118 health strategies that promote healthy aging across diverse populations. Additionally, we
119 also evaluated the fairness of the ML models in predicting healthy aging, ensuring that
120 they do not perpetuate existing disparities and can be applied equitably across different
121 demographic groups. Lastly, we identified the top predictors for healthy aging using
122 SHapley Additive exPlanations (SHAP) values, a well-established explainable ML
123 method, which could inform the development of targeted interventions and policies to
124 support healthy aging by addressing the most influential SDOH.

125 **METHODS**

126 **Data Source and Study Population**

127 We used the registered tier data from the *All of Us* (AoU) Research Program
128 registered tier dataset v7.(29) The AoU was a nationwide program funded by the
129 National Institute of Health, which aimed to provide diverse and comprehensive
130 information among under-represented groups. The database included survey questions
131 (e.g., lifestyle, demographic, and social determinants of health) and electronic health
132 records (EHR).(30) Both survey questions and EHR were standardized and could be
133 mapped utilizing Observational Medical Outcomes Partnership (OMOP) Common Data
134 Model infrastructure.(31) We included individuals aged ≥ 50 years of age who have
135 responded to any of the SDOH survey questions with available EHR data.

136 **Study Outcome**

137 The primary outcome is a dichotomous score of healthy aging, which was
138 measured by a composite score of the status for 1) comorbidities, 2) cognitive
139 conditions, and 3) mobility function. Charlson comorbidity index (CCI) by Quan. et al(32)
140 was used for assessing comorbidity status. We modified the original CCI algorithm to
141 exclude age as a parameter (referred to as modified CCI [mCCI]) since our goal was to
142 predict healthy aging. Secondly, we assessed the cognitive conditions by ICD-9 and -10
143 CM codes with a diagnosis of mild cognitive impairment (MCI). Lastly, to assess the
144 mobility function, we identified individuals in assisted living using CPT/HCPSC codes
145 and records of discharge locations. An individual aged over 75 is classified as
146 experiencing healthy aging if they have a composite score of 0, which includes an mCCI

147 score of 0, no MCI, and are not in assisted living. A composite score greater than 0
148 indicates non-healthy aging.

149 We also defined a secondary cohort as a composite score of 0, with age greater
150 than 85 classified as healthy aging, otherwise as non-healthy aging. Two distinct
151 cohorts were then created for primary outcome and the secondary outcome analysis,
152 respectively. The secondary cohort analysis allows us to examine whether the
153 association between SDOH and healthy aging hold consistent when applying a more
154 stringent definition of healthy aging. Consistent results across both cohorts would
155 reinforce the robustness of our findings across varying definitions of healthy aging.

156 **Study Design**

157 We adopted a retrospective cohort study design and illustrated the cohort
158 selection process in **Figure 1**. Patients aged under 75 with an mCCI score of 0 are
159 excluded from the analysis in primary cohort. For the secondary cohort, this exclusion
160 extends to patients aged under 85 with an mCCI score of 0.

161 **Potential Risk Factors**

162 Potential risk factors (i.e., input features) were SDOH information collected from
163 multifaceted survey questions, including The Basics (demographic information),
164 Lifestyle (smoking, alcohol use, substance use, etc.), Healthcare Access & Utilization
165 (access to and use of health care resources), and Social Factors (neighborhood, social
166 life, stress, etc.). Self-reported race and gender were also recorded and included in the
167 analysis. We reported the counts (percentages) for categorical variables and median
168 (interquartile range) for continuous variables.

169 **Statistical analysis**

170 We aimed to use SDOH features to develop a machine learning model to predict
171 healthy aging. Three machine learning algorithms were applied: logistic regression (LR),
172 multi-layer perceptron (MLP), and Extreme Gradient Boosting(33) (XGBoost).
173 Regularization was employed in both logistic regression (lasso [L1](34), ridge [L2](35),
174 and ElasticNet(36)) and XGBoost (alpha [L1] and lambda [L2]) to reduce overfitting.
175 Following machine learning best practices, we split the data into training and testing
176 with a ratio of 8:2. To account for target class imbalance, we employed both random
177 over-sampling and random under-sampling methods and compared their performance
178 for further analyses.(37) For random over-sampling, we increased the minority class to
179 match the size of the majority class, resulting in a final balanced distribution of 50% for
180 each class. Similarly, for random under-sampling, we reduced the majority class to
181 match the size of the minority class, also achieving a balanced class distribution.

182 After hyperparameters tuning using Bayesian optimization with 5-fold cross-
183 validation over 100 iterations to optimize the area under the receiver operating
184 characteristic curve (AUROC), we reported the performance metrics of the testing set
185 including AUROC, precision, recall, F1 score, and specificity. In addition, we obtained
186 the 95% confidence intervals (CI) of the performance metrics by bootstrap method with
187 50 iterations. The best model is selected based on AUROC and the potential clinical
188 application with the goal of a higher F1 score, showing the balance between precision
189 and recall.

190 We then assessed the fairness of the machine learning model selected by
191 comparing the ratios of metrics such as positive predicted value (PPV), negative
192 predicted value (NPV), false positive rate (FPR), true positive rate (TPR), false negative

193 rate (FNR), and overall accuracy across race and gender. We designated non-Hispanic
194 Whites and females as the privileged groups for race and gender, respectively, and
195 identified Black and males as the protected groups for these categories. Lastly, we
196 adopted SHapley Additive exPlanations (SHAP) values to identify and rank the most
197 important features, with a view to providing explainability and improved clinical decision-
198 making.(38)

199 All analyses were performed in Python (version 3.10 with libraries such as Scikit-
200 learn, Imbalanced-learn). The study followed the STROBE cohort reporting
201 guideline.(39)

202 RESULTS

203 Descriptive Statistics

204 In the primary cohort, 99,936 eligible older adults aged 50 or older who had
205 responded to SDOH survey questions were included, and 44,109 (44%) were identified
206 as healthy aging (age ≥ 75 and a composite condition score of 0, **Table 1**). The mean
207 (SD) age was 74 (9.3) years, with 55,294 (55.3%) females, 41,977 (42.0%) males. Of
208 the cohort, 67,457 (67.5%) were White, 14,612 (14.6%) were Black or African American,
209 11,109 (11.1%) having Hispanic ethnicity. The median (IQR) of the mCCI was 1 (0-2).
210 Most of the individuals lived in their own house (64%), were married (51%), obtained
211 college or advanced degrees (74%), and had Medicare (56.2%).

212 In the secondary cohort, 62,475 participants were included, and 6,648 (10.6%)
213 were identified as healthy aging (i.e., age ≥ 85 and a composite condition score of 0,
214 **Table 1**). The mean (SD) age is 71 (10.6) years, with 36,101 (58%) females, 24,671
215 (40%) males. 38,802 (62%) were White, 11,437 (18%) were Black or African American,
216 8,270 (13%) having Hispanic ethnicity. The median (IQR) of the mCCI was 2 (1-3).
217 Similarly, most of the individuals lived in their own house (57%), were married (46.9%),
218 obtained college or advanced degrees (70%), and had Medicare (46%).

219 Model Performance and Selection

220 Performance metrics on the test dataset and the AUROC for the three models
221 are presented in **Figure 2 and Supplemental table S1**. Bootstrapped performance with
222 95% CI over 50 iterations for the best algorithm are included in **Supplemental table S2**.
223 Overall, all three models achieved decent prediction performance with AUROC > 0.7 .
224 Among them, we found that the XGBoost model with over-sampling adjustments

225 (AUROC: 0.795 and 0.862 for primary and secondary cohort, respectively) shows
226 superior performance. This outperformed both the LR model (AUROC: 0.786 and 0.85)
227 and the MLP model (AUROC: 0.794 and 0.854). Though the AUROC was comparable
228 between XGBoost and MLP, the F1 score and other metrics of MLP classifier are much
229 lower than those of XGBoost classifier. Detailed machine learning algorithm and tuning
230 values are included in **Supplemental table S3**.

231 **Model Fairness Assessment**

232 **Table 2** demonstrates the fairness metrics across gender and race for the best
233 selected model, XGBoost. The ratios of accuracies and PPVs showed no evidence of
234 model biases towards a specific population.

235 **Feature Importance Analysis**

236 **Figure 3** shows the SHAP values to explain the healthy aging prediction of
237 XGBoost model (best performance). In both cohorts, health insurance type (e.g.,
238 Medicare, Medicaid, insurance purchased from a company) is ranked as the most
239 predictive feature (SHAP value: 0.595), followed by employment status (0.233),
240 substance use (0.171), health insurance coverage (yes/no, SHAP value 0.143). The
241 direction of the plots revealed that all top 10 features were positively (red on the right in
242 figure 3) associated with healthy aging.

243 **DISCUSSION**

244 This cohort study leverages the AoU datasets, which not only included diverse
245 populations from historically underrepresented groups and racial/ethnic minority groups,
246 but also provided a rich source of SDOH through standardized OMOP data
247 infrastructure. Our findings suggest that machine learning models could accurately

248 predict healthy aging using SDOH information and highlight the potential for integrating
249 SDOH factors in clinical decision-making to enhance predictive accuracy.

250 It is noteworthy that our work included high-dimensional SDOH with large-scale
251 population and an explainable ML framework. A few of previous studies were available
252 to include SDOH across several domains, such as neighborhood environment,
253 education access, etc..(40,41) This study also added values by the integration of
254 objective measures from EHRs with detailed survey data on SDOH, providing a more
255 comprehensive assessment of healthy aging compared to studies relying solely on self-
256 reported data.(26)

257 Our models showed fairness in predictive parity, where the ratios of both positive
258 (PPV) and negative predictive values (NPV) are close to 1. Ensuring that ML models do
259 not discriminate against different racial or ethnic groups is crucial, as these models must
260 perform equitably independent on sensitive features. In our context, if healthy aging is
261 less accurately identified in disadvantaged groups, it may lead to unnecessary and
262 potentially harmful treatments, thereby increasing their financial burden and causing
263 undue harm. Therefore, maintaining equal PPV and NPV across different demographic
264 groups is imperative to prevent such disparities and ensure equitable and healthcare
265 outcomes.(42)

266 Some machine learning classifiers are notorious as a “black box” where excellent
267 performance is often obtained at the cost of lacking interpretability.(43–47) In the feature
268 importance analysis, health insurance was the strongest positive SDOH factor for
269 predicting healthy aging. Our study identified top SDOH factors from several domains
270 positively associated with healthy aging: health insurance type, employment status,

271 education level, marital status, housing status. These aligned with previous studies
272 indicating that higher socioeconomic status, including higher income and education
273 level, was associated with better health outcomes.(28,48) Although the precise
274 mechanism of marital and housing statuses on healthy aging has not yet been
275 identified, studies have shown that there was an intricate pattern associated with mental
276 health and chronic diseases.(49,50) Our study added the evidence that they could also
277 be related to healthy aging.

278 Past research has suggested that substance use and related drug overdoses
279 may have contributed to lower life expectancy,(51–53) however, we found that
280 substance and alcohol use were positively associated with healthy aging. We suggest
281 the following probable underlying causes of the results. First, it is possible that
282 substance and alcohol users in our cohort be healthier than those in the general
283 population from previous studies since our participants mainly participated AoU
284 voluntarily rather than being randomly recruited into the program. Secondly, some
285 substance and alcohol users may have altered their health behaviors following
286 enrollment in the program, namely, Hawthorne effect, where individuals knowingly adopt
287 a healthier behavior when they were being assessed in a research program.(54,55)
288 Thus, there may be difference in substance and alcohol use status between the
289 baseline period and the follow-up period. Lastly, according to statistics from Centers for
290 Disease Control and Prevention, the drug overdose death rates were higher among
291 groups aged 25-44 (~50 deaths per 100,000 population) compared to those aged over
292 55 (5-35 death per 100,000 population).(56) Thus, our results could be affected by
293 selection bias since we only included participants aged over 50 in our study, whose

294 health behaviors are general substance and alcohol use populations. These factors
295 could contribute to the unexpected positive association between substance use and
296 healthy aging in our study.

297 This study has some limitations. First, to date, a unanimous definition of healthy
298 aging has not yet been reached.(57) Our definition of healthy aging, while based on
299 objective measures, may not capture all constructs of 'healthy aging' such as quality of
300 life, social engagement, and subjective well-being. Also, CCI limited the spectrum of
301 comorbidities. For instance, some may not consider Parkinson's disease as healthy
302 aging, however, it is not covered in CCI.

303 Secondly, while the All of Us Research Program provides a large and diverse
304 dataset, it may not be fully representative of the U.S. population. Participants in All of Us
305 are volunteers who agreed to share their health data, which could introduce selection
306 bias. These individuals may be more health-conscious, have better access to
307 healthcare, or obtained higher educational degrees than the general population,
308 potentially leading to an overestimation of healthy aging in our sample. While the
309 generalizability of our findings is limited to the participants in AoU, it is important to note
310 that the cohort is wide across the nation. Furthermore, the representation of racial and
311 ethnic minority groups has improved in AoU, which enhances the applicability of our
312 results to a more diverse population.(58) However, caution should still be exercised
313 when extrapolating these findings to other populations or other clinical settings.

314 Thirdly, while we attempted to account for a wide range of SDOH factors, there
315 may still be unmeasured confounders. For instance, we did not have data on lifelong
316 health behaviors or early-life exposures that could significantly impact aging trajectories.

317

318 **CONCLUSION**

319 In this cohort study utilizing the AoU database, our machine learning model
320 effectively predicted individuals likely to achieve healthy aging, emphasizing the critical
321 influence of health insurance on this outcome. The findings highlight that access to
322 health insurance is not merely a facilitator of healthcare services but a pivotal
323 determinant of long-term health outcomes in older adults. By addressing the gaps in
324 health insurance, policymakers can contribute to the promotion of healthy aging across
325 diverse populations, ultimately leading to improved quality of life. The integration of
326 health insurance into public health strategies could therefore be a powerful tool in
327 enhancing the overall well-being of aging populations.

328

329 Reference

- 330 1. Menassa M, Stronks K, Khatami F, Díaz ZMR, Espinola OP, Gamba M, et al. Concepts and
331 definitions of healthy ageing: a systematic review and synthesis of theoretical models.
332 eClinicalMedicine [Internet]. 2023 Feb 1 [cited 2024 Jul 6];56. Available from:
333 [https://www.thelancet.com/journals/eclinm/article/PIIS2589-5370\(22\)00550-8/fulltext#](https://www.thelancet.com/journals/eclinm/article/PIIS2589-5370(22)00550-8/fulltext#)
- 334 2. Rudnicka E, Napierała P, Podfigurna A, Męczekalski B, Smolarczyk R, Grymowicz M. The
335 World Health Organization (WHO) approach to healthy ageing. *Maturitas*. 2020 Sep;139:6–
336 11.
- 337 3. Michel JP, Leonardi M, Martin M, Prina M. WHO’s report for the decade of healthy ageing
338 2021–30 sets the stage for globally comparable data on healthy ageing. *The Lancet Healthy*
339 *Longevity*. 2021 Mar 1;2(3):e121–2.
- 340 4. Healthy Aging - PAHO/WHO | Pan American Health Organization [Internet]. 2024 [cited
341 2024 Jul 7]. Available from: <https://www.paho.org/en/healthy-aging>
- 342 5. Beard JR, Officer A, Carvalho IA de, Sadana R, Pot AM, Michel JP, et al. The World report on
343 ageing and health: a policy framework for healthy ageing. *The Lancet*. 2016 May
344 21;387(10033):2145–54.
- 345 6. National Institute on Aging [Internet]. [cited 2024 Jul 7]. Healthy aging. Available from:
346 <https://www.nia.nih.gov/health/healthy-aging>
- 347 7. Jackson EMJ, O’Brien K, McGuire LC, Baumgart M, Gore J, Brandt K, et al. Promoting Healthy
348 Aging: Public Health as a Leader for Reducing Dementia Risk. *Public Policy & Aging Report*.
349 2023 Sep 1;33(3):92–5.
- 350 8. Bureau UC. Census.gov. [cited 2024 Jul 7]. Older Population and Aging. Available from:
351 <https://www.census.gov/topics/population/older-aging.html>
- 352 9. Vespa J, Medina L, Armstrong DM. Population Estimates and Projections.
- 353 10. Noren Hooten N, Pacheco NL, Smith JT, Evans MK. The accelerated aging phenotype: The
354 role of race and social determinants of health on aging. *Ageing Research Reviews*. 2022 Jan
355 1;73:101536.
- 356 11. CDC. About CDC. 2024 [cited 2024 Jul 7]. Social Determinants of Health (SDOH). Available
357 from: <https://www.cdc.gov/about/priorities/why-is-addressing-sdoh-important.html>
- 358 12. Social Determinants of Health - Healthy People 2030 | health.gov [Internet]. Available from:
359 <https://health.gov/healthypeople/priority-areas/social-determinants-health>

- 360 13. Social determinants of health: Key concepts [Internet]. [cited 2024 Jul 7]. Available from:
361 [https://www.who.int/news-room/questions-and-answers/item/social-determinants-of-](https://www.who.int/news-room/questions-and-answers/item/social-determinants-of-health-key-concepts)
362 [health-key-concepts](https://www.who.int/news-room/questions-and-answers/item/social-determinants-of-health-key-concepts)
- 363 14. Rangachari P, Govindarajan A, Mehta R, Seehusen D, Rethemeyer RK. The relationship
364 between Social Determinants of Health (SDoH) and death from cardiovascular disease or
365 opioid use in counties across the United States (2009–2018). *BMC Public Health*. 2022 Feb
366 4;22(1):236.
- 367 15. Tran R, Forman R, Mossialos E, Nasir K, Kulkarni A. Social Determinants of Disparities in
368 Mortality Outcomes in Congenital Heart Disease: A Systematic Review and Meta-Analysis.
369 *Front Cardiovasc Med* [Internet]. 2022 Mar 15 [cited 2024 Jul 7];9. Available from:
370 [https://www.frontiersin.org/journals/cardiovascular-](https://www.frontiersin.org/journals/cardiovascular-medicine/articles/10.3389/fcvm.2022.829902/full)
371 [medicine/articles/10.3389/fcvm.2022.829902/full](https://www.frontiersin.org/journals/cardiovascular-medicine/articles/10.3389/fcvm.2022.829902/full)
- 372 16. Short SE, Mollborn S. Social Determinants and Health Behaviors: Conceptual Frames and
373 Empirical Advances. *Curr Opin Psychol*. 2015 Oct;5:78–84.
- 374 17. Alcántara C, Diaz SV, Cosenzo LG, Loucks EB, Penedo FJ, Williams NJ. Social determinants as
375 moderators of the effectiveness of health behavior change interventions: scientific gaps and
376 opportunities. *Health Psychology Review*. 2020 Jan 2;14(1):132–44.
- 377 18. Ayangunna E, Kalu K, Shah G. Role of Community-level Health Behaviors and Social
378 Determinants of Health in Preventable Hospitalizations. *Journal of the Georgia Public Health*
379 *Association*. 2022 Jan 1;8(3):93–101.
- 380 19. Bundy JD, Mills KT, He H, LaVeist TA, Ferdinand KC, Chen J, et al. Social determinants of
381 health and premature death among adults in the USA from 1999 to 2018: a national cohort
382 study. *The Lancet Public Health*. 2023 Jun 1;8(6):e422–31.
- 383 20. Monroe P, Campbell JA, Harris M, Egede LE. Racial/ethnic differences in social determinants
384 of health and health outcomes among adolescents and youth ages 10–24 years old: a
385 scoping review. *BMC Public Health*. 2023 Mar 1;23(1):410.
- 386 21. Adkins-Jackson PB, George KM, Besser LM, Hyun J, Lamar M, Hill-Jarrett TG, et al. The
387 structural and social determinants of Alzheimer’s disease related dementias. *Alzheimer’s &*
388 *Dementia*. 2023;19(7):3171–85.
- 389 22. Perez FP, Perez CA, Chumbiauca MN. Insights into the Social Determinants of Health in Older
390 Adults. *J Biomed Sci Eng*. 2022 Nov;15(11):261–8.
- 391 23. Llorens-Ortega R, Bertran-Noguer C, Juvinyà-Canals D, Garre-Olmo J, Bosch-Farré C.
392 Influence of social determinants of health in the evolution of the quality of life of older
393 adults in Europe: A comparative analysis between men and women. *Humanit Soc Sci*
394 *Commun*. 2024 Mar 13;11(1):1–13.

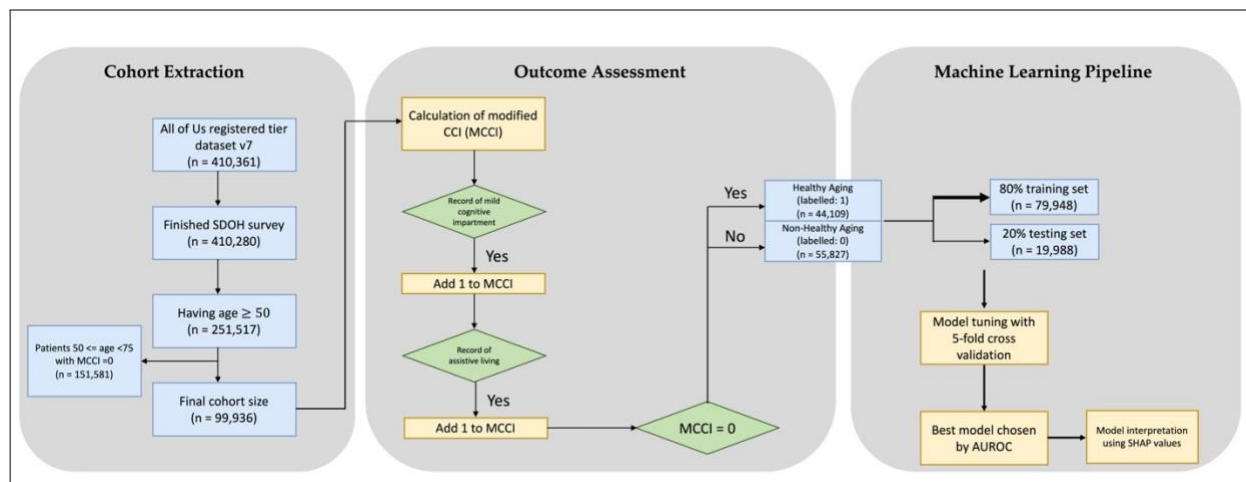
- 395 24. Yearby R. The Social Determinants of Health, Health Disparities, and Health Justice. *J Law*
396 *Med Ethics*. 50(4):641–9.
- 397 25. Abud T, Kounidas G, Martin KR, Werth M, Cooper K, Myint PK. Determinants of healthy
398 ageing: a systematic review of contemporary literature. *Aging Clin Exp Res*. 2022
399 Jun;34(6):1215–23.
- 400 26. Sowa A, Tobiasz-Adamczyk B, Topór-Mądry R, Poscia A, Ia Milia DI. Predictors of healthy
401 ageing: public health policy targets. *BMC Health Serv Res*. 2016 Sep 5;16(5):289.
- 402 27. Wong J, Murray Horwitz M, Zhou L, Toh S. Using Machine Learning to Identify Health
403 Outcomes from Electronic Health Record Data. *Curr Epidemiol Rep*. 2018 Dec 1;5(4):331–42.
- 404 28. Wagg E, Blyth FM, Cumming RG, Khalatbari-Soltani S. Socioeconomic position and healthy
405 ageing: A systematic review of cross-sectional and longitudinal studies. *Ageing Research*
406 *Reviews*. 2021 Aug 1;69:101365.
- 407 29. The “All of Us” Research Program. *New England Journal of Medicine*. 2019 Aug
408 15;381(7):668–76.
- 409 30. Tesfaye S, Cronin RM, Lopez-Class M, Chen Q, Foster CS, Gu CA, et al. Measuring social
410 determinants of health in the All of Us Research Program. *Sci Rep*. 2024 Apr 16;14(1):8815.
- 411 31. Data Standardization – OHDSI [Internet]. [cited 2024 Feb 29]. Available from:
412 <https://www.ohdsi.org/data-standardization/>
- 413 32. Quan H, Li B, Couris CM, Fushimi K, Graham P, Hider P, et al. Updating and validating the
414 Charlson comorbidity index and score for risk adjustment in hospital discharge abstracts
415 using data from 6 countries. *Am J Epidemiol*. 2011 Mar 15;173(6):676–82.
- 416 33. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: *Proceedings of the 22nd*
417 *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* [Internet].
418 2016 [cited 2024 Feb 29]. p. 785–94. Available from: <http://arxiv.org/abs/1603.02754>
- 419 34. Tibshirani R. Regression Shrinkage and Selection Via the Lasso. *Journal of the Royal*
420 *Statistical Society: Series B (Methodological)*. 1996;58(1):267–88.
- 421 35. Hoerl AE, Kennard RW. Ridge Regression: Biased Estimation for Nonorthogonal Problems.
422 *Technometrics*. 1970 Feb 1;12(1):55–67.
- 423 36. Regularization and Variable Selection Via the Elastic Net | *Journal of the Royal Statistical*
424 *Society Series B: Statistical Methodology* | Oxford Academic [Internet]. [cited 2024 Feb 29].
425 Available from: <https://academic.oup.com/jrsssb/article/67/2/301/7109482>
- 426 37. Lemaître G, Nogueira F, Aridas C. Imbalanced-learn: A Python Toolbox to Tackle the Curse of
427 Imbalanced Datasets in Machine Learning. 2016 Sep 21;18.

- 428 38. Lundberg SM, Lee SI. A unified approach to interpreting model predictions. In: Proceedings
429 of the 31st International Conference on Neural Information Processing Systems. Red Hook,
430 NY, USA: Curran Associates Inc.; 2017. p. 4768–77. (NIPS'17).
- 431 39. von Elm E, Altman DG, Egger M, Pocock SJ, Gøtzsche PC, Vandenbroucke JP. Strengthening
432 the reporting of observational studies in epidemiology (STROBE) statement: guidelines for
433 reporting observational studies. *BMJ*. 2007 Oct 20;335(7624):806–8.
- 434 40. de Keijzer C, Bauwelinck M, Dadvand P. Long-Term Exposure to Residential Greenspace and
435 Healthy Ageing: a Systematic Review. *Curr Envir Health Rpt*. 2020 Mar 1;7(1):65–88.
- 436 41. Chen M, Tan X, Padman R. Social determinants of health in electronic health records and
437 their impact on analysis and risk prediction: A systematic review. *Journal of the American*
438 *Medical Informatics Association*. 2020 Nov 1;27(11):1764–73.
- 439 42. Rajkomar A, Hardt M, Howell MD, Corrado G, Chin MH. Ensuring Fairness in Machine
440 Learning to Advance Health Equity. *Ann Intern Med*. 2018 Dec 18;169(12):866–72.
- 441 43. Stiglic G, Kocbek P, Fijacko N, Zitnik M, Verbert K, Cilar L. Interpretability of machine
442 learning-based prediction models in healthcare. *WIREs Data Mining and Knowledge*
443 *Discovery*. 2020;10(5):e1379.
- 444 44. Zihni E, Madai VI, Livne M, Galinovic I, Khalil AA, Fiebach JB, et al. Opening the black box of
445 artificial intelligence for clinical decision support: A study predicting stroke outcome. *PLOS*
446 *ONE*. 2020 Apr 6;15(4):e0231166.
- 447 45. Sajid MR, Khan AA, Albar HM, Muhammad N, Sami W, Bukhari SAC, et al. Exploration of
448 Black Boxes of Supervised Machine Learning Models: A Demonstration on Development of
449 Predictive Heart Risk Score. *Computational Intelligence and Neuroscience*.
450 2022;2022(1):5475313.
- 451 46. Azodi CB, Tang J, Shiu SH. Opening the Black Box: Interpretable Machine Learning for
452 Geneticists. *Trends in Genetics*. 2020 Jun 1;36(6):442–55.
- 453 47. Ratti E, Graves M. Explainable machine learning practices: opening another black box for
454 reliable medical AI. *AI Ethics*. 2022 Nov 1;2(4):801–14.
- 455 48. Wu YT, Daskalopoulou C, Terrera GM, Niubo AS, Rodríguez-Artalejo F, Ayuso-Mateos JL, et al.
456 Education and wealth inequalities in healthy ageing in eight harmonised cohorts in the
457 ATHLOS consortium: a population-based study. *The Lancet Public Health*. 2020 Jul
458 1;5(7):e386–94.
- 459 49. Swope CB, Hernández D. Housing as a determinant of health equity: A conceptual model.
460 *Social Science & Medicine*. 2019 Dec 1;243:112571.

- 461 50. Yannakoulia M, Panagiotakos D, Pitsavos C, Skoumas Y, Stafanadis C. Eating patterns may
462 mediate the association between marital status, body mass index, and blood cholesterol
463 levels in apparently healthy men and women from the ATTICA study. *Social Science &*
464 *Medicine*. 2008 Jun 1;66(11):2230–9.
- 465 51. Rehm J, Probst C. Decreases of Life Expectancy Despite Decreases in Non-Communicable
466 Disease Mortality: The Role of Substance Use and Socioeconomic Status. *European*
467 *Addiction Research*. 2018 Apr 6;24(2):53–9.
- 468 52. Imtiaz S, Probst C, Rehm J. Substance use and population life expectancy in the USA:
469 Interactions with health inequalities and implications for policy. *Drug and Alcohol Review*.
470 2018;37(S1):S263–7.
- 471 53. Gold MS. The Role of Alcohol, Drugs, and Deaths of Despair in the U.S.'s Falling Life
472 Expectancy. *Mo Med*. 2020;117(2):99–101.
- 473 54. Clifford PR, Davis CM, Maisto SA, Stout RL. Alcohol Treatment Research Contributing to
474 Changes in Substance Use Behavior and Related Negative Consequences. *J Stud Alcohol*
475 *Drugs*. 2022 May;83(3):364–73.
- 476 55. Berkhout C, Berbra O, Favre J, Collins C, Calafiore M, Peremans L, et al. Defining and
477 evaluating the Hawthorne effect in primary care, a systematic review and meta-analysis.
478 *Front Med [Internet]*. 2022 Nov 8 [cited 2024 Jul 7];9. Available from:
479 <https://www.frontiersin.org/journals/medicine/articles/10.3389/fmed.2022.1033486/full>
- 480 56. Drug overdose deaths - Health, United States [Internet]. 2023 [cited 2024 Jul 7]. Available
481 from: <https://www.cdc.gov/nchs/hus/topics/drug-overdose-deaths.htm>
- 482 57. Lu W, Pikhart H, Sacker A. Domains and Measurements of Healthy Aging in Epidemiological
483 Studies: A Review. *The Gerontologist*. 2019 Jul 16;59(4):e294–310.
- 484 58. Kathiresan N, Cho SMJ, Bhattacharya R, Truong B, Hornsby W, Natarajan P. Representation of
485 Race and Ethnicity in the Contemporary US Health Cohort All of Us Research Program. *JAMA*
486 *Cardiology*. 2023 Sep 1;8(9):859–64.
- 487
- 488

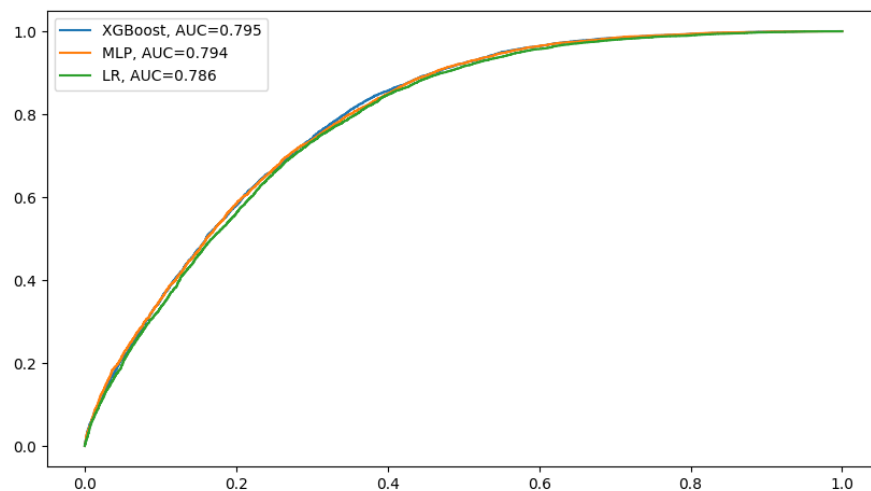
489 **Figures and Tables**

490 **Figure 1.** The overall workflow including participant selection, outcome assessment,
491 and machine learning pipeline.

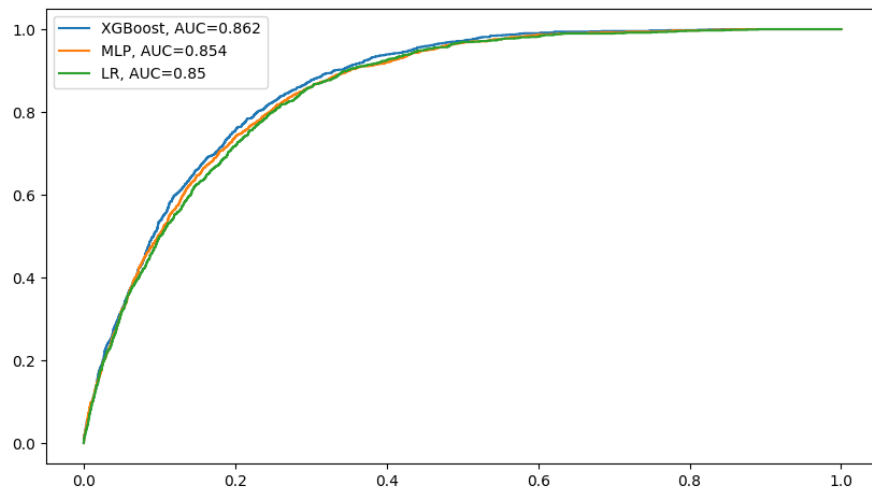


492
493 **Figure 2. Comparison of model performance on test datasets with area under the**
494 **receiver operating characteristic curve**

495 A. Performance on test datasets of the three algorithms in the primary cohort



496
497 B. Performance on test datasets of the three algorithms in the secondary cohort



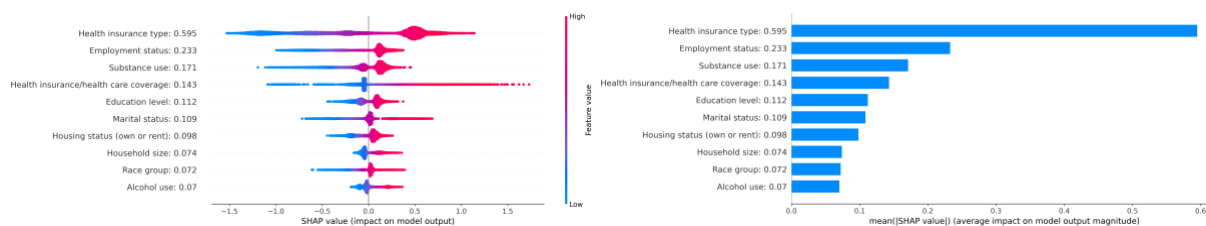
498

499 *XGBoost: extreme gradient boosting, LR: logistic regression, MLP: multilayer

500 perceptron

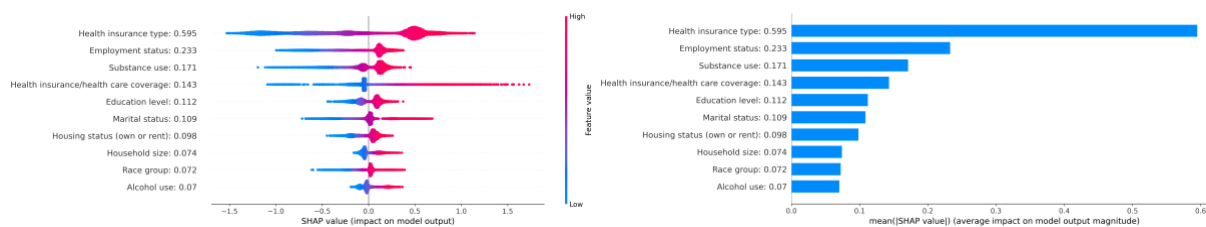
501 **Figure 3.** Distribution of the SHAP values for the top 10 features based on the highest
 502 mean absolute SHAP value (left panels) and their mean absolute contribution of the top
 503 10 features, ranked by their average SHAP value (right panels). Each test sample is
 504 depicted as a point for every feature, with the x-axis indicating whether the feature's
 505 effect on the model's prediction is positive (red on the right) or negative (blue on the
 506 right). The color of each point reflects the feature's value, and this color scale is
 507 adjusted individually according to the value range present in the dataset.

508 A. SHAP values and feature importance for the primary cohort using XGBoost



509

510 B. SHAP values and feature importance for the secondary cohort using XGBoost



511

512

513

514

515

516

517 **Table 1.** Descriptive analysis of participants demographic information

Variables	Primary Cohort, No. (%)			P	Secondary Cohort, No. (%)			P
	Overall, n = 99,936	Healthy aging (75), n = 44,109 (44.1%)	Non-healthy aging, n = 55,827 (55.9%)		Overall, n = 62,475	Healthy aging (85), n = 6,648 (10.6%)	Non-healthy aging, n = 55,827 (89.4%)	
Age, mean (SD)	74.0 (9.3)	79.9 (4.0)	69.4 (9.6)	<.001	71.2 (10.6)	87.2 (1.8)	69.4 (9.6)	<.001
Sex				<.001				<.001
Male	41,977 (42.0%)	20,493 (46.5%)	21,484 (38.5%)		24,671 (39.5%)	3,187 (47.9%)	32,839 (58.8%)	
Female	55,294 (55.3%)	22,455 (50.9%)	32,839 (58.8%)		36,101 (57.8%)	3,262 (49.1%)	21,484 (38.5%)	
Other	2,665 (2.7%)	1,161 (2.6%)	1,504 (2.7%)		1,703 (2.7%)	199 (3.0%)	1,504 (2.7%)	
Race				<.001				<.001
White	67,457 (67.5%)	33,916 (76.9%)	33,541 (60.1%)		38,802 (62.1%)	5,261 (79.1%)	33,541 (60.1%)	
Black or African American	14,612 (14.6%)	3,597 (8.2%)	11,015 (19.7%)		11,437 (18.3%)	422 (6.4%)	11,015 (19.7%)	
Asian	1,775 (1.8%)	941 (2.1%)	834 (1.5%)		954 (1.5%)	120 (1.8%)	834 (1.5%)	
Other/Unknown	16,092 (16.1%)	5,655 (12.8%)	10,437 (18.7%)		11,282 (18.1%)	845 (12.7%)	10,437 (18.7%)	
Modified Charlson comorbidity index, median (IQR)	1 (0-2)	0 (0-0)	2 (1-3)	<.001	2 (1-3)	0 (0-0)	2 (1-3)	<.001
Ethnicity				<.001				<.001
Not Hispanic or Latino	84,332 (84.4%)	38,693 (87.7%)	45,639 (81.8%)		51,477 (82.4%)	5,838 (87.8%)	45,639 (81.8%)	
Hispanic or Latino	11,109 (11.1%)	3,295 (7.5%)	7,814 (14.0%)		8,270 (13.2%)	456 (6.9%)	7,814 (14.0%)	
Other/Unknown	4,495 (4.5%)	2,121 (4.8%)	2,374 (4.2%)		2,728 (4.4%)	354 (5.3%)	2,374 (4.2%)	
Housing Status				<.001				<.001
Own	63,907 (63.9%)	32,924 (74.6%)	30983 (55.5%)		35,729 (57.2%)	4,746 (71.4%)	30,983 (55.5%)	
Rent	27,278 (27.3%)	8,018 (18.2%)	19,260 (34.5%)		20,547 (32.9%)	1,287 (19.4%)	19,260 (34.5%)	
Other/Unknown	8,751 (8.8%)	3,167 (7.2%)	5,584 (10.0%)		6,199 (9.9%)	615 (9.2%)	5,584 (10.0%)	
Current Marital Status				<.001				<.001
Married	51,179 (51.2%)	25,234 (57.2%)	25,945 (46.5%)		29,294 (46.9%)	3,349 (50.4%)	25,945 (46.5%)	
Divorced	17,556 (17.6%)	6,900 (15.6%)	10,656 (19.1%)		11,503 (18.4%)	847 (12.7%)	10,656 (19.1%)	
Widowed	11,711 (11.7%)	6,626 (15.0%)	5,085 (9.1%)		6,864 (11.0%)	1,779 (26.8%)	5,085 (9.1%)	
Never married	11,226 (11.2%)	2,690 (6.1%)	8,536 (15.3%)		8,824 (14.1%)	288 (4.3%)	8,536 (15.3%)	
Other/Unknown	8,264 (8.3%)	2,659 (6.0%)	5,605 (10.0%)		5,990 (9.6%)	385 (5.8%)	5,605 (10.0%)	
Education Level				<.001				<.001
College, graduate or advanced degree	74,320 (74.4%)	35,657 (80.8%)	38663 (69.3%)		43,773 (70.0%)	5,110 (76.9%)	9,948 (17.8%)	
High school or GED	14,895 (14.9%)	4,947 (11.2%)	9,948 (17.8%)		10,864 (17.4%)	916 (13.8%)	38,663 (69.3%)	
Less than high school	7,403 (7.4%)	2,194 (5.0%)	5,209 (9.3%)		5,610 (9.0%)	401 (6.0%)	5,209 (9.3%)	

Other/Unknown	3,318 (3.3%)	1,134 (3.0%)	2,007 (3.6%)		2,228 (3.6%)	221 (3.3%)	2,007 (3.6%)	
Health Insurance*								*
Medicare	56,150 (56.2%)	32,399 (73.5%)	23,751 (42.5%)		28,525 (45.7%)	4,774 (71.8%)	23,751 (42.5%)	
Medicaid	14,553 (14.6%)	2,886 (6.5%)	11,667 (20.9%)		12,082 (19.3%)	415 (6.2%)	11,667 (20.9%)	
Private/Other	48,914 (49%)	22,482 (51%)	26,432 (47.3%)		29,936 (48%)	3,504 (52.7%)	26,432 (47.3%)	
None	143 (0.1%)	32 (0.1%)	111 (0.2%)		115 (0.2%)	<20	111 (0.2%)	
Unknown	1,073 (1.1%)	379 (0.9%)	694 (1.2%)		768 (1.2%)	74 (1.1%)	694 (1.2%)	
Substance use*								*
Marijuana	40,084 (40.1%)	15,215 (34.5%)	24,869 (44.5%)		25,970 (41.6%)	1,101 (16.6%)	24,869 (44.5%)	
Cocaine	13,417 (13.4%)	3,383 (7.7%)	10,034 (18%)		10,181 (16.3%)	147 (2.2%)	10,034 (18.0%)	
Opioids	8,414 (8.4%)	2,380 (5.4%)	6,034 (10.8%)		6,257 (10%)	223 (3.3%)	6,034 (10.8%)	
Other	29,268 (29.3%)	10,306 (23.4%)	18,962 (34%)		19,736 (31.6%)	774 (11.6%)	18,962 (34%)	
None	45,982 (46%)	22,171 (50.3%)	23,811 (42.7%)		28,180 (45.1%)	4,369 (65.7%)	23,811 (42.7%)	
Unknown/skipped	7,799 (7.8%)	4,180 (9.5%)	3,619 (6.5%)		4,355 (7%)	736 (11.1%)	3,619 (6.5%)	

518 *Allow having multiple responses per participant

519 **Table 2. Fairness metrics for XGBoost across gender and race**

Primary Cohort						
	Black	White	Parity	Male	Female	Parity
PPV	0.58	0.67	0.87	0.67	0.66	1.03
TPR	0.34	0.81	0.42	0.77	0.71	1.09
NPV	0.80	0.76	1.06	0.74	0.79	0.94
Accuracy	0.77	0.70	1.10	0.70	0.73	0.96
Secondary Cohort						
	Black	White	Parity	Male	Female	Parity
PPV	0.57	0.67	0.86	0.67	0.66	1.02
TPR	0.33	0.81	0.41	0.77	0.71	1.10
NPV	0.80	0.76	1.06	0.74	0.79	0.94
Accuracy	0.77	0.7	1.10	0.7	0.73	0.96

520 *PPV: positive predicted value, FPR: false positive rate, TPR: true positive rate, FNR:

521 false negative rate, NPV: negative predicted value

522

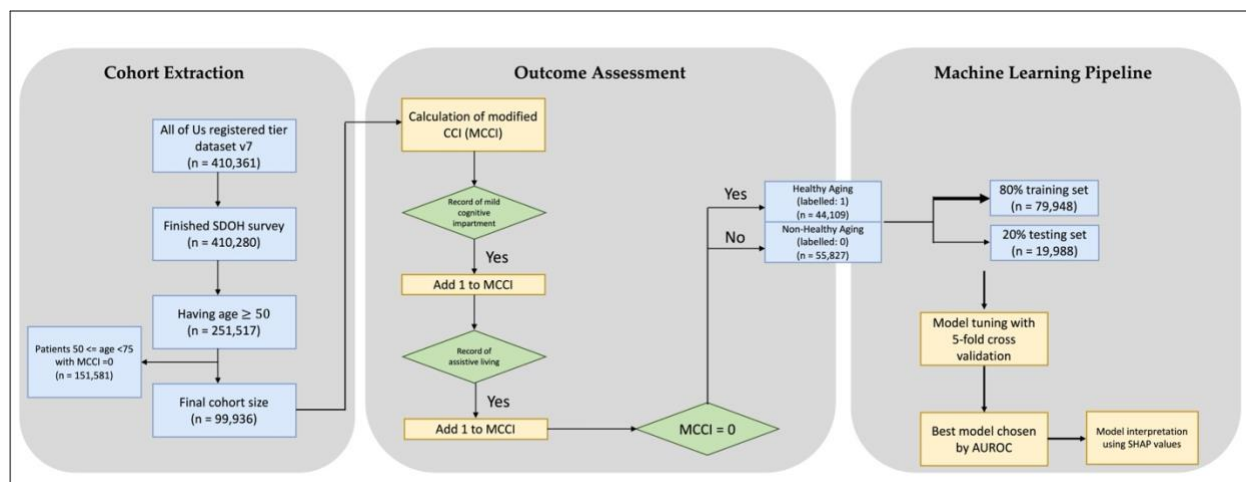
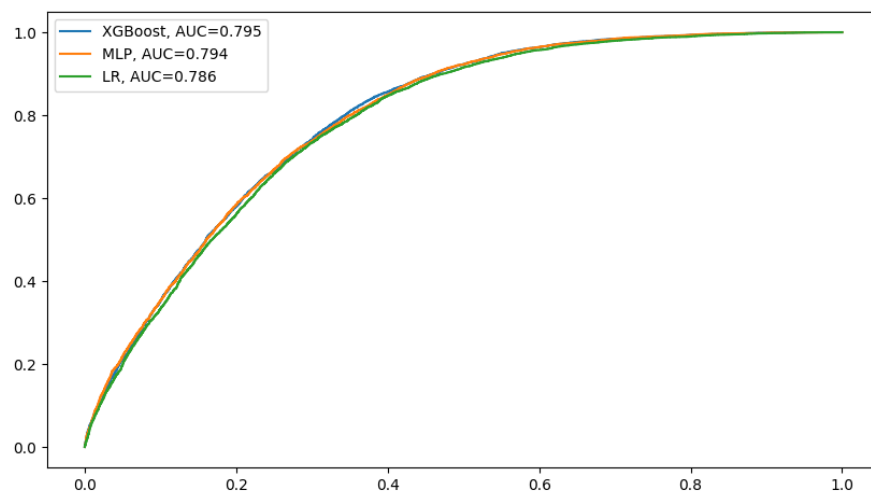
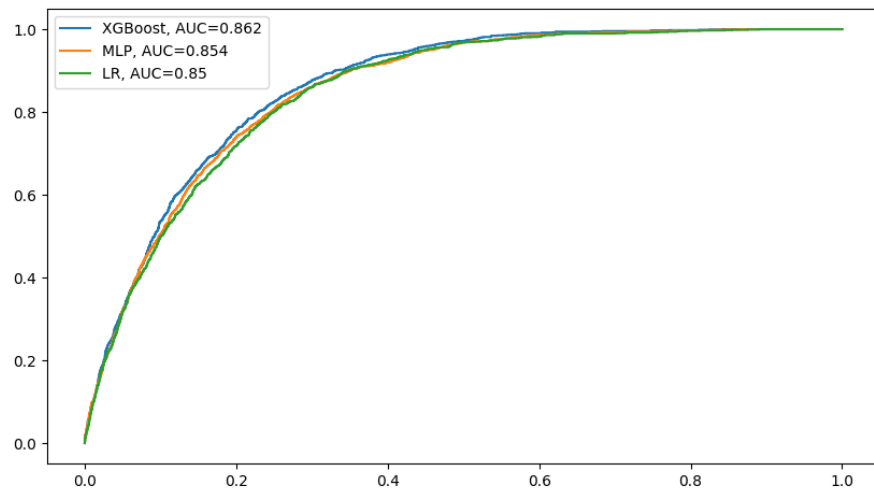


Figure 1. The overall workflow including participant selection, outcome assessment, and machine learning pipeline.

A.



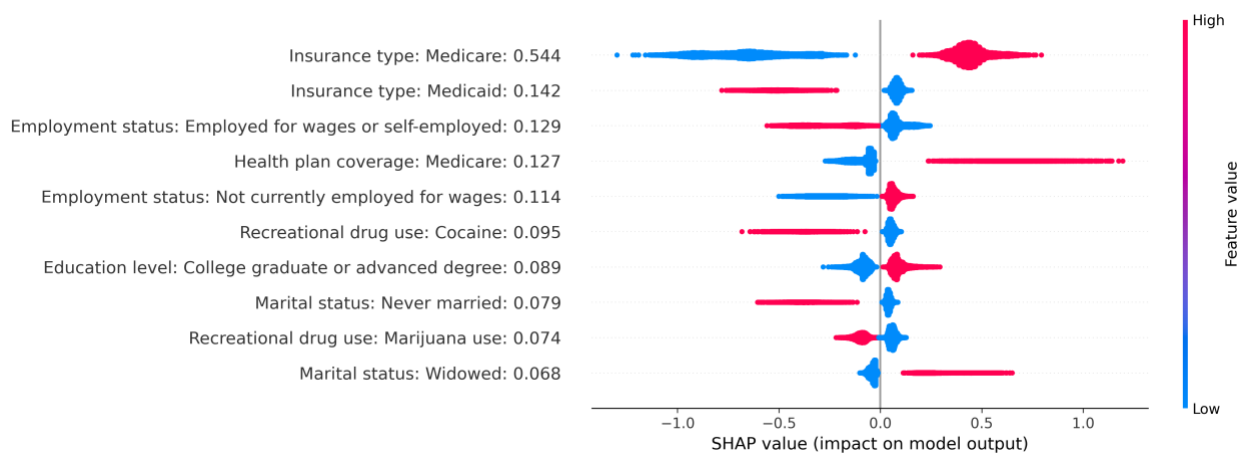
B.



*XGBoost: extreme gradient boosting, LR: logistic regression, MLP: multilayer perception

Figure 3. Distribution of the SHAP values for the top 10 features based on the highest mean absolute SHAP value. Each test sample is depicted as a point for every feature, with the x-axis indicating whether the feature's effect on the model's prediction is positive or negative. The color of each point reflects the feature's value, and this color scale is adjusted individually according to the value range present in the dataset.

A. SHAP values for the primary outcome using XGBoost



B. SHAP values for the secondary outcome using XGBoost

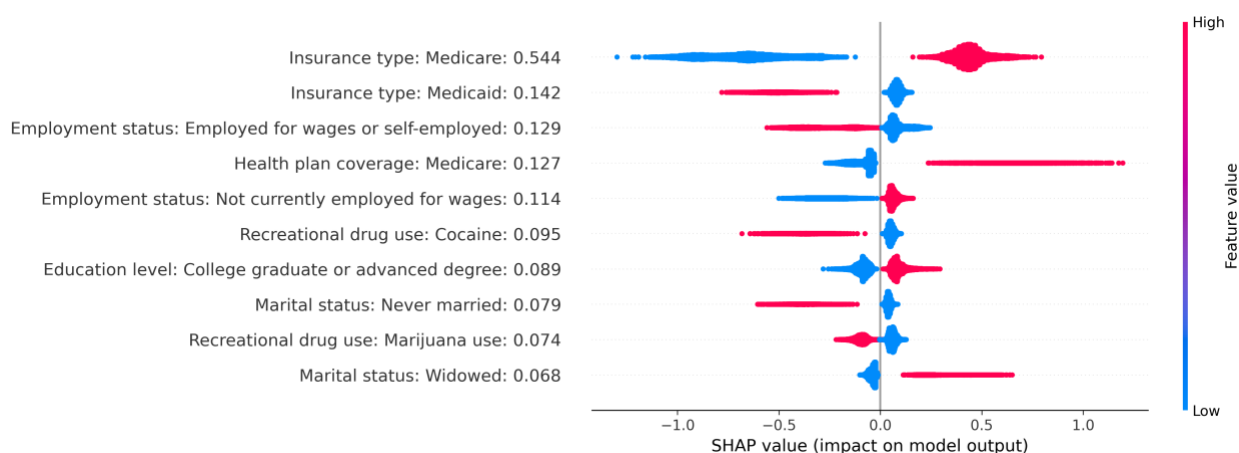
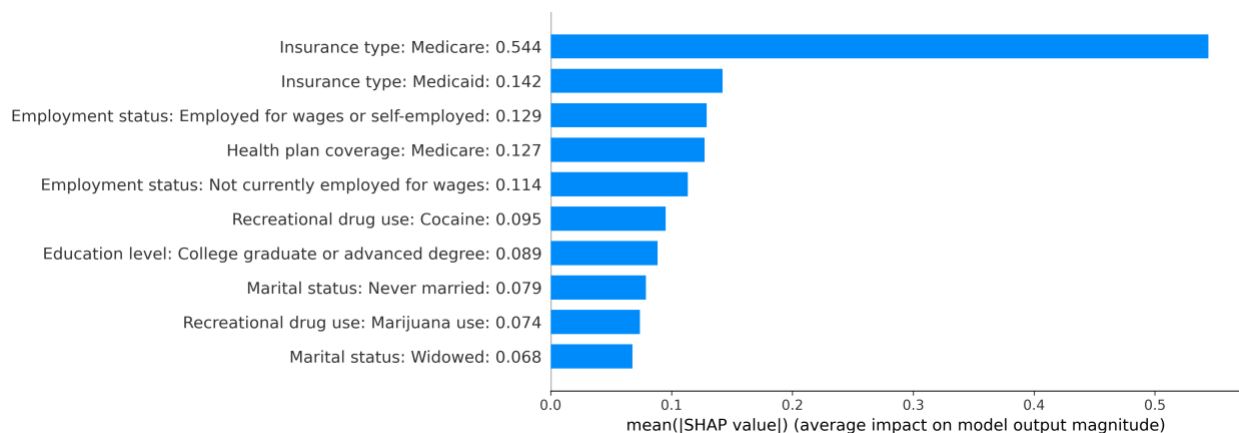


Figure 4. Mean absolute contribution of the top 10 features, ranked by their average SHAP value.

A. Feature importance of predicting the primary outcome, using XGBoost



B. Feature importance of predicting the secondary outcome, using XGBoost

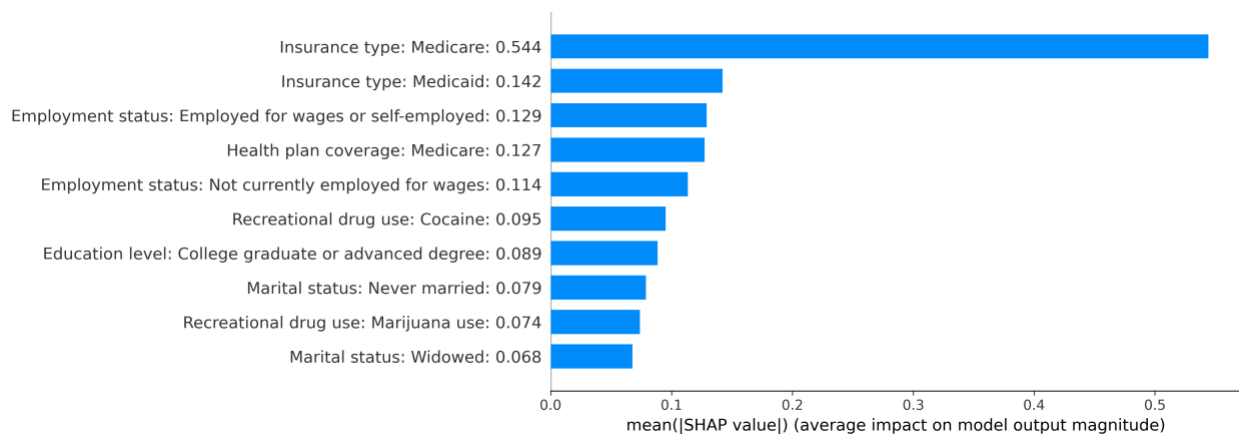


Table 1. Descriptive analysis of participants demographic information

Variables	Group, No. (%)			P	Group, No. (%)			P
	Overall, n = 99,936	Healthy aging (75), n = 44,109	Non- healthy aging (75), n = 55,827		Overall, n = 62,475	Healthy aging (85), n = 6,648	Non- healthy aging (85), n = 55,827	
Age, mean (SD)	74.0 (9.3)	79.9 (4.0)	69.4 (9.6)	<.001	71.2 (10.6)	87.2 (1.8)	69.4 (9.6)	<.001
Sex				<.001				<.001
Male	41,977 (42.0%)	20,493 (46.5%)	21,484 (38.5%)		24,671 (39.5%)	3,187 (47.9%)	32,839 (58.8%)	
Female	55,294 (55.3%)	22,455 (50.9%)	32,839 (58.8%)		36,101 (57.8%)	3,262 (49.1%)	21,484 (38.5%)	
Other	2,665 (2.7%)	1,161 (2.6%)	1,504 (2.7%)		1,703 (2.7%)	199 (3.0%)	1,504 (2.7%)	
Race				<.001				<.001
White	67,457 (67.5%)	33,916 (76.9%)	33,541 (60.1%)		38,802 (62.1%)	5,261 (79.1%)	33,541 (60.1%)	
Black or African American	14,612 (14.6%)	3,597 (8.2%)	11,015 (19.7%)		11,437 (18.3%)	422 (6.4%)	11,015 (19.7%)	
Asian	1,775 (1.8%)	941 (2.1%)	834 (1.5%)		954 (1.5%)	120 (1.8%)	834 (1.5%)	
Other/Unknown	16,092 (16.1%)	5,655 (12.8%)	10,437 (18.7%)		11,282 (18.1%)	845 (12.7%)	10,437 (18.7%)	
Modified Charlson comorbidity index, median (IQR)	1 (0-2)	0 (0-0)	2 (1-3)	<.001	2 (1-3)	0 (0-0)	2 (1-3)	<.001
Ethnicity				<.001				<.001
Not Hispanic or Latino	84,332 (84.4%)	38,693 (87.7%)	45,639 (81.8%)		51,477 (82.4%)	5,838 (87.8%)	45,639 (81.8%)	
Hispanic or Latino	11,109 (11.1%)	3,295 (7.5%)	7,814 (14.0%)		8,270 (13.2%)	456 (6.9%)	7,814 (14.0%)	
Other/Unknown	4,495 (4.5%)	2,121 (4.8%)	2,374 (4.2%)		2,728 (4.4%)	354 (5.3%)	2,374 (4.2%)	

