

1 Identifying probable dementia in undiagnosed Black and White Americans using machine
2 learning in Veterans Health Administration electronic health records

3
4 Yijun Shao^{1,2}, Kaitlin Todd³, Andrew Shutes-David^{3,4}, Steven P. Millard³, Karl Brown³, Amy
5 Thomas^{3,5}, Kathryn Chen,⁶ Katherine Wilson^{3,7}, Qing T. Zeng^{1,2}, Debby W. Tsuang^{3,6,*}

6
7 ¹ Washington DC VA Medical Center, Washington, DC, United States

8 ² George Washington University, Science and Engineering Hall, Washington, DC, United States

9 ³ Geriatric Research, Education, and Clinical Center, VA Puget Sound Health Care System,
10 Seattle, WA, United States

11 ⁴ Mental Illness Research, Education, and Clinical Center, VA Puget Sound Health Care System,
12 Seattle, WA, United States

13 ⁵ Department of Medicine, University of Washington, Seattle, WA, United States

14 ⁶ Department of Psychiatry and Behavioral Sciences, University of Washington, Seattle, WA,
15 United States

16 ⁷ Department of Biostatistics, University of Washington, Seattle, WA, United States

17

18 **RUNNING TITLE**

19 Identifying dementia using machine learning

20

21 **AUTHOR NOTE**

22 Correspondence concerning this article should be addressed to Debby W. Tsuang, Geriatric
23 Research, Education, and Clinical Center, S182 GRECC, VA Puget Sound Health Care System,
24 1660 S. Columbian Way, Seattle, WA 98108, USA; e-mail address: dwt1@uw.edu; phone: (206)
25 277-1333.

26

27 Kaitlin Todd is currently affiliated with Fred Hutchinson Cancer Research Center, 1100 Fairview
28 Ave. N. P.O. Box 19024 Seattle, WA 98109-1024. Kathryn Chen is currently affiliated with the
29 William S. Middleton Memorial Veterans Hospital in Madison, WI, and the University of
30 Wisconsin Department of Psychiatry, 2500 Overlook Terrace, Madison WI, 53705.

31

32

33 **ABSTRACT**

34 The application of machine learning (ML) tools in electronic health records (EHRs) can help
35 reduce the underdiagnosis of dementia, but models that are not designed to reflect minority
36 population may perpetuate that underdiagnosis. To address the underdiagnosis of dementia in
37 both Black Americans (BAs) and white Americans (WAs), we sought to develop and validate
38 ML models that assign race-specific risk scores. These scores were used to identify undiagnosed
39 dementia in BA and WA Veterans in EHRs. More specifically, risk scores were generated
40 separately for BAs (n=10K) and WAs (n=10K) in training samples of cases and controls by
41 performing ML, equivalence mapping, topic modeling, and a support vector-machine (SVM) in
42 structured and unstructured EHR data. Scores were validated via blinded manual chart reviews
43 (n=1.2K) of controls from a separate sample (n=20K). AUCs and negative and positive
44 predictive values (NPVs and PPVs) were calculated to evaluate the models. There was a strong
45 positive relationship between SVM-generated risk scores and undiagnosed dementia. BAs were
46 more likely than WAs to have undiagnosed dementia per chart review, both overall (15.3% vs
47 9.5%) and among Veterans with >90th percentile cutoff scores (25.6% vs 15.3%). With chart
48 reviews as the reference standard and varied cutoff scores, the BA model performed slightly
49 better than the WA model (AUC=0.86 with NPV=0.98 and PPV=0.26 at >90th percentile cutoff
50 vs AUC=0.77 with NPV=0.98 and PPV=0.15 at >90th). The AUCs, NPVs, and PPVs suggest that
51 race-specific ML models can assist in the identification of undiagnosed dementia, particularly in
52 BAs. Future studies should investigate implementing EHR-based risk scores in clinics that serve
53 both BA and WA Veterans.

54

55

56 **KEYWORDS**

57 electronic health record, dementia, machine learning, underdiagnosis, Veterans Health
58 Administration

59

60

61 **1 Introduction**

62 Alzheimer's disease (AD) and related dementias (ARD) are fatal neurodegenerative disorders
63 that account for half of admissions to long-term care facilities (Rice et al., 2001), yet nearly half
64 of those affected by ARD have not been formally diagnosed (Barnes et al., 2020, Amjad et al.,
65 2018). This crisis of underdiagnosis exacerbates existing disparities in health care, as dementia
66 underdiagnosis may disproportionately affect Black Americans (BAs) (Gianattasio et al., 2019).
67 In a large 2019 study of Medicare claims, older BAs with dementia were about two times less
68 likely to be correctly diagnosed with dementia than older White Americans (WAs) with
69 dementia (Gianattasio et al., 2019), and in one of the small handful of studies that examine racial
70 disparity in dementia care within VHA (Sleath et al., 2005, Kalkonde et al., 2009), significantly
71 fewer BA Veterans with suspected dementia underwent neuropsychological testing for the
72 diagnosis of dementia than WA Veterans with suspected dementia (Kalkonde et al., 2009). The
73 underdiagnosis of dementia translates into missed opportunities to treat patients (Cummings et
74 al., 2021), improve quality of life (e.g., through medication management and referrals) (Callahan
75 et al., 1995, Fitten et al., 1995), reduce patient and family burden (Sayegh and Knight, 2013,
76 Hinton et al., 2004), and reduce hospitalization, institutionalization, and health care costs
77 (Rasmussen and Langerman, 2019, Black et al., 2018).

78 We seek to use natural language processing (NLP) and machine learning (ML) tools to
79 address the magnitude of dementia diagnostic disparity in the Veterans Health Administration
80 (VHA) Corporate Data Warehouse (CDW), which is an ideal setting for this work, as it contains
81 comprehensive structured and unstructured data on ~0.4 million BA Veterans who are age 65+
82 and receive care as part of the largest integrated health care system in the nation. ML methods
83 have previously been applied to EHRs (Nadkarni et al., 2011, Gottesman et al., 2013), but we
84 have developed one of the *first* ML models to increase the sensitivity of dementia identification
85 by using *both* structured EHR data (e.g., demographics, diagnoses [ICD codes], procedures
86 [CPTS codes], medications, and clinical note types) and unstructured EHR data (e.g., words in
87 clinical notes) (Shao et al., 2019). In our previous work, we applied topic modeling and logistic
88 regression to develop risk scores for dementia based on the EHRs of older Veterans with
89 (n=1,861, mean age 79.8) and without (n=9,305, mean age 79.5) ICD-9 dementia codes (Shao et
90 al., 2019). Here, we extend this work by building separate predictive models for detecting
91 undiagnosed dementia in BAs and WAs using a larger sample of all VA patients who are 65+
92 years old with and without ICD 9/10 diagnosed dementia. We validate these models by
93 performing chart reviews blinded to dementia risk scores in a new set of patients who lack ICD-
94 9/10 dementia diagnoses and who were not used to build the models; we then compare the chart
95 review diagnoses to the diagnoses based on the model-generated risk scores.

96 **2 Materials and Methods**

97 *2.1 Study population*

98 After receiving IRB approval, we created a cohort of cases (i.e., Veterans with an ICD-9/10
99 dementia code) and controls (Veterans without any ICD-9/10 dementia codes) from the CDW by
100 selecting patients who turned age 65 between 1999 and 2018, lacked a dementia diagnosis at age
101 65, were previously evaluated at a VA clinic, and were identified as BA or WA in their EHRs
102 (top row, Figures 1a and 1b). The selected Veterans were followed until 9/12/2018, until
103 diagnosis (cases), or until censoring due to absence of records (controls).

104 To meet inclusion criteria, cases had to have received at least one ICD-9 or ICD-10 diagnosis
105 of dementia, with the first diagnosis occurring after age 65, and had to have at least 3 years of

106 continuous follow-up (i.e., 2+ documented clinical visits and associated notes during each year)
107 immediately prior to first diagnosis. That is, the one-year-long period in which first diagnosis
108 occurred had to have at least 3 visits (i.e., a diagnosis visit plus 2 previous visits), whereas the
109 other 2 one-year-long periods had to have at least 2 visits. Conversely, controls could not have
110 had any ICD-9/10 dementia codes; could not have filled donepezil, galantamine, rivastigmine, or
111 memantine prescriptions; and needed 3+ years of continuous follow-up (i.e., 2+ documented
112 clinical visits and associated notes during each year) after reaching age 62. We created separate
113 BA and WA cohorts of cases and controls to satisfy these criteria (second row, Figures 1a and
114 1b).

115 All clinical data were collected for a 3-year period that either immediately preceded but did
116 not include the first ICD-9/10 diagnosis of dementia (for cases) or a random visit date that was
117 selected as an index date (for controls). This 3-year period was established to provide adequate
118 structured and unstructured data.

119 The sampling and modeling of the Training and Validation Samples was performed
120 separately for BAs and WAs. We created model Training Samples by randomly sampling 5,000
121 cases and 5,000 controls in each race (total n=20,000). For each control, we randomly chose the
122 index visit among all visits that satisfied the 3-year lookback criterion. We used the Training
123 Samples to build models that produced dementia risk scores. We then created model Validation
124 Samples by randomly sampling 10,000 controls in each race who were not part of the Training
125 Samples (total n=20,000) and used the models to generate scores for these samples. Finally, we
126 sampled 600 Veterans from the Validation Samples for each race to undergo blinded chart
127 reviews (total n=1,200). Veterans were selected for chart review by simple random sampling
128 (n=200) and stratified random sampling (n=400) based on percentiles of the full Validation
129 Sample risk scores, such that 100 Veterans from the >75th – 90th percentiles were included, and
130 30 Veterans in each of the 10 remaining upper percentile ranges (i.e., 30 each from the >90th–
131 91st, >91st–92nd, etc.) were included.

132 2.2 Variable creation

133 2.2.1 Structured data

134 For each Veteran, we aggregated the structured data over the 3-year analysis period, recording
135 the presence/absence of each type of structured data during the 3-year period. Each type of
136 structured data was treated as a candidate binary variable for our model that would produce
137 dementia risk scores, with 0 indicating an absence of the codes/medications/note type and 1
138 indicating their presence.

139 To account for a transition from ICD-9 to ICD-10, we performed equivalence mapping,
140 visualizing the CDC/CMS general equivalence mappings (GEM) as a large bipartite graph that
141 consisted of two disjoint sets of vertices representing all the ICD-9 and ICD-10 codes,
142 respectively, and a number of edges connecting ICD-9 vertices to ICD-10 vertices representing
143 the possible conversions from ICD-9 codes to ICD-10 codes. These mappings allowed us to
144 decompose the GEM, viewed as a large bipartite graph, into a number of smaller disjoint
145 bipartite subgraphs that could not be decomposed into smaller disjoint subgraphs without
146 breaking edges. Then, for each of these minimal equivalence mappings, a new code was defined
147 to represent the group of ICD-9 codes before the transition date and the group of ICD-10 codes
148 after the transition. Variables corresponding to the new codes were defined in the same way as
149 other codes (e.g., CPT codes).

150 2.2.2 Unstructured data

151 Unstructured data were handled using the two-step topic modeling approach previously
152 described in Shao et al. (Shao et al., 2016, Shao et al., 2019). This unsupervised ML method
153 identifies shared topics from a large text corpus. Each topic is defined as a binary variable
154 indicating the presence/absence of that topic, and the proportion of topics within any particular
155 document is calculated. Here, we use the proportion of dementia-related topics observed in
156 excess in cases versus controls to identify dementia-related signs.

157 More specifically, raw topics were identified in clinical notes by running a latent Dirichlet
158 allocation (LDA) algorithm within the Machine Learning for Language Toolkit Java package
159 (Shao et al., 2016, Shao et al., 2019), which includes topic learning and inference functions. The
160 learning function is a time-consuming algorithm that learns the topics from a set of text
161 documents and generates a topic model, whereas the inference function runs much faster and can
162 apply the learned topic model to a new set of text documents and then infer the topic
163 distributions in those documents. For our topic learning subset, we randomly sampled one note
164 per day for each subject from the ~5 million notes collected during the 3-year study period,
165 yielding a sample corpus of 1.8 million notes. We randomly selected 1 million notes from this
166 sample corpus, which allowed for a reduced running time for topic learning while ensuring that
167 main topics were preserved. We then ran LDA topic learning 3 times on the 1 million sampled
168 notes, setting 1,000 as the total number of topics, and applied the 3 resulting models to all of the
169 5 million notes, using the topic inference function to infer the topic proportions in each note.
170 Based on the inferred topic proportions, we calculated the number of words that were associated
171 with each topic in each note by multiplying the topic proportion by the total number of words in
172 the note. Because the “number of words” associated with a topic was not always a whole
173 number, we call it the pseudo word count (PWC).

174 We then applied the stable topic extraction method (Shao et al., 2016, Shao et al., 2019),
175 which yielded 852 stable topics. For each stable topic, there were 3 topics—one from each run—
176 that were very similar to each other, and the stable topic was the “average” of the 3 similar
177 topics. Likewise, the PWC for the stable topic in each note was defined to be the median value of
178 the 3 PWCs corresponding to the 3 topics. By design, topic proportions are always positive
179 numbers, so the PWCs are positive as well. However, because not all of the topics are present in
180 every note, we set a nonzero threshold on the PWCs to indicate whether a topic was present in a
181 note. Empirically, we set the threshold at 2.0, which roughly means that a topic is present in a
182 note only when the $PWC \geq 2.0$. To allow various degrees of topic presence, we defined topic
183 presence to be a function of PWC as follows: (1) presence=0 if $PWC < 2.0$, (2)
184 presence= $PWC/10.0$ if $2.0 \leq PWC \leq 10.0$, and (3) presence=1.0 if $PWC > 10.0$. For the ML model,
185 stable topics were used as variables/features, and the maximum presence value over all the notes
186 of each Veteran was defined as the Veteran’s topic presence value.

187 2.3 Variable selection

188 Separately for BAs and WAs, we selected variables from the structured data that corresponded to
189 the codes/medications/note types that were present in 10+ Veterans in the Training Sample. All
190 of the stable topic variables and two demographic variables (age and sex) were selected. The age
191 variable was normalized so that the value 0 corresponded to 65 years old (minimum age)
192 whereas the value 1 corresponded to 85 years old (maximum age). All other variables were either
193 binary (i.e., values 0 and 1) or continuous (i.e., values between 0 and 1).

194 2.4 Support vector machine (SVM) model

195 Separately for BA and WA Veterans in the Training Sample, we constructed SVM models that
196 used the selected predictor variables to generate dementia “risk” scores. To construct the SVM
197 models, we used the linear SVM model (LinearSVC algorithm) in Python package *scikit-learn*
198 (Pedregosa et al., 2011). The SVM models had only one important hyperparameter: “C,” the cost
199 parameter, which sets the trade-off between misclassification and the simplicity of the decision
200 surface. To determine the best value for C, we performed five-fold cross-validation on the
201 training dataset with various values for C and then selected the value corresponding to the
202 highest predictive area under the receiver operating characteristic (ROC) curve (AUC) in the
203 five-fold cross-validation. The selected C value was used to train the final SVM model on the
204 entire training dataset. The linear SVM model output scores represent the distance to the
205 separation hyperplane in the high-dimensional feature space. The scores have no theoretical
206 limits, and higher scores mean indicate a higher likelihood of having dementia.

207 2.5 Validation of the SVM model

208 We separately generated scores for BA and WA controls in the Validation Sample and then, in a
209 subset of these Veterans, we performed chart reviews in which reviewers were blinded to score.
210 Chart reviews were conducted by experienced cognitive disorder experts (i.e., 2 trained in
211 geriatric psychiatry [DT and KC] and 1 in geriatric medicine [AT]) who achieved interrater
212 reliability on dually reviewed charts (Cohen’s Kappa value of 0.74 [se = 0.25, 95% CI = 0.25 - 1;
213 p = 0.0016]). The reviewers retroactively applied the DSM-V criteria for major neurocognitive
214 disorder (Sachdev et al., 2014) by evaluating memory, apraxia, aphasia, agnosia, executive
215 functioning, and functional domains of ADL and iAD (Katz, 1983) in abstracted notes.
216 Reviewers avoided attributing cognitive or functional deficits due to physical limitations or acute
217 or chronic medical conditions to dementia. When reviewers were uncertain about a Veteran’s
218 dementia status, that Veteran was labeled *uncertain* and then one of the other reviewers
219 adjudicated dementia status independent of the initial reviewer. Dementia status was coded by
220 reviewers as “None,” “Possible,” or “Probable”; a probable or possible dementia code thus
221 indicated that a Veteran had dementia symptoms that had either not been worked up nor
222 previously assigned a dementia diagnosis. Using chart review as the reference standard, we
223 assessed the prevalence of undiagnosed dementia and assessed the sensitivity, specificity,
224 positive predictive value (PPV), negative predictive value (NPV), and AUC by varying the
225 cutoff score for determining when to declare “possible or probable undiagnosed dementia.”
226 Estimates were computed using inverse probability weighting to account for stratified sampling
227 (Alonzo and Pepe, 2005), and confidence intervals were computed using bootstrapping.
228 Demographics, estimates, and confidence intervals were computed using R (R Core Team,
229 2020). We created scatter plots of dementia risks for 3 groups (probable, possible and none) as
230 well as 2 groups (probable/possible combined and none).

232 3 Results

233 3.1 Demographics

234 Among the Veterans who met inclusion/exclusion criteria (see Figures 1a and 1b), the prevalence
235 of dementia was 5.5% for BAs and 4.3% for WAs. Veterans ranged in age from 65 to 84 (see
236 demographics in Table 1). In the Training Sample, cases were older compared to controls (mean
237 [SD]=72.4 [4.8] vs. 69.1 [3.7]), and both cases and controls were overwhelmingly male (97.7 %
238 and 97.2%). BA Veterans were slightly younger than WA Veterans (72.1 [4.8] vs. 72.8 [4.8] for

239 cases; 68.6 [3.5] vs. 69.5 [3.8] for controls). The demographics for controls in the Validation and
240 Training Sample were similar.

241 3.2 Variable selection for the SVM model

242 For the model trained on BA Veterans, a total of 8221 features were selected, including 2
243 demographics, 854 topics, 2229 nondementia ICD code groups, 2561 CPT codes, 686
244 medications, and 1889 note types. For the model trained on WA Veterans, a total of 7716
245 features were selected, including 2 demographics, 854 topics, 2141 nondementia ICD code
246 groups, 2330 CPT codes, 655 medications, and 1734 note types.

247 The most significant topic features are shown in Supplemental Table 1. Note that the terms in
248 a topic can occur in any order or combination, and the presence of a topic in a document does not
249 require that all the terms in a topic be present. Topics that were observed more frequently in
250 cases than in controls were considered dementia related.

251 3.3 Distribution of scores

252 In the Training Sample, cases had higher scores than controls (mean [SD]=0.56 [0.54] vs. -0.50
253 [0.36] for BAs and 0.54 [0.55] vs. -0.47 [0.34] for WAs; Figure 2, Supplemental Figure 1). In the
254 Validation Sample, among Veterans with undiagnosed dementia who underwent chart review,
255 those diagnosed by reviewers with possible/probable dementia had higher scores compared to
256 those diagnosed with no dementia (0.45 [0.38] vs. -0.02 [0.51] for BAs, and 0.38 [0.41] vs. -0.02
257 [0.47] for WAs; Figure 3). For our chart review subsample of the Validation Sample, we
258 oversampled Veterans with higher scores (i.e., Veterans with chart reviews had higher scores
259 compared to all Validation Veterans: 0.05 [0.52] vs. -0.45 [0.41] for BA Veterans, and 0.02
260 [0.48] vs. -0.44 [0.38] for WA Veterans; Supplemental Figure 2), and therefore, we adjusted
261 scores using inverse probability weighting to account for stratified sampling.

262 3.4 Prevalence of undiagnosed dementia and screening test characteristics

263 Of the 1,200 Veterans who underwent chart review, 15.3% (n=92) of BAs and 9.5% (n=57) of
264 WAs were identified with possible/probable dementia. After adjusting for stratified sampling that
265 intentionally oversampled Veterans with higher scores, the estimated prevalence of undiagnosed
266 dementia in the full Validation Sample was 4.1% [3.2, 6.2] for BA Veterans and 3.6% [2.3, 6.3]
267 for WA Veterans. There was a strong positive relationship between risk scores and the
268 prevalence of undiagnosed dementia (Figure 4), and as anticipated, for Veterans with scores
269 below the 90th percentile, the percentages of undiagnosed dementia were low: 3.9% (95% CI
270 [2.1, 7.0]) and 2.9% (95% CI [1.3, 5.8]) for BA and WA Veterans, respectively. Among
271 Veterans with scores above the 90th percentile, we found that a higher percentage of BA
272 Veterans had undiagnosed dementia than WA Veterans: 25.6% (95% CI [20.9, 30.8]) vs. 15.3%
273 (95% CI [11.6, 19.8]).

274 Supplemental Figure 3 shows observed values for sensitivity, specificity, PPV, and NPV of
275 the screening tests that use chart review as the reference standard and vary cutoff score, and
276 Supplemental Table 2 lists values for various score cutoffs. As shown in Supplemental Figure 4,
277 the AUC was moderately high for both BA Veterans (0.86 [0.59, 0.95]) and WA Veterans (0.77
278 [0.59, 0.90]). For score cutoffs above the 50th percentile in the Validation Sample, sensitivity was
279 moderate and specificity was very high for both BA and WA Veterans (e.g., using the 90th
280 percentile as the cutoff, sensitivity and specificity were 0.61 [0.40, 0.76] and 0.92 [0.91, 0.92],
281 respectively, for BA Veterans and 0.43 [0.24, 0.67] and 0.91 [0.91, 0.92], respectively, for WA
282 Veterans). Because of the low prevalence of undiagnosed dementia in the full Validation
283 Samples, as well as the low sensitivity and high specificity of the screening tests, it was
284 unsurprising that PPV was low and NPV was high (Tenny and Hoffman, 2022); using the 90th

285 percentile as the cutoff, PPV was only 0.26 [0.21, 0.30] and 0.15 [0.12, 0.20] for BA and WA
286 Veterans, respectively. In contrast, NPVs remained quite high regardless of the score cutoff.

287 **4 Discussion**

288 *4.1 Significance*

289 We have successfully developed and validated a ML model to identify probable dementia
290 cases in BA and WA Veterans without ICD diagnoses. The dementia risk scores generated by the
291 SVM models were positively correlated with the diagnosis of dementia and achieved a high
292 AUC (0.86 [0.59, 0.95]) for BA Veterans and a satisfactory AUC for WA Veterans (0.77 [0.59,
293 0.90]). Given that BAs are about twice as likely to develop dementia as WAs (Tang et al., 2001,
294 Langa et al., 2017), the good performance of the SVM in this population is particularly
295 important.

296 *4.2 Context*

297 Our preliminary data suggest that BA Veterans have different risk factors for developing
298 dementia than WA Veterans. Using logistic regression to investigate risk factors for incident
299 dementia in all VHA, we identified different risk factors in older BA and WA Veterans (Cheng
300 et al., 2020). For example, among the key baseline characteristics that were significant predictors
301 of dementia in both races, stroke was a significantly stronger predictor among BAs, and Hispanic
302 ethnicity and depression were significantly stronger predictors among WAs ($p < 0.0001$). Those
303 findings motivated the development of the race-specific risk models proposed in the current
304 study, which instead focuses on prediction.

305 Many studies have applied NLP and ML methods in dementia (Chang et al., 2021),
306 particularly in the context of neuroimaging (Popuri et al., 2020, Qiu et al., 2020) or in the use of
307 EHRs to identify cognitive impairment or diagnosed dementia (Amra et al., 2017, Wray et al.,
308 2014), yet few studies have sought to use EHRs as a direct phenotyping tool for undiagnosed
309 dementia. Researchers in the UK developed models (including SVM) to identify patients with
310 dementia (Jammeh et al., 2018), and Kaiser Permanente/UCSF researchers developed the
311 eRADAR tool in research participants and then validated it in two health-care systems (Barnes et
312 al., 2020, Coley et al., 2022); both studies limited their EHR interrogations to structured data and
313 have shown some success in identifying undiagnosed dementia. Likewise, Yadgir et al. used ML
314 to identify structured variables associated with cognitive impairment in ER patients (Yadgir et
315 al., 2022). Conversely, Boustani et al. have developed passive digital signatures for ADRD by
316 searching for predetermined variables in *both* structured and unstructured EHR data, and their
317 work suggests that the combination can improve AUC by up to .11 (Boustani et al., 2020);
318 however, like Barnes et al., Boustani et al. use curated, preselected search terms rather than
319 leveraging the potential of supervised ML to identify topic features associated with dementia.

320 Rather than employing a targeted-word study design like Barnes et al. or Boustani et al., we
321 have sought to improve the identification of dementia by combining supervised ML with an
322 improved clinical standard. More specifically, we have sought to improve upon EHR ICD
323 codes as the basis for ML by incorporating chart reviews by reviewers who have been blinded to
324 the initial ML-derived dementia likelihood scores. We previously published a ML logistic
325 regression model that used this approach on a smaller scale, applying supervised ML to
326 structured and unstructured data from EHRs to identify topics associated with dementia and then
327 identify cases with undiagnosed dementia (Shao et al., 2019). That study included blinded
328 manual reviews in a smaller sample ($n=140$) than our current work and produced a sensitivity of
329 0.825 and a specificity of 0.832. It also had older Veterans (i.e., an average age of 80 vs. 71 in

330 this study); complications with controls in the logistic regression model; and an ad-hoc
331 stratification method for computing sensitivity and specificity, whereas our SVM models avoid
332 these idiosyncrasies in a much larger (n=1,200) and more diverse (600 BA and 600 WA)
333 validation effort.

334 EHR tools and ML models that do not specifically attempt to reflect minoritized
335 communities are more likely to unintentionally generate cycles of exclusion and to thereby
336 perpetuate underdiagnosis in BAs rather than addressing underdiagnosis (Bracic et al., 2022). To
337 our knowledge the present study is the first effort to develop and evaluate a model that
338 specifically focuses on BAs.

339

340 *4.3 Implications*

341 We seek to develop EHR-based dementia risk scores to support future screening of dementia
342 in clinical settings that include both WAs and BAs. Other researchers have noted that PPV and
343 NPV are better at assessing a screening test in clinical practice than sensitivity and specificity
344 (Trevethan, 2017). Our model generated a very high NPV at the 90th percentile for both BA
345 Veterans (0.98 [0.96, 0.99]) and WA Veterans (0.98 [0.94, 0.99]). These findings are similar to
346 the NPVs reported with the eRADAR tool in an EHR sample that was 89% WA (Barnes et al.,
347 2020) but higher than the NPV reported by Yadgir et al (i.e., 0.93) (Yadgir et al., 2022). The
348 PPV in our study was low for both BA Veterans (0.26 [0.21,0.31]) and WA Veterans (0.15
349 [0.12,0.20]) at the 90th percentile cutoff. Practically, this means that at that threshold, about a
350 quarter of the BAs and a seventh of the WAs who were flagged by our model as having potential
351 dementia would actually have dementia according to our manual chart reviews. In contrast,
352 Yadgir et al., achieved PPVs greater than 0.4, but to do so, they applied threshold cutoffs higher
353 than 0.8; this meant that they obtained a high true positive rate at the expense of low sensitivity,
354 which is not optimal as a screening instrument given the high cutoff scores (Yadgir et al., 2022).
355 Our algorithms compare very favorably to the eRADAR tool for dementia, which had a PPV of
356 0.115 in a research setting and 0.020 to 0.048 in patients (Barnes et al., 2020, Coley et al., 2022).
357 Our PPV is similar or superior to the rates of standard screening methods for cancers like
358 mammograms or colonoscopy (reviewed in Barnes et al., 2020). However, cancer screening is
359 often followed by more definitive tests, such as ultrasound and/or biopsies, and thus low PPVs in
360 screening tests may be acceptable. Development of multitier screening and diagnostic tests are
361 therefore necessary prior to the implementation of our SVM model in clinical workflows.

362 *4.4 Limitations and future work*

363 The VA patient population skews heavily toward older males, and our training and test data
364 thus had a low percentage of females; that may limit the generalizability of our final ML models
365 outside VHA, though we expect that the same steps could be applied to generate risk scores
366 within other health care systems with more females. In evaluating the low PPVs in our study, it
367 may be that our standards for the diagnosis of dementia (i.e., manual chart review) are flawed ,
368 as and that due to insufficient information in the charts, we were unable to retrospectively apply
369 the newest AD criteria (NIA-Reagan) are flawed. {Hyman, 1997, 9329452} If signs and
370 symptoms relevant to impairment are not mentioned in clinical notes, reviewers are unable to
371 assign an dementia diagnosis due to insufficient information. Here, this may have led to a low
372 level of dementia prevalence, and a low prevalence of any condition leads to models with high
373 NPVs and low PPVs. It is possible, therefore, that our model may catch signs of dementia that
374 cannot be captured by a manual chart review, which means our model may perform better when

375 compared to more accurate diagnostic standards, like in-person expert diagnoses or
376 neuropathological assessments; this represents a promising area for future research.

377 We recognize that future studies need to assess the portability of the ML models that we have
378 developed. Not all EHRs have notes available to researchers (due to privacy issues), and in those
379 instances, researchers will be unable to leverage the full benefit of our models' ability to search
380 both structured and unstructured data. Future studies should investigate how other ML methods,
381 like deep learning approaches, might improve the detection of undiagnosed dementia; solicit
382 input from BA stakeholders regarding model implementation in clinical processes; and
383 investigate the implementation of our EHR-based risk scores in clinics that serve both BA and
384 WA Veterans.

385 386 **Acknowledgments**

387 This work was supported by NIA R56 AG059739 and was supported in part by the U. S.
388 Department of Veterans Affairs Office of Research and Development Biomedical Laboratory
389 Research Program.

390 391 **Contribution to the field**

392 Up to 50% of dementia is underdiagnosed in primary care settings. This underdiagnosis is
393 especially problematic in elderly Black Americans. Screening all elderly patients (or all elderly
394 Black patients) is time- and resource-intensive, and broad-based screening is not aligned with
395 current clinical guidelines. Thus, other approaches are necessary. This illustrates that cost-
396 effective machine learning algorithms can identify a subset of patients within existing electronic
397 health records who are at high risks for developing dementia. Furthermore, we are one of the first
398 to develop a race-specific algorithm in the context of dementia identification and to thereby
399 leverage machine learning to specifically address dementia-related health-care disparities in
400 Black Americans. That is critical, as models that are not designed to reflect minority population
401 may instead perpetuate underdiagnosis. Moreover, this work may also have tangible financial
402 benefits. Incorporating electronic health records-based algorithms into screening workflows with
403 diagnostic tests as follow-up could focus resources where they will have the most impact in
404 primary care settings, including the prevention of costly health care events that otherwise tend to
405 precede diagnosis.

406 407 **Author contributions:**

408 Study design (DT, QZ, SM); model development (YS, QZ, DT); data analyses (SM, KT, KB,
409 KW); clinical expertise (DT, AT, KC); initial drafting of the paper and literature review (SM,
410 ASD)

411 412 **Dataset restrictions:**

413 The VA EHR data resides in VINCI behind VA firewalls. VA-approved investigators can access
414 the data. SVM algorithms can be made available to interested qualified investigators.

415 416 **Conflict of Interest/Disclosure Statement**

417 The authors have no conflicts of interest to report.
418

419 **References**

- 420 ALONZO, T. & PEPE, M. 2005. Assessing accuracy of a continuous screening test in the
421 presence of verification bias. *J Royal Stat Soc Series C Appl Stats*, 54, 173-190.
- 422 AMJAD, H., ROTH, D. L., SHEEHAN, O. C., LYKETSOS, C. G., WOLFF, J. L. & SAMUS,
423 Q. M. 2018. Underdiagnosis of Dementia: an Observational Study of Patterns in
424 Diagnosis and Awareness in US Older Adults. *J Gen Intern Med*, 33, 1131-1138.
- 425 AMRA, S., O'HORO, J. C., SINGH, T. D., WILSON, G. A., KASHYAP, R., PETERSEN, R.,
426 ROBERTS, R. O., FRYER, J. D., RABINSTEIN, A. A. & GAJIC, O. 2017. Derivation
427 and validation of the automated search algorithms to identify cognitive impairment and
428 dementia in electronic health records. *J Crit Care*, 37, 202-205.
- 429 BARNES, D. E., ZHOU, J., WALKER, R. L., LARSON, E. B., LEE, S. J., BOSCARDIN, W. J.,
430 MARCUM, Z. A. & DUBLIN, S. 2020. Development and validation of eRADAR: A tool
431 using EHR data to detect unrecognized dementia. *J Am Geriatr Soc*, 68, 103-111.
- 432 BLACK, C. M., FILLIT, H., XIE, L., HU, X., KARIBURYO, M. F., AMBEGAONKAR, B. M.,
433 BASER, O., YUCE, H. & KHANDKER, R. K. 2018. Economic burden, mortality, and
434 institutionalization in patients newly diagnosed with AD. *J Alzheimers Dis*, 61, 185-193.
- 435 BOUSTANI, M., PERKINS, A. J., KHANDKER, R. K., DUONG, S., DEXTER, P. R.,
436 LIPTON, R., BLACK, C. M., CHANDRASEKARAN, V., SOLID, C. A. &
437 MONAHAN, P. 2020. Passive digital signature for early identification of Alzheimer's
438 disease and related dementia. *J Am Geriatr Soc*, 68, 511-518.
- 439 BRACIC, A., CALLIER, S. L. & PRICE, W. N., 2ND 2022. Exclusion cycles: Reinforcing
440 disparities in medicine. *Science*, 377, 1158-1160.
- 441 CALLAHAN, C. M., HENDRIE, H. C. & TIERNEY, W. M. 1995. Documentation and
442 evaluation of cognitive impairment in elderly primary care patients. *Ann Intern Med*, 122,
443 422-9.
- 444 CHANG, C. H., LIN, C. H. & LANE, H. Y. 2021. Machine Learning and Novel Biomarkers for
445 the Diagnosis of Alzheimer's Disease. *Int J Mol Sci*, 22.
- 446 CHENG, Y., AHMED, A., ZAMRINI, E., TSUANG, D., SHERIFF, H. & ZENG-TREITLER,
447 Q. 2020. AD and ADRD in older African American and white Veterans. *J Alzheimers*
448 *Dis*, 75, 311-320.
- 449 COLEY, R. Y., SMITH, J. J., KARLINER, L., IDU, A. E., LEE, S. J., FULLER, S., LAM, R.,
450 BARNES, D. E. & DUBLIN, S. 2022. External Validation of the eRADAR Risk Score
451 for Detecting Undiagnosed Dementia in Two Real-World Healthcare Systems. *J Gen*
452 *Intern Med*.
- 453 CUMMINGS, J., AISEN, P., LEMERE, C., ATRI, A., SABBAGH, M. & SALLOWAY, S.
454 2021. Aducanumab produced a clinically meaningful benefit in association with amyloid
455 lowering. *Alzheimers Res Ther*, 13, 98.
- 456 FITTEN, L. J., PERRYMAN, K. M., WILKINSON, C. J., LITTLE, R. J., BURNS, M. M.,
457 PACHANA, N., MERVIS, J. R., MALMGREN, R., SIEMBIEDA, D. W. & GANZELL,
458 S. 1995. Alzheimer and vascular dementias and driving. *JAMA*, 273, 1360-5.
- 459 GIANATTASIO, K. Z., PRATHER, C., GLYMOUR, M. M., CIARLEGLIO, A. & POWER, M.
460 C. 2019. Racial disparities and temporal trends in dementia misdiagnosis risk in the US.
461 *Alzheimers Dement (N Y)*, 5, 891-898.
- 462 GOTTESMAN, O., KUIVANIEMI, H., TROMP, G., FAUCETT, W. A., LI, R., MANOLIO, T.
463 A., SANDERSON, S. C., KANNRY, J., ZINBERG, R., BASFORD, M. A.,
464 BRILLIANT, M., CAREY, D. J., CHISHOLM, R. L., CHUTE, C. G., CONNOLLY, J.

- 465 J., CROSSLIN, D., DENNY, J. C., GALLEGO, C. J., HAINES, J. L., HAKONARSON,
466 H., HARLEY, J., JARVIK, G. P., KOHANE, I., KULLO, I. J., LARSON, E. B.,
467 MCCARTY, C., RITCHIE, M. D., RODEN, D. M., SMITH, M. E., BOTTINGER, E. P.,
468 WILLIAMS, M. S. & E, M. N. 2013. The eMERGE Network. *Genet Med*, 15, 761-71.
- 469 HINTON, L., FRANZ, C. & FRIEND, J. 2004. Pathways to dementia diagnosis: Evidence for
470 cross-ethnic differences. *Alzheimer Dis Assoc Disord*, 18, 134-44.
- 471 HYMAN, B. T. & TROJANOWSKI, J. Q. 1997. Consensus recommendations for the
472 postmortem diagnosis of Alzheimer disease from the National Institute on Aging and the
473 Reagan Institute Working Group on diagnostic criteria for the neuropathological
474 assessment of Alzheimer disease. *J Neuropathol Exp Neurol*, 56, 1095-7.
- 475 JAMMEH, E. A., CARROLL, C. B., PEARSON, S. W., ESCUDERO, J., ANASTASIOU, A.,
476 ZHAO, P., CHENORE, T., ZAJICEK, J. & IFEACHOR, E. 2018. ML based
477 identification of undiagnosed dementia in primary care: A feasibility study. *BJGP Open*,
478 2, [bjgpopen18X101589](https://doi.org/10.1136/bjgpopen18X101589).
- 479 KALKONDE, Y. V., PINTO-PATARROYO, G. P., GOLDMAN, T., STRUTT, A. M., YORK,
480 M. K., KUNIK, M. E. & SCHULZ, P. E. 2009. Ethnic disparities in the treatment of
481 dementia in Veterans. *Dement Geriatr Cogn Disord*, 28, 145-52.
- 482 KATZ, S. 1983. Assessing self-maintenance: activities of daily living, mobility, and instrumental
483 activities of daily living. *J Am Geriatr Soc*, 31, 721-7.
- 484 LANGA, K. M., LARSON, E. B., CRIMMINS, E. M., FAUL, J. D., LEVINE, D. A., KABETO,
485 M. U. & WEIR, D. R. 2017. A comparison of the prevalence of dementia in the US in
486 2000 and 2012. *JAMA Intern Med*, 177, 51-58.
- 487 NADKARNI, P. M., OHNO-MACHADO, L. & CHAPMAN, W. W. 2011. Natural language
488 processing. *J Am Med Inform Assoc*, 18, 544-51.
- 489 PEDREGOSA, F., VAROQUAUX, G., GRAMFORT, A., MICHEL, V., THIRION, B.,
490 GRISEL, O., BLONDEL, M., PRETTENHOFER, P., WEISS, R., DUBOURG, V.,
491 VANDERPLAS, J., PASSOS, A., COURNAPEAU, D., BRUCHER, M., PERROT, M.
492 & DUCHESNAY, E. 2011. Scikit-learn: Machine Learning in Python. *J Mach Learn*
493 *Res*, 12.
- 494 POPURI, K., MA, D., WANG, L. & BEG, M. F. 2020. Using machine learning to quantify
495 structural MRI neurodegeneration patterns of Alzheimer's disease into dementia score:
496 Independent validation on 8,834 images from ADNI, AIBL, OASIS, and MIRIAD
497 databases. *Hum Brain Mapp*, 41, 4127-4147.
- 498 QIU, S., JOSHI, P. S., MILLER, M. I., XUE, C., ZHOU, X., KARJADI, C., CHANG, G. H.,
499 JOSHI, A. S., DWYER, B., ZHU, S., KAKU, M., ZHOU, Y., ALDERAZI, Y. J.,
500 SWAMINATHAN, A., KEDAR, S., SAINT-HILAIRE, M. H., AUERBACH, S. H.,
501 YUAN, J., SARTOR, E. A., AU, R. & KOLACHALAMA, V. B. 2020. Development and
502 validation of an interpretable deep learning framework for Alzheimer's disease
503 classification. *Brain*, 143, 1920-1933.
- 504 R CORE TEAM 2020. The R project for statistical computing. Available from: [https://www.r-](https://www.r-project.org/)
505 [project.org/](https://www.r-project.org/).
- 506 RASMUSSEN, J. & LANGERMAN, H. 2019. AD: Why we need early diagnosis. *Degener*
507 *Neurol Neuromuscul Dis*, 9, 123-130.
- 508 RICE, D. P., FILLIT, H. M., MAX, W., KNOPMAN, D. S., LLOYD, J. R. & DUTTAGUPTA,
509 S. 2001. Prevalence, costs, and treatment of ADRD: a managed care perspective. *Am J*
510 *Manag Care*, 7, 809-18.

- 511 SACHDEV, P. S., BLACKER, D., BLAZER, D. G., GANGULI, M., JESTE, D. V., PAULSEN,
512 J. S. & PETERSEN, R. C. 2014. Classifying neurocognitive disorders: the DSM-5
513 approach. *Nat Rev Neurol*, 10, 634-42.
- 514 SAYEGH, P. & KNIGHT, B. G. 2013. Cross-cultural differences in dementia. *Int Psychogeriatr*,
515 25, 517-30.
- 516 SHAO, Y., MOHANTY, A. F., AHMED, A., WEIR, C. R., BRAY, B. E., SHAH, R. U., REDD,
517 D. & ZENG-TREITLER, Q. 2016. Identification and use of frailty indicators from text to
518 examine associations with clinical outcomes among patients with heart failure. *AMIA*
519 *Annu Symp Proc*, 2016, 1110-1118.
- 520 SHAO, Y., ZENG, Q. T., CHEN, K. K., SHUTES-DAVID, A., THIELKE, S. M. & TSUANG,
521 D. W. 2019. Detection of probable dementia cases in undiagnosed patients using
522 structured and unstructured EHRs. *BMC Med Inform Decis Mak*, 19, 128.
- 523 SLEATH, B., THORPE, J., LANDERMAN, L. R., DOYLE, M. & CLIPP, E. 2005. African-
524 American and white caregivers of older adults with dementia. *J Am Geriatr Soc*, 53, 397-
525 404.
- 526 TANG, M. X., CROSS, P., ANDREWS, H., JACOBS, D. M., SMALL, S., BELL, K.,
527 MERCHANT, C., LANTIGUA, R., COSTA, R., STERN, Y. & MAYEUX, R. 2001.
528 Incidence of AD in African-Americans, Caribbean Hispanics, and Caucasians in northern
529 Manhattan. *Neurology*, 56, 49-56.
- 530 TENNY, S. & HOFFMAN, M. 2022. *StatPearls*, Treasure Island, FL, StatPearls Publishing.
- 531 TREVETHAN, R. 2017. Sensitivity, Specificity, and Predictive Values: Foundations,
532 Pliabilities, and Pitfalls in Research and Practice. *Front Public Health*, 5, 307.
- 533 WRAY, L. O., WADE, M., BEEHLER, G. P., HERSHEY, L. A. & VAIR, C. L. 2014. A
534 program to improve detection of undiagnosed dementia in primary care and its
535 association with healthcare utilization. *Am J Geriatr Psychiatry*, 22, 1282-91.
- 536 YADGIR, S. R., ENGSTROM, C., JACOBSON, G. C., GREEN, R. K., JONES, C. M. C.,
537 CUSHMAN, J. T., CAPRIO, T. V., KIND, A. J. H., LOHMEIER, M., SHAH, M. N. &
538 PATTERSON, B. W. 2022. Machine learning-assisted screening for cognitive
539 impairment in the emergency department. *J Am Geriatr Soc*, 70, 831-837.
- 540

Table 1a. Demographics of the Training Sample by Race (BA: Black American; WA: White American).

Characteristic	Case (n = 10K)			Control (n = 10K)		
	BA (n = 5K)	WA (n = 5K)	Combined (n = 10K)	BA (n = 5K)	WA (n = 5K)	Combined (n = 10K)
Age, mean (SD)	72.1 (4.8)	72.8 (4.8)	72.4 (4.8)	68.6 (3.5)	69.5 (3.8)	69.1 (3.7)
Age category, (%)						
65–69	35.7	29.8	32.8	70.5	60.0	65.2
70–74	34.1	34.4	34.3	22.1	28.4	25.3
75–79	20.9	24.2	22.6	5.8	9.2	7.5
80–84	9.4	11.5	10.4	1.6	2.5	2.0
Gender, % male	97.9	97.5	97.7	96.8	97.6	97.2

Table 1b. Demographics of the Validation Sample by Race (BA: Black American; WA: White American).

Characteristic	Full Validation Sample (n = 20K)			Chart Review (n = 1,200)*					
	BA (n = 10K)	WA (n = 10K)	Combined (n = 20K)	Unweighted			Weighted†		
	BA (n = 10K)	WA (n = 10K)	Combined (n = 20K)	BA (n = 600)	WA (n = 600)	Combined (n = 1,200)	BA (n = 600)	WA (n = 600)	Combined (n = 1,200)
Age, mean (SD)	68.5 (3.4)	69.5 (3.8)	69.0 (3.6)	69.3 (4.2)	70.2 (4.5)	69.8 (4.3)	68.5 (3.4)	69.3 (3.7)	68.9 (3.6)
Age category, (%)									
65–69	70.9	60.3	65.6	64.8	52.8	58.8	70.4	60.9	65.7
70–74	22.3	28.6	25.5	22.3	28.5	25.4	23.2	28.4	25.8
75–79	5.5	8.7	7.1	9.3	13.7	11.5	5.4	8.2	6.7
80–84	1.3	2.4	1.9	3.5	5.0	4.3	1.0	2.5	1.7
Gender, % male	96.1	97.6	96.8	97.2	97.0	97.1	96.7	98.9	97.8

* Patients who underwent chart review were a subset of the full Validation Sample, selected by a combination of random and stratified sampling as described in the text.

[†] Observations were weighted by the inverse probability of being sampled from the full Validation Sample.

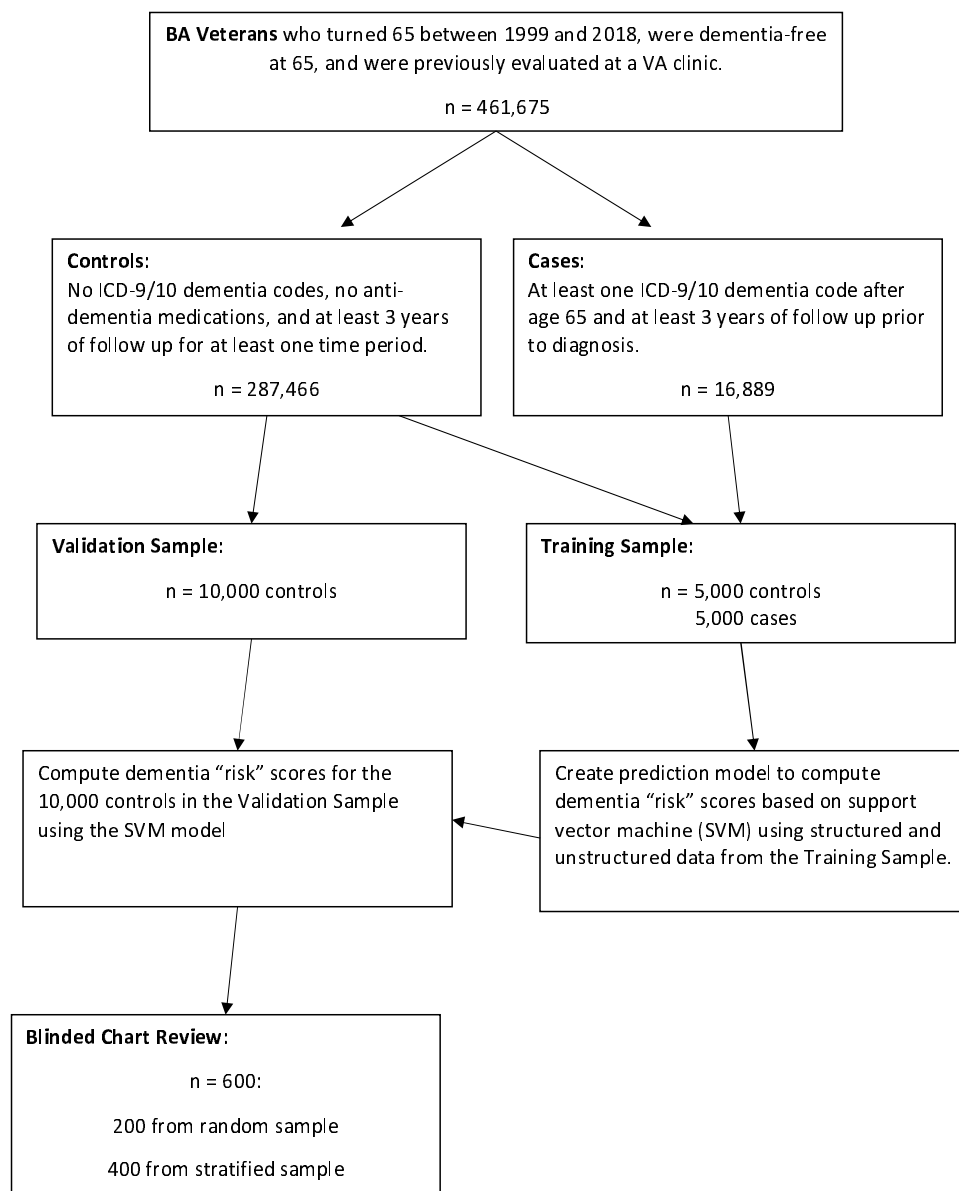


Figure 1a. Study flow diagram for Black American (BA) Veterans. This figure shows the number of BA Veterans available within the Veterans Health Administration (VHA) Corporate Data Warehouse (CDW) for the time period under study who met inclusion/exclusion criteria, as well as the number of Veterans used for model building and validation. Veterans in the Training Sample and Validation Sample were chosen with simple random sampling. Veterans who underwent chart review (blinded to score) were chosen from the 10,000 in the Validation Sample by simple random sampling (n = 200) and stratified random sampling (n = 400), where the strata were based on the scores.

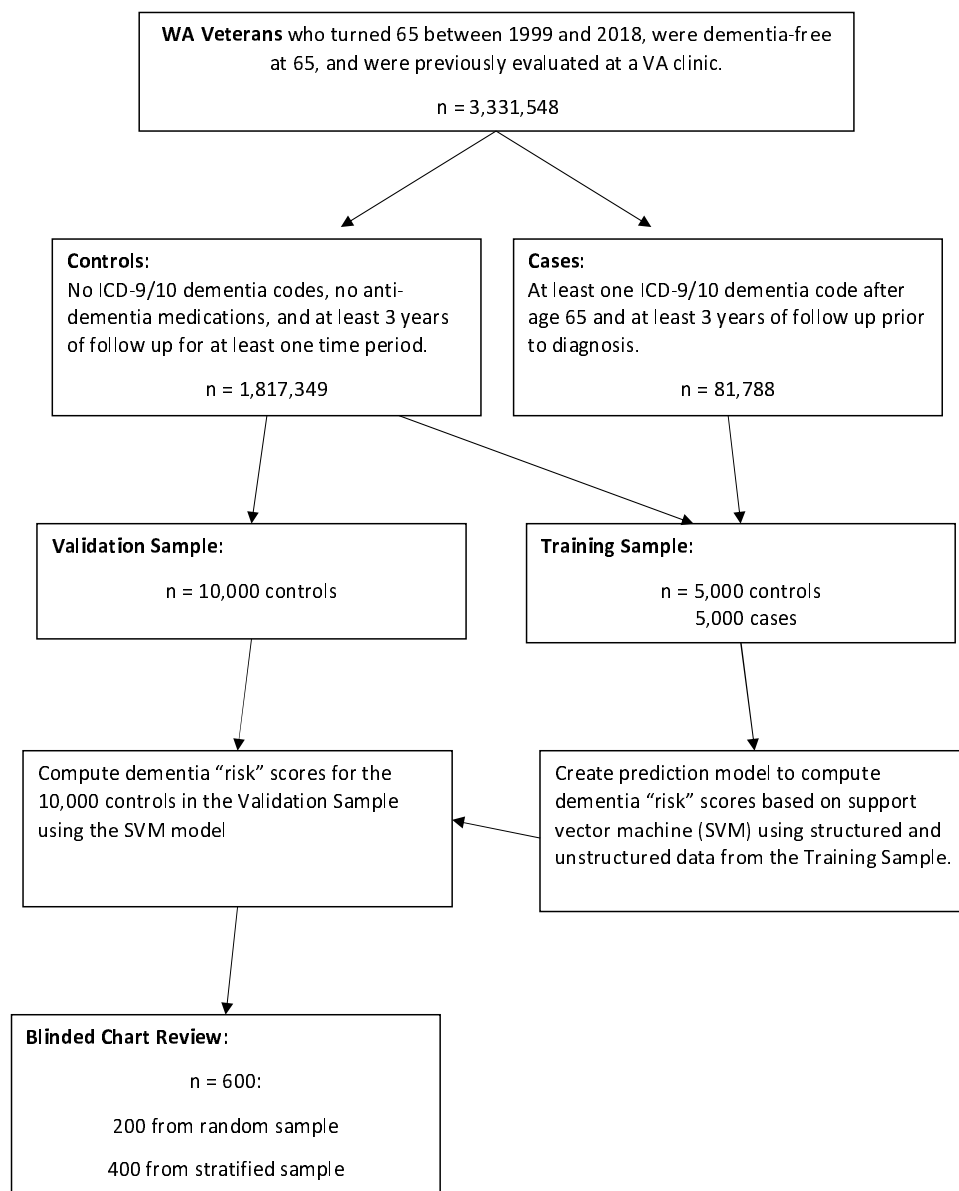
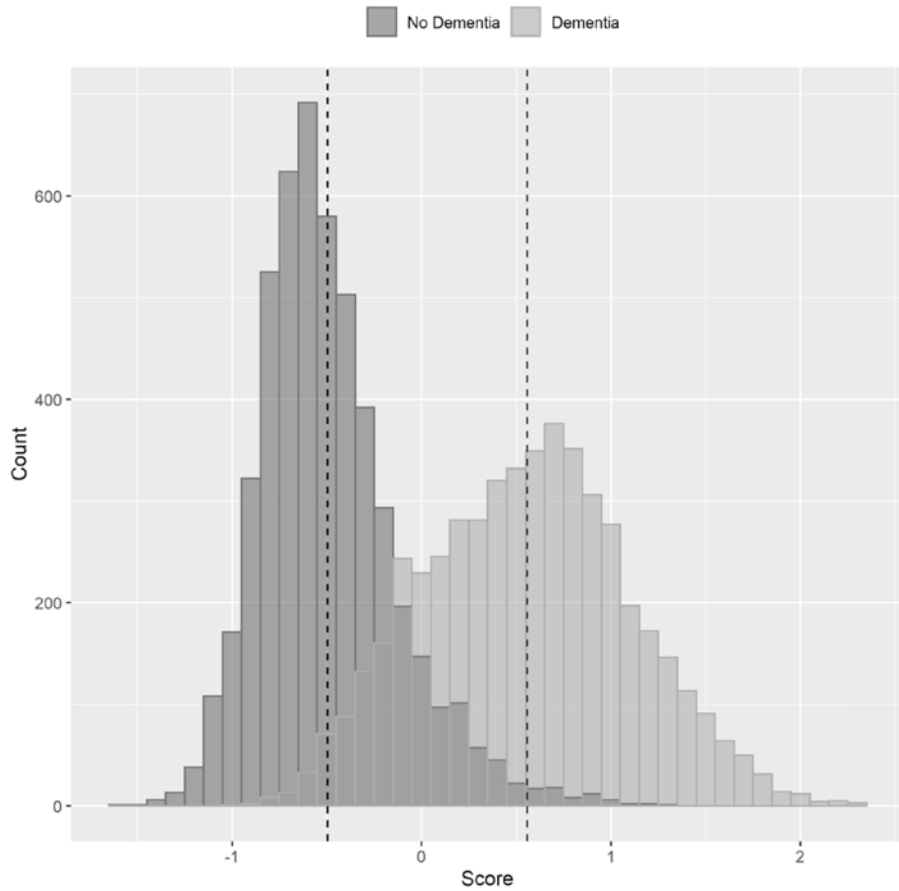


Figure 1b. Study flow diagram for White American (WA) Veterans. This figure shows the number of WA Veterans available within the Veterans Health Administration (VHA) Corporate Data Warehouse (CDW) for the time period under study who met inclusion/exclusion criteria, as well as the number of Veterans used for model building and validation. Veterans in the Training Sample and Validation Sample were chosen with simple random sampling. Veterans who underwent chart review (blinded to score) were chosen from the 10,000 in the Validation Sample by simple random sampling ($n = 200$) and stratified random sampling ($n = 400$), where the strata were based on the scores.

BA



WA

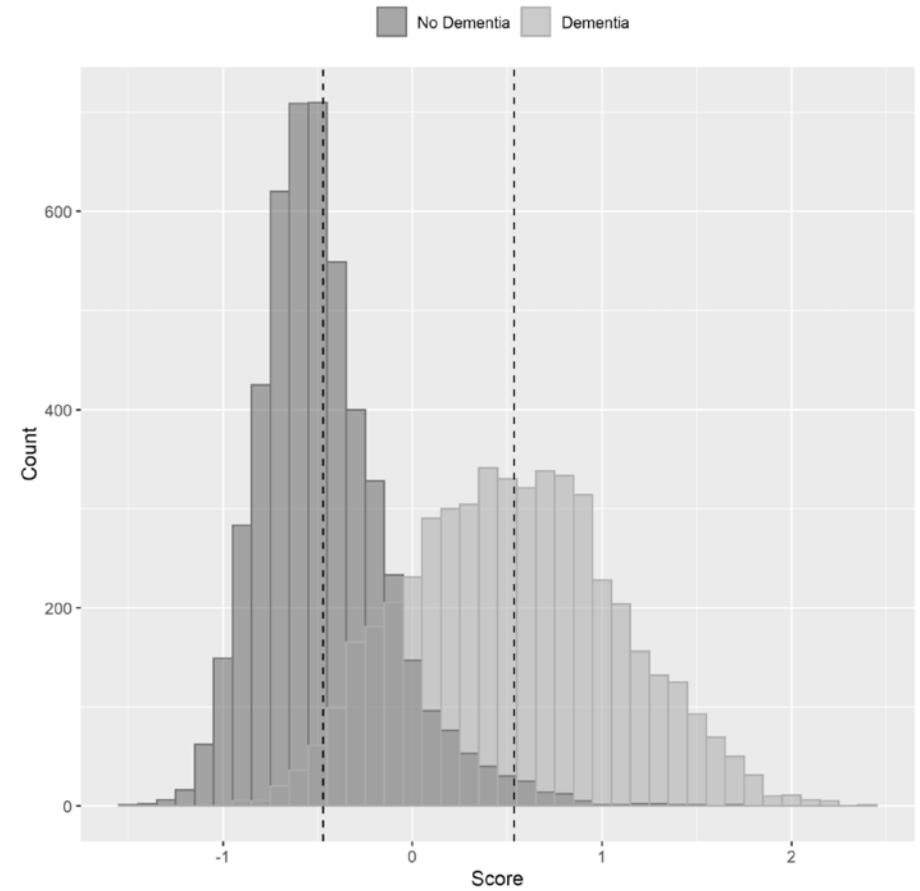


Figure 2. Distribution of scores by dementia status and race (BA: Black American; WA: White American) for Veterans in the Training Sample (n = 5,000 in each dementia status group for each race). Dashed lines represent the means of the distribution.

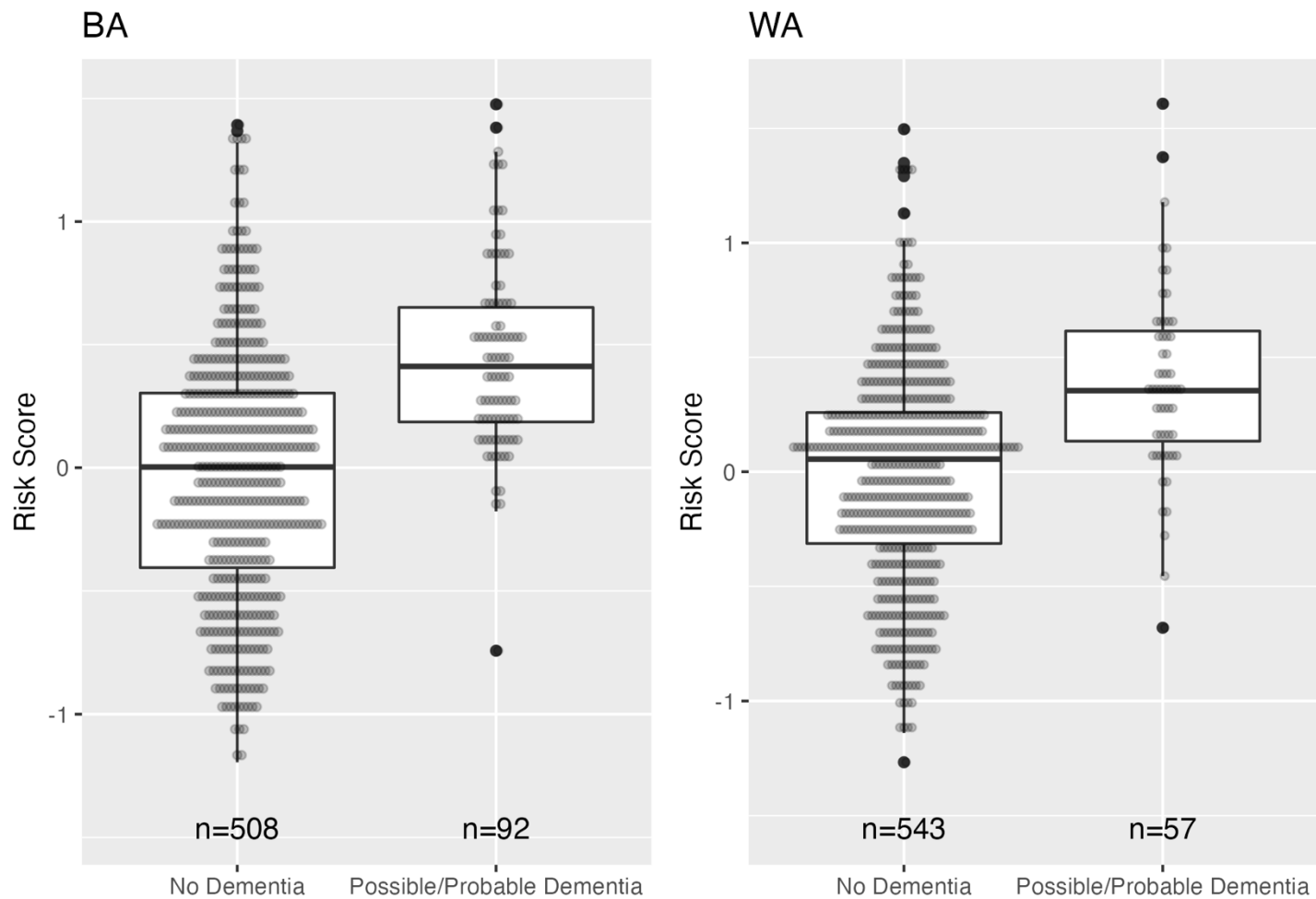


Figure 3a: Distribution of risk scores by dementia status and race (BA: Black American, WA: White American) for Veterans in in the Validation Sample who underwent chart review (n = 600 for each race).

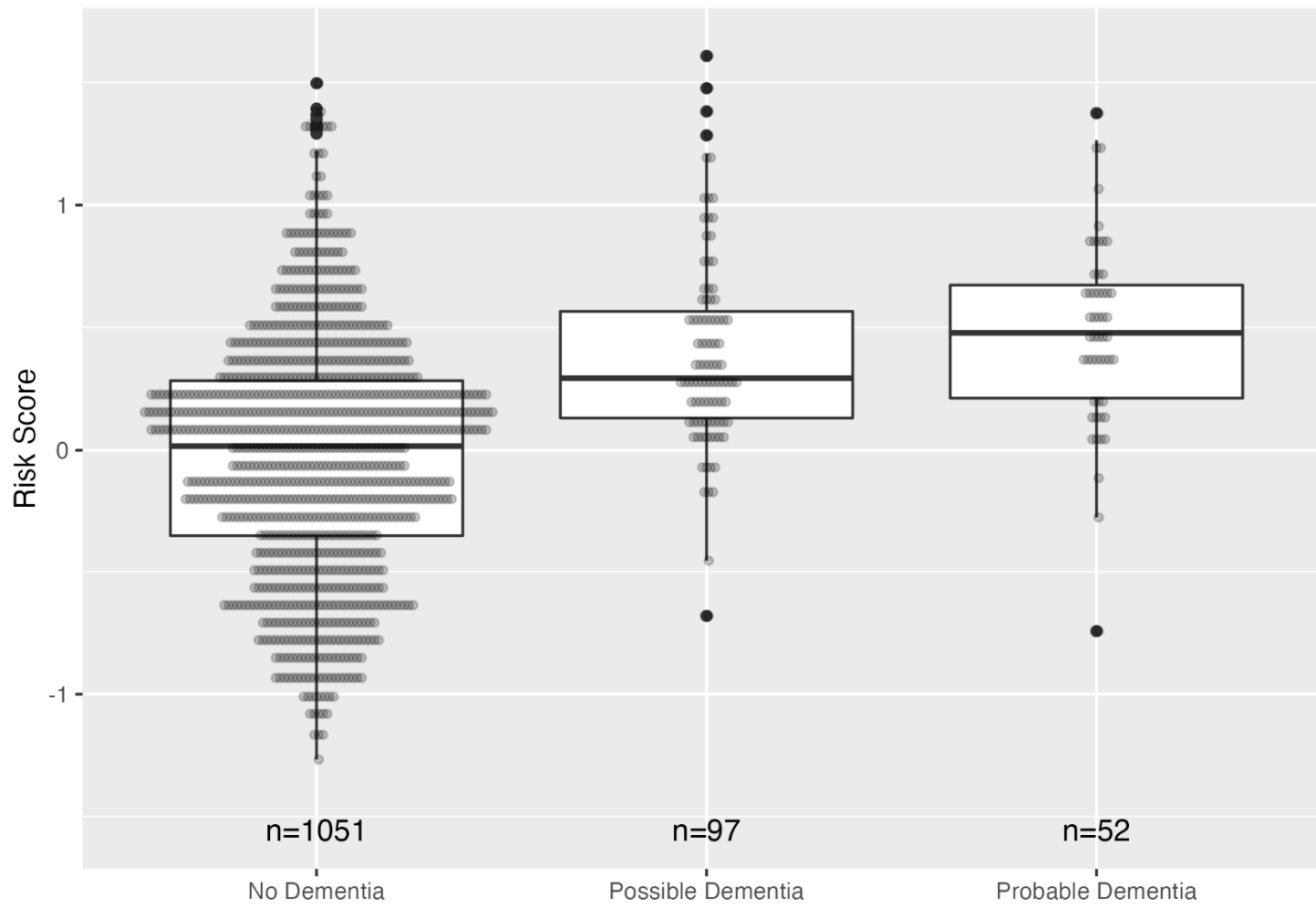


Figure 3b. Distribution of risk scores by dementia status for both Black American and White American Veterans in the Validation Sample who underwent chart review (n = 600 for each race).

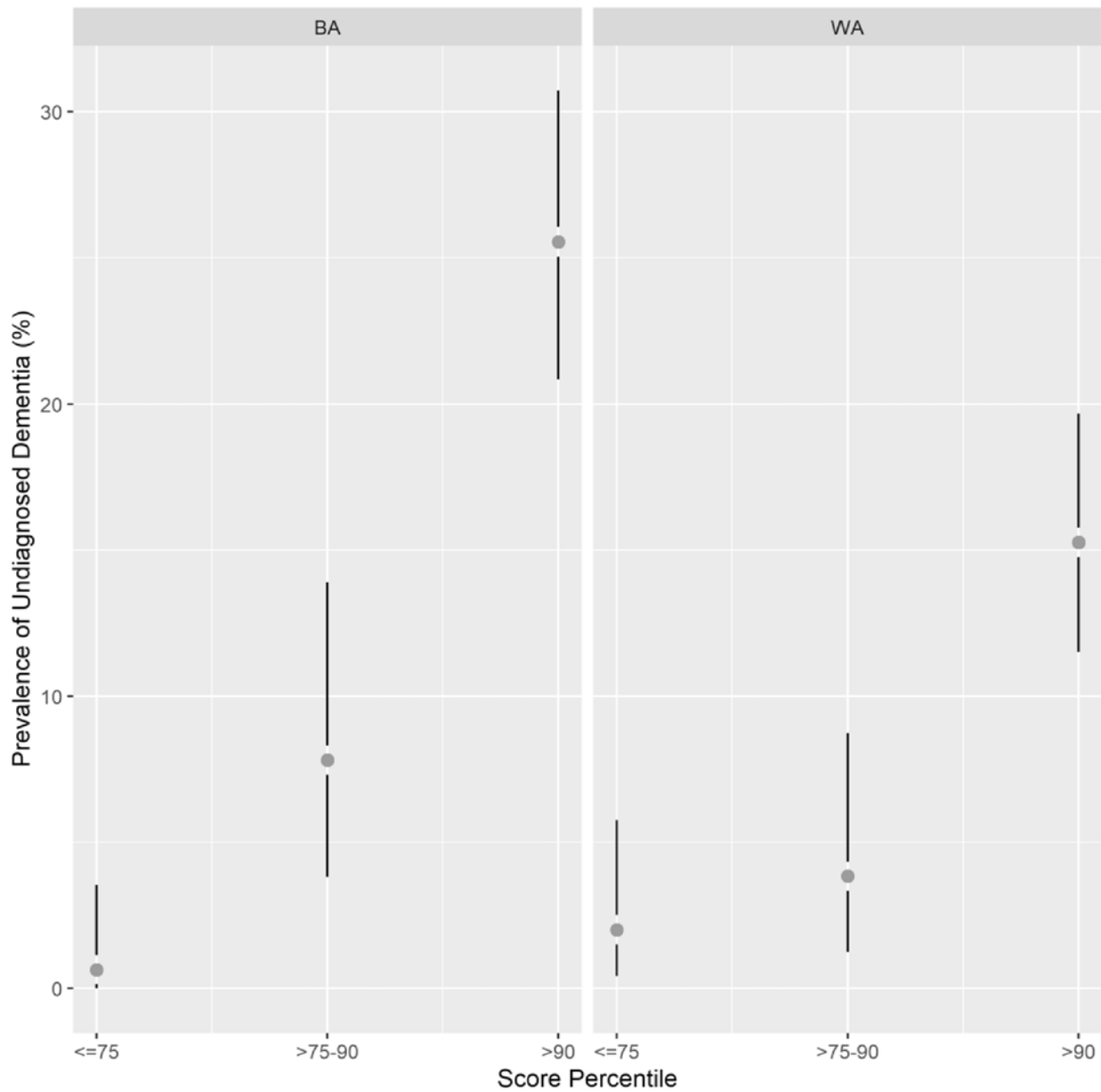


Figure 4. Prevalence of undiagnosed dementia by score percentile stratum and race (BA: Black American; WA: White American) for Veterans who underwent chart review (n = 600 for each race). For each race, score percentiles are based on using the scores from all 10,000 Veterans in the Validation Sample.

